# Introduction to word2vec

## Claire He

## 2023-09-16

```
## [1] "I went on a successful date with someone I felt sympathy and connection with."
## [2] "I was happy when my son got 90% marks in his examination"
## [3] "I went to the gym this morning and did yoga."
## [4] "We had a serious talk with some friends of ours who have been flaky lately. They understood and
## [5] "I went with grandchildren to butterfly display at Crohn Conservatory\n"
## [6] "I meditated last night."
```

```
##   wid  age country gender marital parenthood
## 1   1 37.0     USA      m married          y
## 2   2 29.0     IND      m married          y
## 3   3   25     IND      m  single          n
## 4   4   32     USA      m married          y
## 5   5   29     USA      m married          y
## 6   6   35     IND      m married          y
```

*Question*: Effect of animals on happiness across gender groups.

We are going to use `word2vec` to try to investigate this question. To illustrate this we are going to try to compare use of `word2vec` and `lda`.

First convert the data to a list of characters to input into our model. Choose the dimension of the embedding (arbitrarily chose 15 but this is a tuning parameter you can try a few dimensions and see what works best for you). You can also choose either models based on `cbow` or `skipgram`.

## Overall vocabulary between happiness and animals

```
x <- tolower(cleaned_data$cleaned_hm)
cat(x[1])
```

```
## i went on a successful date with someone i felt sympathy and connection with.
```

Lemmatizing our text and using speech tag (verb, adverb, noun, adjective) will make representation easier (let's say we want to see all adjectives and nouns relative to the topic of animals).

```
## 2023-09-16 21:13:46 Annotating text fragment 1/100392
## 2023-09-16 21:13:46 Annotating text fragment 11/100392
## 2023-09-16 21:13:46 Annotating text fragment 21/100392
## 2023-09-16 21:13:46 Annotating text fragment 31/100392
## 2023-09-16 21:13:46 Annotating text fragment 41/100392
## 2023-09-16 21:13:46 Annotating text fragment 51/100392
```

```
## 2023-09-16 21:13:47 Annotating text fragment 61/100392
## 2023-09-16 21:13:47 Annotating text fragment 71/100392
## 2023-09-16 21:13:47 Annotating text fragment 81/100392
## 2023-09-16 21:13:47 Annotating text fragment 91/100392
## 2023-09-16 21:13:47 Annotating text fragment 101/100392
## 2023-09-16 21:13:47 Annotating text fragment 111/100392
## 2023-09-16 21:13:47 Annotating text fragment 121/100392
## 2023-09-16 21:13:48 Annotating text fragment 131/100392
## 2023-09-16 21:13:48 Annotating text fragment 141/100392
## 2023-09-16 21:13:48 Annotating text fragment 151/100392
## 2023-09-16 21:13:48 Annotating text fragment 161/100392
## 2023-09-16 21:13:49 Annotating text fragment 171/100392
## 2023-09-16 21:13:49 Annotating text fragment 181/100392
## 2023-09-16 21:13:49 Annotating text fragment 191/100392
## 2023-09-16 21:13:49 Annotating text fragment 201/100392
## 2023-09-16 21:13:49 Annotating text fragment 211/100392
## 2023-09-16 21:13:49 Annotating text fragment 221/100392
## 2023-09-16 21:13:49 Annotating text fragment 231/100392
## 2023-09-16 21:13:49 Annotating text fragment 241/100392
## 2023-09-16 21:13:49 Annotating text fragment 251/100392
## 2023-09-16 21:13:49 Annotating text fragment 261/100392
## 2023-09-16 21:13:49 Annotating text fragment 271/100392
## 2023-09-16 21:13:50 Annotating text fragment 281/100392
## 2023-09-16 21:13:50 Annotating text fragment 291/100392
## 2023-09-16 21:13:50 Annotating text fragment 301/100392
## 2023-09-16 21:13:50 Annotating text fragment 311/100392
## 2023-09-16 21:13:50 Annotating text fragment 321/100392
## 2023-09-16 21:13:50 Annotating text fragment 331/100392
## 2023-09-16 21:13:50 Annotating text fragment 341/100392
## 2023-09-16 21:13:50 Annotating text fragment 351/100392
## 2023-09-16 21:13:50 Annotating text fragment 361/100392
## 2023-09-16 21:13:51 Annotating text fragment 371/100392
## 2023-09-16 21:13:51 Annotating text fragment 381/100392
## 2023-09-16 21:13:51 Annotating text fragment 391/100392
## 2023-09-16 21:13:51 Annotating text fragment 401/100392
## 2023-09-16 21:13:51 Annotating text fragment 411/100392
## 2023-09-16 21:13:51 Annotating text fragment 421/100392
## 2023-09-16 21:13:51 Annotating text fragment 431/100392
## 2023-09-16 21:13:51 Annotating text fragment 441/100392
## 2023-09-16 21:13:51 Annotating text fragment 451/100392
## 2023-09-16 21:13:51 Annotating text fragment 461/100392
## 2023-09-16 21:13:51 Annotating text fragment 471/100392
## 2023-09-16 21:13:52 Annotating text fragment 481/100392
## 2023-09-16 21:13:52 Annotating text fragment 491/100392
## 2023-09-16 21:13:52 Annotating text fragment 501/100392
## 2023-09-16 21:13:52 Annotating text fragment 511/100392
## 2023-09-16 21:13:52 Annotating text fragment 521/100392
## 2023-09-16 21:13:52 Annotating text fragment 531/100392
## 2023-09-16 21:13:52 Annotating text fragment 541/100392
## 2023-09-16 21:13:52 Annotating text fragment 551/100392
## 2023-09-16 21:13:52 Annotating text fragment 561/100392
## 2023-09-16 21:13:53 Annotating text fragment 571/100392
## 2023-09-16 21:13:53 Annotating text fragment 581/100392
## 2023-09-16 21:13:53 Annotating text fragment 591/100392
```

```
## 2023-09-16 21:13:53 Annotating text fragment 601/100392
## 2023-09-16 21:13:53 Annotating text fragment 611/100392
## 2023-09-16 21:13:53 Annotating text fragment 621/100392
## 2023-09-16 21:13:53 Annotating text fragment 631/100392
## 2023-09-16 21:13:53 Annotating text fragment 641/100392
## 2023-09-16 21:13:53 Annotating text fragment 651/100392
## 2023-09-16 21:13:53 Annotating text fragment 661/100392
## 2023-09-16 21:13:53 Annotating text fragment 671/100392
## 2023-09-16 21:13:54 Annotating text fragment 681/100392
## 2023-09-16 21:13:54 Annotating text fragment 691/100392
## 2023-09-16 21:13:54 Annotating text fragment 701/100392
## 2023-09-16 21:13:54 Annotating text fragment 711/100392
## 2023-09-16 21:13:54 Annotating text fragment 721/100392
## 2023-09-16 21:13:54 Annotating text fragment 731/100392
## 2023-09-16 21:13:54 Annotating text fragment 741/100392
## 2023-09-16 21:13:54 Annotating text fragment 751/100392
## 2023-09-16 21:13:54 Annotating text fragment 761/100392
## 2023-09-16 21:13:54 Annotating text fragment 771/100392
## 2023-09-16 21:13:55 Annotating text fragment 781/100392
## 2023-09-16 21:13:55 Annotating text fragment 791/100392
## 2023-09-16 21:13:55 Annotating text fragment 801/100392
## 2023-09-16 21:13:55 Annotating text fragment 811/100392
## 2023-09-16 21:13:55 Annotating text fragment 821/100392
## 2023-09-16 21:13:55 Annotating text fragment 831/100392
## 2023-09-16 21:13:55 Annotating text fragment 841/100392
## 2023-09-16 21:13:55 Annotating text fragment 851/100392
## 2023-09-16 21:13:55 Annotating text fragment 861/100392
## 2023-09-16 21:13:55 Annotating text fragment 871/100392
## 2023-09-16 21:13:56 Annotating text fragment 881/100392
## 2023-09-16 21:13:56 Annotating text fragment 891/100392
## 2023-09-16 21:13:56 Annotating text fragment 901/100392
## 2023-09-16 21:13:56 Annotating text fragment 911/100392
## 2023-09-16 21:13:56 Annotating text fragment 921/100392
## 2023-09-16 21:13:56 Annotating text fragment 931/100392
## 2023-09-16 21:13:56 Annotating text fragment 941/100392
## 2023-09-16 21:13:56 Annotating text fragment 951/100392
## 2023-09-16 21:13:56 Annotating text fragment 961/100392
## 2023-09-16 21:13:56 Annotating text fragment 971/100392
## 2023-09-16 21:13:57 Annotating text fragment 981/100392
## 2023-09-16 21:13:57 Annotating text fragment 991/100392
## 2023-09-16 21:13:57 Annotating text fragment 1001/100392
## 2023-09-16 21:13:57 Annotating text fragment 1011/100392
## 2023-09-16 21:13:57 Annotating text fragment 1021/100392
## 2023-09-16 21:13:57 Annotating text fragment 1031/100392
## 2023-09-16 21:13:57 Annotating text fragment 1041/100392
## 2023-09-16 21:13:57 Annotating text fragment 1051/100392
## 2023-09-16 21:13:57 Annotating text fragment 1061/100392
## 2023-09-16 21:13:57 Annotating text fragment 1071/100392
## 2023-09-16 21:13:57 Annotating text fragment 1081/100392
## 2023-09-16 21:13:57 Annotating text fragment 1091/100392
## 2023-09-16 21:13:58 Annotating text fragment 1101/100392
## 2023-09-16 21:13:58 Annotating text fragment 1111/100392
## 2023-09-16 21:13:58 Annotating text fragment 1121/100392
## 2023-09-16 21:13:58 Annotating text fragment 1131/100392
```

```
## 2023-09-16 21:13:58 Annotating text fragment 1141/100392
## 2023-09-16 21:13:58 Annotating text fragment 1151/100392
## 2023-09-16 21:13:58 Annotating text fragment 1161/100392
## 2023-09-16 21:13:58 Annotating text fragment 1171/100392
## 2023-09-16 21:13:58 Annotating text fragment 1181/100392
## 2023-09-16 21:13:58 Annotating text fragment 1191/100392
## 2023-09-16 21:13:58 Annotating text fragment 1201/100392
## 2023-09-16 21:13:58 Annotating text fragment 1211/100392
## 2023-09-16 21:13:58 Annotating text fragment 1221/100392
## 2023-09-16 21:13:59 Annotating text fragment 1231/100392
## 2023-09-16 21:13:59 Annotating text fragment 1241/100392
## 2023-09-16 21:13:59 Annotating text fragment 1251/100392
## 2023-09-16 21:13:59 Annotating text fragment 1261/100392
## 2023-09-16 21:13:59 Annotating text fragment 1271/100392
## 2023-09-16 21:13:59 Annotating text fragment 1281/100392
## 2023-09-16 21:13:59 Annotating text fragment 1291/100392
## 2023-09-16 21:13:59 Annotating text fragment 1301/100392
## 2023-09-16 21:13:59 Annotating text fragment 1311/100392
## 2023-09-16 21:13:59 Annotating text fragment 1321/100392
## 2023-09-16 21:14:00 Annotating text fragment 1331/100392
## 2023-09-16 21:14:00 Annotating text fragment 1341/100392
## 2023-09-16 21:14:00 Annotating text fragment 1351/100392
## 2023-09-16 21:14:00 Annotating text fragment 1361/100392
## 2023-09-16 21:14:00 Annotating text fragment 1371/100392
## 2023-09-16 21:14:00 Annotating text fragment 1381/100392
## 2023-09-16 21:14:00 Annotating text fragment 1391/100392
## 2023-09-16 21:14:00 Annotating text fragment 1401/100392
## 2023-09-16 21:14:01 Annotating text fragment 1411/100392
## 2023-09-16 21:14:01 Annotating text fragment 1421/100392
## 2023-09-16 21:14:02 Annotating text fragment 1431/100392
## 2023-09-16 21:14:02 Annotating text fragment 1441/100392
## 2023-09-16 21:14:02 Annotating text fragment 1451/100392
## 2023-09-16 21:14:02 Annotating text fragment 1461/100392
## 2023-09-16 21:14:02 Annotating text fragment 1471/100392
## 2023-09-16 21:14:02 Annotating text fragment 1481/100392
## 2023-09-16 21:14:02 Annotating text fragment 1491/100392
## 2023-09-16 21:14:02 Annotating text fragment 1501/100392
## 2023-09-16 21:14:02 Annotating text fragment 1511/100392
## 2023-09-16 21:14:02 Annotating text fragment 1521/100392
## 2023-09-16 21:14:02 Annotating text fragment 1531/100392
## 2023-09-16 21:14:03 Annotating text fragment 1541/100392
## 2023-09-16 21:14:03 Annotating text fragment 1551/100392
## 2023-09-16 21:14:03 Annotating text fragment 1561/100392
## 2023-09-16 21:14:03 Annotating text fragment 1571/100392
## 2023-09-16 21:14:03 Annotating text fragment 1581/100392
## 2023-09-16 21:14:03 Annotating text fragment 1591/100392
## 2023-09-16 21:14:03 Annotating text fragment 1601/100392
## 2023-09-16 21:14:03 Annotating text fragment 1611/100392
## 2023-09-16 21:14:03 Annotating text fragment 1621/100392
## 2023-09-16 21:14:03 Annotating text fragment 1631/100392
## 2023-09-16 21:14:03 Annotating text fragment 1641/100392
## 2023-09-16 21:14:03 Annotating text fragment 1651/100392
## 2023-09-16 21:14:04 Annotating text fragment 1661/100392
## 2023-09-16 21:14:04 Annotating text fragment 1671/100392
```

```
## 2023-09-16 21:14:04 Annotating text fragment 1681/100392
## 2023-09-16 21:14:04 Annotating text fragment 1691/100392
## 2023-09-16 21:14:04 Annotating text fragment 1701/100392
## 2023-09-16 21:14:04 Annotating text fragment 1711/100392
## 2023-09-16 21:14:04 Annotating text fragment 1721/100392
## 2023-09-16 21:14:04 Annotating text fragment 1731/100392
## 2023-09-16 21:14:04 Annotating text fragment 1741/100392
## 2023-09-16 21:14:04 Annotating text fragment 1751/100392
## 2023-09-16 21:14:04 Annotating text fragment 1761/100392
## 2023-09-16 21:14:04 Annotating text fragment 1771/100392
## 2023-09-16 21:14:04 Annotating text fragment 1781/100392
## 2023-09-16 21:14:05 Annotating text fragment 1791/100392
## 2023-09-16 21:14:05 Annotating text fragment 1801/100392
## 2023-09-16 21:14:05 Annotating text fragment 1811/100392
## 2023-09-16 21:14:05 Annotating text fragment 1821/100392
## 2023-09-16 21:14:05 Annotating text fragment 1831/100392
## 2023-09-16 21:14:05 Annotating text fragment 1841/100392
## 2023-09-16 21:14:05 Annotating text fragment 1851/100392
## 2023-09-16 21:14:05 Annotating text fragment 1861/100392
## 2023-09-16 21:14:05 Annotating text fragment 1871/100392
## 2023-09-16 21:14:05 Annotating text fragment 1881/100392
## 2023-09-16 21:14:05 Annotating text fragment 1891/100392
## 2023-09-16 21:14:06 Annotating text fragment 1901/100392
## 2023-09-16 21:14:06 Annotating text fragment 1911/100392
## 2023-09-16 21:14:06 Annotating text fragment 1921/100392
## 2023-09-16 21:14:06 Annotating text fragment 1931/100392
## 2023-09-16 21:14:06 Annotating text fragment 1941/100392
## 2023-09-16 21:14:06 Annotating text fragment 1951/100392
## 2023-09-16 21:14:06 Annotating text fragment 1961/100392
## 2023-09-16 21:14:06 Annotating text fragment 1971/100392
## 2023-09-16 21:14:06 Annotating text fragment 1981/100392
## 2023-09-16 21:14:06 Annotating text fragment 1991/100392
## 2023-09-16 21:14:06 Annotating text fragment 2001/100392
## 2023-09-16 21:14:06 Annotating text fragment 2011/100392
## 2023-09-16 21:14:07 Annotating text fragment 2021/100392
## 2023-09-16 21:14:07 Annotating text fragment 2031/100392
## 2023-09-16 21:14:07 Annotating text fragment 2041/100392
## 2023-09-16 21:14:07 Annotating text fragment 2051/100392
## 2023-09-16 21:14:07 Annotating text fragment 2061/100392
## 2023-09-16 21:14:07 Annotating text fragment 2071/100392
## 2023-09-16 21:14:07 Annotating text fragment 2081/100392
## 2023-09-16 21:14:07 Annotating text fragment 2091/100392
## 2023-09-16 21:14:07 Annotating text fragment 2101/100392
## 2023-09-16 21:14:07 Annotating text fragment 2111/100392
## 2023-09-16 21:14:07 Annotating text fragment 2121/100392
## 2023-09-16 21:14:07 Annotating text fragment 2131/100392
## 2023-09-16 21:14:08 Annotating text fragment 2141/100392
## 2023-09-16 21:14:08 Annotating text fragment 2151/100392
## 2023-09-16 21:14:08 Annotating text fragment 2161/100392
## 2023-09-16 21:14:08 Annotating text fragment 2171/100392
## 2023-09-16 21:14:08 Annotating text fragment 2181/100392
## 2023-09-16 21:14:08 Annotating text fragment 2191/100392
## 2023-09-16 21:14:08 Annotating text fragment 2201/100392
## 2023-09-16 21:14:08 Annotating text fragment 2211/100392
```

```
## 2023-09-16 21:14:09 Annotating text fragment 2221/100392
## 2023-09-16 21:14:09 Annotating text fragment 2231/100392
## 2023-09-16 21:14:09 Annotating text fragment 2241/100392
## 2023-09-16 21:14:09 Annotating text fragment 2251/100392
## 2023-09-16 21:14:09 Annotating text fragment 2261/100392
## 2023-09-16 21:14:09 Annotating text fragment 2271/100392
## 2023-09-16 21:14:09 Annotating text fragment 2281/100392
## 2023-09-16 21:14:10 Annotating text fragment 2291/100392
## 2023-09-16 21:14:10 Annotating text fragment 2301/100392
## 2023-09-16 21:14:10 Annotating text fragment 2311/100392
## 2023-09-16 21:14:10 Annotating text fragment 2321/100392
## 2023-09-16 21:14:10 Annotating text fragment 2331/100392
## 2023-09-16 21:14:11 Annotating text fragment 2341/100392
## 2023-09-16 21:14:11 Annotating text fragment 2351/100392
## 2023-09-16 21:14:11 Annotating text fragment 2361/100392
## 2023-09-16 21:14:11 Annotating text fragment 2371/100392
## 2023-09-16 21:14:11 Annotating text fragment 2381/100392
## 2023-09-16 21:14:11 Annotating text fragment 2391/100392
## 2023-09-16 21:14:12 Annotating text fragment 2401/100392
## 2023-09-16 21:14:12 Annotating text fragment 2411/100392
## 2023-09-16 21:14:12 Annotating text fragment 2421/100392
## 2023-09-16 21:14:12 Annotating text fragment 2431/100392
## 2023-09-16 21:14:12 Annotating text fragment 2441/100392
## 2023-09-16 21:14:12 Annotating text fragment 2451/100392
## 2023-09-16 21:14:12 Annotating text fragment 2461/100392
## 2023-09-16 21:14:12 Annotating text fragment 2471/100392
## 2023-09-16 21:14:12 Annotating text fragment 2481/100392
## 2023-09-16 21:14:12 Annotating text fragment 2491/100392
## 2023-09-16 21:14:12 Annotating text fragment 2501/100392
## 2023-09-16 21:14:13 Annotating text fragment 2511/100392
## 2023-09-16 21:14:13 Annotating text fragment 2521/100392
## 2023-09-16 21:14:13 Annotating text fragment 2531/100392
## 2023-09-16 21:14:13 Annotating text fragment 2541/100392
## 2023-09-16 21:14:13 Annotating text fragment 2551/100392
## 2023-09-16 21:14:13 Annotating text fragment 2561/100392
## 2023-09-16 21:14:13 Annotating text fragment 2571/100392
## 2023-09-16 21:14:13 Annotating text fragment 2581/100392
## 2023-09-16 21:14:13 Annotating text fragment 2591/100392
## 2023-09-16 21:14:13 Annotating text fragment 2601/100392
## 2023-09-16 21:14:14 Annotating text fragment 2611/100392
## 2023-09-16 21:14:14 Annotating text fragment 2621/100392
## 2023-09-16 21:14:14 Annotating text fragment 2631/100392
## 2023-09-16 21:14:14 Annotating text fragment 2641/100392
## 2023-09-16 21:14:14 Annotating text fragment 2651/100392
## 2023-09-16 21:14:14 Annotating text fragment 2661/100392
## 2023-09-16 21:14:14 Annotating text fragment 2671/100392
## 2023-09-16 21:14:14 Annotating text fragment 2681/100392
## 2023-09-16 21:14:15 Annotating text fragment 2691/100392
## 2023-09-16 21:14:15 Annotating text fragment 2701/100392
## 2023-09-16 21:14:15 Annotating text fragment 2711/100392
## 2023-09-16 21:14:15 Annotating text fragment 2721/100392
## 2023-09-16 21:14:15 Annotating text fragment 2731/100392
## 2023-09-16 21:14:15 Annotating text fragment 2741/100392
## 2023-09-16 21:14:15 Annotating text fragment 2751/100392
```

```
## 2023-09-16 21:14:16 Annotating text fragment 2761/100392
## 2023-09-16 21:14:16 Annotating text fragment 2771/100392
## 2023-09-16 21:14:16 Annotating text fragment 2781/100392
## 2023-09-16 21:14:16 Annotating text fragment 2791/100392
## 2023-09-16 21:14:16 Annotating text fragment 2801/100392
## 2023-09-16 21:14:16 Annotating text fragment 2811/100392
## 2023-09-16 21:14:16 Annotating text fragment 2821/100392
## 2023-09-16 21:14:16 Annotating text fragment 2831/100392
## 2023-09-16 21:14:17 Annotating text fragment 2841/100392
## 2023-09-16 21:14:17 Annotating text fragment 2851/100392
## 2023-09-16 21:14:17 Annotating text fragment 2861/100392
## 2023-09-16 21:14:17 Annotating text fragment 2871/100392
## 2023-09-16 21:14:17 Annotating text fragment 2881/100392
## 2023-09-16 21:14:17 Annotating text fragment 2891/100392
## 2023-09-16 21:14:17 Annotating text fragment 2901/100392
## 2023-09-16 21:14:17 Annotating text fragment 2911/100392
## 2023-09-16 21:14:18 Annotating text fragment 2921/100392
## 2023-09-16 21:14:18 Annotating text fragment 2931/100392
## 2023-09-16 21:14:18 Annotating text fragment 2941/100392
## 2023-09-16 21:14:18 Annotating text fragment 2951/100392
## 2023-09-16 21:14:18 Annotating text fragment 2961/100392
## 2023-09-16 21:14:18 Annotating text fragment 2971/100392
## 2023-09-16 21:14:18 Annotating text fragment 2981/100392
## 2023-09-16 21:14:18 Annotating text fragment 2991/100392
## 2023-09-16 21:14:18 Annotating text fragment 3001/100392
## 2023-09-16 21:14:19 Annotating text fragment 3011/100392
## 2023-09-16 21:14:19 Annotating text fragment 3021/100392
## 2023-09-16 21:14:19 Annotating text fragment 3031/100392
## 2023-09-16 21:14:19 Annotating text fragment 3041/100392
## 2023-09-16 21:14:19 Annotating text fragment 3051/100392
## 2023-09-16 21:14:19 Annotating text fragment 3061/100392
## 2023-09-16 21:14:19 Annotating text fragment 3071/100392
## 2023-09-16 21:14:20 Annotating text fragment 3081/100392
## 2023-09-16 21:14:20 Annotating text fragment 3091/100392
## 2023-09-16 21:14:20 Annotating text fragment 3101/100392
## 2023-09-16 21:14:20 Annotating text fragment 3111/100392
## 2023-09-16 21:14:20 Annotating text fragment 3121/100392
## 2023-09-16 21:14:20 Annotating text fragment 3131/100392
## 2023-09-16 21:14:20 Annotating text fragment 3141/100392
## 2023-09-16 21:14:20 Annotating text fragment 3151/100392
## 2023-09-16 21:14:20 Annotating text fragment 3161/100392
## 2023-09-16 21:14:20 Annotating text fragment 3171/100392
## 2023-09-16 21:14:21 Annotating text fragment 3181/100392
## 2023-09-16 21:14:21 Annotating text fragment 3191/100392
## 2023-09-16 21:14:21 Annotating text fragment 3201/100392
## 2023-09-16 21:14:21 Annotating text fragment 3211/100392
## 2023-09-16 21:14:21 Annotating text fragment 3221/100392
## 2023-09-16 21:14:21 Annotating text fragment 3231/100392
## 2023-09-16 21:14:21 Annotating text fragment 3241/100392
## 2023-09-16 21:14:21 Annotating text fragment 3251/100392
## 2023-09-16 21:14:22 Annotating text fragment 3261/100392
## 2023-09-16 21:14:22 Annotating text fragment 3271/100392
## 2023-09-16 21:14:22 Annotating text fragment 3281/100392
## 2023-09-16 21:14:22 Annotating text fragment 3291/100392
```

```
## 2023-09-16 21:14:22 Annotating text fragment 3301/100392
## 2023-09-16 21:14:22 Annotating text fragment 3311/100392
## 2023-09-16 21:14:22 Annotating text fragment 3321/100392
## 2023-09-16 21:14:22 Annotating text fragment 3331/100392
## 2023-09-16 21:14:23 Annotating text fragment 3341/100392
## 2023-09-16 21:14:23 Annotating text fragment 3351/100392
## 2023-09-16 21:14:23 Annotating text fragment 3361/100392
## 2023-09-16 21:14:23 Annotating text fragment 3371/100392
## 2023-09-16 21:14:23 Annotating text fragment 3381/100392
## 2023-09-16 21:14:23 Annotating text fragment 3391/100392
## 2023-09-16 21:14:23 Annotating text fragment 3401/100392
## 2023-09-16 21:14:24 Annotating text fragment 3411/100392
## 2023-09-16 21:14:24 Annotating text fragment 3421/100392
## 2023-09-16 21:14:24 Annotating text fragment 3431/100392
## 2023-09-16 21:14:24 Annotating text fragment 3441/100392
## 2023-09-16 21:14:24 Annotating text fragment 3451/100392
## 2023-09-16 21:14:24 Annotating text fragment 3461/100392
## 2023-09-16 21:14:24 Annotating text fragment 3471/100392
## 2023-09-16 21:14:24 Annotating text fragment 3481/100392
## 2023-09-16 21:14:24 Annotating text fragment 3491/100392
## 2023-09-16 21:14:25 Annotating text fragment 3501/100392
## 2023-09-16 21:14:25 Annotating text fragment 3511/100392
## 2023-09-16 21:14:25 Annotating text fragment 3521/100392
## 2023-09-16 21:14:25 Annotating text fragment 3531/100392
## 2023-09-16 21:14:25 Annotating text fragment 3541/100392
## 2023-09-16 21:14:25 Annotating text fragment 3551/100392
## 2023-09-16 21:14:25 Annotating text fragment 3561/100392
## 2023-09-16 21:14:25 Annotating text fragment 3571/100392
## 2023-09-16 21:14:25 Annotating text fragment 3581/100392
## 2023-09-16 21:14:25 Annotating text fragment 3591/100392
## 2023-09-16 21:14:25 Annotating text fragment 3601/100392
## 2023-09-16 21:14:26 Annotating text fragment 3611/100392
## 2023-09-16 21:14:26 Annotating text fragment 3621/100392
## 2023-09-16 21:14:26 Annotating text fragment 3631/100392
## 2023-09-16 21:14:26 Annotating text fragment 3641/100392
## 2023-09-16 21:14:26 Annotating text fragment 3651/100392
## 2023-09-16 21:14:26 Annotating text fragment 3661/100392
## 2023-09-16 21:14:26 Annotating text fragment 3671/100392
## 2023-09-16 21:14:26 Annotating text fragment 3681/100392
## 2023-09-16 21:14:27 Annotating text fragment 3691/100392
## 2023-09-16 21:14:27 Annotating text fragment 3701/100392
## 2023-09-16 21:14:27 Annotating text fragment 3711/100392
## 2023-09-16 21:14:27 Annotating text fragment 3721/100392
## 2023-09-16 21:14:27 Annotating text fragment 3731/100392
## 2023-09-16 21:14:27 Annotating text fragment 3741/100392
## 2023-09-16 21:14:27 Annotating text fragment 3751/100392
## 2023-09-16 21:14:27 Annotating text fragment 3761/100392
## 2023-09-16 21:14:27 Annotating text fragment 3771/100392
## 2023-09-16 21:14:28 Annotating text fragment 3781/100392
## 2023-09-16 21:14:28 Annotating text fragment 3791/100392
## 2023-09-16 21:14:28 Annotating text fragment 3801/100392
## 2023-09-16 21:14:28 Annotating text fragment 3811/100392
## 2023-09-16 21:14:28 Annotating text fragment 3821/100392
## 2023-09-16 21:14:28 Annotating text fragment 3831/100392
```

```
## 2023-09-16 21:14:28 Annotating text fragment 3841/100392
## 2023-09-16 21:14:28 Annotating text fragment 3851/100392
## 2023-09-16 21:14:28 Annotating text fragment 3861/100392
## 2023-09-16 21:14:29 Annotating text fragment 3871/100392
## 2023-09-16 21:14:29 Annotating text fragment 3881/100392
## 2023-09-16 21:14:29 Annotating text fragment 3891/100392
## 2023-09-16 21:14:29 Annotating text fragment 3901/100392
## 2023-09-16 21:14:29 Annotating text fragment 3911/100392
## 2023-09-16 21:14:29 Annotating text fragment 3921/100392
## 2023-09-16 21:14:29 Annotating text fragment 3931/100392
## 2023-09-16 21:14:29 Annotating text fragment 3941/100392
## 2023-09-16 21:14:29 Annotating text fragment 3951/100392
## 2023-09-16 21:14:29 Annotating text fragment 3961/100392
## 2023-09-16 21:14:30 Annotating text fragment 3971/100392
## 2023-09-16 21:14:30 Annotating text fragment 3981/100392
## 2023-09-16 21:14:30 Annotating text fragment 3991/100392
## 2023-09-16 21:14:30 Annotating text fragment 4001/100392
## 2023-09-16 21:14:30 Annotating text fragment 4011/100392
## 2023-09-16 21:14:30 Annotating text fragment 4021/100392
## 2023-09-16 21:14:30 Annotating text fragment 4031/100392
## 2023-09-16 21:14:30 Annotating text fragment 4041/100392
## 2023-09-16 21:14:31 Annotating text fragment 4051/100392
## 2023-09-16 21:14:31 Annotating text fragment 4061/100392
## 2023-09-16 21:14:31 Annotating text fragment 4071/100392
## 2023-09-16 21:14:31 Annotating text fragment 4081/100392
## 2023-09-16 21:14:31 Annotating text fragment 4091/100392
## 2023-09-16 21:14:31 Annotating text fragment 4101/100392
## 2023-09-16 21:14:31 Annotating text fragment 4111/100392
## 2023-09-16 21:14:31 Annotating text fragment 4121/100392
## 2023-09-16 21:14:31 Annotating text fragment 4131/100392
## 2023-09-16 21:14:32 Annotating text fragment 4141/100392
## 2023-09-16 21:14:32 Annotating text fragment 4151/100392
## 2023-09-16 21:14:32 Annotating text fragment 4161/100392
## 2023-09-16 21:14:32 Annotating text fragment 4171/100392
## 2023-09-16 21:14:32 Annotating text fragment 4181/100392
## 2023-09-16 21:14:32 Annotating text fragment 4191/100392
## 2023-09-16 21:14:32 Annotating text fragment 4201/100392
## 2023-09-16 21:14:32 Annotating text fragment 4211/100392
## 2023-09-16 21:14:32 Annotating text fragment 4221/100392
## 2023-09-16 21:14:33 Annotating text fragment 4231/100392
## 2023-09-16 21:14:33 Annotating text fragment 4241/100392
## 2023-09-16 21:14:33 Annotating text fragment 4251/100392
## 2023-09-16 21:14:33 Annotating text fragment 4261/100392
## 2023-09-16 21:14:33 Annotating text fragment 4271/100392
## 2023-09-16 21:14:33 Annotating text fragment 4281/100392
## 2023-09-16 21:14:33 Annotating text fragment 4291/100392
## 2023-09-16 21:14:33 Annotating text fragment 4301/100392
## 2023-09-16 21:14:34 Annotating text fragment 4311/100392
## 2023-09-16 21:14:34 Annotating text fragment 4321/100392
## 2023-09-16 21:14:34 Annotating text fragment 4331/100392
## 2023-09-16 21:14:34 Annotating text fragment 4341/100392
## 2023-09-16 21:14:34 Annotating text fragment 4351/100392
## 2023-09-16 21:14:35 Annotating text fragment 4361/100392
## 2023-09-16 21:14:35 Annotating text fragment 4371/100392
```

```
## 2023-09-16 21:14:35 Annotating text fragment 4381/100392
## 2023-09-16 21:14:35 Annotating text fragment 4391/100392
## 2023-09-16 21:14:35 Annotating text fragment 4401/100392
## 2023-09-16 21:14:35 Annotating text fragment 4411/100392
## 2023-09-16 21:14:35 Annotating text fragment 4421/100392
## 2023-09-16 21:14:35 Annotating text fragment 4431/100392
## 2023-09-16 21:14:35 Annotating text fragment 4441/100392
## 2023-09-16 21:14:35 Annotating text fragment 4451/100392
## 2023-09-16 21:14:36 Annotating text fragment 4461/100392
## 2023-09-16 21:14:36 Annotating text fragment 4471/100392
## 2023-09-16 21:14:36 Annotating text fragment 4481/100392
## 2023-09-16 21:14:36 Annotating text fragment 4491/100392
## 2023-09-16 21:14:36 Annotating text fragment 4501/100392
## 2023-09-16 21:14:36 Annotating text fragment 4511/100392
## 2023-09-16 21:14:36 Annotating text fragment 4521/100392
## 2023-09-16 21:14:36 Annotating text fragment 4531/100392
## 2023-09-16 21:14:36 Annotating text fragment 4541/100392
## 2023-09-16 21:14:36 Annotating text fragment 4551/100392
## 2023-09-16 21:14:37 Annotating text fragment 4561/100392
## 2023-09-16 21:14:37 Annotating text fragment 4571/100392
## 2023-09-16 21:14:37 Annotating text fragment 4581/100392
## 2023-09-16 21:14:37 Annotating text fragment 4591/100392
## 2023-09-16 21:14:37 Annotating text fragment 4601/100392
## 2023-09-16 21:14:37 Annotating text fragment 4611/100392
## 2023-09-16 21:14:37 Annotating text fragment 4621/100392
## 2023-09-16 21:14:37 Annotating text fragment 4631/100392
## 2023-09-16 21:14:37 Annotating text fragment 4641/100392
## 2023-09-16 21:14:37 Annotating text fragment 4651/100392
## 2023-09-16 21:14:38 Annotating text fragment 4661/100392
## 2023-09-16 21:14:38 Annotating text fragment 4671/100392
## 2023-09-16 21:14:38 Annotating text fragment 4681/100392
## 2023-09-16 21:14:38 Annotating text fragment 4691/100392
## 2023-09-16 21:14:38 Annotating text fragment 4701/100392
## 2023-09-16 21:14:38 Annotating text fragment 4711/100392
## 2023-09-16 21:14:38 Annotating text fragment 4721/100392
## 2023-09-16 21:14:38 Annotating text fragment 4731/100392
## 2023-09-16 21:14:39 Annotating text fragment 4741/100392
## 2023-09-16 21:14:39 Annotating text fragment 4751/100392
## 2023-09-16 21:14:39 Annotating text fragment 4761/100392
## 2023-09-16 21:14:39 Annotating text fragment 4771/100392
## 2023-09-16 21:14:39 Annotating text fragment 4781/100392
## 2023-09-16 21:14:39 Annotating text fragment 4791/100392
## 2023-09-16 21:14:39 Annotating text fragment 4801/100392
## 2023-09-16 21:14:39 Annotating text fragment 4811/100392
## 2023-09-16 21:14:39 Annotating text fragment 4821/100392
## 2023-09-16 21:14:40 Annotating text fragment 4831/100392
## 2023-09-16 21:14:40 Annotating text fragment 4841/100392
## 2023-09-16 21:14:40 Annotating text fragment 4851/100392
## 2023-09-16 21:14:40 Annotating text fragment 4861/100392
## 2023-09-16 21:14:40 Annotating text fragment 4871/100392
## 2023-09-16 21:14:40 Annotating text fragment 4881/100392
## 2023-09-16 21:14:40 Annotating text fragment 4891/100392
## 2023-09-16 21:14:40 Annotating text fragment 4901/100392
## 2023-09-16 21:14:40 Annotating text fragment 4911/100392
```

```
## 2023-09-16 21:14:41 Annotating text fragment 4921/100392
## 2023-09-16 21:14:41 Annotating text fragment 4931/100392
## 2023-09-16 21:14:41 Annotating text fragment 4941/100392
## 2023-09-16 21:14:41 Annotating text fragment 4951/100392
## 2023-09-16 21:14:41 Annotating text fragment 4961/100392
## 2023-09-16 21:14:41 Annotating text fragment 4971/100392
## 2023-09-16 21:14:41 Annotating text fragment 4981/100392
## 2023-09-16 21:14:41 Annotating text fragment 4991/100392
## 2023-09-16 21:14:41 Annotating text fragment 5001/100392
## 2023-09-16 21:14:41 Annotating text fragment 5011/100392
## 2023-09-16 21:14:42 Annotating text fragment 5021/100392
## 2023-09-16 21:14:42 Annotating text fragment 5031/100392
## 2023-09-16 21:14:42 Annotating text fragment 5041/100392
## 2023-09-16 21:14:42 Annotating text fragment 5051/100392
## 2023-09-16 21:14:42 Annotating text fragment 5061/100392
## 2023-09-16 21:14:42 Annotating text fragment 5071/100392
## 2023-09-16 21:14:42 Annotating text fragment 5081/100392
## 2023-09-16 21:14:42 Annotating text fragment 5091/100392
## 2023-09-16 21:14:43 Annotating text fragment 5101/100392
## 2023-09-16 21:14:43 Annotating text fragment 5111/100392
## 2023-09-16 21:14:43 Annotating text fragment 5121/100392
## 2023-09-16 21:14:43 Annotating text fragment 5131/100392
## 2023-09-16 21:14:43 Annotating text fragment 5141/100392
## 2023-09-16 21:14:43 Annotating text fragment 5151/100392
## 2023-09-16 21:14:43 Annotating text fragment 5161/100392
## 2023-09-16 21:14:43 Annotating text fragment 5171/100392
## 2023-09-16 21:14:43 Annotating text fragment 5181/100392
## 2023-09-16 21:14:44 Annotating text fragment 5191/100392
## 2023-09-16 21:14:44 Annotating text fragment 5201/100392
## 2023-09-16 21:14:44 Annotating text fragment 5211/100392
## 2023-09-16 21:14:44 Annotating text fragment 5221/100392
## 2023-09-16 21:14:44 Annotating text fragment 5231/100392
## 2023-09-16 21:14:44 Annotating text fragment 5241/100392
## 2023-09-16 21:14:44 Annotating text fragment 5251/100392
## 2023-09-16 21:14:44 Annotating text fragment 5261/100392
## 2023-09-16 21:14:45 Annotating text fragment 5271/100392
## 2023-09-16 21:14:45 Annotating text fragment 5281/100392
## 2023-09-16 21:14:45 Annotating text fragment 5291/100392
## 2023-09-16 21:14:45 Annotating text fragment 5301/100392
## 2023-09-16 21:14:45 Annotating text fragment 5311/100392
## 2023-09-16 21:14:45 Annotating text fragment 5321/100392
## 2023-09-16 21:14:45 Annotating text fragment 5331/100392
## 2023-09-16 21:14:45 Annotating text fragment 5341/100392
## 2023-09-16 21:14:45 Annotating text fragment 5351/100392
## 2023-09-16 21:14:46 Annotating text fragment 5361/100392
## 2023-09-16 21:14:46 Annotating text fragment 5371/100392
## 2023-09-16 21:14:46 Annotating text fragment 5381/100392
## 2023-09-16 21:14:46 Annotating text fragment 5391/100392
## 2023-09-16 21:14:46 Annotating text fragment 5401/100392
## 2023-09-16 21:14:46 Annotating text fragment 5411/100392
## 2023-09-16 21:14:46 Annotating text fragment 5421/100392
## 2023-09-16 21:14:46 Annotating text fragment 5431/100392
## 2023-09-16 21:14:47 Annotating text fragment 5441/100392
## 2023-09-16 21:14:47 Annotating text fragment 5451/100392
```

```
## 2023-09-16 21:14:47 Annotating text fragment 5461/100392
## 2023-09-16 21:14:47 Annotating text fragment 5471/100392
## 2023-09-16 21:14:47 Annotating text fragment 5481/100392
## 2023-09-16 21:14:47 Annotating text fragment 5491/100392
## 2023-09-16 21:14:47 Annotating text fragment 5501/100392
## 2023-09-16 21:14:47 Annotating text fragment 5511/100392
## 2023-09-16 21:14:47 Annotating text fragment 5521/100392
## 2023-09-16 21:14:48 Annotating text fragment 5531/100392
## 2023-09-16 21:14:48 Annotating text fragment 5541/100392
## 2023-09-16 21:14:48 Annotating text fragment 5551/100392
## 2023-09-16 21:14:48 Annotating text fragment 5561/100392
## 2023-09-16 21:14:48 Annotating text fragment 5571/100392
## 2023-09-16 21:14:48 Annotating text fragment 5581/100392
## 2023-09-16 21:14:48 Annotating text fragment 5591/100392
## 2023-09-16 21:14:48 Annotating text fragment 5601/100392
## 2023-09-16 21:14:49 Annotating text fragment 5611/100392
## 2023-09-16 21:14:49 Annotating text fragment 5621/100392
## 2023-09-16 21:14:49 Annotating text fragment 5631/100392
## 2023-09-16 21:14:49 Annotating text fragment 5641/100392
## 2023-09-16 21:14:49 Annotating text fragment 5651/100392
## 2023-09-16 21:14:49 Annotating text fragment 5661/100392
## 2023-09-16 21:14:49 Annotating text fragment 5671/100392
## 2023-09-16 21:14:49 Annotating text fragment 5681/100392
## 2023-09-16 21:14:49 Annotating text fragment 5691/100392
## 2023-09-16 21:14:50 Annotating text fragment 5701/100392
## 2023-09-16 21:14:50 Annotating text fragment 5711/100392
## 2023-09-16 21:14:50 Annotating text fragment 5721/100392
## 2023-09-16 21:14:50 Annotating text fragment 5731/100392
## 2023-09-16 21:14:50 Annotating text fragment 5741/100392
## 2023-09-16 21:14:50 Annotating text fragment 5751/100392
## 2023-09-16 21:14:50 Annotating text fragment 5761/100392
## 2023-09-16 21:14:50 Annotating text fragment 5771/100392
## 2023-09-16 21:14:50 Annotating text fragment 5781/100392
## 2023-09-16 21:14:50 Annotating text fragment 5791/100392
## 2023-09-16 21:14:50 Annotating text fragment 5801/100392
## 2023-09-16 21:14:51 Annotating text fragment 5811/100392
## 2023-09-16 21:14:51 Annotating text fragment 5821/100392
## 2023-09-16 21:14:51 Annotating text fragment 5831/100392
## 2023-09-16 21:14:51 Annotating text fragment 5841/100392
## 2023-09-16 21:14:51 Annotating text fragment 5851/100392
## 2023-09-16 21:14:51 Annotating text fragment 5861/100392
## 2023-09-16 21:14:51 Annotating text fragment 5871/100392
## 2023-09-16 21:14:51 Annotating text fragment 5881/100392
## 2023-09-16 21:14:52 Annotating text fragment 5891/100392
## 2023-09-16 21:14:52 Annotating text fragment 5901/100392
## 2023-09-16 21:14:52 Annotating text fragment 5911/100392
## 2023-09-16 21:14:52 Annotating text fragment 5921/100392
## 2023-09-16 21:14:52 Annotating text fragment 5931/100392
## 2023-09-16 21:14:52 Annotating text fragment 5941/100392
## 2023-09-16 21:14:52 Annotating text fragment 5951/100392
## 2023-09-16 21:14:52 Annotating text fragment 5961/100392
## 2023-09-16 21:14:52 Annotating text fragment 5971/100392
## 2023-09-16 21:14:53 Annotating text fragment 5981/100392
## 2023-09-16 21:14:53 Annotating text fragment 5991/100392
```

```
## 2023-09-16 21:14:53 Annotating text fragment 6001/100392
## 2023-09-16 21:14:53 Annotating text fragment 6011/100392
## 2023-09-16 21:14:53 Annotating text fragment 6021/100392
## 2023-09-16 21:14:53 Annotating text fragment 6031/100392
## 2023-09-16 21:14:53 Annotating text fragment 6041/100392
## 2023-09-16 21:14:53 Annotating text fragment 6051/100392
## 2023-09-16 21:14:53 Annotating text fragment 6061/100392
## 2023-09-16 21:14:54 Annotating text fragment 6071/100392
## 2023-09-16 21:14:54 Annotating text fragment 6081/100392
## 2023-09-16 21:14:54 Annotating text fragment 6091/100392
## 2023-09-16 21:14:54 Annotating text fragment 6101/100392
## 2023-09-16 21:14:54 Annotating text fragment 6111/100392
## 2023-09-16 21:14:54 Annotating text fragment 6121/100392
## 2023-09-16 21:14:54 Annotating text fragment 6131/100392
## 2023-09-16 21:14:54 Annotating text fragment 6141/100392
## 2023-09-16 21:14:54 Annotating text fragment 6151/100392
## 2023-09-16 21:14:55 Annotating text fragment 6161/100392
## 2023-09-16 21:14:55 Annotating text fragment 6171/100392
## 2023-09-16 21:14:55 Annotating text fragment 6181/100392
## 2023-09-16 21:14:55 Annotating text fragment 6191/100392
## 2023-09-16 21:14:55 Annotating text fragment 6201/100392
## 2023-09-16 21:14:55 Annotating text fragment 6211/100392
## 2023-09-16 21:14:55 Annotating text fragment 6221/100392
## 2023-09-16 21:14:55 Annotating text fragment 6231/100392
## 2023-09-16 21:14:55 Annotating text fragment 6241/100392
## 2023-09-16 21:14:55 Annotating text fragment 6251/100392
## 2023-09-16 21:14:56 Annotating text fragment 6261/100392
## 2023-09-16 21:14:56 Annotating text fragment 6271/100392
## 2023-09-16 21:14:56 Annotating text fragment 6281/100392
## 2023-09-16 21:14:56 Annotating text fragment 6291/100392
## 2023-09-16 21:14:56 Annotating text fragment 6301/100392
## 2023-09-16 21:14:56 Annotating text fragment 6311/100392
## 2023-09-16 21:14:57 Annotating text fragment 6321/100392
## 2023-09-16 21:14:57 Annotating text fragment 6331/100392
## 2023-09-16 21:14:57 Annotating text fragment 6341/100392
## 2023-09-16 21:14:57 Annotating text fragment 6351/100392
## 2023-09-16 21:14:57 Annotating text fragment 6361/100392
## 2023-09-16 21:14:57 Annotating text fragment 6371/100392
## 2023-09-16 21:14:58 Annotating text fragment 6381/100392
## 2023-09-16 21:14:58 Annotating text fragment 6391/100392
## 2023-09-16 21:14:58 Annotating text fragment 6401/100392
## 2023-09-16 21:14:58 Annotating text fragment 6411/100392
## 2023-09-16 21:14:58 Annotating text fragment 6421/100392
## 2023-09-16 21:14:58 Annotating text fragment 6431/100392
## 2023-09-16 21:14:58 Annotating text fragment 6441/100392
## 2023-09-16 21:14:58 Annotating text fragment 6451/100392
## 2023-09-16 21:14:59 Annotating text fragment 6461/100392
## 2023-09-16 21:14:59 Annotating text fragment 6471/100392
## 2023-09-16 21:14:59 Annotating text fragment 6481/100392
## 2023-09-16 21:14:59 Annotating text fragment 6491/100392
## 2023-09-16 21:14:59 Annotating text fragment 6501/100392
## 2023-09-16 21:14:59 Annotating text fragment 6511/100392
## 2023-09-16 21:14:59 Annotating text fragment 6521/100392
## 2023-09-16 21:14:59 Annotating text fragment 6531/100392
```

```
## 2023-09-16 21:15:00 Annotating text fragment 6541/100392
## 2023-09-16 21:15:00 Annotating text fragment 6551/100392
## 2023-09-16 21:15:00 Annotating text fragment 6561/100392
## 2023-09-16 21:15:00 Annotating text fragment 6571/100392
## 2023-09-16 21:15:00 Annotating text fragment 6581/100392
## 2023-09-16 21:15:00 Annotating text fragment 6591/100392
## 2023-09-16 21:15:00 Annotating text fragment 6601/100392
## 2023-09-16 21:15:00 Annotating text fragment 6611/100392
## 2023-09-16 21:15:00 Annotating text fragment 6621/100392
## 2023-09-16 21:15:01 Annotating text fragment 6631/100392
## 2023-09-16 21:15:01 Annotating text fragment 6641/100392
## 2023-09-16 21:15:01 Annotating text fragment 6651/100392
## 2023-09-16 21:15:01 Annotating text fragment 6661/100392
## 2023-09-16 21:15:01 Annotating text fragment 6671/100392
## 2023-09-16 21:15:01 Annotating text fragment 6681/100392
## 2023-09-16 21:15:01 Annotating text fragment 6691/100392
## 2023-09-16 21:15:01 Annotating text fragment 6701/100392
## 2023-09-16 21:15:01 Annotating text fragment 6711/100392
## 2023-09-16 21:15:01 Annotating text fragment 6721/100392
## 2023-09-16 21:15:01 Annotating text fragment 6731/100392
## 2023-09-16 21:15:02 Annotating text fragment 6741/100392
## 2023-09-16 21:15:02 Annotating text fragment 6751/100392
## 2023-09-16 21:15:02 Annotating text fragment 6761/100392
## 2023-09-16 21:15:02 Annotating text fragment 6771/100392
## 2023-09-16 21:15:02 Annotating text fragment 6781/100392
## 2023-09-16 21:15:02 Annotating text fragment 6791/100392
## 2023-09-16 21:15:02 Annotating text fragment 6801/100392
## 2023-09-16 21:15:02 Annotating text fragment 6811/100392
## 2023-09-16 21:15:02 Annotating text fragment 6821/100392
## 2023-09-16 21:15:02 Annotating text fragment 6831/100392
## 2023-09-16 21:15:02 Annotating text fragment 6841/100392
## 2023-09-16 21:15:02 Annotating text fragment 6851/100392
## 2023-09-16 21:15:03 Annotating text fragment 6861/100392
## 2023-09-16 21:15:03 Annotating text fragment 6871/100392
## 2023-09-16 21:15:03 Annotating text fragment 6881/100392
## 2023-09-16 21:15:03 Annotating text fragment 6891/100392
## 2023-09-16 21:15:03 Annotating text fragment 6901/100392
## 2023-09-16 21:15:03 Annotating text fragment 6911/100392
## 2023-09-16 21:15:03 Annotating text fragment 6921/100392
## 2023-09-16 21:15:03 Annotating text fragment 6931/100392
## 2023-09-16 21:15:03 Annotating text fragment 6941/100392
## 2023-09-16 21:15:04 Annotating text fragment 6951/100392
## 2023-09-16 21:15:04 Annotating text fragment 6961/100392
## 2023-09-16 21:15:04 Annotating text fragment 6971/100392
## 2023-09-16 21:15:04 Annotating text fragment 6981/100392
## 2023-09-16 21:15:04 Annotating text fragment 6991/100392
## 2023-09-16 21:15:04 Annotating text fragment 7001/100392
## 2023-09-16 21:15:04 Annotating text fragment 7011/100392
## 2023-09-16 21:15:04 Annotating text fragment 7021/100392
## 2023-09-16 21:15:04 Annotating text fragment 7031/100392
## 2023-09-16 21:15:04 Annotating text fragment 7041/100392
## 2023-09-16 21:15:05 Annotating text fragment 7051/100392
## 2023-09-16 21:15:05 Annotating text fragment 7061/100392
## 2023-09-16 21:15:05 Annotating text fragment 7071/100392
```

```
## 2023-09-16 21:15:05 Annotating text fragment 7081/100392
## 2023-09-16 21:15:05 Annotating text fragment 7091/100392
## 2023-09-16 21:15:05 Annotating text fragment 7101/100392
## 2023-09-16 21:15:05 Annotating text fragment 7111/100392
## 2023-09-16 21:15:05 Annotating text fragment 7121/100392
## 2023-09-16 21:15:05 Annotating text fragment 7131/100392
## 2023-09-16 21:15:06 Annotating text fragment 7141/100392
## 2023-09-16 21:15:06 Annotating text fragment 7151/100392
## 2023-09-16 21:15:06 Annotating text fragment 7161/100392
## 2023-09-16 21:15:06 Annotating text fragment 7171/100392
## 2023-09-16 21:15:06 Annotating text fragment 7181/100392
## 2023-09-16 21:15:06 Annotating text fragment 7191/100392
## 2023-09-16 21:15:06 Annotating text fragment 7201/100392
## 2023-09-16 21:15:06 Annotating text fragment 7211/100392
## 2023-09-16 21:15:06 Annotating text fragment 7221/100392
## 2023-09-16 21:15:06 Annotating text fragment 7231/100392
## 2023-09-16 21:15:07 Annotating text fragment 7241/100392
## 2023-09-16 21:15:07 Annotating text fragment 7251/100392
## 2023-09-16 21:15:07 Annotating text fragment 7261/100392
## 2023-09-16 21:15:07 Annotating text fragment 7271/100392
## 2023-09-16 21:15:07 Annotating text fragment 7281/100392
## 2023-09-16 21:15:07 Annotating text fragment 7291/100392
## 2023-09-16 21:15:07 Annotating text fragment 7301/100392
## 2023-09-16 21:15:07 Annotating text fragment 7311/100392
## 2023-09-16 21:15:08 Annotating text fragment 7321/100392
## 2023-09-16 21:15:08 Annotating text fragment 7331/100392
## 2023-09-16 21:15:08 Annotating text fragment 7341/100392
## 2023-09-16 21:15:08 Annotating text fragment 7351/100392
## 2023-09-16 21:15:08 Annotating text fragment 7361/100392
## 2023-09-16 21:15:08 Annotating text fragment 7371/100392
## 2023-09-16 21:15:09 Annotating text fragment 7381/100392
## 2023-09-16 21:15:09 Annotating text fragment 7391/100392
## 2023-09-16 21:15:09 Annotating text fragment 7401/100392
## 2023-09-16 21:15:09 Annotating text fragment 7411/100392
## 2023-09-16 21:15:09 Annotating text fragment 7421/100392
## 2023-09-16 21:15:09 Annotating text fragment 7431/100392
## 2023-09-16 21:15:09 Annotating text fragment 7441/100392
## 2023-09-16 21:15:09 Annotating text fragment 7451/100392
## 2023-09-16 21:15:09 Annotating text fragment 7461/100392
## 2023-09-16 21:15:09 Annotating text fragment 7471/100392
## 2023-09-16 21:15:10 Annotating text fragment 7481/100392
## 2023-09-16 21:15:10 Annotating text fragment 7491/100392
## 2023-09-16 21:15:10 Annotating text fragment 7501/100392
## 2023-09-16 21:15:10 Annotating text fragment 7511/100392
## 2023-09-16 21:15:10 Annotating text fragment 7521/100392
## 2023-09-16 21:15:10 Annotating text fragment 7531/100392
## 2023-09-16 21:15:10 Annotating text fragment 7541/100392
## 2023-09-16 21:15:10 Annotating text fragment 7551/100392
## 2023-09-16 21:15:10 Annotating text fragment 7561/100392
## 2023-09-16 21:15:11 Annotating text fragment 7571/100392
## 2023-09-16 21:15:11 Annotating text fragment 7581/100392
## 2023-09-16 21:15:11 Annotating text fragment 7591/100392
## 2023-09-16 21:15:11 Annotating text fragment 7601/100392
## 2023-09-16 21:15:11 Annotating text fragment 7611/100392
```

```
## 2023-09-16 21:15:11 Annotating text fragment 7621/100392
## 2023-09-16 21:15:11 Annotating text fragment 7631/100392
## 2023-09-16 21:15:11 Annotating text fragment 7641/100392
## 2023-09-16 21:15:11 Annotating text fragment 7651/100392
## 2023-09-16 21:15:11 Annotating text fragment 7661/100392
## 2023-09-16 21:15:11 Annotating text fragment 7671/100392
## 2023-09-16 21:15:11 Annotating text fragment 7681/100392
## 2023-09-16 21:15:11 Annotating text fragment 7691/100392
## 2023-09-16 21:15:12 Annotating text fragment 7701/100392
## 2023-09-16 21:15:12 Annotating text fragment 7711/100392
## 2023-09-16 21:15:12 Annotating text fragment 7721/100392
## 2023-09-16 21:15:12 Annotating text fragment 7731/100392
## 2023-09-16 21:15:12 Annotating text fragment 7741/100392
## 2023-09-16 21:15:12 Annotating text fragment 7751/100392
## 2023-09-16 21:15:12 Annotating text fragment 7761/100392
## 2023-09-16 21:15:12 Annotating text fragment 7771/100392
## 2023-09-16 21:15:12 Annotating text fragment 7781/100392
## 2023-09-16 21:15:12 Annotating text fragment 7791/100392
## 2023-09-16 21:15:12 Annotating text fragment 7801/100392
## 2023-09-16 21:15:12 Annotating text fragment 7811/100392
## 2023-09-16 21:15:13 Annotating text fragment 7821/100392
## 2023-09-16 21:15:13 Annotating text fragment 7831/100392
## 2023-09-16 21:15:13 Annotating text fragment 7841/100392
## 2023-09-16 21:15:13 Annotating text fragment 7851/100392
## 2023-09-16 21:15:13 Annotating text fragment 7861/100392
## 2023-09-16 21:15:13 Annotating text fragment 7871/100392
## 2023-09-16 21:15:13 Annotating text fragment 7881/100392
## 2023-09-16 21:15:13 Annotating text fragment 7891/100392
## 2023-09-16 21:15:13 Annotating text fragment 7901/100392
## 2023-09-16 21:15:13 Annotating text fragment 7911/100392
## 2023-09-16 21:15:14 Annotating text fragment 7921/100392
## 2023-09-16 21:15:14 Annotating text fragment 7931/100392
## 2023-09-16 21:15:14 Annotating text fragment 7941/100392
## 2023-09-16 21:15:14 Annotating text fragment 7951/100392
## 2023-09-16 21:15:14 Annotating text fragment 7961/100392
## 2023-09-16 21:15:14 Annotating text fragment 7971/100392
## 2023-09-16 21:15:14 Annotating text fragment 7981/100392
## 2023-09-16 21:15:14 Annotating text fragment 7991/100392
## 2023-09-16 21:15:14 Annotating text fragment 8001/100392
## 2023-09-16 21:15:14 Annotating text fragment 8011/100392
## 2023-09-16 21:15:15 Annotating text fragment 8021/100392
## 2023-09-16 21:15:15 Annotating text fragment 8031/100392
## 2023-09-16 21:15:15 Annotating text fragment 8041/100392
## 2023-09-16 21:15:15 Annotating text fragment 8051/100392
## 2023-09-16 21:15:15 Annotating text fragment 8061/100392
## 2023-09-16 21:15:15 Annotating text fragment 8071/100392
## 2023-09-16 21:15:15 Annotating text fragment 8081/100392
## 2023-09-16 21:15:15 Annotating text fragment 8091/100392
## 2023-09-16 21:15:15 Annotating text fragment 8101/100392
## 2023-09-16 21:15:15 Annotating text fragment 8111/100392
## 2023-09-16 21:15:15 Annotating text fragment 8121/100392
## 2023-09-16 21:15:16 Annotating text fragment 8131/100392
## 2023-09-16 21:15:16 Annotating text fragment 8141/100392
## 2023-09-16 21:15:16 Annotating text fragment 8151/100392
```

```
## 2023-09-16 21:15:16 Annotating text fragment 8161/100392
## 2023-09-16 21:15:16 Annotating text fragment 8171/100392
## 2023-09-16 21:15:16 Annotating text fragment 8181/100392
## 2023-09-16 21:15:16 Annotating text fragment 8191/100392
## 2023-09-16 21:15:16 Annotating text fragment 8201/100392
## 2023-09-16 21:15:16 Annotating text fragment 8211/100392
## 2023-09-16 21:15:16 Annotating text fragment 8221/100392
## 2023-09-16 21:15:17 Annotating text fragment 8231/100392
## 2023-09-16 21:15:17 Annotating text fragment 8241/100392
## 2023-09-16 21:15:17 Annotating text fragment 8251/100392
## 2023-09-16 21:15:17 Annotating text fragment 8261/100392
## 2023-09-16 21:15:17 Annotating text fragment 8271/100392
## 2023-09-16 21:15:17 Annotating text fragment 8281/100392
## 2023-09-16 21:15:17 Annotating text fragment 8291/100392
## 2023-09-16 21:15:17 Annotating text fragment 8301/100392
## 2023-09-16 21:15:18 Annotating text fragment 8311/100392
## 2023-09-16 21:15:18 Annotating text fragment 8321/100392
## 2023-09-16 21:15:18 Annotating text fragment 8331/100392
## 2023-09-16 21:15:18 Annotating text fragment 8341/100392
## 2023-09-16 21:15:18 Annotating text fragment 8351/100392
## 2023-09-16 21:15:18 Annotating text fragment 8361/100392
## 2023-09-16 21:15:18 Annotating text fragment 8371/100392
## 2023-09-16 21:15:19 Annotating text fragment 8381/100392
## 2023-09-16 21:15:19 Annotating text fragment 8391/100392
## 2023-09-16 21:15:19 Annotating text fragment 8401/100392
## 2023-09-16 21:15:19 Annotating text fragment 8411/100392
## 2023-09-16 21:15:19 Annotating text fragment 8421/100392
## 2023-09-16 21:15:19 Annotating text fragment 8431/100392
## 2023-09-16 21:15:19 Annotating text fragment 8441/100392
## 2023-09-16 21:15:19 Annotating text fragment 8451/100392
## 2023-09-16 21:15:19 Annotating text fragment 8461/100392
## 2023-09-16 21:15:20 Annotating text fragment 8471/100392
## 2023-09-16 21:15:20 Annotating text fragment 8481/100392
## 2023-09-16 21:15:20 Annotating text fragment 8491/100392
## 2023-09-16 21:15:20 Annotating text fragment 8501/100392
## 2023-09-16 21:15:20 Annotating text fragment 8511/100392
## 2023-09-16 21:15:20 Annotating text fragment 8521/100392
## 2023-09-16 21:15:20 Annotating text fragment 8531/100392
## 2023-09-16 21:15:20 Annotating text fragment 8541/100392
## 2023-09-16 21:15:21 Annotating text fragment 8551/100392
## 2023-09-16 21:15:21 Annotating text fragment 8561/100392
## 2023-09-16 21:15:21 Annotating text fragment 8571/100392
## 2023-09-16 21:15:21 Annotating text fragment 8581/100392
## 2023-09-16 21:15:21 Annotating text fragment 8591/100392
## 2023-09-16 21:15:21 Annotating text fragment 8601/100392
## 2023-09-16 21:15:21 Annotating text fragment 8611/100392
## 2023-09-16 21:15:21 Annotating text fragment 8621/100392
## 2023-09-16 21:15:21 Annotating text fragment 8631/100392
## 2023-09-16 21:15:21 Annotating text fragment 8641/100392
## 2023-09-16 21:15:22 Annotating text fragment 8651/100392
## 2023-09-16 21:15:22 Annotating text fragment 8661/100392
## 2023-09-16 21:15:22 Annotating text fragment 8671/100392
## 2023-09-16 21:15:22 Annotating text fragment 8681/100392
## 2023-09-16 21:15:22 Annotating text fragment 8691/100392
```

```
## 2023-09-16 21:15:22 Annotating text fragment 8701/100392
## 2023-09-16 21:15:22 Annotating text fragment 8711/100392
## 2023-09-16 21:15:22 Annotating text fragment 8721/100392
## 2023-09-16 21:15:23 Annotating text fragment 8731/100392
## 2023-09-16 21:15:23 Annotating text fragment 8741/100392
## 2023-09-16 21:15:23 Annotating text fragment 8751/100392
## 2023-09-16 21:15:23 Annotating text fragment 8761/100392
## 2023-09-16 21:15:23 Annotating text fragment 8771/100392
## 2023-09-16 21:15:23 Annotating text fragment 8781/100392
## 2023-09-16 21:15:23 Annotating text fragment 8791/100392
## 2023-09-16 21:15:24 Annotating text fragment 8801/100392
## 2023-09-16 21:15:24 Annotating text fragment 8811/100392
## 2023-09-16 21:15:24 Annotating text fragment 8821/100392
## 2023-09-16 21:15:24 Annotating text fragment 8831/100392
## 2023-09-16 21:15:24 Annotating text fragment 8841/100392
## 2023-09-16 21:15:24 Annotating text fragment 8851/100392
## 2023-09-16 21:15:24 Annotating text fragment 8861/100392
## 2023-09-16 21:15:24 Annotating text fragment 8871/100392
## 2023-09-16 21:15:24 Annotating text fragment 8881/100392
## 2023-09-16 21:15:24 Annotating text fragment 8891/100392
## 2023-09-16 21:15:25 Annotating text fragment 8901/100392
## 2023-09-16 21:15:25 Annotating text fragment 8911/100392
## 2023-09-16 21:15:25 Annotating text fragment 8921/100392
## 2023-09-16 21:15:25 Annotating text fragment 8931/100392
## 2023-09-16 21:15:25 Annotating text fragment 8941/100392
## 2023-09-16 21:15:25 Annotating text fragment 8951/100392
## 2023-09-16 21:15:25 Annotating text fragment 8961/100392
## 2023-09-16 21:15:26 Annotating text fragment 8971/100392
## 2023-09-16 21:15:26 Annotating text fragment 8981/100392
## 2023-09-16 21:15:26 Annotating text fragment 8991/100392
## 2023-09-16 21:15:26 Annotating text fragment 9001/100392
## 2023-09-16 21:15:26 Annotating text fragment 9011/100392
## 2023-09-16 21:15:26 Annotating text fragment 9021/100392
## 2023-09-16 21:15:26 Annotating text fragment 9031/100392
## 2023-09-16 21:15:26 Annotating text fragment 9041/100392
## 2023-09-16 21:15:26 Annotating text fragment 9051/100392
## 2023-09-16 21:15:27 Annotating text fragment 9061/100392
## 2023-09-16 21:15:27 Annotating text fragment 9071/100392
## 2023-09-16 21:15:27 Annotating text fragment 9081/100392
## 2023-09-16 21:15:27 Annotating text fragment 9091/100392
## 2023-09-16 21:15:27 Annotating text fragment 9101/100392
## 2023-09-16 21:15:27 Annotating text fragment 9111/100392
## 2023-09-16 21:15:27 Annotating text fragment 9121/100392
## 2023-09-16 21:15:27 Annotating text fragment 9131/100392
## 2023-09-16 21:15:27 Annotating text fragment 9141/100392
## 2023-09-16 21:15:28 Annotating text fragment 9151/100392
## 2023-09-16 21:15:28 Annotating text fragment 9161/100392
## 2023-09-16 21:15:28 Annotating text fragment 9171/100392
## 2023-09-16 21:15:28 Annotating text fragment 9181/100392
## 2023-09-16 21:15:28 Annotating text fragment 9191/100392
## 2023-09-16 21:15:28 Annotating text fragment 9201/100392
## 2023-09-16 21:15:28 Annotating text fragment 9211/100392
## 2023-09-16 21:15:29 Annotating text fragment 9221/100392
## 2023-09-16 21:15:29 Annotating text fragment 9231/100392
```

```
## 2023-09-16 21:15:29 Annotating text fragment 9241/100392
## 2023-09-16 21:15:29 Annotating text fragment 9251/100392
## 2023-09-16 21:15:29 Annotating text fragment 9261/100392
## 2023-09-16 21:15:29 Annotating text fragment 9271/100392
## 2023-09-16 21:15:29 Annotating text fragment 9281/100392
## 2023-09-16 21:15:29 Annotating text fragment 9291/100392
## 2023-09-16 21:15:29 Annotating text fragment 9301/100392
## 2023-09-16 21:15:30 Annotating text fragment 9311/100392
## 2023-09-16 21:15:30 Annotating text fragment 9321/100392
## 2023-09-16 21:15:30 Annotating text fragment 9331/100392
## 2023-09-16 21:15:30 Annotating text fragment 9341/100392
## 2023-09-16 21:15:30 Annotating text fragment 9351/100392
## 2023-09-16 21:15:30 Annotating text fragment 9361/100392
## 2023-09-16 21:15:30 Annotating text fragment 9371/100392
## 2023-09-16 21:15:31 Annotating text fragment 9381/100392
## 2023-09-16 21:15:31 Annotating text fragment 9391/100392
## 2023-09-16 21:15:31 Annotating text fragment 9401/100392
## 2023-09-16 21:15:31 Annotating text fragment 9411/100392
## 2023-09-16 21:15:31 Annotating text fragment 9421/100392
## 2023-09-16 21:15:31 Annotating text fragment 9431/100392
## 2023-09-16 21:15:31 Annotating text fragment 9441/100392
## 2023-09-16 21:15:31 Annotating text fragment 9451/100392
## 2023-09-16 21:15:31 Annotating text fragment 9461/100392
## 2023-09-16 21:15:31 Annotating text fragment 9471/100392
## 2023-09-16 21:15:31 Annotating text fragment 9481/100392
## 2023-09-16 21:15:32 Annotating text fragment 9491/100392
## 2023-09-16 21:15:32 Annotating text fragment 9501/100392
## 2023-09-16 21:15:32 Annotating text fragment 9511/100392
## 2023-09-16 21:15:32 Annotating text fragment 9521/100392
## 2023-09-16 21:15:32 Annotating text fragment 9531/100392
## 2023-09-16 21:15:33 Annotating text fragment 9541/100392
## 2023-09-16 21:15:33 Annotating text fragment 9551/100392
## 2023-09-16 21:15:33 Annotating text fragment 9561/100392
## 2023-09-16 21:15:33 Annotating text fragment 9571/100392
## 2023-09-16 21:15:33 Annotating text fragment 9581/100392
## 2023-09-16 21:15:33 Annotating text fragment 9591/100392
## 2023-09-16 21:15:33 Annotating text fragment 9601/100392
## 2023-09-16 21:15:33 Annotating text fragment 9611/100392
## 2023-09-16 21:15:33 Annotating text fragment 9621/100392
## 2023-09-16 21:15:34 Annotating text fragment 9631/100392
## 2023-09-16 21:15:34 Annotating text fragment 9641/100392
## 2023-09-16 21:15:34 Annotating text fragment 9651/100392
## 2023-09-16 21:15:34 Annotating text fragment 9661/100392
## 2023-09-16 21:15:34 Annotating text fragment 9671/100392
## 2023-09-16 21:15:34 Annotating text fragment 9681/100392
## 2023-09-16 21:15:34 Annotating text fragment 9691/100392
## 2023-09-16 21:15:34 Annotating text fragment 9701/100392
## 2023-09-16 21:15:35 Annotating text fragment 9711/100392
## 2023-09-16 21:15:35 Annotating text fragment 9721/100392
## 2023-09-16 21:15:35 Annotating text fragment 9731/100392
## 2023-09-16 21:15:35 Annotating text fragment 9741/100392
## 2023-09-16 21:15:35 Annotating text fragment 9751/100392
## 2023-09-16 21:15:35 Annotating text fragment 9761/100392
## 2023-09-16 21:15:35 Annotating text fragment 9771/100392
```

```
## 2023-09-16 21:15:36 Annotating text fragment 9781/100392
## 2023-09-16 21:15:36 Annotating text fragment 9791/100392
## 2023-09-16 21:15:36 Annotating text fragment 9801/100392
## 2023-09-16 21:15:36 Annotating text fragment 9811/100392
## 2023-09-16 21:15:36 Annotating text fragment 9821/100392
## 2023-09-16 21:15:36 Annotating text fragment 9831/100392
## 2023-09-16 21:15:36 Annotating text fragment 9841/100392
## 2023-09-16 21:15:36 Annotating text fragment 9851/100392
## 2023-09-16 21:15:37 Annotating text fragment 9861/100392
## 2023-09-16 21:15:37 Annotating text fragment 9871/100392
## 2023-09-16 21:15:37 Annotating text fragment 9881/100392
## 2023-09-16 21:15:37 Annotating text fragment 9891/100392
## 2023-09-16 21:15:37 Annotating text fragment 9901/100392
## 2023-09-16 21:15:37 Annotating text fragment 9911/100392
## 2023-09-16 21:15:37 Annotating text fragment 9921/100392
## 2023-09-16 21:15:37 Annotating text fragment 9931/100392
## 2023-09-16 21:15:37 Annotating text fragment 9941/100392
## 2023-09-16 21:15:37 Annotating text fragment 9951/100392
## 2023-09-16 21:15:37 Annotating text fragment 9961/100392
## 2023-09-16 21:15:38 Annotating text fragment 9971/100392
## 2023-09-16 21:15:38 Annotating text fragment 9981/100392
## 2023-09-16 21:15:38 Annotating text fragment 9991/100392
## 2023-09-16 21:15:38 Annotating text fragment 10001/100392
## 2023-09-16 21:15:38 Annotating text fragment 10011/100392
## 2023-09-16 21:15:38 Annotating text fragment 10021/100392
## 2023-09-16 21:15:38 Annotating text fragment 10031/100392
## 2023-09-16 21:15:38 Annotating text fragment 10041/100392
## 2023-09-16 21:15:39 Annotating text fragment 10051/100392
## 2023-09-16 21:15:39 Annotating text fragment 10061/100392
## 2023-09-16 21:15:39 Annotating text fragment 10071/100392
## 2023-09-16 21:15:39 Annotating text fragment 10081/100392
## 2023-09-16 21:15:39 Annotating text fragment 10091/100392
## 2023-09-16 21:15:39 Annotating text fragment 10101/100392
## 2023-09-16 21:15:39 Annotating text fragment 10111/100392
## 2023-09-16 21:15:39 Annotating text fragment 10121/100392
## 2023-09-16 21:15:39 Annotating text fragment 10131/100392
## 2023-09-16 21:15:39 Annotating text fragment 10141/100392
## 2023-09-16 21:15:40 Annotating text fragment 10151/100392
## 2023-09-16 21:15:40 Annotating text fragment 10161/100392
## 2023-09-16 21:15:40 Annotating text fragment 10171/100392
## 2023-09-16 21:15:40 Annotating text fragment 10181/100392
## 2023-09-16 21:15:40 Annotating text fragment 10191/100392
## 2023-09-16 21:15:40 Annotating text fragment 10201/100392
## 2023-09-16 21:15:40 Annotating text fragment 10211/100392
## 2023-09-16 21:15:40 Annotating text fragment 10221/100392
## 2023-09-16 21:15:40 Annotating text fragment 10231/100392
## 2023-09-16 21:15:40 Annotating text fragment 10241/100392
## 2023-09-16 21:15:41 Annotating text fragment 10251/100392
## 2023-09-16 21:15:41 Annotating text fragment 10261/100392
## 2023-09-16 21:15:41 Annotating text fragment 10271/100392
## 2023-09-16 21:15:41 Annotating text fragment 10281/100392
## 2023-09-16 21:15:41 Annotating text fragment 10291/100392
## 2023-09-16 21:15:41 Annotating text fragment 10301/100392
## 2023-09-16 21:15:41 Annotating text fragment 10311/100392
```

```
## 2023-09-16 21:15:41 Annotating text fragment 10321/100392
## 2023-09-16 21:15:41 Annotating text fragment 10331/100392
## 2023-09-16 21:15:41 Annotating text fragment 10341/100392
## 2023-09-16 21:15:42 Annotating text fragment 10351/100392
## 2023-09-16 21:15:42 Annotating text fragment 10361/100392
## 2023-09-16 21:15:42 Annotating text fragment 10371/100392
## 2023-09-16 21:15:42 Annotating text fragment 10381/100392
## 2023-09-16 21:15:42 Annotating text fragment 10391/100392
## 2023-09-16 21:15:42 Annotating text fragment 10401/100392
## 2023-09-16 21:15:42 Annotating text fragment 10411/100392
## 2023-09-16 21:15:42 Annotating text fragment 10421/100392
## 2023-09-16 21:15:42 Annotating text fragment 10431/100392
## 2023-09-16 21:15:43 Annotating text fragment 10441/100392
## 2023-09-16 21:15:43 Annotating text fragment 10451/100392
## 2023-09-16 21:15:43 Annotating text fragment 10461/100392
## 2023-09-16 21:15:43 Annotating text fragment 10471/100392
## 2023-09-16 21:15:43 Annotating text fragment 10481/100392
## 2023-09-16 21:15:43 Annotating text fragment 10491/100392
## 2023-09-16 21:15:43 Annotating text fragment 10501/100392
## 2023-09-16 21:15:43 Annotating text fragment 10511/100392
## 2023-09-16 21:15:43 Annotating text fragment 10521/100392
## 2023-09-16 21:15:44 Annotating text fragment 10531/100392
## 2023-09-16 21:15:44 Annotating text fragment 10541/100392
## 2023-09-16 21:15:44 Annotating text fragment 10551/100392
## 2023-09-16 21:15:44 Annotating text fragment 10561/100392
## 2023-09-16 21:15:44 Annotating text fragment 10571/100392
## 2023-09-16 21:15:44 Annotating text fragment 10581/100392
## 2023-09-16 21:15:44 Annotating text fragment 10591/100392
## 2023-09-16 21:15:44 Annotating text fragment 10601/100392
## 2023-09-16 21:15:44 Annotating text fragment 10611/100392
## 2023-09-16 21:15:44 Annotating text fragment 10621/100392
## 2023-09-16 21:15:44 Annotating text fragment 10631/100392
## 2023-09-16 21:15:45 Annotating text fragment 10641/100392
## 2023-09-16 21:15:45 Annotating text fragment 10651/100392
## 2023-09-16 21:15:45 Annotating text fragment 10661/100392
## 2023-09-16 21:15:45 Annotating text fragment 10671/100392
## 2023-09-16 21:15:45 Annotating text fragment 10681/100392
## 2023-09-16 21:15:45 Annotating text fragment 10691/100392
## 2023-09-16 21:15:45 Annotating text fragment 10701/100392
## 2023-09-16 21:15:45 Annotating text fragment 10711/100392
## 2023-09-16 21:15:46 Annotating text fragment 10721/100392
## 2023-09-16 21:15:46 Annotating text fragment 10731/100392
## 2023-09-16 21:15:46 Annotating text fragment 10741/100392
## 2023-09-16 21:15:46 Annotating text fragment 10751/100392
## 2023-09-16 21:15:46 Annotating text fragment 10761/100392
## 2023-09-16 21:15:46 Annotating text fragment 10771/100392
## 2023-09-16 21:15:46 Annotating text fragment 10781/100392
## 2023-09-16 21:15:46 Annotating text fragment 10791/100392
## 2023-09-16 21:15:46 Annotating text fragment 10801/100392
## 2023-09-16 21:15:47 Annotating text fragment 10811/100392
## 2023-09-16 21:15:47 Annotating text fragment 10821/100392
## 2023-09-16 21:15:47 Annotating text fragment 10831/100392
## 2023-09-16 21:15:47 Annotating text fragment 10841/100392
## 2023-09-16 21:15:47 Annotating text fragment 10851/100392
```

```
## 2023-09-16 21:15:47 Annotating text fragment 10861/100392
## 2023-09-16 21:15:47 Annotating text fragment 10871/100392
## 2023-09-16 21:15:47 Annotating text fragment 10881/100392
## 2023-09-16 21:15:47 Annotating text fragment 10891/100392
## 2023-09-16 21:15:47 Annotating text fragment 10901/100392
## 2023-09-16 21:15:48 Annotating text fragment 10911/100392
## 2023-09-16 21:15:48 Annotating text fragment 10921/100392
## 2023-09-16 21:15:48 Annotating text fragment 10931/100392
## 2023-09-16 21:15:48 Annotating text fragment 10941/100392
## 2023-09-16 21:15:48 Annotating text fragment 10951/100392
## 2023-09-16 21:15:48 Annotating text fragment 10961/100392
## 2023-09-16 21:15:48 Annotating text fragment 10971/100392
## 2023-09-16 21:15:48 Annotating text fragment 10981/100392
## 2023-09-16 21:15:48 Annotating text fragment 10991/100392
## 2023-09-16 21:15:48 Annotating text fragment 11001/100392
## 2023-09-16 21:15:48 Annotating text fragment 11011/100392
## 2023-09-16 21:15:49 Annotating text fragment 11021/100392
## 2023-09-16 21:15:49 Annotating text fragment 11031/100392
## 2023-09-16 21:15:49 Annotating text fragment 11041/100392
## 2023-09-16 21:15:49 Annotating text fragment 11051/100392
## 2023-09-16 21:15:49 Annotating text fragment 11061/100392
## 2023-09-16 21:15:49 Annotating text fragment 11071/100392
## 2023-09-16 21:15:49 Annotating text fragment 11081/100392
## 2023-09-16 21:15:49 Annotating text fragment 11091/100392
## 2023-09-16 21:15:49 Annotating text fragment 11101/100392
## 2023-09-16 21:15:50 Annotating text fragment 11111/100392
## 2023-09-16 21:15:50 Annotating text fragment 11121/100392
## 2023-09-16 21:15:50 Annotating text fragment 11131/100392
## 2023-09-16 21:15:50 Annotating text fragment 11141/100392
## 2023-09-16 21:15:50 Annotating text fragment 11151/100392
## 2023-09-16 21:15:50 Annotating text fragment 11161/100392
## 2023-09-16 21:15:50 Annotating text fragment 11171/100392
## 2023-09-16 21:15:50 Annotating text fragment 11181/100392
## 2023-09-16 21:15:50 Annotating text fragment 11191/100392
## 2023-09-16 21:15:50 Annotating text fragment 11201/100392
## 2023-09-16 21:15:50 Annotating text fragment 11211/100392
## 2023-09-16 21:15:51 Annotating text fragment 11221/100392
## 2023-09-16 21:15:51 Annotating text fragment 11231/100392
## 2023-09-16 21:15:51 Annotating text fragment 11241/100392
## 2023-09-16 21:15:51 Annotating text fragment 11251/100392
## 2023-09-16 21:15:51 Annotating text fragment 11261/100392
## 2023-09-16 21:15:51 Annotating text fragment 11271/100392
## 2023-09-16 21:15:51 Annotating text fragment 11281/100392
## 2023-09-16 21:15:51 Annotating text fragment 11291/100392
## 2023-09-16 21:15:51 Annotating text fragment 11301/100392
## 2023-09-16 21:15:51 Annotating text fragment 11311/100392
## 2023-09-16 21:15:51 Annotating text fragment 11321/100392
## 2023-09-16 21:15:52 Annotating text fragment 11331/100392
## 2023-09-16 21:15:52 Annotating text fragment 11341/100392
## 2023-09-16 21:15:52 Annotating text fragment 11351/100392
## 2023-09-16 21:15:52 Annotating text fragment 11361/100392
## 2023-09-16 21:15:52 Annotating text fragment 11371/100392
## 2023-09-16 21:15:52 Annotating text fragment 11381/100392
## 2023-09-16 21:15:52 Annotating text fragment 11391/100392
```

```
## 2023-09-16 21:15:52 Annotating text fragment 11401/100392
## 2023-09-16 21:15:52 Annotating text fragment 11411/100392
## 2023-09-16 21:15:52 Annotating text fragment 11421/100392
## 2023-09-16 21:15:52 Annotating text fragment 11431/100392
## 2023-09-16 21:15:53 Annotating text fragment 11441/100392
## 2023-09-16 21:15:53 Annotating text fragment 11451/100392
## 2023-09-16 21:15:53 Annotating text fragment 11461/100392
## 2023-09-16 21:15:53 Annotating text fragment 11471/100392
## 2023-09-16 21:15:53 Annotating text fragment 11481/100392
## 2023-09-16 21:15:53 Annotating text fragment 11491/100392
## 2023-09-16 21:15:53 Annotating text fragment 11501/100392
## 2023-09-16 21:15:53 Annotating text fragment 11511/100392
## 2023-09-16 21:15:53 Annotating text fragment 11521/100392
## 2023-09-16 21:15:53 Annotating text fragment 11531/100392
## 2023-09-16 21:15:54 Annotating text fragment 11541/100392
## 2023-09-16 21:15:54 Annotating text fragment 11551/100392
## 2023-09-16 21:15:54 Annotating text fragment 11561/100392
## 2023-09-16 21:15:54 Annotating text fragment 11571/100392
## 2023-09-16 21:15:54 Annotating text fragment 11581/100392
## 2023-09-16 21:15:54 Annotating text fragment 11591/100392
## 2023-09-16 21:15:54 Annotating text fragment 11601/100392
## 2023-09-16 21:15:54 Annotating text fragment 11611/100392
## 2023-09-16 21:15:54 Annotating text fragment 11621/100392
## 2023-09-16 21:15:54 Annotating text fragment 11631/100392
## 2023-09-16 21:15:55 Annotating text fragment 11641/100392
## 2023-09-16 21:15:55 Annotating text fragment 11651/100392
## 2023-09-16 21:15:55 Annotating text fragment 11661/100392
## 2023-09-16 21:15:55 Annotating text fragment 11671/100392
## 2023-09-16 21:15:55 Annotating text fragment 11681/100392
## 2023-09-16 21:15:55 Annotating text fragment 11691/100392
## 2023-09-16 21:15:55 Annotating text fragment 11701/100392
## 2023-09-16 21:15:55 Annotating text fragment 11711/100392
## 2023-09-16 21:15:55 Annotating text fragment 11721/100392
## 2023-09-16 21:15:56 Annotating text fragment 11731/100392
## 2023-09-16 21:15:56 Annotating text fragment 11741/100392
## 2023-09-16 21:15:56 Annotating text fragment 11751/100392
## 2023-09-16 21:15:56 Annotating text fragment 11761/100392
## 2023-09-16 21:15:56 Annotating text fragment 11771/100392
## 2023-09-16 21:15:56 Annotating text fragment 11781/100392
## 2023-09-16 21:15:56 Annotating text fragment 11791/100392
## 2023-09-16 21:15:56 Annotating text fragment 11801/100392
## 2023-09-16 21:15:56 Annotating text fragment 11811/100392
## 2023-09-16 21:15:57 Annotating text fragment 11821/100392
## 2023-09-16 21:15:57 Annotating text fragment 11831/100392
## 2023-09-16 21:15:57 Annotating text fragment 11841/100392
## 2023-09-16 21:15:57 Annotating text fragment 11851/100392
## 2023-09-16 21:15:57 Annotating text fragment 11861/100392
## 2023-09-16 21:15:57 Annotating text fragment 11871/100392
## 2023-09-16 21:15:57 Annotating text fragment 11881/100392
## 2023-09-16 21:15:57 Annotating text fragment 11891/100392
## 2023-09-16 21:15:57 Annotating text fragment 11901/100392
## 2023-09-16 21:15:57 Annotating text fragment 11911/100392
## 2023-09-16 21:15:57 Annotating text fragment 11921/100392
## 2023-09-16 21:15:58 Annotating text fragment 11931/100392
```

```
## 2023-09-16 21:15:58 Annotating text fragment 11941/100392
## 2023-09-16 21:15:58 Annotating text fragment 11951/100392
## 2023-09-16 21:15:58 Annotating text fragment 11961/100392
## 2023-09-16 21:15:58 Annotating text fragment 11971/100392
## 2023-09-16 21:15:58 Annotating text fragment 11981/100392
## 2023-09-16 21:15:58 Annotating text fragment 11991/100392
## 2023-09-16 21:15:58 Annotating text fragment 12001/100392
## 2023-09-16 21:15:58 Annotating text fragment 12011/100392
## 2023-09-16 21:15:58 Annotating text fragment 12021/100392
## 2023-09-16 21:15:59 Annotating text fragment 12031/100392
## 2023-09-16 21:15:59 Annotating text fragment 12041/100392
## 2023-09-16 21:15:59 Annotating text fragment 12051/100392
## 2023-09-16 21:15:59 Annotating text fragment 12061/100392
## 2023-09-16 21:15:59 Annotating text fragment 12071/100392
## 2023-09-16 21:15:59 Annotating text fragment 12081/100392
## 2023-09-16 21:15:59 Annotating text fragment 12091/100392
## 2023-09-16 21:15:59 Annotating text fragment 12101/100392
## 2023-09-16 21:15:59 Annotating text fragment 12111/100392
## 2023-09-16 21:15:59 Annotating text fragment 12121/100392
## 2023-09-16 21:16:00 Annotating text fragment 12131/100392
## 2023-09-16 21:16:00 Annotating text fragment 12141/100392
## 2023-09-16 21:16:00 Annotating text fragment 12151/100392
## 2023-09-16 21:16:00 Annotating text fragment 12161/100392
## 2023-09-16 21:16:00 Annotating text fragment 12171/100392
## 2023-09-16 21:16:00 Annotating text fragment 12181/100392
## 2023-09-16 21:16:00 Annotating text fragment 12191/100392
## 2023-09-16 21:16:00 Annotating text fragment 12201/100392
## 2023-09-16 21:16:00 Annotating text fragment 12211/100392
## 2023-09-16 21:16:00 Annotating text fragment 12221/100392
## 2023-09-16 21:16:00 Annotating text fragment 12231/100392
## 2023-09-16 21:16:01 Annotating text fragment 12241/100392
## 2023-09-16 21:16:01 Annotating text fragment 12251/100392
## 2023-09-16 21:16:01 Annotating text fragment 12261/100392
## 2023-09-16 21:16:01 Annotating text fragment 12271/100392
## 2023-09-16 21:16:01 Annotating text fragment 12281/100392
## 2023-09-16 21:16:01 Annotating text fragment 12291/100392
## 2023-09-16 21:16:01 Annotating text fragment 12301/100392
## 2023-09-16 21:16:01 Annotating text fragment 12311/100392
## 2023-09-16 21:16:01 Annotating text fragment 12321/100392
## 2023-09-16 21:16:01 Annotating text fragment 12331/100392
## 2023-09-16 21:16:01 Annotating text fragment 12341/100392
## 2023-09-16 21:16:01 Annotating text fragment 12351/100392
## 2023-09-16 21:16:02 Annotating text fragment 12361/100392
## 2023-09-16 21:16:02 Annotating text fragment 12371/100392
## 2023-09-16 21:16:02 Annotating text fragment 12381/100392
## 2023-09-16 21:16:02 Annotating text fragment 12391/100392
## 2023-09-16 21:16:02 Annotating text fragment 12401/100392
## 2023-09-16 21:16:02 Annotating text fragment 12411/100392
## 2023-09-16 21:16:02 Annotating text fragment 12421/100392
## 2023-09-16 21:16:02 Annotating text fragment 12431/100392
## 2023-09-16 21:16:02 Annotating text fragment 12441/100392
## 2023-09-16 21:16:02 Annotating text fragment 12451/100392
## 2023-09-16 21:16:03 Annotating text fragment 12461/100392
## 2023-09-16 21:16:03 Annotating text fragment 12471/100392
```

```
## 2023-09-16 21:16:03 Annotating text fragment 12481/100392
## 2023-09-16 21:16:03 Annotating text fragment 12491/100392
## 2023-09-16 21:16:03 Annotating text fragment 12501/100392
## 2023-09-16 21:16:03 Annotating text fragment 12511/100392
## 2023-09-16 21:16:03 Annotating text fragment 12521/100392
## 2023-09-16 21:16:03 Annotating text fragment 12531/100392
## 2023-09-16 21:16:03 Annotating text fragment 12541/100392
## 2023-09-16 21:16:03 Annotating text fragment 12551/100392
## 2023-09-16 21:16:04 Annotating text fragment 12561/100392
## 2023-09-16 21:16:04 Annotating text fragment 12571/100392
## 2023-09-16 21:16:04 Annotating text fragment 12581/100392
## 2023-09-16 21:16:04 Annotating text fragment 12591/100392
## 2023-09-16 21:16:04 Annotating text fragment 12601/100392
## 2023-09-16 21:16:04 Annotating text fragment 12611/100392
## 2023-09-16 21:16:04 Annotating text fragment 12621/100392
## 2023-09-16 21:16:04 Annotating text fragment 12631/100392
## 2023-09-16 21:16:04 Annotating text fragment 12641/100392
## 2023-09-16 21:16:05 Annotating text fragment 12651/100392
## 2023-09-16 21:16:05 Annotating text fragment 12661/100392
## 2023-09-16 21:16:05 Annotating text fragment 12671/100392
## 2023-09-16 21:16:05 Annotating text fragment 12681/100392
## 2023-09-16 21:16:05 Annotating text fragment 12691/100392
## 2023-09-16 21:16:05 Annotating text fragment 12701/100392
## 2023-09-16 21:16:05 Annotating text fragment 12711/100392
## 2023-09-16 21:16:05 Annotating text fragment 12721/100392
## 2023-09-16 21:16:05 Annotating text fragment 12731/100392
## 2023-09-16 21:16:06 Annotating text fragment 12741/100392
## 2023-09-16 21:16:06 Annotating text fragment 12751/100392
## 2023-09-16 21:16:06 Annotating text fragment 12761/100392
## 2023-09-16 21:16:06 Annotating text fragment 12771/100392
## 2023-09-16 21:16:06 Annotating text fragment 12781/100392
## 2023-09-16 21:16:06 Annotating text fragment 12791/100392
## 2023-09-16 21:16:06 Annotating text fragment 12801/100392
## 2023-09-16 21:16:06 Annotating text fragment 12811/100392
## 2023-09-16 21:16:06 Annotating text fragment 12821/100392
## 2023-09-16 21:16:06 Annotating text fragment 12831/100392
## 2023-09-16 21:16:07 Annotating text fragment 12841/100392
## 2023-09-16 21:16:07 Annotating text fragment 12851/100392
## 2023-09-16 21:16:07 Annotating text fragment 12861/100392
## 2023-09-16 21:16:07 Annotating text fragment 12871/100392
## 2023-09-16 21:16:07 Annotating text fragment 12881/100392
## 2023-09-16 21:16:07 Annotating text fragment 12891/100392
## 2023-09-16 21:16:07 Annotating text fragment 12901/100392
## 2023-09-16 21:16:07 Annotating text fragment 12911/100392
## 2023-09-16 21:16:07 Annotating text fragment 12921/100392
## 2023-09-16 21:16:08 Annotating text fragment 12931/100392
## 2023-09-16 21:16:08 Annotating text fragment 12941/100392
## 2023-09-16 21:16:08 Annotating text fragment 12951/100392
## 2023-09-16 21:16:08 Annotating text fragment 12961/100392
## 2023-09-16 21:16:08 Annotating text fragment 12971/100392
## 2023-09-16 21:16:08 Annotating text fragment 12981/100392
## 2023-09-16 21:16:08 Annotating text fragment 12991/100392
## 2023-09-16 21:16:08 Annotating text fragment 13001/100392
## 2023-09-16 21:16:08 Annotating text fragment 13011/100392
```

```
## 2023-09-16 21:16:08 Annotating text fragment 13021/100392
## 2023-09-16 21:16:09 Annotating text fragment 13031/100392
## 2023-09-16 21:16:09 Annotating text fragment 13041/100392
## 2023-09-16 21:16:09 Annotating text fragment 13051/100392
## 2023-09-16 21:16:09 Annotating text fragment 13061/100392
## 2023-09-16 21:16:09 Annotating text fragment 13071/100392
## 2023-09-16 21:16:09 Annotating text fragment 13081/100392
## 2023-09-16 21:16:09 Annotating text fragment 13091/100392
## 2023-09-16 21:16:09 Annotating text fragment 13101/100392
## 2023-09-16 21:16:09 Annotating text fragment 13111/100392
## 2023-09-16 21:16:09 Annotating text fragment 13121/100392
## 2023-09-16 21:16:10 Annotating text fragment 13131/100392
## 2023-09-16 21:16:10 Annotating text fragment 13141/100392
## 2023-09-16 21:16:10 Annotating text fragment 13151/100392
## 2023-09-16 21:16:10 Annotating text fragment 13161/100392
## 2023-09-16 21:16:10 Annotating text fragment 13171/100392
## 2023-09-16 21:16:10 Annotating text fragment 13181/100392
## 2023-09-16 21:16:10 Annotating text fragment 13191/100392
## 2023-09-16 21:16:10 Annotating text fragment 13201/100392
## 2023-09-16 21:16:10 Annotating text fragment 13211/100392
## 2023-09-16 21:16:10 Annotating text fragment 13221/100392
## 2023-09-16 21:16:10 Annotating text fragment 13231/100392
## 2023-09-16 21:16:11 Annotating text fragment 13241/100392
## 2023-09-16 21:16:11 Annotating text fragment 13251/100392
## 2023-09-16 21:16:11 Annotating text fragment 13261/100392
## 2023-09-16 21:16:11 Annotating text fragment 13271/100392
## 2023-09-16 21:16:11 Annotating text fragment 13281/100392
## 2023-09-16 21:16:11 Annotating text fragment 13291/100392
## 2023-09-16 21:16:11 Annotating text fragment 13301/100392
## 2023-09-16 21:16:11 Annotating text fragment 13311/100392
## 2023-09-16 21:16:11 Annotating text fragment 13321/100392
## 2023-09-16 21:16:11 Annotating text fragment 13331/100392
## 2023-09-16 21:16:12 Annotating text fragment 13341/100392
## 2023-09-16 21:16:12 Annotating text fragment 13351/100392
## 2023-09-16 21:16:12 Annotating text fragment 13361/100392
## 2023-09-16 21:16:12 Annotating text fragment 13371/100392
## 2023-09-16 21:16:12 Annotating text fragment 13381/100392
## 2023-09-16 21:16:12 Annotating text fragment 13391/100392
## 2023-09-16 21:16:12 Annotating text fragment 13401/100392
## 2023-09-16 21:16:13 Annotating text fragment 13411/100392
## 2023-09-16 21:16:13 Annotating text fragment 13421/100392
## 2023-09-16 21:16:13 Annotating text fragment 13431/100392
## 2023-09-16 21:16:13 Annotating text fragment 13441/100392
## 2023-09-16 21:16:13 Annotating text fragment 13451/100392
## 2023-09-16 21:16:13 Annotating text fragment 13461/100392
## 2023-09-16 21:16:13 Annotating text fragment 13471/100392
## 2023-09-16 21:16:13 Annotating text fragment 13481/100392
## 2023-09-16 21:16:13 Annotating text fragment 13491/100392
## 2023-09-16 21:16:14 Annotating text fragment 13501/100392
## 2023-09-16 21:16:14 Annotating text fragment 13511/100392
## 2023-09-16 21:16:14 Annotating text fragment 13521/100392
## 2023-09-16 21:16:14 Annotating text fragment 13531/100392
## 2023-09-16 21:16:14 Annotating text fragment 13541/100392
## 2023-09-16 21:16:14 Annotating text fragment 13551/100392
```

```
## 2023-09-16 21:16:14 Annotating text fragment 13561/100392
## 2023-09-16 21:16:14 Annotating text fragment 13571/100392
## 2023-09-16 21:16:14 Annotating text fragment 13581/100392
## 2023-09-16 21:16:15 Annotating text fragment 13591/100392
## 2023-09-16 21:16:15 Annotating text fragment 13601/100392
## 2023-09-16 21:16:15 Annotating text fragment 13611/100392
## 2023-09-16 21:16:15 Annotating text fragment 13621/100392
## 2023-09-16 21:16:15 Annotating text fragment 13631/100392
## 2023-09-16 21:16:15 Annotating text fragment 13641/100392
## 2023-09-16 21:16:15 Annotating text fragment 13651/100392
## 2023-09-16 21:16:15 Annotating text fragment 13661/100392
## 2023-09-16 21:16:15 Annotating text fragment 13671/100392
## 2023-09-16 21:16:15 Annotating text fragment 13681/100392
## 2023-09-16 21:16:16 Annotating text fragment 13691/100392
## 2023-09-16 21:16:16 Annotating text fragment 13701/100392
## 2023-09-16 21:16:16 Annotating text fragment 13711/100392
## 2023-09-16 21:16:16 Annotating text fragment 13721/100392
## 2023-09-16 21:16:16 Annotating text fragment 13731/100392
## 2023-09-16 21:16:16 Annotating text fragment 13741/100392
## 2023-09-16 21:16:16 Annotating text fragment 13751/100392
## 2023-09-16 21:16:16 Annotating text fragment 13761/100392
## 2023-09-16 21:16:16 Annotating text fragment 13771/100392
## 2023-09-16 21:16:16 Annotating text fragment 13781/100392
## 2023-09-16 21:16:16 Annotating text fragment 13791/100392
## 2023-09-16 21:16:16 Annotating text fragment 13801/100392
## 2023-09-16 21:16:17 Annotating text fragment 13811/100392
## 2023-09-16 21:16:17 Annotating text fragment 13821/100392
## 2023-09-16 21:16:17 Annotating text fragment 13831/100392
## 2023-09-16 21:16:17 Annotating text fragment 13841/100392
## 2023-09-16 21:16:17 Annotating text fragment 13851/100392
## 2023-09-16 21:16:17 Annotating text fragment 13861/100392
## 2023-09-16 21:16:17 Annotating text fragment 13871/100392
## 2023-09-16 21:16:17 Annotating text fragment 13881/100392
## 2023-09-16 21:16:17 Annotating text fragment 13891/100392
## 2023-09-16 21:16:17 Annotating text fragment 13901/100392
## 2023-09-16 21:16:17 Annotating text fragment 13911/100392
## 2023-09-16 21:16:18 Annotating text fragment 13921/100392
## 2023-09-16 21:16:18 Annotating text fragment 13931/100392
## 2023-09-16 21:16:18 Annotating text fragment 13941/100392
## 2023-09-16 21:16:18 Annotating text fragment 13951/100392
## 2023-09-16 21:16:18 Annotating text fragment 13961/100392
## 2023-09-16 21:16:18 Annotating text fragment 13971/100392
## 2023-09-16 21:16:18 Annotating text fragment 13981/100392
## 2023-09-16 21:16:18 Annotating text fragment 13991/100392
## 2023-09-16 21:16:18 Annotating text fragment 14001/100392
## 2023-09-16 21:16:18 Annotating text fragment 14011/100392
## 2023-09-16 21:16:18 Annotating text fragment 14021/100392
## 2023-09-16 21:16:18 Annotating text fragment 14031/100392
## 2023-09-16 21:16:18 Annotating text fragment 14041/100392
## 2023-09-16 21:16:19 Annotating text fragment 14051/100392
## 2023-09-16 21:16:19 Annotating text fragment 14061/100392
## 2023-09-16 21:16:19 Annotating text fragment 14071/100392
## 2023-09-16 21:16:19 Annotating text fragment 14081/100392
## 2023-09-16 21:16:19 Annotating text fragment 14091/100392
```

```
## 2023-09-16 21:16:19 Annotating text fragment 14101/100392
## 2023-09-16 21:16:19 Annotating text fragment 14111/100392
## 2023-09-16 21:16:19 Annotating text fragment 14121/100392
## 2023-09-16 21:16:19 Annotating text fragment 14131/100392
## 2023-09-16 21:16:19 Annotating text fragment 14141/100392
## 2023-09-16 21:16:20 Annotating text fragment 14151/100392
## 2023-09-16 21:16:20 Annotating text fragment 14161/100392
## 2023-09-16 21:16:20 Annotating text fragment 14171/100392
## 2023-09-16 21:16:20 Annotating text fragment 14181/100392
## 2023-09-16 21:16:20 Annotating text fragment 14191/100392
## 2023-09-16 21:16:20 Annotating text fragment 14201/100392
## 2023-09-16 21:16:20 Annotating text fragment 14211/100392
## 2023-09-16 21:16:20 Annotating text fragment 14221/100392
## 2023-09-16 21:16:20 Annotating text fragment 14231/100392
## 2023-09-16 21:16:20 Annotating text fragment 14241/100392
## 2023-09-16 21:16:21 Annotating text fragment 14251/100392
## 2023-09-16 21:16:21 Annotating text fragment 14261/100392
## 2023-09-16 21:16:21 Annotating text fragment 14271/100392
## 2023-09-16 21:16:21 Annotating text fragment 14281/100392
## 2023-09-16 21:16:21 Annotating text fragment 14291/100392
## 2023-09-16 21:16:21 Annotating text fragment 14301/100392
## 2023-09-16 21:16:21 Annotating text fragment 14311/100392
## 2023-09-16 21:16:22 Annotating text fragment 14321/100392
## 2023-09-16 21:16:22 Annotating text fragment 14331/100392
## 2023-09-16 21:16:22 Annotating text fragment 14341/100392
## 2023-09-16 21:16:22 Annotating text fragment 14351/100392
## 2023-09-16 21:16:22 Annotating text fragment 14361/100392
## 2023-09-16 21:16:22 Annotating text fragment 14371/100392
## 2023-09-16 21:16:22 Annotating text fragment 14381/100392
## 2023-09-16 21:16:22 Annotating text fragment 14391/100392
## 2023-09-16 21:16:22 Annotating text fragment 14401/100392
## 2023-09-16 21:16:22 Annotating text fragment 14411/100392
## 2023-09-16 21:16:23 Annotating text fragment 14421/100392
## 2023-09-16 21:16:23 Annotating text fragment 14431/100392
## 2023-09-16 21:16:23 Annotating text fragment 14441/100392
## 2023-09-16 21:16:23 Annotating text fragment 14451/100392
## 2023-09-16 21:16:23 Annotating text fragment 14461/100392
## 2023-09-16 21:16:23 Annotating text fragment 14471/100392
## 2023-09-16 21:16:23 Annotating text fragment 14481/100392
## 2023-09-16 21:16:23 Annotating text fragment 14491/100392
## 2023-09-16 21:16:23 Annotating text fragment 14501/100392
## 2023-09-16 21:16:23 Annotating text fragment 14511/100392
## 2023-09-16 21:16:24 Annotating text fragment 14521/100392
## 2023-09-16 21:16:24 Annotating text fragment 14531/100392
## 2023-09-16 21:16:24 Annotating text fragment 14541/100392
## 2023-09-16 21:16:24 Annotating text fragment 14551/100392
## 2023-09-16 21:16:24 Annotating text fragment 14561/100392
## 2023-09-16 21:16:24 Annotating text fragment 14571/100392
## 2023-09-16 21:16:24 Annotating text fragment 14581/100392
## 2023-09-16 21:16:24 Annotating text fragment 14591/100392
## 2023-09-16 21:16:24 Annotating text fragment 14601/100392
## 2023-09-16 21:16:24 Annotating text fragment 14611/100392
## 2023-09-16 21:16:24 Annotating text fragment 14621/100392
## 2023-09-16 21:16:25 Annotating text fragment 14631/100392
```

```
## 2023-09-16 21:16:25 Annotating text fragment 14641/100392
## 2023-09-16 21:16:25 Annotating text fragment 14651/100392
## 2023-09-16 21:16:25 Annotating text fragment 14661/100392
## 2023-09-16 21:16:25 Annotating text fragment 14671/100392
## 2023-09-16 21:16:25 Annotating text fragment 14681/100392
## 2023-09-16 21:16:25 Annotating text fragment 14691/100392
## 2023-09-16 21:16:25 Annotating text fragment 14701/100392
## 2023-09-16 21:16:25 Annotating text fragment 14711/100392
## 2023-09-16 21:16:25 Annotating text fragment 14721/100392
## 2023-09-16 21:16:26 Annotating text fragment 14731/100392
## 2023-09-16 21:16:26 Annotating text fragment 14741/100392
## 2023-09-16 21:16:26 Annotating text fragment 14751/100392
## 2023-09-16 21:16:26 Annotating text fragment 14761/100392
## 2023-09-16 21:16:26 Annotating text fragment 14771/100392
## 2023-09-16 21:16:26 Annotating text fragment 14781/100392
## 2023-09-16 21:16:26 Annotating text fragment 14791/100392
## 2023-09-16 21:16:26 Annotating text fragment 14801/100392
## 2023-09-16 21:16:26 Annotating text fragment 14811/100392
## 2023-09-16 21:16:26 Annotating text fragment 14821/100392
## 2023-09-16 21:16:26 Annotating text fragment 14831/100392
## 2023-09-16 21:16:26 Annotating text fragment 14841/100392
## 2023-09-16 21:16:26 Annotating text fragment 14851/100392
## 2023-09-16 21:16:27 Annotating text fragment 14861/100392
## 2023-09-16 21:16:27 Annotating text fragment 14871/100392
## 2023-09-16 21:16:27 Annotating text fragment 14881/100392
## 2023-09-16 21:16:27 Annotating text fragment 14891/100392
## 2023-09-16 21:16:27 Annotating text fragment 14901/100392
## 2023-09-16 21:16:27 Annotating text fragment 14911/100392
## 2023-09-16 21:16:27 Annotating text fragment 14921/100392
## 2023-09-16 21:16:27 Annotating text fragment 14931/100392
## 2023-09-16 21:16:27 Annotating text fragment 14941/100392
## 2023-09-16 21:16:27 Annotating text fragment 14951/100392
## 2023-09-16 21:16:27 Annotating text fragment 14961/100392
## 2023-09-16 21:16:27 Annotating text fragment 14971/100392
## 2023-09-16 21:16:28 Annotating text fragment 14981/100392
## 2023-09-16 21:16:28 Annotating text fragment 14991/100392
## 2023-09-16 21:16:28 Annotating text fragment 15001/100392
## 2023-09-16 21:16:28 Annotating text fragment 15011/100392
## 2023-09-16 21:16:28 Annotating text fragment 15021/100392
## 2023-09-16 21:16:28 Annotating text fragment 15031/100392
## 2023-09-16 21:16:28 Annotating text fragment 15041/100392
## 2023-09-16 21:16:28 Annotating text fragment 15051/100392
## 2023-09-16 21:16:28 Annotating text fragment 15061/100392
## 2023-09-16 21:16:29 Annotating text fragment 15071/100392
## 2023-09-16 21:16:29 Annotating text fragment 15081/100392
## 2023-09-16 21:16:29 Annotating text fragment 15091/100392
## 2023-09-16 21:16:29 Annotating text fragment 15101/100392
## 2023-09-16 21:16:29 Annotating text fragment 15111/100392
## 2023-09-16 21:16:29 Annotating text fragment 15121/100392
## 2023-09-16 21:16:29 Annotating text fragment 15131/100392
## 2023-09-16 21:16:29 Annotating text fragment 15141/100392
## 2023-09-16 21:16:29 Annotating text fragment 15151/100392
## 2023-09-16 21:16:30 Annotating text fragment 15161/100392
## 2023-09-16 21:16:30 Annotating text fragment 15171/100392
```

```
## 2023-09-16 21:16:30 Annotating text fragment 15181/100392
## 2023-09-16 21:16:30 Annotating text fragment 15191/100392
## 2023-09-16 21:16:30 Annotating text fragment 15201/100392
## 2023-09-16 21:16:30 Annotating text fragment 15211/100392
## 2023-09-16 21:16:31 Annotating text fragment 15221/100392
## 2023-09-16 21:16:31 Annotating text fragment 15231/100392
## 2023-09-16 21:16:31 Annotating text fragment 15241/100392
## 2023-09-16 21:16:31 Annotating text fragment 15251/100392
## 2023-09-16 21:16:31 Annotating text fragment 15261/100392
## 2023-09-16 21:16:31 Annotating text fragment 15271/100392
## 2023-09-16 21:16:31 Annotating text fragment 15281/100392
## 2023-09-16 21:16:31 Annotating text fragment 15291/100392
## 2023-09-16 21:16:31 Annotating text fragment 15301/100392
## 2023-09-16 21:16:31 Annotating text fragment 15311/100392
## 2023-09-16 21:16:31 Annotating text fragment 15321/100392
## 2023-09-16 21:16:32 Annotating text fragment 15331/100392
## 2023-09-16 21:16:32 Annotating text fragment 15341/100392
## 2023-09-16 21:16:32 Annotating text fragment 15351/100392
## 2023-09-16 21:16:32 Annotating text fragment 15361/100392
## 2023-09-16 21:16:32 Annotating text fragment 15371/100392
## 2023-09-16 21:16:32 Annotating text fragment 15381/100392
## 2023-09-16 21:16:32 Annotating text fragment 15391/100392
## 2023-09-16 21:16:32 Annotating text fragment 15401/100392
## 2023-09-16 21:16:32 Annotating text fragment 15411/100392
## 2023-09-16 21:16:32 Annotating text fragment 15421/100392
## 2023-09-16 21:16:33 Annotating text fragment 15431/100392
## 2023-09-16 21:16:33 Annotating text fragment 15441/100392
## 2023-09-16 21:16:33 Annotating text fragment 15451/100392
## 2023-09-16 21:16:33 Annotating text fragment 15461/100392
## 2023-09-16 21:16:33 Annotating text fragment 15471/100392
## 2023-09-16 21:16:33 Annotating text fragment 15481/100392
## 2023-09-16 21:16:33 Annotating text fragment 15491/100392
## 2023-09-16 21:16:33 Annotating text fragment 15501/100392
## 2023-09-16 21:16:33 Annotating text fragment 15511/100392
## 2023-09-16 21:16:33 Annotating text fragment 15521/100392
## 2023-09-16 21:16:34 Annotating text fragment 15531/100392
## 2023-09-16 21:16:34 Annotating text fragment 15541/100392
## 2023-09-16 21:16:34 Annotating text fragment 15551/100392
## 2023-09-16 21:16:34 Annotating text fragment 15561/100392
## 2023-09-16 21:16:34 Annotating text fragment 15571/100392
## 2023-09-16 21:16:34 Annotating text fragment 15581/100392
## 2023-09-16 21:16:34 Annotating text fragment 15591/100392
## 2023-09-16 21:16:34 Annotating text fragment 15601/100392
## 2023-09-16 21:16:35 Annotating text fragment 15611/100392
## 2023-09-16 21:16:35 Annotating text fragment 15621/100392
## 2023-09-16 21:16:35 Annotating text fragment 15631/100392
## 2023-09-16 21:16:35 Annotating text fragment 15641/100392
## 2023-09-16 21:16:35 Annotating text fragment 15651/100392
## 2023-09-16 21:16:35 Annotating text fragment 15661/100392
## 2023-09-16 21:16:35 Annotating text fragment 15671/100392
## 2023-09-16 21:16:35 Annotating text fragment 15681/100392
## 2023-09-16 21:16:35 Annotating text fragment 15691/100392
## 2023-09-16 21:16:36 Annotating text fragment 15701/100392
## 2023-09-16 21:16:36 Annotating text fragment 15711/100392
```

```
## 2023-09-16 21:16:36 Annotating text fragment 15721/100392
## 2023-09-16 21:16:36 Annotating text fragment 15731/100392
## 2023-09-16 21:16:36 Annotating text fragment 15741/100392
## 2023-09-16 21:16:36 Annotating text fragment 15751/100392
## 2023-09-16 21:16:36 Annotating text fragment 15761/100392
## 2023-09-16 21:16:36 Annotating text fragment 15771/100392
## 2023-09-16 21:16:37 Annotating text fragment 15781/100392
## 2023-09-16 21:16:37 Annotating text fragment 15791/100392
## 2023-09-16 21:16:37 Annotating text fragment 15801/100392
## 2023-09-16 21:16:37 Annotating text fragment 15811/100392
## 2023-09-16 21:16:37 Annotating text fragment 15821/100392
## 2023-09-16 21:16:37 Annotating text fragment 15831/100392
## 2023-09-16 21:16:37 Annotating text fragment 15841/100392
## 2023-09-16 21:16:37 Annotating text fragment 15851/100392
## 2023-09-16 21:16:37 Annotating text fragment 15861/100392
## 2023-09-16 21:16:37 Annotating text fragment 15871/100392
## 2023-09-16 21:16:37 Annotating text fragment 15881/100392
## 2023-09-16 21:16:37 Annotating text fragment 15891/100392
## 2023-09-16 21:16:38 Annotating text fragment 15901/100392
## 2023-09-16 21:16:38 Annotating text fragment 15911/100392
## 2023-09-16 21:16:38 Annotating text fragment 15921/100392
## 2023-09-16 21:16:38 Annotating text fragment 15931/100392
## 2023-09-16 21:16:38 Annotating text fragment 15941/100392
## 2023-09-16 21:16:38 Annotating text fragment 15951/100392
## 2023-09-16 21:16:38 Annotating text fragment 15961/100392
## 2023-09-16 21:16:38 Annotating text fragment 15971/100392
## 2023-09-16 21:16:38 Annotating text fragment 15981/100392
## 2023-09-16 21:16:38 Annotating text fragment 15991/100392
## 2023-09-16 21:16:38 Annotating text fragment 16001/100392
## 2023-09-16 21:16:39 Annotating text fragment 16011/100392
## 2023-09-16 21:16:39 Annotating text fragment 16021/100392
## 2023-09-16 21:16:39 Annotating text fragment 16031/100392
## 2023-09-16 21:16:39 Annotating text fragment 16041/100392
## 2023-09-16 21:16:39 Annotating text fragment 16051/100392
## 2023-09-16 21:16:39 Annotating text fragment 16061/100392
## 2023-09-16 21:16:39 Annotating text fragment 16071/100392
## 2023-09-16 21:16:39 Annotating text fragment 16081/100392
## 2023-09-16 21:16:39 Annotating text fragment 16091/100392
## 2023-09-16 21:16:39 Annotating text fragment 16101/100392
## 2023-09-16 21:16:40 Annotating text fragment 16111/100392
## 2023-09-16 21:16:40 Annotating text fragment 16121/100392
## 2023-09-16 21:16:40 Annotating text fragment 16131/100392
## 2023-09-16 21:16:40 Annotating text fragment 16141/100392
## 2023-09-16 21:16:40 Annotating text fragment 16151/100392
## 2023-09-16 21:16:40 Annotating text fragment 16161/100392
## 2023-09-16 21:16:40 Annotating text fragment 16171/100392
## 2023-09-16 21:16:40 Annotating text fragment 16181/100392
## 2023-09-16 21:16:40 Annotating text fragment 16191/100392
## 2023-09-16 21:16:40 Annotating text fragment 16201/100392
## 2023-09-16 21:16:40 Annotating text fragment 16211/100392
## 2023-09-16 21:16:41 Annotating text fragment 16221/100392
## 2023-09-16 21:16:41 Annotating text fragment 16231/100392
## 2023-09-16 21:16:41 Annotating text fragment 16241/100392
## 2023-09-16 21:16:41 Annotating text fragment 16251/100392
```

```
## 2023-09-16 21:16:41 Annotating text fragment 16261/100392
## 2023-09-16 21:16:41 Annotating text fragment 16271/100392
## 2023-09-16 21:16:41 Annotating text fragment 16281/100392
## 2023-09-16 21:16:41 Annotating text fragment 16291/100392
## 2023-09-16 21:16:41 Annotating text fragment 16301/100392
## 2023-09-16 21:16:41 Annotating text fragment 16311/100392
## 2023-09-16 21:16:41 Annotating text fragment 16321/100392
## 2023-09-16 21:16:42 Annotating text fragment 16331/100392
## 2023-09-16 21:16:42 Annotating text fragment 16341/100392
## 2023-09-16 21:16:42 Annotating text fragment 16351/100392
## 2023-09-16 21:16:42 Annotating text fragment 16361/100392
## 2023-09-16 21:16:42 Annotating text fragment 16371/100392
## 2023-09-16 21:16:42 Annotating text fragment 16381/100392
## 2023-09-16 21:16:42 Annotating text fragment 16391/100392
## 2023-09-16 21:16:42 Annotating text fragment 16401/100392
## 2023-09-16 21:16:43 Annotating text fragment 16411/100392
## 2023-09-16 21:16:43 Annotating text fragment 16421/100392
## 2023-09-16 21:16:43 Annotating text fragment 16431/100392
## 2023-09-16 21:16:43 Annotating text fragment 16441/100392
## 2023-09-16 21:16:43 Annotating text fragment 16451/100392
## 2023-09-16 21:16:43 Annotating text fragment 16461/100392
## 2023-09-16 21:16:43 Annotating text fragment 16471/100392
## 2023-09-16 21:16:43 Annotating text fragment 16481/100392
## 2023-09-16 21:16:44 Annotating text fragment 16491/100392
## 2023-09-16 21:16:44 Annotating text fragment 16501/100392
## 2023-09-16 21:16:44 Annotating text fragment 16511/100392
## 2023-09-16 21:16:44 Annotating text fragment 16521/100392
## 2023-09-16 21:16:44 Annotating text fragment 16531/100392
## 2023-09-16 21:16:44 Annotating text fragment 16541/100392
## 2023-09-16 21:16:44 Annotating text fragment 16551/100392
## 2023-09-16 21:16:44 Annotating text fragment 16561/100392
## 2023-09-16 21:16:44 Annotating text fragment 16571/100392
## 2023-09-16 21:16:44 Annotating text fragment 16581/100392
## 2023-09-16 21:16:44 Annotating text fragment 16591/100392
## 2023-09-16 21:16:45 Annotating text fragment 16601/100392
## 2023-09-16 21:16:45 Annotating text fragment 16611/100392
## 2023-09-16 21:16:45 Annotating text fragment 16621/100392
## 2023-09-16 21:16:45 Annotating text fragment 16631/100392
## 2023-09-16 21:16:45 Annotating text fragment 16641/100392
## 2023-09-16 21:16:45 Annotating text fragment 16651/100392
## 2023-09-16 21:16:45 Annotating text fragment 16661/100392
## 2023-09-16 21:16:45 Annotating text fragment 16671/100392
## 2023-09-16 21:16:45 Annotating text fragment 16681/100392
## 2023-09-16 21:16:46 Annotating text fragment 16691/100392
## 2023-09-16 21:16:46 Annotating text fragment 16701/100392
## 2023-09-16 21:16:46 Annotating text fragment 16711/100392
## 2023-09-16 21:16:46 Annotating text fragment 16721/100392
## 2023-09-16 21:16:46 Annotating text fragment 16731/100392
## 2023-09-16 21:16:46 Annotating text fragment 16741/100392
## 2023-09-16 21:16:46 Annotating text fragment 16751/100392
## 2023-09-16 21:16:46 Annotating text fragment 16761/100392
## 2023-09-16 21:16:46 Annotating text fragment 16771/100392
## 2023-09-16 21:16:46 Annotating text fragment 16781/100392
## 2023-09-16 21:16:47 Annotating text fragment 16791/100392
```

```
## 2023-09-16 21:16:47 Annotating text fragment 16801/100392
## 2023-09-16 21:16:47 Annotating text fragment 16811/100392
## 2023-09-16 21:16:47 Annotating text fragment 16821/100392
## 2023-09-16 21:16:47 Annotating text fragment 16831/100392
## 2023-09-16 21:16:47 Annotating text fragment 16841/100392
## 2023-09-16 21:16:47 Annotating text fragment 16851/100392
## 2023-09-16 21:16:47 Annotating text fragment 16861/100392
## 2023-09-16 21:16:47 Annotating text fragment 16871/100392
## 2023-09-16 21:16:47 Annotating text fragment 16881/100392
## 2023-09-16 21:16:47 Annotating text fragment 16891/100392
## 2023-09-16 21:16:47 Annotating text fragment 16901/100392
## 2023-09-16 21:16:48 Annotating text fragment 16911/100392
## 2023-09-16 21:16:48 Annotating text fragment 16921/100392
## 2023-09-16 21:16:48 Annotating text fragment 16931/100392
## 2023-09-16 21:16:48 Annotating text fragment 16941/100392
## 2023-09-16 21:16:48 Annotating text fragment 16951/100392
## 2023-09-16 21:16:48 Annotating text fragment 16961/100392
## 2023-09-16 21:16:48 Annotating text fragment 16971/100392
## 2023-09-16 21:16:48 Annotating text fragment 16981/100392
## 2023-09-16 21:16:48 Annotating text fragment 16991/100392
## 2023-09-16 21:16:48 Annotating text fragment 17001/100392
## 2023-09-16 21:16:48 Annotating text fragment 17011/100392
## 2023-09-16 21:16:49 Annotating text fragment 17021/100392
## 2023-09-16 21:16:49 Annotating text fragment 17031/100392
## 2023-09-16 21:16:49 Annotating text fragment 17041/100392
## 2023-09-16 21:16:49 Annotating text fragment 17051/100392
## 2023-09-16 21:16:49 Annotating text fragment 17061/100392
## 2023-09-16 21:16:49 Annotating text fragment 17071/100392
## 2023-09-16 21:16:49 Annotating text fragment 17081/100392
## 2023-09-16 21:16:49 Annotating text fragment 17091/100392
## 2023-09-16 21:16:49 Annotating text fragment 17101/100392
## 2023-09-16 21:16:49 Annotating text fragment 17111/100392
## 2023-09-16 21:16:49 Annotating text fragment 17121/100392
## 2023-09-16 21:16:50 Annotating text fragment 17131/100392
## 2023-09-16 21:16:50 Annotating text fragment 17141/100392
## 2023-09-16 21:16:50 Annotating text fragment 17151/100392
## 2023-09-16 21:16:50 Annotating text fragment 17161/100392
## 2023-09-16 21:16:50 Annotating text fragment 17171/100392
## 2023-09-16 21:16:50 Annotating text fragment 17181/100392
## 2023-09-16 21:16:50 Annotating text fragment 17191/100392
## 2023-09-16 21:16:50 Annotating text fragment 17201/100392
## 2023-09-16 21:16:50 Annotating text fragment 17211/100392
## 2023-09-16 21:16:50 Annotating text fragment 17221/100392
## 2023-09-16 21:16:50 Annotating text fragment 17231/100392
## 2023-09-16 21:16:50 Annotating text fragment 17241/100392
## 2023-09-16 21:16:51 Annotating text fragment 17251/100392
## 2023-09-16 21:16:51 Annotating text fragment 17261/100392
## 2023-09-16 21:16:51 Annotating text fragment 17271/100392
## 2023-09-16 21:16:51 Annotating text fragment 17281/100392
## 2023-09-16 21:16:51 Annotating text fragment 17291/100392
## 2023-09-16 21:16:51 Annotating text fragment 17301/100392
## 2023-09-16 21:16:51 Annotating text fragment 17311/100392
## 2023-09-16 21:16:51 Annotating text fragment 17321/100392
## 2023-09-16 21:16:51 Annotating text fragment 17331/100392
```

```
## 2023-09-16 21:16:52 Annotating text fragment 17341/100392
## 2023-09-16 21:16:52 Annotating text fragment 17351/100392
## 2023-09-16 21:16:52 Annotating text fragment 17361/100392
## 2023-09-16 21:16:52 Annotating text fragment 17371/100392
## 2023-09-16 21:16:52 Annotating text fragment 17381/100392
## 2023-09-16 21:16:52 Annotating text fragment 17391/100392
## 2023-09-16 21:16:52 Annotating text fragment 17401/100392
## 2023-09-16 21:16:52 Annotating text fragment 17411/100392
## 2023-09-16 21:16:52 Annotating text fragment 17421/100392
## 2023-09-16 21:16:53 Annotating text fragment 17431/100392
## 2023-09-16 21:16:53 Annotating text fragment 17441/100392
## 2023-09-16 21:16:53 Annotating text fragment 17451/100392
## 2023-09-16 21:16:53 Annotating text fragment 17461/100392
## 2023-09-16 21:16:53 Annotating text fragment 17471/100392
## 2023-09-16 21:16:53 Annotating text fragment 17481/100392
## 2023-09-16 21:16:53 Annotating text fragment 17491/100392
## 2023-09-16 21:16:53 Annotating text fragment 17501/100392
## 2023-09-16 21:16:53 Annotating text fragment 17511/100392
## 2023-09-16 21:16:53 Annotating text fragment 17521/100392
## 2023-09-16 21:16:54 Annotating text fragment 17531/100392
## 2023-09-16 21:16:54 Annotating text fragment 17541/100392
## 2023-09-16 21:16:54 Annotating text fragment 17551/100392
## 2023-09-16 21:16:54 Annotating text fragment 17561/100392
## 2023-09-16 21:16:54 Annotating text fragment 17571/100392
## 2023-09-16 21:16:54 Annotating text fragment 17581/100392
## 2023-09-16 21:16:54 Annotating text fragment 17591/100392
## 2023-09-16 21:16:54 Annotating text fragment 17601/100392
## 2023-09-16 21:16:54 Annotating text fragment 17611/100392
## 2023-09-16 21:16:54 Annotating text fragment 17621/100392
## 2023-09-16 21:16:54 Annotating text fragment 17631/100392
## 2023-09-16 21:16:55 Annotating text fragment 17641/100392
## 2023-09-16 21:16:55 Annotating text fragment 17651/100392
## 2023-09-16 21:16:55 Annotating text fragment 17661/100392
## 2023-09-16 21:16:55 Annotating text fragment 17671/100392
## 2023-09-16 21:16:55 Annotating text fragment 17681/100392
## 2023-09-16 21:16:55 Annotating text fragment 17691/100392
## 2023-09-16 21:16:55 Annotating text fragment 17701/100392
## 2023-09-16 21:16:55 Annotating text fragment 17711/100392
## 2023-09-16 21:16:55 Annotating text fragment 17721/100392
## 2023-09-16 21:16:55 Annotating text fragment 17731/100392
## 2023-09-16 21:16:56 Annotating text fragment 17741/100392
## 2023-09-16 21:16:56 Annotating text fragment 17751/100392
## 2023-09-16 21:16:56 Annotating text fragment 17761/100392
## 2023-09-16 21:16:56 Annotating text fragment 17771/100392
## 2023-09-16 21:16:56 Annotating text fragment 17781/100392
## 2023-09-16 21:16:56 Annotating text fragment 17791/100392
## 2023-09-16 21:16:56 Annotating text fragment 17801/100392
## 2023-09-16 21:16:56 Annotating text fragment 17811/100392
## 2023-09-16 21:16:56 Annotating text fragment 17821/100392
## 2023-09-16 21:16:56 Annotating text fragment 17831/100392
## 2023-09-16 21:16:56 Annotating text fragment 17841/100392
## 2023-09-16 21:16:57 Annotating text fragment 17851/100392
## 2023-09-16 21:16:57 Annotating text fragment 17861/100392
## 2023-09-16 21:16:57 Annotating text fragment 17871/100392
```

```
## 2023-09-16 21:16:57 Annotating text fragment 17881/100392
## 2023-09-16 21:16:57 Annotating text fragment 17891/100392
## 2023-09-16 21:16:57 Annotating text fragment 17901/100392
## 2023-09-16 21:16:57 Annotating text fragment 17911/100392
## 2023-09-16 21:16:57 Annotating text fragment 17921/100392
## 2023-09-16 21:16:57 Annotating text fragment 17931/100392
## 2023-09-16 21:16:57 Annotating text fragment 17941/100392
## 2023-09-16 21:16:57 Annotating text fragment 17951/100392
## 2023-09-16 21:16:57 Annotating text fragment 17961/100392
## 2023-09-16 21:16:57 Annotating text fragment 17971/100392
## 2023-09-16 21:16:58 Annotating text fragment 17981/100392
## 2023-09-16 21:16:58 Annotating text fragment 17991/100392
## 2023-09-16 21:16:58 Annotating text fragment 18001/100392
## 2023-09-16 21:16:58 Annotating text fragment 18011/100392
## 2023-09-16 21:16:58 Annotating text fragment 18021/100392
## 2023-09-16 21:16:58 Annotating text fragment 18031/100392
## 2023-09-16 21:16:58 Annotating text fragment 18041/100392
## 2023-09-16 21:16:58 Annotating text fragment 18051/100392
## 2023-09-16 21:16:58 Annotating text fragment 18061/100392
## 2023-09-16 21:16:58 Annotating text fragment 18071/100392
## 2023-09-16 21:16:59 Annotating text fragment 18081/100392
## 2023-09-16 21:16:59 Annotating text fragment 18091/100392
## 2023-09-16 21:16:59 Annotating text fragment 18101/100392
## 2023-09-16 21:16:59 Annotating text fragment 18111/100392
## 2023-09-16 21:16:59 Annotating text fragment 18121/100392
## 2023-09-16 21:16:59 Annotating text fragment 18131/100392
## 2023-09-16 21:16:59 Annotating text fragment 18141/100392
## 2023-09-16 21:16:59 Annotating text fragment 18151/100392
## 2023-09-16 21:16:59 Annotating text fragment 18161/100392
## 2023-09-16 21:16:59 Annotating text fragment 18171/100392
## 2023-09-16 21:16:59 Annotating text fragment 18181/100392
## 2023-09-16 21:16:59 Annotating text fragment 18191/100392
## 2023-09-16 21:17:00 Annotating text fragment 18201/100392
## 2023-09-16 21:17:00 Annotating text fragment 18211/100392
## 2023-09-16 21:17:00 Annotating text fragment 18221/100392
## 2023-09-16 21:17:00 Annotating text fragment 18231/100392
## 2023-09-16 21:17:00 Annotating text fragment 18241/100392
## 2023-09-16 21:17:00 Annotating text fragment 18251/100392
## 2023-09-16 21:17:00 Annotating text fragment 18261/100392
## 2023-09-16 21:17:00 Annotating text fragment 18271/100392
## 2023-09-16 21:17:00 Annotating text fragment 18281/100392
## 2023-09-16 21:17:00 Annotating text fragment 18291/100392
## 2023-09-16 21:17:01 Annotating text fragment 18301/100392
## 2023-09-16 21:17:01 Annotating text fragment 18311/100392
## 2023-09-16 21:17:01 Annotating text fragment 18321/100392
## 2023-09-16 21:17:01 Annotating text fragment 18331/100392
## 2023-09-16 21:17:01 Annotating text fragment 18341/100392
## 2023-09-16 21:17:01 Annotating text fragment 18351/100392
## 2023-09-16 21:17:01 Annotating text fragment 18361/100392
## 2023-09-16 21:17:01 Annotating text fragment 18371/100392
## 2023-09-16 21:17:01 Annotating text fragment 18381/100392
## 2023-09-16 21:17:01 Annotating text fragment 18391/100392
## 2023-09-16 21:17:02 Annotating text fragment 18401/100392
## 2023-09-16 21:17:02 Annotating text fragment 18411/100392
```

```
## 2023-09-16 21:17:02 Annotating text fragment 18421/100392
## 2023-09-16 21:17:02 Annotating text fragment 18431/100392
## 2023-09-16 21:17:02 Annotating text fragment 18441/100392
## 2023-09-16 21:17:02 Annotating text fragment 18451/100392
## 2023-09-16 21:17:02 Annotating text fragment 18461/100392
## 2023-09-16 21:17:02 Annotating text fragment 18471/100392
## 2023-09-16 21:17:02 Annotating text fragment 18481/100392
## 2023-09-16 21:17:03 Annotating text fragment 18491/100392
## 2023-09-16 21:17:03 Annotating text fragment 18501/100392
## 2023-09-16 21:17:03 Annotating text fragment 18511/100392
## 2023-09-16 21:17:03 Annotating text fragment 18521/100392
## 2023-09-16 21:17:03 Annotating text fragment 18531/100392
## 2023-09-16 21:17:03 Annotating text fragment 18541/100392
## 2023-09-16 21:17:03 Annotating text fragment 18551/100392
## 2023-09-16 21:17:03 Annotating text fragment 18561/100392
## 2023-09-16 21:17:03 Annotating text fragment 18571/100392
## 2023-09-16 21:17:03 Annotating text fragment 18581/100392
## 2023-09-16 21:17:03 Annotating text fragment 18591/100392
## 2023-09-16 21:17:03 Annotating text fragment 18601/100392
## 2023-09-16 21:17:04 Annotating text fragment 18611/100392
## 2023-09-16 21:17:04 Annotating text fragment 18621/100392
## 2023-09-16 21:17:04 Annotating text fragment 18631/100392
## 2023-09-16 21:17:04 Annotating text fragment 18641/100392
## 2023-09-16 21:17:04 Annotating text fragment 18651/100392
## 2023-09-16 21:17:04 Annotating text fragment 18661/100392
## 2023-09-16 21:17:04 Annotating text fragment 18671/100392
## 2023-09-16 21:17:04 Annotating text fragment 18681/100392
## 2023-09-16 21:17:04 Annotating text fragment 18691/100392
## 2023-09-16 21:17:04 Annotating text fragment 18701/100392
## 2023-09-16 21:17:04 Annotating text fragment 18711/100392
## 2023-09-16 21:17:05 Annotating text fragment 18721/100392
## 2023-09-16 21:17:05 Annotating text fragment 18731/100392
## 2023-09-16 21:17:05 Annotating text fragment 18741/100392
## 2023-09-16 21:17:05 Annotating text fragment 18751/100392
## 2023-09-16 21:17:05 Annotating text fragment 18761/100392
## 2023-09-16 21:17:05 Annotating text fragment 18771/100392
## 2023-09-16 21:17:05 Annotating text fragment 18781/100392
## 2023-09-16 21:17:05 Annotating text fragment 18791/100392
## 2023-09-16 21:17:05 Annotating text fragment 18801/100392
## 2023-09-16 21:17:05 Annotating text fragment 18811/100392
## 2023-09-16 21:17:06 Annotating text fragment 18821/100392
## 2023-09-16 21:17:06 Annotating text fragment 18831/100392
## 2023-09-16 21:17:06 Annotating text fragment 18841/100392
## 2023-09-16 21:17:06 Annotating text fragment 18851/100392
## 2023-09-16 21:17:06 Annotating text fragment 18861/100392
## 2023-09-16 21:17:06 Annotating text fragment 18871/100392
## 2023-09-16 21:17:06 Annotating text fragment 18881/100392
## 2023-09-16 21:17:06 Annotating text fragment 18891/100392
## 2023-09-16 21:17:06 Annotating text fragment 18901/100392
## 2023-09-16 21:17:06 Annotating text fragment 18911/100392
## 2023-09-16 21:17:07 Annotating text fragment 18921/100392
## 2023-09-16 21:17:07 Annotating text fragment 18931/100392
## 2023-09-16 21:17:07 Annotating text fragment 18941/100392
## 2023-09-16 21:17:07 Annotating text fragment 18951/100392
```

```
## 2023-09-16 21:17:07 Annotating text fragment 18961/100392
## 2023-09-16 21:17:07 Annotating text fragment 18971/100392
## 2023-09-16 21:17:07 Annotating text fragment 18981/100392
## 2023-09-16 21:17:07 Annotating text fragment 18991/100392
## 2023-09-16 21:17:07 Annotating text fragment 19001/100392
## 2023-09-16 21:17:07 Annotating text fragment 19011/100392
## 2023-09-16 21:17:08 Annotating text fragment 19021/100392
## 2023-09-16 21:17:08 Annotating text fragment 19031/100392
## 2023-09-16 21:17:08 Annotating text fragment 19041/100392
## 2023-09-16 21:17:08 Annotating text fragment 19051/100392
## 2023-09-16 21:17:08 Annotating text fragment 19061/100392
## 2023-09-16 21:17:08 Annotating text fragment 19071/100392
## 2023-09-16 21:17:08 Annotating text fragment 19081/100392
## 2023-09-16 21:17:09 Annotating text fragment 19091/100392
## 2023-09-16 21:17:09 Annotating text fragment 19101/100392
## 2023-09-16 21:17:09 Annotating text fragment 19111/100392
## 2023-09-16 21:17:09 Annotating text fragment 19121/100392
## 2023-09-16 21:17:09 Annotating text fragment 19131/100392
## 2023-09-16 21:17:09 Annotating text fragment 19141/100392
## 2023-09-16 21:17:09 Annotating text fragment 19151/100392
## 2023-09-16 21:17:09 Annotating text fragment 19161/100392
## 2023-09-16 21:17:09 Annotating text fragment 19171/100392
## 2023-09-16 21:17:09 Annotating text fragment 19181/100392
## 2023-09-16 21:17:09 Annotating text fragment 19191/100392
## 2023-09-16 21:17:10 Annotating text fragment 19201/100392
## 2023-09-16 21:17:10 Annotating text fragment 19211/100392
## 2023-09-16 21:17:10 Annotating text fragment 19221/100392
## 2023-09-16 21:17:10 Annotating text fragment 19231/100392
## 2023-09-16 21:17:10 Annotating text fragment 19241/100392
## 2023-09-16 21:17:10 Annotating text fragment 19251/100392
## 2023-09-16 21:17:10 Annotating text fragment 19261/100392
## 2023-09-16 21:17:10 Annotating text fragment 19271/100392
## 2023-09-16 21:17:10 Annotating text fragment 19281/100392
## 2023-09-16 21:17:10 Annotating text fragment 19291/100392
## 2023-09-16 21:17:10 Annotating text fragment 19301/100392
## 2023-09-16 21:17:11 Annotating text fragment 19311/100392
## 2023-09-16 21:17:11 Annotating text fragment 19321/100392
## 2023-09-16 21:17:11 Annotating text fragment 19331/100392
## 2023-09-16 21:17:11 Annotating text fragment 19341/100392
## 2023-09-16 21:17:11 Annotating text fragment 19351/100392
## 2023-09-16 21:17:11 Annotating text fragment 19361/100392
## 2023-09-16 21:17:11 Annotating text fragment 19371/100392
## 2023-09-16 21:17:11 Annotating text fragment 19381/100392
## 2023-09-16 21:17:11 Annotating text fragment 19391/100392
## 2023-09-16 21:17:12 Annotating text fragment 19401/100392
## 2023-09-16 21:17:12 Annotating text fragment 19411/100392
## 2023-09-16 21:17:12 Annotating text fragment 19421/100392
## 2023-09-16 21:17:12 Annotating text fragment 19431/100392
## 2023-09-16 21:17:12 Annotating text fragment 19441/100392
## 2023-09-16 21:17:12 Annotating text fragment 19451/100392
## 2023-09-16 21:17:12 Annotating text fragment 19461/100392
## 2023-09-16 21:17:12 Annotating text fragment 19471/100392
## 2023-09-16 21:17:12 Annotating text fragment 19481/100392
## 2023-09-16 21:17:12 Annotating text fragment 19491/100392
```

```
## 2023-09-16 21:17:12 Annotating text fragment 19501/100392
## 2023-09-16 21:17:13 Annotating text fragment 19511/100392
## 2023-09-16 21:17:13 Annotating text fragment 19521/100392
## 2023-09-16 21:17:13 Annotating text fragment 19531/100392
## 2023-09-16 21:17:13 Annotating text fragment 19541/100392
## 2023-09-16 21:17:13 Annotating text fragment 19551/100392
## 2023-09-16 21:17:13 Annotating text fragment 19561/100392
## 2023-09-16 21:17:13 Annotating text fragment 19571/100392
## 2023-09-16 21:17:13 Annotating text fragment 19581/100392
## 2023-09-16 21:17:13 Annotating text fragment 19591/100392
## 2023-09-16 21:17:13 Annotating text fragment 19601/100392
## 2023-09-16 21:17:13 Annotating text fragment 19611/100392
## 2023-09-16 21:17:14 Annotating text fragment 19621/100392
## 2023-09-16 21:17:14 Annotating text fragment 19631/100392
## 2023-09-16 21:17:14 Annotating text fragment 19641/100392
## 2023-09-16 21:17:14 Annotating text fragment 19651/100392
## 2023-09-16 21:17:14 Annotating text fragment 19661/100392
## 2023-09-16 21:17:14 Annotating text fragment 19671/100392
## 2023-09-16 21:17:14 Annotating text fragment 19681/100392
## 2023-09-16 21:17:14 Annotating text fragment 19691/100392
## 2023-09-16 21:17:14 Annotating text fragment 19701/100392
## 2023-09-16 21:17:14 Annotating text fragment 19711/100392
## 2023-09-16 21:17:14 Annotating text fragment 19721/100392
## 2023-09-16 21:17:15 Annotating text fragment 19731/100392
## 2023-09-16 21:17:15 Annotating text fragment 19741/100392
## 2023-09-16 21:17:15 Annotating text fragment 19751/100392
## 2023-09-16 21:17:15 Annotating text fragment 19761/100392
## 2023-09-16 21:17:15 Annotating text fragment 19771/100392
## 2023-09-16 21:17:15 Annotating text fragment 19781/100392
## 2023-09-16 21:17:15 Annotating text fragment 19791/100392
## 2023-09-16 21:17:15 Annotating text fragment 19801/100392
## 2023-09-16 21:17:15 Annotating text fragment 19811/100392
## 2023-09-16 21:17:15 Annotating text fragment 19821/100392
## 2023-09-16 21:17:16 Annotating text fragment 19831/100392
## 2023-09-16 21:17:16 Annotating text fragment 19841/100392
## 2023-09-16 21:17:16 Annotating text fragment 19851/100392
## 2023-09-16 21:17:16 Annotating text fragment 19861/100392
## 2023-09-16 21:17:16 Annotating text fragment 19871/100392
## 2023-09-16 21:17:16 Annotating text fragment 19881/100392
## 2023-09-16 21:17:16 Annotating text fragment 19891/100392
## 2023-09-16 21:17:16 Annotating text fragment 19901/100392
## 2023-09-16 21:17:16 Annotating text fragment 19911/100392
## 2023-09-16 21:17:17 Annotating text fragment 19921/100392
## 2023-09-16 21:17:17 Annotating text fragment 19931/100392
## 2023-09-16 21:17:17 Annotating text fragment 19941/100392
## 2023-09-16 21:17:17 Annotating text fragment 19951/100392
## 2023-09-16 21:17:17 Annotating text fragment 19961/100392
## 2023-09-16 21:17:17 Annotating text fragment 19971/100392
## 2023-09-16 21:17:17 Annotating text fragment 19981/100392
## 2023-09-16 21:17:17 Annotating text fragment 19991/100392
## 2023-09-16 21:17:17 Annotating text fragment 20001/100392
## 2023-09-16 21:17:17 Annotating text fragment 20011/100392
## 2023-09-16 21:17:17 Annotating text fragment 20021/100392
## 2023-09-16 21:17:17 Annotating text fragment 20031/100392
```

```
## 2023-09-16 21:17:18 Annotating text fragment 20041/100392
## 2023-09-16 21:17:18 Annotating text fragment 20051/100392
## 2023-09-16 21:17:18 Annotating text fragment 20061/100392
## 2023-09-16 21:17:18 Annotating text fragment 20071/100392
## 2023-09-16 21:17:18 Annotating text fragment 20081/100392
## 2023-09-16 21:17:18 Annotating text fragment 20091/100392
## 2023-09-16 21:17:18 Annotating text fragment 20101/100392
## 2023-09-16 21:17:18 Annotating text fragment 20111/100392
## 2023-09-16 21:17:18 Annotating text fragment 20121/100392
## 2023-09-16 21:17:18 Annotating text fragment 20131/100392
## 2023-09-16 21:17:18 Annotating text fragment 20141/100392
## 2023-09-16 21:17:19 Annotating text fragment 20151/100392
## 2023-09-16 21:17:19 Annotating text fragment 20161/100392
## 2023-09-16 21:17:19 Annotating text fragment 20171/100392
## 2023-09-16 21:17:19 Annotating text fragment 20181/100392
## 2023-09-16 21:17:19 Annotating text fragment 20191/100392
## 2023-09-16 21:17:19 Annotating text fragment 20201/100392
## 2023-09-16 21:17:19 Annotating text fragment 20211/100392
## 2023-09-16 21:17:19 Annotating text fragment 20221/100392
## 2023-09-16 21:17:19 Annotating text fragment 20231/100392
## 2023-09-16 21:17:19 Annotating text fragment 20241/100392
## 2023-09-16 21:17:19 Annotating text fragment 20251/100392
## 2023-09-16 21:17:20 Annotating text fragment 20261/100392
## 2023-09-16 21:17:20 Annotating text fragment 20271/100392
## 2023-09-16 21:17:20 Annotating text fragment 20281/100392
## 2023-09-16 21:17:20 Annotating text fragment 20291/100392
## 2023-09-16 21:17:20 Annotating text fragment 20301/100392
## 2023-09-16 21:17:20 Annotating text fragment 20311/100392
## 2023-09-16 21:17:20 Annotating text fragment 20321/100392
## 2023-09-16 21:17:20 Annotating text fragment 20331/100392
## 2023-09-16 21:17:20 Annotating text fragment 20341/100392
## 2023-09-16 21:17:20 Annotating text fragment 20351/100392
## 2023-09-16 21:17:20 Annotating text fragment 20361/100392
## 2023-09-16 21:17:20 Annotating text fragment 20371/100392
## 2023-09-16 21:17:21 Annotating text fragment 20381/100392
## 2023-09-16 21:17:21 Annotating text fragment 20391/100392
## 2023-09-16 21:17:21 Annotating text fragment 20401/100392
## 2023-09-16 21:17:21 Annotating text fragment 20411/100392
## 2023-09-16 21:17:21 Annotating text fragment 20421/100392
## 2023-09-16 21:17:21 Annotating text fragment 20431/100392
## 2023-09-16 21:17:21 Annotating text fragment 20441/100392
## 2023-09-16 21:17:21 Annotating text fragment 20451/100392
## 2023-09-16 21:17:21 Annotating text fragment 20461/100392
## 2023-09-16 21:17:21 Annotating text fragment 20471/100392
## 2023-09-16 21:17:21 Annotating text fragment 20481/100392
## 2023-09-16 21:17:22 Annotating text fragment 20491/100392
## 2023-09-16 21:17:22 Annotating text fragment 20501/100392
## 2023-09-16 21:17:22 Annotating text fragment 20511/100392
## 2023-09-16 21:17:22 Annotating text fragment 20521/100392
## 2023-09-16 21:17:22 Annotating text fragment 20531/100392
## 2023-09-16 21:17:22 Annotating text fragment 20541/100392
## 2023-09-16 21:17:22 Annotating text fragment 20551/100392
## 2023-09-16 21:17:22 Annotating text fragment 20561/100392
## 2023-09-16 21:17:22 Annotating text fragment 20571/100392
```

```
## 2023-09-16 21:17:22 Annotating text fragment 20581/100392
## 2023-09-16 21:17:23 Annotating text fragment 20591/100392
## 2023-09-16 21:17:23 Annotating text fragment 20601/100392
## 2023-09-16 21:17:23 Annotating text fragment 20611/100392
## 2023-09-16 21:17:23 Annotating text fragment 20621/100392
## 2023-09-16 21:17:23 Annotating text fragment 20631/100392
## 2023-09-16 21:17:23 Annotating text fragment 20641/100392
## 2023-09-16 21:17:23 Annotating text fragment 20651/100392
## 2023-09-16 21:17:23 Annotating text fragment 20661/100392
## 2023-09-16 21:17:23 Annotating text fragment 20671/100392
## 2023-09-16 21:17:23 Annotating text fragment 20681/100392
## 2023-09-16 21:17:24 Annotating text fragment 20691/100392
## 2023-09-16 21:17:24 Annotating text fragment 20701/100392
## 2023-09-16 21:17:24 Annotating text fragment 20711/100392
## 2023-09-16 21:17:24 Annotating text fragment 20721/100392
## 2023-09-16 21:17:24 Annotating text fragment 20731/100392
## 2023-09-16 21:17:24 Annotating text fragment 20741/100392
## 2023-09-16 21:17:24 Annotating text fragment 20751/100392
## 2023-09-16 21:17:24 Annotating text fragment 20761/100392
## 2023-09-16 21:17:24 Annotating text fragment 20771/100392
## 2023-09-16 21:17:24 Annotating text fragment 20781/100392
## 2023-09-16 21:17:25 Annotating text fragment 20791/100392
## 2023-09-16 21:17:25 Annotating text fragment 20801/100392
## 2023-09-16 21:17:25 Annotating text fragment 20811/100392
## 2023-09-16 21:17:25 Annotating text fragment 20821/100392
## 2023-09-16 21:17:25 Annotating text fragment 20831/100392
## 2023-09-16 21:17:25 Annotating text fragment 20841/100392
## 2023-09-16 21:17:25 Annotating text fragment 20851/100392
## 2023-09-16 21:17:25 Annotating text fragment 20861/100392
## 2023-09-16 21:17:25 Annotating text fragment 20871/100392
## 2023-09-16 21:17:25 Annotating text fragment 20881/100392
## 2023-09-16 21:17:25 Annotating text fragment 20891/100392
## 2023-09-16 21:17:25 Annotating text fragment 20901/100392
## 2023-09-16 21:17:26 Annotating text fragment 20911/100392
## 2023-09-16 21:17:26 Annotating text fragment 20921/100392
## 2023-09-16 21:17:26 Annotating text fragment 20931/100392
## 2023-09-16 21:17:26 Annotating text fragment 20941/100392
## 2023-09-16 21:17:26 Annotating text fragment 20951/100392
## 2023-09-16 21:17:26 Annotating text fragment 20961/100392
## 2023-09-16 21:17:26 Annotating text fragment 20971/100392
## 2023-09-16 21:17:26 Annotating text fragment 20981/100392
## 2023-09-16 21:17:26 Annotating text fragment 20991/100392
## 2023-09-16 21:17:26 Annotating text fragment 21001/100392
## 2023-09-16 21:17:26 Annotating text fragment 21011/100392
## 2023-09-16 21:17:27 Annotating text fragment 21021/100392
## 2023-09-16 21:17:27 Annotating text fragment 21031/100392
## 2023-09-16 21:17:27 Annotating text fragment 21041/100392
## 2023-09-16 21:17:27 Annotating text fragment 21051/100392
## 2023-09-16 21:17:27 Annotating text fragment 21061/100392
## 2023-09-16 21:17:27 Annotating text fragment 21071/100392
## 2023-09-16 21:17:27 Annotating text fragment 21081/100392
## 2023-09-16 21:17:27 Annotating text fragment 21091/100392
## 2023-09-16 21:17:27 Annotating text fragment 21101/100392
## 2023-09-16 21:17:27 Annotating text fragment 21111/100392
```

```
## 2023-09-16 21:17:27 Annotating text fragment 21121/100392
## 2023-09-16 21:17:28 Annotating text fragment 21131/100392
## 2023-09-16 21:17:28 Annotating text fragment 21141/100392
## 2023-09-16 21:17:28 Annotating text fragment 21151/100392
## 2023-09-16 21:17:28 Annotating text fragment 21161/100392
## 2023-09-16 21:17:28 Annotating text fragment 21171/100392
## 2023-09-16 21:17:28 Annotating text fragment 21181/100392
## 2023-09-16 21:17:28 Annotating text fragment 21191/100392
## 2023-09-16 21:17:28 Annotating text fragment 21201/100392
## 2023-09-16 21:17:28 Annotating text fragment 21211/100392
## 2023-09-16 21:17:28 Annotating text fragment 21221/100392
## 2023-09-16 21:17:28 Annotating text fragment 21231/100392
## 2023-09-16 21:17:28 Annotating text fragment 21241/100392
## 2023-09-16 21:17:29 Annotating text fragment 21251/100392
## 2023-09-16 21:17:29 Annotating text fragment 21261/100392
## 2023-09-16 21:17:29 Annotating text fragment 21271/100392
## 2023-09-16 21:17:29 Annotating text fragment 21281/100392
## 2023-09-16 21:17:29 Annotating text fragment 21291/100392
## 2023-09-16 21:17:29 Annotating text fragment 21301/100392
## 2023-09-16 21:17:29 Annotating text fragment 21311/100392
## 2023-09-16 21:17:29 Annotating text fragment 21321/100392
## 2023-09-16 21:17:29 Annotating text fragment 21331/100392
## 2023-09-16 21:17:30 Annotating text fragment 21341/100392
## 2023-09-16 21:17:30 Annotating text fragment 21351/100392
## 2023-09-16 21:17:30 Annotating text fragment 21361/100392
## 2023-09-16 21:17:30 Annotating text fragment 21371/100392
## 2023-09-16 21:17:30 Annotating text fragment 21381/100392
## 2023-09-16 21:17:30 Annotating text fragment 21391/100392
## 2023-09-16 21:17:30 Annotating text fragment 21401/100392
## 2023-09-16 21:17:30 Annotating text fragment 21411/100392
## 2023-09-16 21:17:30 Annotating text fragment 21421/100392
## 2023-09-16 21:17:30 Annotating text fragment 21431/100392
## 2023-09-16 21:17:30 Annotating text fragment 21441/100392
## 2023-09-16 21:17:30 Annotating text fragment 21451/100392
## 2023-09-16 21:17:31 Annotating text fragment 21461/100392
## 2023-09-16 21:17:31 Annotating text fragment 21471/100392
## 2023-09-16 21:17:31 Annotating text fragment 21481/100392
## 2023-09-16 21:17:31 Annotating text fragment 21491/100392
## 2023-09-16 21:17:31 Annotating text fragment 21501/100392
## 2023-09-16 21:17:31 Annotating text fragment 21511/100392
## 2023-09-16 21:17:31 Annotating text fragment 21521/100392
## 2023-09-16 21:17:31 Annotating text fragment 21531/100392
## 2023-09-16 21:17:31 Annotating text fragment 21541/100392
## 2023-09-16 21:17:31 Annotating text fragment 21551/100392
## 2023-09-16 21:17:32 Annotating text fragment 21561/100392
## 2023-09-16 21:17:32 Annotating text fragment 21571/100392
## 2023-09-16 21:17:32 Annotating text fragment 21581/100392
## 2023-09-16 21:17:32 Annotating text fragment 21591/100392
## 2023-09-16 21:17:32 Annotating text fragment 21601/100392
## 2023-09-16 21:17:32 Annotating text fragment 21611/100392
## 2023-09-16 21:17:32 Annotating text fragment 21621/100392
## 2023-09-16 21:17:32 Annotating text fragment 21631/100392
## 2023-09-16 21:17:32 Annotating text fragment 21641/100392
## 2023-09-16 21:17:32 Annotating text fragment 21651/100392
```

```
## 2023-09-16 21:17:32 Annotating text fragment 21661/100392
## 2023-09-16 21:17:33 Annotating text fragment 21671/100392
## 2023-09-16 21:17:33 Annotating text fragment 21681/100392
## 2023-09-16 21:17:33 Annotating text fragment 21691/100392
## 2023-09-16 21:17:33 Annotating text fragment 21701/100392
## 2023-09-16 21:17:33 Annotating text fragment 21711/100392
## 2023-09-16 21:17:33 Annotating text fragment 21721/100392
## 2023-09-16 21:17:33 Annotating text fragment 21731/100392
## 2023-09-16 21:17:33 Annotating text fragment 21741/100392
## 2023-09-16 21:17:33 Annotating text fragment 21751/100392
## 2023-09-16 21:17:33 Annotating text fragment 21761/100392
## 2023-09-16 21:17:34 Annotating text fragment 21771/100392
## 2023-09-16 21:17:34 Annotating text fragment 21781/100392
## 2023-09-16 21:17:34 Annotating text fragment 21791/100392
## 2023-09-16 21:17:34 Annotating text fragment 21801/100392
## 2023-09-16 21:17:34 Annotating text fragment 21811/100392
## 2023-09-16 21:17:34 Annotating text fragment 21821/100392
## 2023-09-16 21:17:34 Annotating text fragment 21831/100392
## 2023-09-16 21:17:34 Annotating text fragment 21841/100392
## 2023-09-16 21:17:34 Annotating text fragment 21851/100392
## 2023-09-16 21:17:34 Annotating text fragment 21861/100392
## 2023-09-16 21:17:35 Annotating text fragment 21871/100392
## 2023-09-16 21:17:35 Annotating text fragment 21881/100392
## 2023-09-16 21:17:35 Annotating text fragment 21891/100392
## 2023-09-16 21:17:35 Annotating text fragment 21901/100392
## 2023-09-16 21:17:35 Annotating text fragment 21911/100392
## 2023-09-16 21:17:35 Annotating text fragment 21921/100392
## 2023-09-16 21:17:35 Annotating text fragment 21931/100392
## 2023-09-16 21:17:35 Annotating text fragment 21941/100392
## 2023-09-16 21:17:35 Annotating text fragment 21951/100392
## 2023-09-16 21:17:35 Annotating text fragment 21961/100392
## 2023-09-16 21:17:35 Annotating text fragment 21971/100392
## 2023-09-16 21:17:36 Annotating text fragment 21981/100392
## 2023-09-16 21:17:36 Annotating text fragment 21991/100392
## 2023-09-16 21:17:36 Annotating text fragment 22001/100392
## 2023-09-16 21:17:36 Annotating text fragment 22011/100392
## 2023-09-16 21:17:36 Annotating text fragment 22021/100392
## 2023-09-16 21:17:36 Annotating text fragment 22031/100392
## 2023-09-16 21:17:36 Annotating text fragment 22041/100392
## 2023-09-16 21:17:36 Annotating text fragment 22051/100392
## 2023-09-16 21:17:36 Annotating text fragment 22061/100392
## 2023-09-16 21:17:36 Annotating text fragment 22071/100392
## 2023-09-16 21:17:37 Annotating text fragment 22081/100392
## 2023-09-16 21:17:37 Annotating text fragment 22091/100392
## 2023-09-16 21:17:37 Annotating text fragment 22101/100392
## 2023-09-16 21:17:37 Annotating text fragment 22111/100392
## 2023-09-16 21:17:37 Annotating text fragment 22121/100392
## 2023-09-16 21:17:37 Annotating text fragment 22131/100392
## 2023-09-16 21:17:37 Annotating text fragment 22141/100392
## 2023-09-16 21:17:37 Annotating text fragment 22151/100392
## 2023-09-16 21:17:38 Annotating text fragment 22161/100392
## 2023-09-16 21:17:38 Annotating text fragment 22171/100392
## 2023-09-16 21:17:38 Annotating text fragment 22181/100392
## 2023-09-16 21:17:38 Annotating text fragment 22191/100392
```

```
## 2023-09-16 21:17:38 Annotating text fragment 22201/100392
## 2023-09-16 21:17:38 Annotating text fragment 22211/100392
## 2023-09-16 21:17:38 Annotating text fragment 22221/100392
## 2023-09-16 21:17:38 Annotating text fragment 22231/100392
## 2023-09-16 21:17:38 Annotating text fragment 22241/100392
## 2023-09-16 21:17:38 Annotating text fragment 22251/100392
## 2023-09-16 21:17:38 Annotating text fragment 22261/100392
## 2023-09-16 21:17:39 Annotating text fragment 22271/100392
## 2023-09-16 21:17:39 Annotating text fragment 22281/100392
## 2023-09-16 21:17:39 Annotating text fragment 22291/100392
## 2023-09-16 21:17:39 Annotating text fragment 22301/100392
## 2023-09-16 21:17:39 Annotating text fragment 22311/100392
## 2023-09-16 21:17:39 Annotating text fragment 22321/100392
## 2023-09-16 21:17:39 Annotating text fragment 22331/100392
## 2023-09-16 21:17:39 Annotating text fragment 22341/100392
## 2023-09-16 21:17:39 Annotating text fragment 22351/100392
## 2023-09-16 21:17:39 Annotating text fragment 22361/100392
## 2023-09-16 21:17:39 Annotating text fragment 22371/100392
## 2023-09-16 21:17:40 Annotating text fragment 22381/100392
## 2023-09-16 21:17:40 Annotating text fragment 22391/100392
## 2023-09-16 21:17:40 Annotating text fragment 22401/100392
## 2023-09-16 21:17:40 Annotating text fragment 22411/100392
## 2023-09-16 21:17:40 Annotating text fragment 22421/100392
## 2023-09-16 21:17:40 Annotating text fragment 22431/100392
## 2023-09-16 21:17:40 Annotating text fragment 22441/100392
## 2023-09-16 21:17:40 Annotating text fragment 22451/100392
## 2023-09-16 21:17:40 Annotating text fragment 22461/100392
## 2023-09-16 21:17:40 Annotating text fragment 22471/100392
## 2023-09-16 21:17:40 Annotating text fragment 22481/100392
## 2023-09-16 21:17:40 Annotating text fragment 22491/100392
## 2023-09-16 21:17:41 Annotating text fragment 22501/100392
## 2023-09-16 21:17:41 Annotating text fragment 22511/100392
## 2023-09-16 21:17:41 Annotating text fragment 22521/100392
## 2023-09-16 21:17:41 Annotating text fragment 22531/100392
## 2023-09-16 21:17:41 Annotating text fragment 22541/100392
## 2023-09-16 21:17:41 Annotating text fragment 22551/100392
## 2023-09-16 21:17:41 Annotating text fragment 22561/100392
## 2023-09-16 21:17:41 Annotating text fragment 22571/100392
## 2023-09-16 21:17:41 Annotating text fragment 22581/100392
## 2023-09-16 21:17:41 Annotating text fragment 22591/100392
## 2023-09-16 21:17:41 Annotating text fragment 22601/100392
## 2023-09-16 21:17:41 Annotating text fragment 22611/100392
## 2023-09-16 21:17:42 Annotating text fragment 22621/100392
## 2023-09-16 21:17:42 Annotating text fragment 22631/100392
## 2023-09-16 21:17:42 Annotating text fragment 22641/100392
## 2023-09-16 21:17:42 Annotating text fragment 22651/100392
## 2023-09-16 21:17:42 Annotating text fragment 22661/100392
## 2023-09-16 21:17:42 Annotating text fragment 22671/100392
## 2023-09-16 21:17:42 Annotating text fragment 22681/100392
## 2023-09-16 21:17:42 Annotating text fragment 22691/100392
## 2023-09-16 21:17:42 Annotating text fragment 22701/100392
## 2023-09-16 21:17:42 Annotating text fragment 22711/100392
## 2023-09-16 21:17:43 Annotating text fragment 22721/100392
## 2023-09-16 21:17:43 Annotating text fragment 22731/100392
```

```
## 2023-09-16 21:17:43 Annotating text fragment 22741/100392
## 2023-09-16 21:17:43 Annotating text fragment 22751/100392
## 2023-09-16 21:17:43 Annotating text fragment 22761/100392
## 2023-09-16 21:17:43 Annotating text fragment 22771/100392
## 2023-09-16 21:17:43 Annotating text fragment 22781/100392
## 2023-09-16 21:17:43 Annotating text fragment 22791/100392
## 2023-09-16 21:17:43 Annotating text fragment 22801/100392
## 2023-09-16 21:17:43 Annotating text fragment 22811/100392
## 2023-09-16 21:17:43 Annotating text fragment 22821/100392
## 2023-09-16 21:17:44 Annotating text fragment 22831/100392
## 2023-09-16 21:17:44 Annotating text fragment 22841/100392
## 2023-09-16 21:17:44 Annotating text fragment 22851/100392
## 2023-09-16 21:17:44 Annotating text fragment 22861/100392
## 2023-09-16 21:17:44 Annotating text fragment 22871/100392
## 2023-09-16 21:17:44 Annotating text fragment 22881/100392
## 2023-09-16 21:17:44 Annotating text fragment 22891/100392
## 2023-09-16 21:17:44 Annotating text fragment 22901/100392
## 2023-09-16 21:17:44 Annotating text fragment 22911/100392
## 2023-09-16 21:17:44 Annotating text fragment 22921/100392
## 2023-09-16 21:17:44 Annotating text fragment 22931/100392
## 2023-09-16 21:17:44 Annotating text fragment 22941/100392
## 2023-09-16 21:17:45 Annotating text fragment 22951/100392
## 2023-09-16 21:17:45 Annotating text fragment 22961/100392
## 2023-09-16 21:17:45 Annotating text fragment 22971/100392
## 2023-09-16 21:17:45 Annotating text fragment 22981/100392
## 2023-09-16 21:17:45 Annotating text fragment 22991/100392
## 2023-09-16 21:17:45 Annotating text fragment 23001/100392
## 2023-09-16 21:17:45 Annotating text fragment 23011/100392
## 2023-09-16 21:17:45 Annotating text fragment 23021/100392
## 2023-09-16 21:17:45 Annotating text fragment 23031/100392
## 2023-09-16 21:17:45 Annotating text fragment 23041/100392
## 2023-09-16 21:17:45 Annotating text fragment 23051/100392
## 2023-09-16 21:17:46 Annotating text fragment 23061/100392
## 2023-09-16 21:17:46 Annotating text fragment 23071/100392
## 2023-09-16 21:17:46 Annotating text fragment 23081/100392
## 2023-09-16 21:17:46 Annotating text fragment 23091/100392
## 2023-09-16 21:17:46 Annotating text fragment 23101/100392
## 2023-09-16 21:17:46 Annotating text fragment 23111/100392
## 2023-09-16 21:17:46 Annotating text fragment 23121/100392
## 2023-09-16 21:17:46 Annotating text fragment 23131/100392
## 2023-09-16 21:17:46 Annotating text fragment 23141/100392
## 2023-09-16 21:17:46 Annotating text fragment 23151/100392
## 2023-09-16 21:17:46 Annotating text fragment 23161/100392
## 2023-09-16 21:17:47 Annotating text fragment 23171/100392
## 2023-09-16 21:17:47 Annotating text fragment 23181/100392
## 2023-09-16 21:17:47 Annotating text fragment 23191/100392
## 2023-09-16 21:17:47 Annotating text fragment 23201/100392
## 2023-09-16 21:17:47 Annotating text fragment 23211/100392
## 2023-09-16 21:17:47 Annotating text fragment 23221/100392
## 2023-09-16 21:17:47 Annotating text fragment 23231/100392
## 2023-09-16 21:17:47 Annotating text fragment 23241/100392
## 2023-09-16 21:17:47 Annotating text fragment 23251/100392
## 2023-09-16 21:17:47 Annotating text fragment 23261/100392
## 2023-09-16 21:17:47 Annotating text fragment 23271/100392
```

```
## 2023-09-16 21:17:47 Annotating text fragment 23281/100392
## 2023-09-16 21:17:48 Annotating text fragment 23291/100392
## 2023-09-16 21:17:48 Annotating text fragment 23301/100392
## 2023-09-16 21:17:48 Annotating text fragment 23311/100392
## 2023-09-16 21:17:48 Annotating text fragment 23321/100392
## 2023-09-16 21:17:48 Annotating text fragment 23331/100392
## 2023-09-16 21:17:48 Annotating text fragment 23341/100392
## 2023-09-16 21:17:48 Annotating text fragment 23351/100392
## 2023-09-16 21:17:48 Annotating text fragment 23361/100392
## 2023-09-16 21:17:48 Annotating text fragment 23371/100392
## 2023-09-16 21:17:48 Annotating text fragment 23381/100392
## 2023-09-16 21:17:48 Annotating text fragment 23391/100392
## 2023-09-16 21:17:48 Annotating text fragment 23401/100392
## 2023-09-16 21:17:49 Annotating text fragment 23411/100392
## 2023-09-16 21:17:49 Annotating text fragment 23421/100392
## 2023-09-16 21:17:49 Annotating text fragment 23431/100392
## 2023-09-16 21:17:49 Annotating text fragment 23441/100392
## 2023-09-16 21:17:49 Annotating text fragment 23451/100392
## 2023-09-16 21:17:49 Annotating text fragment 23461/100392
## 2023-09-16 21:17:49 Annotating text fragment 23471/100392
## 2023-09-16 21:17:49 Annotating text fragment 23481/100392
## 2023-09-16 21:17:49 Annotating text fragment 23491/100392
## 2023-09-16 21:17:49 Annotating text fragment 23501/100392
## 2023-09-16 21:17:49 Annotating text fragment 23511/100392
## 2023-09-16 21:17:49 Annotating text fragment 23521/100392
## 2023-09-16 21:17:50 Annotating text fragment 23531/100392
## 2023-09-16 21:17:50 Annotating text fragment 23541/100392
## 2023-09-16 21:17:50 Annotating text fragment 23551/100392
## 2023-09-16 21:17:50 Annotating text fragment 23561/100392
## 2023-09-16 21:17:50 Annotating text fragment 23571/100392
## 2023-09-16 21:17:50 Annotating text fragment 23581/100392
## 2023-09-16 21:17:50 Annotating text fragment 23591/100392
## 2023-09-16 21:17:50 Annotating text fragment 23601/100392
## 2023-09-16 21:17:50 Annotating text fragment 23611/100392
## 2023-09-16 21:17:50 Annotating text fragment 23621/100392
## 2023-09-16 21:17:50 Annotating text fragment 23631/100392
## 2023-09-16 21:17:50 Annotating text fragment 23641/100392
## 2023-09-16 21:17:51 Annotating text fragment 23651/100392
## 2023-09-16 21:17:51 Annotating text fragment 23661/100392
## 2023-09-16 21:17:51 Annotating text fragment 23671/100392
## 2023-09-16 21:17:51 Annotating text fragment 23681/100392
## 2023-09-16 21:17:51 Annotating text fragment 23691/100392
## 2023-09-16 21:17:51 Annotating text fragment 23701/100392
## 2023-09-16 21:17:51 Annotating text fragment 23711/100392
## 2023-09-16 21:17:51 Annotating text fragment 23721/100392
## 2023-09-16 21:17:51 Annotating text fragment 23731/100392
## 2023-09-16 21:17:51 Annotating text fragment 23741/100392
## 2023-09-16 21:17:52 Annotating text fragment 23751/100392
## 2023-09-16 21:17:52 Annotating text fragment 23761/100392
## 2023-09-16 21:17:52 Annotating text fragment 23771/100392
## 2023-09-16 21:17:52 Annotating text fragment 23781/100392
## 2023-09-16 21:17:52 Annotating text fragment 23791/100392
## 2023-09-16 21:17:52 Annotating text fragment 23801/100392
## 2023-09-16 21:17:52 Annotating text fragment 23811/100392
```

```
## 2023-09-16 21:17:52 Annotating text fragment 23821/100392
## 2023-09-16 21:17:52 Annotating text fragment 23831/100392
## 2023-09-16 21:17:52 Annotating text fragment 23841/100392
## 2023-09-16 21:17:53 Annotating text fragment 23851/100392
## 2023-09-16 21:17:53 Annotating text fragment 23861/100392
## 2023-09-16 21:17:53 Annotating text fragment 23871/100392
## 2023-09-16 21:17:53 Annotating text fragment 23881/100392
## 2023-09-16 21:17:53 Annotating text fragment 23891/100392
## 2023-09-16 21:17:53 Annotating text fragment 23901/100392
## 2023-09-16 21:17:53 Annotating text fragment 23911/100392
## 2023-09-16 21:17:53 Annotating text fragment 23921/100392
## 2023-09-16 21:17:53 Annotating text fragment 23931/100392
## 2023-09-16 21:17:53 Annotating text fragment 23941/100392
## 2023-09-16 21:17:53 Annotating text fragment 23951/100392
## 2023-09-16 21:17:53 Annotating text fragment 23961/100392
## 2023-09-16 21:17:54 Annotating text fragment 23971/100392
## 2023-09-16 21:17:54 Annotating text fragment 23981/100392
## 2023-09-16 21:17:54 Annotating text fragment 23991/100392
## 2023-09-16 21:17:54 Annotating text fragment 24001/100392
## 2023-09-16 21:17:54 Annotating text fragment 24011/100392
## 2023-09-16 21:17:54 Annotating text fragment 24021/100392
## 2023-09-16 21:17:54 Annotating text fragment 24031/100392
## 2023-09-16 21:17:54 Annotating text fragment 24041/100392
## 2023-09-16 21:17:54 Annotating text fragment 24051/100392
## 2023-09-16 21:17:54 Annotating text fragment 24061/100392
## 2023-09-16 21:17:54 Annotating text fragment 24071/100392
## 2023-09-16 21:17:55 Annotating text fragment 24081/100392
## 2023-09-16 21:17:55 Annotating text fragment 24091/100392
## 2023-09-16 21:17:55 Annotating text fragment 24101/100392
## 2023-09-16 21:17:55 Annotating text fragment 24111/100392
## 2023-09-16 21:17:55 Annotating text fragment 24121/100392
## 2023-09-16 21:17:55 Annotating text fragment 24131/100392
## 2023-09-16 21:17:55 Annotating text fragment 24141/100392
## 2023-09-16 21:17:55 Annotating text fragment 24151/100392
## 2023-09-16 21:17:55 Annotating text fragment 24161/100392
## 2023-09-16 21:17:55 Annotating text fragment 24171/100392
## 2023-09-16 21:17:55 Annotating text fragment 24181/100392
## 2023-09-16 21:17:56 Annotating text fragment 24191/100392
## 2023-09-16 21:17:56 Annotating text fragment 24201/100392
## 2023-09-16 21:17:56 Annotating text fragment 24211/100392
## 2023-09-16 21:17:56 Annotating text fragment 24221/100392
## 2023-09-16 21:17:56 Annotating text fragment 24231/100392
## 2023-09-16 21:17:56 Annotating text fragment 24241/100392
## 2023-09-16 21:17:56 Annotating text fragment 24251/100392
## 2023-09-16 21:17:56 Annotating text fragment 24261/100392
## 2023-09-16 21:17:56 Annotating text fragment 24271/100392
## 2023-09-16 21:17:56 Annotating text fragment 24281/100392
## 2023-09-16 21:17:57 Annotating text fragment 24291/100392
## 2023-09-16 21:17:57 Annotating text fragment 24301/100392
## 2023-09-16 21:17:57 Annotating text fragment 24311/100392
## 2023-09-16 21:17:57 Annotating text fragment 24321/100392
## 2023-09-16 21:17:57 Annotating text fragment 24331/100392
## 2023-09-16 21:17:57 Annotating text fragment 24341/100392
## 2023-09-16 21:17:57 Annotating text fragment 24351/100392
```

```
## 2023-09-16 21:17:57 Annotating text fragment 24361/100392
## 2023-09-16 21:17:57 Annotating text fragment 24371/100392
## 2023-09-16 21:17:57 Annotating text fragment 24381/100392
## 2023-09-16 21:17:57 Annotating text fragment 24391/100392
## 2023-09-16 21:17:58 Annotating text fragment 24401/100392
## 2023-09-16 21:17:58 Annotating text fragment 24411/100392
## 2023-09-16 21:17:58 Annotating text fragment 24421/100392
## 2023-09-16 21:17:58 Annotating text fragment 24431/100392
## 2023-09-16 21:17:58 Annotating text fragment 24441/100392
## 2023-09-16 21:17:58 Annotating text fragment 24451/100392
## 2023-09-16 21:17:58 Annotating text fragment 24461/100392
## 2023-09-16 21:17:58 Annotating text fragment 24471/100392
## 2023-09-16 21:17:58 Annotating text fragment 24481/100392
## 2023-09-16 21:17:59 Annotating text fragment 24491/100392
## 2023-09-16 21:17:59 Annotating text fragment 24501/100392
## 2023-09-16 21:17:59 Annotating text fragment 24511/100392
## 2023-09-16 21:17:59 Annotating text fragment 24521/100392
## 2023-09-16 21:17:59 Annotating text fragment 24531/100392
## 2023-09-16 21:17:59 Annotating text fragment 24541/100392
## 2023-09-16 21:17:59 Annotating text fragment 24551/100392
## 2023-09-16 21:17:59 Annotating text fragment 24561/100392
## 2023-09-16 21:17:59 Annotating text fragment 24571/100392
## 2023-09-16 21:17:59 Annotating text fragment 24581/100392
## 2023-09-16 21:17:59 Annotating text fragment 24591/100392
## 2023-09-16 21:17:59 Annotating text fragment 24601/100392
## 2023-09-16 21:18:00 Annotating text fragment 24611/100392
## 2023-09-16 21:18:00 Annotating text fragment 24621/100392
## 2023-09-16 21:18:00 Annotating text fragment 24631/100392
## 2023-09-16 21:18:00 Annotating text fragment 24641/100392
## 2023-09-16 21:18:00 Annotating text fragment 24651/100392
## 2023-09-16 21:18:00 Annotating text fragment 24661/100392
## 2023-09-16 21:18:00 Annotating text fragment 24671/100392
## 2023-09-16 21:18:00 Annotating text fragment 24681/100392
## 2023-09-16 21:18:01 Annotating text fragment 24691/100392
## 2023-09-16 21:18:01 Annotating text fragment 24701/100392
## 2023-09-16 21:18:01 Annotating text fragment 24711/100392
## 2023-09-16 21:18:01 Annotating text fragment 24721/100392
## 2023-09-16 21:18:01 Annotating text fragment 24731/100392
## 2023-09-16 21:18:01 Annotating text fragment 24741/100392
## 2023-09-16 21:18:01 Annotating text fragment 24751/100392
## 2023-09-16 21:18:01 Annotating text fragment 24761/100392
## 2023-09-16 21:18:02 Annotating text fragment 24771/100392
## 2023-09-16 21:18:02 Annotating text fragment 24781/100392
## 2023-09-16 21:18:02 Annotating text fragment 24791/100392
## 2023-09-16 21:18:02 Annotating text fragment 24801/100392
## 2023-09-16 21:18:02 Annotating text fragment 24811/100392
## 2023-09-16 21:18:02 Annotating text fragment 24821/100392
## 2023-09-16 21:18:02 Annotating text fragment 24831/100392
## 2023-09-16 21:18:02 Annotating text fragment 24841/100392
## 2023-09-16 21:18:02 Annotating text fragment 24851/100392
## 2023-09-16 21:18:02 Annotating text fragment 24861/100392
## 2023-09-16 21:18:02 Annotating text fragment 24871/100392
## 2023-09-16 21:18:02 Annotating text fragment 24881/100392
## 2023-09-16 21:18:03 Annotating text fragment 24891/100392
```

```
## 2023-09-16 21:18:03 Annotating text fragment 24901/100392
## 2023-09-16 21:18:03 Annotating text fragment 24911/100392
## 2023-09-16 21:18:03 Annotating text fragment 24921/100392
## 2023-09-16 21:18:03 Annotating text fragment 24931/100392
## 2023-09-16 21:18:03 Annotating text fragment 24941/100392
## 2023-09-16 21:18:03 Annotating text fragment 24951/100392
## 2023-09-16 21:18:03 Annotating text fragment 24961/100392
## 2023-09-16 21:18:03 Annotating text fragment 24971/100392
## 2023-09-16 21:18:03 Annotating text fragment 24981/100392
## 2023-09-16 21:18:03 Annotating text fragment 24991/100392
## 2023-09-16 21:18:03 Annotating text fragment 25001/100392
## 2023-09-16 21:18:04 Annotating text fragment 25011/100392
## 2023-09-16 21:18:04 Annotating text fragment 25021/100392
## 2023-09-16 21:18:04 Annotating text fragment 25031/100392
## 2023-09-16 21:18:04 Annotating text fragment 25041/100392
## 2023-09-16 21:18:04 Annotating text fragment 25051/100392
## 2023-09-16 21:18:04 Annotating text fragment 25061/100392
## 2023-09-16 21:18:04 Annotating text fragment 25071/100392
## 2023-09-16 21:18:04 Annotating text fragment 25081/100392
## 2023-09-16 21:18:05 Annotating text fragment 25091/100392
## 2023-09-16 21:18:05 Annotating text fragment 25101/100392
## 2023-09-16 21:18:05 Annotating text fragment 25111/100392
## 2023-09-16 21:18:05 Annotating text fragment 25121/100392
## 2023-09-16 21:18:05 Annotating text fragment 25131/100392
## 2023-09-16 21:18:05 Annotating text fragment 25141/100392
## 2023-09-16 21:18:05 Annotating text fragment 25151/100392
## 2023-09-16 21:18:05 Annotating text fragment 25161/100392
## 2023-09-16 21:18:05 Annotating text fragment 25171/100392
## 2023-09-16 21:18:05 Annotating text fragment 25181/100392
## 2023-09-16 21:18:05 Annotating text fragment 25191/100392
## 2023-09-16 21:18:06 Annotating text fragment 25201/100392
## 2023-09-16 21:18:06 Annotating text fragment 25211/100392
## 2023-09-16 21:18:06 Annotating text fragment 25221/100392
## 2023-09-16 21:18:06 Annotating text fragment 25231/100392
## 2023-09-16 21:18:06 Annotating text fragment 25241/100392
## 2023-09-16 21:18:06 Annotating text fragment 25251/100392
## 2023-09-16 21:18:06 Annotating text fragment 25261/100392
## 2023-09-16 21:18:06 Annotating text fragment 25271/100392
## 2023-09-16 21:18:06 Annotating text fragment 25281/100392
## 2023-09-16 21:18:06 Annotating text fragment 25291/100392
## 2023-09-16 21:18:06 Annotating text fragment 25301/100392
## 2023-09-16 21:18:07 Annotating text fragment 25311/100392
## 2023-09-16 21:18:07 Annotating text fragment 25321/100392
## 2023-09-16 21:18:07 Annotating text fragment 25331/100392
## 2023-09-16 21:18:07 Annotating text fragment 25341/100392
## 2023-09-16 21:18:07 Annotating text fragment 25351/100392
## 2023-09-16 21:18:07 Annotating text fragment 25361/100392
## 2023-09-16 21:18:07 Annotating text fragment 25371/100392
## 2023-09-16 21:18:07 Annotating text fragment 25381/100392
## 2023-09-16 21:18:07 Annotating text fragment 25391/100392
## 2023-09-16 21:18:07 Annotating text fragment 25401/100392
## 2023-09-16 21:18:08 Annotating text fragment 25411/100392
## 2023-09-16 21:18:08 Annotating text fragment 25421/100392
## 2023-09-16 21:18:08 Annotating text fragment 25431/100392
```

```
## 2023-09-16 21:18:08 Annotating text fragment 25441/100392
## 2023-09-16 21:18:08 Annotating text fragment 25451/100392
## 2023-09-16 21:18:08 Annotating text fragment 25461/100392
## 2023-09-16 21:18:08 Annotating text fragment 25471/100392
## 2023-09-16 21:18:08 Annotating text fragment 25481/100392
## 2023-09-16 21:18:08 Annotating text fragment 25491/100392
## 2023-09-16 21:18:08 Annotating text fragment 25501/100392
## 2023-09-16 21:18:08 Annotating text fragment 25511/100392
## 2023-09-16 21:18:08 Annotating text fragment 25521/100392
## 2023-09-16 21:18:09 Annotating text fragment 25531/100392
## 2023-09-16 21:18:09 Annotating text fragment 25541/100392
## 2023-09-16 21:18:09 Annotating text fragment 25551/100392
## 2023-09-16 21:18:09 Annotating text fragment 25561/100392
## 2023-09-16 21:18:09 Annotating text fragment 25571/100392
## 2023-09-16 21:18:09 Annotating text fragment 25581/100392
## 2023-09-16 21:18:09 Annotating text fragment 25591/100392
## 2023-09-16 21:18:09 Annotating text fragment 25601/100392
## 2023-09-16 21:18:09 Annotating text fragment 25611/100392
## 2023-09-16 21:18:09 Annotating text fragment 25621/100392
## 2023-09-16 21:18:10 Annotating text fragment 25631/100392
## 2023-09-16 21:18:10 Annotating text fragment 25641/100392
## 2023-09-16 21:18:10 Annotating text fragment 25651/100392
## 2023-09-16 21:18:10 Annotating text fragment 25661/100392
## 2023-09-16 21:18:10 Annotating text fragment 25671/100392
## 2023-09-16 21:18:10 Annotating text fragment 25681/100392
## 2023-09-16 21:18:10 Annotating text fragment 25691/100392
## 2023-09-16 21:18:10 Annotating text fragment 25701/100392
## 2023-09-16 21:18:10 Annotating text fragment 25711/100392
## 2023-09-16 21:18:11 Annotating text fragment 25721/100392
## 2023-09-16 21:18:11 Annotating text fragment 25731/100392
## 2023-09-16 21:18:11 Annotating text fragment 25741/100392
## 2023-09-16 21:18:11 Annotating text fragment 25751/100392
## 2023-09-16 21:18:11 Annotating text fragment 25761/100392
## 2023-09-16 21:18:11 Annotating text fragment 25771/100392
## 2023-09-16 21:18:11 Annotating text fragment 25781/100392
## 2023-09-16 21:18:11 Annotating text fragment 25791/100392
## 2023-09-16 21:18:12 Annotating text fragment 25801/100392
## 2023-09-16 21:18:12 Annotating text fragment 25811/100392
## 2023-09-16 21:18:12 Annotating text fragment 25821/100392
## 2023-09-16 21:18:12 Annotating text fragment 25831/100392
## 2023-09-16 21:18:12 Annotating text fragment 25841/100392
## 2023-09-16 21:18:12 Annotating text fragment 25851/100392
## 2023-09-16 21:18:12 Annotating text fragment 25861/100392
## 2023-09-16 21:18:12 Annotating text fragment 25871/100392
## 2023-09-16 21:18:12 Annotating text fragment 25881/100392
## 2023-09-16 21:18:12 Annotating text fragment 25891/100392
## 2023-09-16 21:18:12 Annotating text fragment 25901/100392
## 2023-09-16 21:18:12 Annotating text fragment 25911/100392
## 2023-09-16 21:18:13 Annotating text fragment 25921/100392
## 2023-09-16 21:18:13 Annotating text fragment 25931/100392
## 2023-09-16 21:18:13 Annotating text fragment 25941/100392
## 2023-09-16 21:18:13 Annotating text fragment 25951/100392
## 2023-09-16 21:18:13 Annotating text fragment 25961/100392
## 2023-09-16 21:18:13 Annotating text fragment 25971/100392
```

```
## 2023-09-16 21:18:13 Annotating text fragment 25981/100392
## 2023-09-16 21:18:13 Annotating text fragment 25991/100392
## 2023-09-16 21:18:13 Annotating text fragment 26001/100392
## 2023-09-16 21:18:13 Annotating text fragment 26011/100392
## 2023-09-16 21:18:13 Annotating text fragment 26021/100392
## 2023-09-16 21:18:13 Annotating text fragment 26031/100392
## 2023-09-16 21:18:14 Annotating text fragment 26041/100392
## 2023-09-16 21:18:14 Annotating text fragment 26051/100392
## 2023-09-16 21:18:14 Annotating text fragment 26061/100392
## 2023-09-16 21:18:14 Annotating text fragment 26071/100392
## 2023-09-16 21:18:14 Annotating text fragment 26081/100392
## 2023-09-16 21:18:14 Annotating text fragment 26091/100392
## 2023-09-16 21:18:14 Annotating text fragment 26101/100392
## 2023-09-16 21:18:14 Annotating text fragment 26111/100392
## 2023-09-16 21:18:14 Annotating text fragment 26121/100392
## 2023-09-16 21:18:14 Annotating text fragment 26131/100392
## 2023-09-16 21:18:14 Annotating text fragment 26141/100392
## 2023-09-16 21:18:15 Annotating text fragment 26151/100392
## 2023-09-16 21:18:15 Annotating text fragment 26161/100392
## 2023-09-16 21:18:15 Annotating text fragment 26171/100392
## 2023-09-16 21:18:15 Annotating text fragment 26181/100392
## 2023-09-16 21:18:15 Annotating text fragment 26191/100392
## 2023-09-16 21:18:15 Annotating text fragment 26201/100392
## 2023-09-16 21:18:15 Annotating text fragment 26211/100392
## 2023-09-16 21:18:15 Annotating text fragment 26221/100392
## 2023-09-16 21:18:15 Annotating text fragment 26231/100392
## 2023-09-16 21:18:15 Annotating text fragment 26241/100392
## 2023-09-16 21:18:15 Annotating text fragment 26251/100392
## 2023-09-16 21:18:16 Annotating text fragment 26261/100392
## 2023-09-16 21:18:16 Annotating text fragment 26271/100392
## 2023-09-16 21:18:16 Annotating text fragment 26281/100392
## 2023-09-16 21:18:16 Annotating text fragment 26291/100392
## 2023-09-16 21:18:16 Annotating text fragment 26301/100392
## 2023-09-16 21:18:16 Annotating text fragment 26311/100392
## 2023-09-16 21:18:16 Annotating text fragment 26321/100392
## 2023-09-16 21:18:16 Annotating text fragment 26331/100392
## 2023-09-16 21:18:16 Annotating text fragment 26341/100392
## 2023-09-16 21:18:16 Annotating text fragment 26351/100392
## 2023-09-16 21:18:16 Annotating text fragment 26361/100392
## 2023-09-16 21:18:17 Annotating text fragment 26371/100392
## 2023-09-16 21:18:17 Annotating text fragment 26381/100392
## 2023-09-16 21:18:17 Annotating text fragment 26391/100392
## 2023-09-16 21:18:17 Annotating text fragment 26401/100392
## 2023-09-16 21:18:17 Annotating text fragment 26411/100392
## 2023-09-16 21:18:17 Annotating text fragment 26421/100392
## 2023-09-16 21:18:17 Annotating text fragment 26431/100392
## 2023-09-16 21:18:17 Annotating text fragment 26441/100392
## 2023-09-16 21:18:17 Annotating text fragment 26451/100392
## 2023-09-16 21:18:17 Annotating text fragment 26461/100392
## 2023-09-16 21:18:17 Annotating text fragment 26471/100392
## 2023-09-16 21:18:18 Annotating text fragment 26481/100392
## 2023-09-16 21:18:18 Annotating text fragment 26491/100392
## 2023-09-16 21:18:18 Annotating text fragment 26501/100392
## 2023-09-16 21:18:18 Annotating text fragment 26511/100392
```

```
## 2023-09-16 21:18:18 Annotating text fragment 26521/100392
## 2023-09-16 21:18:18 Annotating text fragment 26531/100392
## 2023-09-16 21:18:18 Annotating text fragment 26541/100392
## 2023-09-16 21:18:18 Annotating text fragment 26551/100392
## 2023-09-16 21:18:18 Annotating text fragment 26561/100392
## 2023-09-16 21:18:18 Annotating text fragment 26571/100392
## 2023-09-16 21:18:18 Annotating text fragment 26581/100392
## 2023-09-16 21:18:18 Annotating text fragment 26591/100392
## 2023-09-16 21:18:19 Annotating text fragment 26601/100392
## 2023-09-16 21:18:19 Annotating text fragment 26611/100392
## 2023-09-16 21:18:19 Annotating text fragment 26621/100392
## 2023-09-16 21:18:19 Annotating text fragment 26631/100392
## 2023-09-16 21:18:19 Annotating text fragment 26641/100392
## 2023-09-16 21:18:19 Annotating text fragment 26651/100392
## 2023-09-16 21:18:19 Annotating text fragment 26661/100392
## 2023-09-16 21:18:19 Annotating text fragment 26671/100392
## 2023-09-16 21:18:19 Annotating text fragment 26681/100392
## 2023-09-16 21:18:19 Annotating text fragment 26691/100392
## 2023-09-16 21:18:19 Annotating text fragment 26701/100392
## 2023-09-16 21:18:20 Annotating text fragment 26711/100392
## 2023-09-16 21:18:20 Annotating text fragment 26721/100392
## 2023-09-16 21:18:20 Annotating text fragment 26731/100392
## 2023-09-16 21:18:20 Annotating text fragment 26741/100392
## 2023-09-16 21:18:20 Annotating text fragment 26751/100392
## 2023-09-16 21:18:20 Annotating text fragment 26761/100392
## 2023-09-16 21:18:20 Annotating text fragment 26771/100392
## 2023-09-16 21:18:20 Annotating text fragment 26781/100392
## 2023-09-16 21:18:21 Annotating text fragment 26791/100392
## 2023-09-16 21:18:21 Annotating text fragment 26801/100392
## 2023-09-16 21:18:21 Annotating text fragment 26811/100392
## 2023-09-16 21:18:21 Annotating text fragment 26821/100392
## 2023-09-16 21:18:21 Annotating text fragment 26831/100392
## 2023-09-16 21:18:21 Annotating text fragment 26841/100392
## 2023-09-16 21:18:21 Annotating text fragment 26851/100392
## 2023-09-16 21:18:21 Annotating text fragment 26861/100392
## 2023-09-16 21:18:21 Annotating text fragment 26871/100392
## 2023-09-16 21:18:21 Annotating text fragment 26881/100392
## 2023-09-16 21:18:21 Annotating text fragment 26891/100392
## 2023-09-16 21:18:21 Annotating text fragment 26901/100392
## 2023-09-16 21:18:22 Annotating text fragment 26911/100392
## 2023-09-16 21:18:22 Annotating text fragment 26921/100392
## 2023-09-16 21:18:22 Annotating text fragment 26931/100392
## 2023-09-16 21:18:22 Annotating text fragment 26941/100392
## 2023-09-16 21:18:22 Annotating text fragment 26951/100392
## 2023-09-16 21:18:22 Annotating text fragment 26961/100392
## 2023-09-16 21:18:22 Annotating text fragment 26971/100392
## 2023-09-16 21:18:22 Annotating text fragment 26981/100392
## 2023-09-16 21:18:22 Annotating text fragment 26991/100392
## 2023-09-16 21:18:22 Annotating text fragment 27001/100392
## 2023-09-16 21:18:22 Annotating text fragment 27011/100392
## 2023-09-16 21:18:23 Annotating text fragment 27021/100392
## 2023-09-16 21:18:23 Annotating text fragment 27031/100392
## 2023-09-16 21:18:23 Annotating text fragment 27041/100392
## 2023-09-16 21:18:23 Annotating text fragment 27051/100392
```

```
## 2023-09-16 21:18:23 Annotating text fragment 27061/100392
## 2023-09-16 21:18:23 Annotating text fragment 27071/100392
## 2023-09-16 21:18:23 Annotating text fragment 27081/100392
## 2023-09-16 21:18:23 Annotating text fragment 27091/100392
## 2023-09-16 21:18:23 Annotating text fragment 27101/100392
## 2023-09-16 21:18:23 Annotating text fragment 27111/100392
## 2023-09-16 21:18:24 Annotating text fragment 27121/100392
## 2023-09-16 21:18:24 Annotating text fragment 27131/100392
## 2023-09-16 21:18:24 Annotating text fragment 27141/100392
## 2023-09-16 21:18:24 Annotating text fragment 27151/100392
## 2023-09-16 21:18:24 Annotating text fragment 27161/100392
## 2023-09-16 21:18:24 Annotating text fragment 27171/100392
## 2023-09-16 21:18:24 Annotating text fragment 27181/100392
## 2023-09-16 21:18:24 Annotating text fragment 27191/100392
## 2023-09-16 21:18:25 Annotating text fragment 27201/100392
## 2023-09-16 21:18:25 Annotating text fragment 27211/100392
## 2023-09-16 21:18:25 Annotating text fragment 27221/100392
## 2023-09-16 21:18:25 Annotating text fragment 27231/100392
## 2023-09-16 21:18:25 Annotating text fragment 27241/100392
## 2023-09-16 21:18:25 Annotating text fragment 27251/100392
## 2023-09-16 21:18:25 Annotating text fragment 27261/100392
## 2023-09-16 21:18:25 Annotating text fragment 27271/100392
## 2023-09-16 21:18:25 Annotating text fragment 27281/100392
## 2023-09-16 21:18:25 Annotating text fragment 27291/100392
## 2023-09-16 21:18:26 Annotating text fragment 27301/100392
## 2023-09-16 21:18:26 Annotating text fragment 27311/100392
## 2023-09-16 21:18:26 Annotating text fragment 27321/100392
## 2023-09-16 21:18:26 Annotating text fragment 27331/100392
## 2023-09-16 21:18:26 Annotating text fragment 27341/100392
## 2023-09-16 21:18:26 Annotating text fragment 27351/100392
## 2023-09-16 21:18:26 Annotating text fragment 27361/100392
## 2023-09-16 21:18:26 Annotating text fragment 27371/100392
## 2023-09-16 21:18:26 Annotating text fragment 27381/100392
## 2023-09-16 21:18:26 Annotating text fragment 27391/100392
## 2023-09-16 21:18:26 Annotating text fragment 27401/100392
## 2023-09-16 21:18:27 Annotating text fragment 27411/100392
## 2023-09-16 21:18:27 Annotating text fragment 27421/100392
## 2023-09-16 21:18:27 Annotating text fragment 27431/100392
## 2023-09-16 21:18:27 Annotating text fragment 27441/100392
## 2023-09-16 21:18:27 Annotating text fragment 27451/100392
## 2023-09-16 21:18:27 Annotating text fragment 27461/100392
## 2023-09-16 21:18:27 Annotating text fragment 27471/100392
## 2023-09-16 21:18:27 Annotating text fragment 27481/100392
## 2023-09-16 21:18:28 Annotating text fragment 27491/100392
## 2023-09-16 21:18:28 Annotating text fragment 27501/100392
## 2023-09-16 21:18:28 Annotating text fragment 27511/100392
## 2023-09-16 21:18:28 Annotating text fragment 27521/100392
## 2023-09-16 21:18:28 Annotating text fragment 27531/100392
## 2023-09-16 21:18:28 Annotating text fragment 27541/100392
## 2023-09-16 21:18:28 Annotating text fragment 27551/100392
## 2023-09-16 21:18:28 Annotating text fragment 27561/100392
## 2023-09-16 21:18:28 Annotating text fragment 27571/100392
## 2023-09-16 21:18:28 Annotating text fragment 27581/100392
## 2023-09-16 21:18:28 Annotating text fragment 27591/100392
```

```
## 2023-09-16 21:18:29 Annotating text fragment 27601/100392
## 2023-09-16 21:18:29 Annotating text fragment 27611/100392
## 2023-09-16 21:18:29 Annotating text fragment 27621/100392
## 2023-09-16 21:18:29 Annotating text fragment 27631/100392
## 2023-09-16 21:18:29 Annotating text fragment 27641/100392
## 2023-09-16 21:18:29 Annotating text fragment 27651/100392
## 2023-09-16 21:18:29 Annotating text fragment 27661/100392
## 2023-09-16 21:18:30 Annotating text fragment 27671/100392
## 2023-09-16 21:18:30 Annotating text fragment 27681/100392
## 2023-09-16 21:18:30 Annotating text fragment 27691/100392
## 2023-09-16 21:18:31 Annotating text fragment 27701/100392
## 2023-09-16 21:18:31 Annotating text fragment 27711/100392
## 2023-09-16 21:18:31 Annotating text fragment 27721/100392
## 2023-09-16 21:18:31 Annotating text fragment 27731/100392
## 2023-09-16 21:18:31 Annotating text fragment 27741/100392
## 2023-09-16 21:18:31 Annotating text fragment 27751/100392
## 2023-09-16 21:18:31 Annotating text fragment 27761/100392
## 2023-09-16 21:18:31 Annotating text fragment 27771/100392
## 2023-09-16 21:18:31 Annotating text fragment 27781/100392
## 2023-09-16 21:18:31 Annotating text fragment 27791/100392
## 2023-09-16 21:18:32 Annotating text fragment 27801/100392
## 2023-09-16 21:18:32 Annotating text fragment 27811/100392
## 2023-09-16 21:18:32 Annotating text fragment 27821/100392
## 2023-09-16 21:18:32 Annotating text fragment 27831/100392
## 2023-09-16 21:18:32 Annotating text fragment 27841/100392
## 2023-09-16 21:18:32 Annotating text fragment 27851/100392
## 2023-09-16 21:18:32 Annotating text fragment 27861/100392
## 2023-09-16 21:18:32 Annotating text fragment 27871/100392
## 2023-09-16 21:18:32 Annotating text fragment 27881/100392
## 2023-09-16 21:18:32 Annotating text fragment 27891/100392
## 2023-09-16 21:18:32 Annotating text fragment 27901/100392
## 2023-09-16 21:18:33 Annotating text fragment 27911/100392
## 2023-09-16 21:18:33 Annotating text fragment 27921/100392
## 2023-09-16 21:18:33 Annotating text fragment 27931/100392
## 2023-09-16 21:18:33 Annotating text fragment 27941/100392
## 2023-09-16 21:18:33 Annotating text fragment 27951/100392
## 2023-09-16 21:18:33 Annotating text fragment 27961/100392
## 2023-09-16 21:18:33 Annotating text fragment 27971/100392
## 2023-09-16 21:18:33 Annotating text fragment 27981/100392
## 2023-09-16 21:18:34 Annotating text fragment 27991/100392
## 2023-09-16 21:18:34 Annotating text fragment 28001/100392
## 2023-09-16 21:18:34 Annotating text fragment 28011/100392
## 2023-09-16 21:18:34 Annotating text fragment 28021/100392
## 2023-09-16 21:18:34 Annotating text fragment 28031/100392
## 2023-09-16 21:18:34 Annotating text fragment 28041/100392
## 2023-09-16 21:18:34 Annotating text fragment 28051/100392
## 2023-09-16 21:18:34 Annotating text fragment 28061/100392
## 2023-09-16 21:18:34 Annotating text fragment 28071/100392
## 2023-09-16 21:18:34 Annotating text fragment 28081/100392
## 2023-09-16 21:18:35 Annotating text fragment 28091/100392
## 2023-09-16 21:18:35 Annotating text fragment 28101/100392
## 2023-09-16 21:18:35 Annotating text fragment 28111/100392
## 2023-09-16 21:18:35 Annotating text fragment 28121/100392
## 2023-09-16 21:18:35 Annotating text fragment 28131/100392
```

```
## 2023-09-16 21:18:35 Annotating text fragment 28141/100392
## 2023-09-16 21:18:35 Annotating text fragment 28151/100392
## 2023-09-16 21:18:35 Annotating text fragment 28161/100392
## 2023-09-16 21:18:35 Annotating text fragment 28171/100392
## 2023-09-16 21:18:35 Annotating text fragment 28181/100392
## 2023-09-16 21:18:36 Annotating text fragment 28191/100392
## 2023-09-16 21:18:36 Annotating text fragment 28201/100392
## 2023-09-16 21:18:36 Annotating text fragment 28211/100392
## 2023-09-16 21:18:36 Annotating text fragment 28221/100392
## 2023-09-16 21:18:36 Annotating text fragment 28231/100392
## 2023-09-16 21:18:36 Annotating text fragment 28241/100392
## 2023-09-16 21:18:36 Annotating text fragment 28251/100392
## 2023-09-16 21:18:36 Annotating text fragment 28261/100392
## 2023-09-16 21:18:36 Annotating text fragment 28271/100392
## 2023-09-16 21:18:36 Annotating text fragment 28281/100392
## 2023-09-16 21:18:36 Annotating text fragment 28291/100392
## 2023-09-16 21:18:36 Annotating text fragment 28301/100392
## 2023-09-16 21:18:37 Annotating text fragment 28311/100392
## 2023-09-16 21:18:37 Annotating text fragment 28321/100392
## 2023-09-16 21:18:37 Annotating text fragment 28331/100392
## 2023-09-16 21:18:37 Annotating text fragment 28341/100392
## 2023-09-16 21:18:37 Annotating text fragment 28351/100392
## 2023-09-16 21:18:37 Annotating text fragment 28361/100392
## 2023-09-16 21:18:37 Annotating text fragment 28371/100392
## 2023-09-16 21:18:37 Annotating text fragment 28381/100392
## 2023-09-16 21:18:37 Annotating text fragment 28391/100392
## 2023-09-16 21:18:37 Annotating text fragment 28401/100392
## 2023-09-16 21:18:37 Annotating text fragment 28411/100392
## 2023-09-16 21:18:37 Annotating text fragment 28421/100392
## 2023-09-16 21:18:38 Annotating text fragment 28431/100392
## 2023-09-16 21:18:38 Annotating text fragment 28441/100392
## 2023-09-16 21:18:38 Annotating text fragment 28451/100392
## 2023-09-16 21:18:38 Annotating text fragment 28461/100392
## 2023-09-16 21:18:38 Annotating text fragment 28471/100392
## 2023-09-16 21:18:38 Annotating text fragment 28481/100392
## 2023-09-16 21:18:38 Annotating text fragment 28491/100392
## 2023-09-16 21:18:38 Annotating text fragment 28501/100392
## 2023-09-16 21:18:39 Annotating text fragment 28511/100392
## 2023-09-16 21:18:39 Annotating text fragment 28521/100392
## 2023-09-16 21:18:39 Annotating text fragment 28531/100392
## 2023-09-16 21:18:39 Annotating text fragment 28541/100392
## 2023-09-16 21:18:39 Annotating text fragment 28551/100392
## 2023-09-16 21:18:39 Annotating text fragment 28561/100392
## 2023-09-16 21:18:39 Annotating text fragment 28571/100392
## 2023-09-16 21:18:39 Annotating text fragment 28581/100392
## 2023-09-16 21:18:39 Annotating text fragment 28591/100392
## 2023-09-16 21:18:39 Annotating text fragment 28601/100392
## 2023-09-16 21:18:39 Annotating text fragment 28611/100392
## 2023-09-16 21:18:40 Annotating text fragment 28621/100392
## 2023-09-16 21:18:40 Annotating text fragment 28631/100392
## 2023-09-16 21:18:40 Annotating text fragment 28641/100392
## 2023-09-16 21:18:40 Annotating text fragment 28651/100392
## 2023-09-16 21:18:40 Annotating text fragment 28661/100392
## 2023-09-16 21:18:40 Annotating text fragment 28671/100392
```

```
## 2023-09-16 21:18:40 Annotating text fragment 28681/100392
## 2023-09-16 21:18:40 Annotating text fragment 28691/100392
## 2023-09-16 21:18:40 Annotating text fragment 28701/100392
## 2023-09-16 21:18:40 Annotating text fragment 28711/100392
## 2023-09-16 21:18:40 Annotating text fragment 28721/100392
## 2023-09-16 21:18:40 Annotating text fragment 28731/100392
## 2023-09-16 21:18:41 Annotating text fragment 28741/100392
## 2023-09-16 21:18:41 Annotating text fragment 28751/100392
## 2023-09-16 21:18:41 Annotating text fragment 28761/100392
## 2023-09-16 21:18:41 Annotating text fragment 28771/100392
## 2023-09-16 21:18:41 Annotating text fragment 28781/100392
## 2023-09-16 21:18:41 Annotating text fragment 28791/100392
## 2023-09-16 21:18:41 Annotating text fragment 28801/100392
## 2023-09-16 21:18:41 Annotating text fragment 28811/100392
## 2023-09-16 21:18:41 Annotating text fragment 28821/100392
## 2023-09-16 21:18:41 Annotating text fragment 28831/100392
## 2023-09-16 21:18:41 Annotating text fragment 28841/100392
## 2023-09-16 21:18:42 Annotating text fragment 28851/100392
## 2023-09-16 21:18:42 Annotating text fragment 28861/100392
## 2023-09-16 21:18:42 Annotating text fragment 28871/100392
## 2023-09-16 21:18:42 Annotating text fragment 28881/100392
## 2023-09-16 21:18:42 Annotating text fragment 28891/100392
## 2023-09-16 21:18:42 Annotating text fragment 28901/100392
## 2023-09-16 21:18:42 Annotating text fragment 28911/100392
## 2023-09-16 21:18:42 Annotating text fragment 28921/100392
## 2023-09-16 21:18:42 Annotating text fragment 28931/100392
## 2023-09-16 21:18:43 Annotating text fragment 28941/100392
## 2023-09-16 21:18:43 Annotating text fragment 28951/100392
## 2023-09-16 21:18:43 Annotating text fragment 28961/100392
## 2023-09-16 21:18:43 Annotating text fragment 28971/100392
## 2023-09-16 21:18:43 Annotating text fragment 28981/100392
## 2023-09-16 21:18:43 Annotating text fragment 28991/100392
## 2023-09-16 21:18:43 Annotating text fragment 29001/100392
## 2023-09-16 21:18:43 Annotating text fragment 29011/100392
## 2023-09-16 21:18:43 Annotating text fragment 29021/100392
## 2023-09-16 21:18:43 Annotating text fragment 29031/100392
## 2023-09-16 21:18:43 Annotating text fragment 29041/100392
## 2023-09-16 21:18:43 Annotating text fragment 29051/100392
## 2023-09-16 21:18:44 Annotating text fragment 29061/100392
## 2023-09-16 21:18:44 Annotating text fragment 29071/100392
## 2023-09-16 21:18:44 Annotating text fragment 29081/100392
## 2023-09-16 21:18:44 Annotating text fragment 29091/100392
## 2023-09-16 21:18:44 Annotating text fragment 29101/100392
## 2023-09-16 21:18:44 Annotating text fragment 29111/100392
## 2023-09-16 21:18:44 Annotating text fragment 29121/100392
## 2023-09-16 21:18:44 Annotating text fragment 29131/100392
## 2023-09-16 21:18:44 Annotating text fragment 29141/100392
## 2023-09-16 21:18:45 Annotating text fragment 29151/100392
## 2023-09-16 21:18:45 Annotating text fragment 29161/100392
## 2023-09-16 21:18:45 Annotating text fragment 29171/100392
## 2023-09-16 21:18:45 Annotating text fragment 29181/100392
## 2023-09-16 21:18:45 Annotating text fragment 29191/100392
## 2023-09-16 21:18:45 Annotating text fragment 29201/100392
## 2023-09-16 21:18:45 Annotating text fragment 29211/100392
```

```
## 2023-09-16 21:18:45 Annotating text fragment 29221/100392
## 2023-09-16 21:18:45 Annotating text fragment 29231/100392
## 2023-09-16 21:18:45 Annotating text fragment 29241/100392
## 2023-09-16 21:18:45 Annotating text fragment 29251/100392
## 2023-09-16 21:18:46 Annotating text fragment 29261/100392
## 2023-09-16 21:18:46 Annotating text fragment 29271/100392
## 2023-09-16 21:18:46 Annotating text fragment 29281/100392
## 2023-09-16 21:18:46 Annotating text fragment 29291/100392
## 2023-09-16 21:18:46 Annotating text fragment 29301/100392
## 2023-09-16 21:18:46 Annotating text fragment 29311/100392
## 2023-09-16 21:18:46 Annotating text fragment 29321/100392
## 2023-09-16 21:18:46 Annotating text fragment 29331/100392
## 2023-09-16 21:18:46 Annotating text fragment 29341/100392
## 2023-09-16 21:18:46 Annotating text fragment 29351/100392
## 2023-09-16 21:18:47 Annotating text fragment 29361/100392
## 2023-09-16 21:18:47 Annotating text fragment 29371/100392
## 2023-09-16 21:18:47 Annotating text fragment 29381/100392
## 2023-09-16 21:18:47 Annotating text fragment 29391/100392
## 2023-09-16 21:18:47 Annotating text fragment 29401/100392
## 2023-09-16 21:18:47 Annotating text fragment 29411/100392
## 2023-09-16 21:18:47 Annotating text fragment 29421/100392
## 2023-09-16 21:18:47 Annotating text fragment 29431/100392
## 2023-09-16 21:18:47 Annotating text fragment 29441/100392
## 2023-09-16 21:18:47 Annotating text fragment 29451/100392
## 2023-09-16 21:18:48 Annotating text fragment 29461/100392
## 2023-09-16 21:18:48 Annotating text fragment 29471/100392
## 2023-09-16 21:18:48 Annotating text fragment 29481/100392
## 2023-09-16 21:18:48 Annotating text fragment 29491/100392
## 2023-09-16 21:18:48 Annotating text fragment 29501/100392
## 2023-09-16 21:18:48 Annotating text fragment 29511/100392
## 2023-09-16 21:18:48 Annotating text fragment 29521/100392
## 2023-09-16 21:18:48 Annotating text fragment 29531/100392
## 2023-09-16 21:18:48 Annotating text fragment 29541/100392
## 2023-09-16 21:18:48 Annotating text fragment 29551/100392
## 2023-09-16 21:18:48 Annotating text fragment 29561/100392
## 2023-09-16 21:18:49 Annotating text fragment 29571/100392
## 2023-09-16 21:18:49 Annotating text fragment 29581/100392
## 2023-09-16 21:18:49 Annotating text fragment 29591/100392
## 2023-09-16 21:18:49 Annotating text fragment 29601/100392
## 2023-09-16 21:18:49 Annotating text fragment 29611/100392
## 2023-09-16 21:18:49 Annotating text fragment 29621/100392
## 2023-09-16 21:18:49 Annotating text fragment 29631/100392
## 2023-09-16 21:18:49 Annotating text fragment 29641/100392
## 2023-09-16 21:18:49 Annotating text fragment 29651/100392
## 2023-09-16 21:18:49 Annotating text fragment 29661/100392
## 2023-09-16 21:18:49 Annotating text fragment 29671/100392
## 2023-09-16 21:18:50 Annotating text fragment 29681/100392
## 2023-09-16 21:18:50 Annotating text fragment 29691/100392
## 2023-09-16 21:18:50 Annotating text fragment 29701/100392
## 2023-09-16 21:18:50 Annotating text fragment 29711/100392
## 2023-09-16 21:18:50 Annotating text fragment 29721/100392
## 2023-09-16 21:18:50 Annotating text fragment 29731/100392
## 2023-09-16 21:18:50 Annotating text fragment 29741/100392
## 2023-09-16 21:18:50 Annotating text fragment 29751/100392
```

```
## 2023-09-16 21:18:50 Annotating text fragment 29761/100392
## 2023-09-16 21:18:50 Annotating text fragment 29771/100392
## 2023-09-16 21:18:50 Annotating text fragment 29781/100392
## 2023-09-16 21:18:50 Annotating text fragment 29791/100392
## 2023-09-16 21:18:51 Annotating text fragment 29801/100392
## 2023-09-16 21:18:51 Annotating text fragment 29811/100392
## 2023-09-16 21:18:51 Annotating text fragment 29821/100392
## 2023-09-16 21:18:51 Annotating text fragment 29831/100392
## 2023-09-16 21:18:51 Annotating text fragment 29841/100392
## 2023-09-16 21:18:51 Annotating text fragment 29851/100392
## 2023-09-16 21:18:51 Annotating text fragment 29861/100392
## 2023-09-16 21:18:51 Annotating text fragment 29871/100392
## 2023-09-16 21:18:51 Annotating text fragment 29881/100392
## 2023-09-16 21:18:51 Annotating text fragment 29891/100392
## 2023-09-16 21:18:52 Annotating text fragment 29901/100392
## 2023-09-16 21:18:52 Annotating text fragment 29911/100392
## 2023-09-16 21:18:52 Annotating text fragment 29921/100392
## 2023-09-16 21:18:52 Annotating text fragment 29931/100392
## 2023-09-16 21:18:52 Annotating text fragment 29941/100392
## 2023-09-16 21:18:52 Annotating text fragment 29951/100392
## 2023-09-16 21:18:52 Annotating text fragment 29961/100392
## 2023-09-16 21:18:52 Annotating text fragment 29971/100392
## 2023-09-16 21:18:52 Annotating text fragment 29981/100392
## 2023-09-16 21:18:52 Annotating text fragment 29991/100392
## 2023-09-16 21:18:52 Annotating text fragment 30001/100392
## 2023-09-16 21:18:53 Annotating text fragment 30011/100392
## 2023-09-16 21:18:53 Annotating text fragment 30021/100392
## 2023-09-16 21:18:53 Annotating text fragment 30031/100392
## 2023-09-16 21:18:53 Annotating text fragment 30041/100392
## 2023-09-16 21:18:53 Annotating text fragment 30051/100392
## 2023-09-16 21:18:53 Annotating text fragment 30061/100392
## 2023-09-16 21:18:53 Annotating text fragment 30071/100392
## 2023-09-16 21:18:53 Annotating text fragment 30081/100392
## 2023-09-16 21:18:53 Annotating text fragment 30091/100392
## 2023-09-16 21:18:53 Annotating text fragment 30101/100392
## 2023-09-16 21:18:54 Annotating text fragment 30111/100392
## 2023-09-16 21:18:54 Annotating text fragment 30121/100392
## 2023-09-16 21:18:54 Annotating text fragment 30131/100392
## 2023-09-16 21:18:54 Annotating text fragment 30141/100392
## 2023-09-16 21:18:54 Annotating text fragment 30151/100392
## 2023-09-16 21:18:54 Annotating text fragment 30161/100392
## 2023-09-16 21:18:54 Annotating text fragment 30171/100392
## 2023-09-16 21:18:54 Annotating text fragment 30181/100392
## 2023-09-16 21:18:55 Annotating text fragment 30191/100392
## 2023-09-16 21:18:55 Annotating text fragment 30201/100392
## 2023-09-16 21:18:55 Annotating text fragment 30211/100392
## 2023-09-16 21:18:55 Annotating text fragment 30221/100392
## 2023-09-16 21:18:55 Annotating text fragment 30231/100392
## 2023-09-16 21:18:55 Annotating text fragment 30241/100392
## 2023-09-16 21:18:55 Annotating text fragment 30251/100392
## 2023-09-16 21:18:55 Annotating text fragment 30261/100392
## 2023-09-16 21:18:55 Annotating text fragment 30271/100392
## 2023-09-16 21:18:56 Annotating text fragment 30281/100392
## 2023-09-16 21:18:56 Annotating text fragment 30291/100392
```

```
## 2023-09-16 21:18:56 Annotating text fragment 30301/100392
## 2023-09-16 21:18:56 Annotating text fragment 30311/100392
## 2023-09-16 21:18:56 Annotating text fragment 30321/100392
## 2023-09-16 21:18:56 Annotating text fragment 30331/100392
## 2023-09-16 21:18:56 Annotating text fragment 30341/100392
## 2023-09-16 21:18:56 Annotating text fragment 30351/100392
## 2023-09-16 21:18:57 Annotating text fragment 30361/100392
## 2023-09-16 21:18:57 Annotating text fragment 30371/100392
## 2023-09-16 21:18:57 Annotating text fragment 30381/100392
## 2023-09-16 21:18:57 Annotating text fragment 30391/100392
## 2023-09-16 21:18:57 Annotating text fragment 30401/100392
## 2023-09-16 21:18:57 Annotating text fragment 30411/100392
## 2023-09-16 21:18:57 Annotating text fragment 30421/100392
## 2023-09-16 21:18:57 Annotating text fragment 30431/100392
## 2023-09-16 21:18:57 Annotating text fragment 30441/100392
## 2023-09-16 21:18:57 Annotating text fragment 30451/100392
## 2023-09-16 21:18:58 Annotating text fragment 30461/100392
## 2023-09-16 21:18:58 Annotating text fragment 30471/100392
## 2023-09-16 21:18:58 Annotating text fragment 30481/100392
## 2023-09-16 21:18:58 Annotating text fragment 30491/100392
## 2023-09-16 21:18:58 Annotating text fragment 30501/100392
## 2023-09-16 21:18:58 Annotating text fragment 30511/100392
## 2023-09-16 21:18:58 Annotating text fragment 30521/100392
## 2023-09-16 21:18:59 Annotating text fragment 30531/100392
## 2023-09-16 21:18:59 Annotating text fragment 30541/100392
## 2023-09-16 21:18:59 Annotating text fragment 30551/100392
## 2023-09-16 21:18:59 Annotating text fragment 30561/100392
## 2023-09-16 21:18:59 Annotating text fragment 30571/100392
## 2023-09-16 21:18:59 Annotating text fragment 30581/100392
## 2023-09-16 21:18:59 Annotating text fragment 30591/100392
## 2023-09-16 21:18:59 Annotating text fragment 30601/100392
## 2023-09-16 21:18:59 Annotating text fragment 30611/100392
## 2023-09-16 21:19:00 Annotating text fragment 30621/100392
## 2023-09-16 21:19:00 Annotating text fragment 30631/100392
## 2023-09-16 21:19:00 Annotating text fragment 30641/100392
## 2023-09-16 21:19:00 Annotating text fragment 30651/100392
## 2023-09-16 21:19:00 Annotating text fragment 30661/100392
## 2023-09-16 21:19:00 Annotating text fragment 30671/100392
## 2023-09-16 21:19:00 Annotating text fragment 30681/100392
## 2023-09-16 21:19:00 Annotating text fragment 30691/100392
## 2023-09-16 21:19:00 Annotating text fragment 30701/100392
## 2023-09-16 21:19:01 Annotating text fragment 30711/100392
## 2023-09-16 21:19:01 Annotating text fragment 30721/100392
## 2023-09-16 21:19:01 Annotating text fragment 30731/100392
## 2023-09-16 21:19:01 Annotating text fragment 30741/100392
## 2023-09-16 21:19:01 Annotating text fragment 30751/100392
## 2023-09-16 21:19:01 Annotating text fragment 30761/100392
## 2023-09-16 21:19:01 Annotating text fragment 30771/100392
## 2023-09-16 21:19:01 Annotating text fragment 30781/100392
## 2023-09-16 21:19:01 Annotating text fragment 30791/100392
## 2023-09-16 21:19:01 Annotating text fragment 30801/100392
## 2023-09-16 21:19:02 Annotating text fragment 30811/100392
## 2023-09-16 21:19:02 Annotating text fragment 30821/100392
## 2023-09-16 21:19:02 Annotating text fragment 30831/100392
```

```
## 2023-09-16 21:19:02 Annotating text fragment 30841/100392
## 2023-09-16 21:19:02 Annotating text fragment 30851/100392
## 2023-09-16 21:19:02 Annotating text fragment 30861/100392
## 2023-09-16 21:19:02 Annotating text fragment 30871/100392
## 2023-09-16 21:19:02 Annotating text fragment 30881/100392
## 2023-09-16 21:19:03 Annotating text fragment 30891/100392
## 2023-09-16 21:19:03 Annotating text fragment 30901/100392
## 2023-09-16 21:19:03 Annotating text fragment 30911/100392
## 2023-09-16 21:19:03 Annotating text fragment 30921/100392
## 2023-09-16 21:19:03 Annotating text fragment 30931/100392
## 2023-09-16 21:19:03 Annotating text fragment 30941/100392
## 2023-09-16 21:19:03 Annotating text fragment 30951/100392
## 2023-09-16 21:19:04 Annotating text fragment 30961/100392
## 2023-09-16 21:19:04 Annotating text fragment 30971/100392
## 2023-09-16 21:19:04 Annotating text fragment 30981/100392
## 2023-09-16 21:19:04 Annotating text fragment 30991/100392
## 2023-09-16 21:19:04 Annotating text fragment 31001/100392
## 2023-09-16 21:19:04 Annotating text fragment 31011/100392
## 2023-09-16 21:19:04 Annotating text fragment 31021/100392
## 2023-09-16 21:19:05 Annotating text fragment 31031/100392
## 2023-09-16 21:19:05 Annotating text fragment 31041/100392
## 2023-09-16 21:19:05 Annotating text fragment 31051/100392
## 2023-09-16 21:19:05 Annotating text fragment 31061/100392
## 2023-09-16 21:19:05 Annotating text fragment 31071/100392
## 2023-09-16 21:19:05 Annotating text fragment 31081/100392
## 2023-09-16 21:19:05 Annotating text fragment 31091/100392
## 2023-09-16 21:19:05 Annotating text fragment 31101/100392
## 2023-09-16 21:19:06 Annotating text fragment 31111/100392
## 2023-09-16 21:19:06 Annotating text fragment 31121/100392
## 2023-09-16 21:19:06 Annotating text fragment 31131/100392
## 2023-09-16 21:19:06 Annotating text fragment 31141/100392
## 2023-09-16 21:19:06 Annotating text fragment 31151/100392
## 2023-09-16 21:19:06 Annotating text fragment 31161/100392
## 2023-09-16 21:19:06 Annotating text fragment 31171/100392
## 2023-09-16 21:19:06 Annotating text fragment 31181/100392
## 2023-09-16 21:19:07 Annotating text fragment 31191/100392
## 2023-09-16 21:19:07 Annotating text fragment 31201/100392
## 2023-09-16 21:19:07 Annotating text fragment 31211/100392
## 2023-09-16 21:19:07 Annotating text fragment 31221/100392
## 2023-09-16 21:19:07 Annotating text fragment 31231/100392
## 2023-09-16 21:19:07 Annotating text fragment 31241/100392
## 2023-09-16 21:19:08 Annotating text fragment 31251/100392
## 2023-09-16 21:19:08 Annotating text fragment 31261/100392
## 2023-09-16 21:19:08 Annotating text fragment 31271/100392
## 2023-09-16 21:19:08 Annotating text fragment 31281/100392
## 2023-09-16 21:19:08 Annotating text fragment 31291/100392
## 2023-09-16 21:19:08 Annotating text fragment 31301/100392
## 2023-09-16 21:19:08 Annotating text fragment 31311/100392
## 2023-09-16 21:19:08 Annotating text fragment 31321/100392
## 2023-09-16 21:19:09 Annotating text fragment 31331/100392
## 2023-09-16 21:19:09 Annotating text fragment 31341/100392
## 2023-09-16 21:19:09 Annotating text fragment 31351/100392
## 2023-09-16 21:19:09 Annotating text fragment 31361/100392
## 2023-09-16 21:19:09 Annotating text fragment 31371/100392
```

```
## 2023-09-16 21:19:09 Annotating text fragment 31381/100392
## 2023-09-16 21:19:09 Annotating text fragment 31391/100392
## 2023-09-16 21:19:09 Annotating text fragment 31401/100392
## 2023-09-16 21:19:09 Annotating text fragment 31411/100392
## 2023-09-16 21:19:10 Annotating text fragment 31421/100392
## 2023-09-16 21:19:10 Annotating text fragment 31431/100392
## 2023-09-16 21:19:10 Annotating text fragment 31441/100392
## 2023-09-16 21:19:10 Annotating text fragment 31451/100392
## 2023-09-16 21:19:10 Annotating text fragment 31461/100392
## 2023-09-16 21:19:10 Annotating text fragment 31471/100392
## 2023-09-16 21:19:10 Annotating text fragment 31481/100392
## 2023-09-16 21:19:11 Annotating text fragment 31491/100392
## 2023-09-16 21:19:11 Annotating text fragment 31501/100392
## 2023-09-16 21:19:11 Annotating text fragment 31511/100392
## 2023-09-16 21:19:11 Annotating text fragment 31521/100392
## 2023-09-16 21:19:11 Annotating text fragment 31531/100392
## 2023-09-16 21:19:11 Annotating text fragment 31541/100392
## 2023-09-16 21:19:11 Annotating text fragment 31551/100392
## 2023-09-16 21:19:11 Annotating text fragment 31561/100392
## 2023-09-16 21:19:12 Annotating text fragment 31571/100392
## 2023-09-16 21:19:12 Annotating text fragment 31581/100392
## 2023-09-16 21:19:12 Annotating text fragment 31591/100392
## 2023-09-16 21:19:12 Annotating text fragment 31601/100392
## 2023-09-16 21:19:12 Annotating text fragment 31611/100392
## 2023-09-16 21:19:12 Annotating text fragment 31621/100392
## 2023-09-16 21:19:12 Annotating text fragment 31631/100392
## 2023-09-16 21:19:12 Annotating text fragment 31641/100392
## 2023-09-16 21:19:13 Annotating text fragment 31651/100392
## 2023-09-16 21:19:13 Annotating text fragment 31661/100392
## 2023-09-16 21:19:13 Annotating text fragment 31671/100392
## 2023-09-16 21:19:13 Annotating text fragment 31681/100392
## 2023-09-16 21:19:13 Annotating text fragment 31691/100392
## 2023-09-16 21:19:13 Annotating text fragment 31701/100392
## 2023-09-16 21:19:14 Annotating text fragment 31711/100392
## 2023-09-16 21:19:14 Annotating text fragment 31721/100392
## 2023-09-16 21:19:14 Annotating text fragment 31731/100392
## 2023-09-16 21:19:14 Annotating text fragment 31741/100392
## 2023-09-16 21:19:14 Annotating text fragment 31751/100392
## 2023-09-16 21:19:14 Annotating text fragment 31761/100392
## 2023-09-16 21:19:14 Annotating text fragment 31771/100392
## 2023-09-16 21:19:15 Annotating text fragment 31781/100392
## 2023-09-16 21:19:15 Annotating text fragment 31791/100392
## 2023-09-16 21:19:15 Annotating text fragment 31801/100392
## 2023-09-16 21:19:15 Annotating text fragment 31811/100392
## 2023-09-16 21:19:15 Annotating text fragment 31821/100392
## 2023-09-16 21:19:15 Annotating text fragment 31831/100392
## 2023-09-16 21:19:15 Annotating text fragment 31841/100392
## 2023-09-16 21:19:15 Annotating text fragment 31851/100392
## 2023-09-16 21:19:16 Annotating text fragment 31861/100392
## 2023-09-16 21:19:16 Annotating text fragment 31871/100392
## 2023-09-16 21:19:16 Annotating text fragment 31881/100392
## 2023-09-16 21:19:16 Annotating text fragment 31891/100392
## 2023-09-16 21:19:16 Annotating text fragment 31901/100392
## 2023-09-16 21:19:16 Annotating text fragment 31911/100392
```

```
## 2023-09-16 21:19:16 Annotating text fragment 31921/100392
## 2023-09-16 21:19:16 Annotating text fragment 31931/100392
## 2023-09-16 21:19:17 Annotating text fragment 31941/100392
## 2023-09-16 21:19:17 Annotating text fragment 31951/100392
## 2023-09-16 21:19:17 Annotating text fragment 31961/100392
## 2023-09-16 21:19:17 Annotating text fragment 31971/100392
## 2023-09-16 21:19:17 Annotating text fragment 31981/100392
## 2023-09-16 21:19:17 Annotating text fragment 31991/100392
## 2023-09-16 21:19:17 Annotating text fragment 32001/100392
## 2023-09-16 21:19:17 Annotating text fragment 32011/100392
## 2023-09-16 21:19:18 Annotating text fragment 32021/100392
## 2023-09-16 21:19:18 Annotating text fragment 32031/100392
## 2023-09-16 21:19:18 Annotating text fragment 32041/100392
## 2023-09-16 21:19:18 Annotating text fragment 32051/100392
## 2023-09-16 21:19:18 Annotating text fragment 32061/100392
## 2023-09-16 21:19:18 Annotating text fragment 32071/100392
## 2023-09-16 21:19:18 Annotating text fragment 32081/100392
## 2023-09-16 21:19:18 Annotating text fragment 32091/100392
## 2023-09-16 21:19:18 Annotating text fragment 32101/100392
## 2023-09-16 21:19:18 Annotating text fragment 32111/100392
## 2023-09-16 21:19:19 Annotating text fragment 32121/100392
## 2023-09-16 21:19:19 Annotating text fragment 32131/100392
## 2023-09-16 21:19:19 Annotating text fragment 32141/100392
## 2023-09-16 21:19:19 Annotating text fragment 32151/100392
## 2023-09-16 21:19:19 Annotating text fragment 32161/100392
## 2023-09-16 21:19:19 Annotating text fragment 32171/100392
## 2023-09-16 21:19:19 Annotating text fragment 32181/100392
## 2023-09-16 21:19:20 Annotating text fragment 32191/100392
## 2023-09-16 21:19:20 Annotating text fragment 32201/100392
## 2023-09-16 21:19:20 Annotating text fragment 32211/100392
## 2023-09-16 21:19:20 Annotating text fragment 32221/100392
## 2023-09-16 21:19:20 Annotating text fragment 32231/100392
## 2023-09-16 21:19:20 Annotating text fragment 32241/100392
## 2023-09-16 21:19:20 Annotating text fragment 32251/100392
## 2023-09-16 21:19:21 Annotating text fragment 32261/100392
## 2023-09-16 21:19:21 Annotating text fragment 32271/100392
## 2023-09-16 21:19:21 Annotating text fragment 32281/100392
## 2023-09-16 21:19:21 Annotating text fragment 32291/100392
## 2023-09-16 21:19:21 Annotating text fragment 32301/100392
## 2023-09-16 21:19:21 Annotating text fragment 32311/100392
## 2023-09-16 21:19:21 Annotating text fragment 32321/100392
## 2023-09-16 21:19:21 Annotating text fragment 32331/100392
## 2023-09-16 21:19:22 Annotating text fragment 32341/100392
## 2023-09-16 21:19:22 Annotating text fragment 32351/100392
## 2023-09-16 21:19:22 Annotating text fragment 32361/100392
## 2023-09-16 21:19:22 Annotating text fragment 32371/100392
## 2023-09-16 21:19:22 Annotating text fragment 32381/100392
## 2023-09-16 21:19:22 Annotating text fragment 32391/100392
## 2023-09-16 21:19:22 Annotating text fragment 32401/100392
## 2023-09-16 21:19:23 Annotating text fragment 32411/100392
## 2023-09-16 21:19:23 Annotating text fragment 32421/100392
## 2023-09-16 21:19:23 Annotating text fragment 32431/100392
## 2023-09-16 21:19:23 Annotating text fragment 32441/100392
## 2023-09-16 21:19:23 Annotating text fragment 32451/100392
```

```
## 2023-09-16 21:19:23 Annotating text fragment 32461/100392
## 2023-09-16 21:19:23 Annotating text fragment 32471/100392
## 2023-09-16 21:19:23 Annotating text fragment 32481/100392
## 2023-09-16 21:19:24 Annotating text fragment 32491/100392
## 2023-09-16 21:19:24 Annotating text fragment 32501/100392
## 2023-09-16 21:19:24 Annotating text fragment 32511/100392
## 2023-09-16 21:19:24 Annotating text fragment 32521/100392
## 2023-09-16 21:19:24 Annotating text fragment 32531/100392
## 2023-09-16 21:19:24 Annotating text fragment 32541/100392
## 2023-09-16 21:19:24 Annotating text fragment 32551/100392
## 2023-09-16 21:19:24 Annotating text fragment 32561/100392
## 2023-09-16 21:19:25 Annotating text fragment 32571/100392
## 2023-09-16 21:19:25 Annotating text fragment 32581/100392
## 2023-09-16 21:19:25 Annotating text fragment 32591/100392
## 2023-09-16 21:19:25 Annotating text fragment 32601/100392
## 2023-09-16 21:19:26 Annotating text fragment 32611/100392
## 2023-09-16 21:19:26 Annotating text fragment 32621/100392
## 2023-09-16 21:19:26 Annotating text fragment 32631/100392
## 2023-09-16 21:19:26 Annotating text fragment 32641/100392
## 2023-09-16 21:19:26 Annotating text fragment 32651/100392
## 2023-09-16 21:19:26 Annotating text fragment 32661/100392
## 2023-09-16 21:19:26 Annotating text fragment 32671/100392
## 2023-09-16 21:19:27 Annotating text fragment 32681/100392
## 2023-09-16 21:19:27 Annotating text fragment 32691/100392
## 2023-09-16 21:19:27 Annotating text fragment 32701/100392
## 2023-09-16 21:19:27 Annotating text fragment 32711/100392
## 2023-09-16 21:19:27 Annotating text fragment 32721/100392
## 2023-09-16 21:19:27 Annotating text fragment 32731/100392
## 2023-09-16 21:19:28 Annotating text fragment 32741/100392
## 2023-09-16 21:19:28 Annotating text fragment 32751/100392
## 2023-09-16 21:19:28 Annotating text fragment 32761/100392
## 2023-09-16 21:19:28 Annotating text fragment 32771/100392
## 2023-09-16 21:19:28 Annotating text fragment 32781/100392
## 2023-09-16 21:19:29 Annotating text fragment 32791/100392
## 2023-09-16 21:19:29 Annotating text fragment 32801/100392
## 2023-09-16 21:19:29 Annotating text fragment 32811/100392
## 2023-09-16 21:19:29 Annotating text fragment 32821/100392
## 2023-09-16 21:19:29 Annotating text fragment 32831/100392
## 2023-09-16 21:19:29 Annotating text fragment 32841/100392
## 2023-09-16 21:19:29 Annotating text fragment 32851/100392
## 2023-09-16 21:19:29 Annotating text fragment 32861/100392
## 2023-09-16 21:19:30 Annotating text fragment 32871/100392
## 2023-09-16 21:19:30 Annotating text fragment 32881/100392
## 2023-09-16 21:19:30 Annotating text fragment 32891/100392
## 2023-09-16 21:19:30 Annotating text fragment 32901/100392
## 2023-09-16 21:19:30 Annotating text fragment 32911/100392
## 2023-09-16 21:19:30 Annotating text fragment 32921/100392
## 2023-09-16 21:19:31 Annotating text fragment 32931/100392
## 2023-09-16 21:19:31 Annotating text fragment 32941/100392
## 2023-09-16 21:19:31 Annotating text fragment 32951/100392
## 2023-09-16 21:19:31 Annotating text fragment 32961/100392
## 2023-09-16 21:19:31 Annotating text fragment 32971/100392
## 2023-09-16 21:19:31 Annotating text fragment 32981/100392
## 2023-09-16 21:19:31 Annotating text fragment 32991/100392
```

```
## 2023-09-16 21:19:31 Annotating text fragment 33001/100392
## 2023-09-16 21:19:32 Annotating text fragment 33011/100392
## 2023-09-16 21:19:32 Annotating text fragment 33021/100392
## 2023-09-16 21:19:32 Annotating text fragment 33031/100392
## 2023-09-16 21:19:32 Annotating text fragment 33041/100392
## 2023-09-16 21:19:32 Annotating text fragment 33051/100392
## 2023-09-16 21:19:33 Annotating text fragment 33061/100392
## 2023-09-16 21:19:33 Annotating text fragment 33071/100392
## 2023-09-16 21:19:33 Annotating text fragment 33081/100392
## 2023-09-16 21:19:33 Annotating text fragment 33091/100392
## 2023-09-16 21:19:33 Annotating text fragment 33101/100392
## 2023-09-16 21:19:33 Annotating text fragment 33111/100392
## 2023-09-16 21:19:34 Annotating text fragment 33121/100392
## 2023-09-16 21:19:34 Annotating text fragment 33131/100392
## 2023-09-16 21:19:34 Annotating text fragment 33141/100392
## 2023-09-16 21:19:34 Annotating text fragment 33151/100392
## 2023-09-16 21:19:34 Annotating text fragment 33161/100392
## 2023-09-16 21:19:34 Annotating text fragment 33171/100392
## 2023-09-16 21:19:35 Annotating text fragment 33181/100392
## 2023-09-16 21:19:35 Annotating text fragment 33191/100392
## 2023-09-16 21:19:35 Annotating text fragment 33201/100392
## 2023-09-16 21:19:35 Annotating text fragment 33211/100392
## 2023-09-16 21:19:35 Annotating text fragment 33221/100392
## 2023-09-16 21:19:36 Annotating text fragment 33231/100392
## 2023-09-16 21:19:36 Annotating text fragment 33241/100392
## 2023-09-16 21:19:36 Annotating text fragment 33251/100392
## 2023-09-16 21:19:36 Annotating text fragment 33261/100392
## 2023-09-16 21:19:36 Annotating text fragment 33271/100392
## 2023-09-16 21:19:36 Annotating text fragment 33281/100392
## 2023-09-16 21:19:37 Annotating text fragment 33291/100392
## 2023-09-16 21:19:37 Annotating text fragment 33301/100392
## 2023-09-16 21:19:37 Annotating text fragment 33311/100392
## 2023-09-16 21:19:37 Annotating text fragment 33321/100392
## 2023-09-16 21:19:37 Annotating text fragment 33331/100392
## 2023-09-16 21:19:37 Annotating text fragment 33341/100392
## 2023-09-16 21:19:37 Annotating text fragment 33351/100392
## 2023-09-16 21:19:38 Annotating text fragment 33361/100392
## 2023-09-16 21:19:38 Annotating text fragment 33371/100392
## 2023-09-16 21:19:38 Annotating text fragment 33381/100392
## 2023-09-16 21:19:38 Annotating text fragment 33391/100392
## 2023-09-16 21:19:38 Annotating text fragment 33401/100392
## 2023-09-16 21:19:38 Annotating text fragment 33411/100392
## 2023-09-16 21:19:38 Annotating text fragment 33421/100392
## 2023-09-16 21:19:38 Annotating text fragment 33431/100392
## 2023-09-16 21:19:38 Annotating text fragment 33441/100392
## 2023-09-16 21:19:38 Annotating text fragment 33451/100392
## 2023-09-16 21:19:39 Annotating text fragment 33461/100392
## 2023-09-16 21:19:39 Annotating text fragment 33471/100392
## 2023-09-16 21:19:39 Annotating text fragment 33481/100392
## 2023-09-16 21:19:39 Annotating text fragment 33491/100392
## 2023-09-16 21:19:39 Annotating text fragment 33501/100392
## 2023-09-16 21:19:39 Annotating text fragment 33511/100392
## 2023-09-16 21:19:39 Annotating text fragment 33521/100392
## 2023-09-16 21:19:39 Annotating text fragment 33531/100392
```

```
## 2023-09-16 21:19:40 Annotating text fragment 33541/100392
## 2023-09-16 21:19:40 Annotating text fragment 33551/100392
## 2023-09-16 21:19:41 Annotating text fragment 33561/100392
## 2023-09-16 21:19:43 Annotating text fragment 33571/100392
## 2023-09-16 21:19:43 Annotating text fragment 33581/100392
## 2023-09-16 21:19:43 Annotating text fragment 33591/100392
## 2023-09-16 21:19:44 Annotating text fragment 33601/100392
## 2023-09-16 21:19:44 Annotating text fragment 33611/100392
## 2023-09-16 21:19:44 Annotating text fragment 33621/100392
## 2023-09-16 21:19:44 Annotating text fragment 33631/100392
## 2023-09-16 21:19:44 Annotating text fragment 33641/100392
## 2023-09-16 21:19:44 Annotating text fragment 33651/100392
## 2023-09-16 21:19:44 Annotating text fragment 33661/100392
## 2023-09-16 21:19:44 Annotating text fragment 33671/100392
## 2023-09-16 21:19:44 Annotating text fragment 33681/100392
## 2023-09-16 21:19:44 Annotating text fragment 33691/100392
## 2023-09-16 21:19:45 Annotating text fragment 33701/100392
## 2023-09-16 21:19:45 Annotating text fragment 33711/100392
## 2023-09-16 21:19:45 Annotating text fragment 33721/100392
## 2023-09-16 21:19:45 Annotating text fragment 33731/100392
## 2023-09-16 21:19:45 Annotating text fragment 33741/100392
## 2023-09-16 21:19:45 Annotating text fragment 33751/100392
## 2023-09-16 21:19:45 Annotating text fragment 33761/100392
## 2023-09-16 21:19:45 Annotating text fragment 33771/100392
## 2023-09-16 21:19:45 Annotating text fragment 33781/100392
## 2023-09-16 21:19:45 Annotating text fragment 33791/100392
## 2023-09-16 21:19:46 Annotating text fragment 33801/100392
## 2023-09-16 21:19:46 Annotating text fragment 33811/100392
## 2023-09-16 21:19:46 Annotating text fragment 33821/100392
## 2023-09-16 21:19:46 Annotating text fragment 33831/100392
## 2023-09-16 21:19:46 Annotating text fragment 33841/100392
## 2023-09-16 21:19:46 Annotating text fragment 33851/100392
## 2023-09-16 21:19:46 Annotating text fragment 33861/100392
## 2023-09-16 21:19:47 Annotating text fragment 33871/100392
## 2023-09-16 21:19:47 Annotating text fragment 33881/100392
## 2023-09-16 21:19:47 Annotating text fragment 33891/100392
## 2023-09-16 21:19:47 Annotating text fragment 33901/100392
## 2023-09-16 21:19:47 Annotating text fragment 33911/100392
## 2023-09-16 21:19:47 Annotating text fragment 33921/100392
## 2023-09-16 21:19:47 Annotating text fragment 33931/100392
## 2023-09-16 21:19:48 Annotating text fragment 33941/100392
## 2023-09-16 21:19:48 Annotating text fragment 33951/100392
## 2023-09-16 21:19:48 Annotating text fragment 33961/100392
## 2023-09-16 21:19:48 Annotating text fragment 33971/100392
## 2023-09-16 21:19:48 Annotating text fragment 33981/100392
## 2023-09-16 21:19:48 Annotating text fragment 33991/100392
## 2023-09-16 21:19:48 Annotating text fragment 34001/100392
## 2023-09-16 21:19:49 Annotating text fragment 34011/100392
## 2023-09-16 21:19:49 Annotating text fragment 34021/100392
## 2023-09-16 21:19:49 Annotating text fragment 34031/100392
## 2023-09-16 21:19:49 Annotating text fragment 34041/100392
## 2023-09-16 21:19:49 Annotating text fragment 34051/100392
## 2023-09-16 21:19:49 Annotating text fragment 34061/100392
## 2023-09-16 21:19:49 Annotating text fragment 34071/100392
```

```
## 2023-09-16 21:19:49 Annotating text fragment 34081/100392
## 2023-09-16 21:19:49 Annotating text fragment 34091/100392
## 2023-09-16 21:19:50 Annotating text fragment 34101/100392
## 2023-09-16 21:19:50 Annotating text fragment 34111/100392
## 2023-09-16 21:19:50 Annotating text fragment 34121/100392
## 2023-09-16 21:19:50 Annotating text fragment 34131/100392
## 2023-09-16 21:19:50 Annotating text fragment 34141/100392
## 2023-09-16 21:19:50 Annotating text fragment 34151/100392
## 2023-09-16 21:19:50 Annotating text fragment 34161/100392
## 2023-09-16 21:19:50 Annotating text fragment 34171/100392
## 2023-09-16 21:19:50 Annotating text fragment 34181/100392
## 2023-09-16 21:19:51 Annotating text fragment 34191/100392
## 2023-09-16 21:19:51 Annotating text fragment 34201/100392
## 2023-09-16 21:19:51 Annotating text fragment 34211/100392
## 2023-09-16 21:19:51 Annotating text fragment 34221/100392
## 2023-09-16 21:19:51 Annotating text fragment 34231/100392
## 2023-09-16 21:19:51 Annotating text fragment 34241/100392
## 2023-09-16 21:19:51 Annotating text fragment 34251/100392
## 2023-09-16 21:19:51 Annotating text fragment 34261/100392
## 2023-09-16 21:19:52 Annotating text fragment 34271/100392
## 2023-09-16 21:19:52 Annotating text fragment 34281/100392
## 2023-09-16 21:19:52 Annotating text fragment 34291/100392
## 2023-09-16 21:19:52 Annotating text fragment 34301/100392
## 2023-09-16 21:19:52 Annotating text fragment 34311/100392
## 2023-09-16 21:19:52 Annotating text fragment 34321/100392
## 2023-09-16 21:19:52 Annotating text fragment 34331/100392
## 2023-09-16 21:19:52 Annotating text fragment 34341/100392
## 2023-09-16 21:19:52 Annotating text fragment 34351/100392
## 2023-09-16 21:19:52 Annotating text fragment 34361/100392
## 2023-09-16 21:19:52 Annotating text fragment 34371/100392
## 2023-09-16 21:19:53 Annotating text fragment 34381/100392
## 2023-09-16 21:19:53 Annotating text fragment 34391/100392
## 2023-09-16 21:19:53 Annotating text fragment 34401/100392
## 2023-09-16 21:19:53 Annotating text fragment 34411/100392
## 2023-09-16 21:19:53 Annotating text fragment 34421/100392
## 2023-09-16 21:19:53 Annotating text fragment 34431/100392
## 2023-09-16 21:19:53 Annotating text fragment 34441/100392
## 2023-09-16 21:19:54 Annotating text fragment 34451/100392
## 2023-09-16 21:19:54 Annotating text fragment 34461/100392
## 2023-09-16 21:19:54 Annotating text fragment 34471/100392
## 2023-09-16 21:19:54 Annotating text fragment 34481/100392
## 2023-09-16 21:19:54 Annotating text fragment 34491/100392
## 2023-09-16 21:19:54 Annotating text fragment 34501/100392
## 2023-09-16 21:19:54 Annotating text fragment 34511/100392
## 2023-09-16 21:19:54 Annotating text fragment 34521/100392
## 2023-09-16 21:19:54 Annotating text fragment 34531/100392
## 2023-09-16 21:19:55 Annotating text fragment 34541/100392
## 2023-09-16 21:19:55 Annotating text fragment 34551/100392
## 2023-09-16 21:19:55 Annotating text fragment 34561/100392
## 2023-09-16 21:19:55 Annotating text fragment 34571/100392
## 2023-09-16 21:19:55 Annotating text fragment 34581/100392
## 2023-09-16 21:19:55 Annotating text fragment 34591/100392
## 2023-09-16 21:19:55 Annotating text fragment 34601/100392
## 2023-09-16 21:19:55 Annotating text fragment 34611/100392
```

```
## 2023-09-16 21:19:55 Annotating text fragment 34621/100392
## 2023-09-16 21:19:56 Annotating text fragment 34631/100392
## 2023-09-16 21:19:56 Annotating text fragment 34641/100392
## 2023-09-16 21:19:56 Annotating text fragment 34651/100392
## 2023-09-16 21:19:56 Annotating text fragment 34661/100392
## 2023-09-16 21:19:56 Annotating text fragment 34671/100392
## 2023-09-16 21:19:56 Annotating text fragment 34681/100392
## 2023-09-16 21:19:56 Annotating text fragment 34691/100392
## 2023-09-16 21:19:57 Annotating text fragment 34701/100392
## 2023-09-16 21:19:57 Annotating text fragment 34711/100392
## 2023-09-16 21:19:57 Annotating text fragment 34721/100392
## 2023-09-16 21:19:57 Annotating text fragment 34731/100392
## 2023-09-16 21:19:57 Annotating text fragment 34741/100392
## 2023-09-16 21:19:57 Annotating text fragment 34751/100392
## 2023-09-16 21:19:57 Annotating text fragment 34761/100392
## 2023-09-16 21:19:57 Annotating text fragment 34771/100392
## 2023-09-16 21:19:57 Annotating text fragment 34781/100392
## 2023-09-16 21:19:57 Annotating text fragment 34791/100392
## 2023-09-16 21:19:58 Annotating text fragment 34801/100392
## 2023-09-16 21:19:58 Annotating text fragment 34811/100392
## 2023-09-16 21:19:58 Annotating text fragment 34821/100392
## 2023-09-16 21:19:58 Annotating text fragment 34831/100392
## 2023-09-16 21:19:58 Annotating text fragment 34841/100392
## 2023-09-16 21:19:58 Annotating text fragment 34851/100392
## 2023-09-16 21:19:58 Annotating text fragment 34861/100392
## 2023-09-16 21:19:58 Annotating text fragment 34871/100392
## 2023-09-16 21:19:58 Annotating text fragment 34881/100392
## 2023-09-16 21:19:59 Annotating text fragment 34891/100392
## 2023-09-16 21:19:59 Annotating text fragment 34901/100392
## 2023-09-16 21:19:59 Annotating text fragment 34911/100392
## 2023-09-16 21:19:59 Annotating text fragment 34921/100392
## 2023-09-16 21:19:59 Annotating text fragment 34931/100392
## 2023-09-16 21:19:59 Annotating text fragment 34941/100392
## 2023-09-16 21:19:59 Annotating text fragment 34951/100392
## 2023-09-16 21:20:00 Annotating text fragment 34961/100392
## 2023-09-16 21:20:00 Annotating text fragment 34971/100392
## 2023-09-16 21:20:00 Annotating text fragment 34981/100392
## 2023-09-16 21:20:00 Annotating text fragment 34991/100392
## 2023-09-16 21:20:00 Annotating text fragment 35001/100392
## 2023-09-16 21:20:00 Annotating text fragment 35011/100392
## 2023-09-16 21:20:00 Annotating text fragment 35021/100392
## 2023-09-16 21:20:00 Annotating text fragment 35031/100392
## 2023-09-16 21:20:01 Annotating text fragment 35041/100392
## 2023-09-16 21:20:01 Annotating text fragment 35051/100392
## 2023-09-16 21:20:01 Annotating text fragment 35061/100392
## 2023-09-16 21:20:01 Annotating text fragment 35071/100392
## 2023-09-16 21:20:01 Annotating text fragment 35081/100392
## 2023-09-16 21:20:01 Annotating text fragment 35091/100392
## 2023-09-16 21:20:01 Annotating text fragment 35101/100392
## 2023-09-16 21:20:01 Annotating text fragment 35111/100392
## 2023-09-16 21:20:02 Annotating text fragment 35121/100392
## 2023-09-16 21:20:02 Annotating text fragment 35131/100392
## 2023-09-16 21:20:02 Annotating text fragment 35141/100392
## 2023-09-16 21:20:02 Annotating text fragment 35151/100392
```

```
## 2023-09-16 21:20:02 Annotating text fragment 35161/100392
## 2023-09-16 21:20:02 Annotating text fragment 35171/100392
## 2023-09-16 21:20:02 Annotating text fragment 35181/100392
## 2023-09-16 21:20:03 Annotating text fragment 35191/100392
## 2023-09-16 21:20:03 Annotating text fragment 35201/100392
## 2023-09-16 21:20:03 Annotating text fragment 35211/100392
## 2023-09-16 21:20:03 Annotating text fragment 35221/100392
## 2023-09-16 21:20:03 Annotating text fragment 35231/100392
## 2023-09-16 21:20:03 Annotating text fragment 35241/100392
## 2023-09-16 21:20:03 Annotating text fragment 35251/100392
## 2023-09-16 21:20:03 Annotating text fragment 35261/100392
## 2023-09-16 21:20:04 Annotating text fragment 35271/100392
## 2023-09-16 21:20:04 Annotating text fragment 35281/100392
## 2023-09-16 21:20:04 Annotating text fragment 35291/100392
## 2023-09-16 21:20:04 Annotating text fragment 35301/100392
## 2023-09-16 21:20:04 Annotating text fragment 35311/100392
## 2023-09-16 21:20:04 Annotating text fragment 35321/100392
## 2023-09-16 21:20:04 Annotating text fragment 35331/100392
## 2023-09-16 21:20:04 Annotating text fragment 35341/100392
## 2023-09-16 21:20:04 Annotating text fragment 35351/100392
## 2023-09-16 21:20:05 Annotating text fragment 35361/100392
## 2023-09-16 21:20:05 Annotating text fragment 35371/100392
## 2023-09-16 21:20:05 Annotating text fragment 35381/100392
## 2023-09-16 21:20:05 Annotating text fragment 35391/100392
## 2023-09-16 21:20:05 Annotating text fragment 35401/100392
## 2023-09-16 21:20:05 Annotating text fragment 35411/100392
## 2023-09-16 21:20:05 Annotating text fragment 35421/100392
## 2023-09-16 21:20:05 Annotating text fragment 35431/100392
## 2023-09-16 21:20:06 Annotating text fragment 35441/100392
## 2023-09-16 21:20:06 Annotating text fragment 35451/100392
## 2023-09-16 21:20:06 Annotating text fragment 35461/100392
## 2023-09-16 21:20:06 Annotating text fragment 35471/100392
## 2023-09-16 21:20:06 Annotating text fragment 35481/100392
## 2023-09-16 21:20:06 Annotating text fragment 35491/100392
## 2023-09-16 21:20:06 Annotating text fragment 35501/100392
## 2023-09-16 21:20:06 Annotating text fragment 35511/100392
## 2023-09-16 21:20:06 Annotating text fragment 35521/100392
## 2023-09-16 21:20:07 Annotating text fragment 35531/100392
## 2023-09-16 21:20:07 Annotating text fragment 35541/100392
## 2023-09-16 21:20:07 Annotating text fragment 35551/100392
## 2023-09-16 21:20:07 Annotating text fragment 35561/100392
## 2023-09-16 21:20:07 Annotating text fragment 35571/100392
## 2023-09-16 21:20:07 Annotating text fragment 35581/100392
## 2023-09-16 21:20:07 Annotating text fragment 35591/100392
## 2023-09-16 21:20:07 Annotating text fragment 35601/100392
## 2023-09-16 21:20:08 Annotating text fragment 35611/100392
## 2023-09-16 21:20:08 Annotating text fragment 35621/100392
## 2023-09-16 21:20:08 Annotating text fragment 35631/100392
## 2023-09-16 21:20:08 Annotating text fragment 35641/100392
## 2023-09-16 21:20:08 Annotating text fragment 35651/100392
## 2023-09-16 21:20:08 Annotating text fragment 35661/100392
## 2023-09-16 21:20:08 Annotating text fragment 35671/100392
## 2023-09-16 21:20:08 Annotating text fragment 35681/100392
## 2023-09-16 21:20:09 Annotating text fragment 35691/100392
```

```
## 2023-09-16 21:20:09 Annotating text fragment 35701/100392
## 2023-09-16 21:20:09 Annotating text fragment 35711/100392
## 2023-09-16 21:20:09 Annotating text fragment 35721/100392
## 2023-09-16 21:20:09 Annotating text fragment 35731/100392
## 2023-09-16 21:20:09 Annotating text fragment 35741/100392
## 2023-09-16 21:20:09 Annotating text fragment 35751/100392
## 2023-09-16 21:20:10 Annotating text fragment 35761/100392
## 2023-09-16 21:20:10 Annotating text fragment 35771/100392
## 2023-09-16 21:20:10 Annotating text fragment 35781/100392
## 2023-09-16 21:20:10 Annotating text fragment 35791/100392
## 2023-09-16 21:20:10 Annotating text fragment 35801/100392
## 2023-09-16 21:20:10 Annotating text fragment 35811/100392
## 2023-09-16 21:20:10 Annotating text fragment 35821/100392
## 2023-09-16 21:20:10 Annotating text fragment 35831/100392
## 2023-09-16 21:20:10 Annotating text fragment 35841/100392
## 2023-09-16 21:20:11 Annotating text fragment 35851/100392
## 2023-09-16 21:20:11 Annotating text fragment 35861/100392
## 2023-09-16 21:20:11 Annotating text fragment 35871/100392
## 2023-09-16 21:20:11 Annotating text fragment 35881/100392
## 2023-09-16 21:20:11 Annotating text fragment 35891/100392
## 2023-09-16 21:20:11 Annotating text fragment 35901/100392
## 2023-09-16 21:20:11 Annotating text fragment 35911/100392
## 2023-09-16 21:20:11 Annotating text fragment 35921/100392
## 2023-09-16 21:20:12 Annotating text fragment 35931/100392
## 2023-09-16 21:20:12 Annotating text fragment 35941/100392
## 2023-09-16 21:20:12 Annotating text fragment 35951/100392
## 2023-09-16 21:20:12 Annotating text fragment 35961/100392
## 2023-09-16 21:20:12 Annotating text fragment 35971/100392
## 2023-09-16 21:20:12 Annotating text fragment 35981/100392
## 2023-09-16 21:20:13 Annotating text fragment 35991/100392
## 2023-09-16 21:20:13 Annotating text fragment 36001/100392
## 2023-09-16 21:20:13 Annotating text fragment 36011/100392
## 2023-09-16 21:20:13 Annotating text fragment 36021/100392
## 2023-09-16 21:20:13 Annotating text fragment 36031/100392
## 2023-09-16 21:20:13 Annotating text fragment 36041/100392
## 2023-09-16 21:20:13 Annotating text fragment 36051/100392
## 2023-09-16 21:20:14 Annotating text fragment 36061/100392
## 2023-09-16 21:20:14 Annotating text fragment 36071/100392
## 2023-09-16 21:20:14 Annotating text fragment 36081/100392
## 2023-09-16 21:20:14 Annotating text fragment 36091/100392
## 2023-09-16 21:20:14 Annotating text fragment 36101/100392
## 2023-09-16 21:20:14 Annotating text fragment 36111/100392
## 2023-09-16 21:20:14 Annotating text fragment 36121/100392
## 2023-09-16 21:20:14 Annotating text fragment 36131/100392
## 2023-09-16 21:20:15 Annotating text fragment 36141/100392
## 2023-09-16 21:20:15 Annotating text fragment 36151/100392
## 2023-09-16 21:20:15 Annotating text fragment 36161/100392
## 2023-09-16 21:20:15 Annotating text fragment 36171/100392
## 2023-09-16 21:20:15 Annotating text fragment 36181/100392
## 2023-09-16 21:20:15 Annotating text fragment 36191/100392
## 2023-09-16 21:20:16 Annotating text fragment 36201/100392
## 2023-09-16 21:20:16 Annotating text fragment 36211/100392
## 2023-09-16 21:20:16 Annotating text fragment 36221/100392
## 2023-09-16 21:20:16 Annotating text fragment 36231/100392
```

```
## 2023-09-16 21:20:16 Annotating text fragment 36241/100392
## 2023-09-16 21:20:16 Annotating text fragment 36251/100392
## 2023-09-16 21:20:16 Annotating text fragment 36261/100392
## 2023-09-16 21:20:17 Annotating text fragment 36271/100392
## 2023-09-16 21:20:17 Annotating text fragment 36281/100392
## 2023-09-16 21:20:17 Annotating text fragment 36291/100392
## 2023-09-16 21:20:17 Annotating text fragment 36301/100392
## 2023-09-16 21:20:17 Annotating text fragment 36311/100392
## 2023-09-16 21:20:17 Annotating text fragment 36321/100392
## 2023-09-16 21:20:17 Annotating text fragment 36331/100392
## 2023-09-16 21:20:17 Annotating text fragment 36341/100392
## 2023-09-16 21:20:18 Annotating text fragment 36351/100392
## 2023-09-16 21:20:18 Annotating text fragment 36361/100392
## 2023-09-16 21:20:18 Annotating text fragment 36371/100392
## 2023-09-16 21:20:18 Annotating text fragment 36381/100392
## 2023-09-16 21:20:18 Annotating text fragment 36391/100392
## 2023-09-16 21:20:18 Annotating text fragment 36401/100392
## 2023-09-16 21:20:18 Annotating text fragment 36411/100392
## 2023-09-16 21:20:18 Annotating text fragment 36421/100392
## 2023-09-16 21:20:18 Annotating text fragment 36431/100392
## 2023-09-16 21:20:18 Annotating text fragment 36441/100392
## 2023-09-16 21:20:19 Annotating text fragment 36451/100392
## 2023-09-16 21:20:19 Annotating text fragment 36461/100392
## 2023-09-16 21:20:19 Annotating text fragment 36471/100392
## 2023-09-16 21:20:19 Annotating text fragment 36481/100392
## 2023-09-16 21:20:19 Annotating text fragment 36491/100392
## 2023-09-16 21:20:19 Annotating text fragment 36501/100392
## 2023-09-16 21:20:19 Annotating text fragment 36511/100392
## 2023-09-16 21:20:19 Annotating text fragment 36521/100392
## 2023-09-16 21:20:19 Annotating text fragment 36531/100392
## 2023-09-16 21:20:20 Annotating text fragment 36541/100392
## 2023-09-16 21:20:20 Annotating text fragment 36551/100392
## 2023-09-16 21:20:20 Annotating text fragment 36561/100392
## 2023-09-16 21:20:20 Annotating text fragment 36571/100392
## 2023-09-16 21:20:20 Annotating text fragment 36581/100392
## 2023-09-16 21:20:20 Annotating text fragment 36591/100392
## 2023-09-16 21:20:20 Annotating text fragment 36601/100392
## 2023-09-16 21:20:20 Annotating text fragment 36611/100392
## 2023-09-16 21:20:20 Annotating text fragment 36621/100392
## 2023-09-16 21:20:20 Annotating text fragment 36631/100392
## 2023-09-16 21:20:21 Annotating text fragment 36641/100392
## 2023-09-16 21:20:21 Annotating text fragment 36651/100392
## 2023-09-16 21:20:21 Annotating text fragment 36661/100392
## 2023-09-16 21:20:21 Annotating text fragment 36671/100392
## 2023-09-16 21:20:21 Annotating text fragment 36681/100392
## 2023-09-16 21:20:21 Annotating text fragment 36691/100392
## 2023-09-16 21:20:21 Annotating text fragment 36701/100392
## 2023-09-16 21:20:21 Annotating text fragment 36711/100392
## 2023-09-16 21:20:21 Annotating text fragment 36721/100392
## 2023-09-16 21:20:21 Annotating text fragment 36731/100392
## 2023-09-16 21:20:22 Annotating text fragment 36741/100392
## 2023-09-16 21:20:22 Annotating text fragment 36751/100392
## 2023-09-16 21:20:22 Annotating text fragment 36761/100392
## 2023-09-16 21:20:22 Annotating text fragment 36771/100392
```

```
## 2023-09-16 21:20:22 Annotating text fragment 36781/100392
## 2023-09-16 21:20:22 Annotating text fragment 36791/100392
## 2023-09-16 21:20:22 Annotating text fragment 36801/100392
## 2023-09-16 21:20:22 Annotating text fragment 36811/100392
## 2023-09-16 21:20:23 Annotating text fragment 36821/100392
## 2023-09-16 21:20:23 Annotating text fragment 36831/100392
## 2023-09-16 21:20:23 Annotating text fragment 36841/100392
## 2023-09-16 21:20:23 Annotating text fragment 36851/100392
## 2023-09-16 21:20:23 Annotating text fragment 36861/100392
## 2023-09-16 21:20:23 Annotating text fragment 36871/100392
## 2023-09-16 21:20:23 Annotating text fragment 36881/100392
## 2023-09-16 21:20:23 Annotating text fragment 36891/100392
## 2023-09-16 21:20:24 Annotating text fragment 36901/100392
## 2023-09-16 21:20:24 Annotating text fragment 36911/100392
## 2023-09-16 21:20:24 Annotating text fragment 36921/100392
## 2023-09-16 21:20:24 Annotating text fragment 36931/100392
## 2023-09-16 21:20:24 Annotating text fragment 36941/100392
## 2023-09-16 21:20:24 Annotating text fragment 36951/100392
## 2023-09-16 21:20:24 Annotating text fragment 36961/100392
## 2023-09-16 21:20:24 Annotating text fragment 36971/100392
## 2023-09-16 21:20:24 Annotating text fragment 36981/100392
## 2023-09-16 21:20:25 Annotating text fragment 36991/100392
## 2023-09-16 21:20:25 Annotating text fragment 37001/100392
## 2023-09-16 21:20:25 Annotating text fragment 37011/100392
## 2023-09-16 21:20:25 Annotating text fragment 37021/100392
## 2023-09-16 21:20:25 Annotating text fragment 37031/100392
## 2023-09-16 21:20:25 Annotating text fragment 37041/100392
## 2023-09-16 21:20:25 Annotating text fragment 37051/100392
## 2023-09-16 21:20:25 Annotating text fragment 37061/100392
## 2023-09-16 21:20:25 Annotating text fragment 37071/100392
## 2023-09-16 21:20:26 Annotating text fragment 37081/100392
## 2023-09-16 21:20:26 Annotating text fragment 37091/100392
## 2023-09-16 21:20:26 Annotating text fragment 37101/100392
## 2023-09-16 21:20:26 Annotating text fragment 37111/100392
## 2023-09-16 21:20:26 Annotating text fragment 37121/100392
## 2023-09-16 21:20:26 Annotating text fragment 37131/100392
## 2023-09-16 21:20:26 Annotating text fragment 37141/100392
## 2023-09-16 21:20:26 Annotating text fragment 37151/100392
## 2023-09-16 21:20:26 Annotating text fragment 37161/100392
## 2023-09-16 21:20:27 Annotating text fragment 37171/100392
## 2023-09-16 21:20:27 Annotating text fragment 37181/100392
## 2023-09-16 21:20:27 Annotating text fragment 37191/100392
## 2023-09-16 21:20:27 Annotating text fragment 37201/100392
## 2023-09-16 21:20:27 Annotating text fragment 37211/100392
## 2023-09-16 21:20:27 Annotating text fragment 37221/100392
## 2023-09-16 21:20:27 Annotating text fragment 37231/100392
## 2023-09-16 21:20:27 Annotating text fragment 37241/100392
## 2023-09-16 21:20:28 Annotating text fragment 37251/100392
## 2023-09-16 21:20:28 Annotating text fragment 37261/100392
## 2023-09-16 21:20:28 Annotating text fragment 37271/100392
## 2023-09-16 21:20:28 Annotating text fragment 37281/100392
## 2023-09-16 21:20:28 Annotating text fragment 37291/100392
## 2023-09-16 21:20:28 Annotating text fragment 37301/100392
## 2023-09-16 21:20:28 Annotating text fragment 37311/100392
```

```
## 2023-09-16 21:20:28 Annotating text fragment 37321/100392
## 2023-09-16 21:20:28 Annotating text fragment 37331/100392
## 2023-09-16 21:20:28 Annotating text fragment 37341/100392
## 2023-09-16 21:20:29 Annotating text fragment 37351/100392
## 2023-09-16 21:20:29 Annotating text fragment 37361/100392
## 2023-09-16 21:20:29 Annotating text fragment 37371/100392
## 2023-09-16 21:20:29 Annotating text fragment 37381/100392
## 2023-09-16 21:20:29 Annotating text fragment 37391/100392
## 2023-09-16 21:20:29 Annotating text fragment 37401/100392
## 2023-09-16 21:20:29 Annotating text fragment 37411/100392
## 2023-09-16 21:20:30 Annotating text fragment 37421/100392
## 2023-09-16 21:20:30 Annotating text fragment 37431/100392
## 2023-09-16 21:20:30 Annotating text fragment 37441/100392
## 2023-09-16 21:20:30 Annotating text fragment 37451/100392
## 2023-09-16 21:20:30 Annotating text fragment 37461/100392
## 2023-09-16 21:20:30 Annotating text fragment 37471/100392
## 2023-09-16 21:20:30 Annotating text fragment 37481/100392
## 2023-09-16 21:20:30 Annotating text fragment 37491/100392
## 2023-09-16 21:20:30 Annotating text fragment 37501/100392
## 2023-09-16 21:20:31 Annotating text fragment 37511/100392
## 2023-09-16 21:20:31 Annotating text fragment 37521/100392
## 2023-09-16 21:20:31 Annotating text fragment 37531/100392
## 2023-09-16 21:20:31 Annotating text fragment 37541/100392
## 2023-09-16 21:20:31 Annotating text fragment 37551/100392
## 2023-09-16 21:20:31 Annotating text fragment 37561/100392
## 2023-09-16 21:20:31 Annotating text fragment 37571/100392
## 2023-09-16 21:20:31 Annotating text fragment 37581/100392
## 2023-09-16 21:20:31 Annotating text fragment 37591/100392
## 2023-09-16 21:20:32 Annotating text fragment 37601/100392
## 2023-09-16 21:20:32 Annotating text fragment 37611/100392
## 2023-09-16 21:20:32 Annotating text fragment 37621/100392
## 2023-09-16 21:20:32 Annotating text fragment 37631/100392
## 2023-09-16 21:20:32 Annotating text fragment 37641/100392
## 2023-09-16 21:20:32 Annotating text fragment 37651/100392
## 2023-09-16 21:20:32 Annotating text fragment 37661/100392
## 2023-09-16 21:20:32 Annotating text fragment 37671/100392
## 2023-09-16 21:20:32 Annotating text fragment 37681/100392
## 2023-09-16 21:20:32 Annotating text fragment 37691/100392
## 2023-09-16 21:20:33 Annotating text fragment 37701/100392
## 2023-09-16 21:20:33 Annotating text fragment 37711/100392
## 2023-09-16 21:20:33 Annotating text fragment 37721/100392
## 2023-09-16 21:20:33 Annotating text fragment 37731/100392
## 2023-09-16 21:20:33 Annotating text fragment 37741/100392
## 2023-09-16 21:20:33 Annotating text fragment 37751/100392
## 2023-09-16 21:20:33 Annotating text fragment 37761/100392
## 2023-09-16 21:20:33 Annotating text fragment 37771/100392
## 2023-09-16 21:20:33 Annotating text fragment 37781/100392
## 2023-09-16 21:20:34 Annotating text fragment 37791/100392
## 2023-09-16 21:20:34 Annotating text fragment 37801/100392
## 2023-09-16 21:20:34 Annotating text fragment 37811/100392
## 2023-09-16 21:20:34 Annotating text fragment 37821/100392
## 2023-09-16 21:20:34 Annotating text fragment 37831/100392
## 2023-09-16 21:20:34 Annotating text fragment 37841/100392
## 2023-09-16 21:20:34 Annotating text fragment 37851/100392
```

```
## 2023-09-16 21:20:34 Annotating text fragment 37861/100392
## 2023-09-16 21:20:34 Annotating text fragment 37871/100392
## 2023-09-16 21:20:34 Annotating text fragment 37881/100392
## 2023-09-16 21:20:35 Annotating text fragment 37891/100392
## 2023-09-16 21:20:35 Annotating text fragment 37901/100392
## 2023-09-16 21:20:35 Annotating text fragment 37911/100392
## 2023-09-16 21:20:35 Annotating text fragment 37921/100392
## 2023-09-16 21:20:35 Annotating text fragment 37931/100392
## 2023-09-16 21:20:35 Annotating text fragment 37941/100392
## 2023-09-16 21:20:35 Annotating text fragment 37951/100392
## 2023-09-16 21:20:35 Annotating text fragment 37961/100392
## 2023-09-16 21:20:35 Annotating text fragment 37971/100392
## 2023-09-16 21:20:36 Annotating text fragment 37981/100392
## 2023-09-16 21:20:36 Annotating text fragment 37991/100392
## 2023-09-16 21:20:36 Annotating text fragment 38001/100392
## 2023-09-16 21:20:36 Annotating text fragment 38011/100392
## 2023-09-16 21:20:36 Annotating text fragment 38021/100392
## 2023-09-16 21:20:36 Annotating text fragment 38031/100392
## 2023-09-16 21:20:36 Annotating text fragment 38041/100392
## 2023-09-16 21:20:36 Annotating text fragment 38051/100392
## 2023-09-16 21:20:36 Annotating text fragment 38061/100392
## 2023-09-16 21:20:36 Annotating text fragment 38071/100392
## 2023-09-16 21:20:36 Annotating text fragment 38081/100392
## 2023-09-16 21:20:37 Annotating text fragment 38091/100392
## 2023-09-16 21:20:37 Annotating text fragment 38101/100392
## 2023-09-16 21:20:37 Annotating text fragment 38111/100392
## 2023-09-16 21:20:37 Annotating text fragment 38121/100392
## 2023-09-16 21:20:37 Annotating text fragment 38131/100392
## 2023-09-16 21:20:37 Annotating text fragment 38141/100392
## 2023-09-16 21:20:37 Annotating text fragment 38151/100392
## 2023-09-16 21:20:37 Annotating text fragment 38161/100392
## 2023-09-16 21:20:37 Annotating text fragment 38171/100392
## 2023-09-16 21:20:38 Annotating text fragment 38181/100392
## 2023-09-16 21:20:38 Annotating text fragment 38191/100392
## 2023-09-16 21:20:38 Annotating text fragment 38201/100392
## 2023-09-16 21:20:38 Annotating text fragment 38211/100392
## 2023-09-16 21:20:38 Annotating text fragment 38221/100392
## 2023-09-16 21:20:38 Annotating text fragment 38231/100392
## 2023-09-16 21:20:38 Annotating text fragment 38241/100392
## 2023-09-16 21:20:38 Annotating text fragment 38251/100392
## 2023-09-16 21:20:39 Annotating text fragment 38261/100392
## 2023-09-16 21:20:39 Annotating text fragment 38271/100392
## 2023-09-16 21:20:39 Annotating text fragment 38281/100392
## 2023-09-16 21:20:39 Annotating text fragment 38291/100392
## 2023-09-16 21:20:39 Annotating text fragment 38301/100392
## 2023-09-16 21:20:39 Annotating text fragment 38311/100392
## 2023-09-16 21:20:39 Annotating text fragment 38321/100392
## 2023-09-16 21:20:39 Annotating text fragment 38331/100392
## 2023-09-16 21:20:39 Annotating text fragment 38341/100392
## 2023-09-16 21:20:40 Annotating text fragment 38351/100392
## 2023-09-16 21:20:40 Annotating text fragment 38361/100392
## 2023-09-16 21:20:40 Annotating text fragment 38371/100392
## 2023-09-16 21:20:40 Annotating text fragment 38381/100392
## 2023-09-16 21:20:40 Annotating text fragment 38391/100392
```

```
## 2023-09-16 21:20:40 Annotating text fragment 38401/100392
## 2023-09-16 21:20:40 Annotating text fragment 38411/100392
## 2023-09-16 21:20:40 Annotating text fragment 38421/100392
## 2023-09-16 21:20:40 Annotating text fragment 38431/100392
## 2023-09-16 21:20:41 Annotating text fragment 38441/100392
## 2023-09-16 21:20:41 Annotating text fragment 38451/100392
## 2023-09-16 21:20:41 Annotating text fragment 38461/100392
## 2023-09-16 21:20:41 Annotating text fragment 38471/100392
## 2023-09-16 21:20:41 Annotating text fragment 38481/100392
## 2023-09-16 21:20:41 Annotating text fragment 38491/100392
## 2023-09-16 21:20:41 Annotating text fragment 38501/100392
## 2023-09-16 21:20:41 Annotating text fragment 38511/100392
## 2023-09-16 21:20:41 Annotating text fragment 38521/100392
## 2023-09-16 21:20:41 Annotating text fragment 38531/100392
## 2023-09-16 21:20:42 Annotating text fragment 38541/100392
## 2023-09-16 21:20:42 Annotating text fragment 38551/100392
## 2023-09-16 21:20:42 Annotating text fragment 38561/100392
## 2023-09-16 21:20:42 Annotating text fragment 38571/100392
## 2023-09-16 21:20:42 Annotating text fragment 38581/100392
## 2023-09-16 21:20:42 Annotating text fragment 38591/100392
## 2023-09-16 21:20:43 Annotating text fragment 38601/100392
## 2023-09-16 21:20:43 Annotating text fragment 38611/100392
## 2023-09-16 21:20:43 Annotating text fragment 38621/100392
## 2023-09-16 21:20:43 Annotating text fragment 38631/100392
## 2023-09-16 21:20:43 Annotating text fragment 38641/100392
## 2023-09-16 21:20:43 Annotating text fragment 38651/100392
## 2023-09-16 21:20:43 Annotating text fragment 38661/100392
## 2023-09-16 21:20:43 Annotating text fragment 38671/100392
## 2023-09-16 21:20:44 Annotating text fragment 38681/100392
## 2023-09-16 21:20:44 Annotating text fragment 38691/100392
## 2023-09-16 21:20:44 Annotating text fragment 38701/100392
## 2023-09-16 21:20:44 Annotating text fragment 38711/100392
## 2023-09-16 21:20:44 Annotating text fragment 38721/100392
## 2023-09-16 21:20:44 Annotating text fragment 38731/100392
## 2023-09-16 21:20:44 Annotating text fragment 38741/100392
## 2023-09-16 21:20:45 Annotating text fragment 38751/100392
## 2023-09-16 21:20:45 Annotating text fragment 38761/100392
## 2023-09-16 21:20:45 Annotating text fragment 38771/100392
## 2023-09-16 21:20:45 Annotating text fragment 38781/100392
## 2023-09-16 21:20:45 Annotating text fragment 38791/100392
## 2023-09-16 21:20:45 Annotating text fragment 38801/100392
## 2023-09-16 21:20:45 Annotating text fragment 38811/100392
## 2023-09-16 21:20:45 Annotating text fragment 38821/100392
## 2023-09-16 21:20:45 Annotating text fragment 38831/100392
## 2023-09-16 21:20:45 Annotating text fragment 38841/100392
## 2023-09-16 21:20:46 Annotating text fragment 38851/100392
## 2023-09-16 21:20:46 Annotating text fragment 38861/100392
## 2023-09-16 21:20:46 Annotating text fragment 38871/100392
## 2023-09-16 21:20:46 Annotating text fragment 38881/100392
## 2023-09-16 21:20:46 Annotating text fragment 38891/100392
## 2023-09-16 21:20:46 Annotating text fragment 38901/100392
## 2023-09-16 21:20:46 Annotating text fragment 38911/100392
## 2023-09-16 21:20:46 Annotating text fragment 38921/100392
## 2023-09-16 21:20:46 Annotating text fragment 38931/100392
```

```
## 2023-09-16 21:20:47 Annotating text fragment 38941/100392
## 2023-09-16 21:20:47 Annotating text fragment 38951/100392
## 2023-09-16 21:20:47 Annotating text fragment 38961/100392
## 2023-09-16 21:20:47 Annotating text fragment 38971/100392
## 2023-09-16 21:20:47 Annotating text fragment 38981/100392
## 2023-09-16 21:20:47 Annotating text fragment 38991/100392
## 2023-09-16 21:20:47 Annotating text fragment 39001/100392
## 2023-09-16 21:20:47 Annotating text fragment 39011/100392
## 2023-09-16 21:20:47 Annotating text fragment 39021/100392
## 2023-09-16 21:20:48 Annotating text fragment 39031/100392
## 2023-09-16 21:20:48 Annotating text fragment 39041/100392
## 2023-09-16 21:20:48 Annotating text fragment 39051/100392
## 2023-09-16 21:20:48 Annotating text fragment 39061/100392
## 2023-09-16 21:20:48 Annotating text fragment 39071/100392
## 2023-09-16 21:20:48 Annotating text fragment 39081/100392
## 2023-09-16 21:20:48 Annotating text fragment 39091/100392
## 2023-09-16 21:20:48 Annotating text fragment 39101/100392
## 2023-09-16 21:20:49 Annotating text fragment 39111/100392
## 2023-09-16 21:20:49 Annotating text fragment 39121/100392
## 2023-09-16 21:20:49 Annotating text fragment 39131/100392
## 2023-09-16 21:20:49 Annotating text fragment 39141/100392
## 2023-09-16 21:20:49 Annotating text fragment 39151/100392
## 2023-09-16 21:20:49 Annotating text fragment 39161/100392
## 2023-09-16 21:20:50 Annotating text fragment 39171/100392
## 2023-09-16 21:20:50 Annotating text fragment 39181/100392
## 2023-09-16 21:20:50 Annotating text fragment 39191/100392
## 2023-09-16 21:20:50 Annotating text fragment 39201/100392
## 2023-09-16 21:20:50 Annotating text fragment 39211/100392
## 2023-09-16 21:20:50 Annotating text fragment 39221/100392
## 2023-09-16 21:20:50 Annotating text fragment 39231/100392
## 2023-09-16 21:20:51 Annotating text fragment 39241/100392
## 2023-09-16 21:20:51 Annotating text fragment 39251/100392
## 2023-09-16 21:20:51 Annotating text fragment 39261/100392
## 2023-09-16 21:20:51 Annotating text fragment 39271/100392
## 2023-09-16 21:20:51 Annotating text fragment 39281/100392
## 2023-09-16 21:20:51 Annotating text fragment 39291/100392
## 2023-09-16 21:20:51 Annotating text fragment 39301/100392
## 2023-09-16 21:20:52 Annotating text fragment 39311/100392
## 2023-09-16 21:20:52 Annotating text fragment 39321/100392
## 2023-09-16 21:20:52 Annotating text fragment 39331/100392
## 2023-09-16 21:20:52 Annotating text fragment 39341/100392
## 2023-09-16 21:20:52 Annotating text fragment 39351/100392
## 2023-09-16 21:20:52 Annotating text fragment 39361/100392
## 2023-09-16 21:20:52 Annotating text fragment 39371/100392
## 2023-09-16 21:20:52 Annotating text fragment 39381/100392
## 2023-09-16 21:20:52 Annotating text fragment 39391/100392
## 2023-09-16 21:20:52 Annotating text fragment 39401/100392
## 2023-09-16 21:20:53 Annotating text fragment 39411/100392
## 2023-09-16 21:20:53 Annotating text fragment 39421/100392
## 2023-09-16 21:20:53 Annotating text fragment 39431/100392
## 2023-09-16 21:20:53 Annotating text fragment 39441/100392
## 2023-09-16 21:20:53 Annotating text fragment 39451/100392
## 2023-09-16 21:20:53 Annotating text fragment 39461/100392
## 2023-09-16 21:20:54 Annotating text fragment 39471/100392
```

```
## 2023-09-16 21:20:54 Annotating text fragment 39481/100392
## 2023-09-16 21:20:54 Annotating text fragment 39491/100392
## 2023-09-16 21:20:54 Annotating text fragment 39501/100392
## 2023-09-16 21:20:54 Annotating text fragment 39511/100392
## 2023-09-16 21:20:54 Annotating text fragment 39521/100392
## 2023-09-16 21:20:54 Annotating text fragment 39531/100392
## 2023-09-16 21:20:54 Annotating text fragment 39541/100392
## 2023-09-16 21:20:54 Annotating text fragment 39551/100392
## 2023-09-16 21:20:55 Annotating text fragment 39561/100392
## 2023-09-16 21:20:55 Annotating text fragment 39571/100392
## 2023-09-16 21:20:55 Annotating text fragment 39581/100392
## 2023-09-16 21:20:55 Annotating text fragment 39591/100392
## 2023-09-16 21:20:55 Annotating text fragment 39601/100392
## 2023-09-16 21:20:55 Annotating text fragment 39611/100392
## 2023-09-16 21:20:55 Annotating text fragment 39621/100392
## 2023-09-16 21:20:55 Annotating text fragment 39631/100392
## 2023-09-16 21:20:55 Annotating text fragment 39641/100392
## 2023-09-16 21:20:55 Annotating text fragment 39651/100392
## 2023-09-16 21:20:56 Annotating text fragment 39661/100392
## 2023-09-16 21:20:56 Annotating text fragment 39671/100392
## 2023-09-16 21:20:56 Annotating text fragment 39681/100392
## 2023-09-16 21:20:56 Annotating text fragment 39691/100392
## 2023-09-16 21:20:56 Annotating text fragment 39701/100392
## 2023-09-16 21:20:56 Annotating text fragment 39711/100392
## 2023-09-16 21:20:56 Annotating text fragment 39721/100392
## 2023-09-16 21:20:56 Annotating text fragment 39731/100392
## 2023-09-16 21:20:56 Annotating text fragment 39741/100392
## 2023-09-16 21:20:56 Annotating text fragment 39751/100392
## 2023-09-16 21:20:57 Annotating text fragment 39761/100392
## 2023-09-16 21:20:57 Annotating text fragment 39771/100392
## 2023-09-16 21:20:57 Annotating text fragment 39781/100392
## 2023-09-16 21:20:57 Annotating text fragment 39791/100392
## 2023-09-16 21:20:57 Annotating text fragment 39801/100392
## 2023-09-16 21:20:57 Annotating text fragment 39811/100392
## 2023-09-16 21:20:57 Annotating text fragment 39821/100392
## 2023-09-16 21:20:57 Annotating text fragment 39831/100392
## 2023-09-16 21:20:57 Annotating text fragment 39841/100392
## 2023-09-16 21:20:57 Annotating text fragment 39851/100392
## 2023-09-16 21:20:57 Annotating text fragment 39861/100392
## 2023-09-16 21:20:58 Annotating text fragment 39871/100392
## 2023-09-16 21:20:58 Annotating text fragment 39881/100392
## 2023-09-16 21:20:58 Annotating text fragment 39891/100392
## 2023-09-16 21:20:58 Annotating text fragment 39901/100392
## 2023-09-16 21:20:58 Annotating text fragment 39911/100392
## 2023-09-16 21:20:58 Annotating text fragment 39921/100392
## 2023-09-16 21:20:58 Annotating text fragment 39931/100392
## 2023-09-16 21:20:58 Annotating text fragment 39941/100392
## 2023-09-16 21:20:58 Annotating text fragment 39951/100392
## 2023-09-16 21:20:58 Annotating text fragment 39961/100392
## 2023-09-16 21:20:59 Annotating text fragment 39971/100392
## 2023-09-16 21:20:59 Annotating text fragment 39981/100392
## 2023-09-16 21:20:59 Annotating text fragment 39991/100392
## 2023-09-16 21:20:59 Annotating text fragment 40001/100392
## 2023-09-16 21:20:59 Annotating text fragment 40011/100392
```

```
## 2023-09-16 21:20:59 Annotating text fragment 40021/100392
## 2023-09-16 21:20:59 Annotating text fragment 40031/100392
## 2023-09-16 21:20:59 Annotating text fragment 40041/100392
## 2023-09-16 21:20:59 Annotating text fragment 40051/100392
## 2023-09-16 21:21:00 Annotating text fragment 40061/100392
## 2023-09-16 21:21:00 Annotating text fragment 40071/100392
## 2023-09-16 21:21:00 Annotating text fragment 40081/100392
## 2023-09-16 21:21:00 Annotating text fragment 40091/100392
## 2023-09-16 21:21:00 Annotating text fragment 40101/100392
## 2023-09-16 21:21:00 Annotating text fragment 40111/100392
## 2023-09-16 21:21:00 Annotating text fragment 40121/100392
## 2023-09-16 21:21:00 Annotating text fragment 40131/100392
## 2023-09-16 21:21:01 Annotating text fragment 40141/100392
## 2023-09-16 21:21:01 Annotating text fragment 40151/100392
## 2023-09-16 21:21:01 Annotating text fragment 40161/100392
## 2023-09-16 21:21:01 Annotating text fragment 40171/100392
## 2023-09-16 21:21:01 Annotating text fragment 40181/100392
## 2023-09-16 21:21:01 Annotating text fragment 40191/100392
## 2023-09-16 21:21:01 Annotating text fragment 40201/100392
## 2023-09-16 21:21:01 Annotating text fragment 40211/100392
## 2023-09-16 21:21:01 Annotating text fragment 40221/100392
## 2023-09-16 21:21:02 Annotating text fragment 40231/100392
## 2023-09-16 21:21:02 Annotating text fragment 40241/100392
## 2023-09-16 21:21:02 Annotating text fragment 40251/100392
## 2023-09-16 21:21:02 Annotating text fragment 40261/100392
## 2023-09-16 21:21:02 Annotating text fragment 40271/100392
## 2023-09-16 21:21:02 Annotating text fragment 40281/100392
## 2023-09-16 21:21:02 Annotating text fragment 40291/100392
## 2023-09-16 21:21:02 Annotating text fragment 40301/100392
## 2023-09-16 21:21:02 Annotating text fragment 40311/100392
## 2023-09-16 21:21:03 Annotating text fragment 40321/100392
## 2023-09-16 21:21:03 Annotating text fragment 40331/100392
## 2023-09-16 21:21:03 Annotating text fragment 40341/100392
## 2023-09-16 21:21:03 Annotating text fragment 40351/100392
## 2023-09-16 21:21:03 Annotating text fragment 40361/100392
## 2023-09-16 21:21:03 Annotating text fragment 40371/100392
## 2023-09-16 21:21:03 Annotating text fragment 40381/100392
## 2023-09-16 21:21:03 Annotating text fragment 40391/100392
## 2023-09-16 21:21:03 Annotating text fragment 40401/100392
## 2023-09-16 21:21:04 Annotating text fragment 40411/100392
## 2023-09-16 21:21:04 Annotating text fragment 40421/100392
## 2023-09-16 21:21:04 Annotating text fragment 40431/100392
## 2023-09-16 21:21:04 Annotating text fragment 40441/100392
## 2023-09-16 21:21:04 Annotating text fragment 40451/100392
## 2023-09-16 21:21:04 Annotating text fragment 40461/100392
## 2023-09-16 21:21:04 Annotating text fragment 40471/100392
## 2023-09-16 21:21:04 Annotating text fragment 40481/100392
## 2023-09-16 21:21:05 Annotating text fragment 40491/100392
## 2023-09-16 21:21:05 Annotating text fragment 40501/100392
## 2023-09-16 21:21:05 Annotating text fragment 40511/100392
## 2023-09-16 21:21:05 Annotating text fragment 40521/100392
## 2023-09-16 21:21:05 Annotating text fragment 40531/100392
## 2023-09-16 21:21:05 Annotating text fragment 40541/100392
## 2023-09-16 21:21:05 Annotating text fragment 40551/100392
```

```
## 2023-09-16 21:21:06 Annotating text fragment 40561/100392
## 2023-09-16 21:21:06 Annotating text fragment 40571/100392
## 2023-09-16 21:21:06 Annotating text fragment 40581/100392
## 2023-09-16 21:21:06 Annotating text fragment 40591/100392
## 2023-09-16 21:21:06 Annotating text fragment 40601/100392
## 2023-09-16 21:21:06 Annotating text fragment 40611/100392
## 2023-09-16 21:21:06 Annotating text fragment 40621/100392
## 2023-09-16 21:21:06 Annotating text fragment 40631/100392
## 2023-09-16 21:21:07 Annotating text fragment 40641/100392
## 2023-09-16 21:21:07 Annotating text fragment 40651/100392
## 2023-09-16 21:21:07 Annotating text fragment 40661/100392
## 2023-09-16 21:21:07 Annotating text fragment 40671/100392
## 2023-09-16 21:21:08 Annotating text fragment 40681/100392
## 2023-09-16 21:21:08 Annotating text fragment 40691/100392
## 2023-09-16 21:21:08 Annotating text fragment 40701/100392
## 2023-09-16 21:21:08 Annotating text fragment 40711/100392
## 2023-09-16 21:21:08 Annotating text fragment 40721/100392
## 2023-09-16 21:21:08 Annotating text fragment 40731/100392
## 2023-09-16 21:21:08 Annotating text fragment 40741/100392
## 2023-09-16 21:21:08 Annotating text fragment 40751/100392
## 2023-09-16 21:21:09 Annotating text fragment 40761/100392
## 2023-09-16 21:21:09 Annotating text fragment 40771/100392
## 2023-09-16 21:21:09 Annotating text fragment 40781/100392
## 2023-09-16 21:21:09 Annotating text fragment 40791/100392
## 2023-09-16 21:21:09 Annotating text fragment 40801/100392
## 2023-09-16 21:21:09 Annotating text fragment 40811/100392
## 2023-09-16 21:21:09 Annotating text fragment 40821/100392
## 2023-09-16 21:21:10 Annotating text fragment 40831/100392
## 2023-09-16 21:21:10 Annotating text fragment 40841/100392
## 2023-09-16 21:21:10 Annotating text fragment 40851/100392
## 2023-09-16 21:21:10 Annotating text fragment 40861/100392
## 2023-09-16 21:21:10 Annotating text fragment 40871/100392
## 2023-09-16 21:21:10 Annotating text fragment 40881/100392
## 2023-09-16 21:21:10 Annotating text fragment 40891/100392
## 2023-09-16 21:21:11 Annotating text fragment 40901/100392
## 2023-09-16 21:21:11 Annotating text fragment 40911/100392
## 2023-09-16 21:21:11 Annotating text fragment 40921/100392
## 2023-09-16 21:21:11 Annotating text fragment 40931/100392
## 2023-09-16 21:21:11 Annotating text fragment 40941/100392
## 2023-09-16 21:21:11 Annotating text fragment 40951/100392
## 2023-09-16 21:21:11 Annotating text fragment 40961/100392
## 2023-09-16 21:21:12 Annotating text fragment 40971/100392
## 2023-09-16 21:21:12 Annotating text fragment 40981/100392
## 2023-09-16 21:21:12 Annotating text fragment 40991/100392
## 2023-09-16 21:21:12 Annotating text fragment 41001/100392
## 2023-09-16 21:21:12 Annotating text fragment 41011/100392
## 2023-09-16 21:21:12 Annotating text fragment 41021/100392
## 2023-09-16 21:21:12 Annotating text fragment 41031/100392
## 2023-09-16 21:21:12 Annotating text fragment 41041/100392
## 2023-09-16 21:21:12 Annotating text fragment 41051/100392
## 2023-09-16 21:21:13 Annotating text fragment 41061/100392
## 2023-09-16 21:21:13 Annotating text fragment 41071/100392
## 2023-09-16 21:21:13 Annotating text fragment 41081/100392
## 2023-09-16 21:21:13 Annotating text fragment 41091/100392
```

```
## 2023-09-16 21:21:13 Annotating text fragment 41101/100392
## 2023-09-16 21:21:13 Annotating text fragment 41111/100392
## 2023-09-16 21:21:13 Annotating text fragment 41121/100392
## 2023-09-16 21:21:14 Annotating text fragment 41131/100392
## 2023-09-16 21:21:14 Annotating text fragment 41141/100392
## 2023-09-16 21:21:14 Annotating text fragment 41151/100392
## 2023-09-16 21:21:14 Annotating text fragment 41161/100392
## 2023-09-16 21:21:14 Annotating text fragment 41171/100392
## 2023-09-16 21:21:14 Annotating text fragment 41181/100392
## 2023-09-16 21:21:14 Annotating text fragment 41191/100392
## 2023-09-16 21:21:15 Annotating text fragment 41201/100392
## 2023-09-16 21:21:15 Annotating text fragment 41211/100392
## 2023-09-16 21:21:15 Annotating text fragment 41221/100392
## 2023-09-16 21:21:15 Annotating text fragment 41231/100392
## 2023-09-16 21:21:15 Annotating text fragment 41241/100392
## 2023-09-16 21:21:15 Annotating text fragment 41251/100392
## 2023-09-16 21:21:15 Annotating text fragment 41261/100392
## 2023-09-16 21:21:15 Annotating text fragment 41271/100392
## 2023-09-16 21:21:16 Annotating text fragment 41281/100392
## 2023-09-16 21:21:16 Annotating text fragment 41291/100392
## 2023-09-16 21:21:16 Annotating text fragment 41301/100392
## 2023-09-16 21:21:16 Annotating text fragment 41311/100392
## 2023-09-16 21:21:16 Annotating text fragment 41321/100392
## 2023-09-16 21:21:16 Annotating text fragment 41331/100392
## 2023-09-16 21:21:16 Annotating text fragment 41341/100392
## 2023-09-16 21:21:16 Annotating text fragment 41351/100392
## 2023-09-16 21:21:16 Annotating text fragment 41361/100392
## 2023-09-16 21:21:16 Annotating text fragment 41371/100392
## 2023-09-16 21:21:17 Annotating text fragment 41381/100392
## 2023-09-16 21:21:17 Annotating text fragment 41391/100392
## 2023-09-16 21:21:17 Annotating text fragment 41401/100392
## 2023-09-16 21:21:17 Annotating text fragment 41411/100392
## 2023-09-16 21:21:17 Annotating text fragment 41421/100392
## 2023-09-16 21:21:17 Annotating text fragment 41431/100392
## 2023-09-16 21:21:17 Annotating text fragment 41441/100392
## 2023-09-16 21:21:18 Annotating text fragment 41451/100392
## 2023-09-16 21:21:18 Annotating text fragment 41461/100392
## 2023-09-16 21:21:18 Annotating text fragment 41471/100392
## 2023-09-16 21:21:18 Annotating text fragment 41481/100392
## 2023-09-16 21:21:18 Annotating text fragment 41491/100392
## 2023-09-16 21:21:18 Annotating text fragment 41501/100392
## 2023-09-16 21:21:18 Annotating text fragment 41511/100392
## 2023-09-16 21:21:18 Annotating text fragment 41521/100392
## 2023-09-16 21:21:18 Annotating text fragment 41531/100392
## 2023-09-16 21:21:18 Annotating text fragment 41541/100392
## 2023-09-16 21:21:18 Annotating text fragment 41551/100392
## 2023-09-16 21:21:19 Annotating text fragment 41561/100392
## 2023-09-16 21:21:19 Annotating text fragment 41571/100392
## 2023-09-16 21:21:19 Annotating text fragment 41581/100392
## 2023-09-16 21:21:19 Annotating text fragment 41591/100392
## 2023-09-16 21:21:19 Annotating text fragment 41601/100392
## 2023-09-16 21:21:19 Annotating text fragment 41611/100392
## 2023-09-16 21:21:19 Annotating text fragment 41621/100392
## 2023-09-16 21:21:19 Annotating text fragment 41631/100392
```

```
## 2023-09-16 21:21:19 Annotating text fragment 41641/100392
## 2023-09-16 21:21:19 Annotating text fragment 41651/100392
## 2023-09-16 21:21:20 Annotating text fragment 41661/100392
## 2023-09-16 21:21:20 Annotating text fragment 41671/100392
## 2023-09-16 21:21:20 Annotating text fragment 41681/100392
## 2023-09-16 21:21:20 Annotating text fragment 41691/100392
## 2023-09-16 21:21:20 Annotating text fragment 41701/100392
## 2023-09-16 21:21:20 Annotating text fragment 41711/100392
## 2023-09-16 21:21:20 Annotating text fragment 41721/100392
## 2023-09-16 21:21:20 Annotating text fragment 41731/100392
## 2023-09-16 21:21:20 Annotating text fragment 41741/100392
## 2023-09-16 21:21:20 Annotating text fragment 41751/100392
## 2023-09-16 21:21:21 Annotating text fragment 41761/100392
## 2023-09-16 21:21:21 Annotating text fragment 41771/100392
## 2023-09-16 21:21:21 Annotating text fragment 41781/100392
## 2023-09-16 21:21:21 Annotating text fragment 41791/100392
## 2023-09-16 21:21:21 Annotating text fragment 41801/100392
## 2023-09-16 21:21:21 Annotating text fragment 41811/100392
## 2023-09-16 21:21:21 Annotating text fragment 41821/100392
## 2023-09-16 21:21:21 Annotating text fragment 41831/100392
## 2023-09-16 21:21:21 Annotating text fragment 41841/100392
## 2023-09-16 21:21:21 Annotating text fragment 41851/100392
## 2023-09-16 21:21:21 Annotating text fragment 41861/100392
## 2023-09-16 21:21:21 Annotating text fragment 41871/100392
## 2023-09-16 21:21:22 Annotating text fragment 41881/100392
## 2023-09-16 21:21:22 Annotating text fragment 41891/100392
## 2023-09-16 21:21:22 Annotating text fragment 41901/100392
## 2023-09-16 21:21:22 Annotating text fragment 41911/100392
## 2023-09-16 21:21:22 Annotating text fragment 41921/100392
## 2023-09-16 21:21:22 Annotating text fragment 41931/100392
## 2023-09-16 21:21:22 Annotating text fragment 41941/100392
## 2023-09-16 21:21:22 Annotating text fragment 41951/100392
## 2023-09-16 21:21:22 Annotating text fragment 41961/100392
## 2023-09-16 21:21:22 Annotating text fragment 41971/100392
## 2023-09-16 21:21:23 Annotating text fragment 41981/100392
## 2023-09-16 21:21:23 Annotating text fragment 41991/100392
## 2023-09-16 21:21:23 Annotating text fragment 42001/100392
## 2023-09-16 21:21:23 Annotating text fragment 42011/100392
## 2023-09-16 21:21:23 Annotating text fragment 42021/100392
## 2023-09-16 21:21:23 Annotating text fragment 42031/100392
## 2023-09-16 21:21:23 Annotating text fragment 42041/100392
## 2023-09-16 21:21:23 Annotating text fragment 42051/100392
## 2023-09-16 21:21:23 Annotating text fragment 42061/100392
## 2023-09-16 21:21:24 Annotating text fragment 42071/100392
## 2023-09-16 21:21:24 Annotating text fragment 42081/100392
## 2023-09-16 21:21:24 Annotating text fragment 42091/100392
## 2023-09-16 21:21:24 Annotating text fragment 42101/100392
## 2023-09-16 21:21:24 Annotating text fragment 42111/100392
## 2023-09-16 21:21:24 Annotating text fragment 42121/100392
## 2023-09-16 21:21:24 Annotating text fragment 42131/100392
## 2023-09-16 21:21:24 Annotating text fragment 42141/100392
## 2023-09-16 21:21:24 Annotating text fragment 42151/100392
## 2023-09-16 21:21:24 Annotating text fragment 42161/100392
## 2023-09-16 21:21:24 Annotating text fragment 42171/100392
```

```
## 2023-09-16 21:21:25 Annotating text fragment 42181/100392
## 2023-09-16 21:21:25 Annotating text fragment 42191/100392
## 2023-09-16 21:21:25 Annotating text fragment 42201/100392
## 2023-09-16 21:21:25 Annotating text fragment 42211/100392
## 2023-09-16 21:21:25 Annotating text fragment 42221/100392
## 2023-09-16 21:21:25 Annotating text fragment 42231/100392
## 2023-09-16 21:21:25 Annotating text fragment 42241/100392
## 2023-09-16 21:21:25 Annotating text fragment 42251/100392
## 2023-09-16 21:21:25 Annotating text fragment 42261/100392
## 2023-09-16 21:21:25 Annotating text fragment 42271/100392
## 2023-09-16 21:21:26 Annotating text fragment 42281/100392
## 2023-09-16 21:21:26 Annotating text fragment 42291/100392
## 2023-09-16 21:21:26 Annotating text fragment 42301/100392
## 2023-09-16 21:21:26 Annotating text fragment 42311/100392
## 2023-09-16 21:21:26 Annotating text fragment 42321/100392
## 2023-09-16 21:21:26 Annotating text fragment 42331/100392
## 2023-09-16 21:21:26 Annotating text fragment 42341/100392
## 2023-09-16 21:21:26 Annotating text fragment 42351/100392
## 2023-09-16 21:21:27 Annotating text fragment 42361/100392
## 2023-09-16 21:21:27 Annotating text fragment 42371/100392
## 2023-09-16 21:21:27 Annotating text fragment 42381/100392
## 2023-09-16 21:21:27 Annotating text fragment 42391/100392
## 2023-09-16 21:21:27 Annotating text fragment 42401/100392
## 2023-09-16 21:21:27 Annotating text fragment 42411/100392
## 2023-09-16 21:21:27 Annotating text fragment 42421/100392
## 2023-09-16 21:21:27 Annotating text fragment 42431/100392
## 2023-09-16 21:21:27 Annotating text fragment 42441/100392
## 2023-09-16 21:21:27 Annotating text fragment 42451/100392
## 2023-09-16 21:21:28 Annotating text fragment 42461/100392
## 2023-09-16 21:21:28 Annotating text fragment 42471/100392
## 2023-09-16 21:21:28 Annotating text fragment 42481/100392
## 2023-09-16 21:21:28 Annotating text fragment 42491/100392
## 2023-09-16 21:21:28 Annotating text fragment 42501/100392
## 2023-09-16 21:21:28 Annotating text fragment 42511/100392
## 2023-09-16 21:21:28 Annotating text fragment 42521/100392
## 2023-09-16 21:21:28 Annotating text fragment 42531/100392
## 2023-09-16 21:21:28 Annotating text fragment 42541/100392
## 2023-09-16 21:21:28 Annotating text fragment 42551/100392
## 2023-09-16 21:21:29 Annotating text fragment 42561/100392
## 2023-09-16 21:21:29 Annotating text fragment 42571/100392
## 2023-09-16 21:21:29 Annotating text fragment 42581/100392
## 2023-09-16 21:21:29 Annotating text fragment 42591/100392
## 2023-09-16 21:21:29 Annotating text fragment 42601/100392
## 2023-09-16 21:21:29 Annotating text fragment 42611/100392
## 2023-09-16 21:21:29 Annotating text fragment 42621/100392
## 2023-09-16 21:21:29 Annotating text fragment 42631/100392
## 2023-09-16 21:21:30 Annotating text fragment 42641/100392
## 2023-09-16 21:21:30 Annotating text fragment 42651/100392
## 2023-09-16 21:21:30 Annotating text fragment 42661/100392
## 2023-09-16 21:21:31 Annotating text fragment 42671/100392
## 2023-09-16 21:21:31 Annotating text fragment 42681/100392
## 2023-09-16 21:21:31 Annotating text fragment 42691/100392
## 2023-09-16 21:21:31 Annotating text fragment 42701/100392
## 2023-09-16 21:21:31 Annotating text fragment 42711/100392
```

```
## 2023-09-16 21:21:31 Annotating text fragment 42721/100392
## 2023-09-16 21:21:31 Annotating text fragment 42731/100392
## 2023-09-16 21:21:31 Annotating text fragment 42741/100392
## 2023-09-16 21:21:32 Annotating text fragment 42751/100392
## 2023-09-16 21:21:32 Annotating text fragment 42761/100392
## 2023-09-16 21:21:32 Annotating text fragment 42771/100392
## 2023-09-16 21:21:32 Annotating text fragment 42781/100392
## 2023-09-16 21:21:32 Annotating text fragment 42791/100392
## 2023-09-16 21:21:32 Annotating text fragment 42801/100392
## 2023-09-16 21:21:32 Annotating text fragment 42811/100392
## 2023-09-16 21:21:32 Annotating text fragment 42821/100392
## 2023-09-16 21:21:33 Annotating text fragment 42831/100392
## 2023-09-16 21:21:33 Annotating text fragment 42841/100392
## 2023-09-16 21:21:33 Annotating text fragment 42851/100392
## 2023-09-16 21:21:33 Annotating text fragment 42861/100392
## 2023-09-16 21:21:33 Annotating text fragment 42871/100392
## 2023-09-16 21:21:33 Annotating text fragment 42881/100392
## 2023-09-16 21:21:33 Annotating text fragment 42891/100392
## 2023-09-16 21:21:33 Annotating text fragment 42901/100392
## 2023-09-16 21:21:33 Annotating text fragment 42911/100392
## 2023-09-16 21:21:34 Annotating text fragment 42921/100392
## 2023-09-16 21:21:34 Annotating text fragment 42931/100392
## 2023-09-16 21:21:34 Annotating text fragment 42941/100392
## 2023-09-16 21:21:34 Annotating text fragment 42951/100392
## 2023-09-16 21:21:34 Annotating text fragment 42961/100392
## 2023-09-16 21:21:34 Annotating text fragment 42971/100392
## 2023-09-16 21:21:34 Annotating text fragment 42981/100392
## 2023-09-16 21:21:35 Annotating text fragment 42991/100392
## 2023-09-16 21:21:35 Annotating text fragment 43001/100392
## 2023-09-16 21:21:35 Annotating text fragment 43011/100392
## 2023-09-16 21:21:35 Annotating text fragment 43021/100392
## 2023-09-16 21:21:35 Annotating text fragment 43031/100392
## 2023-09-16 21:21:35 Annotating text fragment 43041/100392
## 2023-09-16 21:21:35 Annotating text fragment 43051/100392
## 2023-09-16 21:21:35 Annotating text fragment 43061/100392
## 2023-09-16 21:21:35 Annotating text fragment 43071/100392
## 2023-09-16 21:21:35 Annotating text fragment 43081/100392
## 2023-09-16 21:21:35 Annotating text fragment 43091/100392
## 2023-09-16 21:21:36 Annotating text fragment 43101/100392
## 2023-09-16 21:21:36 Annotating text fragment 43111/100392
## 2023-09-16 21:21:36 Annotating text fragment 43121/100392
## 2023-09-16 21:21:36 Annotating text fragment 43131/100392
## 2023-09-16 21:21:36 Annotating text fragment 43141/100392
## 2023-09-16 21:21:36 Annotating text fragment 43151/100392
## 2023-09-16 21:21:36 Annotating text fragment 43161/100392
## 2023-09-16 21:21:36 Annotating text fragment 43171/100392
## 2023-09-16 21:21:37 Annotating text fragment 43181/100392
## 2023-09-16 21:21:37 Annotating text fragment 43191/100392
## 2023-09-16 21:21:37 Annotating text fragment 43201/100392
## 2023-09-16 21:21:37 Annotating text fragment 43211/100392
## 2023-09-16 21:21:37 Annotating text fragment 43221/100392
## 2023-09-16 21:21:37 Annotating text fragment 43231/100392
## 2023-09-16 21:21:37 Annotating text fragment 43241/100392
## 2023-09-16 21:21:37 Annotating text fragment 43251/100392
```

```
## 2023-09-16 21:21:38 Annotating text fragment 43261/100392
## 2023-09-16 21:21:38 Annotating text fragment 43271/100392
## 2023-09-16 21:21:38 Annotating text fragment 43281/100392
## 2023-09-16 21:21:38 Annotating text fragment 43291/100392
## 2023-09-16 21:21:38 Annotating text fragment 43301/100392
## 2023-09-16 21:21:38 Annotating text fragment 43311/100392
## 2023-09-16 21:21:38 Annotating text fragment 43321/100392
## 2023-09-16 21:21:38 Annotating text fragment 43331/100392
## 2023-09-16 21:21:38 Annotating text fragment 43341/100392
## 2023-09-16 21:21:39 Annotating text fragment 43351/100392
## 2023-09-16 21:21:39 Annotating text fragment 43361/100392
## 2023-09-16 21:21:39 Annotating text fragment 43371/100392
## 2023-09-16 21:21:39 Annotating text fragment 43381/100392
## 2023-09-16 21:21:39 Annotating text fragment 43391/100392
## 2023-09-16 21:21:39 Annotating text fragment 43401/100392
## 2023-09-16 21:21:39 Annotating text fragment 43411/100392
## 2023-09-16 21:21:39 Annotating text fragment 43421/100392
## 2023-09-16 21:21:39 Annotating text fragment 43431/100392
## 2023-09-16 21:21:40 Annotating text fragment 43441/100392
## 2023-09-16 21:21:40 Annotating text fragment 43451/100392
## 2023-09-16 21:21:40 Annotating text fragment 43461/100392
## 2023-09-16 21:21:40 Annotating text fragment 43471/100392
## 2023-09-16 21:21:40 Annotating text fragment 43481/100392
## 2023-09-16 21:21:40 Annotating text fragment 43491/100392
## 2023-09-16 21:21:40 Annotating text fragment 43501/100392
## 2023-09-16 21:21:40 Annotating text fragment 43511/100392
## 2023-09-16 21:21:40 Annotating text fragment 43521/100392
## 2023-09-16 21:21:41 Annotating text fragment 43531/100392
## 2023-09-16 21:21:41 Annotating text fragment 43541/100392
## 2023-09-16 21:21:41 Annotating text fragment 43551/100392
## 2023-09-16 21:21:41 Annotating text fragment 43561/100392
## 2023-09-16 21:21:41 Annotating text fragment 43571/100392
## 2023-09-16 21:21:41 Annotating text fragment 43581/100392
## 2023-09-16 21:21:41 Annotating text fragment 43591/100392
## 2023-09-16 21:21:41 Annotating text fragment 43601/100392
## 2023-09-16 21:21:41 Annotating text fragment 43611/100392
## 2023-09-16 21:21:41 Annotating text fragment 43621/100392
## 2023-09-16 21:21:42 Annotating text fragment 43631/100392
## 2023-09-16 21:21:42 Annotating text fragment 43641/100392
## 2023-09-16 21:21:42 Annotating text fragment 43651/100392
## 2023-09-16 21:21:42 Annotating text fragment 43661/100392
## 2023-09-16 21:21:42 Annotating text fragment 43671/100392
## 2023-09-16 21:21:42 Annotating text fragment 43681/100392
## 2023-09-16 21:21:42 Annotating text fragment 43691/100392
## 2023-09-16 21:21:42 Annotating text fragment 43701/100392
## 2023-09-16 21:21:42 Annotating text fragment 43711/100392
## 2023-09-16 21:21:42 Annotating text fragment 43721/100392
## 2023-09-16 21:21:43 Annotating text fragment 43731/100392
## 2023-09-16 21:21:43 Annotating text fragment 43741/100392
## 2023-09-16 21:21:43 Annotating text fragment 43751/100392
## 2023-09-16 21:21:43 Annotating text fragment 43761/100392
## 2023-09-16 21:21:43 Annotating text fragment 43771/100392
## 2023-09-16 21:21:43 Annotating text fragment 43781/100392
## 2023-09-16 21:21:43 Annotating text fragment 43791/100392
```

```
## 2023-09-16 21:21:43 Annotating text fragment 43801/100392
## 2023-09-16 21:21:43 Annotating text fragment 43811/100392
## 2023-09-16 21:21:43 Annotating text fragment 43821/100392
## 2023-09-16 21:21:43 Annotating text fragment 43831/100392
## 2023-09-16 21:21:44 Annotating text fragment 43841/100392
## 2023-09-16 21:21:44 Annotating text fragment 43851/100392
## 2023-09-16 21:21:44 Annotating text fragment 43861/100392
## 2023-09-16 21:21:44 Annotating text fragment 43871/100392
## 2023-09-16 21:21:44 Annotating text fragment 43881/100392
## 2023-09-16 21:21:44 Annotating text fragment 43891/100392
## 2023-09-16 21:21:44 Annotating text fragment 43901/100392
## 2023-09-16 21:21:44 Annotating text fragment 43911/100392
## 2023-09-16 21:21:44 Annotating text fragment 43921/100392
## 2023-09-16 21:21:44 Annotating text fragment 43931/100392
## 2023-09-16 21:21:44 Annotating text fragment 43941/100392
## 2023-09-16 21:21:44 Annotating text fragment 43951/100392
## 2023-09-16 21:21:45 Annotating text fragment 43961/100392
## 2023-09-16 21:21:45 Annotating text fragment 43971/100392
## 2023-09-16 21:21:45 Annotating text fragment 43981/100392
## 2023-09-16 21:21:45 Annotating text fragment 43991/100392
## 2023-09-16 21:21:45 Annotating text fragment 44001/100392
## 2023-09-16 21:21:45 Annotating text fragment 44011/100392
## 2023-09-16 21:21:45 Annotating text fragment 44021/100392
## 2023-09-16 21:21:45 Annotating text fragment 44031/100392
## 2023-09-16 21:21:45 Annotating text fragment 44041/100392
## 2023-09-16 21:21:46 Annotating text fragment 44051/100392
## 2023-09-16 21:21:46 Annotating text fragment 44061/100392
## 2023-09-16 21:21:46 Annotating text fragment 44071/100392
## 2023-09-16 21:21:46 Annotating text fragment 44081/100392
## 2023-09-16 21:21:46 Annotating text fragment 44091/100392
## 2023-09-16 21:21:46 Annotating text fragment 44101/100392
## 2023-09-16 21:21:46 Annotating text fragment 44111/100392
## 2023-09-16 21:21:47 Annotating text fragment 44121/100392
## 2023-09-16 21:21:47 Annotating text fragment 44131/100392
## 2023-09-16 21:21:47 Annotating text fragment 44141/100392
## 2023-09-16 21:21:47 Annotating text fragment 44151/100392
## 2023-09-16 21:21:47 Annotating text fragment 44161/100392
## 2023-09-16 21:21:47 Annotating text fragment 44171/100392
## 2023-09-16 21:21:47 Annotating text fragment 44181/100392
## 2023-09-16 21:21:47 Annotating text fragment 44191/100392
## 2023-09-16 21:21:47 Annotating text fragment 44201/100392
## 2023-09-16 21:21:47 Annotating text fragment 44211/100392
## 2023-09-16 21:21:48 Annotating text fragment 44221/100392
## 2023-09-16 21:21:48 Annotating text fragment 44231/100392
## 2023-09-16 21:21:48 Annotating text fragment 44241/100392
## 2023-09-16 21:21:48 Annotating text fragment 44251/100392
## 2023-09-16 21:21:48 Annotating text fragment 44261/100392
## 2023-09-16 21:21:48 Annotating text fragment 44271/100392
## 2023-09-16 21:21:48 Annotating text fragment 44281/100392
## 2023-09-16 21:21:48 Annotating text fragment 44291/100392
## 2023-09-16 21:21:48 Annotating text fragment 44301/100392
## 2023-09-16 21:21:48 Annotating text fragment 44311/100392
## 2023-09-16 21:21:49 Annotating text fragment 44321/100392
## 2023-09-16 21:21:49 Annotating text fragment 44331/100392
```

```
## 2023-09-16 21:21:49 Annotating text fragment 44341/100392
## 2023-09-16 21:21:49 Annotating text fragment 44351/100392
## 2023-09-16 21:21:49 Annotating text fragment 44361/100392
## 2023-09-16 21:21:49 Annotating text fragment 44371/100392
## 2023-09-16 21:21:49 Annotating text fragment 44381/100392
## 2023-09-16 21:21:49 Annotating text fragment 44391/100392
## 2023-09-16 21:21:49 Annotating text fragment 44401/100392
## 2023-09-16 21:21:49 Annotating text fragment 44411/100392
## 2023-09-16 21:21:50 Annotating text fragment 44421/100392
## 2023-09-16 21:21:50 Annotating text fragment 44431/100392
## 2023-09-16 21:21:50 Annotating text fragment 44441/100392
## 2023-09-16 21:21:50 Annotating text fragment 44451/100392
## 2023-09-16 21:21:50 Annotating text fragment 44461/100392
## 2023-09-16 21:21:50 Annotating text fragment 44471/100392
## 2023-09-16 21:21:50 Annotating text fragment 44481/100392
## 2023-09-16 21:21:50 Annotating text fragment 44491/100392
## 2023-09-16 21:21:50 Annotating text fragment 44501/100392
## 2023-09-16 21:21:50 Annotating text fragment 44511/100392
## 2023-09-16 21:21:51 Annotating text fragment 44521/100392
## 2023-09-16 21:21:51 Annotating text fragment 44531/100392
## 2023-09-16 21:21:51 Annotating text fragment 44541/100392
## 2023-09-16 21:21:51 Annotating text fragment 44551/100392
## 2023-09-16 21:21:51 Annotating text fragment 44561/100392
## 2023-09-16 21:21:51 Annotating text fragment 44571/100392
## 2023-09-16 21:21:51 Annotating text fragment 44581/100392
## 2023-09-16 21:21:51 Annotating text fragment 44591/100392
## 2023-09-16 21:21:51 Annotating text fragment 44601/100392
## 2023-09-16 21:21:52 Annotating text fragment 44611/100392
## 2023-09-16 21:21:52 Annotating text fragment 44621/100392
## 2023-09-16 21:21:52 Annotating text fragment 44631/100392
## 2023-09-16 21:21:52 Annotating text fragment 44641/100392
## 2023-09-16 21:21:52 Annotating text fragment 44651/100392
## 2023-09-16 21:21:52 Annotating text fragment 44661/100392
## 2023-09-16 21:21:53 Annotating text fragment 44671/100392
## 2023-09-16 21:21:53 Annotating text fragment 44681/100392
## 2023-09-16 21:21:53 Annotating text fragment 44691/100392
## 2023-09-16 21:21:53 Annotating text fragment 44701/100392
## 2023-09-16 21:21:53 Annotating text fragment 44711/100392
## 2023-09-16 21:21:53 Annotating text fragment 44721/100392
## 2023-09-16 21:21:53 Annotating text fragment 44731/100392
## 2023-09-16 21:21:53 Annotating text fragment 44741/100392
## 2023-09-16 21:21:53 Annotating text fragment 44751/100392
## 2023-09-16 21:21:54 Annotating text fragment 44761/100392
## 2023-09-16 21:21:54 Annotating text fragment 44771/100392
## 2023-09-16 21:21:54 Annotating text fragment 44781/100392
## 2023-09-16 21:21:54 Annotating text fragment 44791/100392
## 2023-09-16 21:21:54 Annotating text fragment 44801/100392
## 2023-09-16 21:21:54 Annotating text fragment 44811/100392
## 2023-09-16 21:21:54 Annotating text fragment 44821/100392
## 2023-09-16 21:21:54 Annotating text fragment 44831/100392
## 2023-09-16 21:21:55 Annotating text fragment 44841/100392
## 2023-09-16 21:21:55 Annotating text fragment 44851/100392
## 2023-09-16 21:21:55 Annotating text fragment 44861/100392
## 2023-09-16 21:21:55 Annotating text fragment 44871/100392
```

```
## 2023-09-16 21:21:55 Annotating text fragment 44881/100392
## 2023-09-16 21:21:55 Annotating text fragment 44891/100392
## 2023-09-16 21:21:55 Annotating text fragment 44901/100392
## 2023-09-16 21:21:55 Annotating text fragment 44911/100392
## 2023-09-16 21:21:55 Annotating text fragment 44921/100392
## 2023-09-16 21:21:56 Annotating text fragment 44931/100392
## 2023-09-16 21:21:56 Annotating text fragment 44941/100392
## 2023-09-16 21:21:56 Annotating text fragment 44951/100392
## 2023-09-16 21:21:56 Annotating text fragment 44961/100392
## 2023-09-16 21:21:56 Annotating text fragment 44971/100392
## 2023-09-16 21:21:56 Annotating text fragment 44981/100392
## 2023-09-16 21:21:56 Annotating text fragment 44991/100392
## 2023-09-16 21:21:56 Annotating text fragment 45001/100392
## 2023-09-16 21:21:57 Annotating text fragment 45011/100392
## 2023-09-16 21:21:57 Annotating text fragment 45021/100392
## 2023-09-16 21:21:57 Annotating text fragment 45031/100392
## 2023-09-16 21:21:57 Annotating text fragment 45041/100392
## 2023-09-16 21:21:57 Annotating text fragment 45051/100392
## 2023-09-16 21:21:57 Annotating text fragment 45061/100392
## 2023-09-16 21:21:57 Annotating text fragment 45071/100392
## 2023-09-16 21:21:57 Annotating text fragment 45081/100392
## 2023-09-16 21:21:57 Annotating text fragment 45091/100392
## 2023-09-16 21:21:58 Annotating text fragment 45101/100392
## 2023-09-16 21:21:58 Annotating text fragment 45111/100392
## 2023-09-16 21:21:58 Annotating text fragment 45121/100392
## 2023-09-16 21:21:58 Annotating text fragment 45131/100392
## 2023-09-16 21:21:58 Annotating text fragment 45141/100392
## 2023-09-16 21:21:58 Annotating text fragment 45151/100392
## 2023-09-16 21:21:58 Annotating text fragment 45161/100392
## 2023-09-16 21:21:58 Annotating text fragment 45171/100392
## 2023-09-16 21:21:59 Annotating text fragment 45181/100392
## 2023-09-16 21:21:59 Annotating text fragment 45191/100392
## 2023-09-16 21:21:59 Annotating text fragment 45201/100392
## 2023-09-16 21:21:59 Annotating text fragment 45211/100392
## 2023-09-16 21:21:59 Annotating text fragment 45221/100392
## 2023-09-16 21:21:59 Annotating text fragment 45231/100392
## 2023-09-16 21:21:59 Annotating text fragment 45241/100392
## 2023-09-16 21:21:59 Annotating text fragment 45251/100392
## 2023-09-16 21:22:00 Annotating text fragment 45261/100392
## 2023-09-16 21:22:00 Annotating text fragment 45271/100392
## 2023-09-16 21:22:00 Annotating text fragment 45281/100392
## 2023-09-16 21:22:00 Annotating text fragment 45291/100392
## 2023-09-16 21:22:00 Annotating text fragment 45301/100392
## 2023-09-16 21:22:00 Annotating text fragment 45311/100392
## 2023-09-16 21:22:00 Annotating text fragment 45321/100392
## 2023-09-16 21:22:00 Annotating text fragment 45331/100392
## 2023-09-16 21:22:00 Annotating text fragment 45341/100392
## 2023-09-16 21:22:01 Annotating text fragment 45351/100392
## 2023-09-16 21:22:01 Annotating text fragment 45361/100392
## 2023-09-16 21:22:01 Annotating text fragment 45371/100392
## 2023-09-16 21:22:01 Annotating text fragment 45381/100392
## 2023-09-16 21:22:01 Annotating text fragment 45391/100392
## 2023-09-16 21:22:01 Annotating text fragment 45401/100392
## 2023-09-16 21:22:01 Annotating text fragment 45411/100392
```

```
## 2023-09-16 21:22:01 Annotating text fragment 45421/100392
## 2023-09-16 21:22:01 Annotating text fragment 45431/100392
## 2023-09-16 21:22:02 Annotating text fragment 45441/100392
## 2023-09-16 21:22:02 Annotating text fragment 45451/100392
## 2023-09-16 21:22:02 Annotating text fragment 45461/100392
## 2023-09-16 21:22:02 Annotating text fragment 45471/100392
## 2023-09-16 21:22:02 Annotating text fragment 45481/100392
## 2023-09-16 21:22:02 Annotating text fragment 45491/100392
## 2023-09-16 21:22:02 Annotating text fragment 45501/100392
## 2023-09-16 21:22:02 Annotating text fragment 45511/100392
## 2023-09-16 21:22:03 Annotating text fragment 45521/100392
## 2023-09-16 21:22:03 Annotating text fragment 45531/100392
## 2023-09-16 21:22:03 Annotating text fragment 45541/100392
## 2023-09-16 21:22:03 Annotating text fragment 45551/100392
## 2023-09-16 21:22:03 Annotating text fragment 45561/100392
## 2023-09-16 21:22:03 Annotating text fragment 45571/100392
## 2023-09-16 21:22:03 Annotating text fragment 45581/100392
## 2023-09-16 21:22:04 Annotating text fragment 45591/100392
## 2023-09-16 21:22:04 Annotating text fragment 45601/100392
## 2023-09-16 21:22:04 Annotating text fragment 45611/100392
## 2023-09-16 21:22:04 Annotating text fragment 45621/100392
## 2023-09-16 21:22:04 Annotating text fragment 45631/100392
## 2023-09-16 21:22:04 Annotating text fragment 45641/100392
## 2023-09-16 21:22:04 Annotating text fragment 45651/100392
## 2023-09-16 21:22:04 Annotating text fragment 45661/100392
## 2023-09-16 21:22:04 Annotating text fragment 45671/100392
## 2023-09-16 21:22:04 Annotating text fragment 45681/100392
## 2023-09-16 21:22:05 Annotating text fragment 45691/100392
## 2023-09-16 21:22:05 Annotating text fragment 45701/100392
## 2023-09-16 21:22:05 Annotating text fragment 45711/100392
## 2023-09-16 21:22:05 Annotating text fragment 45721/100392
## 2023-09-16 21:22:05 Annotating text fragment 45731/100392
## 2023-09-16 21:22:05 Annotating text fragment 45741/100392
## 2023-09-16 21:22:05 Annotating text fragment 45751/100392
## 2023-09-16 21:22:06 Annotating text fragment 45761/100392
## 2023-09-16 21:22:06 Annotating text fragment 45771/100392
## 2023-09-16 21:22:06 Annotating text fragment 45781/100392
## 2023-09-16 21:22:06 Annotating text fragment 45791/100392
## 2023-09-16 21:22:06 Annotating text fragment 45801/100392
## 2023-09-16 21:22:06 Annotating text fragment 45811/100392
## 2023-09-16 21:22:06 Annotating text fragment 45821/100392
## 2023-09-16 21:22:06 Annotating text fragment 45831/100392
## 2023-09-16 21:22:07 Annotating text fragment 45841/100392
## 2023-09-16 21:22:07 Annotating text fragment 45851/100392
## 2023-09-16 21:22:07 Annotating text fragment 45861/100392
## 2023-09-16 21:22:07 Annotating text fragment 45871/100392
## 2023-09-16 21:22:07 Annotating text fragment 45881/100392
## 2023-09-16 21:22:07 Annotating text fragment 45891/100392
## 2023-09-16 21:22:07 Annotating text fragment 45901/100392
## 2023-09-16 21:22:08 Annotating text fragment 45911/100392
## 2023-09-16 21:22:08 Annotating text fragment 45921/100392
## 2023-09-16 21:22:08 Annotating text fragment 45931/100392
## 2023-09-16 21:22:08 Annotating text fragment 45941/100392
## 2023-09-16 21:22:08 Annotating text fragment 45951/100392
```

```
## 2023-09-16 21:22:08 Annotating text fragment 45961/100392
## 2023-09-16 21:22:08 Annotating text fragment 45971/100392
## 2023-09-16 21:22:09 Annotating text fragment 45981/100392
## 2023-09-16 21:22:09 Annotating text fragment 45991/100392
## 2023-09-16 21:22:09 Annotating text fragment 46001/100392
## 2023-09-16 21:22:09 Annotating text fragment 46011/100392
## 2023-09-16 21:22:09 Annotating text fragment 46021/100392
## 2023-09-16 21:22:09 Annotating text fragment 46031/100392
## 2023-09-16 21:22:09 Annotating text fragment 46041/100392
## 2023-09-16 21:22:09 Annotating text fragment 46051/100392
## 2023-09-16 21:22:10 Annotating text fragment 46061/100392
## 2023-09-16 21:22:10 Annotating text fragment 46071/100392
## 2023-09-16 21:22:10 Annotating text fragment 46081/100392
## 2023-09-16 21:22:10 Annotating text fragment 46091/100392
## 2023-09-16 21:22:10 Annotating text fragment 46101/100392
## 2023-09-16 21:22:10 Annotating text fragment 46111/100392
## 2023-09-16 21:22:11 Annotating text fragment 46121/100392
## 2023-09-16 21:22:11 Annotating text fragment 46131/100392
## 2023-09-16 21:22:11 Annotating text fragment 46141/100392
## 2023-09-16 21:22:11 Annotating text fragment 46151/100392
## 2023-09-16 21:22:11 Annotating text fragment 46161/100392
## 2023-09-16 21:22:11 Annotating text fragment 46171/100392
## 2023-09-16 21:22:11 Annotating text fragment 46181/100392
## 2023-09-16 21:22:11 Annotating text fragment 46191/100392
## 2023-09-16 21:22:11 Annotating text fragment 46201/100392
## 2023-09-16 21:22:11 Annotating text fragment 46211/100392
## 2023-09-16 21:22:12 Annotating text fragment 46221/100392
## 2023-09-16 21:22:12 Annotating text fragment 46231/100392
## 2023-09-16 21:22:12 Annotating text fragment 46241/100392
## 2023-09-16 21:22:12 Annotating text fragment 46251/100392
## 2023-09-16 21:22:12 Annotating text fragment 46261/100392
## 2023-09-16 21:22:12 Annotating text fragment 46271/100392
## 2023-09-16 21:22:12 Annotating text fragment 46281/100392
## 2023-09-16 21:22:12 Annotating text fragment 46291/100392
## 2023-09-16 21:22:12 Annotating text fragment 46301/100392
## 2023-09-16 21:22:12 Annotating text fragment 46311/100392
## 2023-09-16 21:22:12 Annotating text fragment 46321/100392
## 2023-09-16 21:22:13 Annotating text fragment 46331/100392
## 2023-09-16 21:22:13 Annotating text fragment 46341/100392
## 2023-09-16 21:22:13 Annotating text fragment 46351/100392
## 2023-09-16 21:22:13 Annotating text fragment 46361/100392
## 2023-09-16 21:22:13 Annotating text fragment 46371/100392
## 2023-09-16 21:22:13 Annotating text fragment 46381/100392
## 2023-09-16 21:22:13 Annotating text fragment 46391/100392
## 2023-09-16 21:22:13 Annotating text fragment 46401/100392
## 2023-09-16 21:22:13 Annotating text fragment 46411/100392
## 2023-09-16 21:22:14 Annotating text fragment 46421/100392
## 2023-09-16 21:22:14 Annotating text fragment 46431/100392
## 2023-09-16 21:22:14 Annotating text fragment 46441/100392
## 2023-09-16 21:22:14 Annotating text fragment 46451/100392
## 2023-09-16 21:22:14 Annotating text fragment 46461/100392
## 2023-09-16 21:22:14 Annotating text fragment 46471/100392
## 2023-09-16 21:22:14 Annotating text fragment 46481/100392
## 2023-09-16 21:22:14 Annotating text fragment 46491/100392
```

```
## 2023-09-16 21:22:14 Annotating text fragment 46501/100392
## 2023-09-16 21:22:14 Annotating text fragment 46511/100392
## 2023-09-16 21:22:15 Annotating text fragment 46521/100392
## 2023-09-16 21:22:15 Annotating text fragment 46531/100392
## 2023-09-16 21:22:15 Annotating text fragment 46541/100392
## 2023-09-16 21:22:15 Annotating text fragment 46551/100392
## 2023-09-16 21:22:15 Annotating text fragment 46561/100392
## 2023-09-16 21:22:15 Annotating text fragment 46571/100392
## 2023-09-16 21:22:15 Annotating text fragment 46581/100392
## 2023-09-16 21:22:15 Annotating text fragment 46591/100392
## 2023-09-16 21:22:15 Annotating text fragment 46601/100392
## 2023-09-16 21:22:16 Annotating text fragment 46611/100392
## 2023-09-16 21:22:16 Annotating text fragment 46621/100392
## 2023-09-16 21:22:16 Annotating text fragment 46631/100392
## 2023-09-16 21:22:16 Annotating text fragment 46641/100392
## 2023-09-16 21:22:16 Annotating text fragment 46651/100392
## 2023-09-16 21:22:16 Annotating text fragment 46661/100392
## 2023-09-16 21:22:16 Annotating text fragment 46671/100392
## 2023-09-16 21:22:16 Annotating text fragment 46681/100392
## 2023-09-16 21:22:17 Annotating text fragment 46691/100392
## 2023-09-16 21:22:17 Annotating text fragment 46701/100392
## 2023-09-16 21:22:17 Annotating text fragment 46711/100392
## 2023-09-16 21:22:17 Annotating text fragment 46721/100392
## 2023-09-16 21:22:17 Annotating text fragment 46731/100392
## 2023-09-16 21:22:17 Annotating text fragment 46741/100392
## 2023-09-16 21:22:17 Annotating text fragment 46751/100392
## 2023-09-16 21:22:18 Annotating text fragment 46761/100392
## 2023-09-16 21:22:18 Annotating text fragment 46771/100392
## 2023-09-16 21:22:18 Annotating text fragment 46781/100392
## 2023-09-16 21:22:18 Annotating text fragment 46791/100392
## 2023-09-16 21:22:18 Annotating text fragment 46801/100392
## 2023-09-16 21:22:18 Annotating text fragment 46811/100392
## 2023-09-16 21:22:18 Annotating text fragment 46821/100392
## 2023-09-16 21:22:18 Annotating text fragment 46831/100392
## 2023-09-16 21:22:19 Annotating text fragment 46841/100392
## 2023-09-16 21:22:19 Annotating text fragment 46851/100392
## 2023-09-16 21:22:19 Annotating text fragment 46861/100392
## 2023-09-16 21:22:19 Annotating text fragment 46871/100392
## 2023-09-16 21:22:19 Annotating text fragment 46881/100392
## 2023-09-16 21:22:19 Annotating text fragment 46891/100392
## 2023-09-16 21:22:19 Annotating text fragment 46901/100392
## 2023-09-16 21:22:19 Annotating text fragment 46911/100392
## 2023-09-16 21:22:20 Annotating text fragment 46921/100392
## 2023-09-16 21:22:20 Annotating text fragment 46931/100392
## 2023-09-16 21:22:20 Annotating text fragment 46941/100392
## 2023-09-16 21:22:20 Annotating text fragment 46951/100392
## 2023-09-16 21:22:20 Annotating text fragment 46961/100392
## 2023-09-16 21:22:20 Annotating text fragment 46971/100392
## 2023-09-16 21:22:20 Annotating text fragment 46981/100392
## 2023-09-16 21:22:20 Annotating text fragment 46991/100392
## 2023-09-16 21:22:21 Annotating text fragment 47001/100392
## 2023-09-16 21:22:21 Annotating text fragment 47011/100392
## 2023-09-16 21:22:21 Annotating text fragment 47021/100392
## 2023-09-16 21:22:21 Annotating text fragment 47031/100392
```

```
## 2023-09-16 21:22:21 Annotating text fragment 47041/100392
## 2023-09-16 21:22:21 Annotating text fragment 47051/100392
## 2023-09-16 21:22:22 Annotating text fragment 47061/100392
## 2023-09-16 21:22:22 Annotating text fragment 47071/100392
## 2023-09-16 21:22:22 Annotating text fragment 47081/100392
## 2023-09-16 21:22:22 Annotating text fragment 47091/100392
## 2023-09-16 21:22:22 Annotating text fragment 47101/100392
## 2023-09-16 21:22:22 Annotating text fragment 47111/100392
## 2023-09-16 21:22:22 Annotating text fragment 47121/100392
## 2023-09-16 21:22:22 Annotating text fragment 47131/100392
## 2023-09-16 21:22:23 Annotating text fragment 47141/100392
## 2023-09-16 21:22:23 Annotating text fragment 47151/100392
## 2023-09-16 21:22:23 Annotating text fragment 47161/100392
## 2023-09-16 21:22:23 Annotating text fragment 47171/100392
## 2023-09-16 21:22:23 Annotating text fragment 47181/100392
## 2023-09-16 21:22:23 Annotating text fragment 47191/100392
## 2023-09-16 21:22:23 Annotating text fragment 47201/100392
## 2023-09-16 21:22:24 Annotating text fragment 47211/100392
## 2023-09-16 21:22:24 Annotating text fragment 47221/100392
## 2023-09-16 21:22:24 Annotating text fragment 47231/100392
## 2023-09-16 21:22:24 Annotating text fragment 47241/100392
## 2023-09-16 21:22:24 Annotating text fragment 47251/100392
## 2023-09-16 21:22:24 Annotating text fragment 47261/100392
## 2023-09-16 21:22:24 Annotating text fragment 47271/100392
## 2023-09-16 21:22:25 Annotating text fragment 47281/100392
## 2023-09-16 21:22:25 Annotating text fragment 47291/100392
## 2023-09-16 21:22:25 Annotating text fragment 47301/100392
## 2023-09-16 21:22:25 Annotating text fragment 47311/100392
## 2023-09-16 21:22:25 Annotating text fragment 47321/100392
## 2023-09-16 21:22:25 Annotating text fragment 47331/100392
## 2023-09-16 21:22:25 Annotating text fragment 47341/100392
## 2023-09-16 21:22:25 Annotating text fragment 47351/100392
## 2023-09-16 21:22:25 Annotating text fragment 47361/100392
## 2023-09-16 21:22:26 Annotating text fragment 47371/100392
## 2023-09-16 21:22:26 Annotating text fragment 47381/100392
## 2023-09-16 21:22:26 Annotating text fragment 47391/100392
## 2023-09-16 21:22:26 Annotating text fragment 47401/100392
## 2023-09-16 21:22:26 Annotating text fragment 47411/100392
## 2023-09-16 21:22:26 Annotating text fragment 47421/100392
## 2023-09-16 21:22:26 Annotating text fragment 47431/100392
## 2023-09-16 21:22:26 Annotating text fragment 47441/100392
## 2023-09-16 21:22:26 Annotating text fragment 47451/100392
## 2023-09-16 21:22:26 Annotating text fragment 47461/100392
## 2023-09-16 21:22:27 Annotating text fragment 47471/100392
## 2023-09-16 21:22:27 Annotating text fragment 47481/100392
## 2023-09-16 21:22:27 Annotating text fragment 47491/100392
## 2023-09-16 21:22:27 Annotating text fragment 47501/100392
## 2023-09-16 21:22:27 Annotating text fragment 47511/100392
## 2023-09-16 21:22:27 Annotating text fragment 47521/100392
## 2023-09-16 21:22:27 Annotating text fragment 47531/100392
## 2023-09-16 21:22:27 Annotating text fragment 47541/100392
## 2023-09-16 21:22:27 Annotating text fragment 47551/100392
## 2023-09-16 21:22:28 Annotating text fragment 47561/100392
## 2023-09-16 21:22:28 Annotating text fragment 47571/100392
```

```
## 2023-09-16 21:22:28 Annotating text fragment 47581/100392
## 2023-09-16 21:22:28 Annotating text fragment 47591/100392
## 2023-09-16 21:22:28 Annotating text fragment 47601/100392
## 2023-09-16 21:22:28 Annotating text fragment 47611/100392
## 2023-09-16 21:22:28 Annotating text fragment 47621/100392
## 2023-09-16 21:22:28 Annotating text fragment 47631/100392
## 2023-09-16 21:22:29 Annotating text fragment 47641/100392
## 2023-09-16 21:22:29 Annotating text fragment 47651/100392
## 2023-09-16 21:22:29 Annotating text fragment 47661/100392
## 2023-09-16 21:22:29 Annotating text fragment 47671/100392
## 2023-09-16 21:22:29 Annotating text fragment 47681/100392
## 2023-09-16 21:22:29 Annotating text fragment 47691/100392
## 2023-09-16 21:22:29 Annotating text fragment 47701/100392
## 2023-09-16 21:22:29 Annotating text fragment 47711/100392
## 2023-09-16 21:22:30 Annotating text fragment 47721/100392
## 2023-09-16 21:22:30 Annotating text fragment 47731/100392
## 2023-09-16 21:22:30 Annotating text fragment 47741/100392
## 2023-09-16 21:22:30 Annotating text fragment 47751/100392
## 2023-09-16 21:22:30 Annotating text fragment 47761/100392
## 2023-09-16 21:22:30 Annotating text fragment 47771/100392
## 2023-09-16 21:22:31 Annotating text fragment 47781/100392
## 2023-09-16 21:22:31 Annotating text fragment 47791/100392
## 2023-09-16 21:22:31 Annotating text fragment 47801/100392
## 2023-09-16 21:22:31 Annotating text fragment 47811/100392
## 2023-09-16 21:22:31 Annotating text fragment 47821/100392
## 2023-09-16 21:22:31 Annotating text fragment 47831/100392
## 2023-09-16 21:22:31 Annotating text fragment 47841/100392
## 2023-09-16 21:22:32 Annotating text fragment 47851/100392
## 2023-09-16 21:22:32 Annotating text fragment 47861/100392
## 2023-09-16 21:22:32 Annotating text fragment 47871/100392
## 2023-09-16 21:22:32 Annotating text fragment 47881/100392
## 2023-09-16 21:22:32 Annotating text fragment 47891/100392
## 2023-09-16 21:22:32 Annotating text fragment 47901/100392
## 2023-09-16 21:22:32 Annotating text fragment 47911/100392
## 2023-09-16 21:22:32 Annotating text fragment 47921/100392
## 2023-09-16 21:22:33 Annotating text fragment 47931/100392
## 2023-09-16 21:22:33 Annotating text fragment 47941/100392
## 2023-09-16 21:22:33 Annotating text fragment 47951/100392
## 2023-09-16 21:22:33 Annotating text fragment 47961/100392
## 2023-09-16 21:22:33 Annotating text fragment 47971/100392
## 2023-09-16 21:22:33 Annotating text fragment 47981/100392
## 2023-09-16 21:22:33 Annotating text fragment 47991/100392
## 2023-09-16 21:22:34 Annotating text fragment 48001/100392
## 2023-09-16 21:22:34 Annotating text fragment 48011/100392
## 2023-09-16 21:22:34 Annotating text fragment 48021/100392
## 2023-09-16 21:22:34 Annotating text fragment 48031/100392
## 2023-09-16 21:22:34 Annotating text fragment 48041/100392
## 2023-09-16 21:22:34 Annotating text fragment 48051/100392
## 2023-09-16 21:22:34 Annotating text fragment 48061/100392
## 2023-09-16 21:22:34 Annotating text fragment 48071/100392
## 2023-09-16 21:22:34 Annotating text fragment 48081/100392
## 2023-09-16 21:22:35 Annotating text fragment 48091/100392
## 2023-09-16 21:22:35 Annotating text fragment 48101/100392
## 2023-09-16 21:22:35 Annotating text fragment 48111/100392
```

```
## 2023-09-16 21:22:35 Annotating text fragment 48121/100392
## 2023-09-16 21:22:35 Annotating text fragment 48131/100392
## 2023-09-16 21:22:35 Annotating text fragment 48141/100392
## 2023-09-16 21:22:35 Annotating text fragment 48151/100392
## 2023-09-16 21:22:35 Annotating text fragment 48161/100392
## 2023-09-16 21:22:35 Annotating text fragment 48171/100392
## 2023-09-16 21:22:35 Annotating text fragment 48181/100392
## 2023-09-16 21:22:35 Annotating text fragment 48191/100392
## 2023-09-16 21:22:36 Annotating text fragment 48201/100392
## 2023-09-16 21:22:36 Annotating text fragment 48211/100392
## 2023-09-16 21:22:36 Annotating text fragment 48221/100392
## 2023-09-16 21:22:36 Annotating text fragment 48231/100392
## 2023-09-16 21:22:36 Annotating text fragment 48241/100392
## 2023-09-16 21:22:36 Annotating text fragment 48251/100392
## 2023-09-16 21:22:36 Annotating text fragment 48261/100392
## 2023-09-16 21:22:36 Annotating text fragment 48271/100392
## 2023-09-16 21:22:36 Annotating text fragment 48281/100392
## 2023-09-16 21:22:36 Annotating text fragment 48291/100392
## 2023-09-16 21:22:36 Annotating text fragment 48301/100392
## 2023-09-16 21:22:36 Annotating text fragment 48311/100392
## 2023-09-16 21:22:37 Annotating text fragment 48321/100392
## 2023-09-16 21:22:37 Annotating text fragment 48331/100392
## 2023-09-16 21:22:37 Annotating text fragment 48341/100392
## 2023-09-16 21:22:37 Annotating text fragment 48351/100392
## 2023-09-16 21:22:37 Annotating text fragment 48361/100392
## 2023-09-16 21:22:37 Annotating text fragment 48371/100392
## 2023-09-16 21:22:37 Annotating text fragment 48381/100392
## 2023-09-16 21:22:37 Annotating text fragment 48391/100392
## 2023-09-16 21:22:37 Annotating text fragment 48401/100392
## 2023-09-16 21:22:37 Annotating text fragment 48411/100392
## 2023-09-16 21:22:38 Annotating text fragment 48421/100392
## 2023-09-16 21:22:38 Annotating text fragment 48431/100392
## 2023-09-16 21:22:38 Annotating text fragment 48441/100392
## 2023-09-16 21:22:38 Annotating text fragment 48451/100392
## 2023-09-16 21:22:38 Annotating text fragment 48461/100392
## 2023-09-16 21:22:38 Annotating text fragment 48471/100392
## 2023-09-16 21:22:38 Annotating text fragment 48481/100392
## 2023-09-16 21:22:38 Annotating text fragment 48491/100392
## 2023-09-16 21:22:38 Annotating text fragment 48501/100392
## 2023-09-16 21:22:38 Annotating text fragment 48511/100392
## 2023-09-16 21:22:39 Annotating text fragment 48521/100392
## 2023-09-16 21:22:39 Annotating text fragment 48531/100392
## 2023-09-16 21:22:39 Annotating text fragment 48541/100392
## 2023-09-16 21:22:39 Annotating text fragment 48551/100392
## 2023-09-16 21:22:39 Annotating text fragment 48561/100392
## 2023-09-16 21:22:39 Annotating text fragment 48571/100392
## 2023-09-16 21:22:39 Annotating text fragment 48581/100392
## 2023-09-16 21:22:39 Annotating text fragment 48591/100392
## 2023-09-16 21:22:39 Annotating text fragment 48601/100392
## 2023-09-16 21:22:39 Annotating text fragment 48611/100392
## 2023-09-16 21:22:40 Annotating text fragment 48621/100392
## 2023-09-16 21:22:40 Annotating text fragment 48631/100392
## 2023-09-16 21:22:40 Annotating text fragment 48641/100392
## 2023-09-16 21:22:40 Annotating text fragment 48651/100392
```

```
## 2023-09-16 21:22:40 Annotating text fragment 48661/100392
## 2023-09-16 21:22:40 Annotating text fragment 48671/100392
## 2023-09-16 21:22:40 Annotating text fragment 48681/100392
## 2023-09-16 21:22:40 Annotating text fragment 48691/100392
## 2023-09-16 21:22:40 Annotating text fragment 48701/100392
## 2023-09-16 21:22:40 Annotating text fragment 48711/100392
## 2023-09-16 21:22:41 Annotating text fragment 48721/100392
## 2023-09-16 21:22:41 Annotating text fragment 48731/100392
## 2023-09-16 21:22:41 Annotating text fragment 48741/100392
## 2023-09-16 21:22:41 Annotating text fragment 48751/100392
## 2023-09-16 21:22:41 Annotating text fragment 48761/100392
## 2023-09-16 21:22:41 Annotating text fragment 48771/100392
## 2023-09-16 21:22:41 Annotating text fragment 48781/100392
## 2023-09-16 21:22:41 Annotating text fragment 48791/100392
## 2023-09-16 21:22:41 Annotating text fragment 48801/100392
## 2023-09-16 21:22:42 Annotating text fragment 48811/100392
## 2023-09-16 21:22:42 Annotating text fragment 48821/100392
## 2023-09-16 21:22:42 Annotating text fragment 48831/100392
## 2023-09-16 21:22:42 Annotating text fragment 48841/100392
## 2023-09-16 21:22:42 Annotating text fragment 48851/100392
## 2023-09-16 21:22:42 Annotating text fragment 48861/100392
## 2023-09-16 21:22:42 Annotating text fragment 48871/100392
## 2023-09-16 21:22:42 Annotating text fragment 48881/100392
## 2023-09-16 21:22:42 Annotating text fragment 48891/100392
## 2023-09-16 21:22:43 Annotating text fragment 48901/100392
## 2023-09-16 21:22:43 Annotating text fragment 48911/100392
## 2023-09-16 21:22:43 Annotating text fragment 48921/100392
## 2023-09-16 21:22:43 Annotating text fragment 48931/100392
## 2023-09-16 21:22:43 Annotating text fragment 48941/100392
## 2023-09-16 21:22:43 Annotating text fragment 48951/100392
## 2023-09-16 21:22:43 Annotating text fragment 48961/100392
## 2023-09-16 21:22:44 Annotating text fragment 48971/100392
## 2023-09-16 21:22:44 Annotating text fragment 48981/100392
## 2023-09-16 21:22:44 Annotating text fragment 48991/100392
## 2023-09-16 21:22:44 Annotating text fragment 49001/100392
## 2023-09-16 21:22:44 Annotating text fragment 49011/100392
## 2023-09-16 21:22:44 Annotating text fragment 49021/100392
## 2023-09-16 21:22:44 Annotating text fragment 49031/100392
## 2023-09-16 21:22:44 Annotating text fragment 49041/100392
## 2023-09-16 21:22:44 Annotating text fragment 49051/100392
## 2023-09-16 21:22:45 Annotating text fragment 49061/100392
## 2023-09-16 21:22:45 Annotating text fragment 49071/100392
## 2023-09-16 21:22:45 Annotating text fragment 49081/100392
## 2023-09-16 21:22:45 Annotating text fragment 49091/100392
## 2023-09-16 21:22:45 Annotating text fragment 49101/100392
## 2023-09-16 21:22:45 Annotating text fragment 49111/100392
## 2023-09-16 21:22:45 Annotating text fragment 49121/100392
## 2023-09-16 21:22:45 Annotating text fragment 49131/100392
## 2023-09-16 21:22:45 Annotating text fragment 49141/100392
## 2023-09-16 21:22:46 Annotating text fragment 49151/100392
## 2023-09-16 21:22:46 Annotating text fragment 49161/100392
## 2023-09-16 21:22:46 Annotating text fragment 49171/100392
## 2023-09-16 21:22:46 Annotating text fragment 49181/100392
## 2023-09-16 21:22:46 Annotating text fragment 49191/100392
```

```
## 2023-09-16 21:22:46 Annotating text fragment 49201/100392
## 2023-09-16 21:22:46 Annotating text fragment 49211/100392
## 2023-09-16 21:22:47 Annotating text fragment 49221/100392
## 2023-09-16 21:22:47 Annotating text fragment 49231/100392
## 2023-09-16 21:22:47 Annotating text fragment 49241/100392
## 2023-09-16 21:22:47 Annotating text fragment 49251/100392
## 2023-09-16 21:22:47 Annotating text fragment 49261/100392
## 2023-09-16 21:22:47 Annotating text fragment 49271/100392
## 2023-09-16 21:22:47 Annotating text fragment 49281/100392
## 2023-09-16 21:22:47 Annotating text fragment 49291/100392
## 2023-09-16 21:22:47 Annotating text fragment 49301/100392
## 2023-09-16 21:22:47 Annotating text fragment 49311/100392
## 2023-09-16 21:22:48 Annotating text fragment 49321/100392
## 2023-09-16 21:22:48 Annotating text fragment 49331/100392
## 2023-09-16 21:22:48 Annotating text fragment 49341/100392
## 2023-09-16 21:22:48 Annotating text fragment 49351/100392
## 2023-09-16 21:22:48 Annotating text fragment 49361/100392
## 2023-09-16 21:22:48 Annotating text fragment 49371/100392
## 2023-09-16 21:22:48 Annotating text fragment 49381/100392
## 2023-09-16 21:22:48 Annotating text fragment 49391/100392
## 2023-09-16 21:22:48 Annotating text fragment 49401/100392
## 2023-09-16 21:22:48 Annotating text fragment 49411/100392
## 2023-09-16 21:22:49 Annotating text fragment 49421/100392
## 2023-09-16 21:22:49 Annotating text fragment 49431/100392
## 2023-09-16 21:22:49 Annotating text fragment 49441/100392
## 2023-09-16 21:22:49 Annotating text fragment 49451/100392
## 2023-09-16 21:22:49 Annotating text fragment 49461/100392
## 2023-09-16 21:22:49 Annotating text fragment 49471/100392
## 2023-09-16 21:22:50 Annotating text fragment 49481/100392
## 2023-09-16 21:22:50 Annotating text fragment 49491/100392
## 2023-09-16 21:22:50 Annotating text fragment 49501/100392
## 2023-09-16 21:22:50 Annotating text fragment 49511/100392
## 2023-09-16 21:22:50 Annotating text fragment 49521/100392
## 2023-09-16 21:22:50 Annotating text fragment 49531/100392
## 2023-09-16 21:22:50 Annotating text fragment 49541/100392
## 2023-09-16 21:22:51 Annotating text fragment 49551/100392
## 2023-09-16 21:22:51 Annotating text fragment 49561/100392
## 2023-09-16 21:22:51 Annotating text fragment 49571/100392
## 2023-09-16 21:22:51 Annotating text fragment 49581/100392
## 2023-09-16 21:22:51 Annotating text fragment 49591/100392
## 2023-09-16 21:22:51 Annotating text fragment 49601/100392
## 2023-09-16 21:22:51 Annotating text fragment 49611/100392
## 2023-09-16 21:22:51 Annotating text fragment 49621/100392
## 2023-09-16 21:22:51 Annotating text fragment 49631/100392
## 2023-09-16 21:22:52 Annotating text fragment 49641/100392
## 2023-09-16 21:22:52 Annotating text fragment 49651/100392
## 2023-09-16 21:22:52 Annotating text fragment 49661/100392
## 2023-09-16 21:22:52 Annotating text fragment 49671/100392
## 2023-09-16 21:22:52 Annotating text fragment 49681/100392
## 2023-09-16 21:22:52 Annotating text fragment 49691/100392
## 2023-09-16 21:22:52 Annotating text fragment 49701/100392
## 2023-09-16 21:22:52 Annotating text fragment 49711/100392
## 2023-09-16 21:22:53 Annotating text fragment 49721/100392
## 2023-09-16 21:22:53 Annotating text fragment 49731/100392
```

```
## 2023-09-16 21:22:53 Annotating text fragment 49741/100392
## 2023-09-16 21:22:53 Annotating text fragment 49751/100392
## 2023-09-16 21:22:53 Annotating text fragment 49761/100392
## 2023-09-16 21:22:53 Annotating text fragment 49771/100392
## 2023-09-16 21:22:54 Annotating text fragment 49781/100392
## 2023-09-16 21:22:54 Annotating text fragment 49791/100392
## 2023-09-16 21:22:54 Annotating text fragment 49801/100392
## 2023-09-16 21:22:54 Annotating text fragment 49811/100392
## 2023-09-16 21:22:54 Annotating text fragment 49821/100392
## 2023-09-16 21:22:54 Annotating text fragment 49831/100392
## 2023-09-16 21:22:54 Annotating text fragment 49841/100392
## 2023-09-16 21:22:54 Annotating text fragment 49851/100392
## 2023-09-16 21:22:55 Annotating text fragment 49861/100392
## 2023-09-16 21:22:55 Annotating text fragment 49871/100392
## 2023-09-16 21:22:55 Annotating text fragment 49881/100392
## 2023-09-16 21:22:55 Annotating text fragment 49891/100392
## 2023-09-16 21:22:55 Annotating text fragment 49901/100392
## 2023-09-16 21:22:55 Annotating text fragment 49911/100392
## 2023-09-16 21:22:55 Annotating text fragment 49921/100392
## 2023-09-16 21:22:56 Annotating text fragment 49931/100392
## 2023-09-16 21:22:56 Annotating text fragment 49941/100392
## 2023-09-16 21:22:56 Annotating text fragment 49951/100392
## 2023-09-16 21:22:56 Annotating text fragment 49961/100392
## 2023-09-16 21:22:56 Annotating text fragment 49971/100392
## 2023-09-16 21:22:56 Annotating text fragment 49981/100392
## 2023-09-16 21:22:56 Annotating text fragment 49991/100392
## 2023-09-16 21:22:57 Annotating text fragment 50001/100392
## 2023-09-16 21:22:57 Annotating text fragment 50011/100392
## 2023-09-16 21:22:57 Annotating text fragment 50021/100392
## 2023-09-16 21:22:57 Annotating text fragment 50031/100392
## 2023-09-16 21:22:57 Annotating text fragment 50041/100392
## 2023-09-16 21:22:57 Annotating text fragment 50051/100392
## 2023-09-16 21:22:57 Annotating text fragment 50061/100392
## 2023-09-16 21:22:57 Annotating text fragment 50071/100392
## 2023-09-16 21:22:57 Annotating text fragment 50081/100392
## 2023-09-16 21:22:58 Annotating text fragment 50091/100392
## 2023-09-16 21:22:58 Annotating text fragment 50101/100392
## 2023-09-16 21:22:58 Annotating text fragment 50111/100392
## 2023-09-16 21:22:58 Annotating text fragment 50121/100392
## 2023-09-16 21:22:58 Annotating text fragment 50131/100392
## 2023-09-16 21:22:58 Annotating text fragment 50141/100392
## 2023-09-16 21:22:58 Annotating text fragment 50151/100392
## 2023-09-16 21:22:58 Annotating text fragment 50161/100392
## 2023-09-16 21:22:58 Annotating text fragment 50171/100392
## 2023-09-16 21:22:58 Annotating text fragment 50181/100392
## 2023-09-16 21:22:58 Annotating text fragment 50191/100392
## 2023-09-16 21:22:59 Annotating text fragment 50201/100392
## 2023-09-16 21:22:59 Annotating text fragment 50211/100392
## 2023-09-16 21:22:59 Annotating text fragment 50221/100392
## 2023-09-16 21:22:59 Annotating text fragment 50231/100392
## 2023-09-16 21:22:59 Annotating text fragment 50241/100392
## 2023-09-16 21:22:59 Annotating text fragment 50251/100392
## 2023-09-16 21:22:59 Annotating text fragment 50261/100392
## 2023-09-16 21:22:59 Annotating text fragment 50271/100392
```

```
## 2023-09-16 21:22:59 Annotating text fragment 50281/100392
## 2023-09-16 21:22:59 Annotating text fragment 50291/100392
## 2023-09-16 21:23:00 Annotating text fragment 50301/100392
## 2023-09-16 21:23:00 Annotating text fragment 50311/100392
## 2023-09-16 21:23:00 Annotating text fragment 50321/100392
## 2023-09-16 21:23:00 Annotating text fragment 50331/100392
## 2023-09-16 21:23:00 Annotating text fragment 50341/100392
## 2023-09-16 21:23:00 Annotating text fragment 50351/100392
## 2023-09-16 21:23:00 Annotating text fragment 50361/100392
## 2023-09-16 21:23:00 Annotating text fragment 50371/100392
## 2023-09-16 21:23:00 Annotating text fragment 50381/100392
## 2023-09-16 21:23:00 Annotating text fragment 50391/100392
## 2023-09-16 21:23:01 Annotating text fragment 50401/100392
## 2023-09-16 21:23:01 Annotating text fragment 50411/100392
## 2023-09-16 21:23:01 Annotating text fragment 50421/100392
## 2023-09-16 21:23:01 Annotating text fragment 50431/100392
## 2023-09-16 21:23:01 Annotating text fragment 50441/100392
## 2023-09-16 21:23:01 Annotating text fragment 50451/100392
## 2023-09-16 21:23:01 Annotating text fragment 50461/100392
## 2023-09-16 21:23:01 Annotating text fragment 50471/100392
## 2023-09-16 21:23:01 Annotating text fragment 50481/100392
## 2023-09-16 21:23:01 Annotating text fragment 50491/100392
## 2023-09-16 21:23:01 Annotating text fragment 50501/100392
## 2023-09-16 21:23:02 Annotating text fragment 50511/100392
## 2023-09-16 21:23:02 Annotating text fragment 50521/100392
## 2023-09-16 21:23:02 Annotating text fragment 50531/100392
## 2023-09-16 21:23:02 Annotating text fragment 50541/100392
## 2023-09-16 21:23:02 Annotating text fragment 50551/100392
## 2023-09-16 21:23:02 Annotating text fragment 50561/100392
## 2023-09-16 21:23:02 Annotating text fragment 50571/100392
## 2023-09-16 21:23:02 Annotating text fragment 50581/100392
## 2023-09-16 21:23:02 Annotating text fragment 50591/100392
## 2023-09-16 21:23:02 Annotating text fragment 50601/100392
## 2023-09-16 21:23:03 Annotating text fragment 50611/100392
## 2023-09-16 21:23:03 Annotating text fragment 50621/100392
## 2023-09-16 21:23:03 Annotating text fragment 50631/100392
## 2023-09-16 21:23:03 Annotating text fragment 50641/100392
## 2023-09-16 21:23:03 Annotating text fragment 50651/100392
## 2023-09-16 21:23:03 Annotating text fragment 50661/100392
## 2023-09-16 21:23:03 Annotating text fragment 50671/100392
## 2023-09-16 21:23:03 Annotating text fragment 50681/100392
## 2023-09-16 21:23:04 Annotating text fragment 50691/100392
## 2023-09-16 21:23:04 Annotating text fragment 50701/100392
## 2023-09-16 21:23:04 Annotating text fragment 50711/100392
## 2023-09-16 21:23:04 Annotating text fragment 50721/100392
## 2023-09-16 21:23:04 Annotating text fragment 50731/100392
## 2023-09-16 21:23:04 Annotating text fragment 50741/100392
## 2023-09-16 21:23:04 Annotating text fragment 50751/100392
## 2023-09-16 21:23:04 Annotating text fragment 50761/100392
## 2023-09-16 21:23:04 Annotating text fragment 50771/100392
## 2023-09-16 21:23:05 Annotating text fragment 50781/100392
## 2023-09-16 21:23:05 Annotating text fragment 50791/100392
## 2023-09-16 21:23:05 Annotating text fragment 50801/100392
## 2023-09-16 21:23:05 Annotating text fragment 50811/100392
```

```
## 2023-09-16 21:23:05 Annotating text fragment 50821/100392
## 2023-09-16 21:23:05 Annotating text fragment 50831/100392
## 2023-09-16 21:23:05 Annotating text fragment 50841/100392
## 2023-09-16 21:23:05 Annotating text fragment 50851/100392
## 2023-09-16 21:23:05 Annotating text fragment 50861/100392
## 2023-09-16 21:23:06 Annotating text fragment 50871/100392
## 2023-09-16 21:23:06 Annotating text fragment 50881/100392
## 2023-09-16 21:23:06 Annotating text fragment 50891/100392
## 2023-09-16 21:23:06 Annotating text fragment 50901/100392
## 2023-09-16 21:23:06 Annotating text fragment 50911/100392
## 2023-09-16 21:23:06 Annotating text fragment 50921/100392
## 2023-09-16 21:23:06 Annotating text fragment 50931/100392
## 2023-09-16 21:23:06 Annotating text fragment 50941/100392
## 2023-09-16 21:23:07 Annotating text fragment 50951/100392
## 2023-09-16 21:23:07 Annotating text fragment 50961/100392
## 2023-09-16 21:23:07 Annotating text fragment 50971/100392
## 2023-09-16 21:23:07 Annotating text fragment 50981/100392
## 2023-09-16 21:23:07 Annotating text fragment 50991/100392
## 2023-09-16 21:23:07 Annotating text fragment 51001/100392
## 2023-09-16 21:23:07 Annotating text fragment 51011/100392
## 2023-09-16 21:23:07 Annotating text fragment 51021/100392
## 2023-09-16 21:23:07 Annotating text fragment 51031/100392
## 2023-09-16 21:23:08 Annotating text fragment 51041/100392
## 2023-09-16 21:23:08 Annotating text fragment 51051/100392
## 2023-09-16 21:23:08 Annotating text fragment 51061/100392
## 2023-09-16 21:23:08 Annotating text fragment 51071/100392
## 2023-09-16 21:23:08 Annotating text fragment 51081/100392
## 2023-09-16 21:23:08 Annotating text fragment 51091/100392
## 2023-09-16 21:23:08 Annotating text fragment 51101/100392
## 2023-09-16 21:23:08 Annotating text fragment 51111/100392
## 2023-09-16 21:23:09 Annotating text fragment 51121/100392
## 2023-09-16 21:23:09 Annotating text fragment 51131/100392
## 2023-09-16 21:23:09 Annotating text fragment 51141/100392
## 2023-09-16 21:23:09 Annotating text fragment 51151/100392
## 2023-09-16 21:23:09 Annotating text fragment 51161/100392
## 2023-09-16 21:23:09 Annotating text fragment 51171/100392
## 2023-09-16 21:23:09 Annotating text fragment 51181/100392
## 2023-09-16 21:23:09 Annotating text fragment 51191/100392
## 2023-09-16 21:23:09 Annotating text fragment 51201/100392
## 2023-09-16 21:23:10 Annotating text fragment 51211/100392
## 2023-09-16 21:23:10 Annotating text fragment 51221/100392
## 2023-09-16 21:23:10 Annotating text fragment 51231/100392
## 2023-09-16 21:23:10 Annotating text fragment 51241/100392
## 2023-09-16 21:23:10 Annotating text fragment 51251/100392
## 2023-09-16 21:23:10 Annotating text fragment 51261/100392
## 2023-09-16 21:23:10 Annotating text fragment 51271/100392
## 2023-09-16 21:23:10 Annotating text fragment 51281/100392
## 2023-09-16 21:23:10 Annotating text fragment 51291/100392
## 2023-09-16 21:23:11 Annotating text fragment 51301/100392
## 2023-09-16 21:23:11 Annotating text fragment 51311/100392
## 2023-09-16 21:23:11 Annotating text fragment 51321/100392
## 2023-09-16 21:23:11 Annotating text fragment 51331/100392
## 2023-09-16 21:23:11 Annotating text fragment 51341/100392
## 2023-09-16 21:23:11 Annotating text fragment 51351/100392
```

```
## 2023-09-16 21:23:11 Annotating text fragment 51361/100392
## 2023-09-16 21:23:11 Annotating text fragment 51371/100392
## 2023-09-16 21:23:11 Annotating text fragment 51381/100392
## 2023-09-16 21:23:11 Annotating text fragment 51391/100392
## 2023-09-16 21:23:12 Annotating text fragment 51401/100392
## 2023-09-16 21:23:12 Annotating text fragment 51411/100392
## 2023-09-16 21:23:12 Annotating text fragment 51421/100392
## 2023-09-16 21:23:12 Annotating text fragment 51431/100392
## 2023-09-16 21:23:12 Annotating text fragment 51441/100392
## 2023-09-16 21:23:12 Annotating text fragment 51451/100392
## 2023-09-16 21:23:12 Annotating text fragment 51461/100392
## 2023-09-16 21:23:12 Annotating text fragment 51471/100392
## 2023-09-16 21:23:12 Annotating text fragment 51481/100392
## 2023-09-16 21:23:12 Annotating text fragment 51491/100392
## 2023-09-16 21:23:13 Annotating text fragment 51501/100392
## 2023-09-16 21:23:13 Annotating text fragment 51511/100392
## 2023-09-16 21:23:13 Annotating text fragment 51521/100392
## 2023-09-16 21:23:13 Annotating text fragment 51531/100392
## 2023-09-16 21:23:13 Annotating text fragment 51541/100392
## 2023-09-16 21:23:13 Annotating text fragment 51551/100392
## 2023-09-16 21:23:13 Annotating text fragment 51561/100392
## 2023-09-16 21:23:13 Annotating text fragment 51571/100392
## 2023-09-16 21:23:13 Annotating text fragment 51581/100392
## 2023-09-16 21:23:14 Annotating text fragment 51591/100392
## 2023-09-16 21:23:14 Annotating text fragment 51601/100392
## 2023-09-16 21:23:14 Annotating text fragment 51611/100392
## 2023-09-16 21:23:14 Annotating text fragment 51621/100392
## 2023-09-16 21:23:14 Annotating text fragment 51631/100392
## 2023-09-16 21:23:14 Annotating text fragment 51641/100392
## 2023-09-16 21:23:14 Annotating text fragment 51651/100392
## 2023-09-16 21:23:14 Annotating text fragment 51661/100392
## 2023-09-16 21:23:14 Annotating text fragment 51671/100392
## 2023-09-16 21:23:15 Annotating text fragment 51681/100392
## 2023-09-16 21:23:15 Annotating text fragment 51691/100392
## 2023-09-16 21:23:15 Annotating text fragment 51701/100392
## 2023-09-16 21:23:15 Annotating text fragment 51711/100392
## 2023-09-16 21:23:15 Annotating text fragment 51721/100392
## 2023-09-16 21:23:15 Annotating text fragment 51731/100392
## 2023-09-16 21:23:15 Annotating text fragment 51741/100392
## 2023-09-16 21:23:15 Annotating text fragment 51751/100392
## 2023-09-16 21:23:15 Annotating text fragment 51761/100392
## 2023-09-16 21:23:15 Annotating text fragment 51771/100392
## 2023-09-16 21:23:15 Annotating text fragment 51781/100392
## 2023-09-16 21:23:16 Annotating text fragment 51791/100392
## 2023-09-16 21:23:16 Annotating text fragment 51801/100392
## 2023-09-16 21:23:16 Annotating text fragment 51811/100392
## 2023-09-16 21:23:16 Annotating text fragment 51821/100392
## 2023-09-16 21:23:16 Annotating text fragment 51831/100392
## 2023-09-16 21:23:16 Annotating text fragment 51841/100392
## 2023-09-16 21:23:16 Annotating text fragment 51851/100392
## 2023-09-16 21:23:16 Annotating text fragment 51861/100392
## 2023-09-16 21:23:16 Annotating text fragment 51871/100392
## 2023-09-16 21:23:17 Annotating text fragment 51881/100392
## 2023-09-16 21:23:17 Annotating text fragment 51891/100392
```

```
## 2023-09-16 21:23:17 Annotating text fragment 51901/100392
## 2023-09-16 21:23:17 Annotating text fragment 51911/100392
## 2023-09-16 21:23:17 Annotating text fragment 51921/100392
## 2023-09-16 21:23:17 Annotating text fragment 51931/100392
## 2023-09-16 21:23:17 Annotating text fragment 51941/100392
## 2023-09-16 21:23:17 Annotating text fragment 51951/100392
## 2023-09-16 21:23:17 Annotating text fragment 51961/100392
## 2023-09-16 21:23:17 Annotating text fragment 51971/100392
## 2023-09-16 21:23:18 Annotating text fragment 51981/100392
## 2023-09-16 21:23:18 Annotating text fragment 51991/100392
## 2023-09-16 21:23:18 Annotating text fragment 52001/100392
## 2023-09-16 21:23:18 Annotating text fragment 52011/100392
## 2023-09-16 21:23:18 Annotating text fragment 52021/100392
## 2023-09-16 21:23:18 Annotating text fragment 52031/100392
## 2023-09-16 21:23:18 Annotating text fragment 52041/100392
## 2023-09-16 21:23:18 Annotating text fragment 52051/100392
## 2023-09-16 21:23:19 Annotating text fragment 52061/100392
## 2023-09-16 21:23:19 Annotating text fragment 52071/100392
## 2023-09-16 21:23:19 Annotating text fragment 52081/100392
## 2023-09-16 21:23:19 Annotating text fragment 52091/100392
## 2023-09-16 21:23:19 Annotating text fragment 52101/100392
## 2023-09-16 21:23:19 Annotating text fragment 52111/100392
## 2023-09-16 21:23:19 Annotating text fragment 52121/100392
## 2023-09-16 21:23:20 Annotating text fragment 52131/100392
## 2023-09-16 21:23:20 Annotating text fragment 52141/100392
## 2023-09-16 21:23:20 Annotating text fragment 52151/100392
## 2023-09-16 21:23:20 Annotating text fragment 52161/100392
## 2023-09-16 21:23:20 Annotating text fragment 52171/100392
## 2023-09-16 21:23:20 Annotating text fragment 52181/100392
## 2023-09-16 21:23:20 Annotating text fragment 52191/100392
## 2023-09-16 21:23:21 Annotating text fragment 52201/100392
## 2023-09-16 21:23:21 Annotating text fragment 52211/100392
## 2023-09-16 21:23:21 Annotating text fragment 52221/100392
## 2023-09-16 21:23:21 Annotating text fragment 52231/100392
## 2023-09-16 21:23:21 Annotating text fragment 52241/100392
## 2023-09-16 21:23:21 Annotating text fragment 52251/100392
## 2023-09-16 21:23:21 Annotating text fragment 52261/100392
## 2023-09-16 21:23:21 Annotating text fragment 52271/100392
## 2023-09-16 21:23:21 Annotating text fragment 52281/100392
## 2023-09-16 21:23:21 Annotating text fragment 52291/100392
## 2023-09-16 21:23:22 Annotating text fragment 52301/100392
## 2023-09-16 21:23:22 Annotating text fragment 52311/100392
## 2023-09-16 21:23:22 Annotating text fragment 52321/100392
## 2023-09-16 21:23:22 Annotating text fragment 52331/100392
## 2023-09-16 21:23:22 Annotating text fragment 52341/100392
## 2023-09-16 21:23:22 Annotating text fragment 52351/100392
## 2023-09-16 21:23:22 Annotating text fragment 52361/100392
## 2023-09-16 21:23:22 Annotating text fragment 52371/100392
## 2023-09-16 21:23:22 Annotating text fragment 52381/100392
## 2023-09-16 21:23:22 Annotating text fragment 52391/100392
## 2023-09-16 21:23:23 Annotating text fragment 52401/100392
## 2023-09-16 21:23:23 Annotating text fragment 52411/100392
## 2023-09-16 21:23:23 Annotating text fragment 52421/100392
## 2023-09-16 21:23:23 Annotating text fragment 52431/100392
```

```
## 2023-09-16 21:23:23 Annotating text fragment 52441/100392
## 2023-09-16 21:23:23 Annotating text fragment 52451/100392
## 2023-09-16 21:23:23 Annotating text fragment 52461/100392
## 2023-09-16 21:23:23 Annotating text fragment 52471/100392
## 2023-09-16 21:23:23 Annotating text fragment 52481/100392
## 2023-09-16 21:23:23 Annotating text fragment 52491/100392
## 2023-09-16 21:23:24 Annotating text fragment 52501/100392
## 2023-09-16 21:23:24 Annotating text fragment 52511/100392
## 2023-09-16 21:23:24 Annotating text fragment 52521/100392
## 2023-09-16 21:23:24 Annotating text fragment 52531/100392
## 2023-09-16 21:23:24 Annotating text fragment 52541/100392
## 2023-09-16 21:23:24 Annotating text fragment 52551/100392
## 2023-09-16 21:23:24 Annotating text fragment 52561/100392
## 2023-09-16 21:23:24 Annotating text fragment 52571/100392
## 2023-09-16 21:23:24 Annotating text fragment 52581/100392
## 2023-09-16 21:23:25 Annotating text fragment 52591/100392
## 2023-09-16 21:23:25 Annotating text fragment 52601/100392
## 2023-09-16 21:23:25 Annotating text fragment 52611/100392
## 2023-09-16 21:23:25 Annotating text fragment 52621/100392
## 2023-09-16 21:23:25 Annotating text fragment 52631/100392
## 2023-09-16 21:23:25 Annotating text fragment 52641/100392
## 2023-09-16 21:23:25 Annotating text fragment 52651/100392
## 2023-09-16 21:23:25 Annotating text fragment 52661/100392
## 2023-09-16 21:23:25 Annotating text fragment 52671/100392
## 2023-09-16 21:23:25 Annotating text fragment 52681/100392
## 2023-09-16 21:23:25 Annotating text fragment 52691/100392
## 2023-09-16 21:23:26 Annotating text fragment 52701/100392
## 2023-09-16 21:23:26 Annotating text fragment 52711/100392
## 2023-09-16 21:23:26 Annotating text fragment 52721/100392
## 2023-09-16 21:23:26 Annotating text fragment 52731/100392
## 2023-09-16 21:23:26 Annotating text fragment 52741/100392
## 2023-09-16 21:23:26 Annotating text fragment 52751/100392
## 2023-09-16 21:23:26 Annotating text fragment 52761/100392
## 2023-09-16 21:23:26 Annotating text fragment 52771/100392
## 2023-09-16 21:23:26 Annotating text fragment 52781/100392
## 2023-09-16 21:23:26 Annotating text fragment 52791/100392
## 2023-09-16 21:23:27 Annotating text fragment 52801/100392
## 2023-09-16 21:23:27 Annotating text fragment 52811/100392
## 2023-09-16 21:23:27 Annotating text fragment 52821/100392
## 2023-09-16 21:23:27 Annotating text fragment 52831/100392
## 2023-09-16 21:23:27 Annotating text fragment 52841/100392
## 2023-09-16 21:23:27 Annotating text fragment 52851/100392
## 2023-09-16 21:23:27 Annotating text fragment 52861/100392
## 2023-09-16 21:23:27 Annotating text fragment 52871/100392
## 2023-09-16 21:23:27 Annotating text fragment 52881/100392
## 2023-09-16 21:23:27 Annotating text fragment 52891/100392
## 2023-09-16 21:23:27 Annotating text fragment 52901/100392
## 2023-09-16 21:23:27 Annotating text fragment 52911/100392
## 2023-09-16 21:23:28 Annotating text fragment 52921/100392
## 2023-09-16 21:23:28 Annotating text fragment 52931/100392
## 2023-09-16 21:23:28 Annotating text fragment 52941/100392
## 2023-09-16 21:23:28 Annotating text fragment 52951/100392
## 2023-09-16 21:23:28 Annotating text fragment 52961/100392
## 2023-09-16 21:23:28 Annotating text fragment 52971/100392
```

```
## 2023-09-16 21:23:28 Annotating text fragment 52981/100392
## 2023-09-16 21:23:28 Annotating text fragment 52991/100392
## 2023-09-16 21:23:28 Annotating text fragment 53001/100392
## 2023-09-16 21:23:28 Annotating text fragment 53011/100392
## 2023-09-16 21:23:29 Annotating text fragment 53021/100392
## 2023-09-16 21:23:29 Annotating text fragment 53031/100392
## 2023-09-16 21:23:29 Annotating text fragment 53041/100392
## 2023-09-16 21:23:29 Annotating text fragment 53051/100392
## 2023-09-16 21:23:29 Annotating text fragment 53061/100392
## 2023-09-16 21:23:29 Annotating text fragment 53071/100392
## 2023-09-16 21:23:29 Annotating text fragment 53081/100392
## 2023-09-16 21:23:29 Annotating text fragment 53091/100392
## 2023-09-16 21:23:29 Annotating text fragment 53101/100392
## 2023-09-16 21:23:29 Annotating text fragment 53111/100392
## 2023-09-16 21:23:30 Annotating text fragment 53121/100392
## 2023-09-16 21:23:30 Annotating text fragment 53131/100392
## 2023-09-16 21:23:30 Annotating text fragment 53141/100392
## 2023-09-16 21:23:30 Annotating text fragment 53151/100392
## 2023-09-16 21:23:30 Annotating text fragment 53161/100392
## 2023-09-16 21:23:30 Annotating text fragment 53171/100392
## 2023-09-16 21:23:30 Annotating text fragment 53181/100392
## 2023-09-16 21:23:30 Annotating text fragment 53191/100392
## 2023-09-16 21:23:30 Annotating text fragment 53201/100392
## 2023-09-16 21:23:31 Annotating text fragment 53211/100392
## 2023-09-16 21:23:31 Annotating text fragment 53221/100392
## 2023-09-16 21:23:31 Annotating text fragment 53231/100392
## 2023-09-16 21:23:31 Annotating text fragment 53241/100392
## 2023-09-16 21:23:31 Annotating text fragment 53251/100392
## 2023-09-16 21:23:31 Annotating text fragment 53261/100392
## 2023-09-16 21:23:31 Annotating text fragment 53271/100392
## 2023-09-16 21:23:31 Annotating text fragment 53281/100392
## 2023-09-16 21:23:32 Annotating text fragment 53291/100392
## 2023-09-16 21:23:32 Annotating text fragment 53301/100392
## 2023-09-16 21:23:32 Annotating text fragment 53311/100392
## 2023-09-16 21:23:32 Annotating text fragment 53321/100392
## 2023-09-16 21:23:32 Annotating text fragment 53331/100392
## 2023-09-16 21:23:32 Annotating text fragment 53341/100392
## 2023-09-16 21:23:32 Annotating text fragment 53351/100392
## 2023-09-16 21:23:32 Annotating text fragment 53361/100392
## 2023-09-16 21:23:32 Annotating text fragment 53371/100392
## 2023-09-16 21:23:32 Annotating text fragment 53381/100392
## 2023-09-16 21:23:33 Annotating text fragment 53391/100392
## 2023-09-16 21:23:33 Annotating text fragment 53401/100392
## 2023-09-16 21:23:33 Annotating text fragment 53411/100392
## 2023-09-16 21:23:33 Annotating text fragment 53421/100392
## 2023-09-16 21:23:33 Annotating text fragment 53431/100392
## 2023-09-16 21:23:33 Annotating text fragment 53441/100392
## 2023-09-16 21:23:33 Annotating text fragment 53451/100392
## 2023-09-16 21:23:33 Annotating text fragment 53461/100392
## 2023-09-16 21:23:33 Annotating text fragment 53471/100392
## 2023-09-16 21:23:34 Annotating text fragment 53481/100392
## 2023-09-16 21:23:34 Annotating text fragment 53491/100392
## 2023-09-16 21:23:34 Annotating text fragment 53501/100392
## 2023-09-16 21:23:34 Annotating text fragment 53511/100392
```

```
## 2023-09-16 21:23:34 Annotating text fragment 53521/100392
## 2023-09-16 21:23:34 Annotating text fragment 53531/100392
## 2023-09-16 21:23:34 Annotating text fragment 53541/100392
## 2023-09-16 21:23:34 Annotating text fragment 53551/100392
## 2023-09-16 21:23:34 Annotating text fragment 53561/100392
## 2023-09-16 21:23:34 Annotating text fragment 53571/100392
## 2023-09-16 21:23:35 Annotating text fragment 53581/100392
## 2023-09-16 21:23:35 Annotating text fragment 53591/100392
## 2023-09-16 21:23:35 Annotating text fragment 53601/100392
## 2023-09-16 21:23:35 Annotating text fragment 53611/100392
## 2023-09-16 21:23:35 Annotating text fragment 53621/100392
## 2023-09-16 21:23:35 Annotating text fragment 53631/100392
## 2023-09-16 21:23:35 Annotating text fragment 53641/100392
## 2023-09-16 21:23:35 Annotating text fragment 53651/100392
## 2023-09-16 21:23:35 Annotating text fragment 53661/100392
## 2023-09-16 21:23:35 Annotating text fragment 53671/100392
## 2023-09-16 21:23:35 Annotating text fragment 53681/100392
## 2023-09-16 21:23:36 Annotating text fragment 53691/100392
## 2023-09-16 21:23:36 Annotating text fragment 53701/100392
## 2023-09-16 21:23:36 Annotating text fragment 53711/100392
## 2023-09-16 21:23:36 Annotating text fragment 53721/100392
## 2023-09-16 21:23:36 Annotating text fragment 53731/100392
## 2023-09-16 21:23:36 Annotating text fragment 53741/100392
## 2023-09-16 21:23:36 Annotating text fragment 53751/100392
## 2023-09-16 21:23:36 Annotating text fragment 53761/100392
## 2023-09-16 21:23:36 Annotating text fragment 53771/100392
## 2023-09-16 21:23:36 Annotating text fragment 53781/100392
## 2023-09-16 21:23:36 Annotating text fragment 53791/100392
## 2023-09-16 21:23:37 Annotating text fragment 53801/100392
## 2023-09-16 21:23:37 Annotating text fragment 53811/100392
## 2023-09-16 21:23:37 Annotating text fragment 53821/100392
## 2023-09-16 21:23:37 Annotating text fragment 53831/100392
## 2023-09-16 21:23:37 Annotating text fragment 53841/100392
## 2023-09-16 21:23:37 Annotating text fragment 53851/100392
## 2023-09-16 21:23:37 Annotating text fragment 53861/100392
## 2023-09-16 21:23:37 Annotating text fragment 53871/100392
## 2023-09-16 21:23:37 Annotating text fragment 53881/100392
## 2023-09-16 21:23:37 Annotating text fragment 53891/100392
## 2023-09-16 21:23:37 Annotating text fragment 53901/100392
## 2023-09-16 21:23:37 Annotating text fragment 53911/100392
## 2023-09-16 21:23:38 Annotating text fragment 53921/100392
## 2023-09-16 21:23:38 Annotating text fragment 53931/100392
## 2023-09-16 21:23:38 Annotating text fragment 53941/100392
## 2023-09-16 21:23:38 Annotating text fragment 53951/100392
## 2023-09-16 21:23:38 Annotating text fragment 53961/100392
## 2023-09-16 21:23:38 Annotating text fragment 53971/100392
## 2023-09-16 21:23:38 Annotating text fragment 53981/100392
## 2023-09-16 21:23:38 Annotating text fragment 53991/100392
## 2023-09-16 21:23:38 Annotating text fragment 54001/100392
## 2023-09-16 21:23:39 Annotating text fragment 54011/100392
## 2023-09-16 21:23:39 Annotating text fragment 54021/100392
## 2023-09-16 21:23:39 Annotating text fragment 54031/100392
## 2023-09-16 21:23:39 Annotating text fragment 54041/100392
## 2023-09-16 21:23:39 Annotating text fragment 54051/100392
```

```
## 2023-09-16 21:23:39 Annotating text fragment 54061/100392
## 2023-09-16 21:23:39 Annotating text fragment 54071/100392
## 2023-09-16 21:23:40 Annotating text fragment 54081/100392
## 2023-09-16 21:23:40 Annotating text fragment 54091/100392
## 2023-09-16 21:23:40 Annotating text fragment 54101/100392
## 2023-09-16 21:23:40 Annotating text fragment 54111/100392
## 2023-09-16 21:23:40 Annotating text fragment 54121/100392
## 2023-09-16 21:23:40 Annotating text fragment 54131/100392
## 2023-09-16 21:23:40 Annotating text fragment 54141/100392
## 2023-09-16 21:23:40 Annotating text fragment 54151/100392
## 2023-09-16 21:23:41 Annotating text fragment 54161/100392
## 2023-09-16 21:23:41 Annotating text fragment 54171/100392
## 2023-09-16 21:23:41 Annotating text fragment 54181/100392
## 2023-09-16 21:23:41 Annotating text fragment 54191/100392
## 2023-09-16 21:23:41 Annotating text fragment 54201/100392
## 2023-09-16 21:23:41 Annotating text fragment 54211/100392
## 2023-09-16 21:23:41 Annotating text fragment 54221/100392
## 2023-09-16 21:23:41 Annotating text fragment 54231/100392
## 2023-09-16 21:23:42 Annotating text fragment 54241/100392
## 2023-09-16 21:23:42 Annotating text fragment 54251/100392
## 2023-09-16 21:23:42 Annotating text fragment 54261/100392
## 2023-09-16 21:23:42 Annotating text fragment 54271/100392
## 2023-09-16 21:23:42 Annotating text fragment 54281/100392
## 2023-09-16 21:23:42 Annotating text fragment 54291/100392
## 2023-09-16 21:23:42 Annotating text fragment 54301/100392
## 2023-09-16 21:23:42 Annotating text fragment 54311/100392
## 2023-09-16 21:23:42 Annotating text fragment 54321/100392
## 2023-09-16 21:23:42 Annotating text fragment 54331/100392
## 2023-09-16 21:23:42 Annotating text fragment 54341/100392
## 2023-09-16 21:23:43 Annotating text fragment 54351/100392
## 2023-09-16 21:23:43 Annotating text fragment 54361/100392
## 2023-09-16 21:23:43 Annotating text fragment 54371/100392
## 2023-09-16 21:23:43 Annotating text fragment 54381/100392
## 2023-09-16 21:23:43 Annotating text fragment 54391/100392
## 2023-09-16 21:23:43 Annotating text fragment 54401/100392
## 2023-09-16 21:23:43 Annotating text fragment 54411/100392
## 2023-09-16 21:23:43 Annotating text fragment 54421/100392
## 2023-09-16 21:23:44 Annotating text fragment 54431/100392
## 2023-09-16 21:23:44 Annotating text fragment 54441/100392
## 2023-09-16 21:23:44 Annotating text fragment 54451/100392
## 2023-09-16 21:23:44 Annotating text fragment 54461/100392
## 2023-09-16 21:23:44 Annotating text fragment 54471/100392
## 2023-09-16 21:23:44 Annotating text fragment 54481/100392
## 2023-09-16 21:23:44 Annotating text fragment 54491/100392
## 2023-09-16 21:23:44 Annotating text fragment 54501/100392
## 2023-09-16 21:23:44 Annotating text fragment 54511/100392
## 2023-09-16 21:23:45 Annotating text fragment 54521/100392
## 2023-09-16 21:23:45 Annotating text fragment 54531/100392
## 2023-09-16 21:23:45 Annotating text fragment 54541/100392
## 2023-09-16 21:23:45 Annotating text fragment 54551/100392
## 2023-09-16 21:23:45 Annotating text fragment 54561/100392
## 2023-09-16 21:23:45 Annotating text fragment 54571/100392
## 2023-09-16 21:23:45 Annotating text fragment 54581/100392
## 2023-09-16 21:23:45 Annotating text fragment 54591/100392
```

```
## 2023-09-16 21:23:45 Annotating text fragment 54601/100392
## 2023-09-16 21:23:45 Annotating text fragment 54611/100392
## 2023-09-16 21:23:46 Annotating text fragment 54621/100392
## 2023-09-16 21:23:46 Annotating text fragment 54631/100392
## 2023-09-16 21:23:46 Annotating text fragment 54641/100392
## 2023-09-16 21:23:46 Annotating text fragment 54651/100392
## 2023-09-16 21:23:46 Annotating text fragment 54661/100392
## 2023-09-16 21:23:46 Annotating text fragment 54671/100392
## 2023-09-16 21:23:46 Annotating text fragment 54681/100392
## 2023-09-16 21:23:46 Annotating text fragment 54691/100392
## 2023-09-16 21:23:46 Annotating text fragment 54701/100392
## 2023-09-16 21:23:46 Annotating text fragment 54711/100392
## 2023-09-16 21:23:47 Annotating text fragment 54721/100392
## 2023-09-16 21:23:47 Annotating text fragment 54731/100392
## 2023-09-16 21:23:47 Annotating text fragment 54741/100392
## 2023-09-16 21:23:47 Annotating text fragment 54751/100392
## 2023-09-16 21:23:47 Annotating text fragment 54761/100392
## 2023-09-16 21:23:47 Annotating text fragment 54771/100392
## 2023-09-16 21:23:47 Annotating text fragment 54781/100392
## 2023-09-16 21:23:47 Annotating text fragment 54791/100392
## 2023-09-16 21:23:47 Annotating text fragment 54801/100392
## 2023-09-16 21:23:47 Annotating text fragment 54811/100392
## 2023-09-16 21:23:47 Annotating text fragment 54821/100392
## 2023-09-16 21:23:48 Annotating text fragment 54831/100392
## 2023-09-16 21:23:48 Annotating text fragment 54841/100392
## 2023-09-16 21:23:48 Annotating text fragment 54851/100392
## 2023-09-16 21:23:48 Annotating text fragment 54861/100392
## 2023-09-16 21:23:48 Annotating text fragment 54871/100392
## 2023-09-16 21:23:48 Annotating text fragment 54881/100392
## 2023-09-16 21:23:48 Annotating text fragment 54891/100392
## 2023-09-16 21:23:48 Annotating text fragment 54901/100392
## 2023-09-16 21:23:49 Annotating text fragment 54911/100392
## 2023-09-16 21:23:49 Annotating text fragment 54921/100392
## 2023-09-16 21:23:49 Annotating text fragment 54931/100392
## 2023-09-16 21:23:49 Annotating text fragment 54941/100392
## 2023-09-16 21:23:49 Annotating text fragment 54951/100392
## 2023-09-16 21:23:49 Annotating text fragment 54961/100392
## 2023-09-16 21:23:49 Annotating text fragment 54971/100392
## 2023-09-16 21:23:49 Annotating text fragment 54981/100392
## 2023-09-16 21:23:49 Annotating text fragment 54991/100392
## 2023-09-16 21:23:50 Annotating text fragment 55001/100392
## 2023-09-16 21:23:50 Annotating text fragment 55011/100392
## 2023-09-16 21:23:50 Annotating text fragment 55021/100392
## 2023-09-16 21:23:50 Annotating text fragment 55031/100392
## 2023-09-16 21:23:50 Annotating text fragment 55041/100392
## 2023-09-16 21:23:50 Annotating text fragment 55051/100392
## 2023-09-16 21:23:50 Annotating text fragment 55061/100392
## 2023-09-16 21:23:50 Annotating text fragment 55071/100392
## 2023-09-16 21:23:50 Annotating text fragment 55081/100392
## 2023-09-16 21:23:51 Annotating text fragment 55091/100392
## 2023-09-16 21:23:51 Annotating text fragment 55101/100392
## 2023-09-16 21:23:51 Annotating text fragment 55111/100392
## 2023-09-16 21:23:51 Annotating text fragment 55121/100392
## 2023-09-16 21:23:51 Annotating text fragment 55131/100392
```

```
## 2023-09-16 21:23:51 Annotating text fragment 55141/100392
## 2023-09-16 21:23:51 Annotating text fragment 55151/100392
## 2023-09-16 21:23:51 Annotating text fragment 55161/100392
## 2023-09-16 21:23:51 Annotating text fragment 55171/100392
## 2023-09-16 21:23:52 Annotating text fragment 55181/100392
## 2023-09-16 21:23:52 Annotating text fragment 55191/100392
## 2023-09-16 21:23:52 Annotating text fragment 55201/100392
## 2023-09-16 21:23:52 Annotating text fragment 55211/100392
## 2023-09-16 21:23:52 Annotating text fragment 55221/100392
## 2023-09-16 21:23:52 Annotating text fragment 55231/100392
## 2023-09-16 21:23:52 Annotating text fragment 55241/100392
## 2023-09-16 21:23:52 Annotating text fragment 55251/100392
## 2023-09-16 21:23:52 Annotating text fragment 55261/100392
## 2023-09-16 21:23:52 Annotating text fragment 55271/100392
## 2023-09-16 21:23:53 Annotating text fragment 55281/100392
## 2023-09-16 21:23:53 Annotating text fragment 55291/100392
## 2023-09-16 21:23:53 Annotating text fragment 55301/100392
## 2023-09-16 21:23:53 Annotating text fragment 55311/100392
## 2023-09-16 21:23:53 Annotating text fragment 55321/100392
## 2023-09-16 21:23:53 Annotating text fragment 55331/100392
## 2023-09-16 21:23:53 Annotating text fragment 55341/100392
## 2023-09-16 21:23:53 Annotating text fragment 55351/100392
## 2023-09-16 21:23:53 Annotating text fragment 55361/100392
## 2023-09-16 21:23:53 Annotating text fragment 55371/100392
## 2023-09-16 21:23:53 Annotating text fragment 55381/100392
## 2023-09-16 21:23:53 Annotating text fragment 55391/100392
## 2023-09-16 21:23:54 Annotating text fragment 55401/100392
## 2023-09-16 21:23:54 Annotating text fragment 55411/100392
## 2023-09-16 21:23:54 Annotating text fragment 55421/100392
## 2023-09-16 21:23:54 Annotating text fragment 55431/100392
## 2023-09-16 21:23:54 Annotating text fragment 55441/100392
## 2023-09-16 21:23:54 Annotating text fragment 55451/100392
## 2023-09-16 21:23:54 Annotating text fragment 55461/100392
## 2023-09-16 21:23:54 Annotating text fragment 55471/100392
## 2023-09-16 21:23:54 Annotating text fragment 55481/100392
## 2023-09-16 21:23:54 Annotating text fragment 55491/100392
## 2023-09-16 21:23:55 Annotating text fragment 55501/100392
## 2023-09-16 21:23:55 Annotating text fragment 55511/100392
## 2023-09-16 21:23:55 Annotating text fragment 55521/100392
## 2023-09-16 21:23:55 Annotating text fragment 55531/100392
## 2023-09-16 21:23:55 Annotating text fragment 55541/100392
## 2023-09-16 21:23:55 Annotating text fragment 55551/100392
## 2023-09-16 21:23:55 Annotating text fragment 55561/100392
## 2023-09-16 21:23:55 Annotating text fragment 55571/100392
## 2023-09-16 21:23:55 Annotating text fragment 55581/100392
## 2023-09-16 21:23:55 Annotating text fragment 55591/100392
## 2023-09-16 21:23:55 Annotating text fragment 55601/100392
## 2023-09-16 21:23:56 Annotating text fragment 55611/100392
## 2023-09-16 21:23:56 Annotating text fragment 55621/100392
## 2023-09-16 21:23:56 Annotating text fragment 55631/100392
## 2023-09-16 21:23:56 Annotating text fragment 55641/100392
## 2023-09-16 21:23:56 Annotating text fragment 55651/100392
## 2023-09-16 21:23:56 Annotating text fragment 55661/100392
## 2023-09-16 21:23:56 Annotating text fragment 55671/100392
```

```
## 2023-09-16 21:23:56 Annotating text fragment 55681/100392
## 2023-09-16 21:23:56 Annotating text fragment 55691/100392
## 2023-09-16 21:23:56 Annotating text fragment 55701/100392
## 2023-09-16 21:23:57 Annotating text fragment 55711/100392
## 2023-09-16 21:23:57 Annotating text fragment 55721/100392
## 2023-09-16 21:23:57 Annotating text fragment 55731/100392
## 2023-09-16 21:23:57 Annotating text fragment 55741/100392
## 2023-09-16 21:23:57 Annotating text fragment 55751/100392
## 2023-09-16 21:23:57 Annotating text fragment 55761/100392
## 2023-09-16 21:23:57 Annotating text fragment 55771/100392
## 2023-09-16 21:23:57 Annotating text fragment 55781/100392
## 2023-09-16 21:23:57 Annotating text fragment 55791/100392
## 2023-09-16 21:23:57 Annotating text fragment 55801/100392
## 2023-09-16 21:23:57 Annotating text fragment 55811/100392
## 2023-09-16 21:23:58 Annotating text fragment 55821/100392
## 2023-09-16 21:23:58 Annotating text fragment 55831/100392
## 2023-09-16 21:23:58 Annotating text fragment 55841/100392
## 2023-09-16 21:23:58 Annotating text fragment 55851/100392
## 2023-09-16 21:23:58 Annotating text fragment 55861/100392
## 2023-09-16 21:23:58 Annotating text fragment 55871/100392
## 2023-09-16 21:23:58 Annotating text fragment 55881/100392
## 2023-09-16 21:23:58 Annotating text fragment 55891/100392
## 2023-09-16 21:23:58 Annotating text fragment 55901/100392
## 2023-09-16 21:23:59 Annotating text fragment 55911/100392
## 2023-09-16 21:23:59 Annotating text fragment 55921/100392
## 2023-09-16 21:23:59 Annotating text fragment 55931/100392
## 2023-09-16 21:23:59 Annotating text fragment 55941/100392
## 2023-09-16 21:23:59 Annotating text fragment 55951/100392
## 2023-09-16 21:23:59 Annotating text fragment 55961/100392
## 2023-09-16 21:23:59 Annotating text fragment 55971/100392
## 2023-09-16 21:23:59 Annotating text fragment 55981/100392
## 2023-09-16 21:23:59 Annotating text fragment 55991/100392
## 2023-09-16 21:24:00 Annotating text fragment 56001/100392
## 2023-09-16 21:24:00 Annotating text fragment 56011/100392
## 2023-09-16 21:24:00 Annotating text fragment 56021/100392
## 2023-09-16 21:24:00 Annotating text fragment 56031/100392
## 2023-09-16 21:24:00 Annotating text fragment 56041/100392
## 2023-09-16 21:24:00 Annotating text fragment 56051/100392
## 2023-09-16 21:24:00 Annotating text fragment 56061/100392
## 2023-09-16 21:24:00 Annotating text fragment 56071/100392
## 2023-09-16 21:24:00 Annotating text fragment 56081/100392
## 2023-09-16 21:24:00 Annotating text fragment 56091/100392
## 2023-09-16 21:24:01 Annotating text fragment 56101/100392
## 2023-09-16 21:24:01 Annotating text fragment 56111/100392
## 2023-09-16 21:24:01 Annotating text fragment 56121/100392
## 2023-09-16 21:24:01 Annotating text fragment 56131/100392
## 2023-09-16 21:24:01 Annotating text fragment 56141/100392
## 2023-09-16 21:24:01 Annotating text fragment 56151/100392
## 2023-09-16 21:24:01 Annotating text fragment 56161/100392
## 2023-09-16 21:24:01 Annotating text fragment 56171/100392
## 2023-09-16 21:24:01 Annotating text fragment 56181/100392
## 2023-09-16 21:24:02 Annotating text fragment 56191/100392
## 2023-09-16 21:24:02 Annotating text fragment 56201/100392
## 2023-09-16 21:24:02 Annotating text fragment 56211/100392
```

```
## 2023-09-16 21:24:02 Annotating text fragment 56221/100392
## 2023-09-16 21:24:02 Annotating text fragment 56231/100392
## 2023-09-16 21:24:02 Annotating text fragment 56241/100392
## 2023-09-16 21:24:02 Annotating text fragment 56251/100392
## 2023-09-16 21:24:02 Annotating text fragment 56261/100392
## 2023-09-16 21:24:02 Annotating text fragment 56271/100392
## 2023-09-16 21:24:02 Annotating text fragment 56281/100392
## 2023-09-16 21:24:02 Annotating text fragment 56291/100392
## 2023-09-16 21:24:02 Annotating text fragment 56301/100392
## 2023-09-16 21:24:03 Annotating text fragment 56311/100392
## 2023-09-16 21:24:03 Annotating text fragment 56321/100392
## 2023-09-16 21:24:03 Annotating text fragment 56331/100392
## 2023-09-16 21:24:03 Annotating text fragment 56341/100392
## 2023-09-16 21:24:03 Annotating text fragment 56351/100392
## 2023-09-16 21:24:03 Annotating text fragment 56361/100392
## 2023-09-16 21:24:03 Annotating text fragment 56371/100392
## 2023-09-16 21:24:03 Annotating text fragment 56381/100392
## 2023-09-16 21:24:03 Annotating text fragment 56391/100392
## 2023-09-16 21:24:04 Annotating text fragment 56401/100392
## 2023-09-16 21:24:04 Annotating text fragment 56411/100392
## 2023-09-16 21:24:04 Annotating text fragment 56421/100392
## 2023-09-16 21:24:04 Annotating text fragment 56431/100392
## 2023-09-16 21:24:04 Annotating text fragment 56441/100392
## 2023-09-16 21:24:04 Annotating text fragment 56451/100392
## 2023-09-16 21:24:04 Annotating text fragment 56461/100392
## 2023-09-16 21:24:04 Annotating text fragment 56471/100392
## 2023-09-16 21:24:04 Annotating text fragment 56481/100392
## 2023-09-16 21:24:04 Annotating text fragment 56491/100392
## 2023-09-16 21:24:04 Annotating text fragment 56501/100392
## 2023-09-16 21:24:04 Annotating text fragment 56511/100392
## 2023-09-16 21:24:05 Annotating text fragment 56521/100392
## 2023-09-16 21:24:05 Annotating text fragment 56531/100392
## 2023-09-16 21:24:05 Annotating text fragment 56541/100392
## 2023-09-16 21:24:05 Annotating text fragment 56551/100392
## 2023-09-16 21:24:05 Annotating text fragment 56561/100392
## 2023-09-16 21:24:05 Annotating text fragment 56571/100392
## 2023-09-16 21:24:05 Annotating text fragment 56581/100392
## 2023-09-16 21:24:05 Annotating text fragment 56591/100392
## 2023-09-16 21:24:05 Annotating text fragment 56601/100392
## 2023-09-16 21:24:05 Annotating text fragment 56611/100392
## 2023-09-16 21:24:06 Annotating text fragment 56621/100392
## 2023-09-16 21:24:06 Annotating text fragment 56631/100392
## 2023-09-16 21:24:06 Annotating text fragment 56641/100392
## 2023-09-16 21:24:06 Annotating text fragment 56651/100392
## 2023-09-16 21:24:06 Annotating text fragment 56661/100392
## 2023-09-16 21:24:06 Annotating text fragment 56671/100392
## 2023-09-16 21:24:06 Annotating text fragment 56681/100392
## 2023-09-16 21:24:06 Annotating text fragment 56691/100392
## 2023-09-16 21:24:06 Annotating text fragment 56701/100392
## 2023-09-16 21:24:06 Annotating text fragment 56711/100392
## 2023-09-16 21:24:06 Annotating text fragment 56721/100392
## 2023-09-16 21:24:07 Annotating text fragment 56731/100392
## 2023-09-16 21:24:07 Annotating text fragment 56741/100392
## 2023-09-16 21:24:07 Annotating text fragment 56751/100392
```

```
## 2023-09-16 21:24:07 Annotating text fragment 56761/100392
## 2023-09-16 21:24:07 Annotating text fragment 56771/100392
## 2023-09-16 21:24:07 Annotating text fragment 56781/100392
## 2023-09-16 21:24:07 Annotating text fragment 56791/100392
## 2023-09-16 21:24:07 Annotating text fragment 56801/100392
## 2023-09-16 21:24:07 Annotating text fragment 56811/100392
## 2023-09-16 21:24:07 Annotating text fragment 56821/100392
## 2023-09-16 21:24:08 Annotating text fragment 56831/100392
## 2023-09-16 21:24:08 Annotating text fragment 56841/100392
## 2023-09-16 21:24:08 Annotating text fragment 56851/100392
## 2023-09-16 21:24:08 Annotating text fragment 56861/100392
## 2023-09-16 21:24:08 Annotating text fragment 56871/100392
## 2023-09-16 21:24:08 Annotating text fragment 56881/100392
## 2023-09-16 21:24:08 Annotating text fragment 56891/100392
## 2023-09-16 21:24:08 Annotating text fragment 56901/100392
## 2023-09-16 21:24:08 Annotating text fragment 56911/100392
## 2023-09-16 21:24:08 Annotating text fragment 56921/100392
## 2023-09-16 21:24:08 Annotating text fragment 56931/100392
## 2023-09-16 21:24:09 Annotating text fragment 56941/100392
## 2023-09-16 21:24:09 Annotating text fragment 56951/100392
## 2023-09-16 21:24:09 Annotating text fragment 56961/100392
## 2023-09-16 21:24:09 Annotating text fragment 56971/100392
## 2023-09-16 21:24:09 Annotating text fragment 56981/100392
## 2023-09-16 21:24:09 Annotating text fragment 56991/100392
## 2023-09-16 21:24:09 Annotating text fragment 57001/100392
## 2023-09-16 21:24:09 Annotating text fragment 57011/100392
## 2023-09-16 21:24:10 Annotating text fragment 57021/100392
## 2023-09-16 21:24:10 Annotating text fragment 57031/100392
## 2023-09-16 21:24:10 Annotating text fragment 57041/100392
## 2023-09-16 21:24:10 Annotating text fragment 57051/100392
## 2023-09-16 21:24:10 Annotating text fragment 57061/100392
## 2023-09-16 21:24:10 Annotating text fragment 57071/100392
## 2023-09-16 21:24:10 Annotating text fragment 57081/100392
## 2023-09-16 21:24:10 Annotating text fragment 57091/100392
## 2023-09-16 21:24:10 Annotating text fragment 57101/100392
## 2023-09-16 21:24:10 Annotating text fragment 57111/100392
## 2023-09-16 21:24:11 Annotating text fragment 57121/100392
## 2023-09-16 21:24:11 Annotating text fragment 57131/100392
## 2023-09-16 21:24:11 Annotating text fragment 57141/100392
## 2023-09-16 21:24:11 Annotating text fragment 57151/100392
## 2023-09-16 21:24:11 Annotating text fragment 57161/100392
## 2023-09-16 21:24:11 Annotating text fragment 57171/100392
## 2023-09-16 21:24:11 Annotating text fragment 57181/100392
## 2023-09-16 21:24:11 Annotating text fragment 57191/100392
## 2023-09-16 21:24:11 Annotating text fragment 57201/100392
## 2023-09-16 21:24:11 Annotating text fragment 57211/100392
## 2023-09-16 21:24:11 Annotating text fragment 57221/100392
## 2023-09-16 21:24:11 Annotating text fragment 57231/100392
## 2023-09-16 21:24:12 Annotating text fragment 57241/100392
## 2023-09-16 21:24:12 Annotating text fragment 57251/100392
## 2023-09-16 21:24:12 Annotating text fragment 57261/100392
## 2023-09-16 21:24:12 Annotating text fragment 57271/100392
## 2023-09-16 21:24:12 Annotating text fragment 57281/100392
## 2023-09-16 21:24:12 Annotating text fragment 57291/100392
```

```
## 2023-09-16 21:24:12 Annotating text fragment 57301/100392
## 2023-09-16 21:24:12 Annotating text fragment 57311/100392
## 2023-09-16 21:24:12 Annotating text fragment 57321/100392
## 2023-09-16 21:24:12 Annotating text fragment 57331/100392
## 2023-09-16 21:24:13 Annotating text fragment 57341/100392
## 2023-09-16 21:24:13 Annotating text fragment 57351/100392
## 2023-09-16 21:24:13 Annotating text fragment 57361/100392
## 2023-09-16 21:24:13 Annotating text fragment 57371/100392
## 2023-09-16 21:24:13 Annotating text fragment 57381/100392
## 2023-09-16 21:24:13 Annotating text fragment 57391/100392
## 2023-09-16 21:24:13 Annotating text fragment 57401/100392
## 2023-09-16 21:24:14 Annotating text fragment 57411/100392
## 2023-09-16 21:24:14 Annotating text fragment 57421/100392
## 2023-09-16 21:24:14 Annotating text fragment 57431/100392
## 2023-09-16 21:24:14 Annotating text fragment 57441/100392
## 2023-09-16 21:24:14 Annotating text fragment 57451/100392
## 2023-09-16 21:24:14 Annotating text fragment 57461/100392
## 2023-09-16 21:24:14 Annotating text fragment 57471/100392
## 2023-09-16 21:24:14 Annotating text fragment 57481/100392
## 2023-09-16 21:24:15 Annotating text fragment 57491/100392
## 2023-09-16 21:24:15 Annotating text fragment 57501/100392
## 2023-09-16 21:24:15 Annotating text fragment 57511/100392
## 2023-09-16 21:24:15 Annotating text fragment 57521/100392
## 2023-09-16 21:24:15 Annotating text fragment 57531/100392
## 2023-09-16 21:24:16 Annotating text fragment 57541/100392
## 2023-09-16 21:24:16 Annotating text fragment 57551/100392
## 2023-09-16 21:24:16 Annotating text fragment 57561/100392
## 2023-09-16 21:24:16 Annotating text fragment 57571/100392
## 2023-09-16 21:24:16 Annotating text fragment 57581/100392
## 2023-09-16 21:24:16 Annotating text fragment 57591/100392
## 2023-09-16 21:24:16 Annotating text fragment 57601/100392
## 2023-09-16 21:24:16 Annotating text fragment 57611/100392
## 2023-09-16 21:24:16 Annotating text fragment 57621/100392
## 2023-09-16 21:24:17 Annotating text fragment 57631/100392
## 2023-09-16 21:24:17 Annotating text fragment 57641/100392
## 2023-09-16 21:24:17 Annotating text fragment 57651/100392
## 2023-09-16 21:24:17 Annotating text fragment 57661/100392
## 2023-09-16 21:24:17 Annotating text fragment 57671/100392
## 2023-09-16 21:24:17 Annotating text fragment 57681/100392
## 2023-09-16 21:24:17 Annotating text fragment 57691/100392
## 2023-09-16 21:24:18 Annotating text fragment 57701/100392
## 2023-09-16 21:24:18 Annotating text fragment 57711/100392
## 2023-09-16 21:24:18 Annotating text fragment 57721/100392
## 2023-09-16 21:24:18 Annotating text fragment 57731/100392
## 2023-09-16 21:24:18 Annotating text fragment 57741/100392
## 2023-09-16 21:24:18 Annotating text fragment 57751/100392
## 2023-09-16 21:24:18 Annotating text fragment 57761/100392
## 2023-09-16 21:24:18 Annotating text fragment 57771/100392
## 2023-09-16 21:24:18 Annotating text fragment 57781/100392
## 2023-09-16 21:24:18 Annotating text fragment 57791/100392
## 2023-09-16 21:24:19 Annotating text fragment 57801/100392
## 2023-09-16 21:24:19 Annotating text fragment 57811/100392
## 2023-09-16 21:24:19 Annotating text fragment 57821/100392
## 2023-09-16 21:24:19 Annotating text fragment 57831/100392
```

```
## 2023-09-16 21:24:19 Annotating text fragment 57841/100392
## 2023-09-16 21:24:19 Annotating text fragment 57851/100392
## 2023-09-16 21:24:19 Annotating text fragment 57861/100392
## 2023-09-16 21:24:19 Annotating text fragment 57871/100392
## 2023-09-16 21:24:19 Annotating text fragment 57881/100392
## 2023-09-16 21:24:20 Annotating text fragment 57891/100392
## 2023-09-16 21:24:20 Annotating text fragment 57901/100392
## 2023-09-16 21:24:20 Annotating text fragment 57911/100392
## 2023-09-16 21:24:20 Annotating text fragment 57921/100392
## 2023-09-16 21:24:20 Annotating text fragment 57931/100392
## 2023-09-16 21:24:20 Annotating text fragment 57941/100392
## 2023-09-16 21:24:20 Annotating text fragment 57951/100392
## 2023-09-16 21:24:20 Annotating text fragment 57961/100392
## 2023-09-16 21:24:20 Annotating text fragment 57971/100392
## 2023-09-16 21:24:21 Annotating text fragment 57981/100392
## 2023-09-16 21:24:21 Annotating text fragment 57991/100392
## 2023-09-16 21:24:21 Annotating text fragment 58001/100392
## 2023-09-16 21:24:21 Annotating text fragment 58011/100392
## 2023-09-16 21:24:21 Annotating text fragment 58021/100392
## 2023-09-16 21:24:21 Annotating text fragment 58031/100392
## 2023-09-16 21:24:21 Annotating text fragment 58041/100392
## 2023-09-16 21:24:21 Annotating text fragment 58051/100392
## 2023-09-16 21:24:21 Annotating text fragment 58061/100392
## 2023-09-16 21:24:22 Annotating text fragment 58071/100392
## 2023-09-16 21:24:22 Annotating text fragment 58081/100392
## 2023-09-16 21:24:22 Annotating text fragment 58091/100392
## 2023-09-16 21:24:22 Annotating text fragment 58101/100392
## 2023-09-16 21:24:22 Annotating text fragment 58111/100392
## 2023-09-16 21:24:22 Annotating text fragment 58121/100392
## 2023-09-16 21:24:22 Annotating text fragment 58131/100392
## 2023-09-16 21:24:22 Annotating text fragment 58141/100392
## 2023-09-16 21:24:22 Annotating text fragment 58151/100392
## 2023-09-16 21:24:23 Annotating text fragment 58161/100392
## 2023-09-16 21:24:23 Annotating text fragment 58171/100392
## 2023-09-16 21:24:23 Annotating text fragment 58181/100392
## 2023-09-16 21:24:23 Annotating text fragment 58191/100392
## 2023-09-16 21:24:23 Annotating text fragment 58201/100392
## 2023-09-16 21:24:23 Annotating text fragment 58211/100392
## 2023-09-16 21:24:24 Annotating text fragment 58221/100392
## 2023-09-16 21:24:24 Annotating text fragment 58231/100392
## 2023-09-16 21:24:24 Annotating text fragment 58241/100392
## 2023-09-16 21:24:24 Annotating text fragment 58251/100392
## 2023-09-16 21:24:24 Annotating text fragment 58261/100392
## 2023-09-16 21:24:24 Annotating text fragment 58271/100392
## 2023-09-16 21:24:24 Annotating text fragment 58281/100392
## 2023-09-16 21:24:24 Annotating text fragment 58291/100392
## 2023-09-16 21:24:24 Annotating text fragment 58301/100392
## 2023-09-16 21:24:25 Annotating text fragment 58311/100392
## 2023-09-16 21:24:25 Annotating text fragment 58321/100392
## 2023-09-16 21:24:25 Annotating text fragment 58331/100392
## 2023-09-16 21:24:25 Annotating text fragment 58341/100392
## 2023-09-16 21:24:25 Annotating text fragment 58351/100392
## 2023-09-16 21:24:25 Annotating text fragment 58361/100392
## 2023-09-16 21:24:25 Annotating text fragment 58371/100392
```

```
## 2023-09-16 21:24:26 Annotating text fragment 58381/100392
## 2023-09-16 21:24:26 Annotating text fragment 58391/100392
## 2023-09-16 21:24:26 Annotating text fragment 58401/100392
## 2023-09-16 21:24:26 Annotating text fragment 58411/100392
## 2023-09-16 21:24:26 Annotating text fragment 58421/100392
## 2023-09-16 21:24:26 Annotating text fragment 58431/100392
## 2023-09-16 21:24:26 Annotating text fragment 58441/100392
## 2023-09-16 21:24:26 Annotating text fragment 58451/100392
## 2023-09-16 21:24:27 Annotating text fragment 58461/100392
## 2023-09-16 21:24:27 Annotating text fragment 58471/100392
## 2023-09-16 21:24:27 Annotating text fragment 58481/100392
## 2023-09-16 21:24:27 Annotating text fragment 58491/100392
## 2023-09-16 21:24:27 Annotating text fragment 58501/100392
## 2023-09-16 21:24:27 Annotating text fragment 58511/100392
## 2023-09-16 21:24:27 Annotating text fragment 58521/100392
## 2023-09-16 21:24:27 Annotating text fragment 58531/100392
## 2023-09-16 21:24:28 Annotating text fragment 58541/100392
## 2023-09-16 21:24:28 Annotating text fragment 58551/100392
## 2023-09-16 21:24:28 Annotating text fragment 58561/100392
## 2023-09-16 21:24:28 Annotating text fragment 58571/100392
## 2023-09-16 21:24:28 Annotating text fragment 58581/100392
## 2023-09-16 21:24:28 Annotating text fragment 58591/100392
## 2023-09-16 21:24:28 Annotating text fragment 58601/100392
## 2023-09-16 21:24:28 Annotating text fragment 58611/100392
## 2023-09-16 21:24:28 Annotating text fragment 58621/100392
## 2023-09-16 21:24:29 Annotating text fragment 58631/100392
## 2023-09-16 21:24:29 Annotating text fragment 58641/100392
## 2023-09-16 21:24:29 Annotating text fragment 58651/100392
## 2023-09-16 21:24:29 Annotating text fragment 58661/100392
## 2023-09-16 21:24:29 Annotating text fragment 58671/100392
## 2023-09-16 21:24:29 Annotating text fragment 58681/100392
## 2023-09-16 21:24:29 Annotating text fragment 58691/100392
## 2023-09-16 21:24:29 Annotating text fragment 58701/100392
## 2023-09-16 21:24:29 Annotating text fragment 58711/100392
## 2023-09-16 21:24:30 Annotating text fragment 58721/100392
## 2023-09-16 21:24:30 Annotating text fragment 58731/100392
## 2023-09-16 21:24:30 Annotating text fragment 58741/100392
## 2023-09-16 21:24:30 Annotating text fragment 58751/100392
## 2023-09-16 21:24:30 Annotating text fragment 58761/100392
## 2023-09-16 21:24:30 Annotating text fragment 58771/100392
## 2023-09-16 21:24:30 Annotating text fragment 58781/100392
## 2023-09-16 21:24:30 Annotating text fragment 58791/100392
## 2023-09-16 21:24:30 Annotating text fragment 58801/100392
## 2023-09-16 21:24:30 Annotating text fragment 58811/100392
## 2023-09-16 21:24:31 Annotating text fragment 58821/100392
## 2023-09-16 21:24:31 Annotating text fragment 58831/100392
## 2023-09-16 21:24:31 Annotating text fragment 58841/100392
## 2023-09-16 21:24:31 Annotating text fragment 58851/100392
## 2023-09-16 21:24:31 Annotating text fragment 58861/100392
## 2023-09-16 21:24:31 Annotating text fragment 58871/100392
## 2023-09-16 21:24:31 Annotating text fragment 58881/100392
## 2023-09-16 21:24:31 Annotating text fragment 58891/100392
## 2023-09-16 21:24:31 Annotating text fragment 58901/100392
## 2023-09-16 21:24:31 Annotating text fragment 58911/100392
```

```
## 2023-09-16 21:24:31 Annotating text fragment 58921/100392
## 2023-09-16 21:24:32 Annotating text fragment 58931/100392
## 2023-09-16 21:24:32 Annotating text fragment 58941/100392
## 2023-09-16 21:24:32 Annotating text fragment 58951/100392
## 2023-09-16 21:24:32 Annotating text fragment 58961/100392
## 2023-09-16 21:24:32 Annotating text fragment 58971/100392
## 2023-09-16 21:24:32 Annotating text fragment 58981/100392
## 2023-09-16 21:24:32 Annotating text fragment 58991/100392
## 2023-09-16 21:24:32 Annotating text fragment 59001/100392
## 2023-09-16 21:24:32 Annotating text fragment 59011/100392
## 2023-09-16 21:24:32 Annotating text fragment 59021/100392
## 2023-09-16 21:24:33 Annotating text fragment 59031/100392
## 2023-09-16 21:24:33 Annotating text fragment 59041/100392
## 2023-09-16 21:24:33 Annotating text fragment 59051/100392
## 2023-09-16 21:24:33 Annotating text fragment 59061/100392
## 2023-09-16 21:24:33 Annotating text fragment 59071/100392
## 2023-09-16 21:24:33 Annotating text fragment 59081/100392
## 2023-09-16 21:24:33 Annotating text fragment 59091/100392
## 2023-09-16 21:24:33 Annotating text fragment 59101/100392
## 2023-09-16 21:24:33 Annotating text fragment 59111/100392
## 2023-09-16 21:24:33 Annotating text fragment 59121/100392
## 2023-09-16 21:24:33 Annotating text fragment 59131/100392
## 2023-09-16 21:24:34 Annotating text fragment 59141/100392
## 2023-09-16 21:24:34 Annotating text fragment 59151/100392
## 2023-09-16 21:24:34 Annotating text fragment 59161/100392
## 2023-09-16 21:24:34 Annotating text fragment 59171/100392
## 2023-09-16 21:24:34 Annotating text fragment 59181/100392
## 2023-09-16 21:24:34 Annotating text fragment 59191/100392
## 2023-09-16 21:24:34 Annotating text fragment 59201/100392
## 2023-09-16 21:24:34 Annotating text fragment 59211/100392
## 2023-09-16 21:24:34 Annotating text fragment 59221/100392
## 2023-09-16 21:24:34 Annotating text fragment 59231/100392
## 2023-09-16 21:24:34 Annotating text fragment 59241/100392
## 2023-09-16 21:24:35 Annotating text fragment 59251/100392
## 2023-09-16 21:24:35 Annotating text fragment 59261/100392
## 2023-09-16 21:24:35 Annotating text fragment 59271/100392
## 2023-09-16 21:24:35 Annotating text fragment 59281/100392
## 2023-09-16 21:24:35 Annotating text fragment 59291/100392
## 2023-09-16 21:24:35 Annotating text fragment 59301/100392
## 2023-09-16 21:24:35 Annotating text fragment 59311/100392
## 2023-09-16 21:24:35 Annotating text fragment 59321/100392
## 2023-09-16 21:24:35 Annotating text fragment 59331/100392
## 2023-09-16 21:24:35 Annotating text fragment 59341/100392
## 2023-09-16 21:24:36 Annotating text fragment 59351/100392
## 2023-09-16 21:24:36 Annotating text fragment 59361/100392
## 2023-09-16 21:24:36 Annotating text fragment 59371/100392
## 2023-09-16 21:24:36 Annotating text fragment 59381/100392
## 2023-09-16 21:24:36 Annotating text fragment 59391/100392
## 2023-09-16 21:24:36 Annotating text fragment 59401/100392
## 2023-09-16 21:24:36 Annotating text fragment 59411/100392
## 2023-09-16 21:24:37 Annotating text fragment 59421/100392
## 2023-09-16 21:24:37 Annotating text fragment 59431/100392
## 2023-09-16 21:24:37 Annotating text fragment 59441/100392
## 2023-09-16 21:24:37 Annotating text fragment 59451/100392
```

```
## 2023-09-16 21:24:37 Annotating text fragment 59461/100392
## 2023-09-16 21:24:37 Annotating text fragment 59471/100392
## 2023-09-16 21:24:37 Annotating text fragment 59481/100392
## 2023-09-16 21:24:37 Annotating text fragment 59491/100392
## 2023-09-16 21:24:37 Annotating text fragment 59501/100392
## 2023-09-16 21:24:37 Annotating text fragment 59511/100392
## 2023-09-16 21:24:38 Annotating text fragment 59521/100392
## 2023-09-16 21:24:38 Annotating text fragment 59531/100392
## 2023-09-16 21:24:38 Annotating text fragment 59541/100392
## 2023-09-16 21:24:38 Annotating text fragment 59551/100392
## 2023-09-16 21:24:38 Annotating text fragment 59561/100392
## 2023-09-16 21:24:38 Annotating text fragment 59571/100392
## 2023-09-16 21:24:38 Annotating text fragment 59581/100392
## 2023-09-16 21:24:38 Annotating text fragment 59591/100392
## 2023-09-16 21:24:39 Annotating text fragment 59601/100392
## 2023-09-16 21:24:39 Annotating text fragment 59611/100392
## 2023-09-16 21:24:39 Annotating text fragment 59621/100392
## 2023-09-16 21:24:39 Annotating text fragment 59631/100392
## 2023-09-16 21:24:39 Annotating text fragment 59641/100392
## 2023-09-16 21:24:39 Annotating text fragment 59651/100392
## 2023-09-16 21:24:39 Annotating text fragment 59661/100392
## 2023-09-16 21:24:39 Annotating text fragment 59671/100392
## 2023-09-16 21:24:39 Annotating text fragment 59681/100392
## 2023-09-16 21:24:40 Annotating text fragment 59691/100392
## 2023-09-16 21:24:40 Annotating text fragment 59701/100392
## 2023-09-16 21:24:40 Annotating text fragment 59711/100392
## 2023-09-16 21:24:40 Annotating text fragment 59721/100392
## 2023-09-16 21:24:40 Annotating text fragment 59731/100392
## 2023-09-16 21:24:40 Annotating text fragment 59741/100392
## 2023-09-16 21:24:40 Annotating text fragment 59751/100392
## 2023-09-16 21:24:40 Annotating text fragment 59761/100392
## 2023-09-16 21:24:40 Annotating text fragment 59771/100392
## 2023-09-16 21:24:40 Annotating text fragment 59781/100392
## 2023-09-16 21:24:41 Annotating text fragment 59791/100392
## 2023-09-16 21:24:41 Annotating text fragment 59801/100392
## 2023-09-16 21:24:41 Annotating text fragment 59811/100392
## 2023-09-16 21:24:41 Annotating text fragment 59821/100392
## 2023-09-16 21:24:41 Annotating text fragment 59831/100392
## 2023-09-16 21:24:41 Annotating text fragment 59841/100392
## 2023-09-16 21:24:41 Annotating text fragment 59851/100392
## 2023-09-16 21:24:41 Annotating text fragment 59861/100392
## 2023-09-16 21:24:41 Annotating text fragment 59871/100392
## 2023-09-16 21:24:41 Annotating text fragment 59881/100392
## 2023-09-16 21:24:42 Annotating text fragment 59891/100392
## 2023-09-16 21:24:42 Annotating text fragment 59901/100392
## 2023-09-16 21:24:42 Annotating text fragment 59911/100392
## 2023-09-16 21:24:42 Annotating text fragment 59921/100392
## 2023-09-16 21:24:42 Annotating text fragment 59931/100392
## 2023-09-16 21:24:42 Annotating text fragment 59941/100392
## 2023-09-16 21:24:42 Annotating text fragment 59951/100392
## 2023-09-16 21:24:42 Annotating text fragment 59961/100392
## 2023-09-16 21:24:42 Annotating text fragment 59971/100392
## 2023-09-16 21:24:42 Annotating text fragment 59981/100392
## 2023-09-16 21:24:43 Annotating text fragment 59991/100392
```

```
## 2023-09-16 21:24:43 Annotating text fragment 60001/100392
## 2023-09-16 21:24:43 Annotating text fragment 60011/100392
## 2023-09-16 21:24:43 Annotating text fragment 60021/100392
## 2023-09-16 21:24:43 Annotating text fragment 60031/100392
## 2023-09-16 21:24:43 Annotating text fragment 60041/100392
## 2023-09-16 21:24:43 Annotating text fragment 60051/100392
## 2023-09-16 21:24:43 Annotating text fragment 60061/100392
## 2023-09-16 21:24:43 Annotating text fragment 60071/100392
## 2023-09-16 21:24:43 Annotating text fragment 60081/100392
## 2023-09-16 21:24:44 Annotating text fragment 60091/100392
## 2023-09-16 21:24:44 Annotating text fragment 60101/100392
## 2023-09-16 21:24:44 Annotating text fragment 60111/100392
## 2023-09-16 21:24:44 Annotating text fragment 60121/100392
## 2023-09-16 21:24:44 Annotating text fragment 60131/100392
## 2023-09-16 21:24:44 Annotating text fragment 60141/100392
## 2023-09-16 21:24:44 Annotating text fragment 60151/100392
## 2023-09-16 21:24:44 Annotating text fragment 60161/100392
## 2023-09-16 21:24:44 Annotating text fragment 60171/100392
## 2023-09-16 21:24:44 Annotating text fragment 60181/100392
## 2023-09-16 21:24:44 Annotating text fragment 60191/100392
## 2023-09-16 21:24:44 Annotating text fragment 60201/100392
## 2023-09-16 21:24:45 Annotating text fragment 60211/100392
## 2023-09-16 21:24:45 Annotating text fragment 60221/100392
## 2023-09-16 21:24:45 Annotating text fragment 60231/100392
## 2023-09-16 21:24:45 Annotating text fragment 60241/100392
## 2023-09-16 21:24:45 Annotating text fragment 60251/100392
## 2023-09-16 21:24:45 Annotating text fragment 60261/100392
## 2023-09-16 21:24:45 Annotating text fragment 60271/100392
## 2023-09-16 21:24:45 Annotating text fragment 60281/100392
## 2023-09-16 21:24:45 Annotating text fragment 60291/100392
## 2023-09-16 21:24:45 Annotating text fragment 60301/100392
## 2023-09-16 21:24:46 Annotating text fragment 60311/100392
## 2023-09-16 21:24:46 Annotating text fragment 60321/100392
## 2023-09-16 21:24:46 Annotating text fragment 60331/100392
## 2023-09-16 21:24:46 Annotating text fragment 60341/100392
## 2023-09-16 21:24:46 Annotating text fragment 60351/100392
## 2023-09-16 21:24:46 Annotating text fragment 60361/100392
## 2023-09-16 21:24:46 Annotating text fragment 60371/100392
## 2023-09-16 21:24:47 Annotating text fragment 60381/100392
## 2023-09-16 21:24:47 Annotating text fragment 60391/100392
## 2023-09-16 21:24:47 Annotating text fragment 60401/100392
## 2023-09-16 21:24:47 Annotating text fragment 60411/100392
## 2023-09-16 21:24:47 Annotating text fragment 60421/100392
## 2023-09-16 21:24:47 Annotating text fragment 60431/100392
## 2023-09-16 21:24:47 Annotating text fragment 60441/100392
## 2023-09-16 21:24:47 Annotating text fragment 60451/100392
## 2023-09-16 21:24:47 Annotating text fragment 60461/100392
## 2023-09-16 21:24:48 Annotating text fragment 60471/100392
## 2023-09-16 21:24:48 Annotating text fragment 60481/100392
## 2023-09-16 21:24:48 Annotating text fragment 60491/100392
## 2023-09-16 21:24:48 Annotating text fragment 60501/100392
## 2023-09-16 21:24:48 Annotating text fragment 60511/100392
## 2023-09-16 21:24:48 Annotating text fragment 60521/100392
## 2023-09-16 21:24:48 Annotating text fragment 60531/100392
```

```
## 2023-09-16 21:24:48 Annotating text fragment 60541/100392
## 2023-09-16 21:24:48 Annotating text fragment 60551/100392
## 2023-09-16 21:24:48 Annotating text fragment 60561/100392
## 2023-09-16 21:24:48 Annotating text fragment 60571/100392
## 2023-09-16 21:24:49 Annotating text fragment 60581/100392
## 2023-09-16 21:24:49 Annotating text fragment 60591/100392
## 2023-09-16 21:24:49 Annotating text fragment 60601/100392
## 2023-09-16 21:24:49 Annotating text fragment 60611/100392
## 2023-09-16 21:24:49 Annotating text fragment 60621/100392
## 2023-09-16 21:24:49 Annotating text fragment 60631/100392
## 2023-09-16 21:24:49 Annotating text fragment 60641/100392
## 2023-09-16 21:24:49 Annotating text fragment 60651/100392
## 2023-09-16 21:24:49 Annotating text fragment 60661/100392
## 2023-09-16 21:24:49 Annotating text fragment 60671/100392
## 2023-09-16 21:24:50 Annotating text fragment 60681/100392
## 2023-09-16 21:24:50 Annotating text fragment 60691/100392
## 2023-09-16 21:24:50 Annotating text fragment 60701/100392
## 2023-09-16 21:24:50 Annotating text fragment 60711/100392
## 2023-09-16 21:24:51 Annotating text fragment 60721/100392
## 2023-09-16 21:24:51 Annotating text fragment 60731/100392
## 2023-09-16 21:24:51 Annotating text fragment 60741/100392
## 2023-09-16 21:24:51 Annotating text fragment 60751/100392
## 2023-09-16 21:24:51 Annotating text fragment 60761/100392
## 2023-09-16 21:24:51 Annotating text fragment 60771/100392
## 2023-09-16 21:24:51 Annotating text fragment 60781/100392
## 2023-09-16 21:24:52 Annotating text fragment 60791/100392
## 2023-09-16 21:24:52 Annotating text fragment 60801/100392
## 2023-09-16 21:24:52 Annotating text fragment 60811/100392
## 2023-09-16 21:24:52 Annotating text fragment 60821/100392
## 2023-09-16 21:24:52 Annotating text fragment 60831/100392
## 2023-09-16 21:24:52 Annotating text fragment 60841/100392
## 2023-09-16 21:24:52 Annotating text fragment 60851/100392
## 2023-09-16 21:24:52 Annotating text fragment 60861/100392
## 2023-09-16 21:24:52 Annotating text fragment 60871/100392
## 2023-09-16 21:24:52 Annotating text fragment 60881/100392
## 2023-09-16 21:24:52 Annotating text fragment 60891/100392
## 2023-09-16 21:24:53 Annotating text fragment 60901/100392
## 2023-09-16 21:24:53 Annotating text fragment 60911/100392
## 2023-09-16 21:24:53 Annotating text fragment 60921/100392
## 2023-09-16 21:24:53 Annotating text fragment 60931/100392
## 2023-09-16 21:24:53 Annotating text fragment 60941/100392
## 2023-09-16 21:24:53 Annotating text fragment 60951/100392
## 2023-09-16 21:24:53 Annotating text fragment 60961/100392
## 2023-09-16 21:24:53 Annotating text fragment 60971/100392
## 2023-09-16 21:24:53 Annotating text fragment 60981/100392
## 2023-09-16 21:24:53 Annotating text fragment 60991/100392
## 2023-09-16 21:24:54 Annotating text fragment 61001/100392
## 2023-09-16 21:24:54 Annotating text fragment 61011/100392
## 2023-09-16 21:24:54 Annotating text fragment 61021/100392
## 2023-09-16 21:24:54 Annotating text fragment 61031/100392
## 2023-09-16 21:24:54 Annotating text fragment 61041/100392
## 2023-09-16 21:24:54 Annotating text fragment 61051/100392
## 2023-09-16 21:24:54 Annotating text fragment 61061/100392
## 2023-09-16 21:24:54 Annotating text fragment 61071/100392
```

```
## 2023-09-16 21:24:54 Annotating text fragment 61081/100392
## 2023-09-16 21:24:54 Annotating text fragment 61091/100392
## 2023-09-16 21:24:55 Annotating text fragment 61101/100392
## 2023-09-16 21:24:55 Annotating text fragment 61111/100392
## 2023-09-16 21:24:55 Annotating text fragment 61121/100392
## 2023-09-16 21:24:55 Annotating text fragment 61131/100392
## 2023-09-16 21:24:55 Annotating text fragment 61141/100392
## 2023-09-16 21:24:55 Annotating text fragment 61151/100392
## 2023-09-16 21:24:56 Annotating text fragment 61161/100392
## 2023-09-16 21:24:56 Annotating text fragment 61171/100392
## 2023-09-16 21:24:56 Annotating text fragment 61181/100392
## 2023-09-16 21:24:56 Annotating text fragment 61191/100392
## 2023-09-16 21:24:56 Annotating text fragment 61201/100392
## 2023-09-16 21:24:56 Annotating text fragment 61211/100392
## 2023-09-16 21:24:56 Annotating text fragment 61221/100392
## 2023-09-16 21:24:56 Annotating text fragment 61231/100392
## 2023-09-16 21:24:56 Annotating text fragment 61241/100392
## 2023-09-16 21:24:56 Annotating text fragment 61251/100392
## 2023-09-16 21:24:56 Annotating text fragment 61261/100392
## 2023-09-16 21:24:57 Annotating text fragment 61271/100392
## 2023-09-16 21:24:57 Annotating text fragment 61281/100392
## 2023-09-16 21:24:57 Annotating text fragment 61291/100392
## 2023-09-16 21:24:57 Annotating text fragment 61301/100392
## 2023-09-16 21:24:57 Annotating text fragment 61311/100392
## 2023-09-16 21:24:57 Annotating text fragment 61321/100392
## 2023-09-16 21:24:57 Annotating text fragment 61331/100392
## 2023-09-16 21:24:57 Annotating text fragment 61341/100392
## 2023-09-16 21:24:57 Annotating text fragment 61351/100392
## 2023-09-16 21:24:57 Annotating text fragment 61361/100392
## 2023-09-16 21:24:57 Annotating text fragment 61371/100392
## 2023-09-16 21:24:57 Annotating text fragment 61381/100392
## 2023-09-16 21:24:58 Annotating text fragment 61391/100392
## 2023-09-16 21:24:58 Annotating text fragment 61401/100392
## 2023-09-16 21:24:58 Annotating text fragment 61411/100392
## 2023-09-16 21:24:58 Annotating text fragment 61421/100392
## 2023-09-16 21:24:58 Annotating text fragment 61431/100392
## 2023-09-16 21:24:58 Annotating text fragment 61441/100392
## 2023-09-16 21:24:58 Annotating text fragment 61451/100392
## 2023-09-16 21:24:58 Annotating text fragment 61461/100392
## 2023-09-16 21:24:58 Annotating text fragment 61471/100392
## 2023-09-16 21:24:58 Annotating text fragment 61481/100392
## 2023-09-16 21:24:59 Annotating text fragment 61491/100392
## 2023-09-16 21:24:59 Annotating text fragment 61501/100392
## 2023-09-16 21:24:59 Annotating text fragment 61511/100392
## 2023-09-16 21:24:59 Annotating text fragment 61521/100392
## 2023-09-16 21:24:59 Annotating text fragment 61531/100392
## 2023-09-16 21:24:59 Annotating text fragment 61541/100392
## 2023-09-16 21:24:59 Annotating text fragment 61551/100392
## 2023-09-16 21:24:59 Annotating text fragment 61561/100392
## 2023-09-16 21:24:59 Annotating text fragment 61571/100392
## 2023-09-16 21:24:59 Annotating text fragment 61581/100392
## 2023-09-16 21:24:59 Annotating text fragment 61591/100392
## 2023-09-16 21:24:59 Annotating text fragment 61601/100392
## 2023-09-16 21:25:00 Annotating text fragment 61611/100392
```

```
## 2023-09-16 21:25:00 Annotating text fragment 61621/100392
## 2023-09-16 21:25:00 Annotating text fragment 61631/100392
## 2023-09-16 21:25:00 Annotating text fragment 61641/100392
## 2023-09-16 21:25:00 Annotating text fragment 61651/100392
## 2023-09-16 21:25:00 Annotating text fragment 61661/100392
## 2023-09-16 21:25:00 Annotating text fragment 61671/100392
## 2023-09-16 21:25:00 Annotating text fragment 61681/100392
## 2023-09-16 21:25:00 Annotating text fragment 61691/100392
## 2023-09-16 21:25:00 Annotating text fragment 61701/100392
## 2023-09-16 21:25:01 Annotating text fragment 61711/100392
## 2023-09-16 21:25:01 Annotating text fragment 61721/100392
## 2023-09-16 21:25:01 Annotating text fragment 61731/100392
## 2023-09-16 21:25:01 Annotating text fragment 61741/100392
## 2023-09-16 21:25:01 Annotating text fragment 61751/100392
## 2023-09-16 21:25:01 Annotating text fragment 61761/100392
## 2023-09-16 21:25:01 Annotating text fragment 61771/100392
## 2023-09-16 21:25:01 Annotating text fragment 61781/100392
## 2023-09-16 21:25:02 Annotating text fragment 61791/100392
## 2023-09-16 21:25:02 Annotating text fragment 61801/100392
## 2023-09-16 21:25:02 Annotating text fragment 61811/100392
## 2023-09-16 21:25:02 Annotating text fragment 61821/100392
## 2023-09-16 21:25:02 Annotating text fragment 61831/100392
## 2023-09-16 21:25:02 Annotating text fragment 61841/100392
## 2023-09-16 21:25:02 Annotating text fragment 61851/100392
## 2023-09-16 21:25:02 Annotating text fragment 61861/100392
## 2023-09-16 21:25:03 Annotating text fragment 61871/100392
## 2023-09-16 21:25:03 Annotating text fragment 61881/100392
## 2023-09-16 21:25:03 Annotating text fragment 61891/100392
## 2023-09-16 21:25:03 Annotating text fragment 61901/100392
## 2023-09-16 21:25:03 Annotating text fragment 61911/100392
## 2023-09-16 21:25:03 Annotating text fragment 61921/100392
## 2023-09-16 21:25:03 Annotating text fragment 61931/100392
## 2023-09-16 21:25:03 Annotating text fragment 61941/100392
## 2023-09-16 21:25:03 Annotating text fragment 61951/100392
## 2023-09-16 21:25:03 Annotating text fragment 61961/100392
## 2023-09-16 21:25:03 Annotating text fragment 61971/100392
## 2023-09-16 21:25:03 Annotating text fragment 61981/100392
## 2023-09-16 21:25:04 Annotating text fragment 61991/100392
## 2023-09-16 21:25:04 Annotating text fragment 62001/100392
## 2023-09-16 21:25:04 Annotating text fragment 62011/100392
## 2023-09-16 21:25:04 Annotating text fragment 62021/100392
## 2023-09-16 21:25:04 Annotating text fragment 62031/100392
## 2023-09-16 21:25:04 Annotating text fragment 62041/100392
## 2023-09-16 21:25:04 Annotating text fragment 62051/100392
## 2023-09-16 21:25:04 Annotating text fragment 62061/100392
## 2023-09-16 21:25:04 Annotating text fragment 62071/100392
## 2023-09-16 21:25:04 Annotating text fragment 62081/100392
## 2023-09-16 21:25:04 Annotating text fragment 62091/100392
## 2023-09-16 21:25:05 Annotating text fragment 62101/100392
## 2023-09-16 21:25:05 Annotating text fragment 62111/100392
## 2023-09-16 21:25:05 Annotating text fragment 62121/100392
## 2023-09-16 21:25:05 Annotating text fragment 62131/100392
## 2023-09-16 21:25:05 Annotating text fragment 62141/100392
## 2023-09-16 21:25:05 Annotating text fragment 62151/100392
```

```
## 2023-09-16 21:25:05 Annotating text fragment 62161/100392
## 2023-09-16 21:25:05 Annotating text fragment 62171/100392
## 2023-09-16 21:25:05 Annotating text fragment 62181/100392
## 2023-09-16 21:25:05 Annotating text fragment 62191/100392
## 2023-09-16 21:25:06 Annotating text fragment 62201/100392
## 2023-09-16 21:25:06 Annotating text fragment 62211/100392
## 2023-09-16 21:25:06 Annotating text fragment 62221/100392
## 2023-09-16 21:25:06 Annotating text fragment 62231/100392
## 2023-09-16 21:25:06 Annotating text fragment 62241/100392
## 2023-09-16 21:25:06 Annotating text fragment 62251/100392
## 2023-09-16 21:25:06 Annotating text fragment 62261/100392
## 2023-09-16 21:25:06 Annotating text fragment 62271/100392
## 2023-09-16 21:25:06 Annotating text fragment 62281/100392
## 2023-09-16 21:25:06 Annotating text fragment 62291/100392
## 2023-09-16 21:25:07 Annotating text fragment 62301/100392
## 2023-09-16 21:25:07 Annotating text fragment 62311/100392
## 2023-09-16 21:25:07 Annotating text fragment 62321/100392
## 2023-09-16 21:25:07 Annotating text fragment 62331/100392
## 2023-09-16 21:25:07 Annotating text fragment 62341/100392
## 2023-09-16 21:25:07 Annotating text fragment 62351/100392
## 2023-09-16 21:25:07 Annotating text fragment 62361/100392
## 2023-09-16 21:25:07 Annotating text fragment 62371/100392
## 2023-09-16 21:25:07 Annotating text fragment 62381/100392
## 2023-09-16 21:25:07 Annotating text fragment 62391/100392
## 2023-09-16 21:25:07 Annotating text fragment 62401/100392
## 2023-09-16 21:25:08 Annotating text fragment 62411/100392
## 2023-09-16 21:25:08 Annotating text fragment 62421/100392
## 2023-09-16 21:25:08 Annotating text fragment 62431/100392
## 2023-09-16 21:25:08 Annotating text fragment 62441/100392
## 2023-09-16 21:25:08 Annotating text fragment 62451/100392
## 2023-09-16 21:25:08 Annotating text fragment 62461/100392
## 2023-09-16 21:25:08 Annotating text fragment 62471/100392
## 2023-09-16 21:25:08 Annotating text fragment 62481/100392
## 2023-09-16 21:25:08 Annotating text fragment 62491/100392
## 2023-09-16 21:25:08 Annotating text fragment 62501/100392
## 2023-09-16 21:25:08 Annotating text fragment 62511/100392
## 2023-09-16 21:25:09 Annotating text fragment 62521/100392
## 2023-09-16 21:25:09 Annotating text fragment 62531/100392
## 2023-09-16 21:25:09 Annotating text fragment 62541/100392
## 2023-09-16 21:25:09 Annotating text fragment 62551/100392
## 2023-09-16 21:25:09 Annotating text fragment 62561/100392
## 2023-09-16 21:25:09 Annotating text fragment 62571/100392
## 2023-09-16 21:25:09 Annotating text fragment 62581/100392
## 2023-09-16 21:25:09 Annotating text fragment 62591/100392
## 2023-09-16 21:25:09 Annotating text fragment 62601/100392
## 2023-09-16 21:25:09 Annotating text fragment 62611/100392
## 2023-09-16 21:25:09 Annotating text fragment 62621/100392
## 2023-09-16 21:25:10 Annotating text fragment 62631/100392
## 2023-09-16 21:25:10 Annotating text fragment 62641/100392
## 2023-09-16 21:25:10 Annotating text fragment 62651/100392
## 2023-09-16 21:25:10 Annotating text fragment 62661/100392
## 2023-09-16 21:25:10 Annotating text fragment 62671/100392
## 2023-09-16 21:25:10 Annotating text fragment 62681/100392
## 2023-09-16 21:25:10 Annotating text fragment 62691/100392
```

```
## 2023-09-16 21:25:10 Annotating text fragment 62701/100392
## 2023-09-16 21:25:11 Annotating text fragment 62711/100392
## 2023-09-16 21:25:11 Annotating text fragment 62721/100392
## 2023-09-16 21:25:11 Annotating text fragment 62731/100392
## 2023-09-16 21:25:11 Annotating text fragment 62741/100392
## 2023-09-16 21:25:11 Annotating text fragment 62751/100392
## 2023-09-16 21:25:11 Annotating text fragment 62761/100392
## 2023-09-16 21:25:11 Annotating text fragment 62771/100392
## 2023-09-16 21:25:11 Annotating text fragment 62781/100392
## 2023-09-16 21:25:11 Annotating text fragment 62791/100392
## 2023-09-16 21:25:12 Annotating text fragment 62801/100392
## 2023-09-16 21:25:12 Annotating text fragment 62811/100392
## 2023-09-16 21:25:12 Annotating text fragment 62821/100392
## 2023-09-16 21:25:12 Annotating text fragment 62831/100392
## 2023-09-16 21:25:12 Annotating text fragment 62841/100392
## 2023-09-16 21:25:12 Annotating text fragment 62851/100392
## 2023-09-16 21:25:12 Annotating text fragment 62861/100392
## 2023-09-16 21:25:12 Annotating text fragment 62871/100392
## 2023-09-16 21:25:12 Annotating text fragment 62881/100392
## 2023-09-16 21:25:12 Annotating text fragment 62891/100392
## 2023-09-16 21:25:13 Annotating text fragment 62901/100392
## 2023-09-16 21:25:13 Annotating text fragment 62911/100392
## 2023-09-16 21:25:13 Annotating text fragment 62921/100392
## 2023-09-16 21:25:13 Annotating text fragment 62931/100392
## 2023-09-16 21:25:13 Annotating text fragment 62941/100392
## 2023-09-16 21:25:13 Annotating text fragment 62951/100392
## 2023-09-16 21:25:13 Annotating text fragment 62961/100392
## 2023-09-16 21:25:13 Annotating text fragment 62971/100392
## 2023-09-16 21:25:13 Annotating text fragment 62981/100392
## 2023-09-16 21:25:13 Annotating text fragment 62991/100392
## 2023-09-16 21:25:13 Annotating text fragment 63001/100392
## 2023-09-16 21:25:14 Annotating text fragment 63011/100392
## 2023-09-16 21:25:14 Annotating text fragment 63021/100392
## 2023-09-16 21:25:14 Annotating text fragment 63031/100392
## 2023-09-16 21:25:14 Annotating text fragment 63041/100392
## 2023-09-16 21:25:14 Annotating text fragment 63051/100392
## 2023-09-16 21:25:14 Annotating text fragment 63061/100392
## 2023-09-16 21:25:14 Annotating text fragment 63071/100392
## 2023-09-16 21:25:14 Annotating text fragment 63081/100392
## 2023-09-16 21:25:14 Annotating text fragment 63091/100392
## 2023-09-16 21:25:14 Annotating text fragment 63101/100392
## 2023-09-16 21:25:14 Annotating text fragment 63111/100392
## 2023-09-16 21:25:15 Annotating text fragment 63121/100392
## 2023-09-16 21:25:15 Annotating text fragment 63131/100392
## 2023-09-16 21:25:15 Annotating text fragment 63141/100392
## 2023-09-16 21:25:15 Annotating text fragment 63151/100392
## 2023-09-16 21:25:15 Annotating text fragment 63161/100392
## 2023-09-16 21:25:15 Annotating text fragment 63171/100392
## 2023-09-16 21:25:15 Annotating text fragment 63181/100392
## 2023-09-16 21:25:15 Annotating text fragment 63191/100392
## 2023-09-16 21:25:15 Annotating text fragment 63201/100392
## 2023-09-16 21:25:15 Annotating text fragment 63211/100392
## 2023-09-16 21:25:15 Annotating text fragment 63221/100392
## 2023-09-16 21:25:15 Annotating text fragment 63231/100392
```

```
## 2023-09-16 21:25:15 Annotating text fragment 63241/100392
## 2023-09-16 21:25:16 Annotating text fragment 63251/100392
## 2023-09-16 21:25:16 Annotating text fragment 63261/100392
## 2023-09-16 21:25:16 Annotating text fragment 63271/100392
## 2023-09-16 21:25:16 Annotating text fragment 63281/100392
## 2023-09-16 21:25:16 Annotating text fragment 63291/100392
## 2023-09-16 21:25:16 Annotating text fragment 63301/100392
## 2023-09-16 21:25:16 Annotating text fragment 63311/100392
## 2023-09-16 21:25:16 Annotating text fragment 63321/100392
## 2023-09-16 21:25:16 Annotating text fragment 63331/100392
## 2023-09-16 21:25:17 Annotating text fragment 63341/100392
## 2023-09-16 21:25:17 Annotating text fragment 63351/100392
## 2023-09-16 21:25:17 Annotating text fragment 63361/100392
## 2023-09-16 21:25:17 Annotating text fragment 63371/100392
## 2023-09-16 21:25:17 Annotating text fragment 63381/100392
## 2023-09-16 21:25:17 Annotating text fragment 63391/100392
## 2023-09-16 21:25:17 Annotating text fragment 63401/100392
## 2023-09-16 21:25:17 Annotating text fragment 63411/100392
## 2023-09-16 21:25:17 Annotating text fragment 63421/100392
## 2023-09-16 21:25:18 Annotating text fragment 63431/100392
## 2023-09-16 21:25:18 Annotating text fragment 63441/100392
## 2023-09-16 21:25:18 Annotating text fragment 63451/100392
## 2023-09-16 21:25:18 Annotating text fragment 63461/100392
## 2023-09-16 21:25:18 Annotating text fragment 63471/100392
## 2023-09-16 21:25:18 Annotating text fragment 63481/100392
## 2023-09-16 21:25:18 Annotating text fragment 63491/100392
## 2023-09-16 21:25:18 Annotating text fragment 63501/100392
## 2023-09-16 21:25:18 Annotating text fragment 63511/100392
## 2023-09-16 21:25:18 Annotating text fragment 63521/100392
## 2023-09-16 21:25:19 Annotating text fragment 63531/100392
## 2023-09-16 21:25:19 Annotating text fragment 63541/100392
## 2023-09-16 21:25:19 Annotating text fragment 63551/100392
## 2023-09-16 21:25:19 Annotating text fragment 63561/100392
## 2023-09-16 21:25:19 Annotating text fragment 63571/100392
## 2023-09-16 21:25:19 Annotating text fragment 63581/100392
## 2023-09-16 21:25:19 Annotating text fragment 63591/100392
## 2023-09-16 21:25:19 Annotating text fragment 63601/100392
## 2023-09-16 21:25:19 Annotating text fragment 63611/100392
## 2023-09-16 21:25:19 Annotating text fragment 63621/100392
## 2023-09-16 21:25:20 Annotating text fragment 63631/100392
## 2023-09-16 21:25:20 Annotating text fragment 63641/100392
## 2023-09-16 21:25:20 Annotating text fragment 63651/100392
## 2023-09-16 21:25:20 Annotating text fragment 63661/100392
## 2023-09-16 21:25:20 Annotating text fragment 63671/100392
## 2023-09-16 21:25:20 Annotating text fragment 63681/100392
## 2023-09-16 21:25:20 Annotating text fragment 63691/100392
## 2023-09-16 21:25:20 Annotating text fragment 63701/100392
## 2023-09-16 21:25:20 Annotating text fragment 63711/100392
## 2023-09-16 21:25:20 Annotating text fragment 63721/100392
## 2023-09-16 21:25:21 Annotating text fragment 63731/100392
## 2023-09-16 21:25:21 Annotating text fragment 63741/100392
## 2023-09-16 21:25:21 Annotating text fragment 63751/100392
## 2023-09-16 21:25:21 Annotating text fragment 63761/100392
## 2023-09-16 21:25:21 Annotating text fragment 63771/100392
```

```
## 2023-09-16 21:25:21 Annotating text fragment 63781/100392
## 2023-09-16 21:25:21 Annotating text fragment 63791/100392
## 2023-09-16 21:25:22 Annotating text fragment 63801/100392
## 2023-09-16 21:25:22 Annotating text fragment 63811/100392
## 2023-09-16 21:25:22 Annotating text fragment 63821/100392
## 2023-09-16 21:25:22 Annotating text fragment 63831/100392
## 2023-09-16 21:25:22 Annotating text fragment 63841/100392
## 2023-09-16 21:25:22 Annotating text fragment 63851/100392
## 2023-09-16 21:25:22 Annotating text fragment 63861/100392
## 2023-09-16 21:25:22 Annotating text fragment 63871/100392
## 2023-09-16 21:25:23 Annotating text fragment 63881/100392
## 2023-09-16 21:25:23 Annotating text fragment 63891/100392
## 2023-09-16 21:25:23 Annotating text fragment 63901/100392
## 2023-09-16 21:25:23 Annotating text fragment 63911/100392
## 2023-09-16 21:25:23 Annotating text fragment 63921/100392
## 2023-09-16 21:25:23 Annotating text fragment 63931/100392
## 2023-09-16 21:25:23 Annotating text fragment 63941/100392
## 2023-09-16 21:25:23 Annotating text fragment 63951/100392
## 2023-09-16 21:25:23 Annotating text fragment 63961/100392
## 2023-09-16 21:25:23 Annotating text fragment 63971/100392
## 2023-09-16 21:25:24 Annotating text fragment 63981/100392
## 2023-09-16 21:25:24 Annotating text fragment 63991/100392
## 2023-09-16 21:25:24 Annotating text fragment 64001/100392
## 2023-09-16 21:25:24 Annotating text fragment 64011/100392
## 2023-09-16 21:25:24 Annotating text fragment 64021/100392
## 2023-09-16 21:25:24 Annotating text fragment 64031/100392
## 2023-09-16 21:25:24 Annotating text fragment 64041/100392
## 2023-09-16 21:25:24 Annotating text fragment 64051/100392
## 2023-09-16 21:25:24 Annotating text fragment 64061/100392
## 2023-09-16 21:25:25 Annotating text fragment 64071/100392
## 2023-09-16 21:25:25 Annotating text fragment 64081/100392
## 2023-09-16 21:25:25 Annotating text fragment 64091/100392
## 2023-09-16 21:25:25 Annotating text fragment 64101/100392
## 2023-09-16 21:25:25 Annotating text fragment 64111/100392
## 2023-09-16 21:25:25 Annotating text fragment 64121/100392
## 2023-09-16 21:25:25 Annotating text fragment 64131/100392
## 2023-09-16 21:25:25 Annotating text fragment 64141/100392
## 2023-09-16 21:25:25 Annotating text fragment 64151/100392
## 2023-09-16 21:25:26 Annotating text fragment 64161/100392
## 2023-09-16 21:25:26 Annotating text fragment 64171/100392
## 2023-09-16 21:25:26 Annotating text fragment 64181/100392
## 2023-09-16 21:25:26 Annotating text fragment 64191/100392
## 2023-09-16 21:25:26 Annotating text fragment 64201/100392
## 2023-09-16 21:25:26 Annotating text fragment 64211/100392
## 2023-09-16 21:25:26 Annotating text fragment 64221/100392
## 2023-09-16 21:25:26 Annotating text fragment 64231/100392
## 2023-09-16 21:25:26 Annotating text fragment 64241/100392
## 2023-09-16 21:25:26 Annotating text fragment 64251/100392
## 2023-09-16 21:25:26 Annotating text fragment 64261/100392
## 2023-09-16 21:25:26 Annotating text fragment 64271/100392
## 2023-09-16 21:25:27 Annotating text fragment 64281/100392
## 2023-09-16 21:25:27 Annotating text fragment 64291/100392
## 2023-09-16 21:25:27 Annotating text fragment 64301/100392
## 2023-09-16 21:25:27 Annotating text fragment 64311/100392
```

```
## 2023-09-16 21:25:27 Annotating text fragment 64321/100392
## 2023-09-16 21:25:27 Annotating text fragment 64331/100392
## 2023-09-16 21:25:27 Annotating text fragment 64341/100392
## 2023-09-16 21:25:27 Annotating text fragment 64351/100392
## 2023-09-16 21:25:27 Annotating text fragment 64361/100392
## 2023-09-16 21:25:27 Annotating text fragment 64371/100392
## 2023-09-16 21:25:28 Annotating text fragment 64381/100392
## 2023-09-16 21:25:28 Annotating text fragment 64391/100392
## 2023-09-16 21:25:28 Annotating text fragment 64401/100392
## 2023-09-16 21:25:28 Annotating text fragment 64411/100392
## 2023-09-16 21:25:28 Annotating text fragment 64421/100392
## 2023-09-16 21:25:28 Annotating text fragment 64431/100392
## 2023-09-16 21:25:28 Annotating text fragment 64441/100392
## 2023-09-16 21:25:28 Annotating text fragment 64451/100392
## 2023-09-16 21:25:28 Annotating text fragment 64461/100392
## 2023-09-16 21:25:28 Annotating text fragment 64471/100392
## 2023-09-16 21:25:28 Annotating text fragment 64481/100392
## 2023-09-16 21:25:29 Annotating text fragment 64491/100392
## 2023-09-16 21:25:29 Annotating text fragment 64501/100392
## 2023-09-16 21:25:29 Annotating text fragment 64511/100392
## 2023-09-16 21:25:29 Annotating text fragment 64521/100392
## 2023-09-16 21:25:29 Annotating text fragment 64531/100392
## 2023-09-16 21:25:29 Annotating text fragment 64541/100392
## 2023-09-16 21:25:29 Annotating text fragment 64551/100392
## 2023-09-16 21:25:29 Annotating text fragment 64561/100392
## 2023-09-16 21:25:29 Annotating text fragment 64571/100392
## 2023-09-16 21:25:30 Annotating text fragment 64581/100392
## 2023-09-16 21:25:30 Annotating text fragment 64591/100392
## 2023-09-16 21:25:30 Annotating text fragment 64601/100392
## 2023-09-16 21:25:30 Annotating text fragment 64611/100392
## 2023-09-16 21:25:30 Annotating text fragment 64621/100392
## 2023-09-16 21:25:30 Annotating text fragment 64631/100392
## 2023-09-16 21:25:30 Annotating text fragment 64641/100392
## 2023-09-16 21:25:30 Annotating text fragment 64651/100392
## 2023-09-16 21:25:30 Annotating text fragment 64661/100392
## 2023-09-16 21:25:30 Annotating text fragment 64671/100392
## 2023-09-16 21:25:30 Annotating text fragment 64681/100392
## 2023-09-16 21:25:31 Annotating text fragment 64691/100392
## 2023-09-16 21:25:31 Annotating text fragment 64701/100392
## 2023-09-16 21:25:31 Annotating text fragment 64711/100392
## 2023-09-16 21:25:31 Annotating text fragment 64721/100392
## 2023-09-16 21:25:31 Annotating text fragment 64731/100392
## 2023-09-16 21:25:31 Annotating text fragment 64741/100392
## 2023-09-16 21:25:31 Annotating text fragment 64751/100392
## 2023-09-16 21:25:31 Annotating text fragment 64761/100392
## 2023-09-16 21:25:31 Annotating text fragment 64771/100392
## 2023-09-16 21:25:32 Annotating text fragment 64781/100392
## 2023-09-16 21:25:32 Annotating text fragment 64791/100392
## 2023-09-16 21:25:32 Annotating text fragment 64801/100392
## 2023-09-16 21:25:32 Annotating text fragment 64811/100392
## 2023-09-16 21:25:32 Annotating text fragment 64821/100392
## 2023-09-16 21:25:32 Annotating text fragment 64831/100392
## 2023-09-16 21:25:32 Annotating text fragment 64841/100392
## 2023-09-16 21:25:32 Annotating text fragment 64851/100392
```

```
## 2023-09-16 21:25:33 Annotating text fragment 64861/100392
## 2023-09-16 21:25:33 Annotating text fragment 64871/100392
## 2023-09-16 21:25:33 Annotating text fragment 64881/100392
## 2023-09-16 21:25:33 Annotating text fragment 64891/100392
## 2023-09-16 21:25:33 Annotating text fragment 64901/100392
## 2023-09-16 21:25:33 Annotating text fragment 64911/100392
## 2023-09-16 21:25:33 Annotating text fragment 64921/100392
## 2023-09-16 21:25:33 Annotating text fragment 64931/100392
## 2023-09-16 21:25:33 Annotating text fragment 64941/100392
## 2023-09-16 21:25:34 Annotating text fragment 64951/100392
## 2023-09-16 21:25:34 Annotating text fragment 64961/100392
## 2023-09-16 21:25:34 Annotating text fragment 64971/100392
## 2023-09-16 21:25:34 Annotating text fragment 64981/100392
## 2023-09-16 21:25:34 Annotating text fragment 64991/100392
## 2023-09-16 21:25:34 Annotating text fragment 65001/100392
## 2023-09-16 21:25:34 Annotating text fragment 65011/100392
## 2023-09-16 21:25:34 Annotating text fragment 65021/100392
## 2023-09-16 21:25:34 Annotating text fragment 65031/100392
## 2023-09-16 21:25:34 Annotating text fragment 65041/100392
## 2023-09-16 21:25:35 Annotating text fragment 65051/100392
## 2023-09-16 21:25:35 Annotating text fragment 65061/100392
## 2023-09-16 21:25:35 Annotating text fragment 65071/100392
## 2023-09-16 21:25:35 Annotating text fragment 65081/100392
## 2023-09-16 21:25:35 Annotating text fragment 65091/100392
## 2023-09-16 21:25:35 Annotating text fragment 65101/100392
## 2023-09-16 21:25:35 Annotating text fragment 65111/100392
## 2023-09-16 21:25:35 Annotating text fragment 65121/100392
## 2023-09-16 21:25:35 Annotating text fragment 65131/100392
## 2023-09-16 21:25:36 Annotating text fragment 65141/100392
## 2023-09-16 21:25:36 Annotating text fragment 65151/100392
## 2023-09-16 21:25:36 Annotating text fragment 65161/100392
## 2023-09-16 21:25:36 Annotating text fragment 65171/100392
## 2023-09-16 21:25:36 Annotating text fragment 65181/100392
## 2023-09-16 21:25:36 Annotating text fragment 65191/100392
## 2023-09-16 21:25:36 Annotating text fragment 65201/100392
## 2023-09-16 21:25:36 Annotating text fragment 65211/100392
## 2023-09-16 21:25:36 Annotating text fragment 65221/100392
## 2023-09-16 21:25:36 Annotating text fragment 65231/100392
## 2023-09-16 21:25:37 Annotating text fragment 65241/100392
## 2023-09-16 21:25:37 Annotating text fragment 65251/100392
## 2023-09-16 21:25:37 Annotating text fragment 65261/100392
## 2023-09-16 21:25:37 Annotating text fragment 65271/100392
## 2023-09-16 21:25:37 Annotating text fragment 65281/100392
## 2023-09-16 21:25:37 Annotating text fragment 65291/100392
## 2023-09-16 21:25:37 Annotating text fragment 65301/100392
## 2023-09-16 21:25:37 Annotating text fragment 65311/100392
## 2023-09-16 21:25:37 Annotating text fragment 65321/100392
## 2023-09-16 21:25:37 Annotating text fragment 65331/100392
## 2023-09-16 21:25:38 Annotating text fragment 65341/100392
## 2023-09-16 21:25:38 Annotating text fragment 65351/100392
## 2023-09-16 21:25:38 Annotating text fragment 65361/100392
## 2023-09-16 21:25:38 Annotating text fragment 65371/100392
## 2023-09-16 21:25:38 Annotating text fragment 65381/100392
## 2023-09-16 21:25:38 Annotating text fragment 65391/100392
```

```
## 2023-09-16 21:25:38 Annotating text fragment 65401/100392
## 2023-09-16 21:25:38 Annotating text fragment 65411/100392
## 2023-09-16 21:25:38 Annotating text fragment 65421/100392
## 2023-09-16 21:25:38 Annotating text fragment 65431/100392
## 2023-09-16 21:25:38 Annotating text fragment 65441/100392
## 2023-09-16 21:25:39 Annotating text fragment 65451/100392
## 2023-09-16 21:25:39 Annotating text fragment 65461/100392
## 2023-09-16 21:25:39 Annotating text fragment 65471/100392
## 2023-09-16 21:25:39 Annotating text fragment 65481/100392
## 2023-09-16 21:25:39 Annotating text fragment 65491/100392
## 2023-09-16 21:25:39 Annotating text fragment 65501/100392
## 2023-09-16 21:25:39 Annotating text fragment 65511/100392
## 2023-09-16 21:25:39 Annotating text fragment 65521/100392
## 2023-09-16 21:25:39 Annotating text fragment 65531/100392
## 2023-09-16 21:25:39 Annotating text fragment 65541/100392
## 2023-09-16 21:25:40 Annotating text fragment 65551/100392
## 2023-09-16 21:25:40 Annotating text fragment 65561/100392
## 2023-09-16 21:25:40 Annotating text fragment 65571/100392
## 2023-09-16 21:25:40 Annotating text fragment 65581/100392
## 2023-09-16 21:25:40 Annotating text fragment 65591/100392
## 2023-09-16 21:25:40 Annotating text fragment 65601/100392
## 2023-09-16 21:25:40 Annotating text fragment 65611/100392
## 2023-09-16 21:25:40 Annotating text fragment 65621/100392
## 2023-09-16 21:25:40 Annotating text fragment 65631/100392
## 2023-09-16 21:25:40 Annotating text fragment 65641/100392
## 2023-09-16 21:25:40 Annotating text fragment 65651/100392
## 2023-09-16 21:25:41 Annotating text fragment 65661/100392
## 2023-09-16 21:25:41 Annotating text fragment 65671/100392
## 2023-09-16 21:25:41 Annotating text fragment 65681/100392
## 2023-09-16 21:25:41 Annotating text fragment 65691/100392
## 2023-09-16 21:25:41 Annotating text fragment 65701/100392
## 2023-09-16 21:25:41 Annotating text fragment 65711/100392
## 2023-09-16 21:25:41 Annotating text fragment 65721/100392
## 2023-09-16 21:25:41 Annotating text fragment 65731/100392
## 2023-09-16 21:25:41 Annotating text fragment 65741/100392
## 2023-09-16 21:25:42 Annotating text fragment 65751/100392
## 2023-09-16 21:25:42 Annotating text fragment 65761/100392
## 2023-09-16 21:25:42 Annotating text fragment 65771/100392
## 2023-09-16 21:25:42 Annotating text fragment 65781/100392
## 2023-09-16 21:25:43 Annotating text fragment 65791/100392
## 2023-09-16 21:25:43 Annotating text fragment 65801/100392
## 2023-09-16 21:25:43 Annotating text fragment 65811/100392
## 2023-09-16 21:25:43 Annotating text fragment 65821/100392
## 2023-09-16 21:25:43 Annotating text fragment 65831/100392
## 2023-09-16 21:25:43 Annotating text fragment 65841/100392
## 2023-09-16 21:25:43 Annotating text fragment 65851/100392
## 2023-09-16 21:25:43 Annotating text fragment 65861/100392
## 2023-09-16 21:25:43 Annotating text fragment 65871/100392
## 2023-09-16 21:25:44 Annotating text fragment 65881/100392
## 2023-09-16 21:25:44 Annotating text fragment 65891/100392
## 2023-09-16 21:25:44 Annotating text fragment 65901/100392
## 2023-09-16 21:25:44 Annotating text fragment 65911/100392
## 2023-09-16 21:25:44 Annotating text fragment 65921/100392
## 2023-09-16 21:25:44 Annotating text fragment 65931/100392
```

```
## 2023-09-16 21:25:44 Annotating text fragment 65941/100392
## 2023-09-16 21:25:44 Annotating text fragment 65951/100392
## 2023-09-16 21:25:44 Annotating text fragment 65961/100392
## 2023-09-16 21:25:44 Annotating text fragment 65971/100392
## 2023-09-16 21:25:44 Annotating text fragment 65981/100392
## 2023-09-16 21:25:45 Annotating text fragment 65991/100392
## 2023-09-16 21:25:45 Annotating text fragment 66001/100392
## 2023-09-16 21:25:45 Annotating text fragment 66011/100392
## 2023-09-16 21:25:45 Annotating text fragment 66021/100392
## 2023-09-16 21:25:45 Annotating text fragment 66031/100392
## 2023-09-16 21:25:45 Annotating text fragment 66041/100392
## 2023-09-16 21:25:45 Annotating text fragment 66051/100392
## 2023-09-16 21:25:45 Annotating text fragment 66061/100392
## 2023-09-16 21:25:45 Annotating text fragment 66071/100392
## 2023-09-16 21:25:45 Annotating text fragment 66081/100392
## 2023-09-16 21:25:45 Annotating text fragment 66091/100392
## 2023-09-16 21:25:46 Annotating text fragment 66101/100392
## 2023-09-16 21:25:46 Annotating text fragment 66111/100392
## 2023-09-16 21:25:46 Annotating text fragment 66121/100392
## 2023-09-16 21:25:46 Annotating text fragment 66131/100392
## 2023-09-16 21:25:46 Annotating text fragment 66141/100392
## 2023-09-16 21:25:46 Annotating text fragment 66151/100392
## 2023-09-16 21:25:46 Annotating text fragment 66161/100392
## 2023-09-16 21:25:46 Annotating text fragment 66171/100392
## 2023-09-16 21:25:46 Annotating text fragment 66181/100392
## 2023-09-16 21:25:46 Annotating text fragment 66191/100392
## 2023-09-16 21:25:46 Annotating text fragment 66201/100392
## 2023-09-16 21:25:47 Annotating text fragment 66211/100392
## 2023-09-16 21:25:47 Annotating text fragment 66221/100392
## 2023-09-16 21:25:47 Annotating text fragment 66231/100392
## 2023-09-16 21:25:47 Annotating text fragment 66241/100392
## 2023-09-16 21:25:47 Annotating text fragment 66251/100392
## 2023-09-16 21:25:47 Annotating text fragment 66261/100392
## 2023-09-16 21:25:47 Annotating text fragment 66271/100392
## 2023-09-16 21:25:47 Annotating text fragment 66281/100392
## 2023-09-16 21:25:47 Annotating text fragment 66291/100392
## 2023-09-16 21:25:48 Annotating text fragment 66301/100392
## 2023-09-16 21:25:48 Annotating text fragment 66311/100392
## 2023-09-16 21:25:48 Annotating text fragment 66321/100392
## 2023-09-16 21:25:48 Annotating text fragment 66331/100392
## 2023-09-16 21:25:48 Annotating text fragment 66341/100392
## 2023-09-16 21:25:48 Annotating text fragment 66351/100392
## 2023-09-16 21:25:48 Annotating text fragment 66361/100392
## 2023-09-16 21:25:48 Annotating text fragment 66371/100392
## 2023-09-16 21:25:48 Annotating text fragment 66381/100392
## 2023-09-16 21:25:49 Annotating text fragment 66391/100392
## 2023-09-16 21:25:49 Annotating text fragment 66401/100392
## 2023-09-16 21:25:49 Annotating text fragment 66411/100392
## 2023-09-16 21:25:49 Annotating text fragment 66421/100392
## 2023-09-16 21:25:49 Annotating text fragment 66431/100392
## 2023-09-16 21:25:49 Annotating text fragment 66441/100392
## 2023-09-16 21:25:49 Annotating text fragment 66451/100392
## 2023-09-16 21:25:49 Annotating text fragment 66461/100392
## 2023-09-16 21:25:49 Annotating text fragment 66471/100392
```

```
## 2023-09-16 21:25:50 Annotating text fragment 66481/100392
## 2023-09-16 21:25:50 Annotating text fragment 66491/100392
## 2023-09-16 21:25:50 Annotating text fragment 66501/100392
## 2023-09-16 21:25:50 Annotating text fragment 66511/100392
## 2023-09-16 21:25:50 Annotating text fragment 66521/100392
## 2023-09-16 21:25:50 Annotating text fragment 66531/100392
## 2023-09-16 21:25:50 Annotating text fragment 66541/100392
## 2023-09-16 21:25:50 Annotating text fragment 66551/100392
## 2023-09-16 21:25:50 Annotating text fragment 66561/100392
## 2023-09-16 21:25:50 Annotating text fragment 66571/100392
## 2023-09-16 21:25:50 Annotating text fragment 66581/100392
## 2023-09-16 21:25:51 Annotating text fragment 66591/100392
## 2023-09-16 21:25:51 Annotating text fragment 66601/100392
## 2023-09-16 21:25:51 Annotating text fragment 66611/100392
## 2023-09-16 21:25:51 Annotating text fragment 66621/100392
## 2023-09-16 21:25:51 Annotating text fragment 66631/100392
## 2023-09-16 21:25:51 Annotating text fragment 66641/100392
## 2023-09-16 21:25:51 Annotating text fragment 66651/100392
## 2023-09-16 21:25:51 Annotating text fragment 66661/100392
## 2023-09-16 21:25:51 Annotating text fragment 66671/100392
## 2023-09-16 21:25:51 Annotating text fragment 66681/100392
## 2023-09-16 21:25:51 Annotating text fragment 66691/100392
## 2023-09-16 21:25:51 Annotating text fragment 66701/100392
## 2023-09-16 21:25:52 Annotating text fragment 66711/100392
## 2023-09-16 21:25:52 Annotating text fragment 66721/100392
## 2023-09-16 21:25:52 Annotating text fragment 66731/100392
## 2023-09-16 21:25:52 Annotating text fragment 66741/100392
## 2023-09-16 21:25:52 Annotating text fragment 66751/100392
## 2023-09-16 21:25:52 Annotating text fragment 66761/100392
## 2023-09-16 21:25:52 Annotating text fragment 66771/100392
## 2023-09-16 21:25:52 Annotating text fragment 66781/100392
## 2023-09-16 21:25:52 Annotating text fragment 66791/100392
## 2023-09-16 21:25:53 Annotating text fragment 66801/100392
## 2023-09-16 21:25:53 Annotating text fragment 66811/100392
## 2023-09-16 21:25:53 Annotating text fragment 66821/100392
## 2023-09-16 21:25:53 Annotating text fragment 66831/100392
## 2023-09-16 21:25:53 Annotating text fragment 66841/100392
## 2023-09-16 21:25:53 Annotating text fragment 66851/100392
## 2023-09-16 21:25:53 Annotating text fragment 66861/100392
## 2023-09-16 21:25:53 Annotating text fragment 66871/100392
## 2023-09-16 21:25:53 Annotating text fragment 66881/100392
## 2023-09-16 21:25:54 Annotating text fragment 66891/100392
## 2023-09-16 21:25:54 Annotating text fragment 66901/100392
## 2023-09-16 21:25:54 Annotating text fragment 66911/100392
## 2023-09-16 21:25:54 Annotating text fragment 66921/100392
## 2023-09-16 21:25:54 Annotating text fragment 66931/100392
## 2023-09-16 21:25:54 Annotating text fragment 66941/100392
## 2023-09-16 21:25:54 Annotating text fragment 66951/100392
## 2023-09-16 21:25:54 Annotating text fragment 66961/100392
## 2023-09-16 21:25:54 Annotating text fragment 66971/100392
## 2023-09-16 21:25:55 Annotating text fragment 66981/100392
## 2023-09-16 21:25:55 Annotating text fragment 66991/100392
## 2023-09-16 21:25:55 Annotating text fragment 67001/100392
## 2023-09-16 21:25:55 Annotating text fragment 67011/100392
```

```
## 2023-09-16 21:25:55 Annotating text fragment 67021/100392
## 2023-09-16 21:25:55 Annotating text fragment 67031/100392
## 2023-09-16 21:25:55 Annotating text fragment 67041/100392
## 2023-09-16 21:25:55 Annotating text fragment 67051/100392
## 2023-09-16 21:25:55 Annotating text fragment 67061/100392
## 2023-09-16 21:25:56 Annotating text fragment 67071/100392
## 2023-09-16 21:25:56 Annotating text fragment 67081/100392
## 2023-09-16 21:25:56 Annotating text fragment 67091/100392
## 2023-09-16 21:25:56 Annotating text fragment 67101/100392
## 2023-09-16 21:25:56 Annotating text fragment 67111/100392
## 2023-09-16 21:25:56 Annotating text fragment 67121/100392
## 2023-09-16 21:25:57 Annotating text fragment 67131/100392
## 2023-09-16 21:25:57 Annotating text fragment 67141/100392
## 2023-09-16 21:25:57 Annotating text fragment 67151/100392
## 2023-09-16 21:25:57 Annotating text fragment 67161/100392
## 2023-09-16 21:25:57 Annotating text fragment 67171/100392
## 2023-09-16 21:25:57 Annotating text fragment 67181/100392
## 2023-09-16 21:25:57 Annotating text fragment 67191/100392
## 2023-09-16 21:25:57 Annotating text fragment 67201/100392
## 2023-09-16 21:25:58 Annotating text fragment 67211/100392
## 2023-09-16 21:25:58 Annotating text fragment 67221/100392
## 2023-09-16 21:25:58 Annotating text fragment 67231/100392
## 2023-09-16 21:25:58 Annotating text fragment 67241/100392
## 2023-09-16 21:25:58 Annotating text fragment 67251/100392
## 2023-09-16 21:25:58 Annotating text fragment 67261/100392
## 2023-09-16 21:25:58 Annotating text fragment 67271/100392
## 2023-09-16 21:25:58 Annotating text fragment 67281/100392
## 2023-09-16 21:25:58 Annotating text fragment 67291/100392
## 2023-09-16 21:25:58 Annotating text fragment 67301/100392
## 2023-09-16 21:25:58 Annotating text fragment 67311/100392
## 2023-09-16 21:25:59 Annotating text fragment 67321/100392
## 2023-09-16 21:25:59 Annotating text fragment 67331/100392
## 2023-09-16 21:25:59 Annotating text fragment 67341/100392
## 2023-09-16 21:25:59 Annotating text fragment 67351/100392
## 2023-09-16 21:25:59 Annotating text fragment 67361/100392
## 2023-09-16 21:25:59 Annotating text fragment 67371/100392
## 2023-09-16 21:25:59 Annotating text fragment 67381/100392
## 2023-09-16 21:25:59 Annotating text fragment 67391/100392
## 2023-09-16 21:25:59 Annotating text fragment 67401/100392
## 2023-09-16 21:25:59 Annotating text fragment 67411/100392
## 2023-09-16 21:26:00 Annotating text fragment 67421/100392
## 2023-09-16 21:26:00 Annotating text fragment 67431/100392
## 2023-09-16 21:26:00 Annotating text fragment 67441/100392
## 2023-09-16 21:26:00 Annotating text fragment 67451/100392
## 2023-09-16 21:26:00 Annotating text fragment 67461/100392
## 2023-09-16 21:26:00 Annotating text fragment 67471/100392
## 2023-09-16 21:26:00 Annotating text fragment 67481/100392
## 2023-09-16 21:26:00 Annotating text fragment 67491/100392
## 2023-09-16 21:26:00 Annotating text fragment 67501/100392
## 2023-09-16 21:26:00 Annotating text fragment 67511/100392
## 2023-09-16 21:26:00 Annotating text fragment 67521/100392
## 2023-09-16 21:26:01 Annotating text fragment 67531/100392
## 2023-09-16 21:26:01 Annotating text fragment 67541/100392
## 2023-09-16 21:26:01 Annotating text fragment 67551/100392
```

```
## 2023-09-16 21:26:01 Annotating text fragment 67561/100392
## 2023-09-16 21:26:01 Annotating text fragment 67571/100392
## 2023-09-16 21:26:01 Annotating text fragment 67581/100392
## 2023-09-16 21:26:01 Annotating text fragment 67591/100392
## 2023-09-16 21:26:01 Annotating text fragment 67601/100392
## 2023-09-16 21:26:01 Annotating text fragment 67611/100392
## 2023-09-16 21:26:01 Annotating text fragment 67621/100392
## 2023-09-16 21:26:01 Annotating text fragment 67631/100392
## 2023-09-16 21:26:01 Annotating text fragment 67641/100392
## 2023-09-16 21:26:01 Annotating text fragment 67651/100392
## 2023-09-16 21:26:02 Annotating text fragment 67661/100392
## 2023-09-16 21:26:02 Annotating text fragment 67671/100392
## 2023-09-16 21:26:02 Annotating text fragment 67681/100392
## 2023-09-16 21:26:02 Annotating text fragment 67691/100392
## 2023-09-16 21:26:02 Annotating text fragment 67701/100392
## 2023-09-16 21:26:02 Annotating text fragment 67711/100392
## 2023-09-16 21:26:02 Annotating text fragment 67721/100392
## 2023-09-16 21:26:02 Annotating text fragment 67731/100392
## 2023-09-16 21:26:02 Annotating text fragment 67741/100392
## 2023-09-16 21:26:03 Annotating text fragment 67751/100392
## 2023-09-16 21:26:03 Annotating text fragment 67761/100392
## 2023-09-16 21:26:03 Annotating text fragment 67771/100392
## 2023-09-16 21:26:03 Annotating text fragment 67781/100392
## 2023-09-16 21:26:03 Annotating text fragment 67791/100392
## 2023-09-16 21:26:03 Annotating text fragment 67801/100392
## 2023-09-16 21:26:03 Annotating text fragment 67811/100392
## 2023-09-16 21:26:03 Annotating text fragment 67821/100392
## 2023-09-16 21:26:03 Annotating text fragment 67831/100392
## 2023-09-16 21:26:03 Annotating text fragment 67841/100392
## 2023-09-16 21:26:04 Annotating text fragment 67851/100392
## 2023-09-16 21:26:04 Annotating text fragment 67861/100392
## 2023-09-16 21:26:04 Annotating text fragment 67871/100392
## 2023-09-16 21:26:04 Annotating text fragment 67881/100392
## 2023-09-16 21:26:04 Annotating text fragment 67891/100392
## 2023-09-16 21:26:04 Annotating text fragment 67901/100392
## 2023-09-16 21:26:04 Annotating text fragment 67911/100392
## 2023-09-16 21:26:04 Annotating text fragment 67921/100392
## 2023-09-16 21:26:04 Annotating text fragment 67931/100392
## 2023-09-16 21:26:04 Annotating text fragment 67941/100392
## 2023-09-16 21:26:05 Annotating text fragment 67951/100392
## 2023-09-16 21:26:05 Annotating text fragment 67961/100392
## 2023-09-16 21:26:05 Annotating text fragment 67971/100392
## 2023-09-16 21:26:05 Annotating text fragment 67981/100392
## 2023-09-16 21:26:05 Annotating text fragment 67991/100392
## 2023-09-16 21:26:05 Annotating text fragment 68001/100392
## 2023-09-16 21:26:05 Annotating text fragment 68011/100392
## 2023-09-16 21:26:05 Annotating text fragment 68021/100392
## 2023-09-16 21:26:05 Annotating text fragment 68031/100392
## 2023-09-16 21:26:05 Annotating text fragment 68041/100392
## 2023-09-16 21:26:06 Annotating text fragment 68051/100392
## 2023-09-16 21:26:06 Annotating text fragment 68061/100392
## 2023-09-16 21:26:06 Annotating text fragment 68071/100392
## 2023-09-16 21:26:06 Annotating text fragment 68081/100392
## 2023-09-16 21:26:06 Annotating text fragment 68091/100392
```

```
## 2023-09-16 21:26:06 Annotating text fragment 68101/100392
## 2023-09-16 21:26:06 Annotating text fragment 68111/100392
## 2023-09-16 21:26:06 Annotating text fragment 68121/100392
## 2023-09-16 21:26:07 Annotating text fragment 68131/100392
## 2023-09-16 21:26:07 Annotating text fragment 68141/100392
## 2023-09-16 21:26:07 Annotating text fragment 68151/100392
## 2023-09-16 21:26:07 Annotating text fragment 68161/100392
## 2023-09-16 21:26:07 Annotating text fragment 68171/100392
## 2023-09-16 21:26:07 Annotating text fragment 68181/100392
## 2023-09-16 21:26:07 Annotating text fragment 68191/100392
## 2023-09-16 21:26:07 Annotating text fragment 68201/100392
## 2023-09-16 21:26:07 Annotating text fragment 68211/100392
## 2023-09-16 21:26:07 Annotating text fragment 68221/100392
## 2023-09-16 21:26:08 Annotating text fragment 68231/100392
## 2023-09-16 21:26:08 Annotating text fragment 68241/100392
## 2023-09-16 21:26:08 Annotating text fragment 68251/100392
## 2023-09-16 21:26:08 Annotating text fragment 68261/100392
## 2023-09-16 21:26:08 Annotating text fragment 68271/100392
## 2023-09-16 21:26:08 Annotating text fragment 68281/100392
## 2023-09-16 21:26:08 Annotating text fragment 68291/100392
## 2023-09-16 21:26:08 Annotating text fragment 68301/100392
## 2023-09-16 21:26:08 Annotating text fragment 68311/100392
## 2023-09-16 21:26:08 Annotating text fragment 68321/100392
## 2023-09-16 21:26:09 Annotating text fragment 68331/100392
## 2023-09-16 21:26:09 Annotating text fragment 68341/100392
## 2023-09-16 21:26:09 Annotating text fragment 68351/100392
## 2023-09-16 21:26:09 Annotating text fragment 68361/100392
## 2023-09-16 21:26:09 Annotating text fragment 68371/100392
## 2023-09-16 21:26:09 Annotating text fragment 68381/100392
## 2023-09-16 21:26:09 Annotating text fragment 68391/100392
## 2023-09-16 21:26:09 Annotating text fragment 68401/100392
## 2023-09-16 21:26:09 Annotating text fragment 68411/100392
## 2023-09-16 21:26:09 Annotating text fragment 68421/100392
## 2023-09-16 21:26:10 Annotating text fragment 68431/100392
## 2023-09-16 21:26:10 Annotating text fragment 68441/100392
## 2023-09-16 21:26:10 Annotating text fragment 68451/100392
## 2023-09-16 21:26:10 Annotating text fragment 68461/100392
## 2023-09-16 21:26:10 Annotating text fragment 68471/100392
## 2023-09-16 21:26:10 Annotating text fragment 68481/100392
## 2023-09-16 21:26:10 Annotating text fragment 68491/100392
## 2023-09-16 21:26:10 Annotating text fragment 68501/100392
## 2023-09-16 21:26:10 Annotating text fragment 68511/100392
## 2023-09-16 21:26:10 Annotating text fragment 68521/100392
## 2023-09-16 21:26:10 Annotating text fragment 68531/100392
## 2023-09-16 21:26:11 Annotating text fragment 68541/100392
## 2023-09-16 21:26:11 Annotating text fragment 68551/100392
## 2023-09-16 21:26:11 Annotating text fragment 68561/100392
## 2023-09-16 21:26:11 Annotating text fragment 68571/100392
## 2023-09-16 21:26:11 Annotating text fragment 68581/100392
## 2023-09-16 21:26:11 Annotating text fragment 68591/100392
## 2023-09-16 21:26:11 Annotating text fragment 68601/100392
## 2023-09-16 21:26:11 Annotating text fragment 68611/100392
## 2023-09-16 21:26:11 Annotating text fragment 68621/100392
## 2023-09-16 21:26:11 Annotating text fragment 68631/100392
```

```
## 2023-09-16 21:26:11 Annotating text fragment 68641/100392
## 2023-09-16 21:26:11 Annotating text fragment 68651/100392
## 2023-09-16 21:26:12 Annotating text fragment 68661/100392
## 2023-09-16 21:26:12 Annotating text fragment 68671/100392
## 2023-09-16 21:26:12 Annotating text fragment 68681/100392
## 2023-09-16 21:26:12 Annotating text fragment 68691/100392
## 2023-09-16 21:26:12 Annotating text fragment 68701/100392
## 2023-09-16 21:26:12 Annotating text fragment 68711/100392
## 2023-09-16 21:26:12 Annotating text fragment 68721/100392
## 2023-09-16 21:26:12 Annotating text fragment 68731/100392
## 2023-09-16 21:26:13 Annotating text fragment 68741/100392
## 2023-09-16 21:26:13 Annotating text fragment 68751/100392
## 2023-09-16 21:26:13 Annotating text fragment 68761/100392
## 2023-09-16 21:26:13 Annotating text fragment 68771/100392
## 2023-09-16 21:26:13 Annotating text fragment 68781/100392
## 2023-09-16 21:26:13 Annotating text fragment 68791/100392
## 2023-09-16 21:26:13 Annotating text fragment 68801/100392
## 2023-09-16 21:26:13 Annotating text fragment 68811/100392
## 2023-09-16 21:26:13 Annotating text fragment 68821/100392
## 2023-09-16 21:26:14 Annotating text fragment 68831/100392
## 2023-09-16 21:26:14 Annotating text fragment 68841/100392
## 2023-09-16 21:26:14 Annotating text fragment 68851/100392
## 2023-09-16 21:26:14 Annotating text fragment 68861/100392
## 2023-09-16 21:26:14 Annotating text fragment 68871/100392
## 2023-09-16 21:26:14 Annotating text fragment 68881/100392
## 2023-09-16 21:26:14 Annotating text fragment 68891/100392
## 2023-09-16 21:26:14 Annotating text fragment 68901/100392
## 2023-09-16 21:26:14 Annotating text fragment 68911/100392
## 2023-09-16 21:26:14 Annotating text fragment 68921/100392
## 2023-09-16 21:26:14 Annotating text fragment 68931/100392
## 2023-09-16 21:26:15 Annotating text fragment 68941/100392
## 2023-09-16 21:26:15 Annotating text fragment 68951/100392
## 2023-09-16 21:26:15 Annotating text fragment 68961/100392
## 2023-09-16 21:26:15 Annotating text fragment 68971/100392
## 2023-09-16 21:26:15 Annotating text fragment 68981/100392
## 2023-09-16 21:26:15 Annotating text fragment 68991/100392
## 2023-09-16 21:26:15 Annotating text fragment 69001/100392
## 2023-09-16 21:26:15 Annotating text fragment 69011/100392
## 2023-09-16 21:26:15 Annotating text fragment 69021/100392
## 2023-09-16 21:26:15 Annotating text fragment 69031/100392
## 2023-09-16 21:26:16 Annotating text fragment 69041/100392
## 2023-09-16 21:26:16 Annotating text fragment 69051/100392
## 2023-09-16 21:26:16 Annotating text fragment 69061/100392
## 2023-09-16 21:26:16 Annotating text fragment 69071/100392
## 2023-09-16 21:26:16 Annotating text fragment 69081/100392
## 2023-09-16 21:26:16 Annotating text fragment 69091/100392
## 2023-09-16 21:26:16 Annotating text fragment 69101/100392
## 2023-09-16 21:26:16 Annotating text fragment 69111/100392
## 2023-09-16 21:26:17 Annotating text fragment 69121/100392
## 2023-09-16 21:26:17 Annotating text fragment 69131/100392
## 2023-09-16 21:26:17 Annotating text fragment 69141/100392
## 2023-09-16 21:26:17 Annotating text fragment 69151/100392
## 2023-09-16 21:26:17 Annotating text fragment 69161/100392
## 2023-09-16 21:26:17 Annotating text fragment 69171/100392
```

```
## 2023-09-16 21:26:17 Annotating text fragment 69181/100392
## 2023-09-16 21:26:17 Annotating text fragment 69191/100392
## 2023-09-16 21:26:17 Annotating text fragment 69201/100392
## 2023-09-16 21:26:17 Annotating text fragment 69211/100392
## 2023-09-16 21:26:18 Annotating text fragment 69221/100392
## 2023-09-16 21:26:18 Annotating text fragment 69231/100392
## 2023-09-16 21:26:18 Annotating text fragment 69241/100392
## 2023-09-16 21:26:18 Annotating text fragment 69251/100392
## 2023-09-16 21:26:18 Annotating text fragment 69261/100392
## 2023-09-16 21:26:18 Annotating text fragment 69271/100392
## 2023-09-16 21:26:18 Annotating text fragment 69281/100392
## 2023-09-16 21:26:18 Annotating text fragment 69291/100392
## 2023-09-16 21:26:18 Annotating text fragment 69301/100392
## 2023-09-16 21:26:19 Annotating text fragment 69311/100392
## 2023-09-16 21:26:19 Annotating text fragment 69321/100392
## 2023-09-16 21:26:19 Annotating text fragment 69331/100392
## 2023-09-16 21:26:19 Annotating text fragment 69341/100392
## 2023-09-16 21:26:19 Annotating text fragment 69351/100392
## 2023-09-16 21:26:19 Annotating text fragment 69361/100392
## 2023-09-16 21:26:19 Annotating text fragment 69371/100392
## 2023-09-16 21:26:19 Annotating text fragment 69381/100392
## 2023-09-16 21:26:19 Annotating text fragment 69391/100392
## 2023-09-16 21:26:19 Annotating text fragment 69401/100392
## 2023-09-16 21:26:19 Annotating text fragment 69411/100392
## 2023-09-16 21:26:19 Annotating text fragment 69421/100392
## 2023-09-16 21:26:19 Annotating text fragment 69431/100392
## 2023-09-16 21:26:20 Annotating text fragment 69441/100392
## 2023-09-16 21:26:20 Annotating text fragment 69451/100392
## 2023-09-16 21:26:20 Annotating text fragment 69461/100392
## 2023-09-16 21:26:20 Annotating text fragment 69471/100392
## 2023-09-16 21:26:20 Annotating text fragment 69481/100392
## 2023-09-16 21:26:20 Annotating text fragment 69491/100392
## 2023-09-16 21:26:20 Annotating text fragment 69501/100392
## 2023-09-16 21:26:20 Annotating text fragment 69511/100392
## 2023-09-16 21:26:20 Annotating text fragment 69521/100392
## 2023-09-16 21:26:21 Annotating text fragment 69531/100392
## 2023-09-16 21:26:21 Annotating text fragment 69541/100392
## 2023-09-16 21:26:21 Annotating text fragment 69551/100392
## 2023-09-16 21:26:21 Annotating text fragment 69561/100392
## 2023-09-16 21:26:21 Annotating text fragment 69571/100392
## 2023-09-16 21:26:21 Annotating text fragment 69581/100392
## 2023-09-16 21:26:21 Annotating text fragment 69591/100392
## 2023-09-16 21:26:21 Annotating text fragment 69601/100392
## 2023-09-16 21:26:21 Annotating text fragment 69611/100392
## 2023-09-16 21:26:22 Annotating text fragment 69621/100392
## 2023-09-16 21:26:22 Annotating text fragment 69631/100392
## 2023-09-16 21:26:22 Annotating text fragment 69641/100392
## 2023-09-16 21:26:22 Annotating text fragment 69651/100392
## 2023-09-16 21:26:22 Annotating text fragment 69661/100392
## 2023-09-16 21:26:22 Annotating text fragment 69671/100392
## 2023-09-16 21:26:22 Annotating text fragment 69681/100392
## 2023-09-16 21:26:22 Annotating text fragment 69691/100392
## 2023-09-16 21:26:22 Annotating text fragment 69701/100392
## 2023-09-16 21:26:23 Annotating text fragment 69711/100392
```

```
## 2023-09-16 21:26:23 Annotating text fragment 69721/100392
## 2023-09-16 21:26:23 Annotating text fragment 69731/100392
## 2023-09-16 21:26:23 Annotating text fragment 69741/100392
## 2023-09-16 21:26:23 Annotating text fragment 69751/100392
## 2023-09-16 21:26:23 Annotating text fragment 69761/100392
## 2023-09-16 21:26:23 Annotating text fragment 69771/100392
## 2023-09-16 21:26:23 Annotating text fragment 69781/100392
## 2023-09-16 21:26:23 Annotating text fragment 69791/100392
## 2023-09-16 21:26:23 Annotating text fragment 69801/100392
## 2023-09-16 21:26:24 Annotating text fragment 69811/100392
## 2023-09-16 21:26:24 Annotating text fragment 69821/100392
## 2023-09-16 21:26:24 Annotating text fragment 69831/100392
## 2023-09-16 21:26:24 Annotating text fragment 69841/100392
## 2023-09-16 21:26:24 Annotating text fragment 69851/100392
## 2023-09-16 21:26:24 Annotating text fragment 69861/100392
## 2023-09-16 21:26:24 Annotating text fragment 69871/100392
## 2023-09-16 21:26:24 Annotating text fragment 69881/100392
## 2023-09-16 21:26:24 Annotating text fragment 69891/100392
## 2023-09-16 21:26:24 Annotating text fragment 69901/100392
## 2023-09-16 21:26:24 Annotating text fragment 69911/100392
## 2023-09-16 21:26:25 Annotating text fragment 69921/100392
## 2023-09-16 21:26:25 Annotating text fragment 69931/100392
## 2023-09-16 21:26:25 Annotating text fragment 69941/100392
## 2023-09-16 21:26:25 Annotating text fragment 69951/100392
## 2023-09-16 21:26:25 Annotating text fragment 69961/100392
## 2023-09-16 21:26:25 Annotating text fragment 69971/100392
## 2023-09-16 21:26:25 Annotating text fragment 69981/100392
## 2023-09-16 21:26:26 Annotating text fragment 69991/100392
## 2023-09-16 21:26:26 Annotating text fragment 70001/100392
## 2023-09-16 21:26:26 Annotating text fragment 70011/100392
## 2023-09-16 21:26:26 Annotating text fragment 70021/100392
## 2023-09-16 21:26:26 Annotating text fragment 70031/100392
## 2023-09-16 21:26:26 Annotating text fragment 70041/100392
## 2023-09-16 21:26:26 Annotating text fragment 70051/100392
## 2023-09-16 21:26:26 Annotating text fragment 70061/100392
## 2023-09-16 21:26:26 Annotating text fragment 70071/100392
## 2023-09-16 21:26:27 Annotating text fragment 70081/100392
## 2023-09-16 21:26:27 Annotating text fragment 70091/100392
## 2023-09-16 21:26:27 Annotating text fragment 70101/100392
## 2023-09-16 21:26:27 Annotating text fragment 70111/100392
## 2023-09-16 21:26:27 Annotating text fragment 70121/100392
## 2023-09-16 21:26:27 Annotating text fragment 70131/100392
## 2023-09-16 21:26:27 Annotating text fragment 70141/100392
## 2023-09-16 21:26:27 Annotating text fragment 70151/100392
## 2023-09-16 21:26:27 Annotating text fragment 70161/100392
## 2023-09-16 21:26:27 Annotating text fragment 70171/100392
## 2023-09-16 21:26:27 Annotating text fragment 70181/100392
## 2023-09-16 21:26:28 Annotating text fragment 70191/100392
## 2023-09-16 21:26:28 Annotating text fragment 70201/100392
## 2023-09-16 21:26:28 Annotating text fragment 70211/100392
## 2023-09-16 21:26:28 Annotating text fragment 70221/100392
## 2023-09-16 21:26:28 Annotating text fragment 70231/100392
## 2023-09-16 21:26:28 Annotating text fragment 70241/100392
## 2023-09-16 21:26:28 Annotating text fragment 70251/100392
```

```
## 2023-09-16 21:26:28 Annotating text fragment 70261/100392
## 2023-09-16 21:26:28 Annotating text fragment 70271/100392
## 2023-09-16 21:26:28 Annotating text fragment 70281/100392
## 2023-09-16 21:26:29 Annotating text fragment 70291/100392
## 2023-09-16 21:26:29 Annotating text fragment 70301/100392
## 2023-09-16 21:26:29 Annotating text fragment 70311/100392
## 2023-09-16 21:26:29 Annotating text fragment 70321/100392
## 2023-09-16 21:26:29 Annotating text fragment 70331/100392
## 2023-09-16 21:26:29 Annotating text fragment 70341/100392
## 2023-09-16 21:26:29 Annotating text fragment 70351/100392
## 2023-09-16 21:26:29 Annotating text fragment 70361/100392
## 2023-09-16 21:26:29 Annotating text fragment 70371/100392
## 2023-09-16 21:26:29 Annotating text fragment 70381/100392
## 2023-09-16 21:26:29 Annotating text fragment 70391/100392
## 2023-09-16 21:26:30 Annotating text fragment 70401/100392
## 2023-09-16 21:26:30 Annotating text fragment 70411/100392
## 2023-09-16 21:26:30 Annotating text fragment 70421/100392
## 2023-09-16 21:26:30 Annotating text fragment 70431/100392
## 2023-09-16 21:26:30 Annotating text fragment 70441/100392
## 2023-09-16 21:26:30 Annotating text fragment 70451/100392
## 2023-09-16 21:26:30 Annotating text fragment 70461/100392
## 2023-09-16 21:26:30 Annotating text fragment 70471/100392
## 2023-09-16 21:26:30 Annotating text fragment 70481/100392
## 2023-09-16 21:26:30 Annotating text fragment 70491/100392
## 2023-09-16 21:26:31 Annotating text fragment 70501/100392
## 2023-09-16 21:26:31 Annotating text fragment 70511/100392
## 2023-09-16 21:26:31 Annotating text fragment 70521/100392
## 2023-09-16 21:26:31 Annotating text fragment 70531/100392
## 2023-09-16 21:26:31 Annotating text fragment 70541/100392
## 2023-09-16 21:26:31 Annotating text fragment 70551/100392
## 2023-09-16 21:26:31 Annotating text fragment 70561/100392
## 2023-09-16 21:26:31 Annotating text fragment 70571/100392
## 2023-09-16 21:26:31 Annotating text fragment 70581/100392
## 2023-09-16 21:26:31 Annotating text fragment 70591/100392
## 2023-09-16 21:26:32 Annotating text fragment 70601/100392
## 2023-09-16 21:26:32 Annotating text fragment 70611/100392
## 2023-09-16 21:26:32 Annotating text fragment 70621/100392
## 2023-09-16 21:26:32 Annotating text fragment 70631/100392
## 2023-09-16 21:26:32 Annotating text fragment 70641/100392
## 2023-09-16 21:26:32 Annotating text fragment 70651/100392
## 2023-09-16 21:26:32 Annotating text fragment 70661/100392
## 2023-09-16 21:26:32 Annotating text fragment 70671/100392
## 2023-09-16 21:26:32 Annotating text fragment 70681/100392
## 2023-09-16 21:26:33 Annotating text fragment 70691/100392
## 2023-09-16 21:26:33 Annotating text fragment 70701/100392
## 2023-09-16 21:26:33 Annotating text fragment 70711/100392
## 2023-09-16 21:26:33 Annotating text fragment 70721/100392
## 2023-09-16 21:26:33 Annotating text fragment 70731/100392
## 2023-09-16 21:26:33 Annotating text fragment 70741/100392
## 2023-09-16 21:26:33 Annotating text fragment 70751/100392
## 2023-09-16 21:26:33 Annotating text fragment 70761/100392
## 2023-09-16 21:26:33 Annotating text fragment 70771/100392
## 2023-09-16 21:26:33 Annotating text fragment 70781/100392
## 2023-09-16 21:26:34 Annotating text fragment 70791/100392
```

```
## 2023-09-16 21:26:34 Annotating text fragment 70801/100392
## 2023-09-16 21:26:34 Annotating text fragment 70811/100392
## 2023-09-16 21:26:34 Annotating text fragment 70821/100392
## 2023-09-16 21:26:34 Annotating text fragment 70831/100392
## 2023-09-16 21:26:34 Annotating text fragment 70841/100392
## 2023-09-16 21:26:34 Annotating text fragment 70851/100392
## 2023-09-16 21:26:34 Annotating text fragment 70861/100392
## 2023-09-16 21:26:34 Annotating text fragment 70871/100392
## 2023-09-16 21:26:34 Annotating text fragment 70881/100392
## 2023-09-16 21:26:34 Annotating text fragment 70891/100392
## 2023-09-16 21:26:34 Annotating text fragment 70901/100392
## 2023-09-16 21:26:35 Annotating text fragment 70911/100392
## 2023-09-16 21:26:35 Annotating text fragment 70921/100392
## 2023-09-16 21:26:35 Annotating text fragment 70931/100392
## 2023-09-16 21:26:35 Annotating text fragment 70941/100392
## 2023-09-16 21:26:35 Annotating text fragment 70951/100392
## 2023-09-16 21:26:35 Annotating text fragment 70961/100392
## 2023-09-16 21:26:35 Annotating text fragment 70971/100392
## 2023-09-16 21:26:35 Annotating text fragment 70981/100392
## 2023-09-16 21:26:36 Annotating text fragment 70991/100392
## 2023-09-16 21:26:36 Annotating text fragment 71001/100392
## 2023-09-16 21:26:36 Annotating text fragment 71011/100392
## 2023-09-16 21:26:36 Annotating text fragment 71021/100392
## 2023-09-16 21:26:36 Annotating text fragment 71031/100392
## 2023-09-16 21:26:36 Annotating text fragment 71041/100392
## 2023-09-16 21:26:36 Annotating text fragment 71051/100392
## 2023-09-16 21:26:36 Annotating text fragment 71061/100392
## 2023-09-16 21:26:36 Annotating text fragment 71071/100392
## 2023-09-16 21:26:37 Annotating text fragment 71081/100392
## 2023-09-16 21:26:37 Annotating text fragment 71091/100392
## 2023-09-16 21:26:37 Annotating text fragment 71101/100392
## 2023-09-16 21:26:37 Annotating text fragment 71111/100392
## 2023-09-16 21:26:37 Annotating text fragment 71121/100392
## 2023-09-16 21:26:37 Annotating text fragment 71131/100392
## 2023-09-16 21:26:37 Annotating text fragment 71141/100392
## 2023-09-16 21:26:37 Annotating text fragment 71151/100392
## 2023-09-16 21:26:37 Annotating text fragment 71161/100392
## 2023-09-16 21:26:37 Annotating text fragment 71171/100392
## 2023-09-16 21:26:37 Annotating text fragment 71181/100392
## 2023-09-16 21:26:37 Annotating text fragment 71191/100392
## 2023-09-16 21:26:37 Annotating text fragment 71201/100392
## 2023-09-16 21:26:38 Annotating text fragment 71211/100392
## 2023-09-16 21:26:38 Annotating text fragment 71221/100392
## 2023-09-16 21:26:38 Annotating text fragment 71231/100392
## 2023-09-16 21:26:38 Annotating text fragment 71241/100392
## 2023-09-16 21:26:38 Annotating text fragment 71251/100392
## 2023-09-16 21:26:38 Annotating text fragment 71261/100392
## 2023-09-16 21:26:38 Annotating text fragment 71271/100392
## 2023-09-16 21:26:38 Annotating text fragment 71281/100392
## 2023-09-16 21:26:38 Annotating text fragment 71291/100392
## 2023-09-16 21:26:38 Annotating text fragment 71301/100392
## 2023-09-16 21:26:39 Annotating text fragment 71311/100392
## 2023-09-16 21:26:39 Annotating text fragment 71321/100392
## 2023-09-16 21:26:39 Annotating text fragment 71331/100392
```

```
## 2023-09-16 21:26:39 Annotating text fragment 71341/100392
## 2023-09-16 21:26:39 Annotating text fragment 71351/100392
## 2023-09-16 21:26:39 Annotating text fragment 71361/100392
## 2023-09-16 21:26:39 Annotating text fragment 71371/100392
## 2023-09-16 21:26:39 Annotating text fragment 71381/100392
## 2023-09-16 21:26:39 Annotating text fragment 71391/100392
## 2023-09-16 21:26:39 Annotating text fragment 71401/100392
## 2023-09-16 21:26:39 Annotating text fragment 71411/100392
## 2023-09-16 21:26:39 Annotating text fragment 71421/100392
## 2023-09-16 21:26:40 Annotating text fragment 71431/100392
## 2023-09-16 21:26:40 Annotating text fragment 71441/100392
## 2023-09-16 21:26:40 Annotating text fragment 71451/100392
## 2023-09-16 21:26:40 Annotating text fragment 71461/100392
## 2023-09-16 21:26:40 Annotating text fragment 71471/100392
## 2023-09-16 21:26:40 Annotating text fragment 71481/100392
## 2023-09-16 21:26:40 Annotating text fragment 71491/100392
## 2023-09-16 21:26:40 Annotating text fragment 71501/100392
## 2023-09-16 21:26:40 Annotating text fragment 71511/100392
## 2023-09-16 21:26:40 Annotating text fragment 71521/100392
## 2023-09-16 21:26:40 Annotating text fragment 71531/100392
## 2023-09-16 21:26:41 Annotating text fragment 71541/100392
## 2023-09-16 21:26:41 Annotating text fragment 71551/100392
## 2023-09-16 21:26:41 Annotating text fragment 71561/100392
## 2023-09-16 21:26:41 Annotating text fragment 71571/100392
## 2023-09-16 21:26:41 Annotating text fragment 71581/100392
## 2023-09-16 21:26:41 Annotating text fragment 71591/100392
## 2023-09-16 21:26:41 Annotating text fragment 71601/100392
## 2023-09-16 21:26:41 Annotating text fragment 71611/100392
## 2023-09-16 21:26:41 Annotating text fragment 71621/100392
## 2023-09-16 21:26:42 Annotating text fragment 71631/100392
## 2023-09-16 21:26:42 Annotating text fragment 71641/100392
## 2023-09-16 21:26:42 Annotating text fragment 71651/100392
## 2023-09-16 21:26:42 Annotating text fragment 71661/100392
## 2023-09-16 21:26:42 Annotating text fragment 71671/100392
## 2023-09-16 21:26:42 Annotating text fragment 71681/100392
## 2023-09-16 21:26:42 Annotating text fragment 71691/100392
## 2023-09-16 21:26:42 Annotating text fragment 71701/100392
## 2023-09-16 21:26:42 Annotating text fragment 71711/100392
## 2023-09-16 21:26:43 Annotating text fragment 71721/100392
## 2023-09-16 21:26:43 Annotating text fragment 71731/100392
## 2023-09-16 21:26:43 Annotating text fragment 71741/100392
## 2023-09-16 21:26:43 Annotating text fragment 71751/100392
## 2023-09-16 21:26:43 Annotating text fragment 71761/100392
## 2023-09-16 21:26:43 Annotating text fragment 71771/100392
## 2023-09-16 21:26:43 Annotating text fragment 71781/100392
## 2023-09-16 21:26:43 Annotating text fragment 71791/100392
## 2023-09-16 21:26:43 Annotating text fragment 71801/100392
## 2023-09-16 21:26:43 Annotating text fragment 71811/100392
## 2023-09-16 21:26:43 Annotating text fragment 71821/100392
## 2023-09-16 21:26:43 Annotating text fragment 71831/100392
## 2023-09-16 21:26:44 Annotating text fragment 71841/100392
## 2023-09-16 21:26:44 Annotating text fragment 71851/100392
## 2023-09-16 21:26:44 Annotating text fragment 71861/100392
## 2023-09-16 21:26:44 Annotating text fragment 71871/100392
```

```
## 2023-09-16 21:26:44 Annotating text fragment 71881/100392
## 2023-09-16 21:26:44 Annotating text fragment 71891/100392
## 2023-09-16 21:26:44 Annotating text fragment 71901/100392
## 2023-09-16 21:26:44 Annotating text fragment 71911/100392
## 2023-09-16 21:26:44 Annotating text fragment 71921/100392
## 2023-09-16 21:26:44 Annotating text fragment 71931/100392
## 2023-09-16 21:26:45 Annotating text fragment 71941/100392
## 2023-09-16 21:26:45 Annotating text fragment 71951/100392
## 2023-09-16 21:26:45 Annotating text fragment 71961/100392
## 2023-09-16 21:26:45 Annotating text fragment 71971/100392
## 2023-09-16 21:26:45 Annotating text fragment 71981/100392
## 2023-09-16 21:26:45 Annotating text fragment 71991/100392
## 2023-09-16 21:26:45 Annotating text fragment 72001/100392
## 2023-09-16 21:26:45 Annotating text fragment 72011/100392
## 2023-09-16 21:26:45 Annotating text fragment 72021/100392
## 2023-09-16 21:26:45 Annotating text fragment 72031/100392
## 2023-09-16 21:26:46 Annotating text fragment 72041/100392
## 2023-09-16 21:26:46 Annotating text fragment 72051/100392
## 2023-09-16 21:26:46 Annotating text fragment 72061/100392
## 2023-09-16 21:26:46 Annotating text fragment 72071/100392
## 2023-09-16 21:26:46 Annotating text fragment 72081/100392
## 2023-09-16 21:26:46 Annotating text fragment 72091/100392
## 2023-09-16 21:26:46 Annotating text fragment 72101/100392
## 2023-09-16 21:26:46 Annotating text fragment 72111/100392
## 2023-09-16 21:26:46 Annotating text fragment 72121/100392
## 2023-09-16 21:26:46 Annotating text fragment 72131/100392
## 2023-09-16 21:26:47 Annotating text fragment 72141/100392
## 2023-09-16 21:26:47 Annotating text fragment 72151/100392
## 2023-09-16 21:26:47 Annotating text fragment 72161/100392
## 2023-09-16 21:26:47 Annotating text fragment 72171/100392
## 2023-09-16 21:26:47 Annotating text fragment 72181/100392
## 2023-09-16 21:26:47 Annotating text fragment 72191/100392
## 2023-09-16 21:26:47 Annotating text fragment 72201/100392
## 2023-09-16 21:26:47 Annotating text fragment 72211/100392
## 2023-09-16 21:26:47 Annotating text fragment 72221/100392
## 2023-09-16 21:26:48 Annotating text fragment 72231/100392
## 2023-09-16 21:26:48 Annotating text fragment 72241/100392
## 2023-09-16 21:26:48 Annotating text fragment 72251/100392
## 2023-09-16 21:26:48 Annotating text fragment 72261/100392
## 2023-09-16 21:26:48 Annotating text fragment 72271/100392
## 2023-09-16 21:26:48 Annotating text fragment 72281/100392
## 2023-09-16 21:26:48 Annotating text fragment 72291/100392
## 2023-09-16 21:26:48 Annotating text fragment 72301/100392
## 2023-09-16 21:26:48 Annotating text fragment 72311/100392
## 2023-09-16 21:26:48 Annotating text fragment 72321/100392
## 2023-09-16 21:26:48 Annotating text fragment 72331/100392
## 2023-09-16 21:26:49 Annotating text fragment 72341/100392
## 2023-09-16 21:26:49 Annotating text fragment 72351/100392
## 2023-09-16 21:26:49 Annotating text fragment 72361/100392
## 2023-09-16 21:26:49 Annotating text fragment 72371/100392
## 2023-09-16 21:26:49 Annotating text fragment 72381/100392
## 2023-09-16 21:26:49 Annotating text fragment 72391/100392
## 2023-09-16 21:26:49 Annotating text fragment 72401/100392
## 2023-09-16 21:26:49 Annotating text fragment 72411/100392
```

```
## 2023-09-16 21:26:49 Annotating text fragment 72421/100392
## 2023-09-16 21:26:50 Annotating text fragment 72431/100392
## 2023-09-16 21:26:50 Annotating text fragment 72441/100392
## 2023-09-16 21:26:50 Annotating text fragment 72451/100392
## 2023-09-16 21:26:50 Annotating text fragment 72461/100392
## 2023-09-16 21:26:50 Annotating text fragment 72471/100392
## 2023-09-16 21:26:50 Annotating text fragment 72481/100392
## 2023-09-16 21:26:50 Annotating text fragment 72491/100392
## 2023-09-16 21:26:50 Annotating text fragment 72501/100392
## 2023-09-16 21:26:50 Annotating text fragment 72511/100392
## 2023-09-16 21:26:50 Annotating text fragment 72521/100392
## 2023-09-16 21:26:51 Annotating text fragment 72531/100392
## 2023-09-16 21:26:51 Annotating text fragment 72541/100392
## 2023-09-16 21:26:51 Annotating text fragment 72551/100392
## 2023-09-16 21:26:51 Annotating text fragment 72561/100392
## 2023-09-16 21:26:51 Annotating text fragment 72571/100392
## 2023-09-16 21:26:51 Annotating text fragment 72581/100392
## 2023-09-16 21:26:51 Annotating text fragment 72591/100392
## 2023-09-16 21:26:51 Annotating text fragment 72601/100392
## 2023-09-16 21:26:51 Annotating text fragment 72611/100392
## 2023-09-16 21:26:51 Annotating text fragment 72621/100392
## 2023-09-16 21:26:51 Annotating text fragment 72631/100392
## 2023-09-16 21:26:52 Annotating text fragment 72641/100392
## 2023-09-16 21:26:52 Annotating text fragment 72651/100392
## 2023-09-16 21:26:52 Annotating text fragment 72661/100392
## 2023-09-16 21:26:52 Annotating text fragment 72671/100392
## 2023-09-16 21:26:52 Annotating text fragment 72681/100392
## 2023-09-16 21:26:52 Annotating text fragment 72691/100392
## 2023-09-16 21:26:52 Annotating text fragment 72701/100392
## 2023-09-16 21:26:52 Annotating text fragment 72711/100392
## 2023-09-16 21:26:52 Annotating text fragment 72721/100392
## 2023-09-16 21:26:53 Annotating text fragment 72731/100392
## 2023-09-16 21:26:53 Annotating text fragment 72741/100392
## 2023-09-16 21:26:53 Annotating text fragment 72751/100392
## 2023-09-16 21:26:53 Annotating text fragment 72761/100392
## 2023-09-16 21:26:53 Annotating text fragment 72771/100392
## 2023-09-16 21:26:53 Annotating text fragment 72781/100392
## 2023-09-16 21:26:53 Annotating text fragment 72791/100392
## 2023-09-16 21:26:53 Annotating text fragment 72801/100392
## 2023-09-16 21:26:53 Annotating text fragment 72811/100392
## 2023-09-16 21:26:53 Annotating text fragment 72821/100392
## 2023-09-16 21:26:54 Annotating text fragment 72831/100392
## 2023-09-16 21:26:54 Annotating text fragment 72841/100392
## 2023-09-16 21:26:54 Annotating text fragment 72851/100392
## 2023-09-16 21:26:54 Annotating text fragment 72861/100392
## 2023-09-16 21:26:54 Annotating text fragment 72871/100392
## 2023-09-16 21:26:54 Annotating text fragment 72881/100392
## 2023-09-16 21:26:54 Annotating text fragment 72891/100392
## 2023-09-16 21:26:54 Annotating text fragment 72901/100392
## 2023-09-16 21:26:55 Annotating text fragment 72911/100392
## 2023-09-16 21:26:55 Annotating text fragment 72921/100392
## 2023-09-16 21:26:55 Annotating text fragment 72931/100392
## 2023-09-16 21:26:55 Annotating text fragment 72941/100392
## 2023-09-16 21:26:55 Annotating text fragment 72951/100392
```

```
## 2023-09-16 21:26:55 Annotating text fragment 72961/100392
## 2023-09-16 21:26:55 Annotating text fragment 72971/100392
## 2023-09-16 21:26:55 Annotating text fragment 72981/100392
## 2023-09-16 21:26:55 Annotating text fragment 72991/100392
## 2023-09-16 21:26:56 Annotating text fragment 73001/100392
## 2023-09-16 21:26:56 Annotating text fragment 73011/100392
## 2023-09-16 21:26:56 Annotating text fragment 73021/100392
## 2023-09-16 21:26:56 Annotating text fragment 73031/100392
## 2023-09-16 21:26:56 Annotating text fragment 73041/100392
## 2023-09-16 21:26:56 Annotating text fragment 73051/100392
## 2023-09-16 21:26:56 Annotating text fragment 73061/100392
## 2023-09-16 21:26:56 Annotating text fragment 73071/100392
## 2023-09-16 21:26:57 Annotating text fragment 73081/100392
## 2023-09-16 21:26:57 Annotating text fragment 73091/100392
## 2023-09-16 21:26:57 Annotating text fragment 73101/100392
## 2023-09-16 21:26:57 Annotating text fragment 73111/100392
## 2023-09-16 21:26:57 Annotating text fragment 73121/100392
## 2023-09-16 21:26:57 Annotating text fragment 73131/100392
## 2023-09-16 21:26:57 Annotating text fragment 73141/100392
## 2023-09-16 21:26:57 Annotating text fragment 73151/100392
## 2023-09-16 21:26:57 Annotating text fragment 73161/100392
## 2023-09-16 21:26:57 Annotating text fragment 73171/100392
## 2023-09-16 21:26:57 Annotating text fragment 73181/100392
## 2023-09-16 21:26:58 Annotating text fragment 73191/100392
## 2023-09-16 21:26:58 Annotating text fragment 73201/100392
## 2023-09-16 21:26:58 Annotating text fragment 73211/100392
## 2023-09-16 21:26:58 Annotating text fragment 73221/100392
## 2023-09-16 21:26:58 Annotating text fragment 73231/100392
## 2023-09-16 21:26:58 Annotating text fragment 73241/100392
## 2023-09-16 21:26:58 Annotating text fragment 73251/100392
## 2023-09-16 21:26:58 Annotating text fragment 73261/100392
## 2023-09-16 21:26:58 Annotating text fragment 73271/100392
## 2023-09-16 21:26:58 Annotating text fragment 73281/100392
## 2023-09-16 21:26:59 Annotating text fragment 73291/100392
## 2023-09-16 21:26:59 Annotating text fragment 73301/100392
## 2023-09-16 21:26:59 Annotating text fragment 73311/100392
## 2023-09-16 21:26:59 Annotating text fragment 73321/100392
## 2023-09-16 21:26:59 Annotating text fragment 73331/100392
## 2023-09-16 21:26:59 Annotating text fragment 73341/100392
## 2023-09-16 21:26:59 Annotating text fragment 73351/100392
## 2023-09-16 21:26:59 Annotating text fragment 73361/100392
## 2023-09-16 21:26:59 Annotating text fragment 73371/100392
## 2023-09-16 21:26:59 Annotating text fragment 73381/100392
## 2023-09-16 21:27:00 Annotating text fragment 73391/100392
## 2023-09-16 21:27:00 Annotating text fragment 73401/100392
## 2023-09-16 21:27:00 Annotating text fragment 73411/100392
## 2023-09-16 21:27:00 Annotating text fragment 73421/100392
## 2023-09-16 21:27:00 Annotating text fragment 73431/100392
## 2023-09-16 21:27:00 Annotating text fragment 73441/100392
## 2023-09-16 21:27:00 Annotating text fragment 73451/100392
## 2023-09-16 21:27:00 Annotating text fragment 73461/100392
## 2023-09-16 21:27:01 Annotating text fragment 73471/100392
## 2023-09-16 21:27:01 Annotating text fragment 73481/100392
## 2023-09-16 21:27:01 Annotating text fragment 73491/100392
```

```
## 2023-09-16 21:27:01 Annotating text fragment 73501/100392
## 2023-09-16 21:27:01 Annotating text fragment 73511/100392
## 2023-09-16 21:27:01 Annotating text fragment 73521/100392
## 2023-09-16 21:27:01 Annotating text fragment 73531/100392
## 2023-09-16 21:27:01 Annotating text fragment 73541/100392
## 2023-09-16 21:27:01 Annotating text fragment 73551/100392
## 2023-09-16 21:27:02 Annotating text fragment 73561/100392
## 2023-09-16 21:27:02 Annotating text fragment 73571/100392
## 2023-09-16 21:27:02 Annotating text fragment 73581/100392
## 2023-09-16 21:27:02 Annotating text fragment 73591/100392
## 2023-09-16 21:27:02 Annotating text fragment 73601/100392
## 2023-09-16 21:27:02 Annotating text fragment 73611/100392
## 2023-09-16 21:27:02 Annotating text fragment 73621/100392
## 2023-09-16 21:27:02 Annotating text fragment 73631/100392
## 2023-09-16 21:27:02 Annotating text fragment 73641/100392
## 2023-09-16 21:27:03 Annotating text fragment 73651/100392
## 2023-09-16 21:27:03 Annotating text fragment 73661/100392
## 2023-09-16 21:27:03 Annotating text fragment 73671/100392
## 2023-09-16 21:27:03 Annotating text fragment 73681/100392
## 2023-09-16 21:27:03 Annotating text fragment 73691/100392
## 2023-09-16 21:27:03 Annotating text fragment 73701/100392
## 2023-09-16 21:27:03 Annotating text fragment 73711/100392
## 2023-09-16 21:27:03 Annotating text fragment 73721/100392
## 2023-09-16 21:27:03 Annotating text fragment 73731/100392
## 2023-09-16 21:27:04 Annotating text fragment 73741/100392
## 2023-09-16 21:27:04 Annotating text fragment 73751/100392
## 2023-09-16 21:27:04 Annotating text fragment 73761/100392
## 2023-09-16 21:27:04 Annotating text fragment 73771/100392
## 2023-09-16 21:27:04 Annotating text fragment 73781/100392
## 2023-09-16 21:27:04 Annotating text fragment 73791/100392
## 2023-09-16 21:27:04 Annotating text fragment 73801/100392
## 2023-09-16 21:27:04 Annotating text fragment 73811/100392
## 2023-09-16 21:27:04 Annotating text fragment 73821/100392
## 2023-09-16 21:27:04 Annotating text fragment 73831/100392
## 2023-09-16 21:27:04 Annotating text fragment 73841/100392
## 2023-09-16 21:27:05 Annotating text fragment 73851/100392
## 2023-09-16 21:27:05 Annotating text fragment 73861/100392
## 2023-09-16 21:27:05 Annotating text fragment 73871/100392
## 2023-09-16 21:27:05 Annotating text fragment 73881/100392
## 2023-09-16 21:27:05 Annotating text fragment 73891/100392
## 2023-09-16 21:27:05 Annotating text fragment 73901/100392
## 2023-09-16 21:27:05 Annotating text fragment 73911/100392
## 2023-09-16 21:27:05 Annotating text fragment 73921/100392
## 2023-09-16 21:27:05 Annotating text fragment 73931/100392
## 2023-09-16 21:27:05 Annotating text fragment 73941/100392
## 2023-09-16 21:27:05 Annotating text fragment 73951/100392
## 2023-09-16 21:27:06 Annotating text fragment 73961/100392
## 2023-09-16 21:27:06 Annotating text fragment 73971/100392
## 2023-09-16 21:27:06 Annotating text fragment 73981/100392
## 2023-09-16 21:27:06 Annotating text fragment 73991/100392
## 2023-09-16 21:27:06 Annotating text fragment 74001/100392
## 2023-09-16 21:27:06 Annotating text fragment 74011/100392
## 2023-09-16 21:27:06 Annotating text fragment 74021/100392
## 2023-09-16 21:27:06 Annotating text fragment 74031/100392
```

```
## 2023-09-16 21:27:07 Annotating text fragment 74041/100392
## 2023-09-16 21:27:07 Annotating text fragment 74051/100392
## 2023-09-16 21:27:07 Annotating text fragment 74061/100392
## 2023-09-16 21:27:07 Annotating text fragment 74071/100392
## 2023-09-16 21:27:07 Annotating text fragment 74081/100392
## 2023-09-16 21:27:07 Annotating text fragment 74091/100392
## 2023-09-16 21:27:07 Annotating text fragment 74101/100392
## 2023-09-16 21:27:07 Annotating text fragment 74111/100392
## 2023-09-16 21:27:07 Annotating text fragment 74121/100392
## 2023-09-16 21:27:07 Annotating text fragment 74131/100392
## 2023-09-16 21:27:08 Annotating text fragment 74141/100392
## 2023-09-16 21:27:08 Annotating text fragment 74151/100392
## 2023-09-16 21:27:08 Annotating text fragment 74161/100392
## 2023-09-16 21:27:08 Annotating text fragment 74171/100392
## 2023-09-16 21:27:08 Annotating text fragment 74181/100392
## 2023-09-16 21:27:08 Annotating text fragment 74191/100392
## 2023-09-16 21:27:08 Annotating text fragment 74201/100392
## 2023-09-16 21:27:08 Annotating text fragment 74211/100392
## 2023-09-16 21:27:08 Annotating text fragment 74221/100392
## 2023-09-16 21:27:09 Annotating text fragment 74231/100392
## 2023-09-16 21:27:09 Annotating text fragment 74241/100392
## 2023-09-16 21:27:09 Annotating text fragment 74251/100392
## 2023-09-16 21:27:09 Annotating text fragment 74261/100392
## 2023-09-16 21:27:09 Annotating text fragment 74271/100392
## 2023-09-16 21:27:09 Annotating text fragment 74281/100392
## 2023-09-16 21:27:09 Annotating text fragment 74291/100392
## 2023-09-16 21:27:09 Annotating text fragment 74301/100392
## 2023-09-16 21:27:09 Annotating text fragment 74311/100392
## 2023-09-16 21:27:10 Annotating text fragment 74321/100392
## 2023-09-16 21:27:10 Annotating text fragment 74331/100392
## 2023-09-16 21:27:10 Annotating text fragment 74341/100392
## 2023-09-16 21:27:10 Annotating text fragment 74351/100392
## 2023-09-16 21:27:10 Annotating text fragment 74361/100392
## 2023-09-16 21:27:10 Annotating text fragment 74371/100392
## 2023-09-16 21:27:10 Annotating text fragment 74381/100392
## 2023-09-16 21:27:10 Annotating text fragment 74391/100392
## 2023-09-16 21:27:10 Annotating text fragment 74401/100392
## 2023-09-16 21:27:10 Annotating text fragment 74411/100392
## 2023-09-16 21:27:10 Annotating text fragment 74421/100392
## 2023-09-16 21:27:10 Annotating text fragment 74431/100392
## 2023-09-16 21:27:11 Annotating text fragment 74441/100392
## 2023-09-16 21:27:11 Annotating text fragment 74451/100392
## 2023-09-16 21:27:11 Annotating text fragment 74461/100392
## 2023-09-16 21:27:11 Annotating text fragment 74471/100392
## 2023-09-16 21:27:11 Annotating text fragment 74481/100392
## 2023-09-16 21:27:11 Annotating text fragment 74491/100392
## 2023-09-16 21:27:11 Annotating text fragment 74501/100392
## 2023-09-16 21:27:11 Annotating text fragment 74511/100392
## 2023-09-16 21:27:11 Annotating text fragment 74521/100392
## 2023-09-16 21:27:11 Annotating text fragment 74531/100392
## 2023-09-16 21:27:11 Annotating text fragment 74541/100392
## 2023-09-16 21:27:11 Annotating text fragment 74551/100392
## 2023-09-16 21:27:12 Annotating text fragment 74561/100392
## 2023-09-16 21:27:12 Annotating text fragment 74571/100392
```

```
## 2023-09-16 21:27:12 Annotating text fragment 74581/100392
## 2023-09-16 21:27:12 Annotating text fragment 74591/100392
## 2023-09-16 21:27:12 Annotating text fragment 74601/100392
## 2023-09-16 21:27:12 Annotating text fragment 74611/100392
## 2023-09-16 21:27:12 Annotating text fragment 74621/100392
## 2023-09-16 21:27:12 Annotating text fragment 74631/100392
## 2023-09-16 21:27:12 Annotating text fragment 74641/100392
## 2023-09-16 21:27:12 Annotating text fragment 74651/100392
## 2023-09-16 21:27:13 Annotating text fragment 74661/100392
## 2023-09-16 21:27:13 Annotating text fragment 74671/100392
## 2023-09-16 21:27:13 Annotating text fragment 74681/100392
## 2023-09-16 21:27:13 Annotating text fragment 74691/100392
## 2023-09-16 21:27:13 Annotating text fragment 74701/100392
## 2023-09-16 21:27:13 Annotating text fragment 74711/100392
## 2023-09-16 21:27:13 Annotating text fragment 74721/100392
## 2023-09-16 21:27:13 Annotating text fragment 74731/100392
## 2023-09-16 21:27:13 Annotating text fragment 74741/100392
## 2023-09-16 21:27:13 Annotating text fragment 74751/100392
## 2023-09-16 21:27:14 Annotating text fragment 74761/100392
## 2023-09-16 21:27:14 Annotating text fragment 74771/100392
## 2023-09-16 21:27:14 Annotating text fragment 74781/100392
## 2023-09-16 21:27:14 Annotating text fragment 74791/100392
## 2023-09-16 21:27:14 Annotating text fragment 74801/100392
## 2023-09-16 21:27:14 Annotating text fragment 74811/100392
## 2023-09-16 21:27:14 Annotating text fragment 74821/100392
## 2023-09-16 21:27:14 Annotating text fragment 74831/100392
## 2023-09-16 21:27:14 Annotating text fragment 74841/100392
## 2023-09-16 21:27:14 Annotating text fragment 74851/100392
## 2023-09-16 21:27:15 Annotating text fragment 74861/100392
## 2023-09-16 21:27:15 Annotating text fragment 74871/100392
## 2023-09-16 21:27:15 Annotating text fragment 74881/100392
## 2023-09-16 21:27:15 Annotating text fragment 74891/100392
## 2023-09-16 21:27:15 Annotating text fragment 74901/100392
## 2023-09-16 21:27:15 Annotating text fragment 74911/100392
## 2023-09-16 21:27:15 Annotating text fragment 74921/100392
## 2023-09-16 21:27:15 Annotating text fragment 74931/100392
## 2023-09-16 21:27:15 Annotating text fragment 74941/100392
## 2023-09-16 21:27:15 Annotating text fragment 74951/100392
## 2023-09-16 21:27:16 Annotating text fragment 74961/100392
## 2023-09-16 21:27:16 Annotating text fragment 74971/100392
## 2023-09-16 21:27:16 Annotating text fragment 74981/100392
## 2023-09-16 21:27:16 Annotating text fragment 74991/100392
## 2023-09-16 21:27:16 Annotating text fragment 75001/100392
## 2023-09-16 21:27:16 Annotating text fragment 75011/100392
## 2023-09-16 21:27:16 Annotating text fragment 75021/100392
## 2023-09-16 21:27:17 Annotating text fragment 75031/100392
## 2023-09-16 21:27:17 Annotating text fragment 75041/100392
## 2023-09-16 21:27:17 Annotating text fragment 75051/100392
## 2023-09-16 21:27:17 Annotating text fragment 75061/100392
## 2023-09-16 21:27:17 Annotating text fragment 75071/100392
## 2023-09-16 21:27:17 Annotating text fragment 75081/100392
## 2023-09-16 21:27:17 Annotating text fragment 75091/100392
## 2023-09-16 21:27:17 Annotating text fragment 75101/100392
## 2023-09-16 21:27:18 Annotating text fragment 75111/100392
```

```
## 2023-09-16 21:27:18 Annotating text fragment 75121/100392
## 2023-09-16 21:27:18 Annotating text fragment 75131/100392
## 2023-09-16 21:27:18 Annotating text fragment 75141/100392
## 2023-09-16 21:27:18 Annotating text fragment 75151/100392
## 2023-09-16 21:27:18 Annotating text fragment 75161/100392
## 2023-09-16 21:27:18 Annotating text fragment 75171/100392
## 2023-09-16 21:27:18 Annotating text fragment 75181/100392
## 2023-09-16 21:27:19 Annotating text fragment 75191/100392
## 2023-09-16 21:27:19 Annotating text fragment 75201/100392
## 2023-09-16 21:27:19 Annotating text fragment 75211/100392
## 2023-09-16 21:27:19 Annotating text fragment 75221/100392
## 2023-09-16 21:27:19 Annotating text fragment 75231/100392
## 2023-09-16 21:27:19 Annotating text fragment 75241/100392
## 2023-09-16 21:27:19 Annotating text fragment 75251/100392
## 2023-09-16 21:27:19 Annotating text fragment 75261/100392
## 2023-09-16 21:27:19 Annotating text fragment 75271/100392
## 2023-09-16 21:27:19 Annotating text fragment 75281/100392
## 2023-09-16 21:27:20 Annotating text fragment 75291/100392
## 2023-09-16 21:27:20 Annotating text fragment 75301/100392
## 2023-09-16 21:27:20 Annotating text fragment 75311/100392
## 2023-09-16 21:27:20 Annotating text fragment 75321/100392
## 2023-09-16 21:27:20 Annotating text fragment 75331/100392
## 2023-09-16 21:27:20 Annotating text fragment 75341/100392
## 2023-09-16 21:27:20 Annotating text fragment 75351/100392
## 2023-09-16 21:27:20 Annotating text fragment 75361/100392
## 2023-09-16 21:27:20 Annotating text fragment 75371/100392
## 2023-09-16 21:27:20 Annotating text fragment 75381/100392
## 2023-09-16 21:27:20 Annotating text fragment 75391/100392
## 2023-09-16 21:27:21 Annotating text fragment 75401/100392
## 2023-09-16 21:27:21 Annotating text fragment 75411/100392
## 2023-09-16 21:27:21 Annotating text fragment 75421/100392
## 2023-09-16 21:27:21 Annotating text fragment 75431/100392
## 2023-09-16 21:27:21 Annotating text fragment 75441/100392
## 2023-09-16 21:27:21 Annotating text fragment 75451/100392
## 2023-09-16 21:27:21 Annotating text fragment 75461/100392
## 2023-09-16 21:27:21 Annotating text fragment 75471/100392
## 2023-09-16 21:27:22 Annotating text fragment 75481/100392
## 2023-09-16 21:27:22 Annotating text fragment 75491/100392
## 2023-09-16 21:27:22 Annotating text fragment 75501/100392
## 2023-09-16 21:27:22 Annotating text fragment 75511/100392
## 2023-09-16 21:27:22 Annotating text fragment 75521/100392
## 2023-09-16 21:27:23 Annotating text fragment 75531/100392
## 2023-09-16 21:27:23 Annotating text fragment 75541/100392
## 2023-09-16 21:27:23 Annotating text fragment 75551/100392
## 2023-09-16 21:27:23 Annotating text fragment 75561/100392
## 2023-09-16 21:27:23 Annotating text fragment 75571/100392
## 2023-09-16 21:27:23 Annotating text fragment 75581/100392
## 2023-09-16 21:27:23 Annotating text fragment 75591/100392
## 2023-09-16 21:27:23 Annotating text fragment 75601/100392
## 2023-09-16 21:27:23 Annotating text fragment 75611/100392
## 2023-09-16 21:27:23 Annotating text fragment 75621/100392
## 2023-09-16 21:27:24 Annotating text fragment 75631/100392
## 2023-09-16 21:27:24 Annotating text fragment 75641/100392
## 2023-09-16 21:27:24 Annotating text fragment 75651/100392
```

```
## 2023-09-16 21:27:24 Annotating text fragment 75661/100392
## 2023-09-16 21:27:24 Annotating text fragment 75671/100392
## 2023-09-16 21:27:24 Annotating text fragment 75681/100392
## 2023-09-16 21:27:24 Annotating text fragment 75691/100392
## 2023-09-16 21:27:24 Annotating text fragment 75701/100392
## 2023-09-16 21:27:24 Annotating text fragment 75711/100392
## 2023-09-16 21:27:25 Annotating text fragment 75721/100392
## 2023-09-16 21:27:25 Annotating text fragment 75731/100392
## 2023-09-16 21:27:25 Annotating text fragment 75741/100392
## 2023-09-16 21:27:25 Annotating text fragment 75751/100392
## 2023-09-16 21:27:25 Annotating text fragment 75761/100392
## 2023-09-16 21:27:25 Annotating text fragment 75771/100392
## 2023-09-16 21:27:25 Annotating text fragment 75781/100392
## 2023-09-16 21:27:25 Annotating text fragment 75791/100392
## 2023-09-16 21:27:25 Annotating text fragment 75801/100392
## 2023-09-16 21:27:26 Annotating text fragment 75811/100392
## 2023-09-16 21:27:26 Annotating text fragment 75821/100392
## 2023-09-16 21:27:26 Annotating text fragment 75831/100392
## 2023-09-16 21:27:26 Annotating text fragment 75841/100392
## 2023-09-16 21:27:26 Annotating text fragment 75851/100392
## 2023-09-16 21:27:26 Annotating text fragment 75861/100392
## 2023-09-16 21:27:26 Annotating text fragment 75871/100392
## 2023-09-16 21:27:27 Annotating text fragment 75881/100392
## 2023-09-16 21:27:27 Annotating text fragment 75891/100392
## 2023-09-16 21:27:27 Annotating text fragment 75901/100392
## 2023-09-16 21:27:27 Annotating text fragment 75911/100392
## 2023-09-16 21:27:27 Annotating text fragment 75921/100392
## 2023-09-16 21:27:27 Annotating text fragment 75931/100392
## 2023-09-16 21:27:27 Annotating text fragment 75941/100392
## 2023-09-16 21:27:27 Annotating text fragment 75951/100392
## 2023-09-16 21:27:27 Annotating text fragment 75961/100392
## 2023-09-16 21:27:28 Annotating text fragment 75971/100392
## 2023-09-16 21:27:28 Annotating text fragment 75981/100392
## 2023-09-16 21:27:28 Annotating text fragment 75991/100392
## 2023-09-16 21:27:28 Annotating text fragment 76001/100392
## 2023-09-16 21:27:28 Annotating text fragment 76011/100392
## 2023-09-16 21:27:28 Annotating text fragment 76021/100392
## 2023-09-16 21:27:28 Annotating text fragment 76031/100392
## 2023-09-16 21:27:28 Annotating text fragment 76041/100392
## 2023-09-16 21:27:29 Annotating text fragment 76051/100392
## 2023-09-16 21:27:29 Annotating text fragment 76061/100392
## 2023-09-16 21:27:29 Annotating text fragment 76071/100392
## 2023-09-16 21:27:29 Annotating text fragment 76081/100392
## 2023-09-16 21:27:29 Annotating text fragment 76091/100392
## 2023-09-16 21:27:29 Annotating text fragment 76101/100392
## 2023-09-16 21:27:29 Annotating text fragment 76111/100392
## 2023-09-16 21:27:30 Annotating text fragment 76121/100392
## 2023-09-16 21:27:30 Annotating text fragment 76131/100392
## 2023-09-16 21:27:30 Annotating text fragment 76141/100392
## 2023-09-16 21:27:30 Annotating text fragment 76151/100392
## 2023-09-16 21:27:30 Annotating text fragment 76161/100392
## 2023-09-16 21:27:30 Annotating text fragment 76171/100392
## 2023-09-16 21:27:30 Annotating text fragment 76181/100392
## 2023-09-16 21:27:30 Annotating text fragment 76191/100392
```

```
## 2023-09-16 21:27:30 Annotating text fragment 76201/100392
## 2023-09-16 21:27:30 Annotating text fragment 76211/100392
## 2023-09-16 21:27:31 Annotating text fragment 76221/100392
## 2023-09-16 21:27:31 Annotating text fragment 76231/100392
## 2023-09-16 21:27:31 Annotating text fragment 76241/100392
## 2023-09-16 21:27:31 Annotating text fragment 76251/100392
## 2023-09-16 21:27:31 Annotating text fragment 76261/100392
## 2023-09-16 21:27:31 Annotating text fragment 76271/100392
## 2023-09-16 21:27:31 Annotating text fragment 76281/100392
## 2023-09-16 21:27:31 Annotating text fragment 76291/100392
## 2023-09-16 21:27:31 Annotating text fragment 76301/100392
## 2023-09-16 21:27:31 Annotating text fragment 76311/100392
## 2023-09-16 21:27:31 Annotating text fragment 76321/100392
## 2023-09-16 21:27:32 Annotating text fragment 76331/100392
## 2023-09-16 21:27:32 Annotating text fragment 76341/100392
## 2023-09-16 21:27:32 Annotating text fragment 76351/100392
## 2023-09-16 21:27:32 Annotating text fragment 76361/100392
## 2023-09-16 21:27:32 Annotating text fragment 76371/100392
## 2023-09-16 21:27:32 Annotating text fragment 76381/100392
## 2023-09-16 21:27:32 Annotating text fragment 76391/100392
## 2023-09-16 21:27:32 Annotating text fragment 76401/100392
## 2023-09-16 21:27:32 Annotating text fragment 76411/100392
## 2023-09-16 21:27:32 Annotating text fragment 76421/100392
## 2023-09-16 21:27:33 Annotating text fragment 76431/100392
## 2023-09-16 21:27:33 Annotating text fragment 76441/100392
## 2023-09-16 21:27:33 Annotating text fragment 76451/100392
## 2023-09-16 21:27:33 Annotating text fragment 76461/100392
## 2023-09-16 21:27:33 Annotating text fragment 76471/100392
## 2023-09-16 21:27:33 Annotating text fragment 76481/100392
## 2023-09-16 21:27:33 Annotating text fragment 76491/100392
## 2023-09-16 21:27:33 Annotating text fragment 76501/100392
## 2023-09-16 21:27:34 Annotating text fragment 76511/100392
## 2023-09-16 21:27:34 Annotating text fragment 76521/100392
## 2023-09-16 21:27:34 Annotating text fragment 76531/100392
## 2023-09-16 21:27:34 Annotating text fragment 76541/100392
## 2023-09-16 21:27:34 Annotating text fragment 76551/100392
## 2023-09-16 21:27:34 Annotating text fragment 76561/100392
## 2023-09-16 21:27:34 Annotating text fragment 76571/100392
## 2023-09-16 21:27:34 Annotating text fragment 76581/100392
## 2023-09-16 21:27:34 Annotating text fragment 76591/100392
## 2023-09-16 21:27:34 Annotating text fragment 76601/100392
## 2023-09-16 21:27:35 Annotating text fragment 76611/100392
## 2023-09-16 21:27:35 Annotating text fragment 76621/100392
## 2023-09-16 21:27:35 Annotating text fragment 76631/100392
## 2023-09-16 21:27:35 Annotating text fragment 76641/100392
## 2023-09-16 21:27:35 Annotating text fragment 76651/100392
## 2023-09-16 21:27:35 Annotating text fragment 76661/100392
## 2023-09-16 21:27:35 Annotating text fragment 76671/100392
## 2023-09-16 21:27:35 Annotating text fragment 76681/100392
## 2023-09-16 21:27:35 Annotating text fragment 76691/100392
## 2023-09-16 21:27:35 Annotating text fragment 76701/100392
## 2023-09-16 21:27:36 Annotating text fragment 76711/100392
## 2023-09-16 21:27:36 Annotating text fragment 76721/100392
## 2023-09-16 21:27:36 Annotating text fragment 76731/100392
```

```
## 2023-09-16 21:27:36 Annotating text fragment 76741/100392
## 2023-09-16 21:27:36 Annotating text fragment 76751/100392
## 2023-09-16 21:27:36 Annotating text fragment 76761/100392
## 2023-09-16 21:27:36 Annotating text fragment 76771/100392
## 2023-09-16 21:27:36 Annotating text fragment 76781/100392
## 2023-09-16 21:27:36 Annotating text fragment 76791/100392
## 2023-09-16 21:27:37 Annotating text fragment 76801/100392
## 2023-09-16 21:27:37 Annotating text fragment 76811/100392
## 2023-09-16 21:27:37 Annotating text fragment 76821/100392
## 2023-09-16 21:27:37 Annotating text fragment 76831/100392
## 2023-09-16 21:27:37 Annotating text fragment 76841/100392
## 2023-09-16 21:27:37 Annotating text fragment 76851/100392
## 2023-09-16 21:27:37 Annotating text fragment 76861/100392
## 2023-09-16 21:27:37 Annotating text fragment 76871/100392
## 2023-09-16 21:27:37 Annotating text fragment 76881/100392
## 2023-09-16 21:27:38 Annotating text fragment 76891/100392
## 2023-09-16 21:27:38 Annotating text fragment 76901/100392
## 2023-09-16 21:27:38 Annotating text fragment 76911/100392
## 2023-09-16 21:27:38 Annotating text fragment 76921/100392
## 2023-09-16 21:27:38 Annotating text fragment 76931/100392
## 2023-09-16 21:27:38 Annotating text fragment 76941/100392
## 2023-09-16 21:27:38 Annotating text fragment 76951/100392
## 2023-09-16 21:27:38 Annotating text fragment 76961/100392
## 2023-09-16 21:27:38 Annotating text fragment 76971/100392
## 2023-09-16 21:27:39 Annotating text fragment 76981/100392
## 2023-09-16 21:27:39 Annotating text fragment 76991/100392
## 2023-09-16 21:27:39 Annotating text fragment 77001/100392
## 2023-09-16 21:27:39 Annotating text fragment 77011/100392
## 2023-09-16 21:27:39 Annotating text fragment 77021/100392
## 2023-09-16 21:27:39 Annotating text fragment 77031/100392
## 2023-09-16 21:27:40 Annotating text fragment 77041/100392
## 2023-09-16 21:27:40 Annotating text fragment 77051/100392
## 2023-09-16 21:27:40 Annotating text fragment 77061/100392
## 2023-09-16 21:27:40 Annotating text fragment 77071/100392
## 2023-09-16 21:27:40 Annotating text fragment 77081/100392
## 2023-09-16 21:27:40 Annotating text fragment 77091/100392
## 2023-09-16 21:27:40 Annotating text fragment 77101/100392
## 2023-09-16 21:27:40 Annotating text fragment 77111/100392
## 2023-09-16 21:27:41 Annotating text fragment 77121/100392
## 2023-09-16 21:27:41 Annotating text fragment 77131/100392
## 2023-09-16 21:27:41 Annotating text fragment 77141/100392
## 2023-09-16 21:27:41 Annotating text fragment 77151/100392
## 2023-09-16 21:27:41 Annotating text fragment 77161/100392
## 2023-09-16 21:27:41 Annotating text fragment 77171/100392
## 2023-09-16 21:27:41 Annotating text fragment 77181/100392
## 2023-09-16 21:27:41 Annotating text fragment 77191/100392
## 2023-09-16 21:27:41 Annotating text fragment 77201/100392
## 2023-09-16 21:27:42 Annotating text fragment 77211/100392
## 2023-09-16 21:27:42 Annotating text fragment 77221/100392
## 2023-09-16 21:27:42 Annotating text fragment 77231/100392
## 2023-09-16 21:27:42 Annotating text fragment 77241/100392
## 2023-09-16 21:27:42 Annotating text fragment 77251/100392
## 2023-09-16 21:27:42 Annotating text fragment 77261/100392
## 2023-09-16 21:27:42 Annotating text fragment 77271/100392
```

```
## 2023-09-16 21:27:42 Annotating text fragment 77281/100392
## 2023-09-16 21:27:42 Annotating text fragment 77291/100392
## 2023-09-16 21:27:42 Annotating text fragment 77301/100392
## 2023-09-16 21:27:42 Annotating text fragment 77311/100392
## 2023-09-16 21:27:43 Annotating text fragment 77321/100392
## 2023-09-16 21:27:43 Annotating text fragment 77331/100392
## 2023-09-16 21:27:43 Annotating text fragment 77341/100392
## 2023-09-16 21:27:43 Annotating text fragment 77351/100392
## 2023-09-16 21:27:43 Annotating text fragment 77361/100392
## 2023-09-16 21:27:43 Annotating text fragment 77371/100392
## 2023-09-16 21:27:43 Annotating text fragment 77381/100392
## 2023-09-16 21:27:43 Annotating text fragment 77391/100392
## 2023-09-16 21:27:43 Annotating text fragment 77401/100392
## 2023-09-16 21:27:43 Annotating text fragment 77411/100392
## 2023-09-16 21:27:43 Annotating text fragment 77421/100392
## 2023-09-16 21:27:44 Annotating text fragment 77431/100392
## 2023-09-16 21:27:44 Annotating text fragment 77441/100392
## 2023-09-16 21:27:44 Annotating text fragment 77451/100392
## 2023-09-16 21:27:44 Annotating text fragment 77461/100392
## 2023-09-16 21:27:44 Annotating text fragment 77471/100392
## 2023-09-16 21:27:44 Annotating text fragment 77481/100392
## 2023-09-16 21:27:44 Annotating text fragment 77491/100392
## 2023-09-16 21:27:45 Annotating text fragment 77501/100392
## 2023-09-16 21:27:45 Annotating text fragment 77511/100392
## 2023-09-16 21:27:45 Annotating text fragment 77521/100392
## 2023-09-16 21:27:45 Annotating text fragment 77531/100392
## 2023-09-16 21:27:45 Annotating text fragment 77541/100392
## 2023-09-16 21:27:45 Annotating text fragment 77551/100392
## 2023-09-16 21:27:45 Annotating text fragment 77561/100392
## 2023-09-16 21:27:45 Annotating text fragment 77571/100392
## 2023-09-16 21:27:45 Annotating text fragment 77581/100392
## 2023-09-16 21:27:45 Annotating text fragment 77591/100392
## 2023-09-16 21:27:46 Annotating text fragment 77601/100392
## 2023-09-16 21:27:46 Annotating text fragment 77611/100392
## 2023-09-16 21:27:46 Annotating text fragment 77621/100392
## 2023-09-16 21:27:46 Annotating text fragment 77631/100392
## 2023-09-16 21:27:46 Annotating text fragment 77641/100392
## 2023-09-16 21:27:46 Annotating text fragment 77651/100392
## 2023-09-16 21:27:46 Annotating text fragment 77661/100392
## 2023-09-16 21:27:46 Annotating text fragment 77671/100392
## 2023-09-16 21:27:46 Annotating text fragment 77681/100392
## 2023-09-16 21:27:47 Annotating text fragment 77691/100392
## 2023-09-16 21:27:47 Annotating text fragment 77701/100392
## 2023-09-16 21:27:47 Annotating text fragment 77711/100392
## 2023-09-16 21:27:47 Annotating text fragment 77721/100392
## 2023-09-16 21:27:47 Annotating text fragment 77731/100392
## 2023-09-16 21:27:47 Annotating text fragment 77741/100392
## 2023-09-16 21:27:47 Annotating text fragment 77751/100392
## 2023-09-16 21:27:47 Annotating text fragment 77761/100392
## 2023-09-16 21:27:47 Annotating text fragment 77771/100392
## 2023-09-16 21:27:48 Annotating text fragment 77781/100392
## 2023-09-16 21:27:48 Annotating text fragment 77791/100392
## 2023-09-16 21:27:48 Annotating text fragment 77801/100392
## 2023-09-16 21:27:48 Annotating text fragment 77811/100392
```

```
## 2023-09-16 21:27:48 Annotating text fragment 77821/100392
## 2023-09-16 21:27:48 Annotating text fragment 77831/100392
## 2023-09-16 21:27:48 Annotating text fragment 77841/100392
## 2023-09-16 21:27:48 Annotating text fragment 77851/100392
## 2023-09-16 21:27:48 Annotating text fragment 77861/100392
## 2023-09-16 21:27:49 Annotating text fragment 77871/100392
## 2023-09-16 21:27:49 Annotating text fragment 77881/100392
## 2023-09-16 21:27:49 Annotating text fragment 77891/100392
## 2023-09-16 21:27:49 Annotating text fragment 77901/100392
## 2023-09-16 21:27:49 Annotating text fragment 77911/100392
## 2023-09-16 21:27:49 Annotating text fragment 77921/100392
## 2023-09-16 21:27:49 Annotating text fragment 77931/100392
## 2023-09-16 21:27:49 Annotating text fragment 77941/100392
## 2023-09-16 21:27:49 Annotating text fragment 77951/100392
## 2023-09-16 21:27:49 Annotating text fragment 77961/100392
## 2023-09-16 21:27:50 Annotating text fragment 77971/100392
## 2023-09-16 21:27:50 Annotating text fragment 77981/100392
## 2023-09-16 21:27:50 Annotating text fragment 77991/100392
## 2023-09-16 21:27:50 Annotating text fragment 78001/100392
## 2023-09-16 21:27:50 Annotating text fragment 78011/100392
## 2023-09-16 21:27:51 Annotating text fragment 78021/100392
## 2023-09-16 21:27:51 Annotating text fragment 78031/100392
## 2023-09-16 21:27:51 Annotating text fragment 78041/100392
## 2023-09-16 21:27:51 Annotating text fragment 78051/100392
## 2023-09-16 21:27:51 Annotating text fragment 78061/100392
## 2023-09-16 21:27:51 Annotating text fragment 78071/100392
## 2023-09-16 21:27:51 Annotating text fragment 78081/100392
## 2023-09-16 21:27:51 Annotating text fragment 78091/100392
## 2023-09-16 21:27:51 Annotating text fragment 78101/100392
## 2023-09-16 21:27:51 Annotating text fragment 78111/100392
## 2023-09-16 21:27:51 Annotating text fragment 78121/100392
## 2023-09-16 21:27:51 Annotating text fragment 78131/100392
## 2023-09-16 21:27:52 Annotating text fragment 78141/100392
## 2023-09-16 21:27:52 Annotating text fragment 78151/100392
## 2023-09-16 21:27:52 Annotating text fragment 78161/100392
## 2023-09-16 21:27:52 Annotating text fragment 78171/100392
## 2023-09-16 21:27:52 Annotating text fragment 78181/100392
## 2023-09-16 21:27:52 Annotating text fragment 78191/100392
## 2023-09-16 21:27:52 Annotating text fragment 78201/100392
## 2023-09-16 21:27:52 Annotating text fragment 78211/100392
## 2023-09-16 21:27:52 Annotating text fragment 78221/100392
## 2023-09-16 21:27:53 Annotating text fragment 78231/100392
## 2023-09-16 21:27:53 Annotating text fragment 78241/100392
## 2023-09-16 21:27:53 Annotating text fragment 78251/100392
## 2023-09-16 21:27:53 Annotating text fragment 78261/100392
## 2023-09-16 21:27:53 Annotating text fragment 78271/100392
## 2023-09-16 21:27:53 Annotating text fragment 78281/100392
## 2023-09-16 21:27:53 Annotating text fragment 78291/100392
## 2023-09-16 21:27:53 Annotating text fragment 78301/100392
## 2023-09-16 21:27:53 Annotating text fragment 78311/100392
## 2023-09-16 21:27:54 Annotating text fragment 78321/100392
## 2023-09-16 21:27:54 Annotating text fragment 78331/100392
## 2023-09-16 21:27:54 Annotating text fragment 78341/100392
## 2023-09-16 21:27:54 Annotating text fragment 78351/100392
```

```
## 2023-09-16 21:27:54 Annotating text fragment 78361/100392
## 2023-09-16 21:27:54 Annotating text fragment 78371/100392
## 2023-09-16 21:27:54 Annotating text fragment 78381/100392
## 2023-09-16 21:27:54 Annotating text fragment 78391/100392
## 2023-09-16 21:27:54 Annotating text fragment 78401/100392
## 2023-09-16 21:27:54 Annotating text fragment 78411/100392
## 2023-09-16 21:27:55 Annotating text fragment 78421/100392
## 2023-09-16 21:27:55 Annotating text fragment 78431/100392
## 2023-09-16 21:27:55 Annotating text fragment 78441/100392
## 2023-09-16 21:27:55 Annotating text fragment 78451/100392
## 2023-09-16 21:27:55 Annotating text fragment 78461/100392
## 2023-09-16 21:27:55 Annotating text fragment 78471/100392
## 2023-09-16 21:27:55 Annotating text fragment 78481/100392
## 2023-09-16 21:27:55 Annotating text fragment 78491/100392
## 2023-09-16 21:27:55 Annotating text fragment 78501/100392
## 2023-09-16 21:27:56 Annotating text fragment 78511/100392
## 2023-09-16 21:27:56 Annotating text fragment 78521/100392
## 2023-09-16 21:27:56 Annotating text fragment 78531/100392
## 2023-09-16 21:27:56 Annotating text fragment 78541/100392
## 2023-09-16 21:27:56 Annotating text fragment 78551/100392
## 2023-09-16 21:27:56 Annotating text fragment 78561/100392
## 2023-09-16 21:27:56 Annotating text fragment 78571/100392
## 2023-09-16 21:27:56 Annotating text fragment 78581/100392
## 2023-09-16 21:27:56 Annotating text fragment 78591/100392
## 2023-09-16 21:27:56 Annotating text fragment 78601/100392
## 2023-09-16 21:27:57 Annotating text fragment 78611/100392
## 2023-09-16 21:27:57 Annotating text fragment 78621/100392
## 2023-09-16 21:27:57 Annotating text fragment 78631/100392
## 2023-09-16 21:27:57 Annotating text fragment 78641/100392
## 2023-09-16 21:27:57 Annotating text fragment 78651/100392
## 2023-09-16 21:27:57 Annotating text fragment 78661/100392
## 2023-09-16 21:27:57 Annotating text fragment 78671/100392
## 2023-09-16 21:27:57 Annotating text fragment 78681/100392
## 2023-09-16 21:27:57 Annotating text fragment 78691/100392
## 2023-09-16 21:27:57 Annotating text fragment 78701/100392
## 2023-09-16 21:27:57 Annotating text fragment 78711/100392
## 2023-09-16 21:27:58 Annotating text fragment 78721/100392
## 2023-09-16 21:27:58 Annotating text fragment 78731/100392
## 2023-09-16 21:27:58 Annotating text fragment 78741/100392
## 2023-09-16 21:27:58 Annotating text fragment 78751/100392
## 2023-09-16 21:27:58 Annotating text fragment 78761/100392
## 2023-09-16 21:27:58 Annotating text fragment 78771/100392
## 2023-09-16 21:27:58 Annotating text fragment 78781/100392
## 2023-09-16 21:27:58 Annotating text fragment 78791/100392
## 2023-09-16 21:27:58 Annotating text fragment 78801/100392
## 2023-09-16 21:27:59 Annotating text fragment 78811/100392
## 2023-09-16 21:27:59 Annotating text fragment 78821/100392
## 2023-09-16 21:27:59 Annotating text fragment 78831/100392
## 2023-09-16 21:27:59 Annotating text fragment 78841/100392
## 2023-09-16 21:27:59 Annotating text fragment 78851/100392
## 2023-09-16 21:27:59 Annotating text fragment 78861/100392
## 2023-09-16 21:27:59 Annotating text fragment 78871/100392
## 2023-09-16 21:27:59 Annotating text fragment 78881/100392
## 2023-09-16 21:27:59 Annotating text fragment 78891/100392
```

```
## 2023-09-16 21:28:00 Annotating text fragment 78901/100392
## 2023-09-16 21:28:00 Annotating text fragment 78911/100392
## 2023-09-16 21:28:00 Annotating text fragment 78921/100392
## 2023-09-16 21:28:00 Annotating text fragment 78931/100392
## 2023-09-16 21:28:00 Annotating text fragment 78941/100392
## 2023-09-16 21:28:00 Annotating text fragment 78951/100392
## 2023-09-16 21:28:00 Annotating text fragment 78961/100392
## 2023-09-16 21:28:00 Annotating text fragment 78971/100392
## 2023-09-16 21:28:00 Annotating text fragment 78981/100392
## 2023-09-16 21:28:00 Annotating text fragment 78991/100392
## 2023-09-16 21:28:00 Annotating text fragment 79001/100392
## 2023-09-16 21:28:01 Annotating text fragment 79011/100392
## 2023-09-16 21:28:01 Annotating text fragment 79021/100392
## 2023-09-16 21:28:01 Annotating text fragment 79031/100392
## 2023-09-16 21:28:01 Annotating text fragment 79041/100392
## 2023-09-16 21:28:01 Annotating text fragment 79051/100392
## 2023-09-16 21:28:01 Annotating text fragment 79061/100392
## 2023-09-16 21:28:01 Annotating text fragment 79071/100392
## 2023-09-16 21:28:01 Annotating text fragment 79081/100392
## 2023-09-16 21:28:01 Annotating text fragment 79091/100392
## 2023-09-16 21:28:01 Annotating text fragment 79101/100392
## 2023-09-16 21:28:02 Annotating text fragment 79111/100392
## 2023-09-16 21:28:02 Annotating text fragment 79121/100392
## 2023-09-16 21:28:02 Annotating text fragment 79131/100392
## 2023-09-16 21:28:02 Annotating text fragment 79141/100392
## 2023-09-16 21:28:02 Annotating text fragment 79151/100392
## 2023-09-16 21:28:02 Annotating text fragment 79161/100392
## 2023-09-16 21:28:02 Annotating text fragment 79171/100392
## 2023-09-16 21:28:02 Annotating text fragment 79181/100392
## 2023-09-16 21:28:02 Annotating text fragment 79191/100392
## 2023-09-16 21:28:02 Annotating text fragment 79201/100392
## 2023-09-16 21:28:02 Annotating text fragment 79211/100392
## 2023-09-16 21:28:02 Annotating text fragment 79221/100392
## 2023-09-16 21:28:03 Annotating text fragment 79231/100392
## 2023-09-16 21:28:03 Annotating text fragment 79241/100392
## 2023-09-16 21:28:03 Annotating text fragment 79251/100392
## 2023-09-16 21:28:03 Annotating text fragment 79261/100392
## 2023-09-16 21:28:03 Annotating text fragment 79271/100392
## 2023-09-16 21:28:03 Annotating text fragment 79281/100392
## 2023-09-16 21:28:03 Annotating text fragment 79291/100392
## 2023-09-16 21:28:03 Annotating text fragment 79301/100392
## 2023-09-16 21:28:04 Annotating text fragment 79311/100392
## 2023-09-16 21:28:04 Annotating text fragment 79321/100392
## 2023-09-16 21:28:04 Annotating text fragment 79331/100392
## 2023-09-16 21:28:04 Annotating text fragment 79341/100392
## 2023-09-16 21:28:04 Annotating text fragment 79351/100392
## 2023-09-16 21:28:04 Annotating text fragment 79361/100392
## 2023-09-16 21:28:04 Annotating text fragment 79371/100392
## 2023-09-16 21:28:04 Annotating text fragment 79381/100392
## 2023-09-16 21:28:04 Annotating text fragment 79391/100392
## 2023-09-16 21:28:04 Annotating text fragment 79401/100392
## 2023-09-16 21:28:05 Annotating text fragment 79411/100392
## 2023-09-16 21:28:05 Annotating text fragment 79421/100392
## 2023-09-16 21:28:05 Annotating text fragment 79431/100392
```

```
## 2023-09-16 21:28:05 Annotating text fragment 79441/100392
## 2023-09-16 21:28:05 Annotating text fragment 79451/100392
## 2023-09-16 21:28:05 Annotating text fragment 79461/100392
## 2023-09-16 21:28:05 Annotating text fragment 79471/100392
## 2023-09-16 21:28:05 Annotating text fragment 79481/100392
## 2023-09-16 21:28:05 Annotating text fragment 79491/100392
## 2023-09-16 21:28:06 Annotating text fragment 79501/100392
## 2023-09-16 21:28:06 Annotating text fragment 79511/100392
## 2023-09-16 21:28:06 Annotating text fragment 79521/100392
## 2023-09-16 21:28:06 Annotating text fragment 79531/100392
## 2023-09-16 21:28:06 Annotating text fragment 79541/100392
## 2023-09-16 21:28:06 Annotating text fragment 79551/100392
## 2023-09-16 21:28:06 Annotating text fragment 79561/100392
## 2023-09-16 21:28:06 Annotating text fragment 79571/100392
## 2023-09-16 21:28:06 Annotating text fragment 79581/100392
## 2023-09-16 21:28:06 Annotating text fragment 79591/100392
## 2023-09-16 21:28:06 Annotating text fragment 79601/100392
## 2023-09-16 21:28:07 Annotating text fragment 79611/100392
## 2023-09-16 21:28:07 Annotating text fragment 79621/100392
## 2023-09-16 21:28:07 Annotating text fragment 79631/100392
## 2023-09-16 21:28:07 Annotating text fragment 79641/100392
## 2023-09-16 21:28:07 Annotating text fragment 79651/100392
## 2023-09-16 21:28:07 Annotating text fragment 79661/100392
## 2023-09-16 21:28:07 Annotating text fragment 79671/100392
## 2023-09-16 21:28:07 Annotating text fragment 79681/100392
## 2023-09-16 21:28:07 Annotating text fragment 79691/100392
## 2023-09-16 21:28:07 Annotating text fragment 79701/100392
## 2023-09-16 21:28:08 Annotating text fragment 79711/100392
## 2023-09-16 21:28:08 Annotating text fragment 79721/100392
## 2023-09-16 21:28:08 Annotating text fragment 79731/100392
## 2023-09-16 21:28:08 Annotating text fragment 79741/100392
## 2023-09-16 21:28:08 Annotating text fragment 79751/100392
## 2023-09-16 21:28:08 Annotating text fragment 79761/100392
## 2023-09-16 21:28:08 Annotating text fragment 79771/100392
## 2023-09-16 21:28:08 Annotating text fragment 79781/100392
## 2023-09-16 21:28:08 Annotating text fragment 79791/100392
## 2023-09-16 21:28:08 Annotating text fragment 79801/100392
## 2023-09-16 21:28:09 Annotating text fragment 79811/100392
## 2023-09-16 21:28:09 Annotating text fragment 79821/100392
## 2023-09-16 21:28:09 Annotating text fragment 79831/100392
## 2023-09-16 21:28:09 Annotating text fragment 79841/100392
## 2023-09-16 21:28:09 Annotating text fragment 79851/100392
## 2023-09-16 21:28:09 Annotating text fragment 79861/100392
## 2023-09-16 21:28:09 Annotating text fragment 79871/100392
## 2023-09-16 21:28:09 Annotating text fragment 79881/100392
## 2023-09-16 21:28:09 Annotating text fragment 79891/100392
## 2023-09-16 21:28:09 Annotating text fragment 79901/100392
## 2023-09-16 21:28:09 Annotating text fragment 79911/100392
## 2023-09-16 21:28:10 Annotating text fragment 79921/100392
## 2023-09-16 21:28:10 Annotating text fragment 79931/100392
## 2023-09-16 21:28:10 Annotating text fragment 79941/100392
## 2023-09-16 21:28:10 Annotating text fragment 79951/100392
## 2023-09-16 21:28:10 Annotating text fragment 79961/100392
## 2023-09-16 21:28:10 Annotating text fragment 79971/100392
```

```
## 2023-09-16 21:28:10 Annotating text fragment 79981/100392
## 2023-09-16 21:28:10 Annotating text fragment 79991/100392
## 2023-09-16 21:28:10 Annotating text fragment 80001/100392
## 2023-09-16 21:28:10 Annotating text fragment 80011/100392
## 2023-09-16 21:28:11 Annotating text fragment 80021/100392
## 2023-09-16 21:28:11 Annotating text fragment 80031/100392
## 2023-09-16 21:28:11 Annotating text fragment 80041/100392
## 2023-09-16 21:28:11 Annotating text fragment 80051/100392
## 2023-09-16 21:28:11 Annotating text fragment 80061/100392
## 2023-09-16 21:28:11 Annotating text fragment 80071/100392
## 2023-09-16 21:28:11 Annotating text fragment 80081/100392
## 2023-09-16 21:28:11 Annotating text fragment 80091/100392
## 2023-09-16 21:28:11 Annotating text fragment 80101/100392
## 2023-09-16 21:28:11 Annotating text fragment 80111/100392
## 2023-09-16 21:28:11 Annotating text fragment 80121/100392
## 2023-09-16 21:28:11 Annotating text fragment 80131/100392
## 2023-09-16 21:28:12 Annotating text fragment 80141/100392
## 2023-09-16 21:28:12 Annotating text fragment 80151/100392
## 2023-09-16 21:28:12 Annotating text fragment 80161/100392
## 2023-09-16 21:28:12 Annotating text fragment 80171/100392
## 2023-09-16 21:28:12 Annotating text fragment 80181/100392
## 2023-09-16 21:28:12 Annotating text fragment 80191/100392
## 2023-09-16 21:28:12 Annotating text fragment 80201/100392
## 2023-09-16 21:28:12 Annotating text fragment 80211/100392
## 2023-09-16 21:28:12 Annotating text fragment 80221/100392
## 2023-09-16 21:28:12 Annotating text fragment 80231/100392
## 2023-09-16 21:28:12 Annotating text fragment 80241/100392
## 2023-09-16 21:28:13 Annotating text fragment 80251/100392
## 2023-09-16 21:28:13 Annotating text fragment 80261/100392
## 2023-09-16 21:28:13 Annotating text fragment 80271/100392
## 2023-09-16 21:28:13 Annotating text fragment 80281/100392
## 2023-09-16 21:28:13 Annotating text fragment 80291/100392
## 2023-09-16 21:28:13 Annotating text fragment 80301/100392
## 2023-09-16 21:28:13 Annotating text fragment 80311/100392
## 2023-09-16 21:28:13 Annotating text fragment 80321/100392
## 2023-09-16 21:28:14 Annotating text fragment 80331/100392
## 2023-09-16 21:28:14 Annotating text fragment 80341/100392
## 2023-09-16 21:28:14 Annotating text fragment 80351/100392
## 2023-09-16 21:28:14 Annotating text fragment 80361/100392
## 2023-09-16 21:28:14 Annotating text fragment 80371/100392
## 2023-09-16 21:28:14 Annotating text fragment 80381/100392
## 2023-09-16 21:28:14 Annotating text fragment 80391/100392
## 2023-09-16 21:28:14 Annotating text fragment 80401/100392
## 2023-09-16 21:28:14 Annotating text fragment 80411/100392
## 2023-09-16 21:28:14 Annotating text fragment 80421/100392
## 2023-09-16 21:28:14 Annotating text fragment 80431/100392
## 2023-09-16 21:28:15 Annotating text fragment 80441/100392
## 2023-09-16 21:28:15 Annotating text fragment 80451/100392
## 2023-09-16 21:28:15 Annotating text fragment 80461/100392
## 2023-09-16 21:28:15 Annotating text fragment 80471/100392
## 2023-09-16 21:28:15 Annotating text fragment 80481/100392
## 2023-09-16 21:28:15 Annotating text fragment 80491/100392
## 2023-09-16 21:28:15 Annotating text fragment 80501/100392
## 2023-09-16 21:28:15 Annotating text fragment 80511/100392
```

```
## 2023-09-16 21:28:16 Annotating text fragment 80521/100392
## 2023-09-16 21:28:16 Annotating text fragment 80531/100392
## 2023-09-16 21:28:16 Annotating text fragment 80541/100392
## 2023-09-16 21:28:16 Annotating text fragment 80551/100392
## 2023-09-16 21:28:16 Annotating text fragment 80561/100392
## 2023-09-16 21:28:16 Annotating text fragment 80571/100392
## 2023-09-16 21:28:16 Annotating text fragment 80581/100392
## 2023-09-16 21:28:16 Annotating text fragment 80591/100392
## 2023-09-16 21:28:16 Annotating text fragment 80601/100392
## 2023-09-16 21:28:16 Annotating text fragment 80611/100392
## 2023-09-16 21:28:17 Annotating text fragment 80621/100392
## 2023-09-16 21:28:17 Annotating text fragment 80631/100392
## 2023-09-16 21:28:17 Annotating text fragment 80641/100392
## 2023-09-16 21:28:17 Annotating text fragment 80651/100392
## 2023-09-16 21:28:17 Annotating text fragment 80661/100392
## 2023-09-16 21:28:17 Annotating text fragment 80671/100392
## 2023-09-16 21:28:17 Annotating text fragment 80681/100392
## 2023-09-16 21:28:17 Annotating text fragment 80691/100392
## 2023-09-16 21:28:17 Annotating text fragment 80701/100392
## 2023-09-16 21:28:17 Annotating text fragment 80711/100392
## 2023-09-16 21:28:17 Annotating text fragment 80721/100392
## 2023-09-16 21:28:18 Annotating text fragment 80731/100392
## 2023-09-16 21:28:18 Annotating text fragment 80741/100392
## 2023-09-16 21:28:18 Annotating text fragment 80751/100392
## 2023-09-16 21:28:18 Annotating text fragment 80761/100392
## 2023-09-16 21:28:18 Annotating text fragment 80771/100392
## 2023-09-16 21:28:18 Annotating text fragment 80781/100392
## 2023-09-16 21:28:18 Annotating text fragment 80791/100392
## 2023-09-16 21:28:18 Annotating text fragment 80801/100392
## 2023-09-16 21:28:18 Annotating text fragment 80811/100392
## 2023-09-16 21:28:19 Annotating text fragment 80821/100392
## 2023-09-16 21:28:19 Annotating text fragment 80831/100392
## 2023-09-16 21:28:19 Annotating text fragment 80841/100392
## 2023-09-16 21:28:19 Annotating text fragment 80851/100392
## 2023-09-16 21:28:19 Annotating text fragment 80861/100392
## 2023-09-16 21:28:19 Annotating text fragment 80871/100392
## 2023-09-16 21:28:19 Annotating text fragment 80881/100392
## 2023-09-16 21:28:19 Annotating text fragment 80891/100392
## 2023-09-16 21:28:19 Annotating text fragment 80901/100392
## 2023-09-16 21:28:19 Annotating text fragment 80911/100392
## 2023-09-16 21:28:20 Annotating text fragment 80921/100392
## 2023-09-16 21:28:20 Annotating text fragment 80931/100392
## 2023-09-16 21:28:20 Annotating text fragment 80941/100392
## 2023-09-16 21:28:20 Annotating text fragment 80951/100392
## 2023-09-16 21:28:20 Annotating text fragment 80961/100392
## 2023-09-16 21:28:20 Annotating text fragment 80971/100392
## 2023-09-16 21:28:20 Annotating text fragment 80981/100392
## 2023-09-16 21:28:20 Annotating text fragment 80991/100392
## 2023-09-16 21:28:20 Annotating text fragment 81001/100392
## 2023-09-16 21:28:20 Annotating text fragment 81011/100392
## 2023-09-16 21:28:21 Annotating text fragment 81021/100392
## 2023-09-16 21:28:21 Annotating text fragment 81031/100392
## 2023-09-16 21:28:21 Annotating text fragment 81041/100392
## 2023-09-16 21:28:21 Annotating text fragment 81051/100392
```

```
## 2023-09-16 21:28:21 Annotating text fragment 81061/100392
## 2023-09-16 21:28:21 Annotating text fragment 81071/100392
## 2023-09-16 21:28:21 Annotating text fragment 81081/100392
## 2023-09-16 21:28:21 Annotating text fragment 81091/100392
## 2023-09-16 21:28:21 Annotating text fragment 81101/100392
## 2023-09-16 21:28:21 Annotating text fragment 81111/100392
## 2023-09-16 21:28:21 Annotating text fragment 81121/100392
## 2023-09-16 21:28:22 Annotating text fragment 81131/100392
## 2023-09-16 21:28:22 Annotating text fragment 81141/100392
## 2023-09-16 21:28:22 Annotating text fragment 81151/100392
## 2023-09-16 21:28:22 Annotating text fragment 81161/100392
## 2023-09-16 21:28:22 Annotating text fragment 81171/100392
## 2023-09-16 21:28:22 Annotating text fragment 81181/100392
## 2023-09-16 21:28:22 Annotating text fragment 81191/100392
## 2023-09-16 21:28:22 Annotating text fragment 81201/100392
## 2023-09-16 21:28:22 Annotating text fragment 81211/100392
## 2023-09-16 21:28:23 Annotating text fragment 81221/100392
## 2023-09-16 21:28:23 Annotating text fragment 81231/100392
## 2023-09-16 21:28:23 Annotating text fragment 81241/100392
## 2023-09-16 21:28:23 Annotating text fragment 81251/100392
## 2023-09-16 21:28:23 Annotating text fragment 81261/100392
## 2023-09-16 21:28:23 Annotating text fragment 81271/100392
## 2023-09-16 21:28:23 Annotating text fragment 81281/100392
## 2023-09-16 21:28:23 Annotating text fragment 81291/100392
## 2023-09-16 21:28:23 Annotating text fragment 81301/100392
## 2023-09-16 21:28:23 Annotating text fragment 81311/100392
## 2023-09-16 21:28:24 Annotating text fragment 81321/100392
## 2023-09-16 21:28:24 Annotating text fragment 81331/100392
## 2023-09-16 21:28:24 Annotating text fragment 81341/100392
## 2023-09-16 21:28:24 Annotating text fragment 81351/100392
## 2023-09-16 21:28:24 Annotating text fragment 81361/100392
## 2023-09-16 21:28:24 Annotating text fragment 81371/100392
## 2023-09-16 21:28:24 Annotating text fragment 81381/100392
## 2023-09-16 21:28:24 Annotating text fragment 81391/100392
## 2023-09-16 21:28:24 Annotating text fragment 81401/100392
## 2023-09-16 21:28:24 Annotating text fragment 81411/100392
## 2023-09-16 21:28:24 Annotating text fragment 81421/100392
## 2023-09-16 21:28:24 Annotating text fragment 81431/100392
## 2023-09-16 21:28:24 Annotating text fragment 81441/100392
## 2023-09-16 21:28:25 Annotating text fragment 81451/100392
## 2023-09-16 21:28:25 Annotating text fragment 81461/100392
## 2023-09-16 21:28:25 Annotating text fragment 81471/100392
## 2023-09-16 21:28:25 Annotating text fragment 81481/100392
## 2023-09-16 21:28:25 Annotating text fragment 81491/100392
## 2023-09-16 21:28:25 Annotating text fragment 81501/100392
## 2023-09-16 21:28:25 Annotating text fragment 81511/100392
## 2023-09-16 21:28:25 Annotating text fragment 81521/100392
## 2023-09-16 21:28:25 Annotating text fragment 81531/100392
## 2023-09-16 21:28:26 Annotating text fragment 81541/100392
## 2023-09-16 21:28:26 Annotating text fragment 81551/100392
## 2023-09-16 21:28:26 Annotating text fragment 81561/100392
## 2023-09-16 21:28:26 Annotating text fragment 81571/100392
## 2023-09-16 21:28:26 Annotating text fragment 81581/100392
## 2023-09-16 21:28:26 Annotating text fragment 81591/100392
```

```
## 2023-09-16 21:28:26 Annotating text fragment 81601/100392
## 2023-09-16 21:28:26 Annotating text fragment 81611/100392
## 2023-09-16 21:28:26 Annotating text fragment 81621/100392
## 2023-09-16 21:28:26 Annotating text fragment 81631/100392
## 2023-09-16 21:28:27 Annotating text fragment 81641/100392
## 2023-09-16 21:28:27 Annotating text fragment 81651/100392
## 2023-09-16 21:28:27 Annotating text fragment 81661/100392
## 2023-09-16 21:28:27 Annotating text fragment 81671/100392
## 2023-09-16 21:28:27 Annotating text fragment 81681/100392
## 2023-09-16 21:28:27 Annotating text fragment 81691/100392
## 2023-09-16 21:28:27 Annotating text fragment 81701/100392
## 2023-09-16 21:28:27 Annotating text fragment 81711/100392
## 2023-09-16 21:28:27 Annotating text fragment 81721/100392
## 2023-09-16 21:28:27 Annotating text fragment 81731/100392
## 2023-09-16 21:28:28 Annotating text fragment 81741/100392
## 2023-09-16 21:28:28 Annotating text fragment 81751/100392
## 2023-09-16 21:28:28 Annotating text fragment 81761/100392
## 2023-09-16 21:28:28 Annotating text fragment 81771/100392
## 2023-09-16 21:28:28 Annotating text fragment 81781/100392
## 2023-09-16 21:28:28 Annotating text fragment 81791/100392
## 2023-09-16 21:28:28 Annotating text fragment 81801/100392
## 2023-09-16 21:28:28 Annotating text fragment 81811/100392
## 2023-09-16 21:28:28 Annotating text fragment 81821/100392
## 2023-09-16 21:28:29 Annotating text fragment 81831/100392
## 2023-09-16 21:28:29 Annotating text fragment 81841/100392
## 2023-09-16 21:28:29 Annotating text fragment 81851/100392
## 2023-09-16 21:28:29 Annotating text fragment 81861/100392
## 2023-09-16 21:28:29 Annotating text fragment 81871/100392
## 2023-09-16 21:28:29 Annotating text fragment 81881/100392
## 2023-09-16 21:28:29 Annotating text fragment 81891/100392
## 2023-09-16 21:28:29 Annotating text fragment 81901/100392
## 2023-09-16 21:28:29 Annotating text fragment 81911/100392
## 2023-09-16 21:28:29 Annotating text fragment 81921/100392
## 2023-09-16 21:28:30 Annotating text fragment 81931/100392
## 2023-09-16 21:28:30 Annotating text fragment 81941/100392
## 2023-09-16 21:28:30 Annotating text fragment 81951/100392
## 2023-09-16 21:28:30 Annotating text fragment 81961/100392
## 2023-09-16 21:28:30 Annotating text fragment 81971/100392
## 2023-09-16 21:28:30 Annotating text fragment 81981/100392
## 2023-09-16 21:28:30 Annotating text fragment 81991/100392
## 2023-09-16 21:28:30 Annotating text fragment 82001/100392
## 2023-09-16 21:28:30 Annotating text fragment 82011/100392
## 2023-09-16 21:28:31 Annotating text fragment 82021/100392
## 2023-09-16 21:28:31 Annotating text fragment 82031/100392
## 2023-09-16 21:28:31 Annotating text fragment 82041/100392
## 2023-09-16 21:28:31 Annotating text fragment 82051/100392
## 2023-09-16 21:28:31 Annotating text fragment 82061/100392
## 2023-09-16 21:28:31 Annotating text fragment 82071/100392
## 2023-09-16 21:28:31 Annotating text fragment 82081/100392
## 2023-09-16 21:28:31 Annotating text fragment 82091/100392
## 2023-09-16 21:28:31 Annotating text fragment 82101/100392
## 2023-09-16 21:28:31 Annotating text fragment 82111/100392
## 2023-09-16 21:28:31 Annotating text fragment 82121/100392
## 2023-09-16 21:28:32 Annotating text fragment 82131/100392
```

```
## 2023-09-16 21:28:32 Annotating text fragment 82141/100392
## 2023-09-16 21:28:32 Annotating text fragment 82151/100392
## 2023-09-16 21:28:32 Annotating text fragment 82161/100392
## 2023-09-16 21:28:32 Annotating text fragment 82171/100392
## 2023-09-16 21:28:32 Annotating text fragment 82181/100392
## 2023-09-16 21:28:32 Annotating text fragment 82191/100392
## 2023-09-16 21:28:32 Annotating text fragment 82201/100392
## 2023-09-16 21:28:32 Annotating text fragment 82211/100392
## 2023-09-16 21:28:33 Annotating text fragment 82221/100392
## 2023-09-16 21:28:33 Annotating text fragment 82231/100392
## 2023-09-16 21:28:33 Annotating text fragment 82241/100392
## 2023-09-16 21:28:33 Annotating text fragment 82251/100392
## 2023-09-16 21:28:33 Annotating text fragment 82261/100392
## 2023-09-16 21:28:33 Annotating text fragment 82271/100392
## 2023-09-16 21:28:33 Annotating text fragment 82281/100392
## 2023-09-16 21:28:33 Annotating text fragment 82291/100392
## 2023-09-16 21:28:33 Annotating text fragment 82301/100392
## 2023-09-16 21:28:33 Annotating text fragment 82311/100392
## 2023-09-16 21:28:33 Annotating text fragment 82321/100392
## 2023-09-16 21:28:34 Annotating text fragment 82331/100392
## 2023-09-16 21:28:34 Annotating text fragment 82341/100392
## 2023-09-16 21:28:34 Annotating text fragment 82351/100392
## 2023-09-16 21:28:34 Annotating text fragment 82361/100392
## 2023-09-16 21:28:34 Annotating text fragment 82371/100392
## 2023-09-16 21:28:34 Annotating text fragment 82381/100392
## 2023-09-16 21:28:34 Annotating text fragment 82391/100392
## 2023-09-16 21:28:34 Annotating text fragment 82401/100392
## 2023-09-16 21:28:34 Annotating text fragment 82411/100392
## 2023-09-16 21:28:34 Annotating text fragment 82421/100392
## 2023-09-16 21:28:34 Annotating text fragment 82431/100392
## 2023-09-16 21:28:34 Annotating text fragment 82441/100392
## 2023-09-16 21:28:34 Annotating text fragment 82451/100392
## 2023-09-16 21:28:35 Annotating text fragment 82461/100392
## 2023-09-16 21:28:35 Annotating text fragment 82471/100392
## 2023-09-16 21:28:35 Annotating text fragment 82481/100392
## 2023-09-16 21:28:35 Annotating text fragment 82491/100392
## 2023-09-16 21:28:35 Annotating text fragment 82501/100392
## 2023-09-16 21:28:35 Annotating text fragment 82511/100392
## 2023-09-16 21:28:35 Annotating text fragment 82521/100392
## 2023-09-16 21:28:35 Annotating text fragment 82531/100392
## 2023-09-16 21:28:35 Annotating text fragment 82541/100392
## 2023-09-16 21:28:35 Annotating text fragment 82551/100392
## 2023-09-16 21:28:35 Annotating text fragment 82561/100392
## 2023-09-16 21:28:36 Annotating text fragment 82571/100392
## 2023-09-16 21:28:36 Annotating text fragment 82581/100392
## 2023-09-16 21:28:36 Annotating text fragment 82591/100392
## 2023-09-16 21:28:36 Annotating text fragment 82601/100392
## 2023-09-16 21:28:36 Annotating text fragment 82611/100392
## 2023-09-16 21:28:36 Annotating text fragment 82621/100392
## 2023-09-16 21:28:36 Annotating text fragment 82631/100392
## 2023-09-16 21:28:36 Annotating text fragment 82641/100392
## 2023-09-16 21:28:36 Annotating text fragment 82651/100392
## 2023-09-16 21:28:37 Annotating text fragment 82661/100392
## 2023-09-16 21:28:37 Annotating text fragment 82671/100392
```

```
## 2023-09-16 21:28:37 Annotating text fragment 82681/100392
## 2023-09-16 21:28:37 Annotating text fragment 82691/100392
## 2023-09-16 21:28:37 Annotating text fragment 82701/100392
## 2023-09-16 21:28:37 Annotating text fragment 82711/100392
## 2023-09-16 21:28:37 Annotating text fragment 82721/100392
## 2023-09-16 21:28:37 Annotating text fragment 82731/100392
## 2023-09-16 21:28:37 Annotating text fragment 82741/100392
## 2023-09-16 21:28:37 Annotating text fragment 82751/100392
## 2023-09-16 21:28:37 Annotating text fragment 82761/100392
## 2023-09-16 21:28:38 Annotating text fragment 82771/100392
## 2023-09-16 21:28:38 Annotating text fragment 82781/100392
## 2023-09-16 21:28:38 Annotating text fragment 82791/100392
## 2023-09-16 21:28:38 Annotating text fragment 82801/100392
## 2023-09-16 21:28:38 Annotating text fragment 82811/100392
## 2023-09-16 21:28:38 Annotating text fragment 82821/100392
## 2023-09-16 21:28:38 Annotating text fragment 82831/100392
## 2023-09-16 21:28:38 Annotating text fragment 82841/100392
## 2023-09-16 21:28:38 Annotating text fragment 82851/100392
## 2023-09-16 21:28:39 Annotating text fragment 82861/100392
## 2023-09-16 21:28:39 Annotating text fragment 82871/100392
## 2023-09-16 21:28:39 Annotating text fragment 82881/100392
## 2023-09-16 21:28:39 Annotating text fragment 82891/100392
## 2023-09-16 21:28:39 Annotating text fragment 82901/100392
## 2023-09-16 21:28:39 Annotating text fragment 82911/100392
## 2023-09-16 21:28:39 Annotating text fragment 82921/100392
## 2023-09-16 21:28:39 Annotating text fragment 82931/100392
## 2023-09-16 21:28:39 Annotating text fragment 82941/100392
## 2023-09-16 21:28:39 Annotating text fragment 82951/100392
## 2023-09-16 21:28:40 Annotating text fragment 82961/100392
## 2023-09-16 21:28:40 Annotating text fragment 82971/100392
## 2023-09-16 21:28:40 Annotating text fragment 82981/100392
## 2023-09-16 21:28:40 Annotating text fragment 82991/100392
## 2023-09-16 21:28:40 Annotating text fragment 83001/100392
## 2023-09-16 21:28:40 Annotating text fragment 83011/100392
## 2023-09-16 21:28:40 Annotating text fragment 83021/100392
## 2023-09-16 21:28:40 Annotating text fragment 83031/100392
## 2023-09-16 21:28:40 Annotating text fragment 83041/100392
## 2023-09-16 21:28:40 Annotating text fragment 83051/100392
## 2023-09-16 21:28:41 Annotating text fragment 83061/100392
## 2023-09-16 21:28:41 Annotating text fragment 83071/100392
## 2023-09-16 21:28:41 Annotating text fragment 83081/100392
## 2023-09-16 21:28:41 Annotating text fragment 83091/100392
## 2023-09-16 21:28:41 Annotating text fragment 83101/100392
## 2023-09-16 21:28:41 Annotating text fragment 83111/100392
## 2023-09-16 21:28:41 Annotating text fragment 83121/100392
## 2023-09-16 21:28:41 Annotating text fragment 83131/100392
## 2023-09-16 21:28:41 Annotating text fragment 83141/100392
## 2023-09-16 21:28:41 Annotating text fragment 83151/100392
## 2023-09-16 21:28:42 Annotating text fragment 83161/100392
## 2023-09-16 21:28:42 Annotating text fragment 83171/100392
## 2023-09-16 21:28:42 Annotating text fragment 83181/100392
## 2023-09-16 21:28:42 Annotating text fragment 83191/100392
## 2023-09-16 21:28:42 Annotating text fragment 83201/100392
## 2023-09-16 21:28:42 Annotating text fragment 83211/100392
```

```
## 2023-09-16 21:28:42 Annotating text fragment 83221/100392
## 2023-09-16 21:28:42 Annotating text fragment 83231/100392
## 2023-09-16 21:28:42 Annotating text fragment 83241/100392
## 2023-09-16 21:28:43 Annotating text fragment 83251/100392
## 2023-09-16 21:28:43 Annotating text fragment 83261/100392
## 2023-09-16 21:28:43 Annotating text fragment 83271/100392
## 2023-09-16 21:28:43 Annotating text fragment 83281/100392
## 2023-09-16 21:28:43 Annotating text fragment 83291/100392
## 2023-09-16 21:28:43 Annotating text fragment 83301/100392
## 2023-09-16 21:28:43 Annotating text fragment 83311/100392
## 2023-09-16 21:28:43 Annotating text fragment 83321/100392
## 2023-09-16 21:28:43 Annotating text fragment 83331/100392
## 2023-09-16 21:28:43 Annotating text fragment 83341/100392
## 2023-09-16 21:28:43 Annotating text fragment 83351/100392
## 2023-09-16 21:28:43 Annotating text fragment 83361/100392
## 2023-09-16 21:28:44 Annotating text fragment 83371/100392
## 2023-09-16 21:28:44 Annotating text fragment 83381/100392
## 2023-09-16 21:28:44 Annotating text fragment 83391/100392
## 2023-09-16 21:28:44 Annotating text fragment 83401/100392
## 2023-09-16 21:28:44 Annotating text fragment 83411/100392
## 2023-09-16 21:28:44 Annotating text fragment 83421/100392
## 2023-09-16 21:28:44 Annotating text fragment 83431/100392
## 2023-09-16 21:28:44 Annotating text fragment 83441/100392
## 2023-09-16 21:28:44 Annotating text fragment 83451/100392
## 2023-09-16 21:28:44 Annotating text fragment 83461/100392
## 2023-09-16 21:28:44 Annotating text fragment 83471/100392
## 2023-09-16 21:28:44 Annotating text fragment 83481/100392
## 2023-09-16 21:28:45 Annotating text fragment 83491/100392
## 2023-09-16 21:28:45 Annotating text fragment 83501/100392
## 2023-09-16 21:28:45 Annotating text fragment 83511/100392
## 2023-09-16 21:28:45 Annotating text fragment 83521/100392
## 2023-09-16 21:28:45 Annotating text fragment 83531/100392
## 2023-09-16 21:28:45 Annotating text fragment 83541/100392
## 2023-09-16 21:28:45 Annotating text fragment 83551/100392
## 2023-09-16 21:28:45 Annotating text fragment 83561/100392
## 2023-09-16 21:28:46 Annotating text fragment 83571/100392
## 2023-09-16 21:28:46 Annotating text fragment 83581/100392
## 2023-09-16 21:28:46 Annotating text fragment 83591/100392
## 2023-09-16 21:28:46 Annotating text fragment 83601/100392
## 2023-09-16 21:28:46 Annotating text fragment 83611/100392
## 2023-09-16 21:28:46 Annotating text fragment 83621/100392
## 2023-09-16 21:28:46 Annotating text fragment 83631/100392
## 2023-09-16 21:28:46 Annotating text fragment 83641/100392
## 2023-09-16 21:28:46 Annotating text fragment 83651/100392
## 2023-09-16 21:28:46 Annotating text fragment 83661/100392
## 2023-09-16 21:28:46 Annotating text fragment 83671/100392
## 2023-09-16 21:28:47 Annotating text fragment 83681/100392
## 2023-09-16 21:28:47 Annotating text fragment 83691/100392
## 2023-09-16 21:28:47 Annotating text fragment 83701/100392
## 2023-09-16 21:28:47 Annotating text fragment 83711/100392
## 2023-09-16 21:28:47 Annotating text fragment 83721/100392
## 2023-09-16 21:28:47 Annotating text fragment 83731/100392
## 2023-09-16 21:28:47 Annotating text fragment 83741/100392
## 2023-09-16 21:28:47 Annotating text fragment 83751/100392
```

```
## 2023-09-16 21:28:47 Annotating text fragment 83761/100392
## 2023-09-16 21:28:47 Annotating text fragment 83771/100392
## 2023-09-16 21:28:47 Annotating text fragment 83781/100392
## 2023-09-16 21:28:48 Annotating text fragment 83791/100392
## 2023-09-16 21:28:48 Annotating text fragment 83801/100392
## 2023-09-16 21:28:48 Annotating text fragment 83811/100392
## 2023-09-16 21:28:48 Annotating text fragment 83821/100392
## 2023-09-16 21:28:48 Annotating text fragment 83831/100392
## 2023-09-16 21:28:48 Annotating text fragment 83841/100392
## 2023-09-16 21:28:48 Annotating text fragment 83851/100392
## 2023-09-16 21:28:48 Annotating text fragment 83861/100392
## 2023-09-16 21:28:48 Annotating text fragment 83871/100392
## 2023-09-16 21:28:48 Annotating text fragment 83881/100392
## 2023-09-16 21:28:48 Annotating text fragment 83891/100392
## 2023-09-16 21:28:49 Annotating text fragment 83901/100392
## 2023-09-16 21:28:49 Annotating text fragment 83911/100392
## 2023-09-16 21:28:49 Annotating text fragment 83921/100392
## 2023-09-16 21:28:49 Annotating text fragment 83931/100392
## 2023-09-16 21:28:49 Annotating text fragment 83941/100392
## 2023-09-16 21:28:49 Annotating text fragment 83951/100392
## 2023-09-16 21:28:49 Annotating text fragment 83961/100392
## 2023-09-16 21:28:49 Annotating text fragment 83971/100392
## 2023-09-16 21:28:49 Annotating text fragment 83981/100392
## 2023-09-16 21:28:49 Annotating text fragment 83991/100392
## 2023-09-16 21:28:50 Annotating text fragment 84001/100392
## 2023-09-16 21:28:50 Annotating text fragment 84011/100392
## 2023-09-16 21:28:50 Annotating text fragment 84021/100392
## 2023-09-16 21:28:50 Annotating text fragment 84031/100392
## 2023-09-16 21:28:50 Annotating text fragment 84041/100392
## 2023-09-16 21:28:50 Annotating text fragment 84051/100392
## 2023-09-16 21:28:50 Annotating text fragment 84061/100392
## 2023-09-16 21:28:50 Annotating text fragment 84071/100392
## 2023-09-16 21:28:50 Annotating text fragment 84081/100392
## 2023-09-16 21:28:50 Annotating text fragment 84091/100392
## 2023-09-16 21:28:50 Annotating text fragment 84101/100392
## 2023-09-16 21:28:50 Annotating text fragment 84111/100392
## 2023-09-16 21:28:51 Annotating text fragment 84121/100392
## 2023-09-16 21:28:51 Annotating text fragment 84131/100392
## 2023-09-16 21:28:51 Annotating text fragment 84141/100392
## 2023-09-16 21:28:51 Annotating text fragment 84151/100392
## 2023-09-16 21:28:51 Annotating text fragment 84161/100392
## 2023-09-16 21:28:51 Annotating text fragment 84171/100392
## 2023-09-16 21:28:51 Annotating text fragment 84181/100392
## 2023-09-16 21:28:51 Annotating text fragment 84191/100392
## 2023-09-16 21:28:51 Annotating text fragment 84201/100392
## 2023-09-16 21:28:52 Annotating text fragment 84211/100392
## 2023-09-16 21:28:52 Annotating text fragment 84221/100392
## 2023-09-16 21:28:52 Annotating text fragment 84231/100392
## 2023-09-16 21:28:52 Annotating text fragment 84241/100392
## 2023-09-16 21:28:52 Annotating text fragment 84251/100392
## 2023-09-16 21:28:52 Annotating text fragment 84261/100392
## 2023-09-16 21:28:52 Annotating text fragment 84271/100392
## 2023-09-16 21:28:52 Annotating text fragment 84281/100392
## 2023-09-16 21:28:52 Annotating text fragment 84291/100392
```

```
## 2023-09-16 21:28:52 Annotating text fragment 84301/100392
## 2023-09-16 21:28:52 Annotating text fragment 84311/100392
## 2023-09-16 21:28:52 Annotating text fragment 84321/100392
## 2023-09-16 21:28:53 Annotating text fragment 84331/100392
## 2023-09-16 21:28:53 Annotating text fragment 84341/100392
## 2023-09-16 21:28:53 Annotating text fragment 84351/100392
## 2023-09-16 21:28:53 Annotating text fragment 84361/100392
## 2023-09-16 21:28:53 Annotating text fragment 84371/100392
## 2023-09-16 21:28:53 Annotating text fragment 84381/100392
## 2023-09-16 21:28:53 Annotating text fragment 84391/100392
## 2023-09-16 21:28:53 Annotating text fragment 84401/100392
## 2023-09-16 21:28:53 Annotating text fragment 84411/100392
## 2023-09-16 21:28:53 Annotating text fragment 84421/100392
## 2023-09-16 21:28:53 Annotating text fragment 84431/100392
## 2023-09-16 21:28:53 Annotating text fragment 84441/100392
## 2023-09-16 21:28:54 Annotating text fragment 84451/100392
## 2023-09-16 21:28:54 Annotating text fragment 84461/100392
## 2023-09-16 21:28:54 Annotating text fragment 84471/100392
## 2023-09-16 21:28:54 Annotating text fragment 84481/100392
## 2023-09-16 21:28:54 Annotating text fragment 84491/100392
## 2023-09-16 21:28:54 Annotating text fragment 84501/100392
## 2023-09-16 21:28:54 Annotating text fragment 84511/100392
## 2023-09-16 21:28:54 Annotating text fragment 84521/100392
## 2023-09-16 21:28:54 Annotating text fragment 84531/100392
## 2023-09-16 21:28:54 Annotating text fragment 84541/100392
## 2023-09-16 21:28:54 Annotating text fragment 84551/100392
## 2023-09-16 21:28:54 Annotating text fragment 84561/100392
## 2023-09-16 21:28:55 Annotating text fragment 84571/100392
## 2023-09-16 21:28:55 Annotating text fragment 84581/100392
## 2023-09-16 21:28:55 Annotating text fragment 84591/100392
## 2023-09-16 21:28:55 Annotating text fragment 84601/100392
## 2023-09-16 21:28:55 Annotating text fragment 84611/100392
## 2023-09-16 21:28:55 Annotating text fragment 84621/100392
## 2023-09-16 21:28:55 Annotating text fragment 84631/100392
## 2023-09-16 21:28:55 Annotating text fragment 84641/100392
## 2023-09-16 21:28:55 Annotating text fragment 84651/100392
## 2023-09-16 21:28:55 Annotating text fragment 84661/100392
## 2023-09-16 21:28:55 Annotating text fragment 84671/100392
## 2023-09-16 21:28:56 Annotating text fragment 84681/100392
## 2023-09-16 21:28:56 Annotating text fragment 84691/100392
## 2023-09-16 21:28:56 Annotating text fragment 84701/100392
## 2023-09-16 21:28:56 Annotating text fragment 84711/100392
## 2023-09-16 21:28:56 Annotating text fragment 84721/100392
## 2023-09-16 21:28:56 Annotating text fragment 84731/100392
## 2023-09-16 21:28:56 Annotating text fragment 84741/100392
## 2023-09-16 21:28:56 Annotating text fragment 84751/100392
## 2023-09-16 21:28:56 Annotating text fragment 84761/100392
## 2023-09-16 21:28:56 Annotating text fragment 84771/100392
## 2023-09-16 21:28:56 Annotating text fragment 84781/100392
## 2023-09-16 21:28:56 Annotating text fragment 84791/100392
## 2023-09-16 21:28:57 Annotating text fragment 84801/100392
## 2023-09-16 21:28:57 Annotating text fragment 84811/100392
## 2023-09-16 21:28:57 Annotating text fragment 84821/100392
## 2023-09-16 21:28:57 Annotating text fragment 84831/100392
```

```
## 2023-09-16 21:28:57 Annotating text fragment 84841/100392
## 2023-09-16 21:28:57 Annotating text fragment 84851/100392
## 2023-09-16 21:28:57 Annotating text fragment 84861/100392
## 2023-09-16 21:28:57 Annotating text fragment 84871/100392
## 2023-09-16 21:28:58 Annotating text fragment 84881/100392
## 2023-09-16 21:28:58 Annotating text fragment 84891/100392
## 2023-09-16 21:28:58 Annotating text fragment 84901/100392
## 2023-09-16 21:28:58 Annotating text fragment 84911/100392
## 2023-09-16 21:28:58 Annotating text fragment 84921/100392
## 2023-09-16 21:28:58 Annotating text fragment 84931/100392
## 2023-09-16 21:28:58 Annotating text fragment 84941/100392
## 2023-09-16 21:28:58 Annotating text fragment 84951/100392
## 2023-09-16 21:28:58 Annotating text fragment 84961/100392
## 2023-09-16 21:28:59 Annotating text fragment 84971/100392
## 2023-09-16 21:28:59 Annotating text fragment 84981/100392
## 2023-09-16 21:28:59 Annotating text fragment 84991/100392
## 2023-09-16 21:28:59 Annotating text fragment 85001/100392
## 2023-09-16 21:28:59 Annotating text fragment 85011/100392
## 2023-09-16 21:28:59 Annotating text fragment 85021/100392
## 2023-09-16 21:28:59 Annotating text fragment 85031/100392
## 2023-09-16 21:28:59 Annotating text fragment 85041/100392
## 2023-09-16 21:28:59 Annotating text fragment 85051/100392
## 2023-09-16 21:28:59 Annotating text fragment 85061/100392
## 2023-09-16 21:29:00 Annotating text fragment 85071/100392
## 2023-09-16 21:29:00 Annotating text fragment 85081/100392
## 2023-09-16 21:29:00 Annotating text fragment 85091/100392
## 2023-09-16 21:29:00 Annotating text fragment 85101/100392
## 2023-09-16 21:29:00 Annotating text fragment 85111/100392
## 2023-09-16 21:29:00 Annotating text fragment 85121/100392
## 2023-09-16 21:29:00 Annotating text fragment 85131/100392
## 2023-09-16 21:29:00 Annotating text fragment 85141/100392
## 2023-09-16 21:29:00 Annotating text fragment 85151/100392
## 2023-09-16 21:29:00 Annotating text fragment 85161/100392
## 2023-09-16 21:29:00 Annotating text fragment 85171/100392
## 2023-09-16 21:29:01 Annotating text fragment 85181/100392
## 2023-09-16 21:29:01 Annotating text fragment 85191/100392
## 2023-09-16 21:29:01 Annotating text fragment 85201/100392
## 2023-09-16 21:29:01 Annotating text fragment 85211/100392
## 2023-09-16 21:29:01 Annotating text fragment 85221/100392
## 2023-09-16 21:29:01 Annotating text fragment 85231/100392
## 2023-09-16 21:29:01 Annotating text fragment 85241/100392
## 2023-09-16 21:29:01 Annotating text fragment 85251/100392
## 2023-09-16 21:29:02 Annotating text fragment 85261/100392
## 2023-09-16 21:29:02 Annotating text fragment 85271/100392
## 2023-09-16 21:29:02 Annotating text fragment 85281/100392
## 2023-09-16 21:29:02 Annotating text fragment 85291/100392
## 2023-09-16 21:29:02 Annotating text fragment 85301/100392
## 2023-09-16 21:29:02 Annotating text fragment 85311/100392
## 2023-09-16 21:29:02 Annotating text fragment 85321/100392
## 2023-09-16 21:29:02 Annotating text fragment 85331/100392
## 2023-09-16 21:29:02 Annotating text fragment 85341/100392
## 2023-09-16 21:29:02 Annotating text fragment 85351/100392
## 2023-09-16 21:29:02 Annotating text fragment 85361/100392
## 2023-09-16 21:29:02 Annotating text fragment 85371/100392
```

```
## 2023-09-16 21:29:03 Annotating text fragment 85381/100392
## 2023-09-16 21:29:03 Annotating text fragment 85391/100392
## 2023-09-16 21:29:03 Annotating text fragment 85401/100392
## 2023-09-16 21:29:03 Annotating text fragment 85411/100392
## 2023-09-16 21:29:03 Annotating text fragment 85421/100392
## 2023-09-16 21:29:03 Annotating text fragment 85431/100392
## 2023-09-16 21:29:03 Annotating text fragment 85441/100392
## 2023-09-16 21:29:03 Annotating text fragment 85451/100392
## 2023-09-16 21:29:03 Annotating text fragment 85461/100392
## 2023-09-16 21:29:03 Annotating text fragment 85471/100392
## 2023-09-16 21:29:03 Annotating text fragment 85481/100392
## 2023-09-16 21:29:03 Annotating text fragment 85491/100392
## 2023-09-16 21:29:03 Annotating text fragment 85501/100392
## 2023-09-16 21:29:03 Annotating text fragment 85511/100392
## 2023-09-16 21:29:04 Annotating text fragment 85521/100392
## 2023-09-16 21:29:04 Annotating text fragment 85531/100392
## 2023-09-16 21:29:04 Annotating text fragment 85541/100392
## 2023-09-16 21:29:04 Annotating text fragment 85551/100392
## 2023-09-16 21:29:04 Annotating text fragment 85561/100392
## 2023-09-16 21:29:04 Annotating text fragment 85571/100392
## 2023-09-16 21:29:04 Annotating text fragment 85581/100392
## 2023-09-16 21:29:04 Annotating text fragment 85591/100392
## 2023-09-16 21:29:04 Annotating text fragment 85601/100392
## 2023-09-16 21:29:04 Annotating text fragment 85611/100392
## 2023-09-16 21:29:04 Annotating text fragment 85621/100392
## 2023-09-16 21:29:04 Annotating text fragment 85631/100392
## 2023-09-16 21:29:05 Annotating text fragment 85641/100392
## 2023-09-16 21:29:05 Annotating text fragment 85651/100392
## 2023-09-16 21:29:05 Annotating text fragment 85661/100392
## 2023-09-16 21:29:05 Annotating text fragment 85671/100392
## 2023-09-16 21:29:05 Annotating text fragment 85681/100392
## 2023-09-16 21:29:05 Annotating text fragment 85691/100392
## 2023-09-16 21:29:05 Annotating text fragment 85701/100392
## 2023-09-16 21:29:05 Annotating text fragment 85711/100392
## 2023-09-16 21:29:05 Annotating text fragment 85721/100392
## 2023-09-16 21:29:05 Annotating text fragment 85731/100392
## 2023-09-16 21:29:05 Annotating text fragment 85741/100392
## 2023-09-16 21:29:05 Annotating text fragment 85751/100392
## 2023-09-16 21:29:06 Annotating text fragment 85761/100392
## 2023-09-16 21:29:06 Annotating text fragment 85771/100392
## 2023-09-16 21:29:06 Annotating text fragment 85781/100392
## 2023-09-16 21:29:06 Annotating text fragment 85791/100392
## 2023-09-16 21:29:06 Annotating text fragment 85801/100392
## 2023-09-16 21:29:06 Annotating text fragment 85811/100392
## 2023-09-16 21:29:06 Annotating text fragment 85821/100392
## 2023-09-16 21:29:06 Annotating text fragment 85831/100392
## 2023-09-16 21:29:06 Annotating text fragment 85841/100392
## 2023-09-16 21:29:06 Annotating text fragment 85851/100392
## 2023-09-16 21:29:07 Annotating text fragment 85861/100392
## 2023-09-16 21:29:07 Annotating text fragment 85871/100392
## 2023-09-16 21:29:07 Annotating text fragment 85881/100392
## 2023-09-16 21:29:07 Annotating text fragment 85891/100392
## 2023-09-16 21:29:07 Annotating text fragment 85901/100392
## 2023-09-16 21:29:07 Annotating text fragment 85911/100392
```

```
## 2023-09-16 21:29:07 Annotating text fragment 85921/100392
## 2023-09-16 21:29:07 Annotating text fragment 85931/100392
## 2023-09-16 21:29:07 Annotating text fragment 85941/100392
## 2023-09-16 21:29:07 Annotating text fragment 85951/100392
## 2023-09-16 21:29:08 Annotating text fragment 85961/100392
## 2023-09-16 21:29:08 Annotating text fragment 85971/100392
## 2023-09-16 21:29:08 Annotating text fragment 85981/100392
## 2023-09-16 21:29:08 Annotating text fragment 85991/100392
## 2023-09-16 21:29:08 Annotating text fragment 86001/100392
## 2023-09-16 21:29:08 Annotating text fragment 86011/100392
## 2023-09-16 21:29:08 Annotating text fragment 86021/100392
## 2023-09-16 21:29:09 Annotating text fragment 86031/100392
## 2023-09-16 21:29:09 Annotating text fragment 86041/100392
## 2023-09-16 21:29:09 Annotating text fragment 86051/100392
## 2023-09-16 21:29:09 Annotating text fragment 86061/100392
## 2023-09-16 21:29:09 Annotating text fragment 86071/100392
## 2023-09-16 21:29:09 Annotating text fragment 86081/100392
## 2023-09-16 21:29:09 Annotating text fragment 86091/100392
## 2023-09-16 21:29:09 Annotating text fragment 86101/100392
## 2023-09-16 21:29:09 Annotating text fragment 86111/100392
## 2023-09-16 21:29:09 Annotating text fragment 86121/100392
## 2023-09-16 21:29:09 Annotating text fragment 86131/100392
## 2023-09-16 21:29:10 Annotating text fragment 86141/100392
## 2023-09-16 21:29:10 Annotating text fragment 86151/100392
## 2023-09-16 21:29:10 Annotating text fragment 86161/100392
## 2023-09-16 21:29:10 Annotating text fragment 86171/100392
## 2023-09-16 21:29:10 Annotating text fragment 86181/100392
## 2023-09-16 21:29:10 Annotating text fragment 86191/100392
## 2023-09-16 21:29:10 Annotating text fragment 86201/100392
## 2023-09-16 21:29:10 Annotating text fragment 86211/100392
## 2023-09-16 21:29:10 Annotating text fragment 86221/100392
## 2023-09-16 21:29:10 Annotating text fragment 86231/100392
## 2023-09-16 21:29:11 Annotating text fragment 86241/100392
## 2023-09-16 21:29:11 Annotating text fragment 86251/100392
## 2023-09-16 21:29:11 Annotating text fragment 86261/100392
## 2023-09-16 21:29:11 Annotating text fragment 86271/100392
## 2023-09-16 21:29:11 Annotating text fragment 86281/100392
## 2023-09-16 21:29:11 Annotating text fragment 86291/100392
## 2023-09-16 21:29:11 Annotating text fragment 86301/100392
## 2023-09-16 21:29:11 Annotating text fragment 86311/100392
## 2023-09-16 21:29:11 Annotating text fragment 86321/100392
## 2023-09-16 21:29:11 Annotating text fragment 86331/100392
## 2023-09-16 21:29:11 Annotating text fragment 86341/100392
## 2023-09-16 21:29:12 Annotating text fragment 86351/100392
## 2023-09-16 21:29:12 Annotating text fragment 86361/100392
## 2023-09-16 21:29:12 Annotating text fragment 86371/100392
## 2023-09-16 21:29:12 Annotating text fragment 86381/100392
## 2023-09-16 21:29:12 Annotating text fragment 86391/100392
## 2023-09-16 21:29:12 Annotating text fragment 86401/100392
## 2023-09-16 21:29:12 Annotating text fragment 86411/100392
## 2023-09-16 21:29:12 Annotating text fragment 86421/100392
## 2023-09-16 21:29:12 Annotating text fragment 86431/100392
## 2023-09-16 21:29:12 Annotating text fragment 86441/100392
## 2023-09-16 21:29:12 Annotating text fragment 86451/100392
```

```
## 2023-09-16 21:29:12 Annotating text fragment 86461/100392
## 2023-09-16 21:29:13 Annotating text fragment 86471/100392
## 2023-09-16 21:29:13 Annotating text fragment 86481/100392
## 2023-09-16 21:29:13 Annotating text fragment 86491/100392
## 2023-09-16 21:29:13 Annotating text fragment 86501/100392
## 2023-09-16 21:29:13 Annotating text fragment 86511/100392
## 2023-09-16 21:29:13 Annotating text fragment 86521/100392
## 2023-09-16 21:29:13 Annotating text fragment 86531/100392
## 2023-09-16 21:29:13 Annotating text fragment 86541/100392
## 2023-09-16 21:29:13 Annotating text fragment 86551/100392
## 2023-09-16 21:29:13 Annotating text fragment 86561/100392
## 2023-09-16 21:29:13 Annotating text fragment 86571/100392
## 2023-09-16 21:29:13 Annotating text fragment 86581/100392
## 2023-09-16 21:29:14 Annotating text fragment 86591/100392
## 2023-09-16 21:29:14 Annotating text fragment 86601/100392
## 2023-09-16 21:29:14 Annotating text fragment 86611/100392
## 2023-09-16 21:29:14 Annotating text fragment 86621/100392
## 2023-09-16 21:29:14 Annotating text fragment 86631/100392
## 2023-09-16 21:29:14 Annotating text fragment 86641/100392
## 2023-09-16 21:29:14 Annotating text fragment 86651/100392
## 2023-09-16 21:29:14 Annotating text fragment 86661/100392
## 2023-09-16 21:29:14 Annotating text fragment 86671/100392
## 2023-09-16 21:29:14 Annotating text fragment 86681/100392
## 2023-09-16 21:29:14 Annotating text fragment 86691/100392
## 2023-09-16 21:29:14 Annotating text fragment 86701/100392
## 2023-09-16 21:29:15 Annotating text fragment 86711/100392
## 2023-09-16 21:29:15 Annotating text fragment 86721/100392
## 2023-09-16 21:29:15 Annotating text fragment 86731/100392
## 2023-09-16 21:29:15 Annotating text fragment 86741/100392
## 2023-09-16 21:29:15 Annotating text fragment 86751/100392
## 2023-09-16 21:29:15 Annotating text fragment 86761/100392
## 2023-09-16 21:29:15 Annotating text fragment 86771/100392
## 2023-09-16 21:29:15 Annotating text fragment 86781/100392
## 2023-09-16 21:29:15 Annotating text fragment 86791/100392
## 2023-09-16 21:29:15 Annotating text fragment 86801/100392
## 2023-09-16 21:29:15 Annotating text fragment 86811/100392
## 2023-09-16 21:29:16 Annotating text fragment 86821/100392
## 2023-09-16 21:29:16 Annotating text fragment 86831/100392
## 2023-09-16 21:29:16 Annotating text fragment 86841/100392
## 2023-09-16 21:29:16 Annotating text fragment 86851/100392
## 2023-09-16 21:29:16 Annotating text fragment 86861/100392
## 2023-09-16 21:29:16 Annotating text fragment 86871/100392
## 2023-09-16 21:29:16 Annotating text fragment 86881/100392
## 2023-09-16 21:29:16 Annotating text fragment 86891/100392
## 2023-09-16 21:29:16 Annotating text fragment 86901/100392
## 2023-09-16 21:29:17 Annotating text fragment 86911/100392
## 2023-09-16 21:29:17 Annotating text fragment 86921/100392
## 2023-09-16 21:29:17 Annotating text fragment 86931/100392
## 2023-09-16 21:29:17 Annotating text fragment 86941/100392
## 2023-09-16 21:29:17 Annotating text fragment 86951/100392
## 2023-09-16 21:29:17 Annotating text fragment 86961/100392
## 2023-09-16 21:29:17 Annotating text fragment 86971/100392
## 2023-09-16 21:29:17 Annotating text fragment 86981/100392
## 2023-09-16 21:29:18 Annotating text fragment 86991/100392
```

```
## 2023-09-16 21:29:18 Annotating text fragment 87001/100392
## 2023-09-16 21:29:18 Annotating text fragment 87011/100392
## 2023-09-16 21:29:18 Annotating text fragment 87021/100392
## 2023-09-16 21:29:18 Annotating text fragment 87031/100392
## 2023-09-16 21:29:18 Annotating text fragment 87041/100392
## 2023-09-16 21:29:18 Annotating text fragment 87051/100392
## 2023-09-16 21:29:18 Annotating text fragment 87061/100392
## 2023-09-16 21:29:18 Annotating text fragment 87071/100392
## 2023-09-16 21:29:18 Annotating text fragment 87081/100392
## 2023-09-16 21:29:19 Annotating text fragment 87091/100392
## 2023-09-16 21:29:19 Annotating text fragment 87101/100392
## 2023-09-16 21:29:19 Annotating text fragment 87111/100392
## 2023-09-16 21:29:19 Annotating text fragment 87121/100392
## 2023-09-16 21:29:19 Annotating text fragment 87131/100392
## 2023-09-16 21:29:19 Annotating text fragment 87141/100392
## 2023-09-16 21:29:19 Annotating text fragment 87151/100392
## 2023-09-16 21:29:19 Annotating text fragment 87161/100392
## 2023-09-16 21:29:19 Annotating text fragment 87171/100392
## 2023-09-16 21:29:20 Annotating text fragment 87181/100392
## 2023-09-16 21:29:20 Annotating text fragment 87191/100392
## 2023-09-16 21:29:20 Annotating text fragment 87201/100392
## 2023-09-16 21:29:20 Annotating text fragment 87211/100392
## 2023-09-16 21:29:20 Annotating text fragment 87221/100392
## 2023-09-16 21:29:20 Annotating text fragment 87231/100392
## 2023-09-16 21:29:20 Annotating text fragment 87241/100392
## 2023-09-16 21:29:20 Annotating text fragment 87251/100392
## 2023-09-16 21:29:20 Annotating text fragment 87261/100392
## 2023-09-16 21:29:20 Annotating text fragment 87271/100392
## 2023-09-16 21:29:20 Annotating text fragment 87281/100392
## 2023-09-16 21:29:20 Annotating text fragment 87291/100392
## 2023-09-16 21:29:20 Annotating text fragment 87301/100392
## 2023-09-16 21:29:21 Annotating text fragment 87311/100392
## 2023-09-16 21:29:21 Annotating text fragment 87321/100392
## 2023-09-16 21:29:21 Annotating text fragment 87331/100392
## 2023-09-16 21:29:21 Annotating text fragment 87341/100392
## 2023-09-16 21:29:21 Annotating text fragment 87351/100392
## 2023-09-16 21:29:21 Annotating text fragment 87361/100392
## 2023-09-16 21:29:21 Annotating text fragment 87371/100392
## 2023-09-16 21:29:21 Annotating text fragment 87381/100392
## 2023-09-16 21:29:21 Annotating text fragment 87391/100392
## 2023-09-16 21:29:21 Annotating text fragment 87401/100392
## 2023-09-16 21:29:21 Annotating text fragment 87411/100392
## 2023-09-16 21:29:22 Annotating text fragment 87421/100392
## 2023-09-16 21:29:22 Annotating text fragment 87431/100392
## 2023-09-16 21:29:22 Annotating text fragment 87441/100392
## 2023-09-16 21:29:22 Annotating text fragment 87451/100392
## 2023-09-16 21:29:22 Annotating text fragment 87461/100392
## 2023-09-16 21:29:22 Annotating text fragment 87471/100392
## 2023-09-16 21:29:22 Annotating text fragment 87481/100392
## 2023-09-16 21:29:22 Annotating text fragment 87491/100392
## 2023-09-16 21:29:22 Annotating text fragment 87501/100392
## 2023-09-16 21:29:22 Annotating text fragment 87511/100392
## 2023-09-16 21:29:22 Annotating text fragment 87521/100392
## 2023-09-16 21:29:23 Annotating text fragment 87531/100392
```

```
## 2023-09-16 21:29:23 Annotating text fragment 87541/100392
## 2023-09-16 21:29:23 Annotating text fragment 87551/100392
## 2023-09-16 21:29:23 Annotating text fragment 87561/100392
## 2023-09-16 21:29:23 Annotating text fragment 87571/100392
## 2023-09-16 21:29:23 Annotating text fragment 87581/100392
## 2023-09-16 21:29:23 Annotating text fragment 87591/100392
## 2023-09-16 21:29:23 Annotating text fragment 87601/100392
## 2023-09-16 21:29:23 Annotating text fragment 87611/100392
## 2023-09-16 21:29:23 Annotating text fragment 87621/100392
## 2023-09-16 21:29:23 Annotating text fragment 87631/100392
## 2023-09-16 21:29:24 Annotating text fragment 87641/100392
## 2023-09-16 21:29:24 Annotating text fragment 87651/100392
## 2023-09-16 21:29:24 Annotating text fragment 87661/100392
## 2023-09-16 21:29:24 Annotating text fragment 87671/100392
## 2023-09-16 21:29:24 Annotating text fragment 87681/100392
## 2023-09-16 21:29:24 Annotating text fragment 87691/100392
## 2023-09-16 21:29:24 Annotating text fragment 87701/100392
## 2023-09-16 21:29:24 Annotating text fragment 87711/100392
## 2023-09-16 21:29:24 Annotating text fragment 87721/100392
## 2023-09-16 21:29:24 Annotating text fragment 87731/100392
## 2023-09-16 21:29:24 Annotating text fragment 87741/100392
## 2023-09-16 21:29:24 Annotating text fragment 87751/100392
## 2023-09-16 21:29:25 Annotating text fragment 87761/100392
## 2023-09-16 21:29:25 Annotating text fragment 87771/100392
## 2023-09-16 21:29:25 Annotating text fragment 87781/100392
## 2023-09-16 21:29:25 Annotating text fragment 87791/100392
## 2023-09-16 21:29:25 Annotating text fragment 87801/100392
## 2023-09-16 21:29:25 Annotating text fragment 87811/100392
## 2023-09-16 21:29:25 Annotating text fragment 87821/100392
## 2023-09-16 21:29:25 Annotating text fragment 87831/100392
## 2023-09-16 21:29:26 Annotating text fragment 87841/100392
## 2023-09-16 21:29:26 Annotating text fragment 87851/100392
## 2023-09-16 21:29:26 Annotating text fragment 87861/100392
## 2023-09-16 21:29:26 Annotating text fragment 87871/100392
## 2023-09-16 21:29:26 Annotating text fragment 87881/100392
## 2023-09-16 21:29:26 Annotating text fragment 87891/100392
## 2023-09-16 21:29:26 Annotating text fragment 87901/100392
## 2023-09-16 21:29:26 Annotating text fragment 87911/100392
## 2023-09-16 21:29:26 Annotating text fragment 87921/100392
## 2023-09-16 21:29:26 Annotating text fragment 87931/100392
## 2023-09-16 21:29:27 Annotating text fragment 87941/100392
## 2023-09-16 21:29:27 Annotating text fragment 87951/100392
## 2023-09-16 21:29:27 Annotating text fragment 87961/100392
## 2023-09-16 21:29:27 Annotating text fragment 87971/100392
## 2023-09-16 21:29:27 Annotating text fragment 87981/100392
## 2023-09-16 21:29:27 Annotating text fragment 87991/100392
## 2023-09-16 21:29:27 Annotating text fragment 88001/100392
## 2023-09-16 21:29:27 Annotating text fragment 88011/100392
## 2023-09-16 21:29:27 Annotating text fragment 88021/100392
## 2023-09-16 21:29:27 Annotating text fragment 88031/100392
## 2023-09-16 21:29:27 Annotating text fragment 88041/100392
## 2023-09-16 21:29:27 Annotating text fragment 88051/100392
## 2023-09-16 21:29:27 Annotating text fragment 88061/100392
## 2023-09-16 21:29:27 Annotating text fragment 88071/100392
```

```
## 2023-09-16 21:29:28 Annotating text fragment 88081/100392
## 2023-09-16 21:29:28 Annotating text fragment 88091/100392
## 2023-09-16 21:29:28 Annotating text fragment 88101/100392
## 2023-09-16 21:29:28 Annotating text fragment 88111/100392
## 2023-09-16 21:29:28 Annotating text fragment 88121/100392
## 2023-09-16 21:29:28 Annotating text fragment 88131/100392
## 2023-09-16 21:29:28 Annotating text fragment 88141/100392
## 2023-09-16 21:29:28 Annotating text fragment 88151/100392
## 2023-09-16 21:29:29 Annotating text fragment 88161/100392
## 2023-09-16 21:29:29 Annotating text fragment 88171/100392
## 2023-09-16 21:29:29 Annotating text fragment 88181/100392
## 2023-09-16 21:29:29 Annotating text fragment 88191/100392
## 2023-09-16 21:29:29 Annotating text fragment 88201/100392
## 2023-09-16 21:29:29 Annotating text fragment 88211/100392
## 2023-09-16 21:29:29 Annotating text fragment 88221/100392
## 2023-09-16 21:29:29 Annotating text fragment 88231/100392
## 2023-09-16 21:29:29 Annotating text fragment 88241/100392
## 2023-09-16 21:29:29 Annotating text fragment 88251/100392
## 2023-09-16 21:29:29 Annotating text fragment 88261/100392
## 2023-09-16 21:29:29 Annotating text fragment 88271/100392
## 2023-09-16 21:29:29 Annotating text fragment 88281/100392
## 2023-09-16 21:29:30 Annotating text fragment 88291/100392
## 2023-09-16 21:29:30 Annotating text fragment 88301/100392
## 2023-09-16 21:29:30 Annotating text fragment 88311/100392
## 2023-09-16 21:29:30 Annotating text fragment 88321/100392
## 2023-09-16 21:29:30 Annotating text fragment 88331/100392
## 2023-09-16 21:29:30 Annotating text fragment 88341/100392
## 2023-09-16 21:29:30 Annotating text fragment 88351/100392
## 2023-09-16 21:29:30 Annotating text fragment 88361/100392
## 2023-09-16 21:29:30 Annotating text fragment 88371/100392
## 2023-09-16 21:29:30 Annotating text fragment 88381/100392
## 2023-09-16 21:29:30 Annotating text fragment 88391/100392
## 2023-09-16 21:29:30 Annotating text fragment 88401/100392
## 2023-09-16 21:29:30 Annotating text fragment 88411/100392
## 2023-09-16 21:29:31 Annotating text fragment 88421/100392
## 2023-09-16 21:29:31 Annotating text fragment 88431/100392
## 2023-09-16 21:29:31 Annotating text fragment 88441/100392
## 2023-09-16 21:29:31 Annotating text fragment 88451/100392
## 2023-09-16 21:29:31 Annotating text fragment 88461/100392
## 2023-09-16 21:29:31 Annotating text fragment 88471/100392
## 2023-09-16 21:29:31 Annotating text fragment 88481/100392
## 2023-09-16 21:29:31 Annotating text fragment 88491/100392
## 2023-09-16 21:29:31 Annotating text fragment 88501/100392
## 2023-09-16 21:29:31 Annotating text fragment 88511/100392
## 2023-09-16 21:29:31 Annotating text fragment 88521/100392
## 2023-09-16 21:29:32 Annotating text fragment 88531/100392
## 2023-09-16 21:29:32 Annotating text fragment 88541/100392
## 2023-09-16 21:29:32 Annotating text fragment 88551/100392
## 2023-09-16 21:29:32 Annotating text fragment 88561/100392
## 2023-09-16 21:29:32 Annotating text fragment 88571/100392
## 2023-09-16 21:29:32 Annotating text fragment 88581/100392
## 2023-09-16 21:29:32 Annotating text fragment 88591/100392
## 2023-09-16 21:29:32 Annotating text fragment 88601/100392
## 2023-09-16 21:29:32 Annotating text fragment 88611/100392
```

```
## 2023-09-16 21:29:32 Annotating text fragment 88621/100392
## 2023-09-16 21:29:32 Annotating text fragment 88631/100392
## 2023-09-16 21:29:32 Annotating text fragment 88641/100392
## 2023-09-16 21:29:33 Annotating text fragment 88651/100392
## 2023-09-16 21:29:33 Annotating text fragment 88661/100392
## 2023-09-16 21:29:33 Annotating text fragment 88671/100392
## 2023-09-16 21:29:33 Annotating text fragment 88681/100392
## 2023-09-16 21:29:33 Annotating text fragment 88691/100392
## 2023-09-16 21:29:33 Annotating text fragment 88701/100392
## 2023-09-16 21:29:33 Annotating text fragment 88711/100392
## 2023-09-16 21:29:33 Annotating text fragment 88721/100392
## 2023-09-16 21:29:33 Annotating text fragment 88731/100392
## 2023-09-16 21:29:33 Annotating text fragment 88741/100392
## 2023-09-16 21:29:33 Annotating text fragment 88751/100392
## 2023-09-16 21:29:33 Annotating text fragment 88761/100392
## 2023-09-16 21:29:34 Annotating text fragment 88771/100392
## 2023-09-16 21:29:34 Annotating text fragment 88781/100392
## 2023-09-16 21:29:34 Annotating text fragment 88791/100392
## 2023-09-16 21:29:34 Annotating text fragment 88801/100392
## 2023-09-16 21:29:34 Annotating text fragment 88811/100392
## 2023-09-16 21:29:34 Annotating text fragment 88821/100392
## 2023-09-16 21:29:34 Annotating text fragment 88831/100392
## 2023-09-16 21:29:34 Annotating text fragment 88841/100392
## 2023-09-16 21:29:34 Annotating text fragment 88851/100392
## 2023-09-16 21:29:34 Annotating text fragment 88861/100392
## 2023-09-16 21:29:35 Annotating text fragment 88871/100392
## 2023-09-16 21:29:35 Annotating text fragment 88881/100392
## 2023-09-16 21:29:35 Annotating text fragment 88891/100392
## 2023-09-16 21:29:35 Annotating text fragment 88901/100392
## 2023-09-16 21:29:35 Annotating text fragment 88911/100392
## 2023-09-16 21:29:35 Annotating text fragment 88921/100392
## 2023-09-16 21:29:36 Annotating text fragment 88931/100392
## 2023-09-16 21:29:36 Annotating text fragment 88941/100392
## 2023-09-16 21:29:36 Annotating text fragment 88951/100392
## 2023-09-16 21:29:36 Annotating text fragment 88961/100392
## 2023-09-16 21:29:36 Annotating text fragment 88971/100392
## 2023-09-16 21:29:36 Annotating text fragment 88981/100392
## 2023-09-16 21:29:36 Annotating text fragment 88991/100392
## 2023-09-16 21:29:36 Annotating text fragment 89001/100392
## 2023-09-16 21:29:36 Annotating text fragment 89011/100392
## 2023-09-16 21:29:36 Annotating text fragment 89021/100392
## 2023-09-16 21:29:36 Annotating text fragment 89031/100392
## 2023-09-16 21:29:36 Annotating text fragment 89041/100392
## 2023-09-16 21:29:36 Annotating text fragment 89051/100392
## 2023-09-16 21:29:36 Annotating text fragment 89061/100392
## 2023-09-16 21:29:37 Annotating text fragment 89071/100392
## 2023-09-16 21:29:37 Annotating text fragment 89081/100392
## 2023-09-16 21:29:37 Annotating text fragment 89091/100392
## 2023-09-16 21:29:37 Annotating text fragment 89101/100392
## 2023-09-16 21:29:37 Annotating text fragment 89111/100392
## 2023-09-16 21:29:37 Annotating text fragment 89121/100392
## 2023-09-16 21:29:37 Annotating text fragment 89131/100392
## 2023-09-16 21:29:37 Annotating text fragment 89141/100392
## 2023-09-16 21:29:38 Annotating text fragment 89151/100392
```

```
## 2023-09-16 21:29:38 Annotating text fragment 89161/100392
## 2023-09-16 21:29:38 Annotating text fragment 89171/100392
## 2023-09-16 21:29:38 Annotating text fragment 89181/100392
## 2023-09-16 21:29:38 Annotating text fragment 89191/100392
## 2023-09-16 21:29:38 Annotating text fragment 89201/100392
## 2023-09-16 21:29:38 Annotating text fragment 89211/100392
## 2023-09-16 21:29:38 Annotating text fragment 89221/100392
## 2023-09-16 21:29:38 Annotating text fragment 89231/100392
## 2023-09-16 21:29:38 Annotating text fragment 89241/100392
## 2023-09-16 21:29:38 Annotating text fragment 89251/100392
## 2023-09-16 21:29:38 Annotating text fragment 89261/100392
## 2023-09-16 21:29:38 Annotating text fragment 89271/100392
## 2023-09-16 21:29:39 Annotating text fragment 89281/100392
## 2023-09-16 21:29:39 Annotating text fragment 89291/100392
## 2023-09-16 21:29:39 Annotating text fragment 89301/100392
## 2023-09-16 21:29:39 Annotating text fragment 89311/100392
## 2023-09-16 21:29:39 Annotating text fragment 89321/100392
## 2023-09-16 21:29:39 Annotating text fragment 89331/100392
## 2023-09-16 21:29:39 Annotating text fragment 89341/100392
## 2023-09-16 21:29:39 Annotating text fragment 89351/100392
## 2023-09-16 21:29:39 Annotating text fragment 89361/100392
## 2023-09-16 21:29:39 Annotating text fragment 89371/100392
## 2023-09-16 21:29:39 Annotating text fragment 89381/100392
## 2023-09-16 21:29:39 Annotating text fragment 89391/100392
## 2023-09-16 21:29:40 Annotating text fragment 89401/100392
## 2023-09-16 21:29:40 Annotating text fragment 89411/100392
## 2023-09-16 21:29:40 Annotating text fragment 89421/100392
## 2023-09-16 21:29:40 Annotating text fragment 89431/100392
## 2023-09-16 21:29:40 Annotating text fragment 89441/100392
## 2023-09-16 21:29:40 Annotating text fragment 89451/100392
## 2023-09-16 21:29:40 Annotating text fragment 89461/100392
## 2023-09-16 21:29:40 Annotating text fragment 89471/100392
## 2023-09-16 21:29:40 Annotating text fragment 89481/100392
## 2023-09-16 21:29:40 Annotating text fragment 89491/100392
## 2023-09-16 21:29:40 Annotating text fragment 89501/100392
## 2023-09-16 21:29:41 Annotating text fragment 89511/100392
## 2023-09-16 21:29:41 Annotating text fragment 89521/100392
## 2023-09-16 21:29:41 Annotating text fragment 89531/100392
## 2023-09-16 21:29:41 Annotating text fragment 89541/100392
## 2023-09-16 21:29:41 Annotating text fragment 89551/100392
## 2023-09-16 21:29:41 Annotating text fragment 89561/100392
## 2023-09-16 21:29:41 Annotating text fragment 89571/100392
## 2023-09-16 21:29:41 Annotating text fragment 89581/100392
## 2023-09-16 21:29:41 Annotating text fragment 89591/100392
## 2023-09-16 21:29:41 Annotating text fragment 89601/100392
## 2023-09-16 21:29:41 Annotating text fragment 89611/100392
## 2023-09-16 21:29:42 Annotating text fragment 89621/100392
## 2023-09-16 21:29:42 Annotating text fragment 89631/100392
## 2023-09-16 21:29:42 Annotating text fragment 89641/100392
## 2023-09-16 21:29:42 Annotating text fragment 89651/100392
## 2023-09-16 21:29:42 Annotating text fragment 89661/100392
## 2023-09-16 21:29:42 Annotating text fragment 89671/100392
## 2023-09-16 21:29:42 Annotating text fragment 89681/100392
## 2023-09-16 21:29:42 Annotating text fragment 89691/100392
```

```
## 2023-09-16 21:29:42 Annotating text fragment 89701/100392
## 2023-09-16 21:29:42 Annotating text fragment 89711/100392
## 2023-09-16 21:29:42 Annotating text fragment 89721/100392
## 2023-09-16 21:29:43 Annotating text fragment 89731/100392
## 2023-09-16 21:29:43 Annotating text fragment 89741/100392
## 2023-09-16 21:29:43 Annotating text fragment 89751/100392
## 2023-09-16 21:29:43 Annotating text fragment 89761/100392
## 2023-09-16 21:29:43 Annotating text fragment 89771/100392
## 2023-09-16 21:29:43 Annotating text fragment 89781/100392
## 2023-09-16 21:29:43 Annotating text fragment 89791/100392
## 2023-09-16 21:29:43 Annotating text fragment 89801/100392
## 2023-09-16 21:29:43 Annotating text fragment 89811/100392
## 2023-09-16 21:29:43 Annotating text fragment 89821/100392
## 2023-09-16 21:29:44 Annotating text fragment 89831/100392
## 2023-09-16 21:29:44 Annotating text fragment 89841/100392
## 2023-09-16 21:29:44 Annotating text fragment 89851/100392
## 2023-09-16 21:29:44 Annotating text fragment 89861/100392
## 2023-09-16 21:29:44 Annotating text fragment 89871/100392
## 2023-09-16 21:29:44 Annotating text fragment 89881/100392
## 2023-09-16 21:29:44 Annotating text fragment 89891/100392
## 2023-09-16 21:29:44 Annotating text fragment 89901/100392
## 2023-09-16 21:29:44 Annotating text fragment 89911/100392
## 2023-09-16 21:29:44 Annotating text fragment 89921/100392
## 2023-09-16 21:29:44 Annotating text fragment 89931/100392
## 2023-09-16 21:29:45 Annotating text fragment 89941/100392
## 2023-09-16 21:29:45 Annotating text fragment 89951/100392
## 2023-09-16 21:29:45 Annotating text fragment 89961/100392
## 2023-09-16 21:29:45 Annotating text fragment 89971/100392
## 2023-09-16 21:29:45 Annotating text fragment 89981/100392
## 2023-09-16 21:29:45 Annotating text fragment 89991/100392
## 2023-09-16 21:29:45 Annotating text fragment 90001/100392
## 2023-09-16 21:29:45 Annotating text fragment 90011/100392
## 2023-09-16 21:29:45 Annotating text fragment 90021/100392
## 2023-09-16 21:29:45 Annotating text fragment 90031/100392
## 2023-09-16 21:29:45 Annotating text fragment 90041/100392
## 2023-09-16 21:29:45 Annotating text fragment 90051/100392
## 2023-09-16 21:29:46 Annotating text fragment 90061/100392
## 2023-09-16 21:29:46 Annotating text fragment 90071/100392
## 2023-09-16 21:29:46 Annotating text fragment 90081/100392
## 2023-09-16 21:29:46 Annotating text fragment 90091/100392
## 2023-09-16 21:29:46 Annotating text fragment 90101/100392
## 2023-09-16 21:29:46 Annotating text fragment 90111/100392
## 2023-09-16 21:29:46 Annotating text fragment 90121/100392
## 2023-09-16 21:29:46 Annotating text fragment 90131/100392
## 2023-09-16 21:29:46 Annotating text fragment 90141/100392
## 2023-09-16 21:29:46 Annotating text fragment 90151/100392
## 2023-09-16 21:29:47 Annotating text fragment 90161/100392
## 2023-09-16 21:29:47 Annotating text fragment 90171/100392
## 2023-09-16 21:29:47 Annotating text fragment 90181/100392
## 2023-09-16 21:29:47 Annotating text fragment 90191/100392
## 2023-09-16 21:29:47 Annotating text fragment 90201/100392
## 2023-09-16 21:29:47 Annotating text fragment 90211/100392
## 2023-09-16 21:29:48 Annotating text fragment 90221/100392
## 2023-09-16 21:29:48 Annotating text fragment 90231/100392
```

```
## 2023-09-16 21:29:48 Annotating text fragment 90241/100392
## 2023-09-16 21:29:48 Annotating text fragment 90251/100392
## 2023-09-16 21:29:48 Annotating text fragment 90261/100392
## 2023-09-16 21:29:48 Annotating text fragment 90271/100392
## 2023-09-16 21:29:48 Annotating text fragment 90281/100392
## 2023-09-16 21:29:49 Annotating text fragment 90291/100392
## 2023-09-16 21:29:49 Annotating text fragment 90301/100392
## 2023-09-16 21:29:49 Annotating text fragment 90311/100392
## 2023-09-16 21:29:49 Annotating text fragment 90321/100392
## 2023-09-16 21:29:49 Annotating text fragment 90331/100392
## 2023-09-16 21:29:49 Annotating text fragment 90341/100392
## 2023-09-16 21:29:49 Annotating text fragment 90351/100392
## 2023-09-16 21:29:49 Annotating text fragment 90361/100392
## 2023-09-16 21:29:49 Annotating text fragment 90371/100392
## 2023-09-16 21:29:49 Annotating text fragment 90381/100392
## 2023-09-16 21:29:49 Annotating text fragment 90391/100392
## 2023-09-16 21:29:49 Annotating text fragment 90401/100392
## 2023-09-16 21:29:50 Annotating text fragment 90411/100392
## 2023-09-16 21:29:50 Annotating text fragment 90421/100392
## 2023-09-16 21:29:50 Annotating text fragment 90431/100392
## 2023-09-16 21:29:50 Annotating text fragment 90441/100392
## 2023-09-16 21:29:50 Annotating text fragment 90451/100392
## 2023-09-16 21:29:50 Annotating text fragment 90461/100392
## 2023-09-16 21:29:50 Annotating text fragment 90471/100392
## 2023-09-16 21:29:50 Annotating text fragment 90481/100392
## 2023-09-16 21:29:50 Annotating text fragment 90491/100392
## 2023-09-16 21:29:50 Annotating text fragment 90501/100392
## 2023-09-16 21:29:51 Annotating text fragment 90511/100392
## 2023-09-16 21:29:51 Annotating text fragment 90521/100392
## 2023-09-16 21:29:51 Annotating text fragment 90531/100392
## 2023-09-16 21:29:51 Annotating text fragment 90541/100392
## 2023-09-16 21:29:51 Annotating text fragment 90551/100392
## 2023-09-16 21:29:51 Annotating text fragment 90561/100392
## 2023-09-16 21:29:51 Annotating text fragment 90571/100392
## 2023-09-16 21:29:51 Annotating text fragment 90581/100392
## 2023-09-16 21:29:51 Annotating text fragment 90591/100392
## 2023-09-16 21:29:52 Annotating text fragment 90601/100392
## 2023-09-16 21:29:52 Annotating text fragment 90611/100392
## 2023-09-16 21:29:52 Annotating text fragment 90621/100392
## 2023-09-16 21:29:52 Annotating text fragment 90631/100392
## 2023-09-16 21:29:52 Annotating text fragment 90641/100392
## 2023-09-16 21:29:52 Annotating text fragment 90651/100392
## 2023-09-16 21:29:52 Annotating text fragment 90661/100392
## 2023-09-16 21:29:52 Annotating text fragment 90671/100392
## 2023-09-16 21:29:52 Annotating text fragment 90681/100392
## 2023-09-16 21:29:52 Annotating text fragment 90691/100392
## 2023-09-16 21:29:53 Annotating text fragment 90701/100392
## 2023-09-16 21:29:53 Annotating text fragment 90711/100392
## 2023-09-16 21:29:53 Annotating text fragment 90721/100392
## 2023-09-16 21:29:53 Annotating text fragment 90731/100392
## 2023-09-16 21:29:53 Annotating text fragment 90741/100392
## 2023-09-16 21:29:53 Annotating text fragment 90751/100392
## 2023-09-16 21:29:53 Annotating text fragment 90761/100392
## 2023-09-16 21:29:53 Annotating text fragment 90771/100392
```

```
## 2023-09-16 21:29:53 Annotating text fragment 90781/100392
## 2023-09-16 21:29:54 Annotating text fragment 90791/100392
## 2023-09-16 21:29:54 Annotating text fragment 90801/100392
## 2023-09-16 21:29:54 Annotating text fragment 90811/100392
## 2023-09-16 21:29:54 Annotating text fragment 90821/100392
## 2023-09-16 21:29:54 Annotating text fragment 90831/100392
## 2023-09-16 21:29:54 Annotating text fragment 90841/100392
## 2023-09-16 21:29:54 Annotating text fragment 90851/100392
## 2023-09-16 21:29:54 Annotating text fragment 90861/100392
## 2023-09-16 21:29:54 Annotating text fragment 90871/100392
## 2023-09-16 21:29:54 Annotating text fragment 90881/100392
## 2023-09-16 21:29:55 Annotating text fragment 90891/100392
## 2023-09-16 21:29:55 Annotating text fragment 90901/100392
## 2023-09-16 21:29:55 Annotating text fragment 90911/100392
## 2023-09-16 21:29:55 Annotating text fragment 90921/100392
## 2023-09-16 21:29:55 Annotating text fragment 90931/100392
## 2023-09-16 21:29:55 Annotating text fragment 90941/100392
## 2023-09-16 21:29:55 Annotating text fragment 90951/100392
## 2023-09-16 21:29:55 Annotating text fragment 90961/100392
## 2023-09-16 21:29:55 Annotating text fragment 90971/100392
## 2023-09-16 21:29:55 Annotating text fragment 90981/100392
## 2023-09-16 21:29:55 Annotating text fragment 90991/100392
## 2023-09-16 21:29:55 Annotating text fragment 91001/100392
## 2023-09-16 21:29:56 Annotating text fragment 91011/100392
## 2023-09-16 21:29:56 Annotating text fragment 91021/100392
## 2023-09-16 21:29:56 Annotating text fragment 91031/100392
## 2023-09-16 21:29:56 Annotating text fragment 91041/100392
## 2023-09-16 21:29:56 Annotating text fragment 91051/100392
## 2023-09-16 21:29:56 Annotating text fragment 91061/100392
## 2023-09-16 21:29:56 Annotating text fragment 91071/100392
## 2023-09-16 21:29:56 Annotating text fragment 91081/100392
## 2023-09-16 21:29:56 Annotating text fragment 91091/100392
## 2023-09-16 21:29:57 Annotating text fragment 91101/100392
## 2023-09-16 21:29:57 Annotating text fragment 91111/100392
## 2023-09-16 21:29:57 Annotating text fragment 91121/100392
## 2023-09-16 21:29:57 Annotating text fragment 91131/100392
## 2023-09-16 21:29:57 Annotating text fragment 91141/100392
## 2023-09-16 21:29:57 Annotating text fragment 91151/100392
## 2023-09-16 21:29:57 Annotating text fragment 91161/100392
## 2023-09-16 21:29:57 Annotating text fragment 91171/100392
## 2023-09-16 21:29:58 Annotating text fragment 91181/100392
## 2023-09-16 21:29:58 Annotating text fragment 91191/100392
## 2023-09-16 21:29:58 Annotating text fragment 91201/100392
## 2023-09-16 21:29:58 Annotating text fragment 91211/100392
## 2023-09-16 21:29:59 Annotating text fragment 91221/100392
## 2023-09-16 21:29:59 Annotating text fragment 91231/100392
## 2023-09-16 21:29:59 Annotating text fragment 91241/100392
## 2023-09-16 21:29:59 Annotating text fragment 91251/100392
## 2023-09-16 21:29:59 Annotating text fragment 91261/100392
## 2023-09-16 21:29:59 Annotating text fragment 91271/100392
## 2023-09-16 21:29:59 Annotating text fragment 91281/100392
## 2023-09-16 21:29:59 Annotating text fragment 91291/100392
## 2023-09-16 21:29:59 Annotating text fragment 91301/100392
## 2023-09-16 21:29:59 Annotating text fragment 91311/100392
```

```
## 2023-09-16 21:29:59 Annotating text fragment 91321/100392
## 2023-09-16 21:30:00 Annotating text fragment 91331/100392
## 2023-09-16 21:30:00 Annotating text fragment 91341/100392
## 2023-09-16 21:30:00 Annotating text fragment 91351/100392
## 2023-09-16 21:30:00 Annotating text fragment 91361/100392
## 2023-09-16 21:30:00 Annotating text fragment 91371/100392
## 2023-09-16 21:30:00 Annotating text fragment 91381/100392
## 2023-09-16 21:30:00 Annotating text fragment 91391/100392
## 2023-09-16 21:30:00 Annotating text fragment 91401/100392
## 2023-09-16 21:30:00 Annotating text fragment 91411/100392
## 2023-09-16 21:30:00 Annotating text fragment 91421/100392
## 2023-09-16 21:30:00 Annotating text fragment 91431/100392
## 2023-09-16 21:30:00 Annotating text fragment 91441/100392
## 2023-09-16 21:30:00 Annotating text fragment 91451/100392
## 2023-09-16 21:30:01 Annotating text fragment 91461/100392
## 2023-09-16 21:30:01 Annotating text fragment 91471/100392
## 2023-09-16 21:30:01 Annotating text fragment 91481/100392
## 2023-09-16 21:30:01 Annotating text fragment 91491/100392
## 2023-09-16 21:30:01 Annotating text fragment 91501/100392
## 2023-09-16 21:30:01 Annotating text fragment 91511/100392
## 2023-09-16 21:30:01 Annotating text fragment 91521/100392
## 2023-09-16 21:30:01 Annotating text fragment 91531/100392
## 2023-09-16 21:30:01 Annotating text fragment 91541/100392
## 2023-09-16 21:30:01 Annotating text fragment 91551/100392
## 2023-09-16 21:30:01 Annotating text fragment 91561/100392
## 2023-09-16 21:30:02 Annotating text fragment 91571/100392
## 2023-09-16 21:30:02 Annotating text fragment 91581/100392
## 2023-09-16 21:30:02 Annotating text fragment 91591/100392
## 2023-09-16 21:30:02 Annotating text fragment 91601/100392
## 2023-09-16 21:30:02 Annotating text fragment 91611/100392
## 2023-09-16 21:30:02 Annotating text fragment 91621/100392
## 2023-09-16 21:30:02 Annotating text fragment 91631/100392
## 2023-09-16 21:30:02 Annotating text fragment 91641/100392
## 2023-09-16 21:30:02 Annotating text fragment 91651/100392
## 2023-09-16 21:30:02 Annotating text fragment 91661/100392
## 2023-09-16 21:30:02 Annotating text fragment 91671/100392
## 2023-09-16 21:30:03 Annotating text fragment 91681/100392
## 2023-09-16 21:30:03 Annotating text fragment 91691/100392
## 2023-09-16 21:30:03 Annotating text fragment 91701/100392
## 2023-09-16 21:30:03 Annotating text fragment 91711/100392
## 2023-09-16 21:30:03 Annotating text fragment 91721/100392
## 2023-09-16 21:30:03 Annotating text fragment 91731/100392
## 2023-09-16 21:30:03 Annotating text fragment 91741/100392
## 2023-09-16 21:30:03 Annotating text fragment 91751/100392
## 2023-09-16 21:30:03 Annotating text fragment 91761/100392
## 2023-09-16 21:30:03 Annotating text fragment 91771/100392
## 2023-09-16 21:30:04 Annotating text fragment 91781/100392
## 2023-09-16 21:30:04 Annotating text fragment 91791/100392
## 2023-09-16 21:30:04 Annotating text fragment 91801/100392
## 2023-09-16 21:30:04 Annotating text fragment 91811/100392
## 2023-09-16 21:30:04 Annotating text fragment 91821/100392
## 2023-09-16 21:30:04 Annotating text fragment 91831/100392
## 2023-09-16 21:30:04 Annotating text fragment 91841/100392
## 2023-09-16 21:30:04 Annotating text fragment 91851/100392
```

```
## 2023-09-16 21:30:04 Annotating text fragment 91861/100392
## 2023-09-16 21:30:04 Annotating text fragment 91871/100392
## 2023-09-16 21:30:04 Annotating text fragment 91881/100392
## 2023-09-16 21:30:04 Annotating text fragment 91891/100392
## 2023-09-16 21:30:05 Annotating text fragment 91901/100392
## 2023-09-16 21:30:05 Annotating text fragment 91911/100392
## 2023-09-16 21:30:05 Annotating text fragment 91921/100392
## 2023-09-16 21:30:05 Annotating text fragment 91931/100392
## 2023-09-16 21:30:05 Annotating text fragment 91941/100392
## 2023-09-16 21:30:05 Annotating text fragment 91951/100392
## 2023-09-16 21:30:05 Annotating text fragment 91961/100392
## 2023-09-16 21:30:05 Annotating text fragment 91971/100392
## 2023-09-16 21:30:05 Annotating text fragment 91981/100392
## 2023-09-16 21:30:05 Annotating text fragment 91991/100392
## 2023-09-16 21:30:05 Annotating text fragment 92001/100392
## 2023-09-16 21:30:05 Annotating text fragment 92011/100392
## 2023-09-16 21:30:05 Annotating text fragment 92021/100392
## 2023-09-16 21:30:06 Annotating text fragment 92031/100392
## 2023-09-16 21:30:06 Annotating text fragment 92041/100392
## 2023-09-16 21:30:06 Annotating text fragment 92051/100392
## 2023-09-16 21:30:06 Annotating text fragment 92061/100392
## 2023-09-16 21:30:06 Annotating text fragment 92071/100392
## 2023-09-16 21:30:06 Annotating text fragment 92081/100392
## 2023-09-16 21:30:06 Annotating text fragment 92091/100392
## 2023-09-16 21:30:06 Annotating text fragment 92101/100392
## 2023-09-16 21:30:06 Annotating text fragment 92111/100392
## 2023-09-16 21:30:07 Annotating text fragment 92121/100392
## 2023-09-16 21:30:07 Annotating text fragment 92131/100392
## 2023-09-16 21:30:07 Annotating text fragment 92141/100392
## 2023-09-16 21:30:07 Annotating text fragment 92151/100392
## 2023-09-16 21:30:07 Annotating text fragment 92161/100392
## 2023-09-16 21:30:08 Annotating text fragment 92171/100392
## 2023-09-16 21:30:08 Annotating text fragment 92181/100392
## 2023-09-16 21:30:08 Annotating text fragment 92191/100392
## 2023-09-16 21:30:08 Annotating text fragment 92201/100392
## 2023-09-16 21:30:08 Annotating text fragment 92211/100392
## 2023-09-16 21:30:08 Annotating text fragment 92221/100392
## 2023-09-16 21:30:09 Annotating text fragment 92231/100392
## 2023-09-16 21:30:09 Annotating text fragment 92241/100392
## 2023-09-16 21:30:09 Annotating text fragment 92251/100392
## 2023-09-16 21:30:09 Annotating text fragment 92261/100392
## 2023-09-16 21:30:09 Annotating text fragment 92271/100392
## 2023-09-16 21:30:09 Annotating text fragment 92281/100392
## 2023-09-16 21:30:09 Annotating text fragment 92291/100392
## 2023-09-16 21:30:09 Annotating text fragment 92301/100392
## 2023-09-16 21:30:09 Annotating text fragment 92311/100392
## 2023-09-16 21:30:09 Annotating text fragment 92321/100392
## 2023-09-16 21:30:09 Annotating text fragment 92331/100392
## 2023-09-16 21:30:09 Annotating text fragment 92341/100392
## 2023-09-16 21:30:09 Annotating text fragment 92351/100392
## 2023-09-16 21:30:10 Annotating text fragment 92361/100392
## 2023-09-16 21:30:10 Annotating text fragment 92371/100392
## 2023-09-16 21:30:10 Annotating text fragment 92381/100392
## 2023-09-16 21:30:10 Annotating text fragment 92391/100392
```

```
## 2023-09-16 21:30:10 Annotating text fragment 92401/100392
## 2023-09-16 21:30:10 Annotating text fragment 92411/100392
## 2023-09-16 21:30:10 Annotating text fragment 92421/100392
## 2023-09-16 21:30:10 Annotating text fragment 92431/100392
## 2023-09-16 21:30:10 Annotating text fragment 92441/100392
## 2023-09-16 21:30:10 Annotating text fragment 92451/100392
## 2023-09-16 21:30:10 Annotating text fragment 92461/100392
## 2023-09-16 21:30:11 Annotating text fragment 92471/100392
## 2023-09-16 21:30:11 Annotating text fragment 92481/100392
## 2023-09-16 21:30:11 Annotating text fragment 92491/100392
## 2023-09-16 21:30:11 Annotating text fragment 92501/100392
## 2023-09-16 21:30:11 Annotating text fragment 92511/100392
## 2023-09-16 21:30:11 Annotating text fragment 92521/100392
## 2023-09-16 21:30:11 Annotating text fragment 92531/100392
## 2023-09-16 21:30:11 Annotating text fragment 92541/100392
## 2023-09-16 21:30:11 Annotating text fragment 92551/100392
## 2023-09-16 21:30:11 Annotating text fragment 92561/100392
## 2023-09-16 21:30:11 Annotating text fragment 92571/100392
## 2023-09-16 21:30:12 Annotating text fragment 92581/100392
## 2023-09-16 21:30:12 Annotating text fragment 92591/100392
## 2023-09-16 21:30:12 Annotating text fragment 92601/100392
## 2023-09-16 21:30:12 Annotating text fragment 92611/100392
## 2023-09-16 21:30:12 Annotating text fragment 92621/100392
## 2023-09-16 21:30:12 Annotating text fragment 92631/100392
## 2023-09-16 21:30:12 Annotating text fragment 92641/100392
## 2023-09-16 21:30:12 Annotating text fragment 92651/100392
## 2023-09-16 21:30:13 Annotating text fragment 92661/100392
## 2023-09-16 21:30:13 Annotating text fragment 92671/100392
## 2023-09-16 21:30:13 Annotating text fragment 92681/100392
## 2023-09-16 21:30:13 Annotating text fragment 92691/100392
## 2023-09-16 21:30:13 Annotating text fragment 92701/100392
## 2023-09-16 21:30:13 Annotating text fragment 92711/100392
## 2023-09-16 21:30:13 Annotating text fragment 92721/100392
## 2023-09-16 21:30:13 Annotating text fragment 92731/100392
## 2023-09-16 21:30:13 Annotating text fragment 92741/100392
## 2023-09-16 21:30:13 Annotating text fragment 92751/100392
## 2023-09-16 21:30:14 Annotating text fragment 92761/100392
## 2023-09-16 21:30:14 Annotating text fragment 92771/100392
## 2023-09-16 21:30:14 Annotating text fragment 92781/100392
## 2023-09-16 21:30:14 Annotating text fragment 92791/100392
## 2023-09-16 21:30:14 Annotating text fragment 92801/100392
## 2023-09-16 21:30:14 Annotating text fragment 92811/100392
## 2023-09-16 21:30:14 Annotating text fragment 92821/100392
## 2023-09-16 21:30:14 Annotating text fragment 92831/100392
## 2023-09-16 21:30:14 Annotating text fragment 92841/100392
## 2023-09-16 21:30:14 Annotating text fragment 92851/100392
## 2023-09-16 21:30:14 Annotating text fragment 92861/100392
## 2023-09-16 21:30:15 Annotating text fragment 92871/100392
## 2023-09-16 21:30:15 Annotating text fragment 92881/100392
## 2023-09-16 21:30:15 Annotating text fragment 92891/100392
## 2023-09-16 21:30:15 Annotating text fragment 92901/100392
## 2023-09-16 21:30:15 Annotating text fragment 92911/100392
## 2023-09-16 21:30:15 Annotating text fragment 92921/100392
## 2023-09-16 21:30:15 Annotating text fragment 92931/100392
```

```
## 2023-09-16 21:30:15 Annotating text fragment 92941/100392
## 2023-09-16 21:30:15 Annotating text fragment 92951/100392
## 2023-09-16 21:30:15 Annotating text fragment 92961/100392
## 2023-09-16 21:30:15 Annotating text fragment 92971/100392
## 2023-09-16 21:30:15 Annotating text fragment 92981/100392
## 2023-09-16 21:30:16 Annotating text fragment 92991/100392
## 2023-09-16 21:30:16 Annotating text fragment 93001/100392
## 2023-09-16 21:30:16 Annotating text fragment 93011/100392
## 2023-09-16 21:30:16 Annotating text fragment 93021/100392
## 2023-09-16 21:30:16 Annotating text fragment 93031/100392
## 2023-09-16 21:30:16 Annotating text fragment 93041/100392
## 2023-09-16 21:30:16 Annotating text fragment 93051/100392
## 2023-09-16 21:30:16 Annotating text fragment 93061/100392
## 2023-09-16 21:30:16 Annotating text fragment 93071/100392
## 2023-09-16 21:30:16 Annotating text fragment 93081/100392
## 2023-09-16 21:30:16 Annotating text fragment 93091/100392
## 2023-09-16 21:30:16 Annotating text fragment 93101/100392
## 2023-09-16 21:30:16 Annotating text fragment 93111/100392
## 2023-09-16 21:30:17 Annotating text fragment 93121/100392
## 2023-09-16 21:30:17 Annotating text fragment 93131/100392
## 2023-09-16 21:30:17 Annotating text fragment 93141/100392
## 2023-09-16 21:30:17 Annotating text fragment 93151/100392
## 2023-09-16 21:30:17 Annotating text fragment 93161/100392
## 2023-09-16 21:30:17 Annotating text fragment 93171/100392
## 2023-09-16 21:30:17 Annotating text fragment 93181/100392
## 2023-09-16 21:30:17 Annotating text fragment 93191/100392
## 2023-09-16 21:30:17 Annotating text fragment 93201/100392
## 2023-09-16 21:30:17 Annotating text fragment 93211/100392
## 2023-09-16 21:30:17 Annotating text fragment 93221/100392
## 2023-09-16 21:30:17 Annotating text fragment 93231/100392
## 2023-09-16 21:30:18 Annotating text fragment 93241/100392
## 2023-09-16 21:30:18 Annotating text fragment 93251/100392
## 2023-09-16 21:30:18 Annotating text fragment 93261/100392
## 2023-09-16 21:30:18 Annotating text fragment 93271/100392
## 2023-09-16 21:30:18 Annotating text fragment 93281/100392
## 2023-09-16 21:30:18 Annotating text fragment 93291/100392
## 2023-09-16 21:30:18 Annotating text fragment 93301/100392
## 2023-09-16 21:30:18 Annotating text fragment 93311/100392
## 2023-09-16 21:30:18 Annotating text fragment 93321/100392
## 2023-09-16 21:30:18 Annotating text fragment 93331/100392
## 2023-09-16 21:30:19 Annotating text fragment 93341/100392
## 2023-09-16 21:30:19 Annotating text fragment 93351/100392
## 2023-09-16 21:30:19 Annotating text fragment 93361/100392
## 2023-09-16 21:30:19 Annotating text fragment 93371/100392
## 2023-09-16 21:30:19 Annotating text fragment 93381/100392
## 2023-09-16 21:30:19 Annotating text fragment 93391/100392
## 2023-09-16 21:30:19 Annotating text fragment 93401/100392
## 2023-09-16 21:30:19 Annotating text fragment 93411/100392
## 2023-09-16 21:30:19 Annotating text fragment 93421/100392
## 2023-09-16 21:30:19 Annotating text fragment 93431/100392
## 2023-09-16 21:30:19 Annotating text fragment 93441/100392
## 2023-09-16 21:30:20 Annotating text fragment 93451/100392
## 2023-09-16 21:30:20 Annotating text fragment 93461/100392
## 2023-09-16 21:30:20 Annotating text fragment 93471/100392
```

```
## 2023-09-16 21:30:20 Annotating text fragment 93481/100392
## 2023-09-16 21:30:20 Annotating text fragment 93491/100392
## 2023-09-16 21:30:20 Annotating text fragment 93501/100392
## 2023-09-16 21:30:20 Annotating text fragment 93511/100392
## 2023-09-16 21:30:20 Annotating text fragment 93521/100392
## 2023-09-16 21:30:20 Annotating text fragment 93531/100392
## 2023-09-16 21:30:20 Annotating text fragment 93541/100392
## 2023-09-16 21:30:21 Annotating text fragment 93551/100392
## 2023-09-16 21:30:21 Annotating text fragment 93561/100392
## 2023-09-16 21:30:21 Annotating text fragment 93571/100392
## 2023-09-16 21:30:21 Annotating text fragment 93581/100392
## 2023-09-16 21:30:21 Annotating text fragment 93591/100392
## 2023-09-16 21:30:21 Annotating text fragment 93601/100392
## 2023-09-16 21:30:21 Annotating text fragment 93611/100392
## 2023-09-16 21:30:21 Annotating text fragment 93621/100392
## 2023-09-16 21:30:21 Annotating text fragment 93631/100392
## 2023-09-16 21:30:22 Annotating text fragment 93641/100392
## 2023-09-16 21:30:22 Annotating text fragment 93651/100392
## 2023-09-16 21:30:22 Annotating text fragment 93661/100392
## 2023-09-16 21:30:22 Annotating text fragment 93671/100392
## 2023-09-16 21:30:22 Annotating text fragment 93681/100392
## 2023-09-16 21:30:23 Annotating text fragment 93691/100392
## 2023-09-16 21:30:23 Annotating text fragment 93701/100392
## 2023-09-16 21:30:23 Annotating text fragment 93711/100392
## 2023-09-16 21:30:23 Annotating text fragment 93721/100392
## 2023-09-16 21:30:23 Annotating text fragment 93731/100392
## 2023-09-16 21:30:23 Annotating text fragment 93741/100392
## 2023-09-16 21:30:23 Annotating text fragment 93751/100392
## 2023-09-16 21:30:23 Annotating text fragment 93761/100392
## 2023-09-16 21:30:23 Annotating text fragment 93771/100392
## 2023-09-16 21:30:23 Annotating text fragment 93781/100392
## 2023-09-16 21:30:24 Annotating text fragment 93791/100392
## 2023-09-16 21:30:24 Annotating text fragment 93801/100392
## 2023-09-16 21:30:24 Annotating text fragment 93811/100392
## 2023-09-16 21:30:24 Annotating text fragment 93821/100392
## 2023-09-16 21:30:24 Annotating text fragment 93831/100392
## 2023-09-16 21:30:24 Annotating text fragment 93841/100392
## 2023-09-16 21:30:24 Annotating text fragment 93851/100392
## 2023-09-16 21:30:24 Annotating text fragment 93861/100392
## 2023-09-16 21:30:24 Annotating text fragment 93871/100392
## 2023-09-16 21:30:24 Annotating text fragment 93881/100392
## 2023-09-16 21:30:24 Annotating text fragment 93891/100392
## 2023-09-16 21:30:25 Annotating text fragment 93901/100392
## 2023-09-16 21:30:25 Annotating text fragment 93911/100392
## 2023-09-16 21:30:25 Annotating text fragment 93921/100392
## 2023-09-16 21:30:25 Annotating text fragment 93931/100392
## 2023-09-16 21:30:25 Annotating text fragment 93941/100392
## 2023-09-16 21:30:25 Annotating text fragment 93951/100392
## 2023-09-16 21:30:25 Annotating text fragment 93961/100392
## 2023-09-16 21:30:25 Annotating text fragment 93971/100392
## 2023-09-16 21:30:25 Annotating text fragment 93981/100392
## 2023-09-16 21:30:25 Annotating text fragment 93991/100392
## 2023-09-16 21:30:25 Annotating text fragment 94001/100392
## 2023-09-16 21:30:26 Annotating text fragment 94011/100392
```

```
## 2023-09-16 21:30:26 Annotating text fragment 94021/100392
## 2023-09-16 21:30:26 Annotating text fragment 94031/100392
## 2023-09-16 21:30:26 Annotating text fragment 94041/100392
## 2023-09-16 21:30:26 Annotating text fragment 94051/100392
## 2023-09-16 21:30:26 Annotating text fragment 94061/100392
## 2023-09-16 21:30:26 Annotating text fragment 94071/100392
## 2023-09-16 21:30:26 Annotating text fragment 94081/100392
## 2023-09-16 21:30:26 Annotating text fragment 94091/100392
## 2023-09-16 21:30:26 Annotating text fragment 94101/100392
## 2023-09-16 21:30:26 Annotating text fragment 94111/100392
## 2023-09-16 21:30:27 Annotating text fragment 94121/100392
## 2023-09-16 21:30:27 Annotating text fragment 94131/100392
## 2023-09-16 21:30:27 Annotating text fragment 94141/100392
## 2023-09-16 21:30:27 Annotating text fragment 94151/100392
## 2023-09-16 21:30:27 Annotating text fragment 94161/100392
## 2023-09-16 21:30:27 Annotating text fragment 94171/100392
## 2023-09-16 21:30:27 Annotating text fragment 94181/100392
## 2023-09-16 21:30:27 Annotating text fragment 94191/100392
## 2023-09-16 21:30:27 Annotating text fragment 94201/100392
## 2023-09-16 21:30:27 Annotating text fragment 94211/100392
## 2023-09-16 21:30:27 Annotating text fragment 94221/100392
## 2023-09-16 21:30:27 Annotating text fragment 94231/100392
## 2023-09-16 21:30:28 Annotating text fragment 94241/100392
## 2023-09-16 21:30:28 Annotating text fragment 94251/100392
## 2023-09-16 21:30:28 Annotating text fragment 94261/100392
## 2023-09-16 21:30:28 Annotating text fragment 94271/100392
## 2023-09-16 21:30:28 Annotating text fragment 94281/100392
## 2023-09-16 21:30:28 Annotating text fragment 94291/100392
## 2023-09-16 21:30:28 Annotating text fragment 94301/100392
## 2023-09-16 21:30:28 Annotating text fragment 94311/100392
## 2023-09-16 21:30:28 Annotating text fragment 94321/100392
## 2023-09-16 21:30:28 Annotating text fragment 94331/100392
## 2023-09-16 21:30:29 Annotating text fragment 94341/100392
## 2023-09-16 21:30:29 Annotating text fragment 94351/100392
## 2023-09-16 21:30:29 Annotating text fragment 94361/100392
## 2023-09-16 21:30:29 Annotating text fragment 94371/100392
## 2023-09-16 21:30:29 Annotating text fragment 94381/100392
## 2023-09-16 21:30:29 Annotating text fragment 94391/100392
## 2023-09-16 21:30:29 Annotating text fragment 94401/100392
## 2023-09-16 21:30:29 Annotating text fragment 94411/100392
## 2023-09-16 21:30:29 Annotating text fragment 94421/100392
## 2023-09-16 21:30:29 Annotating text fragment 94431/100392
## 2023-09-16 21:30:29 Annotating text fragment 94441/100392
## 2023-09-16 21:30:30 Annotating text fragment 94451/100392
## 2023-09-16 21:30:30 Annotating text fragment 94461/100392
## 2023-09-16 21:30:30 Annotating text fragment 94471/100392
## 2023-09-16 21:30:30 Annotating text fragment 94481/100392
## 2023-09-16 21:30:30 Annotating text fragment 94491/100392
## 2023-09-16 21:30:30 Annotating text fragment 94501/100392
## 2023-09-16 21:30:30 Annotating text fragment 94511/100392
## 2023-09-16 21:30:30 Annotating text fragment 94521/100392
## 2023-09-16 21:30:30 Annotating text fragment 94531/100392
## 2023-09-16 21:30:30 Annotating text fragment 94541/100392
## 2023-09-16 21:30:30 Annotating text fragment 94551/100392
```

```
## 2023-09-16 21:30:31 Annotating text fragment 94561/100392
## 2023-09-16 21:30:31 Annotating text fragment 94571/100392
## 2023-09-16 21:30:31 Annotating text fragment 94581/100392
## 2023-09-16 21:30:31 Annotating text fragment 94591/100392
## 2023-09-16 21:30:31 Annotating text fragment 94601/100392
## 2023-09-16 21:30:31 Annotating text fragment 94611/100392
## 2023-09-16 21:30:31 Annotating text fragment 94621/100392
## 2023-09-16 21:30:31 Annotating text fragment 94631/100392
## 2023-09-16 21:30:32 Annotating text fragment 94641/100392
## 2023-09-16 21:30:32 Annotating text fragment 94651/100392
## 2023-09-16 21:30:32 Annotating text fragment 94661/100392
## 2023-09-16 21:30:32 Annotating text fragment 94671/100392
## 2023-09-16 21:30:32 Annotating text fragment 94681/100392
## 2023-09-16 21:30:32 Annotating text fragment 94691/100392
## 2023-09-16 21:30:32 Annotating text fragment 94701/100392
## 2023-09-16 21:30:32 Annotating text fragment 94711/100392
## 2023-09-16 21:30:33 Annotating text fragment 94721/100392
## 2023-09-16 21:30:33 Annotating text fragment 94731/100392
## 2023-09-16 21:30:33 Annotating text fragment 94741/100392
## 2023-09-16 21:30:33 Annotating text fragment 94751/100392
## 2023-09-16 21:30:33 Annotating text fragment 94761/100392
## 2023-09-16 21:30:33 Annotating text fragment 94771/100392
## 2023-09-16 21:30:33 Annotating text fragment 94781/100392
## 2023-09-16 21:30:33 Annotating text fragment 94791/100392
## 2023-09-16 21:30:33 Annotating text fragment 94801/100392
## 2023-09-16 21:30:33 Annotating text fragment 94811/100392
## 2023-09-16 21:30:33 Annotating text fragment 94821/100392
## 2023-09-16 21:30:33 Annotating text fragment 94831/100392
## 2023-09-16 21:30:34 Annotating text fragment 94841/100392
## 2023-09-16 21:30:34 Annotating text fragment 94851/100392
## 2023-09-16 21:30:34 Annotating text fragment 94861/100392
## 2023-09-16 21:30:34 Annotating text fragment 94871/100392
## 2023-09-16 21:30:34 Annotating text fragment 94881/100392
## 2023-09-16 21:30:34 Annotating text fragment 94891/100392
## 2023-09-16 21:30:34 Annotating text fragment 94901/100392
## 2023-09-16 21:30:34 Annotating text fragment 94911/100392
## 2023-09-16 21:30:35 Annotating text fragment 94921/100392
## 2023-09-16 21:30:35 Annotating text fragment 94931/100392
## 2023-09-16 21:30:35 Annotating text fragment 94941/100392
## 2023-09-16 21:30:35 Annotating text fragment 94951/100392
## 2023-09-16 21:30:35 Annotating text fragment 94961/100392
## 2023-09-16 21:30:35 Annotating text fragment 94971/100392
## 2023-09-16 21:30:35 Annotating text fragment 94981/100392
## 2023-09-16 21:30:35 Annotating text fragment 94991/100392
## 2023-09-16 21:30:36 Annotating text fragment 95001/100392
## 2023-09-16 21:30:36 Annotating text fragment 95011/100392
## 2023-09-16 21:30:36 Annotating text fragment 95021/100392
## 2023-09-16 21:30:36 Annotating text fragment 95031/100392
## 2023-09-16 21:30:36 Annotating text fragment 95041/100392
## 2023-09-16 21:30:36 Annotating text fragment 95051/100392
## 2023-09-16 21:30:36 Annotating text fragment 95061/100392
## 2023-09-16 21:30:36 Annotating text fragment 95071/100392
## 2023-09-16 21:30:36 Annotating text fragment 95081/100392
## 2023-09-16 21:30:36 Annotating text fragment 95091/100392
```

```
## 2023-09-16 21:30:37 Annotating text fragment 95101/100392
## 2023-09-16 21:30:37 Annotating text fragment 95111/100392
## 2023-09-16 21:30:37 Annotating text fragment 95121/100392
## 2023-09-16 21:30:37 Annotating text fragment 95131/100392
## 2023-09-16 21:30:37 Annotating text fragment 95141/100392
## 2023-09-16 21:30:37 Annotating text fragment 95151/100392
## 2023-09-16 21:30:37 Annotating text fragment 95161/100392
## 2023-09-16 21:30:37 Annotating text fragment 95171/100392
## 2023-09-16 21:30:37 Annotating text fragment 95181/100392
## 2023-09-16 21:30:37 Annotating text fragment 95191/100392
## 2023-09-16 21:30:37 Annotating text fragment 95201/100392
## 2023-09-16 21:30:37 Annotating text fragment 95211/100392
## 2023-09-16 21:30:37 Annotating text fragment 95221/100392
## 2023-09-16 21:30:38 Annotating text fragment 95231/100392
## 2023-09-16 21:30:38 Annotating text fragment 95241/100392
## 2023-09-16 21:30:38 Annotating text fragment 95251/100392
## 2023-09-16 21:30:38 Annotating text fragment 95261/100392
## 2023-09-16 21:30:38 Annotating text fragment 95271/100392
## 2023-09-16 21:30:38 Annotating text fragment 95281/100392
## 2023-09-16 21:30:38 Annotating text fragment 95291/100392
## 2023-09-16 21:30:38 Annotating text fragment 95301/100392
## 2023-09-16 21:30:38 Annotating text fragment 95311/100392
## 2023-09-16 21:30:38 Annotating text fragment 95321/100392
## 2023-09-16 21:30:39 Annotating text fragment 95331/100392
## 2023-09-16 21:30:39 Annotating text fragment 95341/100392
## 2023-09-16 21:30:39 Annotating text fragment 95351/100392
## 2023-09-16 21:30:39 Annotating text fragment 95361/100392
## 2023-09-16 21:30:39 Annotating text fragment 95371/100392
## 2023-09-16 21:30:39 Annotating text fragment 95381/100392
## 2023-09-16 21:30:39 Annotating text fragment 95391/100392
## 2023-09-16 21:30:39 Annotating text fragment 95401/100392
## 2023-09-16 21:30:39 Annotating text fragment 95411/100392
## 2023-09-16 21:30:39 Annotating text fragment 95421/100392
## 2023-09-16 21:30:39 Annotating text fragment 95431/100392
## 2023-09-16 21:30:39 Annotating text fragment 95441/100392
## 2023-09-16 21:30:40 Annotating text fragment 95451/100392
## 2023-09-16 21:30:40 Annotating text fragment 95461/100392
## 2023-09-16 21:30:40 Annotating text fragment 95471/100392
## 2023-09-16 21:30:40 Annotating text fragment 95481/100392
## 2023-09-16 21:30:40 Annotating text fragment 95491/100392
## 2023-09-16 21:30:40 Annotating text fragment 95501/100392
## 2023-09-16 21:30:40 Annotating text fragment 95511/100392
## 2023-09-16 21:30:40 Annotating text fragment 95521/100392
## 2023-09-16 21:30:40 Annotating text fragment 95531/100392
## 2023-09-16 21:30:40 Annotating text fragment 95541/100392
## 2023-09-16 21:30:40 Annotating text fragment 95551/100392
## 2023-09-16 21:30:41 Annotating text fragment 95561/100392
## 2023-09-16 21:30:41 Annotating text fragment 95571/100392
## 2023-09-16 21:30:41 Annotating text fragment 95581/100392
## 2023-09-16 21:30:41 Annotating text fragment 95591/100392
## 2023-09-16 21:30:41 Annotating text fragment 95601/100392
## 2023-09-16 21:30:41 Annotating text fragment 95611/100392
## 2023-09-16 21:30:41 Annotating text fragment 95621/100392
## 2023-09-16 21:30:41 Annotating text fragment 95631/100392
```

```
## 2023-09-16 21:30:41 Annotating text fragment 95641/100392
## 2023-09-16 21:30:42 Annotating text fragment 95651/100392
## 2023-09-16 21:30:42 Annotating text fragment 95661/100392
## 2023-09-16 21:30:42 Annotating text fragment 95671/100392
## 2023-09-16 21:30:42 Annotating text fragment 95681/100392
## 2023-09-16 21:30:42 Annotating text fragment 95691/100392
## 2023-09-16 21:30:42 Annotating text fragment 95701/100392
## 2023-09-16 21:30:42 Annotating text fragment 95711/100392
## 2023-09-16 21:30:42 Annotating text fragment 95721/100392
## 2023-09-16 21:30:42 Annotating text fragment 95731/100392
## 2023-09-16 21:30:42 Annotating text fragment 95741/100392
## 2023-09-16 21:30:43 Annotating text fragment 95751/100392
## 2023-09-16 21:30:43 Annotating text fragment 95761/100392
## 2023-09-16 21:30:43 Annotating text fragment 95771/100392
## 2023-09-16 21:30:43 Annotating text fragment 95781/100392
## 2023-09-16 21:30:43 Annotating text fragment 95791/100392
## 2023-09-16 21:30:43 Annotating text fragment 95801/100392
## 2023-09-16 21:30:43 Annotating text fragment 95811/100392
## 2023-09-16 21:30:43 Annotating text fragment 95821/100392
## 2023-09-16 21:30:43 Annotating text fragment 95831/100392
## 2023-09-16 21:30:43 Annotating text fragment 95841/100392
## 2023-09-16 21:30:43 Annotating text fragment 95851/100392
## 2023-09-16 21:30:44 Annotating text fragment 95861/100392
## 2023-09-16 21:30:44 Annotating text fragment 95871/100392
## 2023-09-16 21:30:44 Annotating text fragment 95881/100392
## 2023-09-16 21:30:44 Annotating text fragment 95891/100392
## 2023-09-16 21:30:44 Annotating text fragment 95901/100392
## 2023-09-16 21:30:44 Annotating text fragment 95911/100392
## 2023-09-16 21:30:45 Annotating text fragment 95921/100392
## 2023-09-16 21:30:45 Annotating text fragment 95931/100392
## 2023-09-16 21:30:45 Annotating text fragment 95941/100392
## 2023-09-16 21:30:45 Annotating text fragment 95951/100392
## 2023-09-16 21:30:45 Annotating text fragment 95961/100392
## 2023-09-16 21:30:45 Annotating text fragment 95971/100392
## 2023-09-16 21:30:45 Annotating text fragment 95981/100392
## 2023-09-16 21:30:45 Annotating text fragment 95991/100392
## 2023-09-16 21:30:45 Annotating text fragment 96001/100392
## 2023-09-16 21:30:45 Annotating text fragment 96011/100392
## 2023-09-16 21:30:45 Annotating text fragment 96021/100392
## 2023-09-16 21:30:45 Annotating text fragment 96031/100392
## 2023-09-16 21:30:46 Annotating text fragment 96041/100392
## 2023-09-16 21:30:46 Annotating text fragment 96051/100392
## 2023-09-16 21:30:46 Annotating text fragment 96061/100392
## 2023-09-16 21:30:46 Annotating text fragment 96071/100392
## 2023-09-16 21:30:46 Annotating text fragment 96081/100392
## 2023-09-16 21:30:46 Annotating text fragment 96091/100392
## 2023-09-16 21:30:46 Annotating text fragment 96101/100392
## 2023-09-16 21:30:46 Annotating text fragment 96111/100392
## 2023-09-16 21:30:46 Annotating text fragment 96121/100392
## 2023-09-16 21:30:46 Annotating text fragment 96131/100392
## 2023-09-16 21:30:47 Annotating text fragment 96141/100392
## 2023-09-16 21:30:47 Annotating text fragment 96151/100392
## 2023-09-16 21:30:47 Annotating text fragment 96161/100392
## 2023-09-16 21:30:47 Annotating text fragment 96171/100392
```

```
## 2023-09-16 21:30:47 Annotating text fragment 96181/100392
## 2023-09-16 21:30:47 Annotating text fragment 96191/100392
## 2023-09-16 21:30:47 Annotating text fragment 96201/100392
## 2023-09-16 21:30:47 Annotating text fragment 96211/100392
## 2023-09-16 21:30:47 Annotating text fragment 96221/100392
## 2023-09-16 21:30:48 Annotating text fragment 96231/100392
## 2023-09-16 21:30:48 Annotating text fragment 96241/100392
## 2023-09-16 21:30:48 Annotating text fragment 96251/100392
## 2023-09-16 21:30:48 Annotating text fragment 96261/100392
## 2023-09-16 21:30:48 Annotating text fragment 96271/100392
## 2023-09-16 21:30:48 Annotating text fragment 96281/100392
## 2023-09-16 21:30:48 Annotating text fragment 96291/100392
## 2023-09-16 21:30:48 Annotating text fragment 96301/100392
## 2023-09-16 21:30:49 Annotating text fragment 96311/100392
## 2023-09-16 21:30:49 Annotating text fragment 96321/100392
## 2023-09-16 21:30:49 Annotating text fragment 96331/100392
## 2023-09-16 21:30:49 Annotating text fragment 96341/100392
## 2023-09-16 21:30:49 Annotating text fragment 96351/100392
## 2023-09-16 21:30:49 Annotating text fragment 96361/100392
## 2023-09-16 21:30:49 Annotating text fragment 96371/100392
## 2023-09-16 21:30:49 Annotating text fragment 96381/100392
## 2023-09-16 21:30:49 Annotating text fragment 96391/100392
## 2023-09-16 21:30:50 Annotating text fragment 96401/100392
## 2023-09-16 21:30:50 Annotating text fragment 96411/100392
## 2023-09-16 21:30:50 Annotating text fragment 96421/100392
## 2023-09-16 21:30:50 Annotating text fragment 96431/100392
## 2023-09-16 21:30:50 Annotating text fragment 96441/100392
## 2023-09-16 21:30:50 Annotating text fragment 96451/100392
## 2023-09-16 21:30:50 Annotating text fragment 96461/100392
## 2023-09-16 21:30:50 Annotating text fragment 96471/100392
## 2023-09-16 21:30:50 Annotating text fragment 96481/100392
## 2023-09-16 21:30:50 Annotating text fragment 96491/100392
## 2023-09-16 21:30:51 Annotating text fragment 96501/100392
## 2023-09-16 21:30:51 Annotating text fragment 96511/100392
## 2023-09-16 21:30:51 Annotating text fragment 96521/100392
## 2023-09-16 21:30:51 Annotating text fragment 96531/100392
## 2023-09-16 21:30:51 Annotating text fragment 96541/100392
## 2023-09-16 21:30:51 Annotating text fragment 96551/100392
## 2023-09-16 21:30:51 Annotating text fragment 96561/100392
## 2023-09-16 21:30:52 Annotating text fragment 96571/100392
## 2023-09-16 21:30:52 Annotating text fragment 96581/100392
## 2023-09-16 21:30:52 Annotating text fragment 96591/100392
## 2023-09-16 21:30:52 Annotating text fragment 96601/100392
## 2023-09-16 21:30:52 Annotating text fragment 96611/100392
## 2023-09-16 21:30:52 Annotating text fragment 96621/100392
## 2023-09-16 21:30:52 Annotating text fragment 96631/100392
## 2023-09-16 21:30:52 Annotating text fragment 96641/100392
## 2023-09-16 21:30:52 Annotating text fragment 96651/100392
## 2023-09-16 21:30:52 Annotating text fragment 96661/100392
## 2023-09-16 21:30:53 Annotating text fragment 96671/100392
## 2023-09-16 21:30:53 Annotating text fragment 96681/100392
## 2023-09-16 21:30:53 Annotating text fragment 96691/100392
## 2023-09-16 21:30:53 Annotating text fragment 96701/100392
## 2023-09-16 21:30:53 Annotating text fragment 96711/100392
```

```
## 2023-09-16 21:30:53 Annotating text fragment 96721/100392
## 2023-09-16 21:30:53 Annotating text fragment 96731/100392
## 2023-09-16 21:30:53 Annotating text fragment 96741/100392
## 2023-09-16 21:30:53 Annotating text fragment 96751/100392
## 2023-09-16 21:30:53 Annotating text fragment 96761/100392
## 2023-09-16 21:30:54 Annotating text fragment 96771/100392
## 2023-09-16 21:30:54 Annotating text fragment 96781/100392
## 2023-09-16 21:30:54 Annotating text fragment 96791/100392
## 2023-09-16 21:30:54 Annotating text fragment 96801/100392
## 2023-09-16 21:30:54 Annotating text fragment 96811/100392
## 2023-09-16 21:30:55 Annotating text fragment 96821/100392
## 2023-09-16 21:30:55 Annotating text fragment 96831/100392
## 2023-09-16 21:30:55 Annotating text fragment 96841/100392
## 2023-09-16 21:30:55 Annotating text fragment 96851/100392
## 2023-09-16 21:30:55 Annotating text fragment 96861/100392
## 2023-09-16 21:30:55 Annotating text fragment 96871/100392
## 2023-09-16 21:30:55 Annotating text fragment 96881/100392
## 2023-09-16 21:30:55 Annotating text fragment 96891/100392
## 2023-09-16 21:30:55 Annotating text fragment 96901/100392
## 2023-09-16 21:30:55 Annotating text fragment 96911/100392
## 2023-09-16 21:30:56 Annotating text fragment 96921/100392
## 2023-09-16 21:30:56 Annotating text fragment 96931/100392
## 2023-09-16 21:30:56 Annotating text fragment 96941/100392
## 2023-09-16 21:30:56 Annotating text fragment 96951/100392
## 2023-09-16 21:30:56 Annotating text fragment 96961/100392
## 2023-09-16 21:30:56 Annotating text fragment 96971/100392
## 2023-09-16 21:30:56 Annotating text fragment 96981/100392
## 2023-09-16 21:30:56 Annotating text fragment 96991/100392
## 2023-09-16 21:30:56 Annotating text fragment 97001/100392
## 2023-09-16 21:30:56 Annotating text fragment 97011/100392
## 2023-09-16 21:30:56 Annotating text fragment 97021/100392
## 2023-09-16 21:30:56 Annotating text fragment 97031/100392
## 2023-09-16 21:30:56 Annotating text fragment 97041/100392
## 2023-09-16 21:30:57 Annotating text fragment 97051/100392
## 2023-09-16 21:30:57 Annotating text fragment 97061/100392
## 2023-09-16 21:30:57 Annotating text fragment 97071/100392
## 2023-09-16 21:30:57 Annotating text fragment 97081/100392
## 2023-09-16 21:30:57 Annotating text fragment 97091/100392
## 2023-09-16 21:30:57 Annotating text fragment 97101/100392
## 2023-09-16 21:30:57 Annotating text fragment 97111/100392
## 2023-09-16 21:30:57 Annotating text fragment 97121/100392
## 2023-09-16 21:30:57 Annotating text fragment 97131/100392
## 2023-09-16 21:30:57 Annotating text fragment 97141/100392
## 2023-09-16 21:30:57 Annotating text fragment 97151/100392
## 2023-09-16 21:30:58 Annotating text fragment 97161/100392
## 2023-09-16 21:30:58 Annotating text fragment 97171/100392
## 2023-09-16 21:30:58 Annotating text fragment 97181/100392
## 2023-09-16 21:30:58 Annotating text fragment 97191/100392
## 2023-09-16 21:30:58 Annotating text fragment 97201/100392
## 2023-09-16 21:30:58 Annotating text fragment 97211/100392
## 2023-09-16 21:30:58 Annotating text fragment 97221/100392
## 2023-09-16 21:30:58 Annotating text fragment 97231/100392
## 2023-09-16 21:30:58 Annotating text fragment 97241/100392
## 2023-09-16 21:30:58 Annotating text fragment 97251/100392
```

```
## 2023-09-16 21:30:59 Annotating text fragment 97261/100392
## 2023-09-16 21:30:59 Annotating text fragment 97271/100392
## 2023-09-16 21:30:59 Annotating text fragment 97281/100392
## 2023-09-16 21:30:59 Annotating text fragment 97291/100392
## 2023-09-16 21:30:59 Annotating text fragment 97301/100392
## 2023-09-16 21:30:59 Annotating text fragment 97311/100392
## 2023-09-16 21:30:59 Annotating text fragment 97321/100392
## 2023-09-16 21:30:59 Annotating text fragment 97331/100392
## 2023-09-16 21:30:59 Annotating text fragment 97341/100392
## 2023-09-16 21:30:59 Annotating text fragment 97351/100392
## 2023-09-16 21:31:00 Annotating text fragment 97361/100392
## 2023-09-16 21:31:00 Annotating text fragment 97371/100392
## 2023-09-16 21:31:00 Annotating text fragment 97381/100392
## 2023-09-16 21:31:00 Annotating text fragment 97391/100392
## 2023-09-16 21:31:00 Annotating text fragment 97401/100392
## 2023-09-16 21:31:00 Annotating text fragment 97411/100392
## 2023-09-16 21:31:00 Annotating text fragment 97421/100392
## 2023-09-16 21:31:00 Annotating text fragment 97431/100392
## 2023-09-16 21:31:00 Annotating text fragment 97441/100392
## 2023-09-16 21:31:01 Annotating text fragment 97451/100392
## 2023-09-16 21:31:01 Annotating text fragment 97461/100392
## 2023-09-16 21:31:01 Annotating text fragment 97471/100392
## 2023-09-16 21:31:01 Annotating text fragment 97481/100392
## 2023-09-16 21:31:01 Annotating text fragment 97491/100392
## 2023-09-16 21:31:01 Annotating text fragment 97501/100392
## 2023-09-16 21:31:01 Annotating text fragment 97511/100392
## 2023-09-16 21:31:01 Annotating text fragment 97521/100392
## 2023-09-16 21:31:01 Annotating text fragment 97531/100392
## 2023-09-16 21:31:01 Annotating text fragment 97541/100392
## 2023-09-16 21:31:02 Annotating text fragment 97551/100392
## 2023-09-16 21:31:02 Annotating text fragment 97561/100392
## 2023-09-16 21:31:02 Annotating text fragment 97571/100392
## 2023-09-16 21:31:02 Annotating text fragment 97581/100392
## 2023-09-16 21:31:02 Annotating text fragment 97591/100392
## 2023-09-16 21:31:02 Annotating text fragment 97601/100392
## 2023-09-16 21:31:02 Annotating text fragment 97611/100392
## 2023-09-16 21:31:02 Annotating text fragment 97621/100392
## 2023-09-16 21:31:02 Annotating text fragment 97631/100392
## 2023-09-16 21:31:02 Annotating text fragment 97641/100392
## 2023-09-16 21:31:03 Annotating text fragment 97651/100392
## 2023-09-16 21:31:03 Annotating text fragment 97661/100392
## 2023-09-16 21:31:03 Annotating text fragment 97671/100392
## 2023-09-16 21:31:03 Annotating text fragment 97681/100392
## 2023-09-16 21:31:03 Annotating text fragment 97691/100392
## 2023-09-16 21:31:03 Annotating text fragment 97701/100392
## 2023-09-16 21:31:03 Annotating text fragment 97711/100392
## 2023-09-16 21:31:03 Annotating text fragment 97721/100392
## 2023-09-16 21:31:03 Annotating text fragment 97731/100392
## 2023-09-16 21:31:03 Annotating text fragment 97741/100392
## 2023-09-16 21:31:03 Annotating text fragment 97751/100392
## 2023-09-16 21:31:03 Annotating text fragment 97761/100392
## 2023-09-16 21:31:04 Annotating text fragment 97771/100392
## 2023-09-16 21:31:04 Annotating text fragment 97781/100392
## 2023-09-16 21:31:04 Annotating text fragment 97791/100392
```

```
## 2023-09-16 21:31:04 Annotating text fragment 97801/100392
## 2023-09-16 21:31:04 Annotating text fragment 97811/100392
## 2023-09-16 21:31:04 Annotating text fragment 97821/100392
## 2023-09-16 21:31:04 Annotating text fragment 97831/100392
## 2023-09-16 21:31:04 Annotating text fragment 97841/100392
## 2023-09-16 21:31:05 Annotating text fragment 97851/100392
## 2023-09-16 21:31:05 Annotating text fragment 97861/100392
## 2023-09-16 21:31:05 Annotating text fragment 97871/100392
## 2023-09-16 21:31:05 Annotating text fragment 97881/100392
## 2023-09-16 21:31:05 Annotating text fragment 97891/100392
## 2023-09-16 21:31:05 Annotating text fragment 97901/100392
## 2023-09-16 21:31:05 Annotating text fragment 97911/100392
## 2023-09-16 21:31:05 Annotating text fragment 97921/100392
## 2023-09-16 21:31:05 Annotating text fragment 97931/100392
## 2023-09-16 21:31:05 Annotating text fragment 97941/100392
## 2023-09-16 21:31:05 Annotating text fragment 97951/100392
## 2023-09-16 21:31:05 Annotating text fragment 97961/100392
## 2023-09-16 21:31:05 Annotating text fragment 97971/100392
## 2023-09-16 21:31:06 Annotating text fragment 97981/100392
## 2023-09-16 21:31:06 Annotating text fragment 97991/100392
## 2023-09-16 21:31:06 Annotating text fragment 98001/100392
## 2023-09-16 21:31:06 Annotating text fragment 98011/100392
## 2023-09-16 21:31:06 Annotating text fragment 98021/100392
## 2023-09-16 21:31:06 Annotating text fragment 98031/100392
## 2023-09-16 21:31:06 Annotating text fragment 98041/100392
## 2023-09-16 21:31:06 Annotating text fragment 98051/100392
## 2023-09-16 21:31:06 Annotating text fragment 98061/100392
## 2023-09-16 21:31:06 Annotating text fragment 98071/100392
## 2023-09-16 21:31:06 Annotating text fragment 98081/100392
## 2023-09-16 21:31:06 Annotating text fragment 98091/100392
## 2023-09-16 21:31:07 Annotating text fragment 98101/100392
## 2023-09-16 21:31:07 Annotating text fragment 98111/100392
## 2023-09-16 21:31:07 Annotating text fragment 98121/100392
## 2023-09-16 21:31:07 Annotating text fragment 98131/100392
## 2023-09-16 21:31:07 Annotating text fragment 98141/100392
## 2023-09-16 21:31:07 Annotating text fragment 98151/100392
## 2023-09-16 21:31:07 Annotating text fragment 98161/100392
## 2023-09-16 21:31:07 Annotating text fragment 98171/100392
## 2023-09-16 21:31:07 Annotating text fragment 98181/100392
## 2023-09-16 21:31:07 Annotating text fragment 98191/100392
## 2023-09-16 21:31:08 Annotating text fragment 98201/100392
## 2023-09-16 21:31:08 Annotating text fragment 98211/100392
## 2023-09-16 21:31:08 Annotating text fragment 98221/100392
## 2023-09-16 21:31:08 Annotating text fragment 98231/100392
## 2023-09-16 21:31:08 Annotating text fragment 98241/100392
## 2023-09-16 21:31:08 Annotating text fragment 98251/100392
## 2023-09-16 21:31:08 Annotating text fragment 98261/100392
## 2023-09-16 21:31:08 Annotating text fragment 98271/100392
## 2023-09-16 21:31:08 Annotating text fragment 98281/100392
## 2023-09-16 21:31:08 Annotating text fragment 98291/100392
## 2023-09-16 21:31:09 Annotating text fragment 98301/100392
## 2023-09-16 21:31:09 Annotating text fragment 98311/100392
## 2023-09-16 21:31:09 Annotating text fragment 98321/100392
## 2023-09-16 21:31:09 Annotating text fragment 98331/100392
```

```
## 2023-09-16 21:31:09 Annotating text fragment 98341/100392
## 2023-09-16 21:31:09 Annotating text fragment 98351/100392
## 2023-09-16 21:31:09 Annotating text fragment 98361/100392
## 2023-09-16 21:31:09 Annotating text fragment 98371/100392
## 2023-09-16 21:31:09 Annotating text fragment 98381/100392
## 2023-09-16 21:31:09 Annotating text fragment 98391/100392
## 2023-09-16 21:31:10 Annotating text fragment 98401/100392
## 2023-09-16 21:31:10 Annotating text fragment 98411/100392
## 2023-09-16 21:31:10 Annotating text fragment 98421/100392
## 2023-09-16 21:31:10 Annotating text fragment 98431/100392
## 2023-09-16 21:31:10 Annotating text fragment 98441/100392
## 2023-09-16 21:31:10 Annotating text fragment 98451/100392
## 2023-09-16 21:31:10 Annotating text fragment 98461/100392
## 2023-09-16 21:31:10 Annotating text fragment 98471/100392
## 2023-09-16 21:31:10 Annotating text fragment 98481/100392
## 2023-09-16 21:31:10 Annotating text fragment 98491/100392
## 2023-09-16 21:31:11 Annotating text fragment 98501/100392
## 2023-09-16 21:31:11 Annotating text fragment 98511/100392
## 2023-09-16 21:31:11 Annotating text fragment 98521/100392
## 2023-09-16 21:31:11 Annotating text fragment 98531/100392
## 2023-09-16 21:31:11 Annotating text fragment 98541/100392
## 2023-09-16 21:31:11 Annotating text fragment 98551/100392
## 2023-09-16 21:31:11 Annotating text fragment 98561/100392
## 2023-09-16 21:31:11 Annotating text fragment 98571/100392
## 2023-09-16 21:31:11 Annotating text fragment 98581/100392
## 2023-09-16 21:31:11 Annotating text fragment 98591/100392
## 2023-09-16 21:31:12 Annotating text fragment 98601/100392
## 2023-09-16 21:31:12 Annotating text fragment 98611/100392
## 2023-09-16 21:31:12 Annotating text fragment 98621/100392
## 2023-09-16 21:31:12 Annotating text fragment 98631/100392
## 2023-09-16 21:31:12 Annotating text fragment 98641/100392
## 2023-09-16 21:31:12 Annotating text fragment 98651/100392
## 2023-09-16 21:31:12 Annotating text fragment 98661/100392
## 2023-09-16 21:31:12 Annotating text fragment 98671/100392
## 2023-09-16 21:31:12 Annotating text fragment 98681/100392
## 2023-09-16 21:31:12 Annotating text fragment 98691/100392
## 2023-09-16 21:31:12 Annotating text fragment 98701/100392
## 2023-09-16 21:31:12 Annotating text fragment 98711/100392
## 2023-09-16 21:31:13 Annotating text fragment 98721/100392
## 2023-09-16 21:31:13 Annotating text fragment 98731/100392
## 2023-09-16 21:31:13 Annotating text fragment 98741/100392
## 2023-09-16 21:31:13 Annotating text fragment 98751/100392
## 2023-09-16 21:31:13 Annotating text fragment 98761/100392
## 2023-09-16 21:31:13 Annotating text fragment 98771/100392
## 2023-09-16 21:31:13 Annotating text fragment 98781/100392
## 2023-09-16 21:31:13 Annotating text fragment 98791/100392
## 2023-09-16 21:31:14 Annotating text fragment 98801/100392
## 2023-09-16 21:31:14 Annotating text fragment 98811/100392
## 2023-09-16 21:31:14 Annotating text fragment 98821/100392
## 2023-09-16 21:31:14 Annotating text fragment 98831/100392
## 2023-09-16 21:31:14 Annotating text fragment 98841/100392
## 2023-09-16 21:31:14 Annotating text fragment 98851/100392
## 2023-09-16 21:31:14 Annotating text fragment 98861/100392
## 2023-09-16 21:31:14 Annotating text fragment 98871/100392
```

```
## 2023-09-16 21:31:14 Annotating text fragment 98881/100392
## 2023-09-16 21:31:14 Annotating text fragment 98891/100392
## 2023-09-16 21:31:14 Annotating text fragment 98901/100392
## 2023-09-16 21:31:15 Annotating text fragment 98911/100392
## 2023-09-16 21:31:15 Annotating text fragment 98921/100392
## 2023-09-16 21:31:15 Annotating text fragment 98931/100392
## 2023-09-16 21:31:15 Annotating text fragment 98941/100392
## 2023-09-16 21:31:15 Annotating text fragment 98951/100392
## 2023-09-16 21:31:15 Annotating text fragment 98961/100392
## 2023-09-16 21:31:15 Annotating text fragment 98971/100392
## 2023-09-16 21:31:15 Annotating text fragment 98981/100392
## 2023-09-16 21:31:15 Annotating text fragment 98991/100392
## 2023-09-16 21:31:15 Annotating text fragment 99001/100392
## 2023-09-16 21:31:15 Annotating text fragment 99011/100392
## 2023-09-16 21:31:16 Annotating text fragment 99021/100392
## 2023-09-16 21:31:16 Annotating text fragment 99031/100392
## 2023-09-16 21:31:16 Annotating text fragment 99041/100392
## 2023-09-16 21:31:16 Annotating text fragment 99051/100392
## 2023-09-16 21:31:16 Annotating text fragment 99061/100392
## 2023-09-16 21:31:16 Annotating text fragment 99071/100392
## 2023-09-16 21:31:16 Annotating text fragment 99081/100392
## 2023-09-16 21:31:17 Annotating text fragment 99091/100392
## 2023-09-16 21:31:17 Annotating text fragment 99101/100392
## 2023-09-16 21:31:17 Annotating text fragment 99111/100392
## 2023-09-16 21:31:17 Annotating text fragment 99121/100392
## 2023-09-16 21:31:17 Annotating text fragment 99131/100392
## 2023-09-16 21:31:18 Annotating text fragment 99141/100392
## 2023-09-16 21:31:18 Annotating text fragment 99151/100392
## 2023-09-16 21:31:18 Annotating text fragment 99161/100392
## 2023-09-16 21:31:18 Annotating text fragment 99171/100392
## 2023-09-16 21:31:18 Annotating text fragment 99181/100392
## 2023-09-16 21:31:18 Annotating text fragment 99191/100392
## 2023-09-16 21:31:18 Annotating text fragment 99201/100392
## 2023-09-16 21:31:19 Annotating text fragment 99211/100392
## 2023-09-16 21:31:19 Annotating text fragment 99221/100392
## 2023-09-16 21:31:19 Annotating text fragment 99231/100392
## 2023-09-16 21:31:19 Annotating text fragment 99241/100392
## 2023-09-16 21:31:19 Annotating text fragment 99251/100392
## 2023-09-16 21:31:19 Annotating text fragment 99261/100392
## 2023-09-16 21:31:19 Annotating text fragment 99271/100392
## 2023-09-16 21:31:20 Annotating text fragment 99281/100392
## 2023-09-16 21:31:20 Annotating text fragment 99291/100392
## 2023-09-16 21:31:20 Annotating text fragment 99301/100392
## 2023-09-16 21:31:20 Annotating text fragment 99311/100392
## 2023-09-16 21:31:20 Annotating text fragment 99321/100392
## 2023-09-16 21:31:20 Annotating text fragment 99331/100392
## 2023-09-16 21:31:20 Annotating text fragment 99341/100392
## 2023-09-16 21:31:20 Annotating text fragment 99351/100392
## 2023-09-16 21:31:20 Annotating text fragment 99361/100392
## 2023-09-16 21:31:21 Annotating text fragment 99371/100392
## 2023-09-16 21:31:21 Annotating text fragment 99381/100392
## 2023-09-16 21:31:21 Annotating text fragment 99391/100392
## 2023-09-16 21:31:21 Annotating text fragment 99401/100392
## 2023-09-16 21:31:21 Annotating text fragment 99411/100392
```

```
## 2023-09-16 21:31:21 Annotating text fragment 99421/100392
## 2023-09-16 21:31:21 Annotating text fragment 99431/100392
## 2023-09-16 21:31:21 Annotating text fragment 99441/100392
## 2023-09-16 21:31:21 Annotating text fragment 99451/100392
## 2023-09-16 21:31:21 Annotating text fragment 99461/100392
## 2023-09-16 21:31:21 Annotating text fragment 99471/100392
## 2023-09-16 21:31:22 Annotating text fragment 99481/100392
## 2023-09-16 21:31:22 Annotating text fragment 99491/100392
## 2023-09-16 21:31:22 Annotating text fragment 99501/100392
## 2023-09-16 21:31:22 Annotating text fragment 99511/100392
## 2023-09-16 21:31:22 Annotating text fragment 99521/100392
## 2023-09-16 21:31:22 Annotating text fragment 99531/100392
## 2023-09-16 21:31:22 Annotating text fragment 99541/100392
## 2023-09-16 21:31:22 Annotating text fragment 99551/100392
## 2023-09-16 21:31:22 Annotating text fragment 99561/100392
## 2023-09-16 21:31:23 Annotating text fragment 99571/100392
## 2023-09-16 21:31:23 Annotating text fragment 99581/100392
## 2023-09-16 21:31:23 Annotating text fragment 99591/100392
## 2023-09-16 21:31:23 Annotating text fragment 99601/100392
## 2023-09-16 21:31:23 Annotating text fragment 99611/100392
## 2023-09-16 21:31:23 Annotating text fragment 99621/100392
## 2023-09-16 21:31:23 Annotating text fragment 99631/100392
## 2023-09-16 21:31:24 Annotating text fragment 99641/100392
## 2023-09-16 21:31:24 Annotating text fragment 99651/100392
## 2023-09-16 21:31:24 Annotating text fragment 99661/100392
## 2023-09-16 21:31:24 Annotating text fragment 99671/100392
## 2023-09-16 21:31:24 Annotating text fragment 99681/100392
## 2023-09-16 21:31:24 Annotating text fragment 99691/100392
## 2023-09-16 21:31:24 Annotating text fragment 99701/100392
## 2023-09-16 21:31:25 Annotating text fragment 99711/100392
## 2023-09-16 21:31:25 Annotating text fragment 99721/100392
## 2023-09-16 21:31:25 Annotating text fragment 99731/100392
## 2023-09-16 21:31:25 Annotating text fragment 99741/100392
## 2023-09-16 21:31:25 Annotating text fragment 99751/100392
## 2023-09-16 21:31:25 Annotating text fragment 99761/100392
## 2023-09-16 21:31:25 Annotating text fragment 99771/100392
## 2023-09-16 21:31:25 Annotating text fragment 99781/100392
## 2023-09-16 21:31:26 Annotating text fragment 99791/100392
## 2023-09-16 21:31:26 Annotating text fragment 99801/100392
## 2023-09-16 21:31:26 Annotating text fragment 99811/100392
## 2023-09-16 21:31:26 Annotating text fragment 99821/100392
## 2023-09-16 21:31:26 Annotating text fragment 99831/100392
## 2023-09-16 21:31:26 Annotating text fragment 99841/100392
## 2023-09-16 21:31:26 Annotating text fragment 99851/100392
## 2023-09-16 21:31:26 Annotating text fragment 99861/100392
## 2023-09-16 21:31:26 Annotating text fragment 99871/100392
## 2023-09-16 21:31:26 Annotating text fragment 99881/100392
## 2023-09-16 21:31:26 Annotating text fragment 99891/100392
## 2023-09-16 21:31:27 Annotating text fragment 99901/100392
## 2023-09-16 21:31:27 Annotating text fragment 99911/100392
## 2023-09-16 21:31:27 Annotating text fragment 99921/100392
## 2023-09-16 21:31:27 Annotating text fragment 99931/100392
## 2023-09-16 21:31:27 Annotating text fragment 99941/100392
## 2023-09-16 21:31:27 Annotating text fragment 99951/100392
```

```
## 2023-09-16 21:31:27 Annotating text fragment 99961/100392
## 2023-09-16 21:31:27 Annotating text fragment 99971/100392
## 2023-09-16 21:31:28 Annotating text fragment 99981/100392
## 2023-09-16 21:31:28 Annotating text fragment 99991/100392
## 2023-09-16 21:31:28 Annotating text fragment 100001/100392
## 2023-09-16 21:31:28 Annotating text fragment 100011/100392
## 2023-09-16 21:31:28 Annotating text fragment 100021/100392
## 2023-09-16 21:31:28 Annotating text fragment 100031/100392
## 2023-09-16 21:31:28 Annotating text fragment 100041/100392
## 2023-09-16 21:31:28 Annotating text fragment 100051/100392
## 2023-09-16 21:31:28 Annotating text fragment 100061/100392
## 2023-09-16 21:31:29 Annotating text fragment 100071/100392
## 2023-09-16 21:31:29 Annotating text fragment 100081/100392
## 2023-09-16 21:31:29 Annotating text fragment 100091/100392
## 2023-09-16 21:31:29 Annotating text fragment 100101/100392
## 2023-09-16 21:31:29 Annotating text fragment 100111/100392
## 2023-09-16 21:31:29 Annotating text fragment 100121/100392
## 2023-09-16 21:31:30 Annotating text fragment 100131/100392
## 2023-09-16 21:31:30 Annotating text fragment 100141/100392
## 2023-09-16 21:31:30 Annotating text fragment 100151/100392
## 2023-09-16 21:31:30 Annotating text fragment 100161/100392
## 2023-09-16 21:31:30 Annotating text fragment 100171/100392
## 2023-09-16 21:31:30 Annotating text fragment 100181/100392
## 2023-09-16 21:31:30 Annotating text fragment 100191/100392
## 2023-09-16 21:31:31 Annotating text fragment 100201/100392
## 2023-09-16 21:31:31 Annotating text fragment 100211/100392
## 2023-09-16 21:31:31 Annotating text fragment 100221/100392
## 2023-09-16 21:31:31 Annotating text fragment 100231/100392
## 2023-09-16 21:31:31 Annotating text fragment 100241/100392
## 2023-09-16 21:31:31 Annotating text fragment 100251/100392
## 2023-09-16 21:31:31 Annotating text fragment 100261/100392
## 2023-09-16 21:31:31 Annotating text fragment 100271/100392
## 2023-09-16 21:31:32 Annotating text fragment 100281/100392
## 2023-09-16 21:31:32 Annotating text fragment 100291/100392
## 2023-09-16 21:31:32 Annotating text fragment 100301/100392
## 2023-09-16 21:31:32 Annotating text fragment 100311/100392
## 2023-09-16 21:31:32 Annotating text fragment 100321/100392
## 2023-09-16 21:31:32 Annotating text fragment 100331/100392
## 2023-09-16 21:31:32 Annotating text fragment 100341/100392
## 2023-09-16 21:31:32 Annotating text fragment 100351/100392
## 2023-09-16 21:31:32 Annotating text fragment 100361/100392
## 2023-09-16 21:31:32 Annotating text fragment 100371/100392
## 2023-09-16 21:31:33 Annotating text fragment 100381/100392
## 2023-09-16 21:31:33 Annotating text fragment 100391/100392
```

We now want to get the words most similar to animal in the embedding and we compare them to the words most similar to happy.

```
## $title
## [1] "100 most similar words to animal with word2vec - umap"
##
## attr(,"class")
## [1] "labels"
```

## Per gender

We need to combine `demo_data$gender` to `cleaned_data`, join on `wid`.

```
## 2023-09-16 21:33:18 Annotating text fragment 1/42019
## 2023-09-16 21:33:18 Annotating text fragment 11/42019
## 2023-09-16 21:33:18 Annotating text fragment 21/42019
## 2023-09-16 21:33:18 Annotating text fragment 31/42019
## 2023-09-16 21:33:18 Annotating text fragment 41/42019
## 2023-09-16 21:33:18 Annotating text fragment 51/42019
## 2023-09-16 21:33:19 Annotating text fragment 61/42019
## 2023-09-16 21:33:19 Annotating text fragment 71/42019
## 2023-09-16 21:33:19 Annotating text fragment 81/42019
## 2023-09-16 21:33:19 Annotating text fragment 91/42019
## 2023-09-16 21:33:19 Annotating text fragment 101/42019
## 2023-09-16 21:33:19 Annotating text fragment 111/42019
## 2023-09-16 21:33:19 Annotating text fragment 121/42019
## 2023-09-16 21:33:19 Annotating text fragment 131/42019
## 2023-09-16 21:33:19 Annotating text fragment 141/42019
```

```
## 2023-09-16 21:33:19 Annotating text fragment 151/42019
## 2023-09-16 21:33:19 Annotating text fragment 161/42019
## 2023-09-16 21:33:19 Annotating text fragment 171/42019
## 2023-09-16 21:33:20 Annotating text fragment 181/42019
## 2023-09-16 21:33:20 Annotating text fragment 191/42019
## 2023-09-16 21:33:20 Annotating text fragment 201/42019
## 2023-09-16 21:33:20 Annotating text fragment 211/42019
## 2023-09-16 21:33:20 Annotating text fragment 221/42019
## 2023-09-16 21:33:20 Annotating text fragment 231/42019
## 2023-09-16 21:33:20 Annotating text fragment 241/42019
## 2023-09-16 21:33:21 Annotating text fragment 251/42019
## 2023-09-16 21:33:21 Annotating text fragment 261/42019
## 2023-09-16 21:33:21 Annotating text fragment 271/42019
## 2023-09-16 21:33:21 Annotating text fragment 281/42019
## 2023-09-16 21:33:21 Annotating text fragment 291/42019
## 2023-09-16 21:33:21 Annotating text fragment 301/42019
## 2023-09-16 21:33:21 Annotating text fragment 311/42019
## 2023-09-16 21:33:21 Annotating text fragment 321/42019
## 2023-09-16 21:33:21 Annotating text fragment 331/42019
## 2023-09-16 21:33:21 Annotating text fragment 341/42019
## 2023-09-16 21:33:21 Annotating text fragment 351/42019
## 2023-09-16 21:33:21 Annotating text fragment 361/42019
## 2023-09-16 21:33:21 Annotating text fragment 371/42019
## 2023-09-16 21:33:21 Annotating text fragment 381/42019
## 2023-09-16 21:33:21 Annotating text fragment 391/42019
## 2023-09-16 21:33:22 Annotating text fragment 401/42019
## 2023-09-16 21:33:22 Annotating text fragment 411/42019
## 2023-09-16 21:33:22 Annotating text fragment 421/42019
## 2023-09-16 21:33:22 Annotating text fragment 431/42019
## 2023-09-16 21:33:22 Annotating text fragment 441/42019
## 2023-09-16 21:33:22 Annotating text fragment 451/42019
## 2023-09-16 21:33:22 Annotating text fragment 461/42019
## 2023-09-16 21:33:22 Annotating text fragment 471/42019
## 2023-09-16 21:33:23 Annotating text fragment 481/42019
## 2023-09-16 21:33:23 Annotating text fragment 491/42019
## 2023-09-16 21:33:23 Annotating text fragment 501/42019
## 2023-09-16 21:33:23 Annotating text fragment 511/42019
## 2023-09-16 21:33:23 Annotating text fragment 521/42019
## 2023-09-16 21:33:23 Annotating text fragment 531/42019
## 2023-09-16 21:33:23 Annotating text fragment 541/42019
## 2023-09-16 21:33:23 Annotating text fragment 551/42019
## 2023-09-16 21:33:23 Annotating text fragment 561/42019
## 2023-09-16 21:33:23 Annotating text fragment 571/42019
## 2023-09-16 21:33:23 Annotating text fragment 581/42019
## 2023-09-16 21:33:23 Annotating text fragment 591/42019
## 2023-09-16 21:33:24 Annotating text fragment 601/42019
## 2023-09-16 21:33:24 Annotating text fragment 611/42019
## 2023-09-16 21:33:24 Annotating text fragment 621/42019
## 2023-09-16 21:33:24 Annotating text fragment 631/42019
## 2023-09-16 21:33:24 Annotating text fragment 641/42019
## 2023-09-16 21:33:24 Annotating text fragment 651/42019
## 2023-09-16 21:33:24 Annotating text fragment 661/42019
## 2023-09-16 21:33:24 Annotating text fragment 671/42019
## 2023-09-16 21:33:24 Annotating text fragment 681/42019
```

```
## 2023-09-16 21:33:24 Annotating text fragment 691/42019
## 2023-09-16 21:33:24 Annotating text fragment 701/42019
## 2023-09-16 21:33:24 Annotating text fragment 711/42019
## 2023-09-16 21:33:24 Annotating text fragment 721/42019
## 2023-09-16 21:33:24 Annotating text fragment 731/42019
## 2023-09-16 21:33:25 Annotating text fragment 741/42019
## 2023-09-16 21:33:25 Annotating text fragment 751/42019
## 2023-09-16 21:33:25 Annotating text fragment 761/42019
## 2023-09-16 21:33:25 Annotating text fragment 771/42019
## 2023-09-16 21:33:25 Annotating text fragment 781/42019
## 2023-09-16 21:33:26 Annotating text fragment 791/42019
## 2023-09-16 21:33:26 Annotating text fragment 801/42019
## 2023-09-16 21:33:26 Annotating text fragment 811/42019
## 2023-09-16 21:33:26 Annotating text fragment 821/42019
## 2023-09-16 21:33:26 Annotating text fragment 831/42019
## 2023-09-16 21:33:26 Annotating text fragment 841/42019
## 2023-09-16 21:33:26 Annotating text fragment 851/42019
## 2023-09-16 21:33:27 Annotating text fragment 861/42019
## 2023-09-16 21:33:27 Annotating text fragment 871/42019
## 2023-09-16 21:33:27 Annotating text fragment 881/42019
## 2023-09-16 21:33:27 Annotating text fragment 891/42019
## 2023-09-16 21:33:27 Annotating text fragment 901/42019
## 2023-09-16 21:33:27 Annotating text fragment 911/42019
## 2023-09-16 21:33:27 Annotating text fragment 921/42019
## 2023-09-16 21:33:27 Annotating text fragment 931/42019
## 2023-09-16 21:33:27 Annotating text fragment 941/42019
## 2023-09-16 21:33:27 Annotating text fragment 951/42019
## 2023-09-16 21:33:28 Annotating text fragment 961/42019
## 2023-09-16 21:33:28 Annotating text fragment 971/42019
## 2023-09-16 21:33:28 Annotating text fragment 981/42019
## 2023-09-16 21:33:28 Annotating text fragment 991/42019
## 2023-09-16 21:33:28 Annotating text fragment 1001/42019
## 2023-09-16 21:33:29 Annotating text fragment 1011/42019
## 2023-09-16 21:33:29 Annotating text fragment 1021/42019
## 2023-09-16 21:33:29 Annotating text fragment 1031/42019
## 2023-09-16 21:33:29 Annotating text fragment 1041/42019
## 2023-09-16 21:33:29 Annotating text fragment 1051/42019
## 2023-09-16 21:33:30 Annotating text fragment 1061/42019
## 2023-09-16 21:33:30 Annotating text fragment 1071/42019
## 2023-09-16 21:33:30 Annotating text fragment 1081/42019
## 2023-09-16 21:33:30 Annotating text fragment 1091/42019
## 2023-09-16 21:33:30 Annotating text fragment 1101/42019
## 2023-09-16 21:33:30 Annotating text fragment 1111/42019
## 2023-09-16 21:33:30 Annotating text fragment 1121/42019
## 2023-09-16 21:33:30 Annotating text fragment 1131/42019
## 2023-09-16 21:33:31 Annotating text fragment 1141/42019
## 2023-09-16 21:33:31 Annotating text fragment 1151/42019
## 2023-09-16 21:33:31 Annotating text fragment 1161/42019
## 2023-09-16 21:33:31 Annotating text fragment 1171/42019
## 2023-09-16 21:33:31 Annotating text fragment 1181/42019
## 2023-09-16 21:33:32 Annotating text fragment 1191/42019
## 2023-09-16 21:33:32 Annotating text fragment 1201/42019
## 2023-09-16 21:33:32 Annotating text fragment 1211/42019
## 2023-09-16 21:33:32 Annotating text fragment 1221/42019
```

```
## 2023-09-16 21:33:32 Annotating text fragment 1231/42019
## 2023-09-16 21:33:32 Annotating text fragment 1241/42019
## 2023-09-16 21:33:32 Annotating text fragment 1251/42019
## 2023-09-16 21:33:32 Annotating text fragment 1261/42019
## 2023-09-16 21:33:32 Annotating text fragment 1271/42019
## 2023-09-16 21:33:33 Annotating text fragment 1281/42019
## 2023-09-16 21:33:33 Annotating text fragment 1291/42019
## 2023-09-16 21:33:33 Annotating text fragment 1301/42019
## 2023-09-16 21:33:33 Annotating text fragment 1311/42019
## 2023-09-16 21:33:33 Annotating text fragment 1321/42019
## 2023-09-16 21:33:33 Annotating text fragment 1331/42019
## 2023-09-16 21:33:33 Annotating text fragment 1341/42019
## 2023-09-16 21:33:33 Annotating text fragment 1351/42019
## 2023-09-16 21:33:33 Annotating text fragment 1361/42019
## 2023-09-16 21:33:34 Annotating text fragment 1371/42019
## 2023-09-16 21:33:34 Annotating text fragment 1381/42019
## 2023-09-16 21:33:34 Annotating text fragment 1391/42019
## 2023-09-16 21:33:34 Annotating text fragment 1401/42019
## 2023-09-16 21:33:35 Annotating text fragment 1411/42019
## 2023-09-16 21:33:35 Annotating text fragment 1421/42019
## 2023-09-16 21:33:36 Annotating text fragment 1431/42019
## 2023-09-16 21:33:36 Annotating text fragment 1441/42019
## 2023-09-16 21:33:36 Annotating text fragment 1451/42019
## 2023-09-16 21:33:36 Annotating text fragment 1461/42019
## 2023-09-16 21:33:36 Annotating text fragment 1471/42019
## 2023-09-16 21:33:36 Annotating text fragment 1481/42019
## 2023-09-16 21:33:36 Annotating text fragment 1491/42019
## 2023-09-16 21:33:36 Annotating text fragment 1501/42019
## 2023-09-16 21:33:36 Annotating text fragment 1511/42019
## 2023-09-16 21:33:36 Annotating text fragment 1521/42019
## 2023-09-16 21:33:36 Annotating text fragment 1531/42019
## 2023-09-16 21:33:36 Annotating text fragment 1541/42019
## 2023-09-16 21:33:36 Annotating text fragment 1551/42019
## 2023-09-16 21:33:36 Annotating text fragment 1561/42019
## 2023-09-16 21:33:37 Annotating text fragment 1571/42019
## 2023-09-16 21:33:37 Annotating text fragment 1581/42019
## 2023-09-16 21:33:37 Annotating text fragment 1591/42019
## 2023-09-16 21:33:37 Annotating text fragment 1601/42019
## 2023-09-16 21:33:37 Annotating text fragment 1611/42019
## 2023-09-16 21:33:37 Annotating text fragment 1621/42019
## 2023-09-16 21:33:37 Annotating text fragment 1631/42019
## 2023-09-16 21:33:37 Annotating text fragment 1641/42019
## 2023-09-16 21:33:37 Annotating text fragment 1651/42019
## 2023-09-16 21:33:37 Annotating text fragment 1661/42019
## 2023-09-16 21:33:37 Annotating text fragment 1671/42019
## 2023-09-16 21:33:37 Annotating text fragment 1681/42019
## 2023-09-16 21:33:38 Annotating text fragment 1691/42019
## 2023-09-16 21:33:38 Annotating text fragment 1701/42019
## 2023-09-16 21:33:38 Annotating text fragment 1711/42019
## 2023-09-16 21:33:38 Annotating text fragment 1721/42019
## 2023-09-16 21:33:38 Annotating text fragment 1731/42019
## 2023-09-16 21:33:38 Annotating text fragment 1741/42019
## 2023-09-16 21:33:39 Annotating text fragment 1751/42019
## 2023-09-16 21:33:39 Annotating text fragment 1761/42019
```

```
## 2023-09-16 21:33:39 Annotating text fragment 1771/42019
## 2023-09-16 21:33:39 Annotating text fragment 1781/42019
## 2023-09-16 21:33:39 Annotating text fragment 1791/42019
## 2023-09-16 21:33:40 Annotating text fragment 1801/42019
## 2023-09-16 21:33:40 Annotating text fragment 1811/42019
## 2023-09-16 21:33:40 Annotating text fragment 1821/42019
## 2023-09-16 21:33:40 Annotating text fragment 1831/42019
## 2023-09-16 21:33:40 Annotating text fragment 1841/42019
## 2023-09-16 21:33:40 Annotating text fragment 1851/42019
## 2023-09-16 21:33:41 Annotating text fragment 1861/42019
## 2023-09-16 21:33:41 Annotating text fragment 1871/42019
## 2023-09-16 21:33:41 Annotating text fragment 1881/42019
## 2023-09-16 21:33:41 Annotating text fragment 1891/42019
## 2023-09-16 21:33:41 Annotating text fragment 1901/42019
## 2023-09-16 21:33:41 Annotating text fragment 1911/42019
## 2023-09-16 21:33:41 Annotating text fragment 1921/42019
## 2023-09-16 21:33:41 Annotating text fragment 1931/42019
## 2023-09-16 21:33:41 Annotating text fragment 1941/42019
## 2023-09-16 21:33:41 Annotating text fragment 1951/42019
## 2023-09-16 21:33:41 Annotating text fragment 1961/42019
## 2023-09-16 21:33:41 Annotating text fragment 1971/42019
## 2023-09-16 21:33:41 Annotating text fragment 1981/42019
## 2023-09-16 21:33:41 Annotating text fragment 1991/42019
## 2023-09-16 21:33:41 Annotating text fragment 2001/42019
## 2023-09-16 21:33:41 Annotating text fragment 2011/42019
## 2023-09-16 21:33:42 Annotating text fragment 2021/42019
## 2023-09-16 21:33:42 Annotating text fragment 2031/42019
## 2023-09-16 21:33:42 Annotating text fragment 2041/42019
## 2023-09-16 21:33:42 Annotating text fragment 2051/42019
## 2023-09-16 21:33:42 Annotating text fragment 2061/42019
## 2023-09-16 21:33:42 Annotating text fragment 2071/42019
## 2023-09-16 21:33:42 Annotating text fragment 2081/42019
## 2023-09-16 21:33:42 Annotating text fragment 2091/42019
## 2023-09-16 21:33:42 Annotating text fragment 2101/42019
## 2023-09-16 21:33:43 Annotating text fragment 2111/42019
## 2023-09-16 21:33:43 Annotating text fragment 2121/42019
## 2023-09-16 21:33:43 Annotating text fragment 2131/42019
## 2023-09-16 21:33:43 Annotating text fragment 2141/42019
## 2023-09-16 21:33:43 Annotating text fragment 2151/42019
## 2023-09-16 21:33:43 Annotating text fragment 2161/42019
## 2023-09-16 21:33:43 Annotating text fragment 2171/42019
## 2023-09-16 21:33:43 Annotating text fragment 2181/42019
## 2023-09-16 21:33:43 Annotating text fragment 2191/42019
## 2023-09-16 21:33:43 Annotating text fragment 2201/42019
## 2023-09-16 21:33:43 Annotating text fragment 2211/42019
## 2023-09-16 21:33:43 Annotating text fragment 2221/42019
## 2023-09-16 21:33:43 Annotating text fragment 2231/42019
## 2023-09-16 21:33:43 Annotating text fragment 2241/42019
## 2023-09-16 21:33:43 Annotating text fragment 2251/42019
## 2023-09-16 21:33:43 Annotating text fragment 2261/42019
## 2023-09-16 21:33:43 Annotating text fragment 2271/42019
## 2023-09-16 21:33:43 Annotating text fragment 2281/42019
## 2023-09-16 21:33:44 Annotating text fragment 2291/42019
## 2023-09-16 21:33:44 Annotating text fragment 2301/42019
```

```
## 2023-09-16 21:33:44 Annotating text fragment 2311/42019
## 2023-09-16 21:33:44 Annotating text fragment 2321/42019
## 2023-09-16 21:33:44 Annotating text fragment 2331/42019
## 2023-09-16 21:33:44 Annotating text fragment 2341/42019
## 2023-09-16 21:33:44 Annotating text fragment 2351/42019
## 2023-09-16 21:33:44 Annotating text fragment 2361/42019
## 2023-09-16 21:33:44 Annotating text fragment 2371/42019
## 2023-09-16 21:33:44 Annotating text fragment 2381/42019
## 2023-09-16 21:33:44 Annotating text fragment 2391/42019
## 2023-09-16 21:33:44 Annotating text fragment 2401/42019
## 2023-09-16 21:33:44 Annotating text fragment 2411/42019
## 2023-09-16 21:33:44 Annotating text fragment 2421/42019
## 2023-09-16 21:33:44 Annotating text fragment 2431/42019
## 2023-09-16 21:33:45 Annotating text fragment 2441/42019
## 2023-09-16 21:33:45 Annotating text fragment 2451/42019
## 2023-09-16 21:33:45 Annotating text fragment 2461/42019
## 2023-09-16 21:33:45 Annotating text fragment 2471/42019
## 2023-09-16 21:33:45 Annotating text fragment 2481/42019
## 2023-09-16 21:33:45 Annotating text fragment 2491/42019
## 2023-09-16 21:33:45 Annotating text fragment 2501/42019
## 2023-09-16 21:33:45 Annotating text fragment 2511/42019
## 2023-09-16 21:33:46 Annotating text fragment 2521/42019
## 2023-09-16 21:33:46 Annotating text fragment 2531/42019
## 2023-09-16 21:33:46 Annotating text fragment 2541/42019
## 2023-09-16 21:33:46 Annotating text fragment 2551/42019
## 2023-09-16 21:33:47 Annotating text fragment 2561/42019
## 2023-09-16 21:33:47 Annotating text fragment 2571/42019
## 2023-09-16 21:33:47 Annotating text fragment 2581/42019
## 2023-09-16 21:33:47 Annotating text fragment 2591/42019
## 2023-09-16 21:33:47 Annotating text fragment 2601/42019
## 2023-09-16 21:33:47 Annotating text fragment 2611/42019
## 2023-09-16 21:33:47 Annotating text fragment 2621/42019
## 2023-09-16 21:33:47 Annotating text fragment 2631/42019
## 2023-09-16 21:33:47 Annotating text fragment 2641/42019
## 2023-09-16 21:33:47 Annotating text fragment 2651/42019
## 2023-09-16 21:33:47 Annotating text fragment 2661/42019
## 2023-09-16 21:33:48 Annotating text fragment 2671/42019
## 2023-09-16 21:33:48 Annotating text fragment 2681/42019
## 2023-09-16 21:33:48 Annotating text fragment 2691/42019
## 2023-09-16 21:33:48 Annotating text fragment 2701/42019
## 2023-09-16 21:33:48 Annotating text fragment 2711/42019
## 2023-09-16 21:33:48 Annotating text fragment 2721/42019
## 2023-09-16 21:33:48 Annotating text fragment 2731/42019
## 2023-09-16 21:33:48 Annotating text fragment 2741/42019
## 2023-09-16 21:33:48 Annotating text fragment 2751/42019
## 2023-09-16 21:33:48 Annotating text fragment 2761/42019
## 2023-09-16 21:33:48 Annotating text fragment 2771/42019
## 2023-09-16 21:33:48 Annotating text fragment 2781/42019
## 2023-09-16 21:33:48 Annotating text fragment 2791/42019
## 2023-09-16 21:33:49 Annotating text fragment 2801/42019
## 2023-09-16 21:33:49 Annotating text fragment 2811/42019
## 2023-09-16 21:33:49 Annotating text fragment 2821/42019
## 2023-09-16 21:33:49 Annotating text fragment 2831/42019
## 2023-09-16 21:33:49 Annotating text fragment 2841/42019
```

```
## 2023-09-16 21:33:49 Annotating text fragment 2851/42019
## 2023-09-16 21:33:49 Annotating text fragment 2861/42019
## 2023-09-16 21:33:49 Annotating text fragment 2871/42019
## 2023-09-16 21:33:49 Annotating text fragment 2881/42019
## 2023-09-16 21:33:49 Annotating text fragment 2891/42019
## 2023-09-16 21:33:49 Annotating text fragment 2901/42019
## 2023-09-16 21:33:49 Annotating text fragment 2911/42019
## 2023-09-16 21:33:49 Annotating text fragment 2921/42019
## 2023-09-16 21:33:49 Annotating text fragment 2931/42019
## 2023-09-16 21:33:49 Annotating text fragment 2941/42019
## 2023-09-16 21:33:49 Annotating text fragment 2951/42019
## 2023-09-16 21:33:49 Annotating text fragment 2961/42019
## 2023-09-16 21:33:49 Annotating text fragment 2971/42019
## 2023-09-16 21:33:49 Annotating text fragment 2981/42019
## 2023-09-16 21:33:50 Annotating text fragment 2991/42019
## 2023-09-16 21:33:50 Annotating text fragment 3001/42019
## 2023-09-16 21:33:50 Annotating text fragment 3011/42019
## 2023-09-16 21:33:50 Annotating text fragment 3021/42019
## 2023-09-16 21:33:50 Annotating text fragment 3031/42019
## 2023-09-16 21:33:50 Annotating text fragment 3041/42019
## 2023-09-16 21:33:50 Annotating text fragment 3051/42019
## 2023-09-16 21:33:50 Annotating text fragment 3061/42019
## 2023-09-16 21:33:50 Annotating text fragment 3071/42019
## 2023-09-16 21:33:51 Annotating text fragment 3081/42019
## 2023-09-16 21:33:51 Annotating text fragment 3091/42019
## 2023-09-16 21:33:51 Annotating text fragment 3101/42019
## 2023-09-16 21:33:51 Annotating text fragment 3111/42019
## 2023-09-16 21:33:51 Annotating text fragment 3121/42019
## 2023-09-16 21:33:51 Annotating text fragment 3131/42019
## 2023-09-16 21:33:51 Annotating text fragment 3141/42019
## 2023-09-16 21:33:51 Annotating text fragment 3151/42019
## 2023-09-16 21:33:51 Annotating text fragment 3161/42019
## 2023-09-16 21:33:52 Annotating text fragment 3171/42019
## 2023-09-16 21:33:52 Annotating text fragment 3181/42019
## 2023-09-16 21:33:52 Annotating text fragment 3191/42019
## 2023-09-16 21:33:52 Annotating text fragment 3201/42019
## 2023-09-16 21:33:52 Annotating text fragment 3211/42019
## 2023-09-16 21:33:52 Annotating text fragment 3221/42019
## 2023-09-16 21:33:52 Annotating text fragment 3231/42019
## 2023-09-16 21:33:52 Annotating text fragment 3241/42019
## 2023-09-16 21:33:52 Annotating text fragment 3251/42019
## 2023-09-16 21:33:52 Annotating text fragment 3261/42019
## 2023-09-16 21:33:52 Annotating text fragment 3271/42019
## 2023-09-16 21:33:52 Annotating text fragment 3281/42019
## 2023-09-16 21:33:52 Annotating text fragment 3291/42019
## 2023-09-16 21:33:52 Annotating text fragment 3301/42019
## 2023-09-16 21:33:52 Annotating text fragment 3311/42019
## 2023-09-16 21:33:52 Annotating text fragment 3321/42019
## 2023-09-16 21:33:52 Annotating text fragment 3331/42019
## 2023-09-16 21:33:53 Annotating text fragment 3341/42019
## 2023-09-16 21:33:53 Annotating text fragment 3351/42019
## 2023-09-16 21:33:53 Annotating text fragment 3361/42019
## 2023-09-16 21:33:53 Annotating text fragment 3371/42019
## 2023-09-16 21:33:53 Annotating text fragment 3381/42019
```

```
## 2023-09-16 21:33:53 Annotating text fragment 3391/42019
## 2023-09-16 21:33:53 Annotating text fragment 3401/42019
## 2023-09-16 21:33:53 Annotating text fragment 3411/42019
## 2023-09-16 21:33:53 Annotating text fragment 3421/42019
## 2023-09-16 21:33:53 Annotating text fragment 3431/42019
## 2023-09-16 21:33:53 Annotating text fragment 3441/42019
## 2023-09-16 21:33:53 Annotating text fragment 3451/42019
## 2023-09-16 21:33:53 Annotating text fragment 3461/42019
## 2023-09-16 21:33:53 Annotating text fragment 3471/42019
## 2023-09-16 21:33:53 Annotating text fragment 3481/42019
## 2023-09-16 21:33:53 Annotating text fragment 3491/42019
## 2023-09-16 21:33:53 Annotating text fragment 3501/42019
## 2023-09-16 21:33:53 Annotating text fragment 3511/42019
## 2023-09-16 21:33:53 Annotating text fragment 3521/42019
## 2023-09-16 21:33:53 Annotating text fragment 3531/42019
## 2023-09-16 21:33:54 Annotating text fragment 3541/42019
## 2023-09-16 21:33:54 Annotating text fragment 3551/42019
## 2023-09-16 21:33:54 Annotating text fragment 3561/42019
## 2023-09-16 21:33:54 Annotating text fragment 3571/42019
## 2023-09-16 21:33:54 Annotating text fragment 3581/42019
## 2023-09-16 21:33:54 Annotating text fragment 3591/42019
## 2023-09-16 21:33:54 Annotating text fragment 3601/42019
## 2023-09-16 21:33:54 Annotating text fragment 3611/42019
## 2023-09-16 21:33:54 Annotating text fragment 3621/42019
## 2023-09-16 21:33:54 Annotating text fragment 3631/42019
## 2023-09-16 21:33:54 Annotating text fragment 3641/42019
## 2023-09-16 21:33:54 Annotating text fragment 3651/42019
## 2023-09-16 21:33:54 Annotating text fragment 3661/42019
## 2023-09-16 21:33:54 Annotating text fragment 3671/42019
## 2023-09-16 21:33:54 Annotating text fragment 3681/42019
## 2023-09-16 21:33:55 Annotating text fragment 3691/42019
## 2023-09-16 21:33:55 Annotating text fragment 3701/42019
## 2023-09-16 21:33:55 Annotating text fragment 3711/42019
## 2023-09-16 21:33:55 Annotating text fragment 3721/42019
## 2023-09-16 21:33:55 Annotating text fragment 3731/42019
## 2023-09-16 21:33:55 Annotating text fragment 3741/42019
## 2023-09-16 21:33:55 Annotating text fragment 3751/42019
## 2023-09-16 21:33:55 Annotating text fragment 3761/42019
## 2023-09-16 21:33:55 Annotating text fragment 3771/42019
## 2023-09-16 21:33:56 Annotating text fragment 3781/42019
## 2023-09-16 21:33:56 Annotating text fragment 3791/42019
## 2023-09-16 21:33:56 Annotating text fragment 3801/42019
## 2023-09-16 21:33:56 Annotating text fragment 3811/42019
## 2023-09-16 21:33:56 Annotating text fragment 3821/42019
## 2023-09-16 21:33:56 Annotating text fragment 3831/42019
## 2023-09-16 21:33:56 Annotating text fragment 3841/42019
## 2023-09-16 21:33:56 Annotating text fragment 3851/42019
## 2023-09-16 21:33:56 Annotating text fragment 3861/42019
## 2023-09-16 21:33:56 Annotating text fragment 3871/42019
## 2023-09-16 21:33:56 Annotating text fragment 3881/42019
## 2023-09-16 21:33:56 Annotating text fragment 3891/42019
## 2023-09-16 21:33:57 Annotating text fragment 3901/42019
## 2023-09-16 21:33:57 Annotating text fragment 3911/42019
## 2023-09-16 21:33:57 Annotating text fragment 3921/42019
```

```
## 2023-09-16 21:33:57 Annotating text fragment 3931/42019
## 2023-09-16 21:33:57 Annotating text fragment 3941/42019
## 2023-09-16 21:33:58 Annotating text fragment 3951/42019
## 2023-09-16 21:33:58 Annotating text fragment 3961/42019
## 2023-09-16 21:33:58 Annotating text fragment 3971/42019
## 2023-09-16 21:33:58 Annotating text fragment 3981/42019
## 2023-09-16 21:33:58 Annotating text fragment 3991/42019
## 2023-09-16 21:33:58 Annotating text fragment 4001/42019
## 2023-09-16 21:33:58 Annotating text fragment 4011/42019
## 2023-09-16 21:33:58 Annotating text fragment 4021/42019
## 2023-09-16 21:33:58 Annotating text fragment 4031/42019
## 2023-09-16 21:33:58 Annotating text fragment 4041/42019
## 2023-09-16 21:33:58 Annotating text fragment 4051/42019
## 2023-09-16 21:33:58 Annotating text fragment 4061/42019
## 2023-09-16 21:33:59 Annotating text fragment 4071/42019
## 2023-09-16 21:33:59 Annotating text fragment 4081/42019
## 2023-09-16 21:33:59 Annotating text fragment 4091/42019
## 2023-09-16 21:33:59 Annotating text fragment 4101/42019
## 2023-09-16 21:33:59 Annotating text fragment 4111/42019
## 2023-09-16 21:33:59 Annotating text fragment 4121/42019
## 2023-09-16 21:33:59 Annotating text fragment 4131/42019
## 2023-09-16 21:33:59 Annotating text fragment 4141/42019
## 2023-09-16 21:33:59 Annotating text fragment 4151/42019
## 2023-09-16 21:33:59 Annotating text fragment 4161/42019
## 2023-09-16 21:33:59 Annotating text fragment 4171/42019
## 2023-09-16 21:33:59 Annotating text fragment 4181/42019
## 2023-09-16 21:33:59 Annotating text fragment 4191/42019
## 2023-09-16 21:33:59 Annotating text fragment 4201/42019
## 2023-09-16 21:33:59 Annotating text fragment 4211/42019
## 2023-09-16 21:33:59 Annotating text fragment 4221/42019
## 2023-09-16 21:33:59 Annotating text fragment 4231/42019
## 2023-09-16 21:34:00 Annotating text fragment 4241/42019
## 2023-09-16 21:34:00 Annotating text fragment 4251/42019
## 2023-09-16 21:34:00 Annotating text fragment 4261/42019
## 2023-09-16 21:34:00 Annotating text fragment 4271/42019
## 2023-09-16 21:34:00 Annotating text fragment 4281/42019
## 2023-09-16 21:34:00 Annotating text fragment 4291/42019
## 2023-09-16 21:34:00 Annotating text fragment 4301/42019
## 2023-09-16 21:34:00 Annotating text fragment 4311/42019
## 2023-09-16 21:34:00 Annotating text fragment 4321/42019
## 2023-09-16 21:34:00 Annotating text fragment 4331/42019
## 2023-09-16 21:34:00 Annotating text fragment 4341/42019
## 2023-09-16 21:34:00 Annotating text fragment 4351/42019
## 2023-09-16 21:34:01 Annotating text fragment 4361/42019
## 2023-09-16 21:34:01 Annotating text fragment 4371/42019
## 2023-09-16 21:34:01 Annotating text fragment 4381/42019
## 2023-09-16 21:34:01 Annotating text fragment 4391/42019
## 2023-09-16 21:34:01 Annotating text fragment 4401/42019
## 2023-09-16 21:34:01 Annotating text fragment 4411/42019
## 2023-09-16 21:34:01 Annotating text fragment 4421/42019
## 2023-09-16 21:34:01 Annotating text fragment 4431/42019
## 2023-09-16 21:34:01 Annotating text fragment 4441/42019
## 2023-09-16 21:34:01 Annotating text fragment 4451/42019
## 2023-09-16 21:34:01 Annotating text fragment 4461/42019
```

```
## 2023-09-16 21:34:01 Annotating text fragment 4471/42019
## 2023-09-16 21:34:01 Annotating text fragment 4481/42019
## 2023-09-16 21:34:02 Annotating text fragment 4491/42019
## 2023-09-16 21:34:02 Annotating text fragment 4501/42019
## 2023-09-16 21:34:02 Annotating text fragment 4511/42019
## 2023-09-16 21:34:02 Annotating text fragment 4521/42019
## 2023-09-16 21:34:02 Annotating text fragment 4531/42019
## 2023-09-16 21:34:02 Annotating text fragment 4541/42019
## 2023-09-16 21:34:02 Annotating text fragment 4551/42019
## 2023-09-16 21:34:02 Annotating text fragment 4561/42019
## 2023-09-16 21:34:02 Annotating text fragment 4571/42019
## 2023-09-16 21:34:02 Annotating text fragment 4581/42019
## 2023-09-16 21:34:02 Annotating text fragment 4591/42019
## 2023-09-16 21:34:02 Annotating text fragment 4601/42019
## 2023-09-16 21:34:03 Annotating text fragment 4611/42019
## 2023-09-16 21:34:03 Annotating text fragment 4621/42019
## 2023-09-16 21:34:03 Annotating text fragment 4631/42019
## 2023-09-16 21:34:03 Annotating text fragment 4641/42019
## 2023-09-16 21:34:03 Annotating text fragment 4651/42019
## 2023-09-16 21:34:03 Annotating text fragment 4661/42019
## 2023-09-16 21:34:03 Annotating text fragment 4671/42019
## 2023-09-16 21:34:03 Annotating text fragment 4681/42019
## 2023-09-16 21:34:03 Annotating text fragment 4691/42019
## 2023-09-16 21:34:04 Annotating text fragment 4701/42019
## 2023-09-16 21:34:04 Annotating text fragment 4711/42019
## 2023-09-16 21:34:04 Annotating text fragment 4721/42019
## 2023-09-16 21:34:04 Annotating text fragment 4731/42019
## 2023-09-16 21:34:04 Annotating text fragment 4741/42019
## 2023-09-16 21:34:04 Annotating text fragment 4751/42019
## 2023-09-16 21:34:04 Annotating text fragment 4761/42019
## 2023-09-16 21:34:04 Annotating text fragment 4771/42019
## 2023-09-16 21:34:04 Annotating text fragment 4781/42019
## 2023-09-16 21:34:04 Annotating text fragment 4791/42019
## 2023-09-16 21:34:05 Annotating text fragment 4801/42019
## 2023-09-16 21:34:05 Annotating text fragment 4811/42019
## 2023-09-16 21:34:05 Annotating text fragment 4821/42019
## 2023-09-16 21:34:05 Annotating text fragment 4831/42019
## 2023-09-16 21:34:05 Annotating text fragment 4841/42019
## 2023-09-16 21:34:05 Annotating text fragment 4851/42019
## 2023-09-16 21:34:05 Annotating text fragment 4861/42019
## 2023-09-16 21:34:05 Annotating text fragment 4871/42019
## 2023-09-16 21:34:05 Annotating text fragment 4881/42019
## 2023-09-16 21:34:05 Annotating text fragment 4891/42019
## 2023-09-16 21:34:05 Annotating text fragment 4901/42019
## 2023-09-16 21:34:05 Annotating text fragment 4911/42019
## 2023-09-16 21:34:05 Annotating text fragment 4921/42019
## 2023-09-16 21:34:06 Annotating text fragment 4931/42019
## 2023-09-16 21:34:06 Annotating text fragment 4941/42019
## 2023-09-16 21:34:06 Annotating text fragment 4951/42019
## 2023-09-16 21:34:06 Annotating text fragment 4961/42019
## 2023-09-16 21:34:07 Annotating text fragment 4971/42019
## 2023-09-16 21:34:07 Annotating text fragment 4981/42019
## 2023-09-16 21:34:07 Annotating text fragment 4991/42019
## 2023-09-16 21:34:07 Annotating text fragment 5001/42019
```

```
## 2023-09-16 21:34:07 Annotating text fragment 5011/42019
## 2023-09-16 21:34:07 Annotating text fragment 5021/42019
## 2023-09-16 21:34:07 Annotating text fragment 5031/42019
## 2023-09-16 21:34:07 Annotating text fragment 5041/42019
## 2023-09-16 21:34:07 Annotating text fragment 5051/42019
## 2023-09-16 21:34:07 Annotating text fragment 5061/42019
## 2023-09-16 21:34:08 Annotating text fragment 5071/42019
## 2023-09-16 21:34:08 Annotating text fragment 5081/42019
## 2023-09-16 21:34:08 Annotating text fragment 5091/42019
## 2023-09-16 21:34:08 Annotating text fragment 5101/42019
## 2023-09-16 21:34:08 Annotating text fragment 5111/42019
## 2023-09-16 21:34:08 Annotating text fragment 5121/42019
## 2023-09-16 21:34:08 Annotating text fragment 5131/42019
## 2023-09-16 21:34:08 Annotating text fragment 5141/42019
## 2023-09-16 21:34:08 Annotating text fragment 5151/42019
## 2023-09-16 21:34:08 Annotating text fragment 5161/42019
## 2023-09-16 21:34:08 Annotating text fragment 5171/42019
## 2023-09-16 21:34:09 Annotating text fragment 5181/42019
## 2023-09-16 21:34:09 Annotating text fragment 5191/42019
## 2023-09-16 21:34:09 Annotating text fragment 5201/42019
## 2023-09-16 21:34:09 Annotating text fragment 5211/42019
## 2023-09-16 21:34:09 Annotating text fragment 5221/42019
## 2023-09-16 21:34:09 Annotating text fragment 5231/42019
## 2023-09-16 21:34:09 Annotating text fragment 5241/42019
## 2023-09-16 21:34:09 Annotating text fragment 5251/42019
## 2023-09-16 21:34:09 Annotating text fragment 5261/42019
## 2023-09-16 21:34:10 Annotating text fragment 5271/42019
## 2023-09-16 21:34:10 Annotating text fragment 5281/42019
## 2023-09-16 21:34:10 Annotating text fragment 5291/42019
## 2023-09-16 21:34:10 Annotating text fragment 5301/42019
## 2023-09-16 21:34:10 Annotating text fragment 5311/42019
## 2023-09-16 21:34:10 Annotating text fragment 5321/42019
## 2023-09-16 21:34:10 Annotating text fragment 5331/42019
## 2023-09-16 21:34:10 Annotating text fragment 5341/42019
## 2023-09-16 21:34:11 Annotating text fragment 5351/42019
## 2023-09-16 21:34:11 Annotating text fragment 5361/42019
## 2023-09-16 21:34:11 Annotating text fragment 5371/42019
## 2023-09-16 21:34:11 Annotating text fragment 5381/42019
## 2023-09-16 21:34:11 Annotating text fragment 5391/42019
## 2023-09-16 21:34:11 Annotating text fragment 5401/42019
## 2023-09-16 21:34:11 Annotating text fragment 5411/42019
## 2023-09-16 21:34:11 Annotating text fragment 5421/42019
## 2023-09-16 21:34:11 Annotating text fragment 5431/42019
## 2023-09-16 21:34:11 Annotating text fragment 5441/42019
## 2023-09-16 21:34:12 Annotating text fragment 5451/42019
## 2023-09-16 21:34:12 Annotating text fragment 5461/42019
## 2023-09-16 21:34:12 Annotating text fragment 5471/42019
## 2023-09-16 21:34:12 Annotating text fragment 5481/42019
## 2023-09-16 21:34:12 Annotating text fragment 5491/42019
## 2023-09-16 21:34:12 Annotating text fragment 5501/42019
## 2023-09-16 21:34:12 Annotating text fragment 5511/42019
## 2023-09-16 21:34:12 Annotating text fragment 5521/42019
## 2023-09-16 21:34:12 Annotating text fragment 5531/42019
## 2023-09-16 21:34:13 Annotating text fragment 5541/42019
```

```
## 2023-09-16 21:34:13 Annotating text fragment 5551/42019
## 2023-09-16 21:34:13 Annotating text fragment 5561/42019
## 2023-09-16 21:34:13 Annotating text fragment 5571/42019
## 2023-09-16 21:34:13 Annotating text fragment 5581/42019
## 2023-09-16 21:34:13 Annotating text fragment 5591/42019
## 2023-09-16 21:34:13 Annotating text fragment 5601/42019
## 2023-09-16 21:34:13 Annotating text fragment 5611/42019
## 2023-09-16 21:34:13 Annotating text fragment 5621/42019
## 2023-09-16 21:34:13 Annotating text fragment 5631/42019
## 2023-09-16 21:34:13 Annotating text fragment 5641/42019
## 2023-09-16 21:34:13 Annotating text fragment 5651/42019
## 2023-09-16 21:34:13 Annotating text fragment 5661/42019
## 2023-09-16 21:34:13 Annotating text fragment 5671/42019
## 2023-09-16 21:34:13 Annotating text fragment 5681/42019
## 2023-09-16 21:34:14 Annotating text fragment 5691/42019
## 2023-09-16 21:34:14 Annotating text fragment 5701/42019
## 2023-09-16 21:34:14 Annotating text fragment 5711/42019
## 2023-09-16 21:34:14 Annotating text fragment 5721/42019
## 2023-09-16 21:34:14 Annotating text fragment 5731/42019
## 2023-09-16 21:34:14 Annotating text fragment 5741/42019
## 2023-09-16 21:34:15 Annotating text fragment 5751/42019
## 2023-09-16 21:34:15 Annotating text fragment 5761/42019
## 2023-09-16 21:34:15 Annotating text fragment 5771/42019
## 2023-09-16 21:34:15 Annotating text fragment 5781/42019
## 2023-09-16 21:34:15 Annotating text fragment 5791/42019
## 2023-09-16 21:34:15 Annotating text fragment 5801/42019
## 2023-09-16 21:34:15 Annotating text fragment 5811/42019
## 2023-09-16 21:34:15 Annotating text fragment 5821/42019
## 2023-09-16 21:34:15 Annotating text fragment 5831/42019
## 2023-09-16 21:34:15 Annotating text fragment 5841/42019
## 2023-09-16 21:34:15 Annotating text fragment 5851/42019
## 2023-09-16 21:34:15 Annotating text fragment 5861/42019
## 2023-09-16 21:34:16 Annotating text fragment 5871/42019
## 2023-09-16 21:34:16 Annotating text fragment 5881/42019
## 2023-09-16 21:34:16 Annotating text fragment 5891/42019
## 2023-09-16 21:34:16 Annotating text fragment 5901/42019
## 2023-09-16 21:34:16 Annotating text fragment 5911/42019
## 2023-09-16 21:34:16 Annotating text fragment 5921/42019
## 2023-09-16 21:34:16 Annotating text fragment 5931/42019
## 2023-09-16 21:34:16 Annotating text fragment 5941/42019
## 2023-09-16 21:34:16 Annotating text fragment 5951/42019
## 2023-09-16 21:34:16 Annotating text fragment 5961/42019
## 2023-09-16 21:34:16 Annotating text fragment 5971/42019
## 2023-09-16 21:34:16 Annotating text fragment 5981/42019
## 2023-09-16 21:34:16 Annotating text fragment 5991/42019
## 2023-09-16 21:34:16 Annotating text fragment 6001/42019
## 2023-09-16 21:34:16 Annotating text fragment 6011/42019
## 2023-09-16 21:34:17 Annotating text fragment 6021/42019
## 2023-09-16 21:34:17 Annotating text fragment 6031/42019
## 2023-09-16 21:34:17 Annotating text fragment 6041/42019
## 2023-09-16 21:34:17 Annotating text fragment 6051/42019
## 2023-09-16 21:34:17 Annotating text fragment 6061/42019
## 2023-09-16 21:34:17 Annotating text fragment 6071/42019
## 2023-09-16 21:34:17 Annotating text fragment 6081/42019
```

```
## 2023-09-16 21:34:17 Annotating text fragment 6091/42019
## 2023-09-16 21:34:17 Annotating text fragment 6101/42019
## 2023-09-16 21:34:18 Annotating text fragment 6111/42019
## 2023-09-16 21:34:18 Annotating text fragment 6121/42019
## 2023-09-16 21:34:18 Annotating text fragment 6131/42019
## 2023-09-16 21:34:18 Annotating text fragment 6141/42019
## 2023-09-16 21:34:18 Annotating text fragment 6151/42019
## 2023-09-16 21:34:18 Annotating text fragment 6161/42019
## 2023-09-16 21:34:18 Annotating text fragment 6171/42019
## 2023-09-16 21:34:18 Annotating text fragment 6181/42019
## 2023-09-16 21:34:19 Annotating text fragment 6191/42019
## 2023-09-16 21:34:19 Annotating text fragment 6201/42019
## 2023-09-16 21:34:19 Annotating text fragment 6211/42019
## 2023-09-16 21:34:19 Annotating text fragment 6221/42019
## 2023-09-16 21:34:19 Annotating text fragment 6231/42019
## 2023-09-16 21:34:19 Annotating text fragment 6241/42019
## 2023-09-16 21:34:19 Annotating text fragment 6251/42019
## 2023-09-16 21:34:19 Annotating text fragment 6261/42019
## 2023-09-16 21:34:19 Annotating text fragment 6271/42019
## 2023-09-16 21:34:20 Annotating text fragment 6281/42019
## 2023-09-16 21:34:20 Annotating text fragment 6291/42019
## 2023-09-16 21:34:20 Annotating text fragment 6301/42019
## 2023-09-16 21:34:20 Annotating text fragment 6311/42019
## 2023-09-16 21:34:21 Annotating text fragment 6321/42019
## 2023-09-16 21:34:21 Annotating text fragment 6331/42019
## 2023-09-16 21:34:21 Annotating text fragment 6341/42019
## 2023-09-16 21:34:21 Annotating text fragment 6351/42019
## 2023-09-16 21:34:22 Annotating text fragment 6361/42019
## 2023-09-16 21:34:22 Annotating text fragment 6371/42019
## 2023-09-16 21:34:22 Annotating text fragment 6381/42019
## 2023-09-16 21:34:22 Annotating text fragment 6391/42019
## 2023-09-16 21:34:22 Annotating text fragment 6401/42019
## 2023-09-16 21:34:22 Annotating text fragment 6411/42019
## 2023-09-16 21:34:22 Annotating text fragment 6421/42019
## 2023-09-16 21:34:22 Annotating text fragment 6431/42019
## 2023-09-16 21:34:22 Annotating text fragment 6441/42019
## 2023-09-16 21:34:23 Annotating text fragment 6451/42019
## 2023-09-16 21:34:23 Annotating text fragment 6461/42019
## 2023-09-16 21:34:23 Annotating text fragment 6471/42019
## 2023-09-16 21:34:23 Annotating text fragment 6481/42019
## 2023-09-16 21:34:23 Annotating text fragment 6491/42019
## 2023-09-16 21:34:23 Annotating text fragment 6501/42019
## 2023-09-16 21:34:23 Annotating text fragment 6511/42019
## 2023-09-16 21:34:23 Annotating text fragment 6521/42019
## 2023-09-16 21:34:23 Annotating text fragment 6531/42019
## 2023-09-16 21:34:23 Annotating text fragment 6541/42019
## 2023-09-16 21:34:24 Annotating text fragment 6551/42019
## 2023-09-16 21:34:24 Annotating text fragment 6561/42019
## 2023-09-16 21:34:24 Annotating text fragment 6571/42019
## 2023-09-16 21:34:24 Annotating text fragment 6581/42019
## 2023-09-16 21:34:24 Annotating text fragment 6591/42019
## 2023-09-16 21:34:24 Annotating text fragment 6601/42019
## 2023-09-16 21:34:24 Annotating text fragment 6611/42019
## 2023-09-16 21:34:24 Annotating text fragment 6621/42019
```

```
## 2023-09-16 21:34:24 Annotating text fragment 6631/42019
## 2023-09-16 21:34:24 Annotating text fragment 6641/42019
## 2023-09-16 21:34:25 Annotating text fragment 6651/42019
## 2023-09-16 21:34:25 Annotating text fragment 6661/42019
## 2023-09-16 21:34:25 Annotating text fragment 6671/42019
## 2023-09-16 21:34:25 Annotating text fragment 6681/42019
## 2023-09-16 21:34:25 Annotating text fragment 6691/42019
## 2023-09-16 21:34:26 Annotating text fragment 6701/42019
## 2023-09-16 21:34:26 Annotating text fragment 6711/42019
## 2023-09-16 21:34:26 Annotating text fragment 6721/42019
## 2023-09-16 21:34:26 Annotating text fragment 6731/42019
## 2023-09-16 21:34:26 Annotating text fragment 6741/42019
## 2023-09-16 21:34:26 Annotating text fragment 6751/42019
## 2023-09-16 21:34:26 Annotating text fragment 6761/42019
## 2023-09-16 21:34:27 Annotating text fragment 6771/42019
## 2023-09-16 21:34:27 Annotating text fragment 6781/42019
## 2023-09-16 21:34:27 Annotating text fragment 6791/42019
## 2023-09-16 21:34:27 Annotating text fragment 6801/42019
## 2023-09-16 21:34:27 Annotating text fragment 6811/42019
## 2023-09-16 21:34:27 Annotating text fragment 6821/42019
## 2023-09-16 21:34:27 Annotating text fragment 6831/42019
## 2023-09-16 21:34:27 Annotating text fragment 6841/42019
## 2023-09-16 21:34:27 Annotating text fragment 6851/42019
## 2023-09-16 21:34:27 Annotating text fragment 6861/42019
## 2023-09-16 21:34:27 Annotating text fragment 6871/42019
## 2023-09-16 21:34:27 Annotating text fragment 6881/42019
## 2023-09-16 21:34:27 Annotating text fragment 6891/42019
## 2023-09-16 21:34:28 Annotating text fragment 6901/42019
## 2023-09-16 21:34:28 Annotating text fragment 6911/42019
## 2023-09-16 21:34:28 Annotating text fragment 6921/42019
## 2023-09-16 21:34:28 Annotating text fragment 6931/42019
## 2023-09-16 21:34:28 Annotating text fragment 6941/42019
## 2023-09-16 21:34:28 Annotating text fragment 6951/42019
## 2023-09-16 21:34:28 Annotating text fragment 6961/42019
## 2023-09-16 21:34:28 Annotating text fragment 6971/42019
## 2023-09-16 21:34:28 Annotating text fragment 6981/42019
## 2023-09-16 21:34:28 Annotating text fragment 6991/42019
## 2023-09-16 21:34:29 Annotating text fragment 7001/42019
## 2023-09-16 21:34:29 Annotating text fragment 7011/42019
## 2023-09-16 21:34:29 Annotating text fragment 7021/42019
## 2023-09-16 21:34:29 Annotating text fragment 7031/42019
## 2023-09-16 21:34:29 Annotating text fragment 7041/42019
## 2023-09-16 21:34:29 Annotating text fragment 7051/42019
## 2023-09-16 21:34:29 Annotating text fragment 7061/42019
## 2023-09-16 21:34:30 Annotating text fragment 7071/42019
## 2023-09-16 21:34:30 Annotating text fragment 7081/42019
## 2023-09-16 21:34:30 Annotating text fragment 7091/42019
## 2023-09-16 21:34:30 Annotating text fragment 7101/42019
## 2023-09-16 21:34:30 Annotating text fragment 7111/42019
## 2023-09-16 21:34:30 Annotating text fragment 7121/42019
## 2023-09-16 21:34:30 Annotating text fragment 7131/42019
## 2023-09-16 21:34:30 Annotating text fragment 7141/42019
## 2023-09-16 21:34:30 Annotating text fragment 7151/42019
## 2023-09-16 21:34:30 Annotating text fragment 7161/42019
```

```
## 2023-09-16 21:34:31 Annotating text fragment 7171/42019
## 2023-09-16 21:34:31 Annotating text fragment 7181/42019
## 2023-09-16 21:34:31 Annotating text fragment 7191/42019
## 2023-09-16 21:34:31 Annotating text fragment 7201/42019
## 2023-09-16 21:34:31 Annotating text fragment 7211/42019
## 2023-09-16 21:34:32 Annotating text fragment 7221/42019
## 2023-09-16 21:34:32 Annotating text fragment 7231/42019
## 2023-09-16 21:34:32 Annotating text fragment 7241/42019
## 2023-09-16 21:34:32 Annotating text fragment 7251/42019
## 2023-09-16 21:34:32 Annotating text fragment 7261/42019
## 2023-09-16 21:34:32 Annotating text fragment 7271/42019
## 2023-09-16 21:34:33 Annotating text fragment 7281/42019
## 2023-09-16 21:34:33 Annotating text fragment 7291/42019
## 2023-09-16 21:34:33 Annotating text fragment 7301/42019
## 2023-09-16 21:34:33 Annotating text fragment 7311/42019
## 2023-09-16 21:34:33 Annotating text fragment 7321/42019
## 2023-09-16 21:34:33 Annotating text fragment 7331/42019
## 2023-09-16 21:34:33 Annotating text fragment 7341/42019
## 2023-09-16 21:34:33 Annotating text fragment 7351/42019
## 2023-09-16 21:34:34 Annotating text fragment 7361/42019
## 2023-09-16 21:34:34 Annotating text fragment 7371/42019
## 2023-09-16 21:34:34 Annotating text fragment 7381/42019
## 2023-09-16 21:34:34 Annotating text fragment 7391/42019
## 2023-09-16 21:34:34 Annotating text fragment 7401/42019
## 2023-09-16 21:34:34 Annotating text fragment 7411/42019
## 2023-09-16 21:34:34 Annotating text fragment 7421/42019
## 2023-09-16 21:34:35 Annotating text fragment 7431/42019
## 2023-09-16 21:34:35 Annotating text fragment 7441/42019
## 2023-09-16 21:34:35 Annotating text fragment 7451/42019
## 2023-09-16 21:34:35 Annotating text fragment 7461/42019
## 2023-09-16 21:34:35 Annotating text fragment 7471/42019
## 2023-09-16 21:34:35 Annotating text fragment 7481/42019
## 2023-09-16 21:34:35 Annotating text fragment 7491/42019
## 2023-09-16 21:34:35 Annotating text fragment 7501/42019
## 2023-09-16 21:34:35 Annotating text fragment 7511/42019
## 2023-09-16 21:34:35 Annotating text fragment 7521/42019
## 2023-09-16 21:34:36 Annotating text fragment 7531/42019
## 2023-09-16 21:34:36 Annotating text fragment 7541/42019
## 2023-09-16 21:34:36 Annotating text fragment 7551/42019
## 2023-09-16 21:34:36 Annotating text fragment 7561/42019
## 2023-09-16 21:34:36 Annotating text fragment 7571/42019
## 2023-09-16 21:34:36 Annotating text fragment 7581/42019
## 2023-09-16 21:34:36 Annotating text fragment 7591/42019
## 2023-09-16 21:34:36 Annotating text fragment 7601/42019
## 2023-09-16 21:34:36 Annotating text fragment 7611/42019
## 2023-09-16 21:34:36 Annotating text fragment 7621/42019
## 2023-09-16 21:34:36 Annotating text fragment 7631/42019
## 2023-09-16 21:34:36 Annotating text fragment 7641/42019
## 2023-09-16 21:34:37 Annotating text fragment 7651/42019
## 2023-09-16 21:34:37 Annotating text fragment 7661/42019
## 2023-09-16 21:34:37 Annotating text fragment 7671/42019
## 2023-09-16 21:34:37 Annotating text fragment 7681/42019
## 2023-09-16 21:34:37 Annotating text fragment 7691/42019
## 2023-09-16 21:34:37 Annotating text fragment 7701/42019
```

```
## 2023-09-16 21:34:37 Annotating text fragment 7711/42019
## 2023-09-16 21:34:37 Annotating text fragment 7721/42019
## 2023-09-16 21:34:37 Annotating text fragment 7731/42019
## 2023-09-16 21:34:37 Annotating text fragment 7741/42019
## 2023-09-16 21:34:37 Annotating text fragment 7751/42019
## 2023-09-16 21:34:38 Annotating text fragment 7761/42019
## 2023-09-16 21:34:38 Annotating text fragment 7771/42019
## 2023-09-16 21:34:38 Annotating text fragment 7781/42019
## 2023-09-16 21:34:38 Annotating text fragment 7791/42019
## 2023-09-16 21:34:38 Annotating text fragment 7801/42019
## 2023-09-16 21:34:38 Annotating text fragment 7811/42019
## 2023-09-16 21:34:38 Annotating text fragment 7821/42019
## 2023-09-16 21:34:39 Annotating text fragment 7831/42019
## 2023-09-16 21:34:39 Annotating text fragment 7841/42019
## 2023-09-16 21:34:39 Annotating text fragment 7851/42019
## 2023-09-16 21:34:39 Annotating text fragment 7861/42019
## 2023-09-16 21:34:39 Annotating text fragment 7871/42019
## 2023-09-16 21:34:39 Annotating text fragment 7881/42019
## 2023-09-16 21:34:39 Annotating text fragment 7891/42019
## 2023-09-16 21:34:39 Annotating text fragment 7901/42019
## 2023-09-16 21:34:39 Annotating text fragment 7911/42019
## 2023-09-16 21:34:39 Annotating text fragment 7921/42019
## 2023-09-16 21:34:39 Annotating text fragment 7931/42019
## 2023-09-16 21:34:39 Annotating text fragment 7941/42019
## 2023-09-16 21:34:39 Annotating text fragment 7951/42019
## 2023-09-16 21:34:39 Annotating text fragment 7961/42019
## 2023-09-16 21:34:40 Annotating text fragment 7971/42019
## 2023-09-16 21:34:40 Annotating text fragment 7981/42019
## 2023-09-16 21:34:40 Annotating text fragment 7991/42019
## 2023-09-16 21:34:40 Annotating text fragment 8001/42019
## 2023-09-16 21:34:40 Annotating text fragment 8011/42019
## 2023-09-16 21:34:40 Annotating text fragment 8021/42019
## 2023-09-16 21:34:40 Annotating text fragment 8031/42019
## 2023-09-16 21:34:40 Annotating text fragment 8041/42019
## 2023-09-16 21:34:40 Annotating text fragment 8051/42019
## 2023-09-16 21:34:40 Annotating text fragment 8061/42019
## 2023-09-16 21:34:40 Annotating text fragment 8071/42019
## 2023-09-16 21:34:40 Annotating text fragment 8081/42019
## 2023-09-16 21:34:40 Annotating text fragment 8091/42019
## 2023-09-16 21:34:40 Annotating text fragment 8101/42019
## 2023-09-16 21:34:41 Annotating text fragment 8111/42019
## 2023-09-16 21:34:41 Annotating text fragment 8121/42019
## 2023-09-16 21:34:41 Annotating text fragment 8131/42019
## 2023-09-16 21:34:41 Annotating text fragment 8141/42019
## 2023-09-16 21:34:41 Annotating text fragment 8151/42019
## 2023-09-16 21:34:41 Annotating text fragment 8161/42019
## 2023-09-16 21:34:41 Annotating text fragment 8171/42019
## 2023-09-16 21:34:41 Annotating text fragment 8181/42019
## 2023-09-16 21:34:42 Annotating text fragment 8191/42019
## 2023-09-16 21:34:42 Annotating text fragment 8201/42019
## 2023-09-16 21:34:42 Annotating text fragment 8211/42019
## 2023-09-16 21:34:42 Annotating text fragment 8221/42019
## 2023-09-16 21:34:42 Annotating text fragment 8231/42019
## 2023-09-16 21:34:42 Annotating text fragment 8241/42019
```

```
## 2023-09-16 21:34:42 Annotating text fragment 8251/42019
## 2023-09-16 21:34:42 Annotating text fragment 8261/42019
## 2023-09-16 21:34:42 Annotating text fragment 8271/42019
## 2023-09-16 21:34:42 Annotating text fragment 8281/42019
## 2023-09-16 21:34:42 Annotating text fragment 8291/42019
## 2023-09-16 21:34:42 Annotating text fragment 8301/42019
## 2023-09-16 21:34:42 Annotating text fragment 8311/42019
## 2023-09-16 21:34:43 Annotating text fragment 8321/42019
## 2023-09-16 21:34:43 Annotating text fragment 8331/42019
## 2023-09-16 21:34:43 Annotating text fragment 8341/42019
## 2023-09-16 21:34:43 Annotating text fragment 8351/42019
## 2023-09-16 21:34:43 Annotating text fragment 8361/42019
## 2023-09-16 21:34:43 Annotating text fragment 8371/42019
## 2023-09-16 21:34:43 Annotating text fragment 8381/42019
## 2023-09-16 21:34:43 Annotating text fragment 8391/42019
## 2023-09-16 21:34:43 Annotating text fragment 8401/42019
## 2023-09-16 21:34:43 Annotating text fragment 8411/42019
## 2023-09-16 21:34:43 Annotating text fragment 8421/42019
## 2023-09-16 21:34:43 Annotating text fragment 8431/42019
## 2023-09-16 21:34:44 Annotating text fragment 8441/42019
## 2023-09-16 21:34:44 Annotating text fragment 8451/42019
## 2023-09-16 21:34:44 Annotating text fragment 8461/42019
## 2023-09-16 21:34:44 Annotating text fragment 8471/42019
## 2023-09-16 21:34:44 Annotating text fragment 8481/42019
## 2023-09-16 21:34:44 Annotating text fragment 8491/42019
## 2023-09-16 21:34:45 Annotating text fragment 8501/42019
## 2023-09-16 21:34:45 Annotating text fragment 8511/42019
## 2023-09-16 21:34:45 Annotating text fragment 8521/42019
## 2023-09-16 21:34:45 Annotating text fragment 8531/42019
## 2023-09-16 21:34:45 Annotating text fragment 8541/42019
## 2023-09-16 21:34:45 Annotating text fragment 8551/42019
## 2023-09-16 21:34:45 Annotating text fragment 8561/42019
## 2023-09-16 21:34:45 Annotating text fragment 8571/42019
## 2023-09-16 21:34:45 Annotating text fragment 8581/42019
## 2023-09-16 21:34:45 Annotating text fragment 8591/42019
## 2023-09-16 21:34:45 Annotating text fragment 8601/42019
## 2023-09-16 21:34:46 Annotating text fragment 8611/42019
## 2023-09-16 21:34:46 Annotating text fragment 8621/42019
## 2023-09-16 21:34:46 Annotating text fragment 8631/42019
## 2023-09-16 21:34:46 Annotating text fragment 8641/42019
## 2023-09-16 21:34:46 Annotating text fragment 8651/42019
## 2023-09-16 21:34:46 Annotating text fragment 8661/42019
## 2023-09-16 21:34:46 Annotating text fragment 8671/42019
## 2023-09-16 21:34:46 Annotating text fragment 8681/42019
## 2023-09-16 21:34:46 Annotating text fragment 8691/42019
## 2023-09-16 21:34:46 Annotating text fragment 8701/42019
## 2023-09-16 21:34:46 Annotating text fragment 8711/42019
## 2023-09-16 21:34:47 Annotating text fragment 8721/42019
## 2023-09-16 21:34:47 Annotating text fragment 8731/42019
## 2023-09-16 21:34:47 Annotating text fragment 8741/42019
## 2023-09-16 21:34:47 Annotating text fragment 8751/42019
## 2023-09-16 21:34:47 Annotating text fragment 8761/42019
## 2023-09-16 21:34:47 Annotating text fragment 8771/42019
## 2023-09-16 21:34:47 Annotating text fragment 8781/42019
```

```
## 2023-09-16 21:34:47 Annotating text fragment 8791/42019
## 2023-09-16 21:34:47 Annotating text fragment 8801/42019
## 2023-09-16 21:34:48 Annotating text fragment 8811/42019
## 2023-09-16 21:34:48 Annotating text fragment 8821/42019
## 2023-09-16 21:34:48 Annotating text fragment 8831/42019
## 2023-09-16 21:34:48 Annotating text fragment 8841/42019
## 2023-09-16 21:34:48 Annotating text fragment 8851/42019
## 2023-09-16 21:34:48 Annotating text fragment 8861/42019
## 2023-09-16 21:34:48 Annotating text fragment 8871/42019
## 2023-09-16 21:34:48 Annotating text fragment 8881/42019
## 2023-09-16 21:34:48 Annotating text fragment 8891/42019
## 2023-09-16 21:34:48 Annotating text fragment 8901/42019
## 2023-09-16 21:34:48 Annotating text fragment 8911/42019
## 2023-09-16 21:34:48 Annotating text fragment 8921/42019
## 2023-09-16 21:34:49 Annotating text fragment 8931/42019
## 2023-09-16 21:34:49 Annotating text fragment 8941/42019
## 2023-09-16 21:34:49 Annotating text fragment 8951/42019
## 2023-09-16 21:34:49 Annotating text fragment 8961/42019
## 2023-09-16 21:34:49 Annotating text fragment 8971/42019
## 2023-09-16 21:34:49 Annotating text fragment 8981/42019
## 2023-09-16 21:34:49 Annotating text fragment 8991/42019
## 2023-09-16 21:34:49 Annotating text fragment 9001/42019
## 2023-09-16 21:34:49 Annotating text fragment 9011/42019
## 2023-09-16 21:34:49 Annotating text fragment 9021/42019
## 2023-09-16 21:34:49 Annotating text fragment 9031/42019
## 2023-09-16 21:34:50 Annotating text fragment 9041/42019
## 2023-09-16 21:34:50 Annotating text fragment 9051/42019
## 2023-09-16 21:34:50 Annotating text fragment 9061/42019
## 2023-09-16 21:34:50 Annotating text fragment 9071/42019
## 2023-09-16 21:34:50 Annotating text fragment 9081/42019
## 2023-09-16 21:34:50 Annotating text fragment 9091/42019
## 2023-09-16 21:34:50 Annotating text fragment 9101/42019
## 2023-09-16 21:34:50 Annotating text fragment 9111/42019
## 2023-09-16 21:34:50 Annotating text fragment 9121/42019
## 2023-09-16 21:34:50 Annotating text fragment 9131/42019
## 2023-09-16 21:34:50 Annotating text fragment 9141/42019
## 2023-09-16 21:34:50 Annotating text fragment 9151/42019
## 2023-09-16 21:34:51 Annotating text fragment 9161/42019
## 2023-09-16 21:34:51 Annotating text fragment 9171/42019
## 2023-09-16 21:34:51 Annotating text fragment 9181/42019
## 2023-09-16 21:34:51 Annotating text fragment 9191/42019
## 2023-09-16 21:34:51 Annotating text fragment 9201/42019
## 2023-09-16 21:34:51 Annotating text fragment 9211/42019
## 2023-09-16 21:34:51 Annotating text fragment 9221/42019
## 2023-09-16 21:34:51 Annotating text fragment 9231/42019
## 2023-09-16 21:34:51 Annotating text fragment 9241/42019
## 2023-09-16 21:34:51 Annotating text fragment 9251/42019
## 2023-09-16 21:34:52 Annotating text fragment 9261/42019
## 2023-09-16 21:34:52 Annotating text fragment 9271/42019
## 2023-09-16 21:34:52 Annotating text fragment 9281/42019
## 2023-09-16 21:34:52 Annotating text fragment 9291/42019
## 2023-09-16 21:34:52 Annotating text fragment 9301/42019
## 2023-09-16 21:34:52 Annotating text fragment 9311/42019
## 2023-09-16 21:34:52 Annotating text fragment 9321/42019
```

```
## 2023-09-16 21:34:52 Annotating text fragment 9331/42019
## 2023-09-16 21:34:52 Annotating text fragment 9341/42019
## 2023-09-16 21:34:52 Annotating text fragment 9351/42019
## 2023-09-16 21:34:52 Annotating text fragment 9361/42019
## 2023-09-16 21:34:52 Annotating text fragment 9371/42019
## 2023-09-16 21:34:52 Annotating text fragment 9381/42019
## 2023-09-16 21:34:53 Annotating text fragment 9391/42019
## 2023-09-16 21:34:53 Annotating text fragment 9401/42019
## 2023-09-16 21:34:53 Annotating text fragment 9411/42019
## 2023-09-16 21:34:53 Annotating text fragment 9421/42019
## 2023-09-16 21:34:53 Annotating text fragment 9431/42019
## 2023-09-16 21:34:53 Annotating text fragment 9441/42019
## 2023-09-16 21:34:53 Annotating text fragment 9451/42019
## 2023-09-16 21:34:53 Annotating text fragment 9461/42019
## 2023-09-16 21:34:54 Annotating text fragment 9471/42019
## 2023-09-16 21:34:54 Annotating text fragment 9481/42019
## 2023-09-16 21:34:54 Annotating text fragment 9491/42019
## 2023-09-16 21:34:54 Annotating text fragment 9501/42019
## 2023-09-16 21:34:54 Annotating text fragment 9511/42019
## 2023-09-16 21:34:54 Annotating text fragment 9521/42019
## 2023-09-16 21:34:54 Annotating text fragment 9531/42019
## 2023-09-16 21:34:54 Annotating text fragment 9541/42019
## 2023-09-16 21:34:54 Annotating text fragment 9551/42019
## 2023-09-16 21:34:54 Annotating text fragment 9561/42019
## 2023-09-16 21:34:54 Annotating text fragment 9571/42019
## 2023-09-16 21:34:55 Annotating text fragment 9581/42019
## 2023-09-16 21:34:55 Annotating text fragment 9591/42019
## 2023-09-16 21:34:56 Annotating text fragment 9601/42019
## 2023-09-16 21:34:56 Annotating text fragment 9611/42019
## 2023-09-16 21:34:56 Annotating text fragment 9621/42019
## 2023-09-16 21:34:56 Annotating text fragment 9631/42019
## 2023-09-16 21:34:57 Annotating text fragment 9641/42019
## 2023-09-16 21:34:57 Annotating text fragment 9651/42019
## 2023-09-16 21:34:57 Annotating text fragment 9661/42019
## 2023-09-16 21:34:57 Annotating text fragment 9671/42019
## 2023-09-16 21:34:57 Annotating text fragment 9681/42019
## 2023-09-16 21:34:57 Annotating text fragment 9691/42019
## 2023-09-16 21:34:57 Annotating text fragment 9701/42019
## 2023-09-16 21:34:57 Annotating text fragment 9711/42019
## 2023-09-16 21:34:57 Annotating text fragment 9721/42019
## 2023-09-16 21:34:57 Annotating text fragment 9731/42019
## 2023-09-16 21:34:57 Annotating text fragment 9741/42019
## 2023-09-16 21:34:57 Annotating text fragment 9751/42019
## 2023-09-16 21:34:57 Annotating text fragment 9761/42019
## 2023-09-16 21:34:57 Annotating text fragment 9771/42019
## 2023-09-16 21:34:57 Annotating text fragment 9781/42019
## 2023-09-16 21:34:57 Annotating text fragment 9791/42019
## 2023-09-16 21:34:57 Annotating text fragment 9801/42019
## 2023-09-16 21:34:57 Annotating text fragment 9811/42019
## 2023-09-16 21:34:58 Annotating text fragment 9821/42019
## 2023-09-16 21:34:58 Annotating text fragment 9831/42019
## 2023-09-16 21:34:58 Annotating text fragment 9841/42019
## 2023-09-16 21:34:58 Annotating text fragment 9851/42019
## 2023-09-16 21:34:58 Annotating text fragment 9861/42019
```

```
## 2023-09-16 21:34:58 Annotating text fragment 9871/42019
## 2023-09-16 21:34:58 Annotating text fragment 9881/42019
## 2023-09-16 21:34:58 Annotating text fragment 9891/42019
## 2023-09-16 21:34:58 Annotating text fragment 9901/42019
## 2023-09-16 21:34:58 Annotating text fragment 9911/42019
## 2023-09-16 21:34:59 Annotating text fragment 9921/42019
## 2023-09-16 21:34:59 Annotating text fragment 9931/42019
## 2023-09-16 21:34:59 Annotating text fragment 9941/42019
## 2023-09-16 21:34:59 Annotating text fragment 9951/42019
## 2023-09-16 21:34:59 Annotating text fragment 9961/42019
## 2023-09-16 21:34:59 Annotating text fragment 9971/42019
## 2023-09-16 21:34:59 Annotating text fragment 9981/42019
## 2023-09-16 21:35:00 Annotating text fragment 9991/42019
## 2023-09-16 21:35:00 Annotating text fragment 10001/42019
## 2023-09-16 21:35:00 Annotating text fragment 10011/42019
## 2023-09-16 21:35:00 Annotating text fragment 10021/42019
## 2023-09-16 21:35:00 Annotating text fragment 10031/42019
## 2023-09-16 21:35:00 Annotating text fragment 10041/42019
## 2023-09-16 21:35:00 Annotating text fragment 10051/42019
## 2023-09-16 21:35:00 Annotating text fragment 10061/42019
## 2023-09-16 21:35:00 Annotating text fragment 10071/42019
## 2023-09-16 21:35:00 Annotating text fragment 10081/42019
## 2023-09-16 21:35:00 Annotating text fragment 10091/42019
## 2023-09-16 21:35:00 Annotating text fragment 10101/42019
## 2023-09-16 21:35:00 Annotating text fragment 10111/42019
## 2023-09-16 21:35:00 Annotating text fragment 10121/42019
## 2023-09-16 21:35:00 Annotating text fragment 10131/42019
## 2023-09-16 21:35:01 Annotating text fragment 10141/42019
## 2023-09-16 21:35:01 Annotating text fragment 10151/42019
## 2023-09-16 21:35:02 Annotating text fragment 10161/42019
## 2023-09-16 21:35:02 Annotating text fragment 10171/42019
## 2023-09-16 21:35:02 Annotating text fragment 10181/42019
## 2023-09-16 21:35:02 Annotating text fragment 10191/42019
## 2023-09-16 21:35:02 Annotating text fragment 10201/42019
## 2023-09-16 21:35:02 Annotating text fragment 10211/42019
## 2023-09-16 21:35:02 Annotating text fragment 10221/42019
## 2023-09-16 21:35:02 Annotating text fragment 10231/42019
## 2023-09-16 21:35:02 Annotating text fragment 10241/42019
## 2023-09-16 21:35:02 Annotating text fragment 10251/42019
## 2023-09-16 21:35:02 Annotating text fragment 10261/42019
## 2023-09-16 21:35:03 Annotating text fragment 10271/42019
## 2023-09-16 21:35:03 Annotating text fragment 10281/42019
## 2023-09-16 21:35:03 Annotating text fragment 10291/42019
## 2023-09-16 21:35:03 Annotating text fragment 10301/42019
## 2023-09-16 21:35:03 Annotating text fragment 10311/42019
## 2023-09-16 21:35:03 Annotating text fragment 10321/42019
## 2023-09-16 21:35:03 Annotating text fragment 10331/42019
## 2023-09-16 21:35:03 Annotating text fragment 10341/42019
## 2023-09-16 21:35:03 Annotating text fragment 10351/42019
## 2023-09-16 21:35:04 Annotating text fragment 10361/42019
## 2023-09-16 21:35:04 Annotating text fragment 10371/42019
## 2023-09-16 21:35:04 Annotating text fragment 10381/42019
## 2023-09-16 21:35:04 Annotating text fragment 10391/42019
## 2023-09-16 21:35:04 Annotating text fragment 10401/42019
```

```
## 2023-09-16 21:35:04 Annotating text fragment 10411/42019
## 2023-09-16 21:35:04 Annotating text fragment 10421/42019
## 2023-09-16 21:35:04 Annotating text fragment 10431/42019
## 2023-09-16 21:35:04 Annotating text fragment 10441/42019
## 2023-09-16 21:35:04 Annotating text fragment 10451/42019
## 2023-09-16 21:35:04 Annotating text fragment 10461/42019
## 2023-09-16 21:35:04 Annotating text fragment 10471/42019
## 2023-09-16 21:35:05 Annotating text fragment 10481/42019
## 2023-09-16 21:35:05 Annotating text fragment 10491/42019
## 2023-09-16 21:35:05 Annotating text fragment 10501/42019
## 2023-09-16 21:35:05 Annotating text fragment 10511/42019
## 2023-09-16 21:35:05 Annotating text fragment 10521/42019
## 2023-09-16 21:35:05 Annotating text fragment 10531/42019
## 2023-09-16 21:35:05 Annotating text fragment 10541/42019
## 2023-09-16 21:35:05 Annotating text fragment 10551/42019
## 2023-09-16 21:35:05 Annotating text fragment 10561/42019
## 2023-09-16 21:35:05 Annotating text fragment 10571/42019
## 2023-09-16 21:35:05 Annotating text fragment 10581/42019
## 2023-09-16 21:35:06 Annotating text fragment 10591/42019
## 2023-09-16 21:35:06 Annotating text fragment 10601/42019
## 2023-09-16 21:35:06 Annotating text fragment 10611/42019
## 2023-09-16 21:35:06 Annotating text fragment 10621/42019
## 2023-09-16 21:35:06 Annotating text fragment 10631/42019
## 2023-09-16 21:35:06 Annotating text fragment 10641/42019
## 2023-09-16 21:35:06 Annotating text fragment 10651/42019
## 2023-09-16 21:35:06 Annotating text fragment 10661/42019
## 2023-09-16 21:35:06 Annotating text fragment 10671/42019
## 2023-09-16 21:35:07 Annotating text fragment 10681/42019
## 2023-09-16 21:35:07 Annotating text fragment 10691/42019
## 2023-09-16 21:35:07 Annotating text fragment 10701/42019
## 2023-09-16 21:35:07 Annotating text fragment 10711/42019
## 2023-09-16 21:35:07 Annotating text fragment 10721/42019
## 2023-09-16 21:35:07 Annotating text fragment 10731/42019
## 2023-09-16 21:35:07 Annotating text fragment 10741/42019
## 2023-09-16 21:35:07 Annotating text fragment 10751/42019
## 2023-09-16 21:35:07 Annotating text fragment 10761/42019
## 2023-09-16 21:35:07 Annotating text fragment 10771/42019
## 2023-09-16 21:35:07 Annotating text fragment 10781/42019
## 2023-09-16 21:35:07 Annotating text fragment 10791/42019
## 2023-09-16 21:35:07 Annotating text fragment 10801/42019
## 2023-09-16 21:35:07 Annotating text fragment 10811/42019
## 2023-09-16 21:35:07 Annotating text fragment 10821/42019
## 2023-09-16 21:35:07 Annotating text fragment 10831/42019
## 2023-09-16 21:35:08 Annotating text fragment 10841/42019
## 2023-09-16 21:35:08 Annotating text fragment 10851/42019
## 2023-09-16 21:35:08 Annotating text fragment 10861/42019
## 2023-09-16 21:35:08 Annotating text fragment 10871/42019
## 2023-09-16 21:35:08 Annotating text fragment 10881/42019
## 2023-09-16 21:35:08 Annotating text fragment 10891/42019
## 2023-09-16 21:35:08 Annotating text fragment 10901/42019
## 2023-09-16 21:35:08 Annotating text fragment 10911/42019
## 2023-09-16 21:35:08 Annotating text fragment 10921/42019
## 2023-09-16 21:35:08 Annotating text fragment 10931/42019
## 2023-09-16 21:35:08 Annotating text fragment 10941/42019
```

```
## 2023-09-16 21:35:08 Annotating text fragment 10951/42019
## 2023-09-16 21:35:08 Annotating text fragment 10961/42019
## 2023-09-16 21:35:08 Annotating text fragment 10971/42019
## 2023-09-16 21:35:08 Annotating text fragment 10981/42019
## 2023-09-16 21:35:08 Annotating text fragment 10991/42019
## 2023-09-16 21:35:08 Annotating text fragment 11001/42019
## 2023-09-16 21:35:08 Annotating text fragment 11011/42019
## 2023-09-16 21:35:09 Annotating text fragment 11021/42019
## 2023-09-16 21:35:09 Annotating text fragment 11031/42019
## 2023-09-16 21:35:09 Annotating text fragment 11041/42019
## 2023-09-16 21:35:09 Annotating text fragment 11051/42019
## 2023-09-16 21:35:09 Annotating text fragment 11061/42019
## 2023-09-16 21:35:09 Annotating text fragment 11071/42019
## 2023-09-16 21:35:09 Annotating text fragment 11081/42019
## 2023-09-16 21:35:09 Annotating text fragment 11091/42019
## 2023-09-16 21:35:09 Annotating text fragment 11101/42019
## 2023-09-16 21:35:09 Annotating text fragment 11111/42019
## 2023-09-16 21:35:09 Annotating text fragment 11121/42019
## 2023-09-16 21:35:10 Annotating text fragment 11131/42019
## 2023-09-16 21:35:10 Annotating text fragment 11141/42019
## 2023-09-16 21:35:10 Annotating text fragment 11151/42019
## 2023-09-16 21:35:10 Annotating text fragment 11161/42019
## 2023-09-16 21:35:10 Annotating text fragment 11171/42019
## 2023-09-16 21:35:10 Annotating text fragment 11181/42019
## 2023-09-16 21:35:10 Annotating text fragment 11191/42019
## 2023-09-16 21:35:10 Annotating text fragment 11201/42019
## 2023-09-16 21:35:10 Annotating text fragment 11211/42019
## 2023-09-16 21:35:11 Annotating text fragment 11221/42019
## 2023-09-16 21:35:11 Annotating text fragment 11231/42019
## 2023-09-16 21:35:11 Annotating text fragment 11241/42019
## 2023-09-16 21:35:11 Annotating text fragment 11251/42019
## 2023-09-16 21:35:11 Annotating text fragment 11261/42019
## 2023-09-16 21:35:11 Annotating text fragment 11271/42019
## 2023-09-16 21:35:11 Annotating text fragment 11281/42019
## 2023-09-16 21:35:11 Annotating text fragment 11291/42019
## 2023-09-16 21:35:11 Annotating text fragment 11301/42019
## 2023-09-16 21:35:11 Annotating text fragment 11311/42019
## 2023-09-16 21:35:12 Annotating text fragment 11321/42019
## 2023-09-16 21:35:12 Annotating text fragment 11331/42019
## 2023-09-16 21:35:12 Annotating text fragment 11341/42019
## 2023-09-16 21:35:12 Annotating text fragment 11351/42019
## 2023-09-16 21:35:12 Annotating text fragment 11361/42019
## 2023-09-16 21:35:12 Annotating text fragment 11371/42019
## 2023-09-16 21:35:12 Annotating text fragment 11381/42019
## 2023-09-16 21:35:12 Annotating text fragment 11391/42019
## 2023-09-16 21:35:12 Annotating text fragment 11401/42019
## 2023-09-16 21:35:12 Annotating text fragment 11411/42019
## 2023-09-16 21:35:13 Annotating text fragment 11421/42019
## 2023-09-16 21:35:13 Annotating text fragment 11431/42019
## 2023-09-16 21:35:13 Annotating text fragment 11441/42019
## 2023-09-16 21:35:13 Annotating text fragment 11451/42019
## 2023-09-16 21:35:13 Annotating text fragment 11461/42019
## 2023-09-16 21:35:13 Annotating text fragment 11471/42019
## 2023-09-16 21:35:13 Annotating text fragment 11481/42019
```

```
## 2023-09-16 21:35:13 Annotating text fragment 11491/42019
## 2023-09-16 21:35:13 Annotating text fragment 11501/42019
## 2023-09-16 21:35:13 Annotating text fragment 11511/42019
## 2023-09-16 21:35:14 Annotating text fragment 11521/42019
## 2023-09-16 21:35:14 Annotating text fragment 11531/42019
## 2023-09-16 21:35:15 Annotating text fragment 11541/42019
## 2023-09-16 21:35:15 Annotating text fragment 11551/42019
## 2023-09-16 21:35:16 Annotating text fragment 11561/42019
## 2023-09-16 21:35:16 Annotating text fragment 11571/42019
## 2023-09-16 21:35:16 Annotating text fragment 11581/42019
## 2023-09-16 21:35:16 Annotating text fragment 11591/42019
## 2023-09-16 21:35:16 Annotating text fragment 11601/42019
## 2023-09-16 21:35:16 Annotating text fragment 11611/42019
## 2023-09-16 21:35:16 Annotating text fragment 11621/42019
## 2023-09-16 21:35:16 Annotating text fragment 11631/42019
## 2023-09-16 21:35:17 Annotating text fragment 11641/42019
## 2023-09-16 21:35:17 Annotating text fragment 11651/42019
## 2023-09-16 21:35:17 Annotating text fragment 11661/42019
## 2023-09-16 21:35:17 Annotating text fragment 11671/42019
## 2023-09-16 21:35:17 Annotating text fragment 11681/42019
## 2023-09-16 21:35:17 Annotating text fragment 11691/42019
## 2023-09-16 21:35:17 Annotating text fragment 11701/42019
## 2023-09-16 21:35:17 Annotating text fragment 11711/42019
## 2023-09-16 21:35:17 Annotating text fragment 11721/42019
## 2023-09-16 21:35:17 Annotating text fragment 11731/42019
## 2023-09-16 21:35:17 Annotating text fragment 11741/42019
## 2023-09-16 21:35:17 Annotating text fragment 11751/42019
## 2023-09-16 21:35:17 Annotating text fragment 11761/42019
## 2023-09-16 21:35:17 Annotating text fragment 11771/42019
## 2023-09-16 21:35:18 Annotating text fragment 11781/42019
## 2023-09-16 21:35:18 Annotating text fragment 11791/42019
## 2023-09-16 21:35:18 Annotating text fragment 11801/42019
## 2023-09-16 21:35:18 Annotating text fragment 11811/42019
## 2023-09-16 21:35:18 Annotating text fragment 11821/42019
## 2023-09-16 21:35:18 Annotating text fragment 11831/42019
## 2023-09-16 21:35:18 Annotating text fragment 11841/42019
## 2023-09-16 21:35:18 Annotating text fragment 11851/42019
## 2023-09-16 21:35:18 Annotating text fragment 11861/42019
## 2023-09-16 21:35:18 Annotating text fragment 11871/42019
## 2023-09-16 21:35:18 Annotating text fragment 11881/42019
## 2023-09-16 21:35:18 Annotating text fragment 11891/42019
## 2023-09-16 21:35:19 Annotating text fragment 11901/42019
## 2023-09-16 21:35:19 Annotating text fragment 11911/42019
## 2023-09-16 21:35:19 Annotating text fragment 11921/42019
## 2023-09-16 21:35:19 Annotating text fragment 11931/42019
## 2023-09-16 21:35:19 Annotating text fragment 11941/42019
## 2023-09-16 21:35:19 Annotating text fragment 11951/42019
## 2023-09-16 21:35:19 Annotating text fragment 11961/42019
## 2023-09-16 21:35:19 Annotating text fragment 11971/42019
## 2023-09-16 21:35:19 Annotating text fragment 11981/42019
## 2023-09-16 21:35:19 Annotating text fragment 11991/42019
## 2023-09-16 21:35:19 Annotating text fragment 12001/42019
## 2023-09-16 21:35:20 Annotating text fragment 12011/42019
## 2023-09-16 21:35:20 Annotating text fragment 12021/42019
```

```
## 2023-09-16 21:35:20 Annotating text fragment 12031/42019
## 2023-09-16 21:35:20 Annotating text fragment 12041/42019
## 2023-09-16 21:35:20 Annotating text fragment 12051/42019
## 2023-09-16 21:35:20 Annotating text fragment 12061/42019
## 2023-09-16 21:35:20 Annotating text fragment 12071/42019
## 2023-09-16 21:35:20 Annotating text fragment 12081/42019
## 2023-09-16 21:35:20 Annotating text fragment 12091/42019
## 2023-09-16 21:35:20 Annotating text fragment 12101/42019
## 2023-09-16 21:35:20 Annotating text fragment 12111/42019
## 2023-09-16 21:35:20 Annotating text fragment 12121/42019
## 2023-09-16 21:35:20 Annotating text fragment 12131/42019
## 2023-09-16 21:35:20 Annotating text fragment 12141/42019
## 2023-09-16 21:35:20 Annotating text fragment 12151/42019
## 2023-09-16 21:35:20 Annotating text fragment 12161/42019
## 2023-09-16 21:35:21 Annotating text fragment 12171/42019
## 2023-09-16 21:35:21 Annotating text fragment 12181/42019
## 2023-09-16 21:35:21 Annotating text fragment 12191/42019
## 2023-09-16 21:35:21 Annotating text fragment 12201/42019
## 2023-09-16 21:35:21 Annotating text fragment 12211/42019
## 2023-09-16 21:35:21 Annotating text fragment 12221/42019
## 2023-09-16 21:35:22 Annotating text fragment 12231/42019
## 2023-09-16 21:35:22 Annotating text fragment 12241/42019
## 2023-09-16 21:35:22 Annotating text fragment 12251/42019
## 2023-09-16 21:35:22 Annotating text fragment 12261/42019
## 2023-09-16 21:35:22 Annotating text fragment 12271/42019
## 2023-09-16 21:35:22 Annotating text fragment 12281/42019
## 2023-09-16 21:35:22 Annotating text fragment 12291/42019
## 2023-09-16 21:35:22 Annotating text fragment 12301/42019
## 2023-09-16 21:35:22 Annotating text fragment 12311/42019
## 2023-09-16 21:35:22 Annotating text fragment 12321/42019
## 2023-09-16 21:35:22 Annotating text fragment 12331/42019
## 2023-09-16 21:35:23 Annotating text fragment 12341/42019
## 2023-09-16 21:35:23 Annotating text fragment 12351/42019
## 2023-09-16 21:35:23 Annotating text fragment 12361/42019
## 2023-09-16 21:35:23 Annotating text fragment 12371/42019
## 2023-09-16 21:35:23 Annotating text fragment 12381/42019
## 2023-09-16 21:35:23 Annotating text fragment 12391/42019
## 2023-09-16 21:35:23 Annotating text fragment 12401/42019
## 2023-09-16 21:35:23 Annotating text fragment 12411/42019
## 2023-09-16 21:35:23 Annotating text fragment 12421/42019
## 2023-09-16 21:35:23 Annotating text fragment 12431/42019
## 2023-09-16 21:35:23 Annotating text fragment 12441/42019
## 2023-09-16 21:35:23 Annotating text fragment 12451/42019
## 2023-09-16 21:35:23 Annotating text fragment 12461/42019
## 2023-09-16 21:35:23 Annotating text fragment 12471/42019
## 2023-09-16 21:35:23 Annotating text fragment 12481/42019
## 2023-09-16 21:35:23 Annotating text fragment 12491/42019
## 2023-09-16 21:35:24 Annotating text fragment 12501/42019
## 2023-09-16 21:35:24 Annotating text fragment 12511/42019
## 2023-09-16 21:35:24 Annotating text fragment 12521/42019
## 2023-09-16 21:35:24 Annotating text fragment 12531/42019
## 2023-09-16 21:35:24 Annotating text fragment 12541/42019
## 2023-09-16 21:35:24 Annotating text fragment 12551/42019
## 2023-09-16 21:35:24 Annotating text fragment 12561/42019
```

```
## 2023-09-16 21:35:24 Annotating text fragment 12571/42019
## 2023-09-16 21:35:24 Annotating text fragment 12581/42019
## 2023-09-16 21:35:24 Annotating text fragment 12591/42019
## 2023-09-16 21:35:25 Annotating text fragment 12601/42019
## 2023-09-16 21:35:25 Annotating text fragment 12611/42019
## 2023-09-16 21:35:25 Annotating text fragment 12621/42019
## 2023-09-16 21:35:25 Annotating text fragment 12631/42019
## 2023-09-16 21:35:25 Annotating text fragment 12641/42019
## 2023-09-16 21:35:25 Annotating text fragment 12651/42019
## 2023-09-16 21:35:25 Annotating text fragment 12661/42019
## 2023-09-16 21:35:25 Annotating text fragment 12671/42019
## 2023-09-16 21:35:25 Annotating text fragment 12681/42019
## 2023-09-16 21:35:26 Annotating text fragment 12691/42019
## 2023-09-16 21:35:26 Annotating text fragment 12701/42019
## 2023-09-16 21:35:26 Annotating text fragment 12711/42019
## 2023-09-16 21:35:26 Annotating text fragment 12721/42019
## 2023-09-16 21:35:26 Annotating text fragment 12731/42019
## 2023-09-16 21:35:26 Annotating text fragment 12741/42019
## 2023-09-16 21:35:26 Annotating text fragment 12751/42019
## 2023-09-16 21:35:26 Annotating text fragment 12761/42019
## 2023-09-16 21:35:26 Annotating text fragment 12771/42019
## 2023-09-16 21:35:26 Annotating text fragment 12781/42019
## 2023-09-16 21:35:27 Annotating text fragment 12791/42019
## 2023-09-16 21:35:27 Annotating text fragment 12801/42019
## 2023-09-16 21:35:27 Annotating text fragment 12811/42019
## 2023-09-16 21:35:27 Annotating text fragment 12821/42019
## 2023-09-16 21:35:27 Annotating text fragment 12831/42019
## 2023-09-16 21:35:27 Annotating text fragment 12841/42019
## 2023-09-16 21:35:27 Annotating text fragment 12851/42019
## 2023-09-16 21:35:27 Annotating text fragment 12861/42019
## 2023-09-16 21:35:27 Annotating text fragment 12871/42019
## 2023-09-16 21:35:27 Annotating text fragment 12881/42019
## 2023-09-16 21:35:27 Annotating text fragment 12891/42019
## 2023-09-16 21:35:27 Annotating text fragment 12901/42019
## 2023-09-16 21:35:27 Annotating text fragment 12911/42019
## 2023-09-16 21:35:28 Annotating text fragment 12921/42019
## 2023-09-16 21:35:28 Annotating text fragment 12931/42019
## 2023-09-16 21:35:28 Annotating text fragment 12941/42019
## 2023-09-16 21:35:28 Annotating text fragment 12951/42019
## 2023-09-16 21:35:28 Annotating text fragment 12961/42019
## 2023-09-16 21:35:28 Annotating text fragment 12971/42019
## 2023-09-16 21:35:28 Annotating text fragment 12981/42019
## 2023-09-16 21:35:28 Annotating text fragment 12991/42019
## 2023-09-16 21:35:28 Annotating text fragment 13001/42019
## 2023-09-16 21:35:28 Annotating text fragment 13011/42019
## 2023-09-16 21:35:28 Annotating text fragment 13021/42019
## 2023-09-16 21:35:28 Annotating text fragment 13031/42019
## 2023-09-16 21:35:28 Annotating text fragment 13041/42019
## 2023-09-16 21:35:28 Annotating text fragment 13051/42019
## 2023-09-16 21:35:28 Annotating text fragment 13061/42019
## 2023-09-16 21:35:29 Annotating text fragment 13071/42019
## 2023-09-16 21:35:29 Annotating text fragment 13081/42019
## 2023-09-16 21:35:29 Annotating text fragment 13091/42019
## 2023-09-16 21:35:29 Annotating text fragment 13101/42019
```

```
## 2023-09-16 21:35:29 Annotating text fragment 13111/42019
## 2023-09-16 21:35:29 Annotating text fragment 13121/42019
## 2023-09-16 21:35:29 Annotating text fragment 13131/42019
## 2023-09-16 21:35:29 Annotating text fragment 13141/42019
## 2023-09-16 21:35:29 Annotating text fragment 13151/42019
## 2023-09-16 21:35:29 Annotating text fragment 13161/42019
## 2023-09-16 21:35:29 Annotating text fragment 13171/42019
## 2023-09-16 21:35:29 Annotating text fragment 13181/42019
## 2023-09-16 21:35:29 Annotating text fragment 13191/42019
## 2023-09-16 21:35:29 Annotating text fragment 13201/42019
## 2023-09-16 21:35:30 Annotating text fragment 13211/42019
## 2023-09-16 21:35:30 Annotating text fragment 13221/42019
## 2023-09-16 21:35:30 Annotating text fragment 13231/42019
## 2023-09-16 21:35:30 Annotating text fragment 13241/42019
## 2023-09-16 21:35:30 Annotating text fragment 13251/42019
## 2023-09-16 21:35:30 Annotating text fragment 13261/42019
## 2023-09-16 21:35:31 Annotating text fragment 13271/42019
## 2023-09-16 21:35:31 Annotating text fragment 13281/42019
## 2023-09-16 21:35:31 Annotating text fragment 13291/42019
## 2023-09-16 21:35:31 Annotating text fragment 13301/42019
## 2023-09-16 21:35:31 Annotating text fragment 13311/42019
## 2023-09-16 21:35:31 Annotating text fragment 13321/42019
## 2023-09-16 21:35:32 Annotating text fragment 13331/42019
## 2023-09-16 21:35:32 Annotating text fragment 13341/42019
## 2023-09-16 21:35:32 Annotating text fragment 13351/42019
## 2023-09-16 21:35:32 Annotating text fragment 13361/42019
## 2023-09-16 21:35:32 Annotating text fragment 13371/42019
## 2023-09-16 21:35:32 Annotating text fragment 13381/42019
## 2023-09-16 21:35:32 Annotating text fragment 13391/42019
## 2023-09-16 21:35:32 Annotating text fragment 13401/42019
## 2023-09-16 21:35:32 Annotating text fragment 13411/42019
## 2023-09-16 21:35:32 Annotating text fragment 13421/42019
## 2023-09-16 21:35:32 Annotating text fragment 13431/42019
## 2023-09-16 21:35:33 Annotating text fragment 13441/42019
## 2023-09-16 21:35:33 Annotating text fragment 13451/42019
## 2023-09-16 21:35:33 Annotating text fragment 13461/42019
## 2023-09-16 21:35:33 Annotating text fragment 13471/42019
## 2023-09-16 21:35:33 Annotating text fragment 13481/42019
## 2023-09-16 21:35:33 Annotating text fragment 13491/42019
## 2023-09-16 21:35:33 Annotating text fragment 13501/42019
## 2023-09-16 21:35:33 Annotating text fragment 13511/42019
## 2023-09-16 21:35:34 Annotating text fragment 13521/42019
## 2023-09-16 21:35:34 Annotating text fragment 13531/42019
## 2023-09-16 21:35:34 Annotating text fragment 13541/42019
## 2023-09-16 21:35:35 Annotating text fragment 13551/42019
## 2023-09-16 21:35:35 Annotating text fragment 13561/42019
## 2023-09-16 21:35:35 Annotating text fragment 13571/42019
## 2023-09-16 21:35:35 Annotating text fragment 13581/42019
## 2023-09-16 21:35:35 Annotating text fragment 13591/42019
## 2023-09-16 21:35:35 Annotating text fragment 13601/42019
## 2023-09-16 21:35:35 Annotating text fragment 13611/42019
## 2023-09-16 21:35:35 Annotating text fragment 13621/42019
## 2023-09-16 21:35:35 Annotating text fragment 13631/42019
## 2023-09-16 21:35:36 Annotating text fragment 13641/42019
```

```
## 2023-09-16 21:35:36 Annotating text fragment 13651/42019
## 2023-09-16 21:35:36 Annotating text fragment 13661/42019
## 2023-09-16 21:35:36 Annotating text fragment 13671/42019
## 2023-09-16 21:35:36 Annotating text fragment 13681/42019
## 2023-09-16 21:35:36 Annotating text fragment 13691/42019
## 2023-09-16 21:35:36 Annotating text fragment 13701/42019
## 2023-09-16 21:35:36 Annotating text fragment 13711/42019
## 2023-09-16 21:35:36 Annotating text fragment 13721/42019
## 2023-09-16 21:35:36 Annotating text fragment 13731/42019
## 2023-09-16 21:35:37 Annotating text fragment 13741/42019
## 2023-09-16 21:35:37 Annotating text fragment 13751/42019
## 2023-09-16 21:35:37 Annotating text fragment 13761/42019
## 2023-09-16 21:35:37 Annotating text fragment 13771/42019
## 2023-09-16 21:35:37 Annotating text fragment 13781/42019
## 2023-09-16 21:35:37 Annotating text fragment 13791/42019
## 2023-09-16 21:35:37 Annotating text fragment 13801/42019
## 2023-09-16 21:35:37 Annotating text fragment 13811/42019
## 2023-09-16 21:35:37 Annotating text fragment 13821/42019
## 2023-09-16 21:35:37 Annotating text fragment 13831/42019
## 2023-09-16 21:35:37 Annotating text fragment 13841/42019
## 2023-09-16 21:35:38 Annotating text fragment 13851/42019
## 2023-09-16 21:35:38 Annotating text fragment 13861/42019
## 2023-09-16 21:35:38 Annotating text fragment 13871/42019
## 2023-09-16 21:35:38 Annotating text fragment 13881/42019
## 2023-09-16 21:35:38 Annotating text fragment 13891/42019
## 2023-09-16 21:35:38 Annotating text fragment 13901/42019
## 2023-09-16 21:35:38 Annotating text fragment 13911/42019
## 2023-09-16 21:35:39 Annotating text fragment 13921/42019
## 2023-09-16 21:35:39 Annotating text fragment 13931/42019
## 2023-09-16 21:35:39 Annotating text fragment 13941/42019
## 2023-09-16 21:35:39 Annotating text fragment 13951/42019
## 2023-09-16 21:35:39 Annotating text fragment 13961/42019
## 2023-09-16 21:35:39 Annotating text fragment 13971/42019
## 2023-09-16 21:35:39 Annotating text fragment 13981/42019
## 2023-09-16 21:35:39 Annotating text fragment 13991/42019
## 2023-09-16 21:35:39 Annotating text fragment 14001/42019
## 2023-09-16 21:35:39 Annotating text fragment 14011/42019
## 2023-09-16 21:35:39 Annotating text fragment 14021/42019
## 2023-09-16 21:35:39 Annotating text fragment 14031/42019
## 2023-09-16 21:35:39 Annotating text fragment 14041/42019
## 2023-09-16 21:35:39 Annotating text fragment 14051/42019
## 2023-09-16 21:35:40 Annotating text fragment 14061/42019
## 2023-09-16 21:35:40 Annotating text fragment 14071/42019
## 2023-09-16 21:35:40 Annotating text fragment 14081/42019
## 2023-09-16 21:35:40 Annotating text fragment 14091/42019
## 2023-09-16 21:35:40 Annotating text fragment 14101/42019
## 2023-09-16 21:35:40 Annotating text fragment 14111/42019
## 2023-09-16 21:35:40 Annotating text fragment 14121/42019
## 2023-09-16 21:35:40 Annotating text fragment 14131/42019
## 2023-09-16 21:35:40 Annotating text fragment 14141/42019
## 2023-09-16 21:35:40 Annotating text fragment 14151/42019
## 2023-09-16 21:35:40 Annotating text fragment 14161/42019
## 2023-09-16 21:35:40 Annotating text fragment 14171/42019
## 2023-09-16 21:35:40 Annotating text fragment 14181/42019
```

```
## 2023-09-16 21:35:40 Annotating text fragment 14191/42019
## 2023-09-16 21:35:40 Annotating text fragment 14201/42019
## 2023-09-16 21:35:41 Annotating text fragment 14211/42019
## 2023-09-16 21:35:41 Annotating text fragment 14221/42019
## 2023-09-16 21:35:41 Annotating text fragment 14231/42019
## 2023-09-16 21:35:41 Annotating text fragment 14241/42019
## 2023-09-16 21:35:41 Annotating text fragment 14251/42019
## 2023-09-16 21:35:41 Annotating text fragment 14261/42019
## 2023-09-16 21:35:41 Annotating text fragment 14271/42019
## 2023-09-16 21:35:41 Annotating text fragment 14281/42019
## 2023-09-16 21:35:41 Annotating text fragment 14291/42019
## 2023-09-16 21:35:41 Annotating text fragment 14301/42019
## 2023-09-16 21:35:41 Annotating text fragment 14311/42019
## 2023-09-16 21:35:41 Annotating text fragment 14321/42019
## 2023-09-16 21:35:41 Annotating text fragment 14331/42019
## 2023-09-16 21:35:41 Annotating text fragment 14341/42019
## 2023-09-16 21:35:41 Annotating text fragment 14351/42019
## 2023-09-16 21:35:41 Annotating text fragment 14361/42019
## 2023-09-16 21:35:41 Annotating text fragment 14371/42019
## 2023-09-16 21:35:41 Annotating text fragment 14381/42019
## 2023-09-16 21:35:42 Annotating text fragment 14391/42019
## 2023-09-16 21:35:42 Annotating text fragment 14401/42019
## 2023-09-16 21:35:42 Annotating text fragment 14411/42019
## 2023-09-16 21:35:42 Annotating text fragment 14421/42019
## 2023-09-16 21:35:42 Annotating text fragment 14431/42019
## 2023-09-16 21:35:42 Annotating text fragment 14441/42019
## 2023-09-16 21:35:42 Annotating text fragment 14451/42019
## 2023-09-16 21:35:42 Annotating text fragment 14461/42019
## 2023-09-16 21:35:42 Annotating text fragment 14471/42019
## 2023-09-16 21:35:42 Annotating text fragment 14481/42019
## 2023-09-16 21:35:42 Annotating text fragment 14491/42019
## 2023-09-16 21:35:42 Annotating text fragment 14501/42019
## 2023-09-16 21:35:42 Annotating text fragment 14511/42019
## 2023-09-16 21:35:42 Annotating text fragment 14521/42019
## 2023-09-16 21:35:42 Annotating text fragment 14531/42019
## 2023-09-16 21:35:43 Annotating text fragment 14541/42019
## 2023-09-16 21:35:43 Annotating text fragment 14551/42019
## 2023-09-16 21:35:43 Annotating text fragment 14561/42019
## 2023-09-16 21:35:43 Annotating text fragment 14571/42019
## 2023-09-16 21:35:43 Annotating text fragment 14581/42019
## 2023-09-16 21:35:43 Annotating text fragment 14591/42019
## 2023-09-16 21:35:43 Annotating text fragment 14601/42019
## 2023-09-16 21:35:43 Annotating text fragment 14611/42019
## 2023-09-16 21:35:43 Annotating text fragment 14621/42019
## 2023-09-16 21:35:43 Annotating text fragment 14631/42019
## 2023-09-16 21:35:43 Annotating text fragment 14641/42019
## 2023-09-16 21:35:43 Annotating text fragment 14651/42019
## 2023-09-16 21:35:43 Annotating text fragment 14661/42019
## 2023-09-16 21:35:44 Annotating text fragment 14671/42019
## 2023-09-16 21:35:44 Annotating text fragment 14681/42019
## 2023-09-16 21:35:44 Annotating text fragment 14691/42019
## 2023-09-16 21:35:44 Annotating text fragment 14701/42019
## 2023-09-16 21:35:44 Annotating text fragment 14711/42019
## 2023-09-16 21:35:44 Annotating text fragment 14721/42019
```

```
## 2023-09-16 21:35:44 Annotating text fragment 14731/42019
## 2023-09-16 21:35:44 Annotating text fragment 14741/42019
## 2023-09-16 21:35:44 Annotating text fragment 14751/42019
## 2023-09-16 21:35:44 Annotating text fragment 14761/42019
## 2023-09-16 21:35:44 Annotating text fragment 14771/42019
## 2023-09-16 21:35:44 Annotating text fragment 14781/42019
## 2023-09-16 21:35:44 Annotating text fragment 14791/42019
## 2023-09-16 21:35:44 Annotating text fragment 14801/42019
## 2023-09-16 21:35:45 Annotating text fragment 14811/42019
## 2023-09-16 21:35:45 Annotating text fragment 14821/42019
## 2023-09-16 21:35:45 Annotating text fragment 14831/42019
## 2023-09-16 21:35:45 Annotating text fragment 14841/42019
## 2023-09-16 21:35:45 Annotating text fragment 14851/42019
## 2023-09-16 21:35:45 Annotating text fragment 14861/42019
## 2023-09-16 21:35:45 Annotating text fragment 14871/42019
## 2023-09-16 21:35:45 Annotating text fragment 14881/42019
## 2023-09-16 21:35:45 Annotating text fragment 14891/42019
## 2023-09-16 21:35:45 Annotating text fragment 14901/42019
## 2023-09-16 21:35:45 Annotating text fragment 14911/42019
## 2023-09-16 21:35:45 Annotating text fragment 14921/42019
## 2023-09-16 21:35:45 Annotating text fragment 14931/42019
## 2023-09-16 21:35:46 Annotating text fragment 14941/42019
## 2023-09-16 21:35:46 Annotating text fragment 14951/42019
## 2023-09-16 21:35:46 Annotating text fragment 14961/42019
## 2023-09-16 21:35:46 Annotating text fragment 14971/42019
## 2023-09-16 21:35:46 Annotating text fragment 14981/42019
## 2023-09-16 21:35:46 Annotating text fragment 14991/42019
## 2023-09-16 21:35:46 Annotating text fragment 15001/42019
## 2023-09-16 21:35:46 Annotating text fragment 15011/42019
## 2023-09-16 21:35:46 Annotating text fragment 15021/42019
## 2023-09-16 21:35:46 Annotating text fragment 15031/42019
## 2023-09-16 21:35:46 Annotating text fragment 15041/42019
## 2023-09-16 21:35:46 Annotating text fragment 15051/42019
## 2023-09-16 21:35:46 Annotating text fragment 15061/42019
## 2023-09-16 21:35:46 Annotating text fragment 15071/42019
## 2023-09-16 21:35:47 Annotating text fragment 15081/42019
## 2023-09-16 21:35:47 Annotating text fragment 15091/42019
## 2023-09-16 21:35:47 Annotating text fragment 15101/42019
## 2023-09-16 21:35:47 Annotating text fragment 15111/42019
## 2023-09-16 21:35:47 Annotating text fragment 15121/42019
## 2023-09-16 21:35:47 Annotating text fragment 15131/42019
## 2023-09-16 21:35:47 Annotating text fragment 15141/42019
## 2023-09-16 21:35:47 Annotating text fragment 15151/42019
## 2023-09-16 21:35:47 Annotating text fragment 15161/42019
## 2023-09-16 21:35:47 Annotating text fragment 15171/42019
## 2023-09-16 21:35:47 Annotating text fragment 15181/42019
## 2023-09-16 21:35:47 Annotating text fragment 15191/42019
## 2023-09-16 21:35:47 Annotating text fragment 15201/42019
## 2023-09-16 21:35:47 Annotating text fragment 15211/42019
## 2023-09-16 21:35:48 Annotating text fragment 15221/42019
## 2023-09-16 21:35:48 Annotating text fragment 15231/42019
## 2023-09-16 21:35:48 Annotating text fragment 15241/42019
## 2023-09-16 21:35:48 Annotating text fragment 15251/42019
## 2023-09-16 21:35:48 Annotating text fragment 15261/42019
```

```
## 2023-09-16 21:35:48 Annotating text fragment 15271/42019
## 2023-09-16 21:35:48 Annotating text fragment 15281/42019
## 2023-09-16 21:35:48 Annotating text fragment 15291/42019
## 2023-09-16 21:35:48 Annotating text fragment 15301/42019
## 2023-09-16 21:35:48 Annotating text fragment 15311/42019
## 2023-09-16 21:35:48 Annotating text fragment 15321/42019
## 2023-09-16 21:35:48 Annotating text fragment 15331/42019
## 2023-09-16 21:35:48 Annotating text fragment 15341/42019
## 2023-09-16 21:35:48 Annotating text fragment 15351/42019
## 2023-09-16 21:35:48 Annotating text fragment 15361/42019
## 2023-09-16 21:35:48 Annotating text fragment 15371/42019
## 2023-09-16 21:35:48 Annotating text fragment 15381/42019
## 2023-09-16 21:35:48 Annotating text fragment 15391/42019
## 2023-09-16 21:35:49 Annotating text fragment 15401/42019
## 2023-09-16 21:35:49 Annotating text fragment 15411/42019
## 2023-09-16 21:35:49 Annotating text fragment 15421/42019
## 2023-09-16 21:35:49 Annotating text fragment 15431/42019
## 2023-09-16 21:35:49 Annotating text fragment 15441/42019
## 2023-09-16 21:35:49 Annotating text fragment 15451/42019
## 2023-09-16 21:35:49 Annotating text fragment 15461/42019
## 2023-09-16 21:35:49 Annotating text fragment 15471/42019
## 2023-09-16 21:35:49 Annotating text fragment 15481/42019
## 2023-09-16 21:35:49 Annotating text fragment 15491/42019
## 2023-09-16 21:35:49 Annotating text fragment 15501/42019
## 2023-09-16 21:35:49 Annotating text fragment 15511/42019
## 2023-09-16 21:35:49 Annotating text fragment 15521/42019
## 2023-09-16 21:35:50 Annotating text fragment 15531/42019
## 2023-09-16 21:35:50 Annotating text fragment 15541/42019
## 2023-09-16 21:35:50 Annotating text fragment 15551/42019
## 2023-09-16 21:35:50 Annotating text fragment 15561/42019
## 2023-09-16 21:35:50 Annotating text fragment 15571/42019
## 2023-09-16 21:35:50 Annotating text fragment 15581/42019
## 2023-09-16 21:35:50 Annotating text fragment 15591/42019
## 2023-09-16 21:35:50 Annotating text fragment 15601/42019
## 2023-09-16 21:35:50 Annotating text fragment 15611/42019
## 2023-09-16 21:35:50 Annotating text fragment 15621/42019
## 2023-09-16 21:35:50 Annotating text fragment 15631/42019
## 2023-09-16 21:35:50 Annotating text fragment 15641/42019
## 2023-09-16 21:35:50 Annotating text fragment 15651/42019
## 2023-09-16 21:35:50 Annotating text fragment 15661/42019
## 2023-09-16 21:35:51 Annotating text fragment 15671/42019
## 2023-09-16 21:35:51 Annotating text fragment 15681/42019
## 2023-09-16 21:35:51 Annotating text fragment 15691/42019
## 2023-09-16 21:35:51 Annotating text fragment 15701/42019
## 2023-09-16 21:35:51 Annotating text fragment 15711/42019
## 2023-09-16 21:35:51 Annotating text fragment 15721/42019
## 2023-09-16 21:35:51 Annotating text fragment 15731/42019
## 2023-09-16 21:35:51 Annotating text fragment 15741/42019
## 2023-09-16 21:35:51 Annotating text fragment 15751/42019
## 2023-09-16 21:35:51 Annotating text fragment 15761/42019
## 2023-09-16 21:35:51 Annotating text fragment 15771/42019
## 2023-09-16 21:35:51 Annotating text fragment 15781/42019
## 2023-09-16 21:35:52 Annotating text fragment 15791/42019
## 2023-09-16 21:35:52 Annotating text fragment 15801/42019
```

```
## 2023-09-16 21:35:52 Annotating text fragment 15811/42019
## 2023-09-16 21:35:52 Annotating text fragment 15821/42019
## 2023-09-16 21:35:52 Annotating text fragment 15831/42019
## 2023-09-16 21:35:52 Annotating text fragment 15841/42019
## 2023-09-16 21:35:52 Annotating text fragment 15851/42019
## 2023-09-16 21:35:53 Annotating text fragment 15861/42019
## 2023-09-16 21:35:53 Annotating text fragment 15871/42019
## 2023-09-16 21:35:53 Annotating text fragment 15881/42019
## 2023-09-16 21:35:53 Annotating text fragment 15891/42019
## 2023-09-16 21:35:53 Annotating text fragment 15901/42019
## 2023-09-16 21:35:53 Annotating text fragment 15911/42019
## 2023-09-16 21:35:53 Annotating text fragment 15921/42019
## 2023-09-16 21:35:53 Annotating text fragment 15931/42019
## 2023-09-16 21:35:54 Annotating text fragment 15941/42019
## 2023-09-16 21:35:54 Annotating text fragment 15951/42019
## 2023-09-16 21:35:54 Annotating text fragment 15961/42019
## 2023-09-16 21:35:54 Annotating text fragment 15971/42019
## 2023-09-16 21:35:54 Annotating text fragment 15981/42019
## 2023-09-16 21:35:54 Annotating text fragment 15991/42019
## 2023-09-16 21:35:54 Annotating text fragment 16001/42019
## 2023-09-16 21:35:54 Annotating text fragment 16011/42019
## 2023-09-16 21:35:54 Annotating text fragment 16021/42019
## 2023-09-16 21:35:54 Annotating text fragment 16031/42019
## 2023-09-16 21:35:54 Annotating text fragment 16041/42019
## 2023-09-16 21:35:54 Annotating text fragment 16051/42019
## 2023-09-16 21:35:55 Annotating text fragment 16061/42019
## 2023-09-16 21:35:55 Annotating text fragment 16071/42019
## 2023-09-16 21:35:55 Annotating text fragment 16081/42019
## 2023-09-16 21:35:55 Annotating text fragment 16091/42019
## 2023-09-16 21:35:55 Annotating text fragment 16101/42019
## 2023-09-16 21:35:55 Annotating text fragment 16111/42019
## 2023-09-16 21:35:55 Annotating text fragment 16121/42019
## 2023-09-16 21:35:55 Annotating text fragment 16131/42019
## 2023-09-16 21:35:55 Annotating text fragment 16141/42019
## 2023-09-16 21:35:55 Annotating text fragment 16151/42019
## 2023-09-16 21:35:55 Annotating text fragment 16161/42019
## 2023-09-16 21:35:55 Annotating text fragment 16171/42019
## 2023-09-16 21:35:55 Annotating text fragment 16181/42019
## 2023-09-16 21:35:55 Annotating text fragment 16191/42019
## 2023-09-16 21:35:55 Annotating text fragment 16201/42019
## 2023-09-16 21:35:55 Annotating text fragment 16211/42019
## 2023-09-16 21:35:56 Annotating text fragment 16221/42019
## 2023-09-16 21:35:56 Annotating text fragment 16231/42019
## 2023-09-16 21:35:56 Annotating text fragment 16241/42019
## 2023-09-16 21:35:56 Annotating text fragment 16251/42019
## 2023-09-16 21:35:56 Annotating text fragment 16261/42019
## 2023-09-16 21:35:56 Annotating text fragment 16271/42019
## 2023-09-16 21:35:56 Annotating text fragment 16281/42019
## 2023-09-16 21:35:56 Annotating text fragment 16291/42019
## 2023-09-16 21:35:56 Annotating text fragment 16301/42019
## 2023-09-16 21:35:57 Annotating text fragment 16311/42019
## 2023-09-16 21:35:57 Annotating text fragment 16321/42019
## 2023-09-16 21:35:57 Annotating text fragment 16331/42019
## 2023-09-16 21:35:57 Annotating text fragment 16341/42019
```

```
## 2023-09-16 21:35:57 Annotating text fragment 16351/42019
## 2023-09-16 21:35:57 Annotating text fragment 16361/42019
## 2023-09-16 21:35:57 Annotating text fragment 16371/42019
## 2023-09-16 21:35:57 Annotating text fragment 16381/42019
## 2023-09-16 21:35:57 Annotating text fragment 16391/42019
## 2023-09-16 21:35:57 Annotating text fragment 16401/42019
## 2023-09-16 21:35:57 Annotating text fragment 16411/42019
## 2023-09-16 21:35:57 Annotating text fragment 16421/42019
## 2023-09-16 21:35:57 Annotating text fragment 16431/42019
## 2023-09-16 21:35:57 Annotating text fragment 16441/42019
## 2023-09-16 21:35:57 Annotating text fragment 16451/42019
## 2023-09-16 21:35:57 Annotating text fragment 16461/42019
## 2023-09-16 21:35:57 Annotating text fragment 16471/42019
## 2023-09-16 21:35:57 Annotating text fragment 16481/42019
## 2023-09-16 21:35:58 Annotating text fragment 16491/42019
## 2023-09-16 21:35:58 Annotating text fragment 16501/42019
## 2023-09-16 21:35:58 Annotating text fragment 16511/42019
## 2023-09-16 21:35:58 Annotating text fragment 16521/42019
## 2023-09-16 21:35:58 Annotating text fragment 16531/42019
## 2023-09-16 21:35:58 Annotating text fragment 16541/42019
## 2023-09-16 21:35:58 Annotating text fragment 16551/42019
## 2023-09-16 21:35:58 Annotating text fragment 16561/42019
## 2023-09-16 21:35:58 Annotating text fragment 16571/42019
## 2023-09-16 21:35:59 Annotating text fragment 16581/42019
## 2023-09-16 21:35:59 Annotating text fragment 16591/42019
## 2023-09-16 21:35:59 Annotating text fragment 16601/42019
## 2023-09-16 21:35:59 Annotating text fragment 16611/42019
## 2023-09-16 21:35:59 Annotating text fragment 16621/42019
## 2023-09-16 21:35:59 Annotating text fragment 16631/42019
## 2023-09-16 21:35:59 Annotating text fragment 16641/42019
## 2023-09-16 21:35:59 Annotating text fragment 16651/42019
## 2023-09-16 21:35:59 Annotating text fragment 16661/42019
## 2023-09-16 21:35:59 Annotating text fragment 16671/42019
## 2023-09-16 21:35:59 Annotating text fragment 16681/42019
## 2023-09-16 21:35:59 Annotating text fragment 16691/42019
## 2023-09-16 21:35:59 Annotating text fragment 16701/42019
## 2023-09-16 21:35:59 Annotating text fragment 16711/42019
## 2023-09-16 21:36:00 Annotating text fragment 16721/42019
## 2023-09-16 21:36:00 Annotating text fragment 16731/42019
## 2023-09-16 21:36:00 Annotating text fragment 16741/42019
## 2023-09-16 21:36:00 Annotating text fragment 16751/42019
## 2023-09-16 21:36:00 Annotating text fragment 16761/42019
## 2023-09-16 21:36:00 Annotating text fragment 16771/42019
## 2023-09-16 21:36:00 Annotating text fragment 16781/42019
## 2023-09-16 21:36:00 Annotating text fragment 16791/42019
## 2023-09-16 21:36:00 Annotating text fragment 16801/42019
## 2023-09-16 21:36:00 Annotating text fragment 16811/42019
## 2023-09-16 21:36:00 Annotating text fragment 16821/42019
## 2023-09-16 21:36:00 Annotating text fragment 16831/42019
## 2023-09-16 21:36:00 Annotating text fragment 16841/42019
## 2023-09-16 21:36:00 Annotating text fragment 16851/42019
## 2023-09-16 21:36:00 Annotating text fragment 16861/42019
## 2023-09-16 21:36:00 Annotating text fragment 16871/42019
## 2023-09-16 21:36:00 Annotating text fragment 16881/42019
```

```
## 2023-09-16 21:36:00 Annotating text fragment 16891/42019
## 2023-09-16 21:36:01 Annotating text fragment 16901/42019
## 2023-09-16 21:36:01 Annotating text fragment 16911/42019
## 2023-09-16 21:36:01 Annotating text fragment 16921/42019
## 2023-09-16 21:36:01 Annotating text fragment 16931/42019
## 2023-09-16 21:36:01 Annotating text fragment 16941/42019
## 2023-09-16 21:36:01 Annotating text fragment 16951/42019
## 2023-09-16 21:36:01 Annotating text fragment 16961/42019
## 2023-09-16 21:36:01 Annotating text fragment 16971/42019
## 2023-09-16 21:36:01 Annotating text fragment 16981/42019
## 2023-09-16 21:36:01 Annotating text fragment 16991/42019
## 2023-09-16 21:36:02 Annotating text fragment 17001/42019
## 2023-09-16 21:36:02 Annotating text fragment 17011/42019
## 2023-09-16 21:36:02 Annotating text fragment 17021/42019
## 2023-09-16 21:36:02 Annotating text fragment 17031/42019
## 2023-09-16 21:36:02 Annotating text fragment 17041/42019
## 2023-09-16 21:36:02 Annotating text fragment 17051/42019
## 2023-09-16 21:36:02 Annotating text fragment 17061/42019
## 2023-09-16 21:36:02 Annotating text fragment 17071/42019
## 2023-09-16 21:36:02 Annotating text fragment 17081/42019
## 2023-09-16 21:36:02 Annotating text fragment 17091/42019
## 2023-09-16 21:36:02 Annotating text fragment 17101/42019
## 2023-09-16 21:36:03 Annotating text fragment 17111/42019
## 2023-09-16 21:36:03 Annotating text fragment 17121/42019
## 2023-09-16 21:36:03 Annotating text fragment 17131/42019
## 2023-09-16 21:36:03 Annotating text fragment 17141/42019
## 2023-09-16 21:36:03 Annotating text fragment 17151/42019
## 2023-09-16 21:36:03 Annotating text fragment 17161/42019
## 2023-09-16 21:36:03 Annotating text fragment 17171/42019
## 2023-09-16 21:36:04 Annotating text fragment 17181/42019
## 2023-09-16 21:36:04 Annotating text fragment 17191/42019
## 2023-09-16 21:36:05 Annotating text fragment 17201/42019
## 2023-09-16 21:36:05 Annotating text fragment 17211/42019
## 2023-09-16 21:36:05 Annotating text fragment 17221/42019
## 2023-09-16 21:36:05 Annotating text fragment 17231/42019
## 2023-09-16 21:36:05 Annotating text fragment 17241/42019
## 2023-09-16 21:36:05 Annotating text fragment 17251/42019
## 2023-09-16 21:36:05 Annotating text fragment 17261/42019
## 2023-09-16 21:36:05 Annotating text fragment 17271/42019
## 2023-09-16 21:36:05 Annotating text fragment 17281/42019
## 2023-09-16 21:36:05 Annotating text fragment 17291/42019
## 2023-09-16 21:36:05 Annotating text fragment 17301/42019
## 2023-09-16 21:36:06 Annotating text fragment 17311/42019
## 2023-09-16 21:36:06 Annotating text fragment 17321/42019
## 2023-09-16 21:36:06 Annotating text fragment 17331/42019
## 2023-09-16 21:36:06 Annotating text fragment 17341/42019
## 2023-09-16 21:36:06 Annotating text fragment 17351/42019
## 2023-09-16 21:36:06 Annotating text fragment 17361/42019
## 2023-09-16 21:36:06 Annotating text fragment 17371/42019
## 2023-09-16 21:36:06 Annotating text fragment 17381/42019
## 2023-09-16 21:36:06 Annotating text fragment 17391/42019
## 2023-09-16 21:36:07 Annotating text fragment 17401/42019
## 2023-09-16 21:36:07 Annotating text fragment 17411/42019
## 2023-09-16 21:36:07 Annotating text fragment 17421/42019
```

```
## 2023-09-16 21:36:07 Annotating text fragment 17431/42019
## 2023-09-16 21:36:07 Annotating text fragment 17441/42019
## 2023-09-16 21:36:07 Annotating text fragment 17451/42019
## 2023-09-16 21:36:07 Annotating text fragment 17461/42019
## 2023-09-16 21:36:07 Annotating text fragment 17471/42019
## 2023-09-16 21:36:08 Annotating text fragment 17481/42019
## 2023-09-16 21:36:08 Annotating text fragment 17491/42019
## 2023-09-16 21:36:08 Annotating text fragment 17501/42019
## 2023-09-16 21:36:08 Annotating text fragment 17511/42019
## 2023-09-16 21:36:08 Annotating text fragment 17521/42019
## 2023-09-16 21:36:08 Annotating text fragment 17531/42019
## 2023-09-16 21:36:08 Annotating text fragment 17541/42019
## 2023-09-16 21:36:08 Annotating text fragment 17551/42019
## 2023-09-16 21:36:08 Annotating text fragment 17561/42019
## 2023-09-16 21:36:08 Annotating text fragment 17571/42019
## 2023-09-16 21:36:08 Annotating text fragment 17581/42019
## 2023-09-16 21:36:08 Annotating text fragment 17591/42019
## 2023-09-16 21:36:08 Annotating text fragment 17601/42019
## 2023-09-16 21:36:09 Annotating text fragment 17611/42019
## 2023-09-16 21:36:09 Annotating text fragment 17621/42019
## 2023-09-16 21:36:09 Annotating text fragment 17631/42019
## 2023-09-16 21:36:09 Annotating text fragment 17641/42019
## 2023-09-16 21:36:09 Annotating text fragment 17651/42019
## 2023-09-16 21:36:09 Annotating text fragment 17661/42019
## 2023-09-16 21:36:09 Annotating text fragment 17671/42019
## 2023-09-16 21:36:09 Annotating text fragment 17681/42019
## 2023-09-16 21:36:09 Annotating text fragment 17691/42019
## 2023-09-16 21:36:09 Annotating text fragment 17701/42019
## 2023-09-16 21:36:10 Annotating text fragment 17711/42019
## 2023-09-16 21:36:10 Annotating text fragment 17721/42019
## 2023-09-16 21:36:10 Annotating text fragment 17731/42019
## 2023-09-16 21:36:10 Annotating text fragment 17741/42019
## 2023-09-16 21:36:10 Annotating text fragment 17751/42019
## 2023-09-16 21:36:10 Annotating text fragment 17761/42019
## 2023-09-16 21:36:10 Annotating text fragment 17771/42019
## 2023-09-16 21:36:10 Annotating text fragment 17781/42019
## 2023-09-16 21:36:10 Annotating text fragment 17791/42019
## 2023-09-16 21:36:10 Annotating text fragment 17801/42019
## 2023-09-16 21:36:11 Annotating text fragment 17811/42019
## 2023-09-16 21:36:11 Annotating text fragment 17821/42019
## 2023-09-16 21:36:11 Annotating text fragment 17831/42019
## 2023-09-16 21:36:11 Annotating text fragment 17841/42019
## 2023-09-16 21:36:11 Annotating text fragment 17851/42019
## 2023-09-16 21:36:11 Annotating text fragment 17861/42019
## 2023-09-16 21:36:11 Annotating text fragment 17871/42019
## 2023-09-16 21:36:11 Annotating text fragment 17881/42019
## 2023-09-16 21:36:11 Annotating text fragment 17891/42019
## 2023-09-16 21:36:11 Annotating text fragment 17901/42019
## 2023-09-16 21:36:11 Annotating text fragment 17911/42019
## 2023-09-16 21:36:11 Annotating text fragment 17921/42019
## 2023-09-16 21:36:11 Annotating text fragment 17931/42019
## 2023-09-16 21:36:11 Annotating text fragment 17941/42019
## 2023-09-16 21:36:11 Annotating text fragment 17951/42019
## 2023-09-16 21:36:12 Annotating text fragment 17961/42019
```

```
## 2023-09-16 21:36:12 Annotating text fragment 17971/42019
## 2023-09-16 21:36:12 Annotating text fragment 17981/42019
## 2023-09-16 21:36:12 Annotating text fragment 17991/42019
## 2023-09-16 21:36:12 Annotating text fragment 18001/42019
## 2023-09-16 21:36:12 Annotating text fragment 18011/42019
## 2023-09-16 21:36:12 Annotating text fragment 18021/42019
## 2023-09-16 21:36:13 Annotating text fragment 18031/42019
## 2023-09-16 21:36:13 Annotating text fragment 18041/42019
## 2023-09-16 21:36:13 Annotating text fragment 18051/42019
## 2023-09-16 21:36:13 Annotating text fragment 18061/42019
## 2023-09-16 21:36:13 Annotating text fragment 18071/42019
## 2023-09-16 21:36:13 Annotating text fragment 18081/42019
## 2023-09-16 21:36:13 Annotating text fragment 18091/42019
## 2023-09-16 21:36:13 Annotating text fragment 18101/42019
## 2023-09-16 21:36:13 Annotating text fragment 18111/42019
## 2023-09-16 21:36:13 Annotating text fragment 18121/42019
## 2023-09-16 21:36:13 Annotating text fragment 18131/42019
## 2023-09-16 21:36:13 Annotating text fragment 18141/42019
## 2023-09-16 21:36:13 Annotating text fragment 18151/42019
## 2023-09-16 21:36:13 Annotating text fragment 18161/42019
## 2023-09-16 21:36:13 Annotating text fragment 18171/42019
## 2023-09-16 21:36:14 Annotating text fragment 18181/42019
## 2023-09-16 21:36:14 Annotating text fragment 18191/42019
## 2023-09-16 21:36:14 Annotating text fragment 18201/42019
## 2023-09-16 21:36:14 Annotating text fragment 18211/42019
## 2023-09-16 21:36:14 Annotating text fragment 18221/42019
## 2023-09-16 21:36:14 Annotating text fragment 18231/42019
## 2023-09-16 21:36:14 Annotating text fragment 18241/42019
## 2023-09-16 21:36:14 Annotating text fragment 18251/42019
## 2023-09-16 21:36:14 Annotating text fragment 18261/42019
## 2023-09-16 21:36:14 Annotating text fragment 18271/42019
## 2023-09-16 21:36:14 Annotating text fragment 18281/42019
## 2023-09-16 21:36:15 Annotating text fragment 18291/42019
## 2023-09-16 21:36:15 Annotating text fragment 18301/42019
## 2023-09-16 21:36:15 Annotating text fragment 18311/42019
## 2023-09-16 21:36:15 Annotating text fragment 18321/42019
## 2023-09-16 21:36:15 Annotating text fragment 18331/42019
## 2023-09-16 21:36:15 Annotating text fragment 18341/42019
## 2023-09-16 21:36:15 Annotating text fragment 18351/42019
## 2023-09-16 21:36:15 Annotating text fragment 18361/42019
## 2023-09-16 21:36:15 Annotating text fragment 18371/42019
## 2023-09-16 21:36:15 Annotating text fragment 18381/42019
## 2023-09-16 21:36:15 Annotating text fragment 18391/42019
## 2023-09-16 21:36:15 Annotating text fragment 18401/42019
## 2023-09-16 21:36:15 Annotating text fragment 18411/42019
## 2023-09-16 21:36:15 Annotating text fragment 18421/42019
## 2023-09-16 21:36:15 Annotating text fragment 18431/42019
## 2023-09-16 21:36:15 Annotating text fragment 18441/42019
## 2023-09-16 21:36:15 Annotating text fragment 18451/42019
## 2023-09-16 21:36:16 Annotating text fragment 18461/42019
## 2023-09-16 21:36:16 Annotating text fragment 18471/42019
## 2023-09-16 21:36:16 Annotating text fragment 18481/42019
## 2023-09-16 21:36:16 Annotating text fragment 18491/42019
## 2023-09-16 21:36:16 Annotating text fragment 18501/42019
```

```
## 2023-09-16 21:36:16 Annotating text fragment 18511/42019
## 2023-09-16 21:36:16 Annotating text fragment 18521/42019
## 2023-09-16 21:36:16 Annotating text fragment 18531/42019
## 2023-09-16 21:36:16 Annotating text fragment 18541/42019
## 2023-09-16 21:36:16 Annotating text fragment 18551/42019
## 2023-09-16 21:36:16 Annotating text fragment 18561/42019
## 2023-09-16 21:36:16 Annotating text fragment 18571/42019
## 2023-09-16 21:36:16 Annotating text fragment 18581/42019
## 2023-09-16 21:36:16 Annotating text fragment 18591/42019
## 2023-09-16 21:36:17 Annotating text fragment 18601/42019
## 2023-09-16 21:36:17 Annotating text fragment 18611/42019
## 2023-09-16 21:36:17 Annotating text fragment 18621/42019
## 2023-09-16 21:36:17 Annotating text fragment 18631/42019
## 2023-09-16 21:36:17 Annotating text fragment 18641/42019
## 2023-09-16 21:36:17 Annotating text fragment 18651/42019
## 2023-09-16 21:36:17 Annotating text fragment 18661/42019
## 2023-09-16 21:36:17 Annotating text fragment 18671/42019
## 2023-09-16 21:36:17 Annotating text fragment 18681/42019
## 2023-09-16 21:36:17 Annotating text fragment 18691/42019
## 2023-09-16 21:36:17 Annotating text fragment 18701/42019
## 2023-09-16 21:36:17 Annotating text fragment 18711/42019
## 2023-09-16 21:36:17 Annotating text fragment 18721/42019
## 2023-09-16 21:36:18 Annotating text fragment 18731/42019
## 2023-09-16 21:36:18 Annotating text fragment 18741/42019
## 2023-09-16 21:36:18 Annotating text fragment 18751/42019
## 2023-09-16 21:36:18 Annotating text fragment 18761/42019
## 2023-09-16 21:36:18 Annotating text fragment 18771/42019
## 2023-09-16 21:36:18 Annotating text fragment 18781/42019
## 2023-09-16 21:36:18 Annotating text fragment 18791/42019
## 2023-09-16 21:36:18 Annotating text fragment 18801/42019
## 2023-09-16 21:36:18 Annotating text fragment 18811/42019
## 2023-09-16 21:36:18 Annotating text fragment 18821/42019
## 2023-09-16 21:36:18 Annotating text fragment 18831/42019
## 2023-09-16 21:36:18 Annotating text fragment 18841/42019
## 2023-09-16 21:36:18 Annotating text fragment 18851/42019
## 2023-09-16 21:36:19 Annotating text fragment 18861/42019
## 2023-09-16 21:36:19 Annotating text fragment 18871/42019
## 2023-09-16 21:36:19 Annotating text fragment 18881/42019
## 2023-09-16 21:36:19 Annotating text fragment 18891/42019
## 2023-09-16 21:36:19 Annotating text fragment 18901/42019
## 2023-09-16 21:36:19 Annotating text fragment 18911/42019
## 2023-09-16 21:36:19 Annotating text fragment 18921/42019
## 2023-09-16 21:36:19 Annotating text fragment 18931/42019
## 2023-09-16 21:36:19 Annotating text fragment 18941/42019
## 2023-09-16 21:36:19 Annotating text fragment 18951/42019
## 2023-09-16 21:36:20 Annotating text fragment 18961/42019
## 2023-09-16 21:36:20 Annotating text fragment 18971/42019
## 2023-09-16 21:36:20 Annotating text fragment 18981/42019
## 2023-09-16 21:36:20 Annotating text fragment 18991/42019
## 2023-09-16 21:36:20 Annotating text fragment 19001/42019
## 2023-09-16 21:36:20 Annotating text fragment 19011/42019
## 2023-09-16 21:36:20 Annotating text fragment 19021/42019
## 2023-09-16 21:36:20 Annotating text fragment 19031/42019
## 2023-09-16 21:36:20 Annotating text fragment 19041/42019
```

```
## 2023-09-16 21:36:20 Annotating text fragment 19051/42019
## 2023-09-16 21:36:20 Annotating text fragment 19061/42019
## 2023-09-16 21:36:20 Annotating text fragment 19071/42019
## 2023-09-16 21:36:20 Annotating text fragment 19081/42019
## 2023-09-16 21:36:20 Annotating text fragment 19091/42019
## 2023-09-16 21:36:20 Annotating text fragment 19101/42019
## 2023-09-16 21:36:20 Annotating text fragment 19111/42019
## 2023-09-16 21:36:21 Annotating text fragment 19121/42019
## 2023-09-16 21:36:21 Annotating text fragment 19131/42019
## 2023-09-16 21:36:21 Annotating text fragment 19141/42019
## 2023-09-16 21:36:21 Annotating text fragment 19151/42019
## 2023-09-16 21:36:21 Annotating text fragment 19161/42019
## 2023-09-16 21:36:21 Annotating text fragment 19171/42019
## 2023-09-16 21:36:21 Annotating text fragment 19181/42019
## 2023-09-16 21:36:21 Annotating text fragment 19191/42019
## 2023-09-16 21:36:21 Annotating text fragment 19201/42019
## 2023-09-16 21:36:21 Annotating text fragment 19211/42019
## 2023-09-16 21:36:22 Annotating text fragment 19221/42019
## 2023-09-16 21:36:22 Annotating text fragment 19231/42019
## 2023-09-16 21:36:22 Annotating text fragment 19241/42019
## 2023-09-16 21:36:22 Annotating text fragment 19251/42019
## 2023-09-16 21:36:22 Annotating text fragment 19261/42019
## 2023-09-16 21:36:22 Annotating text fragment 19271/42019
## 2023-09-16 21:36:22 Annotating text fragment 19281/42019
## 2023-09-16 21:36:22 Annotating text fragment 19291/42019
## 2023-09-16 21:36:22 Annotating text fragment 19301/42019
## 2023-09-16 21:36:22 Annotating text fragment 19311/42019
## 2023-09-16 21:36:22 Annotating text fragment 19321/42019
## 2023-09-16 21:36:22 Annotating text fragment 19331/42019
## 2023-09-16 21:36:22 Annotating text fragment 19341/42019
## 2023-09-16 21:36:23 Annotating text fragment 19351/42019
## 2023-09-16 21:36:23 Annotating text fragment 19361/42019
## 2023-09-16 21:36:23 Annotating text fragment 19371/42019
## 2023-09-16 21:36:24 Annotating text fragment 19381/42019
## 2023-09-16 21:36:24 Annotating text fragment 19391/42019
## 2023-09-16 21:36:24 Annotating text fragment 19401/42019
## 2023-09-16 21:36:24 Annotating text fragment 19411/42019
## 2023-09-16 21:36:24 Annotating text fragment 19421/42019
## 2023-09-16 21:36:24 Annotating text fragment 19431/42019
## 2023-09-16 21:36:24 Annotating text fragment 19441/42019
## 2023-09-16 21:36:24 Annotating text fragment 19451/42019
## 2023-09-16 21:36:24 Annotating text fragment 19461/42019
## 2023-09-16 21:36:24 Annotating text fragment 19471/42019
## 2023-09-16 21:36:24 Annotating text fragment 19481/42019
## 2023-09-16 21:36:24 Annotating text fragment 19491/42019
## 2023-09-16 21:36:25 Annotating text fragment 19501/42019
## 2023-09-16 21:36:25 Annotating text fragment 19511/42019
## 2023-09-16 21:36:25 Annotating text fragment 19521/42019
## 2023-09-16 21:36:25 Annotating text fragment 19531/42019
## 2023-09-16 21:36:25 Annotating text fragment 19541/42019
## 2023-09-16 21:36:25 Annotating text fragment 19551/42019
## 2023-09-16 21:36:25 Annotating text fragment 19561/42019
## 2023-09-16 21:36:25 Annotating text fragment 19571/42019
## 2023-09-16 21:36:25 Annotating text fragment 19581/42019
```

```
## 2023-09-16 21:36:25 Annotating text fragment 19591/42019
## 2023-09-16 21:36:26 Annotating text fragment 19601/42019
## 2023-09-16 21:36:26 Annotating text fragment 19611/42019
## 2023-09-16 21:36:26 Annotating text fragment 19621/42019
## 2023-09-16 21:36:26 Annotating text fragment 19631/42019
## 2023-09-16 21:36:26 Annotating text fragment 19641/42019
## 2023-09-16 21:36:26 Annotating text fragment 19651/42019
## 2023-09-16 21:36:26 Annotating text fragment 19661/42019
## 2023-09-16 21:36:26 Annotating text fragment 19671/42019
## 2023-09-16 21:36:27 Annotating text fragment 19681/42019
## 2023-09-16 21:36:27 Annotating text fragment 19691/42019
## 2023-09-16 21:36:27 Annotating text fragment 19701/42019
## 2023-09-16 21:36:27 Annotating text fragment 19711/42019
## 2023-09-16 21:36:27 Annotating text fragment 19721/42019
## 2023-09-16 21:36:27 Annotating text fragment 19731/42019
## 2023-09-16 21:36:27 Annotating text fragment 19741/42019
## 2023-09-16 21:36:27 Annotating text fragment 19751/42019
## 2023-09-16 21:36:27 Annotating text fragment 19761/42019
## 2023-09-16 21:36:27 Annotating text fragment 19771/42019
## 2023-09-16 21:36:28 Annotating text fragment 19781/42019
## 2023-09-16 21:36:28 Annotating text fragment 19791/42019
## 2023-09-16 21:36:28 Annotating text fragment 19801/42019
## 2023-09-16 21:36:28 Annotating text fragment 19811/42019
## 2023-09-16 21:36:28 Annotating text fragment 19821/42019
## 2023-09-16 21:36:28 Annotating text fragment 19831/42019
## 2023-09-16 21:36:28 Annotating text fragment 19841/42019
## 2023-09-16 21:36:28 Annotating text fragment 19851/42019
## 2023-09-16 21:36:28 Annotating text fragment 19861/42019
## 2023-09-16 21:36:29 Annotating text fragment 19871/42019
## 2023-09-16 21:36:29 Annotating text fragment 19881/42019
## 2023-09-16 21:36:29 Annotating text fragment 19891/42019
## 2023-09-16 21:36:29 Annotating text fragment 19901/42019
## 2023-09-16 21:36:29 Annotating text fragment 19911/42019
## 2023-09-16 21:36:29 Annotating text fragment 19921/42019
## 2023-09-16 21:36:29 Annotating text fragment 19931/42019
## 2023-09-16 21:36:29 Annotating text fragment 19941/42019
## 2023-09-16 21:36:29 Annotating text fragment 19951/42019
## 2023-09-16 21:36:29 Annotating text fragment 19961/42019
## 2023-09-16 21:36:29 Annotating text fragment 19971/42019
## 2023-09-16 21:36:29 Annotating text fragment 19981/42019
## 2023-09-16 21:36:30 Annotating text fragment 19991/42019
## 2023-09-16 21:36:30 Annotating text fragment 20001/42019
## 2023-09-16 21:36:30 Annotating text fragment 20011/42019
## 2023-09-16 21:36:30 Annotating text fragment 20021/42019
## 2023-09-16 21:36:30 Annotating text fragment 20031/42019
## 2023-09-16 21:36:30 Annotating text fragment 20041/42019
## 2023-09-16 21:36:30 Annotating text fragment 20051/42019
## 2023-09-16 21:36:30 Annotating text fragment 20061/42019
## 2023-09-16 21:36:30 Annotating text fragment 20071/42019
## 2023-09-16 21:36:31 Annotating text fragment 20081/42019
## 2023-09-16 21:36:31 Annotating text fragment 20091/42019
## 2023-09-16 21:36:31 Annotating text fragment 20101/42019
## 2023-09-16 21:36:31 Annotating text fragment 20111/42019
## 2023-09-16 21:36:31 Annotating text fragment 20121/42019
```

```
## 2023-09-16 21:36:31 Annotating text fragment 20131/42019
## 2023-09-16 21:36:31 Annotating text fragment 20141/42019
## 2023-09-16 21:36:31 Annotating text fragment 20151/42019
## 2023-09-16 21:36:31 Annotating text fragment 20161/42019
## 2023-09-16 21:36:31 Annotating text fragment 20171/42019
## 2023-09-16 21:36:31 Annotating text fragment 20181/42019
## 2023-09-16 21:36:31 Annotating text fragment 20191/42019
## 2023-09-16 21:36:32 Annotating text fragment 20201/42019
## 2023-09-16 21:36:32 Annotating text fragment 20211/42019
## 2023-09-16 21:36:32 Annotating text fragment 20221/42019
## 2023-09-16 21:36:32 Annotating text fragment 20231/42019
## 2023-09-16 21:36:32 Annotating text fragment 20241/42019
## 2023-09-16 21:36:32 Annotating text fragment 20251/42019
## 2023-09-16 21:36:32 Annotating text fragment 20261/42019
## 2023-09-16 21:36:32 Annotating text fragment 20271/42019
## 2023-09-16 21:36:32 Annotating text fragment 20281/42019
## 2023-09-16 21:36:32 Annotating text fragment 20291/42019
## 2023-09-16 21:36:32 Annotating text fragment 20301/42019
## 2023-09-16 21:36:32 Annotating text fragment 20311/42019
## 2023-09-16 21:36:32 Annotating text fragment 20321/42019
## 2023-09-16 21:36:33 Annotating text fragment 20331/42019
## 2023-09-16 21:36:33 Annotating text fragment 20341/42019
## 2023-09-16 21:36:33 Annotating text fragment 20351/42019
## 2023-09-16 21:36:33 Annotating text fragment 20361/42019
## 2023-09-16 21:36:33 Annotating text fragment 20371/42019
## 2023-09-16 21:36:33 Annotating text fragment 20381/42019
## 2023-09-16 21:36:33 Annotating text fragment 20391/42019
## 2023-09-16 21:36:33 Annotating text fragment 20401/42019
## 2023-09-16 21:36:33 Annotating text fragment 20411/42019
## 2023-09-16 21:36:33 Annotating text fragment 20421/42019
## 2023-09-16 21:36:33 Annotating text fragment 20431/42019
## 2023-09-16 21:36:33 Annotating text fragment 20441/42019
## 2023-09-16 21:36:33 Annotating text fragment 20451/42019
## 2023-09-16 21:36:34 Annotating text fragment 20461/42019
## 2023-09-16 21:36:34 Annotating text fragment 20471/42019
## 2023-09-16 21:36:34 Annotating text fragment 20481/42019
## 2023-09-16 21:36:34 Annotating text fragment 20491/42019
## 2023-09-16 21:36:34 Annotating text fragment 20501/42019
## 2023-09-16 21:36:34 Annotating text fragment 20511/42019
## 2023-09-16 21:36:34 Annotating text fragment 20521/42019
## 2023-09-16 21:36:35 Annotating text fragment 20531/42019
## 2023-09-16 21:36:35 Annotating text fragment 20541/42019
## 2023-09-16 21:36:35 Annotating text fragment 20551/42019
## 2023-09-16 21:36:35 Annotating text fragment 20561/42019
## 2023-09-16 21:36:35 Annotating text fragment 20571/42019
## 2023-09-16 21:36:35 Annotating text fragment 20581/42019
## 2023-09-16 21:36:35 Annotating text fragment 20591/42019
## 2023-09-16 21:36:35 Annotating text fragment 20601/42019
## 2023-09-16 21:36:35 Annotating text fragment 20611/42019
## 2023-09-16 21:36:35 Annotating text fragment 20621/42019
## 2023-09-16 21:36:35 Annotating text fragment 20631/42019
## 2023-09-16 21:36:36 Annotating text fragment 20641/42019
## 2023-09-16 21:36:36 Annotating text fragment 20651/42019
## 2023-09-16 21:36:36 Annotating text fragment 20661/42019
```

```
## 2023-09-16 21:36:36 Annotating text fragment 20671/42019
## 2023-09-16 21:36:36 Annotating text fragment 20681/42019
## 2023-09-16 21:36:36 Annotating text fragment 20691/42019
## 2023-09-16 21:36:36 Annotating text fragment 20701/42019
## 2023-09-16 21:36:36 Annotating text fragment 20711/42019
## 2023-09-16 21:36:36 Annotating text fragment 20721/42019
## 2023-09-16 21:36:36 Annotating text fragment 20731/42019
## 2023-09-16 21:36:37 Annotating text fragment 20741/42019
## 2023-09-16 21:36:37 Annotating text fragment 20751/42019
## 2023-09-16 21:36:37 Annotating text fragment 20761/42019
## 2023-09-16 21:36:37 Annotating text fragment 20771/42019
## 2023-09-16 21:36:37 Annotating text fragment 20781/42019
## 2023-09-16 21:36:37 Annotating text fragment 20791/42019
## 2023-09-16 21:36:37 Annotating text fragment 20801/42019
## 2023-09-16 21:36:37 Annotating text fragment 20811/42019
## 2023-09-16 21:36:37 Annotating text fragment 20821/42019
## 2023-09-16 21:36:38 Annotating text fragment 20831/42019
## 2023-09-16 21:36:38 Annotating text fragment 20841/42019
## 2023-09-16 21:36:38 Annotating text fragment 20851/42019
## 2023-09-16 21:36:38 Annotating text fragment 20861/42019
## 2023-09-16 21:36:38 Annotating text fragment 20871/42019
## 2023-09-16 21:36:38 Annotating text fragment 20881/42019
## 2023-09-16 21:36:38 Annotating text fragment 20891/42019
## 2023-09-16 21:36:38 Annotating text fragment 20901/42019
## 2023-09-16 21:36:38 Annotating text fragment 20911/42019
## 2023-09-16 21:36:38 Annotating text fragment 20921/42019
## 2023-09-16 21:36:38 Annotating text fragment 20931/42019
## 2023-09-16 21:36:39 Annotating text fragment 20941/42019
## 2023-09-16 21:36:39 Annotating text fragment 20951/42019
## 2023-09-16 21:36:39 Annotating text fragment 20961/42019
## 2023-09-16 21:36:39 Annotating text fragment 20971/42019
## 2023-09-16 21:36:39 Annotating text fragment 20981/42019
## 2023-09-16 21:36:39 Annotating text fragment 20991/42019
## 2023-09-16 21:36:39 Annotating text fragment 21001/42019
## 2023-09-16 21:36:39 Annotating text fragment 21011/42019
## 2023-09-16 21:36:39 Annotating text fragment 21021/42019
## 2023-09-16 21:36:39 Annotating text fragment 21031/42019
## 2023-09-16 21:36:40 Annotating text fragment 21041/42019
## 2023-09-16 21:36:40 Annotating text fragment 21051/42019
## 2023-09-16 21:36:41 Annotating text fragment 21061/42019
## 2023-09-16 21:36:41 Annotating text fragment 21071/42019
## 2023-09-16 21:36:41 Annotating text fragment 21081/42019
## 2023-09-16 21:36:41 Annotating text fragment 21091/42019
## 2023-09-16 21:36:41 Annotating text fragment 21101/42019
## 2023-09-16 21:36:41 Annotating text fragment 21111/42019
## 2023-09-16 21:36:41 Annotating text fragment 21121/42019
## 2023-09-16 21:36:41 Annotating text fragment 21131/42019
## 2023-09-16 21:36:41 Annotating text fragment 21141/42019
## 2023-09-16 21:36:41 Annotating text fragment 21151/42019
## 2023-09-16 21:36:42 Annotating text fragment 21161/42019
## 2023-09-16 21:36:42 Annotating text fragment 21171/42019
## 2023-09-16 21:36:42 Annotating text fragment 21181/42019
## 2023-09-16 21:36:42 Annotating text fragment 21191/42019
## 2023-09-16 21:36:42 Annotating text fragment 21201/42019
```

```
## 2023-09-16 21:36:42 Annotating text fragment 21211/42019
## 2023-09-16 21:36:42 Annotating text fragment 21221/42019
## 2023-09-16 21:36:42 Annotating text fragment 21231/42019
## 2023-09-16 21:36:42 Annotating text fragment 21241/42019
## 2023-09-16 21:36:42 Annotating text fragment 21251/42019
## 2023-09-16 21:36:42 Annotating text fragment 21261/42019
## 2023-09-16 21:36:42 Annotating text fragment 21271/42019
## 2023-09-16 21:36:43 Annotating text fragment 21281/42019
## 2023-09-16 21:36:43 Annotating text fragment 21291/42019
## 2023-09-16 21:36:43 Annotating text fragment 21301/42019
## 2023-09-16 21:36:43 Annotating text fragment 21311/42019
## 2023-09-16 21:36:43 Annotating text fragment 21321/42019
## 2023-09-16 21:36:44 Annotating text fragment 21331/42019
## 2023-09-16 21:36:44 Annotating text fragment 21341/42019
## 2023-09-16 21:36:44 Annotating text fragment 21351/42019
## 2023-09-16 21:36:44 Annotating text fragment 21361/42019
## 2023-09-16 21:36:44 Annotating text fragment 21371/42019
## 2023-09-16 21:36:44 Annotating text fragment 21381/42019
## 2023-09-16 21:36:44 Annotating text fragment 21391/42019
## 2023-09-16 21:36:44 Annotating text fragment 21401/42019
## 2023-09-16 21:36:44 Annotating text fragment 21411/42019
## 2023-09-16 21:36:44 Annotating text fragment 21421/42019
## 2023-09-16 21:36:45 Annotating text fragment 21431/42019
## 2023-09-16 21:36:45 Annotating text fragment 21441/42019
## 2023-09-16 21:36:45 Annotating text fragment 21451/42019
## 2023-09-16 21:36:45 Annotating text fragment 21461/42019
## 2023-09-16 21:36:45 Annotating text fragment 21471/42019
## 2023-09-16 21:36:45 Annotating text fragment 21481/42019
## 2023-09-16 21:36:45 Annotating text fragment 21491/42019
## 2023-09-16 21:36:45 Annotating text fragment 21501/42019
## 2023-09-16 21:36:45 Annotating text fragment 21511/42019
## 2023-09-16 21:36:45 Annotating text fragment 21521/42019
## 2023-09-16 21:36:45 Annotating text fragment 21531/42019
## 2023-09-16 21:36:46 Annotating text fragment 21541/42019
## 2023-09-16 21:36:46 Annotating text fragment 21551/42019
## 2023-09-16 21:36:46 Annotating text fragment 21561/42019
## 2023-09-16 21:36:46 Annotating text fragment 21571/42019
## 2023-09-16 21:36:46 Annotating text fragment 21581/42019
## 2023-09-16 21:36:46 Annotating text fragment 21591/42019
## 2023-09-16 21:36:46 Annotating text fragment 21601/42019
## 2023-09-16 21:36:46 Annotating text fragment 21611/42019
## 2023-09-16 21:36:46 Annotating text fragment 21621/42019
## 2023-09-16 21:36:46 Annotating text fragment 21631/42019
## 2023-09-16 21:36:46 Annotating text fragment 21641/42019
## 2023-09-16 21:36:46 Annotating text fragment 21651/42019
## 2023-09-16 21:36:47 Annotating text fragment 21661/42019
## 2023-09-16 21:36:47 Annotating text fragment 21671/42019
## 2023-09-16 21:36:47 Annotating text fragment 21681/42019
## 2023-09-16 21:36:47 Annotating text fragment 21691/42019
## 2023-09-16 21:36:47 Annotating text fragment 21701/42019
## 2023-09-16 21:36:47 Annotating text fragment 21711/42019
## 2023-09-16 21:36:47 Annotating text fragment 21721/42019
## 2023-09-16 21:36:47 Annotating text fragment 21731/42019
## 2023-09-16 21:36:48 Annotating text fragment 21741/42019
```

```
## 2023-09-16 21:36:48 Annotating text fragment 21751/42019
## 2023-09-16 21:36:48 Annotating text fragment 21761/42019
## 2023-09-16 21:36:48 Annotating text fragment 21771/42019
## 2023-09-16 21:36:48 Annotating text fragment 21781/42019
## 2023-09-16 21:36:48 Annotating text fragment 21791/42019
## 2023-09-16 21:36:48 Annotating text fragment 21801/42019
## 2023-09-16 21:36:49 Annotating text fragment 21811/42019
## 2023-09-16 21:36:49 Annotating text fragment 21821/42019
## 2023-09-16 21:36:49 Annotating text fragment 21831/42019
## 2023-09-16 21:36:49 Annotating text fragment 21841/42019
## 2023-09-16 21:36:49 Annotating text fragment 21851/42019
## 2023-09-16 21:36:49 Annotating text fragment 21861/42019
## 2023-09-16 21:36:49 Annotating text fragment 21871/42019
## 2023-09-16 21:36:49 Annotating text fragment 21881/42019
## 2023-09-16 21:36:49 Annotating text fragment 21891/42019
## 2023-09-16 21:36:49 Annotating text fragment 21901/42019
## 2023-09-16 21:36:50 Annotating text fragment 21911/42019
## 2023-09-16 21:36:50 Annotating text fragment 21921/42019
## 2023-09-16 21:36:50 Annotating text fragment 21931/42019
## 2023-09-16 21:36:50 Annotating text fragment 21941/42019
## 2023-09-16 21:36:50 Annotating text fragment 21951/42019
## 2023-09-16 21:36:50 Annotating text fragment 21961/42019
## 2023-09-16 21:36:50 Annotating text fragment 21971/42019
## 2023-09-16 21:36:51 Annotating text fragment 21981/42019
## 2023-09-16 21:36:51 Annotating text fragment 21991/42019
## 2023-09-16 21:36:51 Annotating text fragment 22001/42019
## 2023-09-16 21:36:51 Annotating text fragment 22011/42019
## 2023-09-16 21:36:51 Annotating text fragment 22021/42019
## 2023-09-16 21:36:51 Annotating text fragment 22031/42019
## 2023-09-16 21:36:51 Annotating text fragment 22041/42019
## 2023-09-16 21:36:51 Annotating text fragment 22051/42019
## 2023-09-16 21:36:51 Annotating text fragment 22061/42019
## 2023-09-16 21:36:51 Annotating text fragment 22071/42019
## 2023-09-16 21:36:52 Annotating text fragment 22081/42019
## 2023-09-16 21:36:52 Annotating text fragment 22091/42019
## 2023-09-16 21:36:52 Annotating text fragment 22101/42019
## 2023-09-16 21:36:52 Annotating text fragment 22111/42019
## 2023-09-16 21:36:52 Annotating text fragment 22121/42019
## 2023-09-16 21:36:52 Annotating text fragment 22131/42019
## 2023-09-16 21:36:52 Annotating text fragment 22141/42019
## 2023-09-16 21:36:52 Annotating text fragment 22151/42019
## 2023-09-16 21:36:52 Annotating text fragment 22161/42019
## 2023-09-16 21:36:52 Annotating text fragment 22171/42019
## 2023-09-16 21:36:52 Annotating text fragment 22181/42019
## 2023-09-16 21:36:52 Annotating text fragment 22191/42019
## 2023-09-16 21:36:53 Annotating text fragment 22201/42019
## 2023-09-16 21:36:53 Annotating text fragment 22211/42019
## 2023-09-16 21:36:53 Annotating text fragment 22221/42019
## 2023-09-16 21:36:53 Annotating text fragment 22231/42019
## 2023-09-16 21:36:53 Annotating text fragment 22241/42019
## 2023-09-16 21:36:53 Annotating text fragment 22251/42019
## 2023-09-16 21:36:53 Annotating text fragment 22261/42019
## 2023-09-16 21:36:54 Annotating text fragment 22271/42019
## 2023-09-16 21:36:54 Annotating text fragment 22281/42019
```

```
## 2023-09-16 21:36:54 Annotating text fragment 22291/42019
## 2023-09-16 21:36:54 Annotating text fragment 22301/42019
## 2023-09-16 21:36:54 Annotating text fragment 22311/42019
## 2023-09-16 21:36:54 Annotating text fragment 22321/42019
## 2023-09-16 21:36:54 Annotating text fragment 22331/42019
## 2023-09-16 21:36:54 Annotating text fragment 22341/42019
## 2023-09-16 21:36:54 Annotating text fragment 22351/42019
## 2023-09-16 21:36:54 Annotating text fragment 22361/42019
## 2023-09-16 21:36:54 Annotating text fragment 22371/42019
## 2023-09-16 21:36:54 Annotating text fragment 22381/42019
## 2023-09-16 21:36:54 Annotating text fragment 22391/42019
## 2023-09-16 21:36:55 Annotating text fragment 22401/42019
## 2023-09-16 21:36:55 Annotating text fragment 22411/42019
## 2023-09-16 21:36:55 Annotating text fragment 22421/42019
## 2023-09-16 21:36:55 Annotating text fragment 22431/42019
## 2023-09-16 21:36:55 Annotating text fragment 22441/42019
## 2023-09-16 21:36:55 Annotating text fragment 22451/42019
## 2023-09-16 21:36:55 Annotating text fragment 22461/42019
## 2023-09-16 21:36:55 Annotating text fragment 22471/42019
## 2023-09-16 21:36:55 Annotating text fragment 22481/42019
## 2023-09-16 21:36:55 Annotating text fragment 22491/42019
## 2023-09-16 21:36:55 Annotating text fragment 22501/42019
## 2023-09-16 21:36:55 Annotating text fragment 22511/42019
## 2023-09-16 21:36:56 Annotating text fragment 22521/42019
## 2023-09-16 21:36:56 Annotating text fragment 22531/42019
## 2023-09-16 21:36:56 Annotating text fragment 22541/42019
## 2023-09-16 21:36:56 Annotating text fragment 22551/42019
## 2023-09-16 21:36:56 Annotating text fragment 22561/42019
## 2023-09-16 21:36:56 Annotating text fragment 22571/42019
## 2023-09-16 21:36:56 Annotating text fragment 22581/42019
## 2023-09-16 21:36:56 Annotating text fragment 22591/42019
## 2023-09-16 21:36:57 Annotating text fragment 22601/42019
## 2023-09-16 21:36:57 Annotating text fragment 22611/42019
## 2023-09-16 21:36:57 Annotating text fragment 22621/42019
## 2023-09-16 21:36:57 Annotating text fragment 22631/42019
## 2023-09-16 21:36:57 Annotating text fragment 22641/42019
## 2023-09-16 21:36:57 Annotating text fragment 22651/42019
## 2023-09-16 21:36:57 Annotating text fragment 22661/42019
## 2023-09-16 21:36:57 Annotating text fragment 22671/42019
## 2023-09-16 21:36:57 Annotating text fragment 22681/42019
## 2023-09-16 21:36:57 Annotating text fragment 22691/42019
## 2023-09-16 21:36:57 Annotating text fragment 22701/42019
## 2023-09-16 21:36:57 Annotating text fragment 22711/42019
## 2023-09-16 21:36:57 Annotating text fragment 22721/42019
## 2023-09-16 21:36:58 Annotating text fragment 22731/42019
## 2023-09-16 21:36:58 Annotating text fragment 22741/42019
## 2023-09-16 21:36:58 Annotating text fragment 22751/42019
## 2023-09-16 21:36:58 Annotating text fragment 22761/42019
## 2023-09-16 21:36:58 Annotating text fragment 22771/42019
## 2023-09-16 21:36:58 Annotating text fragment 22781/42019
## 2023-09-16 21:36:58 Annotating text fragment 22791/42019
## 2023-09-16 21:36:58 Annotating text fragment 22801/42019
## 2023-09-16 21:36:58 Annotating text fragment 22811/42019
## 2023-09-16 21:36:58 Annotating text fragment 22821/42019
```

```
## 2023-09-16 21:36:58 Annotating text fragment 22831/42019
## 2023-09-16 21:36:59 Annotating text fragment 22841/42019
## 2023-09-16 21:36:59 Annotating text fragment 22851/42019
## 2023-09-16 21:36:59 Annotating text fragment 22861/42019
## 2023-09-16 21:36:59 Annotating text fragment 22871/42019
## 2023-09-16 21:36:59 Annotating text fragment 22881/42019
## 2023-09-16 21:36:59 Annotating text fragment 22891/42019
## 2023-09-16 21:36:59 Annotating text fragment 22901/42019
## 2023-09-16 21:36:59 Annotating text fragment 22911/42019
## 2023-09-16 21:36:59 Annotating text fragment 22921/42019
## 2023-09-16 21:36:59 Annotating text fragment 22931/42019
## 2023-09-16 21:37:00 Annotating text fragment 22941/42019
## 2023-09-16 21:37:00 Annotating text fragment 22951/42019
## 2023-09-16 21:37:00 Annotating text fragment 22961/42019
## 2023-09-16 21:37:00 Annotating text fragment 22971/42019
## 2023-09-16 21:37:00 Annotating text fragment 22981/42019
## 2023-09-16 21:37:00 Annotating text fragment 22991/42019
## 2023-09-16 21:37:00 Annotating text fragment 23001/42019
## 2023-09-16 21:37:00 Annotating text fragment 23011/42019
## 2023-09-16 21:37:00 Annotating text fragment 23021/42019
## 2023-09-16 21:37:00 Annotating text fragment 23031/42019
## 2023-09-16 21:37:00 Annotating text fragment 23041/42019
## 2023-09-16 21:37:00 Annotating text fragment 23051/42019
## 2023-09-16 21:37:01 Annotating text fragment 23061/42019
## 2023-09-16 21:37:01 Annotating text fragment 23071/42019
## 2023-09-16 21:37:01 Annotating text fragment 23081/42019
## 2023-09-16 21:37:01 Annotating text fragment 23091/42019
## 2023-09-16 21:37:01 Annotating text fragment 23101/42019
## 2023-09-16 21:37:01 Annotating text fragment 23111/42019
## 2023-09-16 21:37:01 Annotating text fragment 23121/42019
## 2023-09-16 21:37:01 Annotating text fragment 23131/42019
## 2023-09-16 21:37:01 Annotating text fragment 23141/42019
## 2023-09-16 21:37:01 Annotating text fragment 23151/42019
## 2023-09-16 21:37:02 Annotating text fragment 23161/42019
## 2023-09-16 21:37:02 Annotating text fragment 23171/42019
## 2023-09-16 21:37:02 Annotating text fragment 23181/42019
## 2023-09-16 21:37:02 Annotating text fragment 23191/42019
## 2023-09-16 21:37:02 Annotating text fragment 23201/42019
## 2023-09-16 21:37:02 Annotating text fragment 23211/42019
## 2023-09-16 21:37:02 Annotating text fragment 23221/42019
## 2023-09-16 21:37:02 Annotating text fragment 23231/42019
## 2023-09-16 21:37:02 Annotating text fragment 23241/42019
## 2023-09-16 21:37:02 Annotating text fragment 23251/42019
## 2023-09-16 21:37:02 Annotating text fragment 23261/42019
## 2023-09-16 21:37:03 Annotating text fragment 23271/42019
## 2023-09-16 21:37:03 Annotating text fragment 23281/42019
## 2023-09-16 21:37:03 Annotating text fragment 23291/42019
## 2023-09-16 21:37:03 Annotating text fragment 23301/42019
## 2023-09-16 21:37:03 Annotating text fragment 23311/42019
## 2023-09-16 21:37:03 Annotating text fragment 23321/42019
## 2023-09-16 21:37:03 Annotating text fragment 23331/42019
## 2023-09-16 21:37:03 Annotating text fragment 23341/42019
## 2023-09-16 21:37:04 Annotating text fragment 23351/42019
## 2023-09-16 21:37:04 Annotating text fragment 23361/42019
```

```
## 2023-09-16 21:37:04 Annotating text fragment 23371/42019
## 2023-09-16 21:37:04 Annotating text fragment 23381/42019
## 2023-09-16 21:37:04 Annotating text fragment 23391/42019
## 2023-09-16 21:37:04 Annotating text fragment 23401/42019
## 2023-09-16 21:37:04 Annotating text fragment 23411/42019
## 2023-09-16 21:37:04 Annotating text fragment 23421/42019
## 2023-09-16 21:37:04 Annotating text fragment 23431/42019
## 2023-09-16 21:37:04 Annotating text fragment 23441/42019
## 2023-09-16 21:37:04 Annotating text fragment 23451/42019
## 2023-09-16 21:37:05 Annotating text fragment 23461/42019
## 2023-09-16 21:37:05 Annotating text fragment 23471/42019
## 2023-09-16 21:37:05 Annotating text fragment 23481/42019
## 2023-09-16 21:37:05 Annotating text fragment 23491/42019
## 2023-09-16 21:37:05 Annotating text fragment 23501/42019
## 2023-09-16 21:37:05 Annotating text fragment 23511/42019
## 2023-09-16 21:37:05 Annotating text fragment 23521/42019
## 2023-09-16 21:37:05 Annotating text fragment 23531/42019
## 2023-09-16 21:37:05 Annotating text fragment 23541/42019
## 2023-09-16 21:37:06 Annotating text fragment 23551/42019
## 2023-09-16 21:37:06 Annotating text fragment 23561/42019
## 2023-09-16 21:37:06 Annotating text fragment 23571/42019
## 2023-09-16 21:37:06 Annotating text fragment 23581/42019
## 2023-09-16 21:37:06 Annotating text fragment 23591/42019
## 2023-09-16 21:37:06 Annotating text fragment 23601/42019
## 2023-09-16 21:37:06 Annotating text fragment 23611/42019
## 2023-09-16 21:37:06 Annotating text fragment 23621/42019
## 2023-09-16 21:37:06 Annotating text fragment 23631/42019
## 2023-09-16 21:37:06 Annotating text fragment 23641/42019
## 2023-09-16 21:37:06 Annotating text fragment 23651/42019
## 2023-09-16 21:37:07 Annotating text fragment 23661/42019
## 2023-09-16 21:37:07 Annotating text fragment 23671/42019
## 2023-09-16 21:37:07 Annotating text fragment 23681/42019
## 2023-09-16 21:37:07 Annotating text fragment 23691/42019
## 2023-09-16 21:37:07 Annotating text fragment 23701/42019
## 2023-09-16 21:37:07 Annotating text fragment 23711/42019
## 2023-09-16 21:37:07 Annotating text fragment 23721/42019
## 2023-09-16 21:37:07 Annotating text fragment 23731/42019
## 2023-09-16 21:37:07 Annotating text fragment 23741/42019
## 2023-09-16 21:37:07 Annotating text fragment 23751/42019
## 2023-09-16 21:37:07 Annotating text fragment 23761/42019
## 2023-09-16 21:37:08 Annotating text fragment 23771/42019
## 2023-09-16 21:37:08 Annotating text fragment 23781/42019
## 2023-09-16 21:37:08 Annotating text fragment 23791/42019
## 2023-09-16 21:37:08 Annotating text fragment 23801/42019
## 2023-09-16 21:37:08 Annotating text fragment 23811/42019
## 2023-09-16 21:37:08 Annotating text fragment 23821/42019
## 2023-09-16 21:37:08 Annotating text fragment 23831/42019
## 2023-09-16 21:37:08 Annotating text fragment 23841/42019
## 2023-09-16 21:37:08 Annotating text fragment 23851/42019
## 2023-09-16 21:37:08 Annotating text fragment 23861/42019
## 2023-09-16 21:37:08 Annotating text fragment 23871/42019
## 2023-09-16 21:37:08 Annotating text fragment 23881/42019
## 2023-09-16 21:37:09 Annotating text fragment 23891/42019
## 2023-09-16 21:37:09 Annotating text fragment 23901/42019
```

```
## 2023-09-16 21:37:09 Annotating text fragment 23911/42019
## 2023-09-16 21:37:09 Annotating text fragment 23921/42019
## 2023-09-16 21:37:09 Annotating text fragment 23931/42019
## 2023-09-16 21:37:09 Annotating text fragment 23941/42019
## 2023-09-16 21:37:09 Annotating text fragment 23951/42019
## 2023-09-16 21:37:09 Annotating text fragment 23961/42019
## 2023-09-16 21:37:09 Annotating text fragment 23971/42019
## 2023-09-16 21:37:09 Annotating text fragment 23981/42019
## 2023-09-16 21:37:09 Annotating text fragment 23991/42019
## 2023-09-16 21:37:10 Annotating text fragment 24001/42019
## 2023-09-16 21:37:10 Annotating text fragment 24011/42019
## 2023-09-16 21:37:10 Annotating text fragment 24021/42019
## 2023-09-16 21:37:10 Annotating text fragment 24031/42019
## 2023-09-16 21:37:10 Annotating text fragment 24041/42019
## 2023-09-16 21:37:10 Annotating text fragment 24051/42019
## 2023-09-16 21:37:10 Annotating text fragment 24061/42019
## 2023-09-16 21:37:10 Annotating text fragment 24071/42019
## 2023-09-16 21:37:10 Annotating text fragment 24081/42019
## 2023-09-16 21:37:10 Annotating text fragment 24091/42019
## 2023-09-16 21:37:10 Annotating text fragment 24101/42019
## 2023-09-16 21:37:10 Annotating text fragment 24111/42019
## 2023-09-16 21:37:10 Annotating text fragment 24121/42019
## 2023-09-16 21:37:10 Annotating text fragment 24131/42019
## 2023-09-16 21:37:11 Annotating text fragment 24141/42019
## 2023-09-16 21:37:11 Annotating text fragment 24151/42019
## 2023-09-16 21:37:11 Annotating text fragment 24161/42019
## 2023-09-16 21:37:11 Annotating text fragment 24171/42019
## 2023-09-16 21:37:11 Annotating text fragment 24181/42019
## 2023-09-16 21:37:11 Annotating text fragment 24191/42019
## 2023-09-16 21:37:11 Annotating text fragment 24201/42019
## 2023-09-16 21:37:11 Annotating text fragment 24211/42019
## 2023-09-16 21:37:11 Annotating text fragment 24221/42019
## 2023-09-16 21:37:11 Annotating text fragment 24231/42019
## 2023-09-16 21:37:11 Annotating text fragment 24241/42019
## 2023-09-16 21:37:11 Annotating text fragment 24251/42019
## 2023-09-16 21:37:11 Annotating text fragment 24261/42019
## 2023-09-16 21:37:12 Annotating text fragment 24271/42019
## 2023-09-16 21:37:12 Annotating text fragment 24281/42019
## 2023-09-16 21:37:12 Annotating text fragment 24291/42019
## 2023-09-16 21:37:12 Annotating text fragment 24301/42019
## 2023-09-16 21:37:12 Annotating text fragment 24311/42019
## 2023-09-16 21:37:12 Annotating text fragment 24321/42019
## 2023-09-16 21:37:12 Annotating text fragment 24331/42019
## 2023-09-16 21:37:12 Annotating text fragment 24341/42019
## 2023-09-16 21:37:12 Annotating text fragment 24351/42019
## 2023-09-16 21:37:12 Annotating text fragment 24361/42019
## 2023-09-16 21:37:13 Annotating text fragment 24371/42019
## 2023-09-16 21:37:13 Annotating text fragment 24381/42019
## 2023-09-16 21:37:13 Annotating text fragment 24391/42019
## 2023-09-16 21:37:13 Annotating text fragment 24401/42019
## 2023-09-16 21:37:13 Annotating text fragment 24411/42019
## 2023-09-16 21:37:13 Annotating text fragment 24421/42019
## 2023-09-16 21:37:13 Annotating text fragment 24431/42019
## 2023-09-16 21:37:13 Annotating text fragment 24441/42019
```

```
## 2023-09-16 21:37:13 Annotating text fragment 24451/42019
## 2023-09-16 21:37:13 Annotating text fragment 24461/42019
## 2023-09-16 21:37:13 Annotating text fragment 24471/42019
## 2023-09-16 21:37:13 Annotating text fragment 24481/42019
## 2023-09-16 21:37:14 Annotating text fragment 24491/42019
## 2023-09-16 21:37:14 Annotating text fragment 24501/42019
## 2023-09-16 21:37:14 Annotating text fragment 24511/42019
## 2023-09-16 21:37:14 Annotating text fragment 24521/42019
## 2023-09-16 21:37:14 Annotating text fragment 24531/42019
## 2023-09-16 21:37:14 Annotating text fragment 24541/42019
## 2023-09-16 21:37:14 Annotating text fragment 24551/42019
## 2023-09-16 21:37:14 Annotating text fragment 24561/42019
## 2023-09-16 21:37:14 Annotating text fragment 24571/42019
## 2023-09-16 21:37:14 Annotating text fragment 24581/42019
## 2023-09-16 21:37:14 Annotating text fragment 24591/42019
## 2023-09-16 21:37:14 Annotating text fragment 24601/42019
## 2023-09-16 21:37:15 Annotating text fragment 24611/42019
## 2023-09-16 21:37:15 Annotating text fragment 24621/42019
## 2023-09-16 21:37:15 Annotating text fragment 24631/42019
## 2023-09-16 21:37:15 Annotating text fragment 24641/42019
## 2023-09-16 21:37:15 Annotating text fragment 24651/42019
## 2023-09-16 21:37:15 Annotating text fragment 24661/42019
## 2023-09-16 21:37:15 Annotating text fragment 24671/42019
## 2023-09-16 21:37:15 Annotating text fragment 24681/42019
## 2023-09-16 21:37:15 Annotating text fragment 24691/42019
## 2023-09-16 21:37:15 Annotating text fragment 24701/42019
## 2023-09-16 21:37:16 Annotating text fragment 24711/42019
## 2023-09-16 21:37:16 Annotating text fragment 24721/42019
## 2023-09-16 21:37:16 Annotating text fragment 24731/42019
## 2023-09-16 21:37:16 Annotating text fragment 24741/42019
## 2023-09-16 21:37:16 Annotating text fragment 24751/42019
## 2023-09-16 21:37:16 Annotating text fragment 24761/42019
## 2023-09-16 21:37:16 Annotating text fragment 24771/42019
## 2023-09-16 21:37:16 Annotating text fragment 24781/42019
## 2023-09-16 21:37:16 Annotating text fragment 24791/42019
## 2023-09-16 21:37:16 Annotating text fragment 24801/42019
## 2023-09-16 21:37:17 Annotating text fragment 24811/42019
## 2023-09-16 21:37:17 Annotating text fragment 24821/42019
## 2023-09-16 21:37:17 Annotating text fragment 24831/42019
## 2023-09-16 21:37:17 Annotating text fragment 24841/42019
## 2023-09-16 21:37:17 Annotating text fragment 24851/42019
## 2023-09-16 21:37:17 Annotating text fragment 24861/42019
## 2023-09-16 21:37:17 Annotating text fragment 24871/42019
## 2023-09-16 21:37:17 Annotating text fragment 24881/42019
## 2023-09-16 21:37:17 Annotating text fragment 24891/42019
## 2023-09-16 21:37:17 Annotating text fragment 24901/42019
## 2023-09-16 21:37:17 Annotating text fragment 24911/42019
## 2023-09-16 21:37:17 Annotating text fragment 24921/42019
## 2023-09-16 21:37:17 Annotating text fragment 24931/42019
## 2023-09-16 21:37:18 Annotating text fragment 24941/42019
## 2023-09-16 21:37:18 Annotating text fragment 24951/42019
## 2023-09-16 21:37:18 Annotating text fragment 24961/42019
## 2023-09-16 21:37:18 Annotating text fragment 24971/42019
## 2023-09-16 21:37:18 Annotating text fragment 24981/42019
```

```
## 2023-09-16 21:37:18 Annotating text fragment 24991/42019
## 2023-09-16 21:37:18 Annotating text fragment 25001/42019
## 2023-09-16 21:37:18 Annotating text fragment 25011/42019
## 2023-09-16 21:37:18 Annotating text fragment 25021/42019
## 2023-09-16 21:37:19 Annotating text fragment 25031/42019
## 2023-09-16 21:37:19 Annotating text fragment 25041/42019
## 2023-09-16 21:37:19 Annotating text fragment 25051/42019
## 2023-09-16 21:37:19 Annotating text fragment 25061/42019
## 2023-09-16 21:37:19 Annotating text fragment 25071/42019
## 2023-09-16 21:37:19 Annotating text fragment 25081/42019
## 2023-09-16 21:37:19 Annotating text fragment 25091/42019
## 2023-09-16 21:37:19 Annotating text fragment 25101/42019
## 2023-09-16 21:37:19 Annotating text fragment 25111/42019
## 2023-09-16 21:37:19 Annotating text fragment 25121/42019
## 2023-09-16 21:37:20 Annotating text fragment 25131/42019
## 2023-09-16 21:37:20 Annotating text fragment 25141/42019
## 2023-09-16 21:37:20 Annotating text fragment 25151/42019
## 2023-09-16 21:37:20 Annotating text fragment 25161/42019
## 2023-09-16 21:37:20 Annotating text fragment 25171/42019
## 2023-09-16 21:37:20 Annotating text fragment 25181/42019
## 2023-09-16 21:37:21 Annotating text fragment 25191/42019
## 2023-09-16 21:37:21 Annotating text fragment 25201/42019
## 2023-09-16 21:37:21 Annotating text fragment 25211/42019
## 2023-09-16 21:37:21 Annotating text fragment 25221/42019
## 2023-09-16 21:37:21 Annotating text fragment 25231/42019
## 2023-09-16 21:37:21 Annotating text fragment 25241/42019
## 2023-09-16 21:37:21 Annotating text fragment 25251/42019
## 2023-09-16 21:37:22 Annotating text fragment 25261/42019
## 2023-09-16 21:37:22 Annotating text fragment 25271/42019
## 2023-09-16 21:37:22 Annotating text fragment 25281/42019
## 2023-09-16 21:37:22 Annotating text fragment 25291/42019
## 2023-09-16 21:37:22 Annotating text fragment 25301/42019
## 2023-09-16 21:37:22 Annotating text fragment 25311/42019
## 2023-09-16 21:37:22 Annotating text fragment 25321/42019
## 2023-09-16 21:37:22 Annotating text fragment 25331/42019
## 2023-09-16 21:37:23 Annotating text fragment 25341/42019
## 2023-09-16 21:37:23 Annotating text fragment 25351/42019
## 2023-09-16 21:37:23 Annotating text fragment 25361/42019
## 2023-09-16 21:37:23 Annotating text fragment 25371/42019
## 2023-09-16 21:37:23 Annotating text fragment 25381/42019
## 2023-09-16 21:37:23 Annotating text fragment 25391/42019
## 2023-09-16 21:37:23 Annotating text fragment 25401/42019
## 2023-09-16 21:37:23 Annotating text fragment 25411/42019
## 2023-09-16 21:37:23 Annotating text fragment 25421/42019
## 2023-09-16 21:37:23 Annotating text fragment 25431/42019
## 2023-09-16 21:37:23 Annotating text fragment 25441/42019
## 2023-09-16 21:37:24 Annotating text fragment 25451/42019
## 2023-09-16 21:37:24 Annotating text fragment 25461/42019
## 2023-09-16 21:37:24 Annotating text fragment 25471/42019
## 2023-09-16 21:37:24 Annotating text fragment 25481/42019
## 2023-09-16 21:37:24 Annotating text fragment 25491/42019
## 2023-09-16 21:37:24 Annotating text fragment 25501/42019
## 2023-09-16 21:37:24 Annotating text fragment 25511/42019
## 2023-09-16 21:37:24 Annotating text fragment 25521/42019
```

```
## 2023-09-16 21:37:25 Annotating text fragment 25531/42019
## 2023-09-16 21:37:25 Annotating text fragment 25541/42019
## 2023-09-16 21:37:25 Annotating text fragment 25551/42019
## 2023-09-16 21:37:25 Annotating text fragment 25561/42019
## 2023-09-16 21:37:25 Annotating text fragment 25571/42019
## 2023-09-16 21:37:25 Annotating text fragment 25581/42019
## 2023-09-16 21:37:26 Annotating text fragment 25591/42019
## 2023-09-16 21:37:26 Annotating text fragment 25601/42019
## 2023-09-16 21:37:26 Annotating text fragment 25611/42019
## 2023-09-16 21:37:26 Annotating text fragment 25621/42019
## 2023-09-16 21:37:26 Annotating text fragment 25631/42019
## 2023-09-16 21:37:26 Annotating text fragment 25641/42019
## 2023-09-16 21:37:26 Annotating text fragment 25651/42019
## 2023-09-16 21:37:26 Annotating text fragment 25661/42019
## 2023-09-16 21:37:27 Annotating text fragment 25671/42019
## 2023-09-16 21:37:27 Annotating text fragment 25681/42019
## 2023-09-16 21:37:27 Annotating text fragment 25691/42019
## 2023-09-16 21:37:27 Annotating text fragment 25701/42019
## 2023-09-16 21:37:27 Annotating text fragment 25711/42019
## 2023-09-16 21:37:27 Annotating text fragment 25721/42019
## 2023-09-16 21:37:27 Annotating text fragment 25731/42019
## 2023-09-16 21:37:27 Annotating text fragment 25741/42019
## 2023-09-16 21:37:27 Annotating text fragment 25751/42019
## 2023-09-16 21:37:27 Annotating text fragment 25761/42019
## 2023-09-16 21:37:27 Annotating text fragment 25771/42019
## 2023-09-16 21:37:28 Annotating text fragment 25781/42019
## 2023-09-16 21:37:28 Annotating text fragment 25791/42019
## 2023-09-16 21:37:28 Annotating text fragment 25801/42019
## 2023-09-16 21:37:28 Annotating text fragment 25811/42019
## 2023-09-16 21:37:28 Annotating text fragment 25821/42019
## 2023-09-16 21:37:28 Annotating text fragment 25831/42019
## 2023-09-16 21:37:28 Annotating text fragment 25841/42019
## 2023-09-16 21:37:28 Annotating text fragment 25851/42019
## 2023-09-16 21:37:28 Annotating text fragment 25861/42019
## 2023-09-16 21:37:28 Annotating text fragment 25871/42019
## 2023-09-16 21:37:29 Annotating text fragment 25881/42019
## 2023-09-16 21:37:29 Annotating text fragment 25891/42019
## 2023-09-16 21:37:29 Annotating text fragment 25901/42019
## 2023-09-16 21:37:29 Annotating text fragment 25911/42019
## 2023-09-16 21:37:29 Annotating text fragment 25921/42019
## 2023-09-16 21:37:29 Annotating text fragment 25931/42019
## 2023-09-16 21:37:29 Annotating text fragment 25941/42019
## 2023-09-16 21:37:29 Annotating text fragment 25951/42019
## 2023-09-16 21:37:29 Annotating text fragment 25961/42019
## 2023-09-16 21:37:29 Annotating text fragment 25971/42019
## 2023-09-16 21:37:30 Annotating text fragment 25981/42019
## 2023-09-16 21:37:30 Annotating text fragment 25991/42019
## 2023-09-16 21:37:30 Annotating text fragment 26001/42019
## 2023-09-16 21:37:30 Annotating text fragment 26011/42019
## 2023-09-16 21:37:30 Annotating text fragment 26021/42019
## 2023-09-16 21:37:30 Annotating text fragment 26031/42019
## 2023-09-16 21:37:30 Annotating text fragment 26041/42019
## 2023-09-16 21:37:30 Annotating text fragment 26051/42019
## 2023-09-16 21:37:30 Annotating text fragment 26061/42019
```

```
## 2023-09-16 21:37:30 Annotating text fragment 26071/42019
## 2023-09-16 21:37:30 Annotating text fragment 26081/42019
## 2023-09-16 21:37:31 Annotating text fragment 26091/42019
## 2023-09-16 21:37:31 Annotating text fragment 26101/42019
## 2023-09-16 21:37:31 Annotating text fragment 26111/42019
## 2023-09-16 21:37:31 Annotating text fragment 26121/42019
## 2023-09-16 21:37:31 Annotating text fragment 26131/42019
## 2023-09-16 21:37:32 Annotating text fragment 26141/42019
## 2023-09-16 21:37:32 Annotating text fragment 26151/42019
## 2023-09-16 21:37:32 Annotating text fragment 26161/42019
## 2023-09-16 21:37:32 Annotating text fragment 26171/42019
## 2023-09-16 21:37:32 Annotating text fragment 26181/42019
## 2023-09-16 21:37:32 Annotating text fragment 26191/42019
## 2023-09-16 21:37:32 Annotating text fragment 26201/42019
## 2023-09-16 21:37:33 Annotating text fragment 26211/42019
## 2023-09-16 21:37:33 Annotating text fragment 26221/42019
## 2023-09-16 21:37:33 Annotating text fragment 26231/42019
## 2023-09-16 21:37:33 Annotating text fragment 26241/42019
## 2023-09-16 21:37:34 Annotating text fragment 26251/42019
## 2023-09-16 21:37:34 Annotating text fragment 26261/42019
## 2023-09-16 21:37:34 Annotating text fragment 26271/42019
## 2023-09-16 21:37:34 Annotating text fragment 26281/42019
## 2023-09-16 21:37:34 Annotating text fragment 26291/42019
## 2023-09-16 21:37:34 Annotating text fragment 26301/42019
## 2023-09-16 21:37:34 Annotating text fragment 26311/42019
## 2023-09-16 21:37:34 Annotating text fragment 26321/42019
## 2023-09-16 21:37:35 Annotating text fragment 26331/42019
## 2023-09-16 21:37:35 Annotating text fragment 26341/42019
## 2023-09-16 21:37:35 Annotating text fragment 26351/42019
## 2023-09-16 21:37:35 Annotating text fragment 26361/42019
## 2023-09-16 21:37:35 Annotating text fragment 26371/42019
## 2023-09-16 21:37:35 Annotating text fragment 26381/42019
## 2023-09-16 21:37:35 Annotating text fragment 26391/42019
## 2023-09-16 21:37:35 Annotating text fragment 26401/42019
## 2023-09-16 21:37:35 Annotating text fragment 26411/42019
## 2023-09-16 21:37:35 Annotating text fragment 26421/42019
## 2023-09-16 21:37:35 Annotating text fragment 26431/42019
## 2023-09-16 21:37:36 Annotating text fragment 26441/42019
## 2023-09-16 21:37:36 Annotating text fragment 26451/42019
## 2023-09-16 21:37:36 Annotating text fragment 26461/42019
## 2023-09-16 21:37:36 Annotating text fragment 26471/42019
## 2023-09-16 21:37:36 Annotating text fragment 26481/42019
## 2023-09-16 21:37:36 Annotating text fragment 26491/42019
## 2023-09-16 21:37:36 Annotating text fragment 26501/42019
## 2023-09-16 21:37:36 Annotating text fragment 26511/42019
## 2023-09-16 21:37:36 Annotating text fragment 26521/42019
## 2023-09-16 21:37:37 Annotating text fragment 26531/42019
## 2023-09-16 21:37:37 Annotating text fragment 26541/42019
## 2023-09-16 21:37:37 Annotating text fragment 26551/42019
## 2023-09-16 21:37:37 Annotating text fragment 26561/42019
## 2023-09-16 21:37:37 Annotating text fragment 26571/42019
## 2023-09-16 21:37:37 Annotating text fragment 26581/42019
## 2023-09-16 21:37:38 Annotating text fragment 26591/42019
## 2023-09-16 21:37:38 Annotating text fragment 26601/42019
```

```
## 2023-09-16 21:37:38 Annotating text fragment 26611/42019
## 2023-09-16 21:37:38 Annotating text fragment 26621/42019
## 2023-09-16 21:37:38 Annotating text fragment 26631/42019
## 2023-09-16 21:37:38 Annotating text fragment 26641/42019
## 2023-09-16 21:37:38 Annotating text fragment 26651/42019
## 2023-09-16 21:37:38 Annotating text fragment 26661/42019
## 2023-09-16 21:37:38 Annotating text fragment 26671/42019
## 2023-09-16 21:37:39 Annotating text fragment 26681/42019
## 2023-09-16 21:37:39 Annotating text fragment 26691/42019
## 2023-09-16 21:37:39 Annotating text fragment 26701/42019
## 2023-09-16 21:37:39 Annotating text fragment 26711/42019
## 2023-09-16 21:37:39 Annotating text fragment 26721/42019
## 2023-09-16 21:37:39 Annotating text fragment 26731/42019
## 2023-09-16 21:37:39 Annotating text fragment 26741/42019
## 2023-09-16 21:37:39 Annotating text fragment 26751/42019
## 2023-09-16 21:37:40 Annotating text fragment 26761/42019
## 2023-09-16 21:37:40 Annotating text fragment 26771/42019
## 2023-09-16 21:37:40 Annotating text fragment 26781/42019
## 2023-09-16 21:37:40 Annotating text fragment 26791/42019
## 2023-09-16 21:37:40 Annotating text fragment 26801/42019
## 2023-09-16 21:37:40 Annotating text fragment 26811/42019
## 2023-09-16 21:37:41 Annotating text fragment 26821/42019
## 2023-09-16 21:37:41 Annotating text fragment 26831/42019
## 2023-09-16 21:37:41 Annotating text fragment 26841/42019
## 2023-09-16 21:37:41 Annotating text fragment 26851/42019
## 2023-09-16 21:37:41 Annotating text fragment 26861/42019
## 2023-09-16 21:37:41 Annotating text fragment 26871/42019
## 2023-09-16 21:37:41 Annotating text fragment 26881/42019
## 2023-09-16 21:37:41 Annotating text fragment 26891/42019
## 2023-09-16 21:37:41 Annotating text fragment 26901/42019
## 2023-09-16 21:37:42 Annotating text fragment 26911/42019
## 2023-09-16 21:37:42 Annotating text fragment 26921/42019
## 2023-09-16 21:37:42 Annotating text fragment 26931/42019
## 2023-09-16 21:37:42 Annotating text fragment 26941/42019
## 2023-09-16 21:37:42 Annotating text fragment 26951/42019
## 2023-09-16 21:37:42 Annotating text fragment 26961/42019
## 2023-09-16 21:37:42 Annotating text fragment 26971/42019
## 2023-09-16 21:37:42 Annotating text fragment 26981/42019
## 2023-09-16 21:37:43 Annotating text fragment 26991/42019
## 2023-09-16 21:37:43 Annotating text fragment 27001/42019
## 2023-09-16 21:37:43 Annotating text fragment 27011/42019
## 2023-09-16 21:37:43 Annotating text fragment 27021/42019
## 2023-09-16 21:37:43 Annotating text fragment 27031/42019
## 2023-09-16 21:37:43 Annotating text fragment 27041/42019
## 2023-09-16 21:37:43 Annotating text fragment 27051/42019
## 2023-09-16 21:37:43 Annotating text fragment 27061/42019
## 2023-09-16 21:37:44 Annotating text fragment 27071/42019
## 2023-09-16 21:37:44 Annotating text fragment 27081/42019
## 2023-09-16 21:37:44 Annotating text fragment 27091/42019
## 2023-09-16 21:37:44 Annotating text fragment 27101/42019
## 2023-09-16 21:37:44 Annotating text fragment 27111/42019
## 2023-09-16 21:37:44 Annotating text fragment 27121/42019
## 2023-09-16 21:37:44 Annotating text fragment 27131/42019
## 2023-09-16 21:37:45 Annotating text fragment 27141/42019
```

```
## 2023-09-16 21:37:45 Annotating text fragment 27151/42019
## 2023-09-16 21:37:45 Annotating text fragment 27161/42019
## 2023-09-16 21:37:45 Annotating text fragment 27171/42019
## 2023-09-16 21:37:45 Annotating text fragment 27181/42019
## 2023-09-16 21:37:45 Annotating text fragment 27191/42019
## 2023-09-16 21:37:45 Annotating text fragment 27201/42019
## 2023-09-16 21:37:45 Annotating text fragment 27211/42019
## 2023-09-16 21:37:45 Annotating text fragment 27221/42019
## 2023-09-16 21:37:45 Annotating text fragment 27231/42019
## 2023-09-16 21:37:46 Annotating text fragment 27241/42019
## 2023-09-16 21:37:46 Annotating text fragment 27251/42019
## 2023-09-16 21:37:46 Annotating text fragment 27261/42019
## 2023-09-16 21:37:46 Annotating text fragment 27271/42019
## 2023-09-16 21:37:46 Annotating text fragment 27281/42019
## 2023-09-16 21:37:46 Annotating text fragment 27291/42019
## 2023-09-16 21:37:46 Annotating text fragment 27301/42019
## 2023-09-16 21:37:46 Annotating text fragment 27311/42019
## 2023-09-16 21:37:47 Annotating text fragment 27321/42019
## 2023-09-16 21:37:47 Annotating text fragment 27331/42019
## 2023-09-16 21:37:47 Annotating text fragment 27341/42019
## 2023-09-16 21:37:47 Annotating text fragment 27351/42019
## 2023-09-16 21:37:47 Annotating text fragment 27361/42019
## 2023-09-16 21:37:47 Annotating text fragment 27371/42019
## 2023-09-16 21:37:47 Annotating text fragment 27381/42019
## 2023-09-16 21:37:47 Annotating text fragment 27391/42019
## 2023-09-16 21:37:47 Annotating text fragment 27401/42019
## 2023-09-16 21:37:47 Annotating text fragment 27411/42019
## 2023-09-16 21:37:47 Annotating text fragment 27421/42019
## 2023-09-16 21:37:47 Annotating text fragment 27431/42019
## 2023-09-16 21:37:48 Annotating text fragment 27441/42019
## 2023-09-16 21:37:48 Annotating text fragment 27451/42019
## 2023-09-16 21:37:48 Annotating text fragment 27461/42019
## 2023-09-16 21:37:48 Annotating text fragment 27471/42019
## 2023-09-16 21:37:48 Annotating text fragment 27481/42019
## 2023-09-16 21:37:48 Annotating text fragment 27491/42019
## 2023-09-16 21:37:48 Annotating text fragment 27501/42019
## 2023-09-16 21:37:48 Annotating text fragment 27511/42019
## 2023-09-16 21:37:49 Annotating text fragment 27521/42019
## 2023-09-16 21:37:49 Annotating text fragment 27531/42019
## 2023-09-16 21:37:49 Annotating text fragment 27541/42019
## 2023-09-16 21:37:49 Annotating text fragment 27551/42019
## 2023-09-16 21:37:49 Annotating text fragment 27561/42019
## 2023-09-16 21:37:49 Annotating text fragment 27571/42019
## 2023-09-16 21:37:49 Annotating text fragment 27581/42019
## 2023-09-16 21:37:49 Annotating text fragment 27591/42019
## 2023-09-16 21:37:49 Annotating text fragment 27601/42019
## 2023-09-16 21:37:49 Annotating text fragment 27611/42019
## 2023-09-16 21:37:50 Annotating text fragment 27621/42019
## 2023-09-16 21:37:50 Annotating text fragment 27631/42019
## 2023-09-16 21:37:50 Annotating text fragment 27641/42019
## 2023-09-16 21:37:50 Annotating text fragment 27651/42019
## 2023-09-16 21:37:50 Annotating text fragment 27661/42019
## 2023-09-16 21:37:50 Annotating text fragment 27671/42019
## 2023-09-16 21:37:50 Annotating text fragment 27681/42019
```

```
## 2023-09-16 21:37:50 Annotating text fragment 27691/42019
## 2023-09-16 21:37:50 Annotating text fragment 27701/42019
## 2023-09-16 21:37:50 Annotating text fragment 27711/42019
## 2023-09-16 21:37:51 Annotating text fragment 27721/42019
## 2023-09-16 21:37:51 Annotating text fragment 27731/42019
## 2023-09-16 21:37:51 Annotating text fragment 27741/42019
## 2023-09-16 21:37:51 Annotating text fragment 27751/42019
## 2023-09-16 21:37:51 Annotating text fragment 27761/42019
## 2023-09-16 21:37:51 Annotating text fragment 27771/42019
## 2023-09-16 21:37:51 Annotating text fragment 27781/42019
## 2023-09-16 21:37:51 Annotating text fragment 27791/42019
## 2023-09-16 21:37:51 Annotating text fragment 27801/42019
## 2023-09-16 21:37:52 Annotating text fragment 27811/42019
## 2023-09-16 21:37:52 Annotating text fragment 27821/42019
## 2023-09-16 21:37:52 Annotating text fragment 27831/42019
## 2023-09-16 21:37:52 Annotating text fragment 27841/42019
## 2023-09-16 21:37:52 Annotating text fragment 27851/42019
## 2023-09-16 21:37:52 Annotating text fragment 27861/42019
## 2023-09-16 21:37:52 Annotating text fragment 27871/42019
## 2023-09-16 21:37:53 Annotating text fragment 27881/42019
## 2023-09-16 21:37:53 Annotating text fragment 27891/42019
## 2023-09-16 21:37:53 Annotating text fragment 27901/42019
## 2023-09-16 21:37:53 Annotating text fragment 27911/42019
## 2023-09-16 21:37:53 Annotating text fragment 27921/42019
## 2023-09-16 21:37:53 Annotating text fragment 27931/42019
## 2023-09-16 21:37:53 Annotating text fragment 27941/42019
## 2023-09-16 21:37:53 Annotating text fragment 27951/42019
## 2023-09-16 21:37:53 Annotating text fragment 27961/42019
## 2023-09-16 21:37:53 Annotating text fragment 27971/42019
## 2023-09-16 21:37:53 Annotating text fragment 27981/42019
## 2023-09-16 21:37:54 Annotating text fragment 27991/42019
## 2023-09-16 21:37:54 Annotating text fragment 28001/42019
## 2023-09-16 21:37:54 Annotating text fragment 28011/42019
## 2023-09-16 21:37:54 Annotating text fragment 28021/42019
## 2023-09-16 21:37:54 Annotating text fragment 28031/42019
## 2023-09-16 21:37:54 Annotating text fragment 28041/42019
## 2023-09-16 21:37:54 Annotating text fragment 28051/42019
## 2023-09-16 21:37:55 Annotating text fragment 28061/42019
## 2023-09-16 21:37:55 Annotating text fragment 28071/42019
## 2023-09-16 21:37:55 Annotating text fragment 28081/42019
## 2023-09-16 21:37:55 Annotating text fragment 28091/42019
## 2023-09-16 21:37:55 Annotating text fragment 28101/42019
## 2023-09-16 21:37:55 Annotating text fragment 28111/42019
## 2023-09-16 21:37:55 Annotating text fragment 28121/42019
## 2023-09-16 21:37:55 Annotating text fragment 28131/42019
## 2023-09-16 21:37:55 Annotating text fragment 28141/42019
## 2023-09-16 21:37:55 Annotating text fragment 28151/42019
## 2023-09-16 21:37:56 Annotating text fragment 28161/42019
## 2023-09-16 21:37:56 Annotating text fragment 28171/42019
## 2023-09-16 21:37:56 Annotating text fragment 28181/42019
## 2023-09-16 21:37:56 Annotating text fragment 28191/42019
## 2023-09-16 21:37:56 Annotating text fragment 28201/42019
## 2023-09-16 21:37:56 Annotating text fragment 28211/42019
## 2023-09-16 21:37:56 Annotating text fragment 28221/42019
```

```
## 2023-09-16 21:37:56 Annotating text fragment 28231/42019
## 2023-09-16 21:37:57 Annotating text fragment 28241/42019
## 2023-09-16 21:37:57 Annotating text fragment 28251/42019
## 2023-09-16 21:37:57 Annotating text fragment 28261/42019
## 2023-09-16 21:37:57 Annotating text fragment 28271/42019
## 2023-09-16 21:37:57 Annotating text fragment 28281/42019
## 2023-09-16 21:37:57 Annotating text fragment 28291/42019
## 2023-09-16 21:37:57 Annotating text fragment 28301/42019
## 2023-09-16 21:37:58 Annotating text fragment 28311/42019
## 2023-09-16 21:37:58 Annotating text fragment 28321/42019
## 2023-09-16 21:37:58 Annotating text fragment 28331/42019
## 2023-09-16 21:37:58 Annotating text fragment 28341/42019
## 2023-09-16 21:37:58 Annotating text fragment 28351/42019
## 2023-09-16 21:37:58 Annotating text fragment 28361/42019
## 2023-09-16 21:37:58 Annotating text fragment 28371/42019
## 2023-09-16 21:37:58 Annotating text fragment 28381/42019
## 2023-09-16 21:37:59 Annotating text fragment 28391/42019
## 2023-09-16 21:37:59 Annotating text fragment 28401/42019
## 2023-09-16 21:37:59 Annotating text fragment 28411/42019
## 2023-09-16 21:37:59 Annotating text fragment 28421/42019
## 2023-09-16 21:37:59 Annotating text fragment 28431/42019
## 2023-09-16 21:37:59 Annotating text fragment 28441/42019
## 2023-09-16 21:37:59 Annotating text fragment 28451/42019
## 2023-09-16 21:37:59 Annotating text fragment 28461/42019
## 2023-09-16 21:38:00 Annotating text fragment 28471/42019
## 2023-09-16 21:38:00 Annotating text fragment 28481/42019
## 2023-09-16 21:38:00 Annotating text fragment 28491/42019
## 2023-09-16 21:38:00 Annotating text fragment 28501/42019
## 2023-09-16 21:38:00 Annotating text fragment 28511/42019
## 2023-09-16 21:38:00 Annotating text fragment 28521/42019
## 2023-09-16 21:38:00 Annotating text fragment 28531/42019
## 2023-09-16 21:38:00 Annotating text fragment 28541/42019
## 2023-09-16 21:38:00 Annotating text fragment 28551/42019
## 2023-09-16 21:38:00 Annotating text fragment 28561/42019
## 2023-09-16 21:38:01 Annotating text fragment 28571/42019
## 2023-09-16 21:38:01 Annotating text fragment 28581/42019
## 2023-09-16 21:38:01 Annotating text fragment 28591/42019
## 2023-09-16 21:38:01 Annotating text fragment 28601/42019
## 2023-09-16 21:38:01 Annotating text fragment 28611/42019
## 2023-09-16 21:38:01 Annotating text fragment 28621/42019
## 2023-09-16 21:38:01 Annotating text fragment 28631/42019
## 2023-09-16 21:38:01 Annotating text fragment 28641/42019
## 2023-09-16 21:38:02 Annotating text fragment 28651/42019
## 2023-09-16 21:38:02 Annotating text fragment 28661/42019
## 2023-09-16 21:38:02 Annotating text fragment 28671/42019
## 2023-09-16 21:38:02 Annotating text fragment 28681/42019
## 2023-09-16 21:38:02 Annotating text fragment 28691/42019
## 2023-09-16 21:38:02 Annotating text fragment 28701/42019
## 2023-09-16 21:38:02 Annotating text fragment 28711/42019
## 2023-09-16 21:38:02 Annotating text fragment 28721/42019
## 2023-09-16 21:38:02 Annotating text fragment 28731/42019
## 2023-09-16 21:38:03 Annotating text fragment 28741/42019
## 2023-09-16 21:38:03 Annotating text fragment 28751/42019
## 2023-09-16 21:38:03 Annotating text fragment 28761/42019
```

```
## 2023-09-16 21:38:03 Annotating text fragment 28771/42019
## 2023-09-16 21:38:03 Annotating text fragment 28781/42019
## 2023-09-16 21:38:03 Annotating text fragment 28791/42019
## 2023-09-16 21:38:04 Annotating text fragment 28801/42019
## 2023-09-16 21:38:04 Annotating text fragment 28811/42019
## 2023-09-16 21:38:04 Annotating text fragment 28821/42019
## 2023-09-16 21:38:04 Annotating text fragment 28831/42019
## 2023-09-16 21:38:04 Annotating text fragment 28841/42019
## 2023-09-16 21:38:04 Annotating text fragment 28851/42019
## 2023-09-16 21:38:04 Annotating text fragment 28861/42019
## 2023-09-16 21:38:04 Annotating text fragment 28871/42019
## 2023-09-16 21:38:04 Annotating text fragment 28881/42019
## 2023-09-16 21:38:04 Annotating text fragment 28891/42019
## 2023-09-16 21:38:05 Annotating text fragment 28901/42019
## 2023-09-16 21:38:05 Annotating text fragment 28911/42019
## 2023-09-16 21:38:05 Annotating text fragment 28921/42019
## 2023-09-16 21:38:05 Annotating text fragment 28931/42019
## 2023-09-16 21:38:05 Annotating text fragment 28941/42019
## 2023-09-16 21:38:05 Annotating text fragment 28951/42019
## 2023-09-16 21:38:05 Annotating text fragment 28961/42019
## 2023-09-16 21:38:05 Annotating text fragment 28971/42019
## 2023-09-16 21:38:05 Annotating text fragment 28981/42019
## 2023-09-16 21:38:05 Annotating text fragment 28991/42019
## 2023-09-16 21:38:06 Annotating text fragment 29001/42019
## 2023-09-16 21:38:06 Annotating text fragment 29011/42019
## 2023-09-16 21:38:06 Annotating text fragment 29021/42019
## 2023-09-16 21:38:06 Annotating text fragment 29031/42019
## 2023-09-16 21:38:06 Annotating text fragment 29041/42019
## 2023-09-16 21:38:06 Annotating text fragment 29051/42019
## 2023-09-16 21:38:06 Annotating text fragment 29061/42019
## 2023-09-16 21:38:07 Annotating text fragment 29071/42019
## 2023-09-16 21:38:07 Annotating text fragment 29081/42019
## 2023-09-16 21:38:07 Annotating text fragment 29091/42019
## 2023-09-16 21:38:07 Annotating text fragment 29101/42019
## 2023-09-16 21:38:07 Annotating text fragment 29111/42019
## 2023-09-16 21:38:07 Annotating text fragment 29121/42019
## 2023-09-16 21:38:07 Annotating text fragment 29131/42019
## 2023-09-16 21:38:07 Annotating text fragment 29141/42019
## 2023-09-16 21:38:08 Annotating text fragment 29151/42019
## 2023-09-16 21:38:08 Annotating text fragment 29161/42019
## 2023-09-16 21:38:08 Annotating text fragment 29171/42019
## 2023-09-16 21:38:08 Annotating text fragment 29181/42019
## 2023-09-16 21:38:08 Annotating text fragment 29191/42019
## 2023-09-16 21:38:08 Annotating text fragment 29201/42019
## 2023-09-16 21:38:08 Annotating text fragment 29211/42019
## 2023-09-16 21:38:09 Annotating text fragment 29221/42019
## 2023-09-16 21:38:09 Annotating text fragment 29231/42019
## 2023-09-16 21:38:09 Annotating text fragment 29241/42019
## 2023-09-16 21:38:09 Annotating text fragment 29251/42019
## 2023-09-16 21:38:09 Annotating text fragment 29261/42019
## 2023-09-16 21:38:09 Annotating text fragment 29271/42019
## 2023-09-16 21:38:09 Annotating text fragment 29281/42019
## 2023-09-16 21:38:09 Annotating text fragment 29291/42019
## 2023-09-16 21:38:09 Annotating text fragment 29301/42019
```

```
## 2023-09-16 21:38:10 Annotating text fragment 29311/42019
## 2023-09-16 21:38:10 Annotating text fragment 29321/42019
## 2023-09-16 21:38:10 Annotating text fragment 29331/42019
## 2023-09-16 21:38:10 Annotating text fragment 29341/42019
## 2023-09-16 21:38:10 Annotating text fragment 29351/42019
## 2023-09-16 21:38:10 Annotating text fragment 29361/42019
## 2023-09-16 21:38:10 Annotating text fragment 29371/42019
## 2023-09-16 21:38:10 Annotating text fragment 29381/42019
## 2023-09-16 21:38:10 Annotating text fragment 29391/42019
## 2023-09-16 21:38:10 Annotating text fragment 29401/42019
## 2023-09-16 21:38:10 Annotating text fragment 29411/42019
## 2023-09-16 21:38:11 Annotating text fragment 29421/42019
## 2023-09-16 21:38:11 Annotating text fragment 29431/42019
## 2023-09-16 21:38:11 Annotating text fragment 29441/42019
## 2023-09-16 21:38:11 Annotating text fragment 29451/42019
## 2023-09-16 21:38:11 Annotating text fragment 29461/42019
## 2023-09-16 21:38:11 Annotating text fragment 29471/42019
## 2023-09-16 21:38:11 Annotating text fragment 29481/42019
## 2023-09-16 21:38:11 Annotating text fragment 29491/42019
## 2023-09-16 21:38:11 Annotating text fragment 29501/42019
## 2023-09-16 21:38:11 Annotating text fragment 29511/42019
## 2023-09-16 21:38:12 Annotating text fragment 29521/42019
## 2023-09-16 21:38:12 Annotating text fragment 29531/42019
## 2023-09-16 21:38:12 Annotating text fragment 29541/42019
## 2023-09-16 21:38:12 Annotating text fragment 29551/42019
## 2023-09-16 21:38:12 Annotating text fragment 29561/42019
## 2023-09-16 21:38:12 Annotating text fragment 29571/42019
## 2023-09-16 21:38:12 Annotating text fragment 29581/42019
## 2023-09-16 21:38:12 Annotating text fragment 29591/42019
## 2023-09-16 21:38:13 Annotating text fragment 29601/42019
## 2023-09-16 21:38:13 Annotating text fragment 29611/42019
## 2023-09-16 21:38:13 Annotating text fragment 29621/42019
## 2023-09-16 21:38:13 Annotating text fragment 29631/42019
## 2023-09-16 21:38:13 Annotating text fragment 29641/42019
## 2023-09-16 21:38:13 Annotating text fragment 29651/42019
## 2023-09-16 21:38:13 Annotating text fragment 29661/42019
## 2023-09-16 21:38:13 Annotating text fragment 29671/42019
## 2023-09-16 21:38:13 Annotating text fragment 29681/42019
## 2023-09-16 21:38:13 Annotating text fragment 29691/42019
## 2023-09-16 21:38:14 Annotating text fragment 29701/42019
## 2023-09-16 21:38:14 Annotating text fragment 29711/42019
## 2023-09-16 21:38:14 Annotating text fragment 29721/42019
## 2023-09-16 21:38:14 Annotating text fragment 29731/42019
## 2023-09-16 21:38:14 Annotating text fragment 29741/42019
## 2023-09-16 21:38:14 Annotating text fragment 29751/42019
## 2023-09-16 21:38:15 Annotating text fragment 29761/42019
## 2023-09-16 21:38:15 Annotating text fragment 29771/42019
## 2023-09-16 21:38:15 Annotating text fragment 29781/42019
## 2023-09-16 21:38:15 Annotating text fragment 29791/42019
## 2023-09-16 21:38:15 Annotating text fragment 29801/42019
## 2023-09-16 21:38:15 Annotating text fragment 29811/42019
## 2023-09-16 21:38:15 Annotating text fragment 29821/42019
## 2023-09-16 21:38:15 Annotating text fragment 29831/42019
## 2023-09-16 21:38:15 Annotating text fragment 29841/42019
```

```
## 2023-09-16 21:38:15 Annotating text fragment 29851/42019
## 2023-09-16 21:38:15 Annotating text fragment 29861/42019
## 2023-09-16 21:38:15 Annotating text fragment 29871/42019
## 2023-09-16 21:38:16 Annotating text fragment 29881/42019
## 2023-09-16 21:38:16 Annotating text fragment 29891/42019
## 2023-09-16 21:38:16 Annotating text fragment 29901/42019
## 2023-09-16 21:38:16 Annotating text fragment 29911/42019
## 2023-09-16 21:38:16 Annotating text fragment 29921/42019
## 2023-09-16 21:38:16 Annotating text fragment 29931/42019
## 2023-09-16 21:38:16 Annotating text fragment 29941/42019
## 2023-09-16 21:38:16 Annotating text fragment 29951/42019
## 2023-09-16 21:38:17 Annotating text fragment 29961/42019
## 2023-09-16 21:38:17 Annotating text fragment 29971/42019
## 2023-09-16 21:38:17 Annotating text fragment 29981/42019
## 2023-09-16 21:38:17 Annotating text fragment 29991/42019
## 2023-09-16 21:38:17 Annotating text fragment 30001/42019
## 2023-09-16 21:38:17 Annotating text fragment 30011/42019
## 2023-09-16 21:38:17 Annotating text fragment 30021/42019
## 2023-09-16 21:38:17 Annotating text fragment 30031/42019
## 2023-09-16 21:38:17 Annotating text fragment 30041/42019
## 2023-09-16 21:38:17 Annotating text fragment 30051/42019
## 2023-09-16 21:38:18 Annotating text fragment 30061/42019
## 2023-09-16 21:38:18 Annotating text fragment 30071/42019
## 2023-09-16 21:38:18 Annotating text fragment 30081/42019
## 2023-09-16 21:38:18 Annotating text fragment 30091/42019
## 2023-09-16 21:38:18 Annotating text fragment 30101/42019
## 2023-09-16 21:38:18 Annotating text fragment 30111/42019
## 2023-09-16 21:38:18 Annotating text fragment 30121/42019
## 2023-09-16 21:38:18 Annotating text fragment 30131/42019
## 2023-09-16 21:38:18 Annotating text fragment 30141/42019
## 2023-09-16 21:38:18 Annotating text fragment 30151/42019
## 2023-09-16 21:38:18 Annotating text fragment 30161/42019
## 2023-09-16 21:38:19 Annotating text fragment 30171/42019
## 2023-09-16 21:38:19 Annotating text fragment 30181/42019
## 2023-09-16 21:38:19 Annotating text fragment 30191/42019
## 2023-09-16 21:38:19 Annotating text fragment 30201/42019
## 2023-09-16 21:38:19 Annotating text fragment 30211/42019
## 2023-09-16 21:38:19 Annotating text fragment 30221/42019
## 2023-09-16 21:38:19 Annotating text fragment 30231/42019
## 2023-09-16 21:38:19 Annotating text fragment 30241/42019
## 2023-09-16 21:38:19 Annotating text fragment 30251/42019
## 2023-09-16 21:38:19 Annotating text fragment 30261/42019
## 2023-09-16 21:38:20 Annotating text fragment 30271/42019
## 2023-09-16 21:38:20 Annotating text fragment 30281/42019
## 2023-09-16 21:38:20 Annotating text fragment 30291/42019
## 2023-09-16 21:38:20 Annotating text fragment 30301/42019
## 2023-09-16 21:38:20 Annotating text fragment 30311/42019
## 2023-09-16 21:38:20 Annotating text fragment 30321/42019
## 2023-09-16 21:38:20 Annotating text fragment 30331/42019
## 2023-09-16 21:38:21 Annotating text fragment 30341/42019
## 2023-09-16 21:38:21 Annotating text fragment 30351/42019
## 2023-09-16 21:38:21 Annotating text fragment 30361/42019
## 2023-09-16 21:38:21 Annotating text fragment 30371/42019
## 2023-09-16 21:38:21 Annotating text fragment 30381/42019
```

```
## 2023-09-16 21:38:21 Annotating text fragment 30391/42019
## 2023-09-16 21:38:21 Annotating text fragment 30401/42019
## 2023-09-16 21:38:22 Annotating text fragment 30411/42019
## 2023-09-16 21:38:22 Annotating text fragment 30421/42019
## 2023-09-16 21:38:22 Annotating text fragment 30431/42019
## 2023-09-16 21:38:22 Annotating text fragment 30441/42019
## 2023-09-16 21:38:22 Annotating text fragment 30451/42019
## 2023-09-16 21:38:22 Annotating text fragment 30461/42019
## 2023-09-16 21:38:22 Annotating text fragment 30471/42019
## 2023-09-16 21:38:23 Annotating text fragment 30481/42019
## 2023-09-16 21:38:23 Annotating text fragment 30491/42019
## 2023-09-16 21:38:23 Annotating text fragment 30501/42019
## 2023-09-16 21:38:23 Annotating text fragment 30511/42019
## 2023-09-16 21:38:23 Annotating text fragment 30521/42019
## 2023-09-16 21:38:23 Annotating text fragment 30531/42019
## 2023-09-16 21:38:23 Annotating text fragment 30541/42019
## 2023-09-16 21:38:24 Annotating text fragment 30551/42019
## 2023-09-16 21:38:24 Annotating text fragment 30561/42019
## 2023-09-16 21:38:24 Annotating text fragment 30571/42019
## 2023-09-16 21:38:24 Annotating text fragment 30581/42019
## 2023-09-16 21:38:24 Annotating text fragment 30591/42019
## 2023-09-16 21:38:24 Annotating text fragment 30601/42019
## 2023-09-16 21:38:24 Annotating text fragment 30611/42019
## 2023-09-16 21:38:24 Annotating text fragment 30621/42019
## 2023-09-16 21:38:25 Annotating text fragment 30631/42019
## 2023-09-16 21:38:25 Annotating text fragment 30641/42019
## 2023-09-16 21:38:25 Annotating text fragment 30651/42019
## 2023-09-16 21:38:25 Annotating text fragment 30661/42019
## 2023-09-16 21:38:25 Annotating text fragment 30671/42019
## 2023-09-16 21:38:25 Annotating text fragment 30681/42019
## 2023-09-16 21:38:25 Annotating text fragment 30691/42019
## 2023-09-16 21:38:25 Annotating text fragment 30701/42019
## 2023-09-16 21:38:25 Annotating text fragment 30711/42019
## 2023-09-16 21:38:25 Annotating text fragment 30721/42019
## 2023-09-16 21:38:25 Annotating text fragment 30731/42019
## 2023-09-16 21:38:26 Annotating text fragment 30741/42019
## 2023-09-16 21:38:26 Annotating text fragment 30751/42019
## 2023-09-16 21:38:26 Annotating text fragment 30761/42019
## 2023-09-16 21:38:26 Annotating text fragment 30771/42019
## 2023-09-16 21:38:26 Annotating text fragment 30781/42019
## 2023-09-16 21:38:26 Annotating text fragment 30791/42019
## 2023-09-16 21:38:26 Annotating text fragment 30801/42019
## 2023-09-16 21:38:26 Annotating text fragment 30811/42019
## 2023-09-16 21:38:27 Annotating text fragment 30821/42019
## 2023-09-16 21:38:27 Annotating text fragment 30831/42019
## 2023-09-16 21:38:27 Annotating text fragment 30841/42019
## 2023-09-16 21:38:27 Annotating text fragment 30851/42019
## 2023-09-16 21:38:27 Annotating text fragment 30861/42019
## 2023-09-16 21:38:27 Annotating text fragment 30871/42019
## 2023-09-16 21:38:27 Annotating text fragment 30881/42019
## 2023-09-16 21:38:27 Annotating text fragment 30891/42019
## 2023-09-16 21:38:28 Annotating text fragment 30901/42019
## 2023-09-16 21:38:28 Annotating text fragment 30911/42019
## 2023-09-16 21:38:28 Annotating text fragment 30921/42019
```

```
## 2023-09-16 21:38:28 Annotating text fragment 30931/42019
## 2023-09-16 21:38:28 Annotating text fragment 30941/42019
## 2023-09-16 21:38:28 Annotating text fragment 30951/42019
## 2023-09-16 21:38:28 Annotating text fragment 30961/42019
## 2023-09-16 21:38:28 Annotating text fragment 30971/42019
## 2023-09-16 21:38:28 Annotating text fragment 30981/42019
## 2023-09-16 21:38:29 Annotating text fragment 30991/42019
## 2023-09-16 21:38:29 Annotating text fragment 31001/42019
## 2023-09-16 21:38:29 Annotating text fragment 31011/42019
## 2023-09-16 21:38:29 Annotating text fragment 31021/42019
## 2023-09-16 21:38:29 Annotating text fragment 31031/42019
## 2023-09-16 21:38:29 Annotating text fragment 31041/42019
## 2023-09-16 21:38:29 Annotating text fragment 31051/42019
## 2023-09-16 21:38:30 Annotating text fragment 31061/42019
## 2023-09-16 21:38:30 Annotating text fragment 31071/42019
## 2023-09-16 21:38:30 Annotating text fragment 31081/42019
## 2023-09-16 21:38:30 Annotating text fragment 31091/42019
## 2023-09-16 21:38:30 Annotating text fragment 31101/42019
## 2023-09-16 21:38:30 Annotating text fragment 31111/42019
## 2023-09-16 21:38:30 Annotating text fragment 31121/42019
## 2023-09-16 21:38:30 Annotating text fragment 31131/42019
## 2023-09-16 21:38:30 Annotating text fragment 31141/42019
## 2023-09-16 21:38:31 Annotating text fragment 31151/42019
## 2023-09-16 21:38:31 Annotating text fragment 31161/42019
## 2023-09-16 21:38:31 Annotating text fragment 31171/42019
## 2023-09-16 21:38:31 Annotating text fragment 31181/42019
## 2023-09-16 21:38:31 Annotating text fragment 31191/42019
## 2023-09-16 21:38:31 Annotating text fragment 31201/42019
## 2023-09-16 21:38:32 Annotating text fragment 31211/42019
## 2023-09-16 21:38:32 Annotating text fragment 31221/42019
## 2023-09-16 21:38:32 Annotating text fragment 31231/42019
## 2023-09-16 21:38:32 Annotating text fragment 31241/42019
## 2023-09-16 21:38:32 Annotating text fragment 31251/42019
## 2023-09-16 21:38:32 Annotating text fragment 31261/42019
## 2023-09-16 21:38:32 Annotating text fragment 31271/42019
## 2023-09-16 21:38:33 Annotating text fragment 31281/42019
## 2023-09-16 21:38:33 Annotating text fragment 31291/42019
## 2023-09-16 21:38:33 Annotating text fragment 31301/42019
## 2023-09-16 21:38:33 Annotating text fragment 31311/42019
## 2023-09-16 21:38:33 Annotating text fragment 31321/42019
## 2023-09-16 21:38:33 Annotating text fragment 31331/42019
## 2023-09-16 21:38:33 Annotating text fragment 31341/42019
## 2023-09-16 21:38:33 Annotating text fragment 31351/42019
## 2023-09-16 21:38:33 Annotating text fragment 31361/42019
## 2023-09-16 21:38:33 Annotating text fragment 31371/42019
## 2023-09-16 21:38:34 Annotating text fragment 31381/42019
## 2023-09-16 21:38:34 Annotating text fragment 31391/42019
## 2023-09-16 21:38:34 Annotating text fragment 31401/42019
## 2023-09-16 21:38:34 Annotating text fragment 31411/42019
## 2023-09-16 21:38:34 Annotating text fragment 31421/42019
## 2023-09-16 21:38:34 Annotating text fragment 31431/42019
## 2023-09-16 21:38:34 Annotating text fragment 31441/42019
## 2023-09-16 21:38:34 Annotating text fragment 31451/42019
## 2023-09-16 21:38:35 Annotating text fragment 31461/42019
```

```
## 2023-09-16 21:38:35 Annotating text fragment 31471/42019
## 2023-09-16 21:38:35 Annotating text fragment 31481/42019
## 2023-09-16 21:38:35 Annotating text fragment 31491/42019
## 2023-09-16 21:38:35 Annotating text fragment 31501/42019
## 2023-09-16 21:38:35 Annotating text fragment 31511/42019
## 2023-09-16 21:38:35 Annotating text fragment 31521/42019
## 2023-09-16 21:38:36 Annotating text fragment 31531/42019
## 2023-09-16 21:38:36 Annotating text fragment 31541/42019
## 2023-09-16 21:38:36 Annotating text fragment 31551/42019
## 2023-09-16 21:38:36 Annotating text fragment 31561/42019
## 2023-09-16 21:38:36 Annotating text fragment 31571/42019
## 2023-09-16 21:38:36 Annotating text fragment 31581/42019
## 2023-09-16 21:38:36 Annotating text fragment 31591/42019
## 2023-09-16 21:38:36 Annotating text fragment 31601/42019
## 2023-09-16 21:38:36 Annotating text fragment 31611/42019
## 2023-09-16 21:38:37 Annotating text fragment 31621/42019
## 2023-09-16 21:38:37 Annotating text fragment 31631/42019
## 2023-09-16 21:38:37 Annotating text fragment 31641/42019
## 2023-09-16 21:38:37 Annotating text fragment 31651/42019
## 2023-09-16 21:38:37 Annotating text fragment 31661/42019
## 2023-09-16 21:38:37 Annotating text fragment 31671/42019
## 2023-09-16 21:38:37 Annotating text fragment 31681/42019
## 2023-09-16 21:38:38 Annotating text fragment 31691/42019
## 2023-09-16 21:38:38 Annotating text fragment 31701/42019
## 2023-09-16 21:38:38 Annotating text fragment 31711/42019
## 2023-09-16 21:38:38 Annotating text fragment 31721/42019
## 2023-09-16 21:38:38 Annotating text fragment 31731/42019
## 2023-09-16 21:38:38 Annotating text fragment 31741/42019
## 2023-09-16 21:38:38 Annotating text fragment 31751/42019
## 2023-09-16 21:38:39 Annotating text fragment 31761/42019
## 2023-09-16 21:38:39 Annotating text fragment 31771/42019
## 2023-09-16 21:38:39 Annotating text fragment 31781/42019
## 2023-09-16 21:38:39 Annotating text fragment 31791/42019
## 2023-09-16 21:38:39 Annotating text fragment 31801/42019
## 2023-09-16 21:38:39 Annotating text fragment 31811/42019
## 2023-09-16 21:38:39 Annotating text fragment 31821/42019
## 2023-09-16 21:38:39 Annotating text fragment 31831/42019
## 2023-09-16 21:38:39 Annotating text fragment 31841/42019
## 2023-09-16 21:38:40 Annotating text fragment 31851/42019
## 2023-09-16 21:38:40 Annotating text fragment 31861/42019
## 2023-09-16 21:38:40 Annotating text fragment 31871/42019
## 2023-09-16 21:38:40 Annotating text fragment 31881/42019
## 2023-09-16 21:38:40 Annotating text fragment 31891/42019
## 2023-09-16 21:38:40 Annotating text fragment 31901/42019
## 2023-09-16 21:38:41 Annotating text fragment 31911/42019
## 2023-09-16 21:38:41 Annotating text fragment 31921/42019
## 2023-09-16 21:38:41 Annotating text fragment 31931/42019
## 2023-09-16 21:38:41 Annotating text fragment 31941/42019
## 2023-09-16 21:38:41 Annotating text fragment 31951/42019
## 2023-09-16 21:38:41 Annotating text fragment 31961/42019
## 2023-09-16 21:38:41 Annotating text fragment 31971/42019
## 2023-09-16 21:38:42 Annotating text fragment 31981/42019
## 2023-09-16 21:38:42 Annotating text fragment 31991/42019
## 2023-09-16 21:38:42 Annotating text fragment 32001/42019
```

```
## 2023-09-16 21:38:42 Annotating text fragment 32011/42019
## 2023-09-16 21:38:42 Annotating text fragment 32021/42019
## 2023-09-16 21:38:42 Annotating text fragment 32031/42019
## 2023-09-16 21:38:42 Annotating text fragment 32041/42019
## 2023-09-16 21:38:42 Annotating text fragment 32051/42019
## 2023-09-16 21:38:42 Annotating text fragment 32061/42019
## 2023-09-16 21:38:42 Annotating text fragment 32071/42019
## 2023-09-16 21:38:43 Annotating text fragment 32081/42019
## 2023-09-16 21:38:43 Annotating text fragment 32091/42019
## 2023-09-16 21:38:43 Annotating text fragment 32101/42019
## 2023-09-16 21:38:43 Annotating text fragment 32111/42019
## 2023-09-16 21:38:43 Annotating text fragment 32121/42019
## 2023-09-16 21:38:43 Annotating text fragment 32131/42019
## 2023-09-16 21:38:43 Annotating text fragment 32141/42019
## 2023-09-16 21:38:43 Annotating text fragment 32151/42019
## 2023-09-16 21:38:43 Annotating text fragment 32161/42019
## 2023-09-16 21:38:43 Annotating text fragment 32171/42019
## 2023-09-16 21:38:44 Annotating text fragment 32181/42019
## 2023-09-16 21:38:44 Annotating text fragment 32191/42019
## 2023-09-16 21:38:44 Annotating text fragment 32201/42019
## 2023-09-16 21:38:44 Annotating text fragment 32211/42019
## 2023-09-16 21:38:44 Annotating text fragment 32221/42019
## 2023-09-16 21:38:44 Annotating text fragment 32231/42019
## 2023-09-16 21:38:44 Annotating text fragment 32241/42019
## 2023-09-16 21:38:44 Annotating text fragment 32251/42019
## 2023-09-16 21:38:44 Annotating text fragment 32261/42019
## 2023-09-16 21:38:44 Annotating text fragment 32271/42019
## 2023-09-16 21:38:45 Annotating text fragment 32281/42019
## 2023-09-16 21:38:45 Annotating text fragment 32291/42019
## 2023-09-16 21:38:45 Annotating text fragment 32301/42019
## 2023-09-16 21:38:45 Annotating text fragment 32311/42019
## 2023-09-16 21:38:45 Annotating text fragment 32321/42019
## 2023-09-16 21:38:45 Annotating text fragment 32331/42019
## 2023-09-16 21:38:45 Annotating text fragment 32341/42019
## 2023-09-16 21:38:45 Annotating text fragment 32351/42019
## 2023-09-16 21:38:46 Annotating text fragment 32361/42019
## 2023-09-16 21:38:46 Annotating text fragment 32371/42019
## 2023-09-16 21:38:46 Annotating text fragment 32381/42019
## 2023-09-16 21:38:46 Annotating text fragment 32391/42019
## 2023-09-16 21:38:46 Annotating text fragment 32401/42019
## 2023-09-16 21:38:46 Annotating text fragment 32411/42019
## 2023-09-16 21:38:46 Annotating text fragment 32421/42019
## 2023-09-16 21:38:46 Annotating text fragment 32431/42019
## 2023-09-16 21:38:47 Annotating text fragment 32441/42019
## 2023-09-16 21:38:47 Annotating text fragment 32451/42019
## 2023-09-16 21:38:47 Annotating text fragment 32461/42019
## 2023-09-16 21:38:47 Annotating text fragment 32471/42019
## 2023-09-16 21:38:47 Annotating text fragment 32481/42019
## 2023-09-16 21:38:47 Annotating text fragment 32491/42019
## 2023-09-16 21:38:47 Annotating text fragment 32501/42019
## 2023-09-16 21:38:47 Annotating text fragment 32511/42019
## 2023-09-16 21:38:47 Annotating text fragment 32521/42019
## 2023-09-16 21:38:48 Annotating text fragment 32531/42019
## 2023-09-16 21:38:48 Annotating text fragment 32541/42019
```

```
## 2023-09-16 21:38:48 Annotating text fragment 32551/42019
## 2023-09-16 21:38:48 Annotating text fragment 32561/42019
## 2023-09-16 21:38:48 Annotating text fragment 32571/42019
## 2023-09-16 21:38:48 Annotating text fragment 32581/42019
## 2023-09-16 21:38:49 Annotating text fragment 32591/42019
## 2023-09-16 21:38:49 Annotating text fragment 32601/42019
## 2023-09-16 21:38:49 Annotating text fragment 32611/42019
## 2023-09-16 21:38:49 Annotating text fragment 32621/42019
## 2023-09-16 21:38:49 Annotating text fragment 32631/42019
## 2023-09-16 21:38:49 Annotating text fragment 32641/42019
## 2023-09-16 21:38:49 Annotating text fragment 32651/42019
## 2023-09-16 21:38:49 Annotating text fragment 32661/42019
## 2023-09-16 21:38:50 Annotating text fragment 32671/42019
## 2023-09-16 21:38:50 Annotating text fragment 32681/42019
## 2023-09-16 21:38:50 Annotating text fragment 32691/42019
## 2023-09-16 21:38:50 Annotating text fragment 32701/42019
## 2023-09-16 21:38:50 Annotating text fragment 32711/42019
## 2023-09-16 21:38:50 Annotating text fragment 32721/42019
## 2023-09-16 21:38:51 Annotating text fragment 32731/42019
## 2023-09-16 21:38:51 Annotating text fragment 32741/42019
## 2023-09-16 21:38:51 Annotating text fragment 32751/42019
## 2023-09-16 21:38:51 Annotating text fragment 32761/42019
## 2023-09-16 21:38:51 Annotating text fragment 32771/42019
## 2023-09-16 21:38:51 Annotating text fragment 32781/42019
## 2023-09-16 21:38:51 Annotating text fragment 32791/42019
## 2023-09-16 21:38:52 Annotating text fragment 32801/42019
## 2023-09-16 21:38:52 Annotating text fragment 32811/42019
## 2023-09-16 21:38:52 Annotating text fragment 32821/42019
## 2023-09-16 21:38:52 Annotating text fragment 32831/42019
## 2023-09-16 21:38:52 Annotating text fragment 32841/42019
## 2023-09-16 21:38:52 Annotating text fragment 32851/42019
## 2023-09-16 21:38:52 Annotating text fragment 32861/42019
## 2023-09-16 21:38:53 Annotating text fragment 32871/42019
## 2023-09-16 21:38:53 Annotating text fragment 32881/42019
## 2023-09-16 21:38:53 Annotating text fragment 32891/42019
## 2023-09-16 21:38:53 Annotating text fragment 32901/42019
## 2023-09-16 21:38:53 Annotating text fragment 32911/42019
## 2023-09-16 21:38:53 Annotating text fragment 32921/42019
## 2023-09-16 21:38:53 Annotating text fragment 32931/42019
## 2023-09-16 21:38:53 Annotating text fragment 32941/42019
## 2023-09-16 21:38:53 Annotating text fragment 32951/42019
## 2023-09-16 21:38:54 Annotating text fragment 32961/42019
## 2023-09-16 21:38:54 Annotating text fragment 32971/42019
## 2023-09-16 21:38:54 Annotating text fragment 32981/42019
## 2023-09-16 21:38:54 Annotating text fragment 32991/42019
## 2023-09-16 21:38:54 Annotating text fragment 33001/42019
## 2023-09-16 21:38:54 Annotating text fragment 33011/42019
## 2023-09-16 21:38:54 Annotating text fragment 33021/42019
## 2023-09-16 21:38:54 Annotating text fragment 33031/42019
## 2023-09-16 21:38:55 Annotating text fragment 33041/42019
## 2023-09-16 21:38:55 Annotating text fragment 33051/42019
## 2023-09-16 21:38:55 Annotating text fragment 33061/42019
## 2023-09-16 21:38:55 Annotating text fragment 33071/42019
## 2023-09-16 21:38:55 Annotating text fragment 33081/42019
```

```
## 2023-09-16 21:38:55 Annotating text fragment 33091/42019
## 2023-09-16 21:38:55 Annotating text fragment 33101/42019
## 2023-09-16 21:38:55 Annotating text fragment 33111/42019
## 2023-09-16 21:38:56 Annotating text fragment 33121/42019
## 2023-09-16 21:38:56 Annotating text fragment 33131/42019
## 2023-09-16 21:38:56 Annotating text fragment 33141/42019
## 2023-09-16 21:38:56 Annotating text fragment 33151/42019
## 2023-09-16 21:38:56 Annotating text fragment 33161/42019
## 2023-09-16 21:38:56 Annotating text fragment 33171/42019
## 2023-09-16 21:38:56 Annotating text fragment 33181/42019
## 2023-09-16 21:38:56 Annotating text fragment 33191/42019
## 2023-09-16 21:38:56 Annotating text fragment 33201/42019
## 2023-09-16 21:38:56 Annotating text fragment 33211/42019
## 2023-09-16 21:38:57 Annotating text fragment 33221/42019
## 2023-09-16 21:38:57 Annotating text fragment 33231/42019
## 2023-09-16 21:38:57 Annotating text fragment 33241/42019
## 2023-09-16 21:38:57 Annotating text fragment 33251/42019
## 2023-09-16 21:38:57 Annotating text fragment 33261/42019
## 2023-09-16 21:38:57 Annotating text fragment 33271/42019
## 2023-09-16 21:38:57 Annotating text fragment 33281/42019
## 2023-09-16 21:38:58 Annotating text fragment 33291/42019
## 2023-09-16 21:38:58 Annotating text fragment 33301/42019
## 2023-09-16 21:38:58 Annotating text fragment 33311/42019
## 2023-09-16 21:38:58 Annotating text fragment 33321/42019
## 2023-09-16 21:38:58 Annotating text fragment 33331/42019
## 2023-09-16 21:38:58 Annotating text fragment 33341/42019
## 2023-09-16 21:38:58 Annotating text fragment 33351/42019
## 2023-09-16 21:38:58 Annotating text fragment 33361/42019
## 2023-09-16 21:38:58 Annotating text fragment 33371/42019
## 2023-09-16 21:38:58 Annotating text fragment 33381/42019
## 2023-09-16 21:38:59 Annotating text fragment 33391/42019
## 2023-09-16 21:38:59 Annotating text fragment 33401/42019
## 2023-09-16 21:38:59 Annotating text fragment 33411/42019
## 2023-09-16 21:38:59 Annotating text fragment 33421/42019
## 2023-09-16 21:38:59 Annotating text fragment 33431/42019
## 2023-09-16 21:39:00 Annotating text fragment 33441/42019
## 2023-09-16 21:39:00 Annotating text fragment 33451/42019
## 2023-09-16 21:39:00 Annotating text fragment 33461/42019
## 2023-09-16 21:39:00 Annotating text fragment 33471/42019
## 2023-09-16 21:39:00 Annotating text fragment 33481/42019
## 2023-09-16 21:39:00 Annotating text fragment 33491/42019
## 2023-09-16 21:39:00 Annotating text fragment 33501/42019
## 2023-09-16 21:39:00 Annotating text fragment 33511/42019
## 2023-09-16 21:39:00 Annotating text fragment 33521/42019
## 2023-09-16 21:39:01 Annotating text fragment 33531/42019
## 2023-09-16 21:39:01 Annotating text fragment 33541/42019
## 2023-09-16 21:39:01 Annotating text fragment 33551/42019
## 2023-09-16 21:39:01 Annotating text fragment 33561/42019
## 2023-09-16 21:39:01 Annotating text fragment 33571/42019
## 2023-09-16 21:39:01 Annotating text fragment 33581/42019
## 2023-09-16 21:39:02 Annotating text fragment 33591/42019
## 2023-09-16 21:39:02 Annotating text fragment 33601/42019
## 2023-09-16 21:39:02 Annotating text fragment 33611/42019
## 2023-09-16 21:39:02 Annotating text fragment 33621/42019
```

```
## 2023-09-16 21:39:02 Annotating text fragment 33631/42019
## 2023-09-16 21:39:02 Annotating text fragment 33641/42019
## 2023-09-16 21:39:02 Annotating text fragment 33651/42019
## 2023-09-16 21:39:03 Annotating text fragment 33661/42019
## 2023-09-16 21:39:03 Annotating text fragment 33671/42019
## 2023-09-16 21:39:03 Annotating text fragment 33681/42019
## 2023-09-16 21:39:03 Annotating text fragment 33691/42019
## 2023-09-16 21:39:03 Annotating text fragment 33701/42019
## 2023-09-16 21:39:03 Annotating text fragment 33711/42019
## 2023-09-16 21:39:03 Annotating text fragment 33721/42019
## 2023-09-16 21:39:03 Annotating text fragment 33731/42019
## 2023-09-16 21:39:04 Annotating text fragment 33741/42019
## 2023-09-16 21:39:04 Annotating text fragment 33751/42019
## 2023-09-16 21:39:04 Annotating text fragment 33761/42019
## 2023-09-16 21:39:04 Annotating text fragment 33771/42019
## 2023-09-16 21:39:04 Annotating text fragment 33781/42019
## 2023-09-16 21:39:04 Annotating text fragment 33791/42019
## 2023-09-16 21:39:04 Annotating text fragment 33801/42019
## 2023-09-16 21:39:04 Annotating text fragment 33811/42019
## 2023-09-16 21:39:04 Annotating text fragment 33821/42019
## 2023-09-16 21:39:05 Annotating text fragment 33831/42019
## 2023-09-16 21:39:05 Annotating text fragment 33841/42019
## 2023-09-16 21:39:05 Annotating text fragment 33851/42019
## 2023-09-16 21:39:05 Annotating text fragment 33861/42019
## 2023-09-16 21:39:05 Annotating text fragment 33871/42019
## 2023-09-16 21:39:05 Annotating text fragment 33881/42019
## 2023-09-16 21:39:05 Annotating text fragment 33891/42019
## 2023-09-16 21:39:05 Annotating text fragment 33901/42019
## 2023-09-16 21:39:06 Annotating text fragment 33911/42019
## 2023-09-16 21:39:06 Annotating text fragment 33921/42019
## 2023-09-16 21:39:06 Annotating text fragment 33931/42019
## 2023-09-16 21:39:06 Annotating text fragment 33941/42019
## 2023-09-16 21:39:06 Annotating text fragment 33951/42019
## 2023-09-16 21:39:06 Annotating text fragment 33961/42019
## 2023-09-16 21:39:06 Annotating text fragment 33971/42019
## 2023-09-16 21:39:06 Annotating text fragment 33981/42019
## 2023-09-16 21:39:07 Annotating text fragment 33991/42019
## 2023-09-16 21:39:07 Annotating text fragment 34001/42019
## 2023-09-16 21:39:07 Annotating text fragment 34011/42019
## 2023-09-16 21:39:07 Annotating text fragment 34021/42019
## 2023-09-16 21:39:07 Annotating text fragment 34031/42019
## 2023-09-16 21:39:07 Annotating text fragment 34041/42019
## 2023-09-16 21:39:07 Annotating text fragment 34051/42019
## 2023-09-16 21:39:08 Annotating text fragment 34061/42019
## 2023-09-16 21:39:08 Annotating text fragment 34071/42019
## 2023-09-16 21:39:08 Annotating text fragment 34081/42019
## 2023-09-16 21:39:08 Annotating text fragment 34091/42019
## 2023-09-16 21:39:08 Annotating text fragment 34101/42019
## 2023-09-16 21:39:08 Annotating text fragment 34111/42019
## 2023-09-16 21:39:08 Annotating text fragment 34121/42019
## 2023-09-16 21:39:09 Annotating text fragment 34131/42019
## 2023-09-16 21:39:09 Annotating text fragment 34141/42019
## 2023-09-16 21:39:09 Annotating text fragment 34151/42019
## 2023-09-16 21:39:09 Annotating text fragment 34161/42019
```

```
## 2023-09-16 21:39:09 Annotating text fragment 34171/42019
## 2023-09-16 21:39:09 Annotating text fragment 34181/42019
## 2023-09-16 21:39:09 Annotating text fragment 34191/42019
## 2023-09-16 21:39:09 Annotating text fragment 34201/42019
## 2023-09-16 21:39:09 Annotating text fragment 34211/42019
## 2023-09-16 21:39:10 Annotating text fragment 34221/42019
## 2023-09-16 21:39:10 Annotating text fragment 34231/42019
## 2023-09-16 21:39:10 Annotating text fragment 34241/42019
## 2023-09-16 21:39:10 Annotating text fragment 34251/42019
## 2023-09-16 21:39:10 Annotating text fragment 34261/42019
## 2023-09-16 21:39:10 Annotating text fragment 34271/42019
## 2023-09-16 21:39:10 Annotating text fragment 34281/42019
## 2023-09-16 21:39:10 Annotating text fragment 34291/42019
## 2023-09-16 21:39:10 Annotating text fragment 34301/42019
## 2023-09-16 21:39:10 Annotating text fragment 34311/42019
## 2023-09-16 21:39:11 Annotating text fragment 34321/42019
## 2023-09-16 21:39:11 Annotating text fragment 34331/42019
## 2023-09-16 21:39:11 Annotating text fragment 34341/42019
## 2023-09-16 21:39:11 Annotating text fragment 34351/42019
## 2023-09-16 21:39:11 Annotating text fragment 34361/42019
## 2023-09-16 21:39:11 Annotating text fragment 34371/42019
## 2023-09-16 21:39:11 Annotating text fragment 34381/42019
## 2023-09-16 21:39:12 Annotating text fragment 34391/42019
## 2023-09-16 21:39:12 Annotating text fragment 34401/42019
## 2023-09-16 21:39:12 Annotating text fragment 34411/42019
## 2023-09-16 21:39:12 Annotating text fragment 34421/42019
## 2023-09-16 21:39:12 Annotating text fragment 34431/42019
## 2023-09-16 21:39:12 Annotating text fragment 34441/42019
## 2023-09-16 21:39:12 Annotating text fragment 34451/42019
## 2023-09-16 21:39:13 Annotating text fragment 34461/42019
## 2023-09-16 21:39:13 Annotating text fragment 34471/42019
## 2023-09-16 21:39:13 Annotating text fragment 34481/42019
## 2023-09-16 21:39:13 Annotating text fragment 34491/42019
## 2023-09-16 21:39:13 Annotating text fragment 34501/42019
## 2023-09-16 21:39:13 Annotating text fragment 34511/42019
## 2023-09-16 21:39:13 Annotating text fragment 34521/42019
## 2023-09-16 21:39:14 Annotating text fragment 34531/42019
## 2023-09-16 21:39:14 Annotating text fragment 34541/42019
## 2023-09-16 21:39:14 Annotating text fragment 34551/42019
## 2023-09-16 21:39:14 Annotating text fragment 34561/42019
## 2023-09-16 21:39:14 Annotating text fragment 34571/42019
## 2023-09-16 21:39:15 Annotating text fragment 34581/42019
## 2023-09-16 21:39:15 Annotating text fragment 34591/42019
## 2023-09-16 21:39:15 Annotating text fragment 34601/42019
## 2023-09-16 21:39:15 Annotating text fragment 34611/42019
## 2023-09-16 21:39:15 Annotating text fragment 34621/42019
## 2023-09-16 21:39:15 Annotating text fragment 34631/42019
## 2023-09-16 21:39:15 Annotating text fragment 34641/42019
## 2023-09-16 21:39:16 Annotating text fragment 34651/42019
## 2023-09-16 21:39:16 Annotating text fragment 34661/42019
## 2023-09-16 21:39:16 Annotating text fragment 34671/42019
## 2023-09-16 21:39:16 Annotating text fragment 34681/42019
## 2023-09-16 21:39:16 Annotating text fragment 34691/42019
## 2023-09-16 21:39:16 Annotating text fragment 34701/42019
```

```
## 2023-09-16 21:39:16 Annotating text fragment 34711/42019
## 2023-09-16 21:39:16 Annotating text fragment 34721/42019
## 2023-09-16 21:39:16 Annotating text fragment 34731/42019
## 2023-09-16 21:39:16 Annotating text fragment 34741/42019
## 2023-09-16 21:39:17 Annotating text fragment 34751/42019
## 2023-09-16 21:39:17 Annotating text fragment 34761/42019
## 2023-09-16 21:39:17 Annotating text fragment 34771/42019
## 2023-09-16 21:39:17 Annotating text fragment 34781/42019
## 2023-09-16 21:39:17 Annotating text fragment 34791/42019
## 2023-09-16 21:39:18 Annotating text fragment 34801/42019
## 2023-09-16 21:39:18 Annotating text fragment 34811/42019
## 2023-09-16 21:39:18 Annotating text fragment 34821/42019
## 2023-09-16 21:39:18 Annotating text fragment 34831/42019
## 2023-09-16 21:39:18 Annotating text fragment 34841/42019
## 2023-09-16 21:39:18 Annotating text fragment 34851/42019
## 2023-09-16 21:39:18 Annotating text fragment 34861/42019
## 2023-09-16 21:39:18 Annotating text fragment 34871/42019
## 2023-09-16 21:39:18 Annotating text fragment 34881/42019
## 2023-09-16 21:39:19 Annotating text fragment 34891/42019
## 2023-09-16 21:39:19 Annotating text fragment 34901/42019
## 2023-09-16 21:39:19 Annotating text fragment 34911/42019
## 2023-09-16 21:39:19 Annotating text fragment 34921/42019
## 2023-09-16 21:39:19 Annotating text fragment 34931/42019
## 2023-09-16 21:39:19 Annotating text fragment 34941/42019
## 2023-09-16 21:39:20 Annotating text fragment 34951/42019
## 2023-09-16 21:39:20 Annotating text fragment 34961/42019
## 2023-09-16 21:39:20 Annotating text fragment 34971/42019
## 2023-09-16 21:39:20 Annotating text fragment 34981/42019
## 2023-09-16 21:39:20 Annotating text fragment 34991/42019
## 2023-09-16 21:39:20 Annotating text fragment 35001/42019
## 2023-09-16 21:39:20 Annotating text fragment 35011/42019
## 2023-09-16 21:39:20 Annotating text fragment 35021/42019
## 2023-09-16 21:39:21 Annotating text fragment 35031/42019
## 2023-09-16 21:39:21 Annotating text fragment 35041/42019
## 2023-09-16 21:39:21 Annotating text fragment 35051/42019
## 2023-09-16 21:39:21 Annotating text fragment 35061/42019
## 2023-09-16 21:39:21 Annotating text fragment 35071/42019
## 2023-09-16 21:39:21 Annotating text fragment 35081/42019
## 2023-09-16 21:39:21 Annotating text fragment 35091/42019
## 2023-09-16 21:39:21 Annotating text fragment 35101/42019
## 2023-09-16 21:39:21 Annotating text fragment 35111/42019
## 2023-09-16 21:39:22 Annotating text fragment 35121/42019
## 2023-09-16 21:39:22 Annotating text fragment 35131/42019
## 2023-09-16 21:39:22 Annotating text fragment 35141/42019
## 2023-09-16 21:39:22 Annotating text fragment 35151/42019
## 2023-09-16 21:39:22 Annotating text fragment 35161/42019
## 2023-09-16 21:39:22 Annotating text fragment 35171/42019
## 2023-09-16 21:39:22 Annotating text fragment 35181/42019
## 2023-09-16 21:39:22 Annotating text fragment 35191/42019
## 2023-09-16 21:39:22 Annotating text fragment 35201/42019
## 2023-09-16 21:39:23 Annotating text fragment 35211/42019
## 2023-09-16 21:39:23 Annotating text fragment 35221/42019
## 2023-09-16 21:39:23 Annotating text fragment 35231/42019
## 2023-09-16 21:39:23 Annotating text fragment 35241/42019
```

```
## 2023-09-16 21:39:23 Annotating text fragment 35251/42019
## 2023-09-16 21:39:23 Annotating text fragment 35261/42019
## 2023-09-16 21:39:23 Annotating text fragment 35271/42019
## 2023-09-16 21:39:23 Annotating text fragment 35281/42019
## 2023-09-16 21:39:24 Annotating text fragment 35291/42019
## 2023-09-16 21:39:24 Annotating text fragment 35301/42019
## 2023-09-16 21:39:24 Annotating text fragment 35311/42019
## 2023-09-16 21:39:24 Annotating text fragment 35321/42019
## 2023-09-16 21:39:24 Annotating text fragment 35331/42019
## 2023-09-16 21:39:24 Annotating text fragment 35341/42019
## 2023-09-16 21:39:24 Annotating text fragment 35351/42019
## 2023-09-16 21:39:24 Annotating text fragment 35361/42019
## 2023-09-16 21:39:25 Annotating text fragment 35371/42019
## 2023-09-16 21:39:25 Annotating text fragment 35381/42019
## 2023-09-16 21:39:25 Annotating text fragment 35391/42019
## 2023-09-16 21:39:25 Annotating text fragment 35401/42019
## 2023-09-16 21:39:25 Annotating text fragment 35411/42019
## 2023-09-16 21:39:25 Annotating text fragment 35421/42019
## 2023-09-16 21:39:25 Annotating text fragment 35431/42019
## 2023-09-16 21:39:25 Annotating text fragment 35441/42019
## 2023-09-16 21:39:26 Annotating text fragment 35451/42019
## 2023-09-16 21:39:26 Annotating text fragment 35461/42019
## 2023-09-16 21:39:26 Annotating text fragment 35471/42019
## 2023-09-16 21:39:26 Annotating text fragment 35481/42019
## 2023-09-16 21:39:26 Annotating text fragment 35491/42019
## 2023-09-16 21:39:26 Annotating text fragment 35501/42019
## 2023-09-16 21:39:26 Annotating text fragment 35511/42019
## 2023-09-16 21:39:27 Annotating text fragment 35521/42019
## 2023-09-16 21:39:27 Annotating text fragment 35531/42019
## 2023-09-16 21:39:27 Annotating text fragment 35541/42019
## 2023-09-16 21:39:27 Annotating text fragment 35551/42019
## 2023-09-16 21:39:27 Annotating text fragment 35561/42019
## 2023-09-16 21:39:27 Annotating text fragment 35571/42019
## 2023-09-16 21:39:27 Annotating text fragment 35581/42019
## 2023-09-16 21:39:27 Annotating text fragment 35591/42019
## 2023-09-16 21:39:27 Annotating text fragment 35601/42019
## 2023-09-16 21:39:27 Annotating text fragment 35611/42019
## 2023-09-16 21:39:28 Annotating text fragment 35621/42019
## 2023-09-16 21:39:28 Annotating text fragment 35631/42019
## 2023-09-16 21:39:28 Annotating text fragment 35641/42019
## 2023-09-16 21:39:28 Annotating text fragment 35651/42019
## 2023-09-16 21:39:28 Annotating text fragment 35661/42019
## 2023-09-16 21:39:28 Annotating text fragment 35671/42019
## 2023-09-16 21:39:28 Annotating text fragment 35681/42019
## 2023-09-16 21:39:28 Annotating text fragment 35691/42019
## 2023-09-16 21:39:29 Annotating text fragment 35701/42019
## 2023-09-16 21:39:29 Annotating text fragment 35711/42019
## 2023-09-16 21:39:29 Annotating text fragment 35721/42019
## 2023-09-16 21:39:29 Annotating text fragment 35731/42019
## 2023-09-16 21:39:29 Annotating text fragment 35741/42019
## 2023-09-16 21:39:29 Annotating text fragment 35751/42019
## 2023-09-16 21:39:29 Annotating text fragment 35761/42019
## 2023-09-16 21:39:29 Annotating text fragment 35771/42019
## 2023-09-16 21:39:30 Annotating text fragment 35781/42019
```

```
## 2023-09-16 21:39:30 Annotating text fragment 35791/42019
## 2023-09-16 21:39:30 Annotating text fragment 35801/42019
## 2023-09-16 21:39:30 Annotating text fragment 35811/42019
## 2023-09-16 21:39:30 Annotating text fragment 35821/42019
## 2023-09-16 21:39:30 Annotating text fragment 35831/42019
## 2023-09-16 21:39:30 Annotating text fragment 35841/42019
## 2023-09-16 21:39:30 Annotating text fragment 35851/42019
## 2023-09-16 21:39:31 Annotating text fragment 35861/42019
## 2023-09-16 21:39:31 Annotating text fragment 35871/42019
## 2023-09-16 21:39:31 Annotating text fragment 35881/42019
## 2023-09-16 21:39:31 Annotating text fragment 35891/42019
## 2023-09-16 21:39:31 Annotating text fragment 35901/42019
## 2023-09-16 21:39:32 Annotating text fragment 35911/42019
## 2023-09-16 21:39:32 Annotating text fragment 35921/42019
## 2023-09-16 21:39:32 Annotating text fragment 35931/42019
## 2023-09-16 21:39:32 Annotating text fragment 35941/42019
## 2023-09-16 21:39:32 Annotating text fragment 35951/42019
## 2023-09-16 21:39:32 Annotating text fragment 35961/42019
## 2023-09-16 21:39:32 Annotating text fragment 35971/42019
## 2023-09-16 21:39:32 Annotating text fragment 35981/42019
## 2023-09-16 21:39:33 Annotating text fragment 35991/42019
## 2023-09-16 21:39:33 Annotating text fragment 36001/42019
## 2023-09-16 21:39:33 Annotating text fragment 36011/42019
## 2023-09-16 21:39:33 Annotating text fragment 36021/42019
## 2023-09-16 21:39:33 Annotating text fragment 36031/42019
## 2023-09-16 21:39:33 Annotating text fragment 36041/42019
## 2023-09-16 21:39:33 Annotating text fragment 36051/42019
## 2023-09-16 21:39:34 Annotating text fragment 36061/42019
## 2023-09-16 21:39:34 Annotating text fragment 36071/42019
## 2023-09-16 21:39:34 Annotating text fragment 36081/42019
## 2023-09-16 21:39:34 Annotating text fragment 36091/42019
## 2023-09-16 21:39:34 Annotating text fragment 36101/42019
## 2023-09-16 21:39:34 Annotating text fragment 36111/42019
## 2023-09-16 21:39:34 Annotating text fragment 36121/42019
## 2023-09-16 21:39:34 Annotating text fragment 36131/42019
## 2023-09-16 21:39:34 Annotating text fragment 36141/42019
## 2023-09-16 21:39:34 Annotating text fragment 36151/42019
## 2023-09-16 21:39:35 Annotating text fragment 36161/42019
## 2023-09-16 21:39:35 Annotating text fragment 36171/42019
## 2023-09-16 21:39:35 Annotating text fragment 36181/42019
## 2023-09-16 21:39:35 Annotating text fragment 36191/42019
## 2023-09-16 21:39:35 Annotating text fragment 36201/42019
## 2023-09-16 21:39:35 Annotating text fragment 36211/42019
## 2023-09-16 21:39:35 Annotating text fragment 36221/42019
## 2023-09-16 21:39:35 Annotating text fragment 36231/42019
## 2023-09-16 21:39:36 Annotating text fragment 36241/42019
## 2023-09-16 21:39:36 Annotating text fragment 36251/42019
## 2023-09-16 21:39:36 Annotating text fragment 36261/42019
## 2023-09-16 21:39:36 Annotating text fragment 36271/42019
## 2023-09-16 21:39:36 Annotating text fragment 36281/42019
## 2023-09-16 21:39:36 Annotating text fragment 36291/42019
## 2023-09-16 21:39:36 Annotating text fragment 36301/42019
## 2023-09-16 21:39:36 Annotating text fragment 36311/42019
## 2023-09-16 21:39:37 Annotating text fragment 36321/42019
```

```
## 2023-09-16 21:39:37 Annotating text fragment 36331/42019
## 2023-09-16 21:39:37 Annotating text fragment 36341/42019
## 2023-09-16 21:39:37 Annotating text fragment 36351/42019
## 2023-09-16 21:39:37 Annotating text fragment 36361/42019
## 2023-09-16 21:39:37 Annotating text fragment 36371/42019
## 2023-09-16 21:39:37 Annotating text fragment 36381/42019
## 2023-09-16 21:39:37 Annotating text fragment 36391/42019
## 2023-09-16 21:39:37 Annotating text fragment 36401/42019
## 2023-09-16 21:39:37 Annotating text fragment 36411/42019
## 2023-09-16 21:39:38 Annotating text fragment 36421/42019
## 2023-09-16 21:39:38 Annotating text fragment 36431/42019
## 2023-09-16 21:39:38 Annotating text fragment 36441/42019
## 2023-09-16 21:39:38 Annotating text fragment 36451/42019
## 2023-09-16 21:39:38 Annotating text fragment 36461/42019
## 2023-09-16 21:39:38 Annotating text fragment 36471/42019
## 2023-09-16 21:39:38 Annotating text fragment 36481/42019
## 2023-09-16 21:39:38 Annotating text fragment 36491/42019
## 2023-09-16 21:39:38 Annotating text fragment 36501/42019
## 2023-09-16 21:39:39 Annotating text fragment 36511/42019
## 2023-09-16 21:39:39 Annotating text fragment 36521/42019
## 2023-09-16 21:39:39 Annotating text fragment 36531/42019
## 2023-09-16 21:39:39 Annotating text fragment 36541/42019
## 2023-09-16 21:39:39 Annotating text fragment 36551/42019
## 2023-09-16 21:39:39 Annotating text fragment 36561/42019
## 2023-09-16 21:39:39 Annotating text fragment 36571/42019
## 2023-09-16 21:39:39 Annotating text fragment 36581/42019
## 2023-09-16 21:39:39 Annotating text fragment 36591/42019
## 2023-09-16 21:39:40 Annotating text fragment 36601/42019
## 2023-09-16 21:39:40 Annotating text fragment 36611/42019
## 2023-09-16 21:39:40 Annotating text fragment 36621/42019
## 2023-09-16 21:39:40 Annotating text fragment 36631/42019
## 2023-09-16 21:39:40 Annotating text fragment 36641/42019
## 2023-09-16 21:39:40 Annotating text fragment 36651/42019
## 2023-09-16 21:39:40 Annotating text fragment 36661/42019
## 2023-09-16 21:39:40 Annotating text fragment 36671/42019
## 2023-09-16 21:39:40 Annotating text fragment 36681/42019
## 2023-09-16 21:39:41 Annotating text fragment 36691/42019
## 2023-09-16 21:39:41 Annotating text fragment 36701/42019
## 2023-09-16 21:39:41 Annotating text fragment 36711/42019
## 2023-09-16 21:39:41 Annotating text fragment 36721/42019
## 2023-09-16 21:39:41 Annotating text fragment 36731/42019
## 2023-09-16 21:39:41 Annotating text fragment 36741/42019
## 2023-09-16 21:39:41 Annotating text fragment 36751/42019
## 2023-09-16 21:39:41 Annotating text fragment 36761/42019
## 2023-09-16 21:39:42 Annotating text fragment 36771/42019
## 2023-09-16 21:39:42 Annotating text fragment 36781/42019
## 2023-09-16 21:39:42 Annotating text fragment 36791/42019
## 2023-09-16 21:39:42 Annotating text fragment 36801/42019
## 2023-09-16 21:39:42 Annotating text fragment 36811/42019
## 2023-09-16 21:39:42 Annotating text fragment 36821/42019
## 2023-09-16 21:39:42 Annotating text fragment 36831/42019
## 2023-09-16 21:39:43 Annotating text fragment 36841/42019
## 2023-09-16 21:39:43 Annotating text fragment 36851/42019
## 2023-09-16 21:39:43 Annotating text fragment 36861/42019
```

```
## 2023-09-16 21:39:43 Annotating text fragment 36871/42019
## 2023-09-16 21:39:43 Annotating text fragment 36881/42019
## 2023-09-16 21:39:43 Annotating text fragment 36891/42019
## 2023-09-16 21:39:43 Annotating text fragment 36901/42019
## 2023-09-16 21:39:43 Annotating text fragment 36911/42019
## 2023-09-16 21:39:43 Annotating text fragment 36921/42019
## 2023-09-16 21:39:44 Annotating text fragment 36931/42019
## 2023-09-16 21:39:44 Annotating text fragment 36941/42019
## 2023-09-16 21:39:44 Annotating text fragment 36951/42019
## 2023-09-16 21:39:44 Annotating text fragment 36961/42019
## 2023-09-16 21:39:44 Annotating text fragment 36971/42019
## 2023-09-16 21:39:44 Annotating text fragment 36981/42019
## 2023-09-16 21:39:44 Annotating text fragment 36991/42019
## 2023-09-16 21:39:44 Annotating text fragment 37001/42019
## 2023-09-16 21:39:45 Annotating text fragment 37011/42019
## 2023-09-16 21:39:45 Annotating text fragment 37021/42019
## 2023-09-16 21:39:45 Annotating text fragment 37031/42019
## 2023-09-16 21:39:45 Annotating text fragment 37041/42019
## 2023-09-16 21:39:45 Annotating text fragment 37051/42019
## 2023-09-16 21:39:45 Annotating text fragment 37061/42019
## 2023-09-16 21:39:45 Annotating text fragment 37071/42019
## 2023-09-16 21:39:45 Annotating text fragment 37081/42019
## 2023-09-16 21:39:45 Annotating text fragment 37091/42019
## 2023-09-16 21:39:45 Annotating text fragment 37101/42019
## 2023-09-16 21:39:46 Annotating text fragment 37111/42019
## 2023-09-16 21:39:46 Annotating text fragment 37121/42019
## 2023-09-16 21:39:46 Annotating text fragment 37131/42019
## 2023-09-16 21:39:46 Annotating text fragment 37141/42019
## 2023-09-16 21:39:46 Annotating text fragment 37151/42019
## 2023-09-16 21:39:46 Annotating text fragment 37161/42019
## 2023-09-16 21:39:46 Annotating text fragment 37171/42019
## 2023-09-16 21:39:46 Annotating text fragment 37181/42019
## 2023-09-16 21:39:46 Annotating text fragment 37191/42019
## 2023-09-16 21:39:47 Annotating text fragment 37201/42019
## 2023-09-16 21:39:47 Annotating text fragment 37211/42019
## 2023-09-16 21:39:47 Annotating text fragment 37221/42019
## 2023-09-16 21:39:47 Annotating text fragment 37231/42019
## 2023-09-16 21:39:47 Annotating text fragment 37241/42019
## 2023-09-16 21:39:47 Annotating text fragment 37251/42019
## 2023-09-16 21:39:47 Annotating text fragment 37261/42019
## 2023-09-16 21:39:47 Annotating text fragment 37271/42019
## 2023-09-16 21:39:47 Annotating text fragment 37281/42019
## 2023-09-16 21:39:47 Annotating text fragment 37291/42019
## 2023-09-16 21:39:48 Annotating text fragment 37301/42019
## 2023-09-16 21:39:48 Annotating text fragment 37311/42019
## 2023-09-16 21:39:48 Annotating text fragment 37321/42019
## 2023-09-16 21:39:48 Annotating text fragment 37331/42019
## 2023-09-16 21:39:48 Annotating text fragment 37341/42019
## 2023-09-16 21:39:48 Annotating text fragment 37351/42019
## 2023-09-16 21:39:48 Annotating text fragment 37361/42019
## 2023-09-16 21:39:48 Annotating text fragment 37371/42019
## 2023-09-16 21:39:48 Annotating text fragment 37381/42019
## 2023-09-16 21:39:48 Annotating text fragment 37391/42019
## 2023-09-16 21:39:49 Annotating text fragment 37401/42019
```

```
## 2023-09-16 21:39:49 Annotating text fragment 37411/42019
## 2023-09-16 21:39:49 Annotating text fragment 37421/42019
## 2023-09-16 21:39:49 Annotating text fragment 37431/42019
## 2023-09-16 21:39:49 Annotating text fragment 37441/42019
## 2023-09-16 21:39:49 Annotating text fragment 37451/42019
## 2023-09-16 21:39:49 Annotating text fragment 37461/42019
## 2023-09-16 21:39:49 Annotating text fragment 37471/42019
## 2023-09-16 21:39:50 Annotating text fragment 37481/42019
## 2023-09-16 21:39:50 Annotating text fragment 37491/42019
## 2023-09-16 21:39:50 Annotating text fragment 37501/42019
## 2023-09-16 21:39:50 Annotating text fragment 37511/42019
## 2023-09-16 21:39:50 Annotating text fragment 37521/42019
## 2023-09-16 21:39:50 Annotating text fragment 37531/42019
## 2023-09-16 21:39:50 Annotating text fragment 37541/42019
## 2023-09-16 21:39:50 Annotating text fragment 37551/42019
## 2023-09-16 21:39:50 Annotating text fragment 37561/42019
## 2023-09-16 21:39:51 Annotating text fragment 37571/42019
## 2023-09-16 21:39:51 Annotating text fragment 37581/42019
## 2023-09-16 21:39:51 Annotating text fragment 37591/42019
## 2023-09-16 21:39:51 Annotating text fragment 37601/42019
## 2023-09-16 21:39:51 Annotating text fragment 37611/42019
## 2023-09-16 21:39:51 Annotating text fragment 37621/42019
## 2023-09-16 21:39:51 Annotating text fragment 37631/42019
## 2023-09-16 21:39:51 Annotating text fragment 37641/42019
## 2023-09-16 21:39:52 Annotating text fragment 37651/42019
## 2023-09-16 21:39:52 Annotating text fragment 37661/42019
## 2023-09-16 21:39:52 Annotating text fragment 37671/42019
## 2023-09-16 21:39:52 Annotating text fragment 37681/42019
## 2023-09-16 21:39:52 Annotating text fragment 37691/42019
## 2023-09-16 21:39:52 Annotating text fragment 37701/42019
## 2023-09-16 21:39:52 Annotating text fragment 37711/42019
## 2023-09-16 21:39:53 Annotating text fragment 37721/42019
## 2023-09-16 21:39:53 Annotating text fragment 37731/42019
## 2023-09-16 21:39:53 Annotating text fragment 37741/42019
## 2023-09-16 21:39:53 Annotating text fragment 37751/42019
## 2023-09-16 21:39:53 Annotating text fragment 37761/42019
## 2023-09-16 21:39:53 Annotating text fragment 37771/42019
## 2023-09-16 21:39:54 Annotating text fragment 37781/42019
## 2023-09-16 21:39:54 Annotating text fragment 37791/42019
## 2023-09-16 21:39:54 Annotating text fragment 37801/42019
## 2023-09-16 21:39:54 Annotating text fragment 37811/42019
## 2023-09-16 21:39:54 Annotating text fragment 37821/42019
## 2023-09-16 21:39:54 Annotating text fragment 37831/42019
## 2023-09-16 21:39:54 Annotating text fragment 37841/42019
## 2023-09-16 21:39:55 Annotating text fragment 37851/42019
## 2023-09-16 21:39:55 Annotating text fragment 37861/42019
## 2023-09-16 21:39:55 Annotating text fragment 37871/42019
## 2023-09-16 21:39:55 Annotating text fragment 37881/42019
## 2023-09-16 21:39:55 Annotating text fragment 37891/42019
## 2023-09-16 21:39:55 Annotating text fragment 37901/42019
## 2023-09-16 21:39:55 Annotating text fragment 37911/42019
## 2023-09-16 21:39:55 Annotating text fragment 37921/42019
## 2023-09-16 21:39:55 Annotating text fragment 37931/42019
## 2023-09-16 21:39:55 Annotating text fragment 37941/42019
```

```
## 2023-09-16 21:39:56 Annotating text fragment 37951/42019
## 2023-09-16 21:39:56 Annotating text fragment 37961/42019
## 2023-09-16 21:39:56 Annotating text fragment 37971/42019
## 2023-09-16 21:39:56 Annotating text fragment 37981/42019
## 2023-09-16 21:39:56 Annotating text fragment 37991/42019
## 2023-09-16 21:39:57 Annotating text fragment 38001/42019
## 2023-09-16 21:39:57 Annotating text fragment 38011/42019
## 2023-09-16 21:39:57 Annotating text fragment 38021/42019
## 2023-09-16 21:39:57 Annotating text fragment 38031/42019
## 2023-09-16 21:39:57 Annotating text fragment 38041/42019
## 2023-09-16 21:39:57 Annotating text fragment 38051/42019
## 2023-09-16 21:39:57 Annotating text fragment 38061/42019
## 2023-09-16 21:39:57 Annotating text fragment 38071/42019
## 2023-09-16 21:39:57 Annotating text fragment 38081/42019
## 2023-09-16 21:39:58 Annotating text fragment 38091/42019
## 2023-09-16 21:39:58 Annotating text fragment 38101/42019
## 2023-09-16 21:39:58 Annotating text fragment 38111/42019
## 2023-09-16 21:39:58 Annotating text fragment 38121/42019
## 2023-09-16 21:39:58 Annotating text fragment 38131/42019
## 2023-09-16 21:39:58 Annotating text fragment 38141/42019
## 2023-09-16 21:39:58 Annotating text fragment 38151/42019
## 2023-09-16 21:39:58 Annotating text fragment 38161/42019
## 2023-09-16 21:39:59 Annotating text fragment 38171/42019
## 2023-09-16 21:39:59 Annotating text fragment 38181/42019
## 2023-09-16 21:39:59 Annotating text fragment 38191/42019
## 2023-09-16 21:39:59 Annotating text fragment 38201/42019
## 2023-09-16 21:39:59 Annotating text fragment 38211/42019
## 2023-09-16 21:39:59 Annotating text fragment 38221/42019
## 2023-09-16 21:40:00 Annotating text fragment 38231/42019
## 2023-09-16 21:40:00 Annotating text fragment 38241/42019
## 2023-09-16 21:40:00 Annotating text fragment 38251/42019
## 2023-09-16 21:40:00 Annotating text fragment 38261/42019
## 2023-09-16 21:40:00 Annotating text fragment 38271/42019
## 2023-09-16 21:40:00 Annotating text fragment 38281/42019
## 2023-09-16 21:40:00 Annotating text fragment 38291/42019
## 2023-09-16 21:40:00 Annotating text fragment 38301/42019
## 2023-09-16 21:40:01 Annotating text fragment 38311/42019
## 2023-09-16 21:40:01 Annotating text fragment 38321/42019
## 2023-09-16 21:40:01 Annotating text fragment 38331/42019
## 2023-09-16 21:40:01 Annotating text fragment 38341/42019
## 2023-09-16 21:40:01 Annotating text fragment 38351/42019
## 2023-09-16 21:40:01 Annotating text fragment 38361/42019
## 2023-09-16 21:40:01 Annotating text fragment 38371/42019
## 2023-09-16 21:40:01 Annotating text fragment 38381/42019
## 2023-09-16 21:40:01 Annotating text fragment 38391/42019
## 2023-09-16 21:40:02 Annotating text fragment 38401/42019
## 2023-09-16 21:40:02 Annotating text fragment 38411/42019
## 2023-09-16 21:40:02 Annotating text fragment 38421/42019
## 2023-09-16 21:40:02 Annotating text fragment 38431/42019
## 2023-09-16 21:40:02 Annotating text fragment 38441/42019
## 2023-09-16 21:40:02 Annotating text fragment 38451/42019
## 2023-09-16 21:40:02 Annotating text fragment 38461/42019
## 2023-09-16 21:40:02 Annotating text fragment 38471/42019
## 2023-09-16 21:40:02 Annotating text fragment 38481/42019
```

```
## 2023-09-16 21:40:03 Annotating text fragment 38491/42019
## 2023-09-16 21:40:03 Annotating text fragment 38501/42019
## 2023-09-16 21:40:03 Annotating text fragment 38511/42019
## 2023-09-16 21:40:03 Annotating text fragment 38521/42019
## 2023-09-16 21:40:03 Annotating text fragment 38531/42019
## 2023-09-16 21:40:03 Annotating text fragment 38541/42019
## 2023-09-16 21:40:03 Annotating text fragment 38551/42019
## 2023-09-16 21:40:04 Annotating text fragment 38561/42019
## 2023-09-16 21:40:04 Annotating text fragment 38571/42019
## 2023-09-16 21:40:04 Annotating text fragment 38581/42019
## 2023-09-16 21:40:04 Annotating text fragment 38591/42019
## 2023-09-16 21:40:04 Annotating text fragment 38601/42019
## 2023-09-16 21:40:04 Annotating text fragment 38611/42019
## 2023-09-16 21:40:04 Annotating text fragment 38621/42019
## 2023-09-16 21:40:04 Annotating text fragment 38631/42019
## 2023-09-16 21:40:04 Annotating text fragment 38641/42019
## 2023-09-16 21:40:05 Annotating text fragment 38651/42019
## 2023-09-16 21:40:05 Annotating text fragment 38661/42019
## 2023-09-16 21:40:05 Annotating text fragment 38671/42019
## 2023-09-16 21:40:05 Annotating text fragment 38681/42019
## 2023-09-16 21:40:05 Annotating text fragment 38691/42019
## 2023-09-16 21:40:05 Annotating text fragment 38701/42019
## 2023-09-16 21:40:06 Annotating text fragment 38711/42019
## 2023-09-16 21:40:06 Annotating text fragment 38721/42019
## 2023-09-16 21:40:06 Annotating text fragment 38731/42019
## 2023-09-16 21:40:06 Annotating text fragment 38741/42019
## 2023-09-16 21:40:06 Annotating text fragment 38751/42019
## 2023-09-16 21:40:06 Annotating text fragment 38761/42019
## 2023-09-16 21:40:06 Annotating text fragment 38771/42019
## 2023-09-16 21:40:06 Annotating text fragment 38781/42019
## 2023-09-16 21:40:06 Annotating text fragment 38791/42019
## 2023-09-16 21:40:07 Annotating text fragment 38801/42019
## 2023-09-16 21:40:07 Annotating text fragment 38811/42019
## 2023-09-16 21:40:07 Annotating text fragment 38821/42019
## 2023-09-16 21:40:07 Annotating text fragment 38831/42019
## 2023-09-16 21:40:07 Annotating text fragment 38841/42019
## 2023-09-16 21:40:07 Annotating text fragment 38851/42019
## 2023-09-16 21:40:07 Annotating text fragment 38861/42019
## 2023-09-16 21:40:07 Annotating text fragment 38871/42019
## 2023-09-16 21:40:07 Annotating text fragment 38881/42019
## 2023-09-16 21:40:08 Annotating text fragment 38891/42019
## 2023-09-16 21:40:08 Annotating text fragment 38901/42019
## 2023-09-16 21:40:08 Annotating text fragment 38911/42019
## 2023-09-16 21:40:08 Annotating text fragment 38921/42019
## 2023-09-16 21:40:08 Annotating text fragment 38931/42019
## 2023-09-16 21:40:08 Annotating text fragment 38941/42019
## 2023-09-16 21:40:08 Annotating text fragment 38951/42019
## 2023-09-16 21:40:08 Annotating text fragment 38961/42019
## 2023-09-16 21:40:08 Annotating text fragment 38971/42019
## 2023-09-16 21:40:08 Annotating text fragment 38981/42019
## 2023-09-16 21:40:09 Annotating text fragment 38991/42019
## 2023-09-16 21:40:09 Annotating text fragment 39001/42019
## 2023-09-16 21:40:09 Annotating text fragment 39011/42019
## 2023-09-16 21:40:09 Annotating text fragment 39021/42019
```

```
## 2023-09-16 21:40:09 Annotating text fragment 39031/42019
## 2023-09-16 21:40:09 Annotating text fragment 39041/42019
## 2023-09-16 21:40:09 Annotating text fragment 39051/42019
## 2023-09-16 21:40:09 Annotating text fragment 39061/42019
## 2023-09-16 21:40:10 Annotating text fragment 39071/42019
## 2023-09-16 21:40:10 Annotating text fragment 39081/42019
## 2023-09-16 21:40:10 Annotating text fragment 39091/42019
## 2023-09-16 21:40:10 Annotating text fragment 39101/42019
## 2023-09-16 21:40:10 Annotating text fragment 39111/42019
## 2023-09-16 21:40:10 Annotating text fragment 39121/42019
## 2023-09-16 21:40:10 Annotating text fragment 39131/42019
## 2023-09-16 21:40:10 Annotating text fragment 39141/42019
## 2023-09-16 21:40:11 Annotating text fragment 39151/42019
## 2023-09-16 21:40:11 Annotating text fragment 39161/42019
## 2023-09-16 21:40:11 Annotating text fragment 39171/42019
## 2023-09-16 21:40:11 Annotating text fragment 39181/42019
## 2023-09-16 21:40:11 Annotating text fragment 39191/42019
## 2023-09-16 21:40:11 Annotating text fragment 39201/42019
## 2023-09-16 21:40:11 Annotating text fragment 39211/42019
## 2023-09-16 21:40:12 Annotating text fragment 39221/42019
## 2023-09-16 21:40:12 Annotating text fragment 39231/42019
## 2023-09-16 21:40:12 Annotating text fragment 39241/42019
## 2023-09-16 21:40:12 Annotating text fragment 39251/42019
## 2023-09-16 21:40:12 Annotating text fragment 39261/42019
## 2023-09-16 21:40:12 Annotating text fragment 39271/42019
## 2023-09-16 21:40:13 Annotating text fragment 39281/42019
## 2023-09-16 21:40:13 Annotating text fragment 39291/42019
## 2023-09-16 21:40:13 Annotating text fragment 39301/42019
## 2023-09-16 21:40:13 Annotating text fragment 39311/42019
## 2023-09-16 21:40:13 Annotating text fragment 39321/42019
## 2023-09-16 21:40:13 Annotating text fragment 39331/42019
## 2023-09-16 21:40:13 Annotating text fragment 39341/42019
## 2023-09-16 21:40:13 Annotating text fragment 39351/42019
## 2023-09-16 21:40:13 Annotating text fragment 39361/42019
## 2023-09-16 21:40:13 Annotating text fragment 39371/42019
## 2023-09-16 21:40:14 Annotating text fragment 39381/42019
## 2023-09-16 21:40:14 Annotating text fragment 39391/42019
## 2023-09-16 21:40:14 Annotating text fragment 39401/42019
## 2023-09-16 21:40:14 Annotating text fragment 39411/42019
## 2023-09-16 21:40:14 Annotating text fragment 39421/42019
## 2023-09-16 21:40:14 Annotating text fragment 39431/42019
## 2023-09-16 21:40:14 Annotating text fragment 39441/42019
## 2023-09-16 21:40:14 Annotating text fragment 39451/42019
## 2023-09-16 21:40:14 Annotating text fragment 39461/42019
## 2023-09-16 21:40:14 Annotating text fragment 39471/42019
## 2023-09-16 21:40:15 Annotating text fragment 39481/42019
## 2023-09-16 21:40:15 Annotating text fragment 39491/42019
## 2023-09-16 21:40:15 Annotating text fragment 39501/42019
## 2023-09-16 21:40:15 Annotating text fragment 39511/42019
## 2023-09-16 21:40:15 Annotating text fragment 39521/42019
## 2023-09-16 21:40:15 Annotating text fragment 39531/42019
## 2023-09-16 21:40:15 Annotating text fragment 39541/42019
## 2023-09-16 21:40:15 Annotating text fragment 39551/42019
## 2023-09-16 21:40:16 Annotating text fragment 39561/42019
```

```
## 2023-09-16 21:40:16 Annotating text fragment 39571/42019
## 2023-09-16 21:40:16 Annotating text fragment 39581/42019
## 2023-09-16 21:40:16 Annotating text fragment 39591/42019
## 2023-09-16 21:40:16 Annotating text fragment 39601/42019
## 2023-09-16 21:40:16 Annotating text fragment 39611/42019
## 2023-09-16 21:40:16 Annotating text fragment 39621/42019
## 2023-09-16 21:40:17 Annotating text fragment 39631/42019
## 2023-09-16 21:40:17 Annotating text fragment 39641/42019
## 2023-09-16 21:40:17 Annotating text fragment 39651/42019
## 2023-09-16 21:40:17 Annotating text fragment 39661/42019
## 2023-09-16 21:40:17 Annotating text fragment 39671/42019
## 2023-09-16 21:40:17 Annotating text fragment 39681/42019
## 2023-09-16 21:40:17 Annotating text fragment 39691/42019
## 2023-09-16 21:40:17 Annotating text fragment 39701/42019
## 2023-09-16 21:40:17 Annotating text fragment 39711/42019
## 2023-09-16 21:40:17 Annotating text fragment 39721/42019
## 2023-09-16 21:40:18 Annotating text fragment 39731/42019
## 2023-09-16 21:40:18 Annotating text fragment 39741/42019
## 2023-09-16 21:40:18 Annotating text fragment 39751/42019
## 2023-09-16 21:40:18 Annotating text fragment 39761/42019
## 2023-09-16 21:40:18 Annotating text fragment 39771/42019
## 2023-09-16 21:40:18 Annotating text fragment 39781/42019
## 2023-09-16 21:40:19 Annotating text fragment 39791/42019
## 2023-09-16 21:40:19 Annotating text fragment 39801/42019
## 2023-09-16 21:40:19 Annotating text fragment 39811/42019
## 2023-09-16 21:40:19 Annotating text fragment 39821/42019
## 2023-09-16 21:40:19 Annotating text fragment 39831/42019
## 2023-09-16 21:40:19 Annotating text fragment 39841/42019
## 2023-09-16 21:40:20 Annotating text fragment 39851/42019
## 2023-09-16 21:40:20 Annotating text fragment 39861/42019
## 2023-09-16 21:40:20 Annotating text fragment 39871/42019
## 2023-09-16 21:40:20 Annotating text fragment 39881/42019
## 2023-09-16 21:40:20 Annotating text fragment 39891/42019
## 2023-09-16 21:40:20 Annotating text fragment 39901/42019
## 2023-09-16 21:40:20 Annotating text fragment 39911/42019
## 2023-09-16 21:40:20 Annotating text fragment 39921/42019
## 2023-09-16 21:40:21 Annotating text fragment 39931/42019
## 2023-09-16 21:40:21 Annotating text fragment 39941/42019
## 2023-09-16 21:40:21 Annotating text fragment 39951/42019
## 2023-09-16 21:40:21 Annotating text fragment 39961/42019
## 2023-09-16 21:40:21 Annotating text fragment 39971/42019
## 2023-09-16 21:40:21 Annotating text fragment 39981/42019
## 2023-09-16 21:40:21 Annotating text fragment 39991/42019
## 2023-09-16 21:40:22 Annotating text fragment 40001/42019
## 2023-09-16 21:40:22 Annotating text fragment 40011/42019
## 2023-09-16 21:40:22 Annotating text fragment 40021/42019
## 2023-09-16 21:40:22 Annotating text fragment 40031/42019
## 2023-09-16 21:40:22 Annotating text fragment 40041/42019
## 2023-09-16 21:40:22 Annotating text fragment 40051/42019
## 2023-09-16 21:40:22 Annotating text fragment 40061/42019
## 2023-09-16 21:40:23 Annotating text fragment 40071/42019
## 2023-09-16 21:40:23 Annotating text fragment 40081/42019
## 2023-09-16 21:40:23 Annotating text fragment 40091/42019
## 2023-09-16 21:40:23 Annotating text fragment 40101/42019
```

```
## 2023-09-16 21:40:23 Annotating text fragment 40111/42019
## 2023-09-16 21:40:23 Annotating text fragment 40121/42019
## 2023-09-16 21:40:23 Annotating text fragment 40131/42019
## 2023-09-16 21:40:24 Annotating text fragment 40141/42019
## 2023-09-16 21:40:24 Annotating text fragment 40151/42019
## 2023-09-16 21:40:24 Annotating text fragment 40161/42019
## 2023-09-16 21:40:24 Annotating text fragment 40171/42019
## 2023-09-16 21:40:24 Annotating text fragment 40181/42019
## 2023-09-16 21:40:24 Annotating text fragment 40191/42019
## 2023-09-16 21:40:24 Annotating text fragment 40201/42019
## 2023-09-16 21:40:25 Annotating text fragment 40211/42019
## 2023-09-16 21:40:25 Annotating text fragment 40221/42019
## 2023-09-16 21:40:25 Annotating text fragment 40231/42019
## 2023-09-16 21:40:25 Annotating text fragment 40241/42019
## 2023-09-16 21:40:25 Annotating text fragment 40251/42019
## 2023-09-16 21:40:25 Annotating text fragment 40261/42019
## 2023-09-16 21:40:25 Annotating text fragment 40271/42019
## 2023-09-16 21:40:26 Annotating text fragment 40281/42019
## 2023-09-16 21:40:26 Annotating text fragment 40291/42019
## 2023-09-16 21:40:26 Annotating text fragment 40301/42019
## 2023-09-16 21:40:26 Annotating text fragment 40311/42019
## 2023-09-16 21:40:26 Annotating text fragment 40321/42019
## 2023-09-16 21:40:26 Annotating text fragment 40331/42019
## 2023-09-16 21:40:27 Annotating text fragment 40341/42019
## 2023-09-16 21:40:27 Annotating text fragment 40351/42019
## 2023-09-16 21:40:27 Annotating text fragment 40361/42019
## 2023-09-16 21:40:27 Annotating text fragment 40371/42019
## 2023-09-16 21:40:27 Annotating text fragment 40381/42019
## 2023-09-16 21:40:27 Annotating text fragment 40391/42019
## 2023-09-16 21:40:27 Annotating text fragment 40401/42019
## 2023-09-16 21:40:28 Annotating text fragment 40411/42019
## 2023-09-16 21:40:28 Annotating text fragment 40421/42019
## 2023-09-16 21:40:28 Annotating text fragment 40431/42019
## 2023-09-16 21:40:28 Annotating text fragment 40441/42019
## 2023-09-16 21:40:28 Annotating text fragment 40451/42019
## 2023-09-16 21:40:28 Annotating text fragment 40461/42019
## 2023-09-16 21:40:28 Annotating text fragment 40471/42019
## 2023-09-16 21:40:28 Annotating text fragment 40481/42019
## 2023-09-16 21:40:28 Annotating text fragment 40491/42019
## 2023-09-16 21:40:29 Annotating text fragment 40501/42019
## 2023-09-16 21:40:29 Annotating text fragment 40511/42019
## 2023-09-16 21:40:29 Annotating text fragment 40521/42019
## 2023-09-16 21:40:29 Annotating text fragment 40531/42019
## 2023-09-16 21:40:29 Annotating text fragment 40541/42019
## 2023-09-16 21:40:29 Annotating text fragment 40551/42019
## 2023-09-16 21:40:30 Annotating text fragment 40561/42019
## 2023-09-16 21:40:30 Annotating text fragment 40571/42019
## 2023-09-16 21:40:30 Annotating text fragment 40581/42019
## 2023-09-16 21:40:30 Annotating text fragment 40591/42019
## 2023-09-16 21:40:30 Annotating text fragment 40601/42019
## 2023-09-16 21:40:30 Annotating text fragment 40611/42019
## 2023-09-16 21:40:30 Annotating text fragment 40621/42019
## 2023-09-16 21:40:30 Annotating text fragment 40631/42019
## 2023-09-16 21:40:31 Annotating text fragment 40641/42019
```

```
## 2023-09-16 21:40:31 Annotating text fragment 40651/42019
## 2023-09-16 21:40:31 Annotating text fragment 40661/42019
## 2023-09-16 21:40:31 Annotating text fragment 40671/42019
## 2023-09-16 21:40:31 Annotating text fragment 40681/42019
## 2023-09-16 21:40:31 Annotating text fragment 40691/42019
## 2023-09-16 21:40:32 Annotating text fragment 40701/42019
## 2023-09-16 21:40:32 Annotating text fragment 40711/42019
## 2023-09-16 21:40:32 Annotating text fragment 40721/42019
## 2023-09-16 21:40:32 Annotating text fragment 40731/42019
## 2023-09-16 21:40:32 Annotating text fragment 40741/42019
## 2023-09-16 21:40:32 Annotating text fragment 40751/42019
## 2023-09-16 21:40:32 Annotating text fragment 40761/42019
## 2023-09-16 21:40:32 Annotating text fragment 40771/42019
## 2023-09-16 21:40:32 Annotating text fragment 40781/42019
## 2023-09-16 21:40:32 Annotating text fragment 40791/42019
## 2023-09-16 21:40:33 Annotating text fragment 40801/42019
## 2023-09-16 21:40:33 Annotating text fragment 40811/42019
## 2023-09-16 21:40:33 Annotating text fragment 40821/42019
## 2023-09-16 21:40:33 Annotating text fragment 40831/42019
## 2023-09-16 21:40:33 Annotating text fragment 40841/42019
## 2023-09-16 21:40:33 Annotating text fragment 40851/42019
## 2023-09-16 21:40:33 Annotating text fragment 40861/42019
## 2023-09-16 21:40:33 Annotating text fragment 40871/42019
## 2023-09-16 21:40:33 Annotating text fragment 40881/42019
## 2023-09-16 21:40:34 Annotating text fragment 40891/42019
## 2023-09-16 21:40:34 Annotating text fragment 40901/42019
## 2023-09-16 21:40:34 Annotating text fragment 40911/42019
## 2023-09-16 21:40:34 Annotating text fragment 40921/42019
## 2023-09-16 21:40:34 Annotating text fragment 40931/42019
## 2023-09-16 21:40:34 Annotating text fragment 40941/42019
## 2023-09-16 21:40:34 Annotating text fragment 40951/42019
## 2023-09-16 21:40:34 Annotating text fragment 40961/42019
## 2023-09-16 21:40:35 Annotating text fragment 40971/42019
## 2023-09-16 21:40:35 Annotating text fragment 40981/42019
## 2023-09-16 21:40:35 Annotating text fragment 40991/42019
## 2023-09-16 21:40:35 Annotating text fragment 41001/42019
## 2023-09-16 21:40:35 Annotating text fragment 41011/42019
## 2023-09-16 21:40:35 Annotating text fragment 41021/42019
## 2023-09-16 21:40:35 Annotating text fragment 41031/42019
## 2023-09-16 21:40:35 Annotating text fragment 41041/42019
## 2023-09-16 21:40:36 Annotating text fragment 41051/42019
## 2023-09-16 21:40:36 Annotating text fragment 41061/42019
## 2023-09-16 21:40:36 Annotating text fragment 41071/42019
## 2023-09-16 21:40:36 Annotating text fragment 41081/42019
## 2023-09-16 21:40:36 Annotating text fragment 41091/42019
## 2023-09-16 21:40:36 Annotating text fragment 41101/42019
## 2023-09-16 21:40:37 Annotating text fragment 41111/42019
## 2023-09-16 21:40:37 Annotating text fragment 41121/42019
## 2023-09-16 21:40:37 Annotating text fragment 41131/42019
## 2023-09-16 21:40:37 Annotating text fragment 41141/42019
## 2023-09-16 21:40:38 Annotating text fragment 41151/42019
## 2023-09-16 21:40:38 Annotating text fragment 41161/42019
## 2023-09-16 21:40:38 Annotating text fragment 41171/42019
## 2023-09-16 21:40:38 Annotating text fragment 41181/42019
```

```
## 2023-09-16 21:40:38 Annotating text fragment 41191/42019
## 2023-09-16 21:40:38 Annotating text fragment 41201/42019
## 2023-09-16 21:40:38 Annotating text fragment 41211/42019
## 2023-09-16 21:40:38 Annotating text fragment 41221/42019
## 2023-09-16 21:40:39 Annotating text fragment 41231/42019
## 2023-09-16 21:40:39 Annotating text fragment 41241/42019
## 2023-09-16 21:40:39 Annotating text fragment 41251/42019
## 2023-09-16 21:40:39 Annotating text fragment 41261/42019
## 2023-09-16 21:40:39 Annotating text fragment 41271/42019
## 2023-09-16 21:40:39 Annotating text fragment 41281/42019
## 2023-09-16 21:40:39 Annotating text fragment 41291/42019
## 2023-09-16 21:40:39 Annotating text fragment 41301/42019
## 2023-09-16 21:40:40 Annotating text fragment 41311/42019
## 2023-09-16 21:40:40 Annotating text fragment 41321/42019
## 2023-09-16 21:40:40 Annotating text fragment 41331/42019
## 2023-09-16 21:40:40 Annotating text fragment 41341/42019
## 2023-09-16 21:40:40 Annotating text fragment 41351/42019
## 2023-09-16 21:40:40 Annotating text fragment 41361/42019
## 2023-09-16 21:40:40 Annotating text fragment 41371/42019
## 2023-09-16 21:40:41 Annotating text fragment 41381/42019
## 2023-09-16 21:40:41 Annotating text fragment 41391/42019
## 2023-09-16 21:40:41 Annotating text fragment 41401/42019
## 2023-09-16 21:40:41 Annotating text fragment 41411/42019
## 2023-09-16 21:40:41 Annotating text fragment 41421/42019
## 2023-09-16 21:40:41 Annotating text fragment 41431/42019
## 2023-09-16 21:40:41 Annotating text fragment 41441/42019
## 2023-09-16 21:40:41 Annotating text fragment 41451/42019
## 2023-09-16 21:40:42 Annotating text fragment 41461/42019
## 2023-09-16 21:40:42 Annotating text fragment 41471/42019
## 2023-09-16 21:40:42 Annotating text fragment 41481/42019
## 2023-09-16 21:40:42 Annotating text fragment 41491/42019
## 2023-09-16 21:40:42 Annotating text fragment 41501/42019
## 2023-09-16 21:40:43 Annotating text fragment 41511/42019
## 2023-09-16 21:40:43 Annotating text fragment 41521/42019
## 2023-09-16 21:40:43 Annotating text fragment 41531/42019
## 2023-09-16 21:40:43 Annotating text fragment 41541/42019
## 2023-09-16 21:40:43 Annotating text fragment 41551/42019
## 2023-09-16 21:40:43 Annotating text fragment 41561/42019
## 2023-09-16 21:40:43 Annotating text fragment 41571/42019
## 2023-09-16 21:40:43 Annotating text fragment 41581/42019
## 2023-09-16 21:40:44 Annotating text fragment 41591/42019
## 2023-09-16 21:40:44 Annotating text fragment 41601/42019
## 2023-09-16 21:40:44 Annotating text fragment 41611/42019
## 2023-09-16 21:40:44 Annotating text fragment 41621/42019
## 2023-09-16 21:40:44 Annotating text fragment 41631/42019
## 2023-09-16 21:40:44 Annotating text fragment 41641/42019
## 2023-09-16 21:40:44 Annotating text fragment 41651/42019
## 2023-09-16 21:40:44 Annotating text fragment 41661/42019
## 2023-09-16 21:40:45 Annotating text fragment 41671/42019
## 2023-09-16 21:40:45 Annotating text fragment 41681/42019
## 2023-09-16 21:40:45 Annotating text fragment 41691/42019
## 2023-09-16 21:40:45 Annotating text fragment 41701/42019
## 2023-09-16 21:40:45 Annotating text fragment 41711/42019
## 2023-09-16 21:40:45 Annotating text fragment 41721/42019
```

```
## 2023-09-16 21:40:45 Annotating text fragment 41731/42019
## 2023-09-16 21:40:46 Annotating text fragment 41741/42019
## 2023-09-16 21:40:46 Annotating text fragment 41751/42019
## 2023-09-16 21:40:46 Annotating text fragment 41761/42019
## 2023-09-16 21:40:46 Annotating text fragment 41771/42019
## 2023-09-16 21:40:46 Annotating text fragment 41781/42019
## 2023-09-16 21:40:46 Annotating text fragment 41791/42019
## 2023-09-16 21:40:46 Annotating text fragment 41801/42019
## 2023-09-16 21:40:47 Annotating text fragment 41811/42019
## 2023-09-16 21:40:47 Annotating text fragment 41821/42019
## 2023-09-16 21:40:47 Annotating text fragment 41831/42019
## 2023-09-16 21:40:47 Annotating text fragment 41841/42019
## 2023-09-16 21:40:47 Annotating text fragment 41851/42019
## 2023-09-16 21:40:47 Annotating text fragment 41861/42019
## 2023-09-16 21:40:47 Annotating text fragment 41871/42019
## 2023-09-16 21:40:47 Annotating text fragment 41881/42019
## 2023-09-16 21:40:48 Annotating text fragment 41891/42019
## 2023-09-16 21:40:48 Annotating text fragment 41901/42019
## 2023-09-16 21:40:48 Annotating text fragment 41911/42019
## 2023-09-16 21:40:48 Annotating text fragment 41921/42019
## 2023-09-16 21:40:48 Annotating text fragment 41931/42019
## 2023-09-16 21:40:48 Annotating text fragment 41941/42019
## 2023-09-16 21:40:48 Annotating text fragment 41951/42019
## 2023-09-16 21:40:48 Annotating text fragment 41961/42019
## 2023-09-16 21:40:49 Annotating text fragment 41971/42019
## 2023-09-16 21:40:49 Annotating text fragment 41981/42019
## 2023-09-16 21:40:49 Annotating text fragment 41991/42019
## 2023-09-16 21:40:49 Annotating text fragment 42001/42019
## 2023-09-16 21:40:49 Annotating text fragment 42011/42019


## 2023-09-16 21:41:39 Annotating text fragment 1/57597
## 2023-09-16 21:41:39 Annotating text fragment 11/57597
## 2023-09-16 21:41:39 Annotating text fragment 21/57597
## 2023-09-16 21:41:39 Annotating text fragment 31/57597
## 2023-09-16 21:41:39 Annotating text fragment 41/57597
## 2023-09-16 21:41:39 Annotating text fragment 51/57597
## 2023-09-16 21:41:39 Annotating text fragment 61/57597
## 2023-09-16 21:41:39 Annotating text fragment 71/57597
## 2023-09-16 21:41:39 Annotating text fragment 81/57597
## 2023-09-16 21:41:40 Annotating text fragment 91/57597
## 2023-09-16 21:41:40 Annotating text fragment 101/57597
## 2023-09-16 21:41:40 Annotating text fragment 111/57597
## 2023-09-16 21:41:40 Annotating text fragment 121/57597
## 2023-09-16 21:41:40 Annotating text fragment 131/57597
## 2023-09-16 21:41:40 Annotating text fragment 141/57597
## 2023-09-16 21:41:40 Annotating text fragment 151/57597
## 2023-09-16 21:41:40 Annotating text fragment 161/57597
## 2023-09-16 21:41:40 Annotating text fragment 171/57597
## 2023-09-16 21:41:40 Annotating text fragment 181/57597
## 2023-09-16 21:41:40 Annotating text fragment 191/57597
## 2023-09-16 21:41:41 Annotating text fragment 201/57597
## 2023-09-16 21:41:41 Annotating text fragment 211/57597
## 2023-09-16 21:41:41 Annotating text fragment 221/57597
## 2023-09-16 21:41:41 Annotating text fragment 231/57597
```

```
## 2023-09-16 21:41:41 Annotating text fragment 241/57597
## 2023-09-16 21:41:41 Annotating text fragment 251/57597
## 2023-09-16 21:41:41 Annotating text fragment 261/57597
## 2023-09-16 21:41:41 Annotating text fragment 271/57597
## 2023-09-16 21:41:41 Annotating text fragment 281/57597
## 2023-09-16 21:41:41 Annotating text fragment 291/57597
## 2023-09-16 21:41:41 Annotating text fragment 301/57597
## 2023-09-16 21:41:41 Annotating text fragment 311/57597
## 2023-09-16 21:41:42 Annotating text fragment 321/57597
## 2023-09-16 21:41:42 Annotating text fragment 331/57597
## 2023-09-16 21:41:42 Annotating text fragment 341/57597
## 2023-09-16 21:41:42 Annotating text fragment 351/57597
## 2023-09-16 21:41:42 Annotating text fragment 361/57597
## 2023-09-16 21:41:42 Annotating text fragment 371/57597
## 2023-09-16 21:41:42 Annotating text fragment 381/57597
## 2023-09-16 21:41:42 Annotating text fragment 391/57597
## 2023-09-16 21:41:42 Annotating text fragment 401/57597
## 2023-09-16 21:41:42 Annotating text fragment 411/57597
## 2023-09-16 21:41:42 Annotating text fragment 421/57597
## 2023-09-16 21:41:43 Annotating text fragment 431/57597
## 2023-09-16 21:41:43 Annotating text fragment 441/57597
## 2023-09-16 21:41:43 Annotating text fragment 451/57597
## 2023-09-16 21:41:43 Annotating text fragment 461/57597
## 2023-09-16 21:41:43 Annotating text fragment 471/57597
## 2023-09-16 21:41:43 Annotating text fragment 481/57597
## 2023-09-16 21:41:43 Annotating text fragment 491/57597
## 2023-09-16 21:41:43 Annotating text fragment 501/57597
## 2023-09-16 21:41:44 Annotating text fragment 511/57597
## 2023-09-16 21:41:44 Annotating text fragment 521/57597
## 2023-09-16 21:41:44 Annotating text fragment 531/57597
## 2023-09-16 21:41:44 Annotating text fragment 541/57597
## 2023-09-16 21:41:44 Annotating text fragment 551/57597
## 2023-09-16 21:41:44 Annotating text fragment 561/57597
## 2023-09-16 21:41:45 Annotating text fragment 571/57597
## 2023-09-16 21:41:45 Annotating text fragment 581/57597
## 2023-09-16 21:41:45 Annotating text fragment 591/57597
## 2023-09-16 21:41:45 Annotating text fragment 601/57597
## 2023-09-16 21:41:45 Annotating text fragment 611/57597
## 2023-09-16 21:41:45 Annotating text fragment 621/57597
## 2023-09-16 21:41:45 Annotating text fragment 631/57597
## 2023-09-16 21:41:45 Annotating text fragment 641/57597
## 2023-09-16 21:41:45 Annotating text fragment 651/57597
## 2023-09-16 21:41:45 Annotating text fragment 661/57597
## 2023-09-16 21:41:45 Annotating text fragment 671/57597
## 2023-09-16 21:41:45 Annotating text fragment 681/57597
## 2023-09-16 21:41:46 Annotating text fragment 691/57597
## 2023-09-16 21:41:46 Annotating text fragment 701/57597
## 2023-09-16 21:41:46 Annotating text fragment 711/57597
## 2023-09-16 21:41:46 Annotating text fragment 721/57597
## 2023-09-16 21:41:46 Annotating text fragment 731/57597
## 2023-09-16 21:41:46 Annotating text fragment 741/57597
## 2023-09-16 21:41:46 Annotating text fragment 751/57597
## 2023-09-16 21:41:47 Annotating text fragment 761/57597
## 2023-09-16 21:41:47 Annotating text fragment 771/57597
```

```
## 2023-09-16 21:41:47 Annotating text fragment 781/57597
## 2023-09-16 21:41:47 Annotating text fragment 791/57597
## 2023-09-16 21:41:47 Annotating text fragment 801/57597
## 2023-09-16 21:41:47 Annotating text fragment 811/57597
## 2023-09-16 21:41:47 Annotating text fragment 821/57597
## 2023-09-16 21:41:47 Annotating text fragment 831/57597
## 2023-09-16 21:41:47 Annotating text fragment 841/57597
## 2023-09-16 21:41:47 Annotating text fragment 851/57597
## 2023-09-16 21:41:47 Annotating text fragment 861/57597
## 2023-09-16 21:41:47 Annotating text fragment 871/57597
## 2023-09-16 21:41:47 Annotating text fragment 881/57597
## 2023-09-16 21:41:47 Annotating text fragment 891/57597
## 2023-09-16 21:41:48 Annotating text fragment 901/57597
## 2023-09-16 21:41:48 Annotating text fragment 911/57597
## 2023-09-16 21:41:48 Annotating text fragment 921/57597
## 2023-09-16 21:41:48 Annotating text fragment 931/57597
## 2023-09-16 21:41:48 Annotating text fragment 941/57597
## 2023-09-16 21:41:48 Annotating text fragment 951/57597
## 2023-09-16 21:41:48 Annotating text fragment 961/57597
## 2023-09-16 21:41:48 Annotating text fragment 971/57597
## 2023-09-16 21:41:48 Annotating text fragment 981/57597
## 2023-09-16 21:41:48 Annotating text fragment 991/57597
## 2023-09-16 21:41:49 Annotating text fragment 1001/57597
## 2023-09-16 21:41:49 Annotating text fragment 1011/57597
## 2023-09-16 21:41:49 Annotating text fragment 1021/57597
## 2023-09-16 21:41:49 Annotating text fragment 1031/57597
## 2023-09-16 21:41:49 Annotating text fragment 1041/57597
## 2023-09-16 21:41:49 Annotating text fragment 1051/57597
## 2023-09-16 21:41:49 Annotating text fragment 1061/57597
## 2023-09-16 21:41:49 Annotating text fragment 1071/57597
## 2023-09-16 21:41:49 Annotating text fragment 1081/57597
## 2023-09-16 21:41:50 Annotating text fragment 1091/57597
## 2023-09-16 21:41:50 Annotating text fragment 1101/57597
## 2023-09-16 21:41:50 Annotating text fragment 1111/57597
## 2023-09-16 21:41:50 Annotating text fragment 1121/57597
## 2023-09-16 21:41:50 Annotating text fragment 1131/57597
## 2023-09-16 21:41:50 Annotating text fragment 1141/57597
## 2023-09-16 21:41:50 Annotating text fragment 1151/57597
## 2023-09-16 21:41:50 Annotating text fragment 1161/57597
## 2023-09-16 21:41:50 Annotating text fragment 1171/57597
## 2023-09-16 21:41:50 Annotating text fragment 1181/57597
## 2023-09-16 21:41:50 Annotating text fragment 1191/57597
## 2023-09-16 21:41:50 Annotating text fragment 1201/57597
## 2023-09-16 21:41:50 Annotating text fragment 1211/57597
## 2023-09-16 21:41:50 Annotating text fragment 1221/57597
## 2023-09-16 21:41:50 Annotating text fragment 1231/57597
## 2023-09-16 21:41:51 Annotating text fragment 1241/57597
## 2023-09-16 21:41:51 Annotating text fragment 1251/57597
## 2023-09-16 21:41:51 Annotating text fragment 1261/57597
## 2023-09-16 21:41:51 Annotating text fragment 1271/57597
## 2023-09-16 21:41:51 Annotating text fragment 1281/57597
## 2023-09-16 21:41:51 Annotating text fragment 1291/57597
## 2023-09-16 21:41:51 Annotating text fragment 1301/57597
## 2023-09-16 21:41:51 Annotating text fragment 1311/57597
```

```
## 2023-09-16 21:41:51 Annotating text fragment 1321/57597
## 2023-09-16 21:41:51 Annotating text fragment 1331/57597
## 2023-09-16 21:41:51 Annotating text fragment 1341/57597
## 2023-09-16 21:41:51 Annotating text fragment 1351/57597
## 2023-09-16 21:41:51 Annotating text fragment 1361/57597
## 2023-09-16 21:41:51 Annotating text fragment 1371/57597
## 2023-09-16 21:41:51 Annotating text fragment 1381/57597
## 2023-09-16 21:41:52 Annotating text fragment 1391/57597
## 2023-09-16 21:41:52 Annotating text fragment 1401/57597
## 2023-09-16 21:41:52 Annotating text fragment 1411/57597
## 2023-09-16 21:41:52 Annotating text fragment 1421/57597
## 2023-09-16 21:41:52 Annotating text fragment 1431/57597
## 2023-09-16 21:41:52 Annotating text fragment 1441/57597
## 2023-09-16 21:41:52 Annotating text fragment 1451/57597
## 2023-09-16 21:41:52 Annotating text fragment 1461/57597
## 2023-09-16 21:41:52 Annotating text fragment 1471/57597
## 2023-09-16 21:41:52 Annotating text fragment 1481/57597
## 2023-09-16 21:41:52 Annotating text fragment 1491/57597
## 2023-09-16 21:41:52 Annotating text fragment 1501/57597
## 2023-09-16 21:41:52 Annotating text fragment 1511/57597
## 2023-09-16 21:41:53 Annotating text fragment 1521/57597
## 2023-09-16 21:41:53 Annotating text fragment 1531/57597
## 2023-09-16 21:41:53 Annotating text fragment 1541/57597
## 2023-09-16 21:41:53 Annotating text fragment 1551/57597
## 2023-09-16 21:41:53 Annotating text fragment 1561/57597
## 2023-09-16 21:41:53 Annotating text fragment 1571/57597
## 2023-09-16 21:41:53 Annotating text fragment 1581/57597
## 2023-09-16 21:41:53 Annotating text fragment 1591/57597
## 2023-09-16 21:41:54 Annotating text fragment 1601/57597
## 2023-09-16 21:41:54 Annotating text fragment 1611/57597
## 2023-09-16 21:41:54 Annotating text fragment 1621/57597
## 2023-09-16 21:41:54 Annotating text fragment 1631/57597
## 2023-09-16 21:41:54 Annotating text fragment 1641/57597
## 2023-09-16 21:41:54 Annotating text fragment 1651/57597
## 2023-09-16 21:41:54 Annotating text fragment 1661/57597
## 2023-09-16 21:41:54 Annotating text fragment 1671/57597
## 2023-09-16 21:41:54 Annotating text fragment 1681/57597
## 2023-09-16 21:41:54 Annotating text fragment 1691/57597
## 2023-09-16 21:41:54 Annotating text fragment 1701/57597
## 2023-09-16 21:41:54 Annotating text fragment 1711/57597
## 2023-09-16 21:41:54 Annotating text fragment 1721/57597
## 2023-09-16 21:41:55 Annotating text fragment 1731/57597
## 2023-09-16 21:41:55 Annotating text fragment 1741/57597
## 2023-09-16 21:41:55 Annotating text fragment 1751/57597
## 2023-09-16 21:41:55 Annotating text fragment 1761/57597
## 2023-09-16 21:41:55 Annotating text fragment 1771/57597
## 2023-09-16 21:41:55 Annotating text fragment 1781/57597
## 2023-09-16 21:41:55 Annotating text fragment 1791/57597
## 2023-09-16 21:41:55 Annotating text fragment 1801/57597
## 2023-09-16 21:41:55 Annotating text fragment 1811/57597
## 2023-09-16 21:41:55 Annotating text fragment 1821/57597
## 2023-09-16 21:41:55 Annotating text fragment 1831/57597
## 2023-09-16 21:41:55 Annotating text fragment 1841/57597
## 2023-09-16 21:41:55 Annotating text fragment 1851/57597
```

```
## 2023-09-16 21:41:55 Annotating text fragment 1861/57597
## 2023-09-16 21:41:55 Annotating text fragment 1871/57597
## 2023-09-16 21:41:55 Annotating text fragment 1881/57597
## 2023-09-16 21:41:55 Annotating text fragment 1891/57597
## 2023-09-16 21:41:55 Annotating text fragment 1901/57597
## 2023-09-16 21:41:56 Annotating text fragment 1911/57597
## 2023-09-16 21:41:56 Annotating text fragment 1921/57597
## 2023-09-16 21:41:56 Annotating text fragment 1931/57597
## 2023-09-16 21:41:56 Annotating text fragment 1941/57597
## 2023-09-16 21:41:56 Annotating text fragment 1951/57597
## 2023-09-16 21:41:56 Annotating text fragment 1961/57597
## 2023-09-16 21:41:56 Annotating text fragment 1971/57597
## 2023-09-16 21:41:56 Annotating text fragment 1981/57597
## 2023-09-16 21:41:56 Annotating text fragment 1991/57597
## 2023-09-16 21:41:56 Annotating text fragment 2001/57597
## 2023-09-16 21:41:56 Annotating text fragment 2011/57597
## 2023-09-16 21:41:56 Annotating text fragment 2021/57597
## 2023-09-16 21:41:56 Annotating text fragment 2031/57597
## 2023-09-16 21:41:56 Annotating text fragment 2041/57597
## 2023-09-16 21:41:57 Annotating text fragment 2051/57597
## 2023-09-16 21:41:57 Annotating text fragment 2061/57597
## 2023-09-16 21:41:57 Annotating text fragment 2071/57597
## 2023-09-16 21:41:57 Annotating text fragment 2081/57597
## 2023-09-16 21:41:57 Annotating text fragment 2091/57597
## 2023-09-16 21:41:57 Annotating text fragment 2101/57597
## 2023-09-16 21:41:57 Annotating text fragment 2111/57597
## 2023-09-16 21:41:57 Annotating text fragment 2121/57597
## 2023-09-16 21:41:57 Annotating text fragment 2131/57597
## 2023-09-16 21:41:57 Annotating text fragment 2141/57597
## 2023-09-16 21:41:57 Annotating text fragment 2151/57597
## 2023-09-16 21:41:57 Annotating text fragment 2161/57597
## 2023-09-16 21:41:57 Annotating text fragment 2171/57597
## 2023-09-16 21:41:57 Annotating text fragment 2181/57597
## 2023-09-16 21:41:57 Annotating text fragment 2191/57597
## 2023-09-16 21:41:57 Annotating text fragment 2201/57597
## 2023-09-16 21:41:57 Annotating text fragment 2211/57597
## 2023-09-16 21:41:57 Annotating text fragment 2221/57597
## 2023-09-16 21:41:58 Annotating text fragment 2231/57597
## 2023-09-16 21:41:58 Annotating text fragment 2241/57597
## 2023-09-16 21:41:58 Annotating text fragment 2251/57597
## 2023-09-16 21:41:58 Annotating text fragment 2261/57597
## 2023-09-16 21:41:58 Annotating text fragment 2271/57597
## 2023-09-16 21:41:58 Annotating text fragment 2281/57597
## 2023-09-16 21:41:58 Annotating text fragment 2291/57597
## 2023-09-16 21:41:58 Annotating text fragment 2301/57597
## 2023-09-16 21:41:58 Annotating text fragment 2311/57597
## 2023-09-16 21:41:58 Annotating text fragment 2321/57597
## 2023-09-16 21:41:58 Annotating text fragment 2331/57597
## 2023-09-16 21:41:58 Annotating text fragment 2341/57597
## 2023-09-16 21:41:58 Annotating text fragment 2351/57597
## 2023-09-16 21:41:58 Annotating text fragment 2361/57597
## 2023-09-16 21:41:58 Annotating text fragment 2371/57597
## 2023-09-16 21:41:59 Annotating text fragment 2381/57597
## 2023-09-16 21:41:59 Annotating text fragment 2391/57597
```

```
## 2023-09-16 21:41:59 Annotating text fragment 2401/57597
## 2023-09-16 21:41:59 Annotating text fragment 2411/57597
## 2023-09-16 21:41:59 Annotating text fragment 2421/57597
## 2023-09-16 21:41:59 Annotating text fragment 2431/57597
## 2023-09-16 21:41:59 Annotating text fragment 2441/57597
## 2023-09-16 21:41:59 Annotating text fragment 2451/57597
## 2023-09-16 21:41:59 Annotating text fragment 2461/57597
## 2023-09-16 21:41:59 Annotating text fragment 2471/57597
## 2023-09-16 21:41:59 Annotating text fragment 2481/57597
## 2023-09-16 21:41:59 Annotating text fragment 2491/57597
## 2023-09-16 21:41:59 Annotating text fragment 2501/57597
## 2023-09-16 21:42:00 Annotating text fragment 2511/57597
## 2023-09-16 21:42:00 Annotating text fragment 2521/57597
## 2023-09-16 21:42:00 Annotating text fragment 2531/57597
## 2023-09-16 21:42:00 Annotating text fragment 2541/57597
## 2023-09-16 21:42:01 Annotating text fragment 2551/57597
## 2023-09-16 21:42:01 Annotating text fragment 2561/57597
## 2023-09-16 21:42:01 Annotating text fragment 2571/57597
## 2023-09-16 21:42:01 Annotating text fragment 2581/57597
## 2023-09-16 21:42:01 Annotating text fragment 2591/57597
## 2023-09-16 21:42:02 Annotating text fragment 2601/57597
## 2023-09-16 21:42:03 Annotating text fragment 2611/57597
## 2023-09-16 21:42:03 Annotating text fragment 2621/57597
## 2023-09-16 21:42:04 Annotating text fragment 2631/57597
## 2023-09-16 21:42:05 Annotating text fragment 2641/57597
## 2023-09-16 21:42:05 Annotating text fragment 2651/57597
## 2023-09-16 21:42:06 Annotating text fragment 2661/57597
## 2023-09-16 21:42:06 Annotating text fragment 2671/57597
## 2023-09-16 21:42:07 Annotating text fragment 2681/57597
## 2023-09-16 21:42:07 Annotating text fragment 2691/57597
## 2023-09-16 21:42:07 Annotating text fragment 2701/57597
## 2023-09-16 21:42:07 Annotating text fragment 2711/57597
## 2023-09-16 21:42:07 Annotating text fragment 2721/57597
## 2023-09-16 21:42:07 Annotating text fragment 2731/57597
## 2023-09-16 21:42:07 Annotating text fragment 2741/57597
## 2023-09-16 21:42:07 Annotating text fragment 2751/57597
## 2023-09-16 21:42:07 Annotating text fragment 2761/57597
## 2023-09-16 21:42:08 Annotating text fragment 2771/57597
## 2023-09-16 21:42:08 Annotating text fragment 2781/57597
## 2023-09-16 21:42:08 Annotating text fragment 2791/57597
## 2023-09-16 21:42:08 Annotating text fragment 2801/57597
## 2023-09-16 21:42:08 Annotating text fragment 2811/57597
## 2023-09-16 21:42:08 Annotating text fragment 2821/57597
## 2023-09-16 21:42:08 Annotating text fragment 2831/57597
## 2023-09-16 21:42:08 Annotating text fragment 2841/57597
## 2023-09-16 21:42:08 Annotating text fragment 2851/57597
## 2023-09-16 21:42:08 Annotating text fragment 2861/57597
## 2023-09-16 21:42:09 Annotating text fragment 2871/57597
## 2023-09-16 21:42:09 Annotating text fragment 2881/57597
## 2023-09-16 21:42:09 Annotating text fragment 2891/57597
## 2023-09-16 21:42:09 Annotating text fragment 2901/57597
## 2023-09-16 21:42:09 Annotating text fragment 2911/57597
## 2023-09-16 21:42:09 Annotating text fragment 2921/57597
## 2023-09-16 21:42:09 Annotating text fragment 2931/57597
```

```
## 2023-09-16 21:42:09 Annotating text fragment 2941/57597
## 2023-09-16 21:42:09 Annotating text fragment 2951/57597
## 2023-09-16 21:42:09 Annotating text fragment 2961/57597
## 2023-09-16 21:42:09 Annotating text fragment 2971/57597
## 2023-09-16 21:42:09 Annotating text fragment 2981/57597
## 2023-09-16 21:42:09 Annotating text fragment 2991/57597
## 2023-09-16 21:42:10 Annotating text fragment 3001/57597
## 2023-09-16 21:42:10 Annotating text fragment 3011/57597
## 2023-09-16 21:42:10 Annotating text fragment 3021/57597
## 2023-09-16 21:42:10 Annotating text fragment 3031/57597
## 2023-09-16 21:42:10 Annotating text fragment 3041/57597
## 2023-09-16 21:42:10 Annotating text fragment 3051/57597
## 2023-09-16 21:42:10 Annotating text fragment 3061/57597
## 2023-09-16 21:42:10 Annotating text fragment 3071/57597
## 2023-09-16 21:42:10 Annotating text fragment 3081/57597
## 2023-09-16 21:42:10 Annotating text fragment 3091/57597
## 2023-09-16 21:42:10 Annotating text fragment 3101/57597
## 2023-09-16 21:42:11 Annotating text fragment 3111/57597
## 2023-09-16 21:42:11 Annotating text fragment 3121/57597
## 2023-09-16 21:42:11 Annotating text fragment 3131/57597
## 2023-09-16 21:42:11 Annotating text fragment 3141/57597
## 2023-09-16 21:42:11 Annotating text fragment 3151/57597
## 2023-09-16 21:42:11 Annotating text fragment 3161/57597
## 2023-09-16 21:42:11 Annotating text fragment 3171/57597
## 2023-09-16 21:42:11 Annotating text fragment 3181/57597
## 2023-09-16 21:42:11 Annotating text fragment 3191/57597
## 2023-09-16 21:42:11 Annotating text fragment 3201/57597
## 2023-09-16 21:42:11 Annotating text fragment 3211/57597
## 2023-09-16 21:42:11 Annotating text fragment 3221/57597
## 2023-09-16 21:42:11 Annotating text fragment 3231/57597
## 2023-09-16 21:42:11 Annotating text fragment 3241/57597
## 2023-09-16 21:42:12 Annotating text fragment 3251/57597
## 2023-09-16 21:42:12 Annotating text fragment 3261/57597
## 2023-09-16 21:42:12 Annotating text fragment 3271/57597
## 2023-09-16 21:42:12 Annotating text fragment 3281/57597
## 2023-09-16 21:42:12 Annotating text fragment 3291/57597
## 2023-09-16 21:42:12 Annotating text fragment 3301/57597
## 2023-09-16 21:42:12 Annotating text fragment 3311/57597
## 2023-09-16 21:42:12 Annotating text fragment 3321/57597
## 2023-09-16 21:42:12 Annotating text fragment 3331/57597
## 2023-09-16 21:42:12 Annotating text fragment 3341/57597
## 2023-09-16 21:42:12 Annotating text fragment 3351/57597
## 2023-09-16 21:42:12 Annotating text fragment 3361/57597
## 2023-09-16 21:42:13 Annotating text fragment 3371/57597
## 2023-09-16 21:42:13 Annotating text fragment 3381/57597
## 2023-09-16 21:42:14 Annotating text fragment 3391/57597
## 2023-09-16 21:42:14 Annotating text fragment 3401/57597
## 2023-09-16 21:42:14 Annotating text fragment 3411/57597
## 2023-09-16 21:42:15 Annotating text fragment 3421/57597
## 2023-09-16 21:42:15 Annotating text fragment 3431/57597
## 2023-09-16 21:42:15 Annotating text fragment 3441/57597
## 2023-09-16 21:42:15 Annotating text fragment 3451/57597
## 2023-09-16 21:42:15 Annotating text fragment 3461/57597
## 2023-09-16 21:42:16 Annotating text fragment 3471/57597
```

```
## 2023-09-16 21:42:16 Annotating text fragment 3481/57597
## 2023-09-16 21:42:16 Annotating text fragment 3491/57597
## 2023-09-16 21:42:16 Annotating text fragment 3501/57597
## 2023-09-16 21:42:16 Annotating text fragment 3511/57597
## 2023-09-16 21:42:16 Annotating text fragment 3521/57597
## 2023-09-16 21:42:16 Annotating text fragment 3531/57597
## 2023-09-16 21:42:16 Annotating text fragment 3541/57597
## 2023-09-16 21:42:16 Annotating text fragment 3551/57597
## 2023-09-16 21:42:16 Annotating text fragment 3561/57597
## 2023-09-16 21:42:16 Annotating text fragment 3571/57597
## 2023-09-16 21:42:16 Annotating text fragment 3581/57597
## 2023-09-16 21:42:16 Annotating text fragment 3591/57597
## 2023-09-16 21:42:16 Annotating text fragment 3601/57597
## 2023-09-16 21:42:16 Annotating text fragment 3611/57597
## 2023-09-16 21:42:17 Annotating text fragment 3621/57597
## 2023-09-16 21:42:17 Annotating text fragment 3631/57597
## 2023-09-16 21:42:17 Annotating text fragment 3641/57597
## 2023-09-16 21:42:17 Annotating text fragment 3651/57597
## 2023-09-16 21:42:17 Annotating text fragment 3661/57597
## 2023-09-16 21:42:17 Annotating text fragment 3671/57597
## 2023-09-16 21:42:17 Annotating text fragment 3681/57597
## 2023-09-16 21:42:17 Annotating text fragment 3691/57597
## 2023-09-16 21:42:17 Annotating text fragment 3701/57597
## 2023-09-16 21:42:17 Annotating text fragment 3711/57597
## 2023-09-16 21:42:17 Annotating text fragment 3721/57597
## 2023-09-16 21:42:17 Annotating text fragment 3731/57597
## 2023-09-16 21:42:18 Annotating text fragment 3741/57597
## 2023-09-16 21:42:18 Annotating text fragment 3751/57597
## 2023-09-16 21:42:18 Annotating text fragment 3761/57597
## 2023-09-16 21:42:18 Annotating text fragment 3771/57597
## 2023-09-16 21:42:18 Annotating text fragment 3781/57597
## 2023-09-16 21:42:18 Annotating text fragment 3791/57597
## 2023-09-16 21:42:18 Annotating text fragment 3801/57597
## 2023-09-16 21:42:19 Annotating text fragment 3811/57597
## 2023-09-16 21:42:19 Annotating text fragment 3821/57597
## 2023-09-16 21:42:19 Annotating text fragment 3831/57597
## 2023-09-16 21:42:19 Annotating text fragment 3841/57597
## 2023-09-16 21:42:19 Annotating text fragment 3851/57597
## 2023-09-16 21:42:19 Annotating text fragment 3861/57597
## 2023-09-16 21:42:19 Annotating text fragment 3871/57597
## 2023-09-16 21:42:19 Annotating text fragment 3881/57597
## 2023-09-16 21:42:19 Annotating text fragment 3891/57597
## 2023-09-16 21:42:19 Annotating text fragment 3901/57597
## 2023-09-16 21:42:19 Annotating text fragment 3911/57597
## 2023-09-16 21:42:19 Annotating text fragment 3921/57597
## 2023-09-16 21:42:19 Annotating text fragment 3931/57597
## 2023-09-16 21:42:19 Annotating text fragment 3941/57597
## 2023-09-16 21:42:19 Annotating text fragment 3951/57597
## 2023-09-16 21:42:19 Annotating text fragment 3961/57597
## 2023-09-16 21:42:19 Annotating text fragment 3971/57597
## 2023-09-16 21:42:20 Annotating text fragment 3981/57597
## 2023-09-16 21:42:20 Annotating text fragment 3991/57597
## 2023-09-16 21:42:20 Annotating text fragment 4001/57597
## 2023-09-16 21:42:20 Annotating text fragment 4011/57597
```

```
## 2023-09-16 21:42:21 Annotating text fragment 4021/57597
## 2023-09-16 21:42:21 Annotating text fragment 4031/57597
## 2023-09-16 21:42:21 Annotating text fragment 4041/57597
## 2023-09-16 21:42:21 Annotating text fragment 4051/57597
## 2023-09-16 21:42:21 Annotating text fragment 4061/57597
## 2023-09-16 21:42:21 Annotating text fragment 4071/57597
## 2023-09-16 21:42:22 Annotating text fragment 4081/57597
## 2023-09-16 21:42:22 Annotating text fragment 4091/57597
## 2023-09-16 21:42:22 Annotating text fragment 4101/57597
## 2023-09-16 21:42:22 Annotating text fragment 4111/57597
## 2023-09-16 21:42:22 Annotating text fragment 4121/57597
## 2023-09-16 21:42:22 Annotating text fragment 4131/57597
## 2023-09-16 21:42:22 Annotating text fragment 4141/57597
## 2023-09-16 21:42:22 Annotating text fragment 4151/57597
## 2023-09-16 21:42:22 Annotating text fragment 4161/57597
## 2023-09-16 21:42:22 Annotating text fragment 4171/57597
## 2023-09-16 21:42:23 Annotating text fragment 4181/57597
## 2023-09-16 21:42:23 Annotating text fragment 4191/57597
## 2023-09-16 21:42:23 Annotating text fragment 4201/57597
## 2023-09-16 21:42:23 Annotating text fragment 4211/57597
## 2023-09-16 21:42:23 Annotating text fragment 4221/57597
## 2023-09-16 21:42:23 Annotating text fragment 4231/57597
## 2023-09-16 21:42:23 Annotating text fragment 4241/57597
## 2023-09-16 21:42:23 Annotating text fragment 4251/57597
## 2023-09-16 21:42:23 Annotating text fragment 4261/57597
## 2023-09-16 21:42:24 Annotating text fragment 4271/57597
## 2023-09-16 21:42:24 Annotating text fragment 4281/57597
## 2023-09-16 21:42:24 Annotating text fragment 4291/57597
## 2023-09-16 21:42:24 Annotating text fragment 4301/57597
## 2023-09-16 21:42:24 Annotating text fragment 4311/57597
## 2023-09-16 21:42:24 Annotating text fragment 4321/57597
## 2023-09-16 21:42:24 Annotating text fragment 4331/57597
## 2023-09-16 21:42:24 Annotating text fragment 4341/57597
## 2023-09-16 21:42:24 Annotating text fragment 4351/57597
## 2023-09-16 21:42:24 Annotating text fragment 4361/57597
## 2023-09-16 21:42:25 Annotating text fragment 4371/57597
## 2023-09-16 21:42:25 Annotating text fragment 4381/57597
## 2023-09-16 21:42:25 Annotating text fragment 4391/57597
## 2023-09-16 21:42:25 Annotating text fragment 4401/57597
## 2023-09-16 21:42:25 Annotating text fragment 4411/57597
## 2023-09-16 21:42:26 Annotating text fragment 4421/57597
## 2023-09-16 21:42:26 Annotating text fragment 4431/57597
## 2023-09-16 21:42:26 Annotating text fragment 4441/57597
## 2023-09-16 21:42:26 Annotating text fragment 4451/57597
## 2023-09-16 21:42:26 Annotating text fragment 4461/57597
## 2023-09-16 21:42:26 Annotating text fragment 4471/57597
## 2023-09-16 21:42:26 Annotating text fragment 4481/57597
## 2023-09-16 21:42:27 Annotating text fragment 4491/57597
## 2023-09-16 21:42:27 Annotating text fragment 4501/57597
## 2023-09-16 21:42:27 Annotating text fragment 4511/57597
## 2023-09-16 21:42:27 Annotating text fragment 4521/57597
## 2023-09-16 21:42:27 Annotating text fragment 4531/57597
## 2023-09-16 21:42:27 Annotating text fragment 4541/57597
## 2023-09-16 21:42:27 Annotating text fragment 4551/57597
```

```
## 2023-09-16 21:42:27 Annotating text fragment 4561/57597
## 2023-09-16 21:42:27 Annotating text fragment 4571/57597
## 2023-09-16 21:42:27 Annotating text fragment 4581/57597
## 2023-09-16 21:42:27 Annotating text fragment 4591/57597
## 2023-09-16 21:42:27 Annotating text fragment 4601/57597
## 2023-09-16 21:42:27 Annotating text fragment 4611/57597
## 2023-09-16 21:42:27 Annotating text fragment 4621/57597
## 2023-09-16 21:42:27 Annotating text fragment 4631/57597
## 2023-09-16 21:42:27 Annotating text fragment 4641/57597
## 2023-09-16 21:42:27 Annotating text fragment 4651/57597
## 2023-09-16 21:42:27 Annotating text fragment 4661/57597
## 2023-09-16 21:42:27 Annotating text fragment 4671/57597
## 2023-09-16 21:42:28 Annotating text fragment 4681/57597
## 2023-09-16 21:42:28 Annotating text fragment 4691/57597
## 2023-09-16 21:42:28 Annotating text fragment 4701/57597
## 2023-09-16 21:42:28 Annotating text fragment 4711/57597
## 2023-09-16 21:42:28 Annotating text fragment 4721/57597
## 2023-09-16 21:42:28 Annotating text fragment 4731/57597
## 2023-09-16 21:42:28 Annotating text fragment 4741/57597
## 2023-09-16 21:42:28 Annotating text fragment 4751/57597
## 2023-09-16 21:42:28 Annotating text fragment 4761/57597
## 2023-09-16 21:42:28 Annotating text fragment 4771/57597
## 2023-09-16 21:42:28 Annotating text fragment 4781/57597
## 2023-09-16 21:42:28 Annotating text fragment 4791/57597
## 2023-09-16 21:42:29 Annotating text fragment 4801/57597
## 2023-09-16 21:42:29 Annotating text fragment 4811/57597
## 2023-09-16 21:42:29 Annotating text fragment 4821/57597
## 2023-09-16 21:42:29 Annotating text fragment 4831/57597
## 2023-09-16 21:42:29 Annotating text fragment 4841/57597
## 2023-09-16 21:42:29 Annotating text fragment 4851/57597
## 2023-09-16 21:42:29 Annotating text fragment 4861/57597
## 2023-09-16 21:42:29 Annotating text fragment 4871/57597
## 2023-09-16 21:42:29 Annotating text fragment 4881/57597
## 2023-09-16 21:42:29 Annotating text fragment 4891/57597
## 2023-09-16 21:42:29 Annotating text fragment 4901/57597
## 2023-09-16 21:42:30 Annotating text fragment 4911/57597
## 2023-09-16 21:42:30 Annotating text fragment 4921/57597
## 2023-09-16 21:42:30 Annotating text fragment 4931/57597
## 2023-09-16 21:42:30 Annotating text fragment 4941/57597
## 2023-09-16 21:42:31 Annotating text fragment 4951/57597
## 2023-09-16 21:42:31 Annotating text fragment 4961/57597
## 2023-09-16 21:42:32 Annotating text fragment 4971/57597
## 2023-09-16 21:42:32 Annotating text fragment 4981/57597
## 2023-09-16 21:42:32 Annotating text fragment 4991/57597
## 2023-09-16 21:42:33 Annotating text fragment 5001/57597
## 2023-09-16 21:42:33 Annotating text fragment 5011/57597
## 2023-09-16 21:42:33 Annotating text fragment 5021/57597
## 2023-09-16 21:42:33 Annotating text fragment 5031/57597
## 2023-09-16 21:42:33 Annotating text fragment 5041/57597
## 2023-09-16 21:42:33 Annotating text fragment 5051/57597
## 2023-09-16 21:42:33 Annotating text fragment 5061/57597
## 2023-09-16 21:42:33 Annotating text fragment 5071/57597
## 2023-09-16 21:42:33 Annotating text fragment 5081/57597
## 2023-09-16 21:42:33 Annotating text fragment 5091/57597
```

```
## 2023-09-16 21:42:33 Annotating text fragment 5101/57597
## 2023-09-16 21:42:33 Annotating text fragment 5111/57597
## 2023-09-16 21:42:33 Annotating text fragment 5121/57597
## 2023-09-16 21:42:33 Annotating text fragment 5131/57597
## 2023-09-16 21:42:33 Annotating text fragment 5141/57597
## 2023-09-16 21:42:33 Annotating text fragment 5151/57597
## 2023-09-16 21:42:33 Annotating text fragment 5161/57597
## 2023-09-16 21:42:33 Annotating text fragment 5171/57597
## 2023-09-16 21:42:33 Annotating text fragment 5181/57597
## 2023-09-16 21:42:33 Annotating text fragment 5191/57597
## 2023-09-16 21:42:33 Annotating text fragment 5201/57597
## 2023-09-16 21:42:34 Annotating text fragment 5211/57597
## 2023-09-16 21:42:34 Annotating text fragment 5221/57597
## 2023-09-16 21:42:34 Annotating text fragment 5231/57597
## 2023-09-16 21:42:34 Annotating text fragment 5241/57597
## 2023-09-16 21:42:34 Annotating text fragment 5251/57597
## 2023-09-16 21:42:34 Annotating text fragment 5261/57597
## 2023-09-16 21:42:34 Annotating text fragment 5271/57597
## 2023-09-16 21:42:34 Annotating text fragment 5281/57597
## 2023-09-16 21:42:34 Annotating text fragment 5291/57597
## 2023-09-16 21:42:34 Annotating text fragment 5301/57597
## 2023-09-16 21:42:34 Annotating text fragment 5311/57597
## 2023-09-16 21:42:34 Annotating text fragment 5321/57597
## 2023-09-16 21:42:34 Annotating text fragment 5331/57597
## 2023-09-16 21:42:34 Annotating text fragment 5341/57597
## 2023-09-16 21:42:34 Annotating text fragment 5351/57597
## 2023-09-16 21:42:34 Annotating text fragment 5361/57597
## 2023-09-16 21:42:34 Annotating text fragment 5371/57597
## 2023-09-16 21:42:34 Annotating text fragment 5381/57597
## 2023-09-16 21:42:35 Annotating text fragment 5391/57597
## 2023-09-16 21:42:35 Annotating text fragment 5401/57597
## 2023-09-16 21:42:35 Annotating text fragment 5411/57597
## 2023-09-16 21:42:35 Annotating text fragment 5421/57597
## 2023-09-16 21:42:35 Annotating text fragment 5431/57597
## 2023-09-16 21:42:35 Annotating text fragment 5441/57597
## 2023-09-16 21:42:35 Annotating text fragment 5451/57597
## 2023-09-16 21:42:35 Annotating text fragment 5461/57597
## 2023-09-16 21:42:35 Annotating text fragment 5471/57597
## 2023-09-16 21:42:35 Annotating text fragment 5481/57597
## 2023-09-16 21:42:35 Annotating text fragment 5491/57597
## 2023-09-16 21:42:35 Annotating text fragment 5501/57597
## 2023-09-16 21:42:35 Annotating text fragment 5511/57597
## 2023-09-16 21:42:36 Annotating text fragment 5521/57597
## 2023-09-16 21:42:36 Annotating text fragment 5531/57597
## 2023-09-16 21:42:36 Annotating text fragment 5541/57597
## 2023-09-16 21:42:36 Annotating text fragment 5551/57597
## 2023-09-16 21:42:36 Annotating text fragment 5561/57597
## 2023-09-16 21:42:36 Annotating text fragment 5571/57597
## 2023-09-16 21:42:36 Annotating text fragment 5581/57597
## 2023-09-16 21:42:36 Annotating text fragment 5591/57597
## 2023-09-16 21:42:36 Annotating text fragment 5601/57597
## 2023-09-16 21:42:36 Annotating text fragment 5611/57597
## 2023-09-16 21:42:36 Annotating text fragment 5621/57597
## 2023-09-16 21:42:36 Annotating text fragment 5631/57597
```

```
## 2023-09-16 21:42:36 Annotating text fragment 5641/57597
## 2023-09-16 21:42:36 Annotating text fragment 5651/57597
## 2023-09-16 21:42:36 Annotating text fragment 5661/57597
## 2023-09-16 21:42:36 Annotating text fragment 5671/57597
## 2023-09-16 21:42:36 Annotating text fragment 5681/57597
## 2023-09-16 21:42:36 Annotating text fragment 5691/57597
## 2023-09-16 21:42:36 Annotating text fragment 5701/57597
## 2023-09-16 21:42:36 Annotating text fragment 5711/57597
## 2023-09-16 21:42:37 Annotating text fragment 5721/57597
## 2023-09-16 21:42:37 Annotating text fragment 5731/57597
## 2023-09-16 21:42:37 Annotating text fragment 5741/57597
## 2023-09-16 21:42:37 Annotating text fragment 5751/57597
## 2023-09-16 21:42:37 Annotating text fragment 5761/57597
## 2023-09-16 21:42:37 Annotating text fragment 5771/57597
## 2023-09-16 21:42:37 Annotating text fragment 5781/57597
## 2023-09-16 21:42:37 Annotating text fragment 5791/57597
## 2023-09-16 21:42:37 Annotating text fragment 5801/57597
## 2023-09-16 21:42:37 Annotating text fragment 5811/57597
## 2023-09-16 21:42:37 Annotating text fragment 5821/57597
## 2023-09-16 21:42:37 Annotating text fragment 5831/57597
## 2023-09-16 21:42:37 Annotating text fragment 5841/57597
## 2023-09-16 21:42:37 Annotating text fragment 5851/57597
## 2023-09-16 21:42:37 Annotating text fragment 5861/57597
## 2023-09-16 21:42:37 Annotating text fragment 5871/57597
## 2023-09-16 21:42:37 Annotating text fragment 5881/57597
## 2023-09-16 21:42:37 Annotating text fragment 5891/57597
## 2023-09-16 21:42:37 Annotating text fragment 5901/57597
## 2023-09-16 21:42:38 Annotating text fragment 5911/57597
## 2023-09-16 21:42:38 Annotating text fragment 5921/57597
## 2023-09-16 21:42:38 Annotating text fragment 5931/57597
## 2023-09-16 21:42:38 Annotating text fragment 5941/57597
## 2023-09-16 21:42:38 Annotating text fragment 5951/57597
## 2023-09-16 21:42:38 Annotating text fragment 5961/57597
## 2023-09-16 21:42:38 Annotating text fragment 5971/57597
## 2023-09-16 21:42:39 Annotating text fragment 5981/57597
## 2023-09-16 21:42:39 Annotating text fragment 5991/57597
## 2023-09-16 21:42:39 Annotating text fragment 6001/57597
## 2023-09-16 21:42:39 Annotating text fragment 6011/57597
## 2023-09-16 21:42:40 Annotating text fragment 6021/57597
## 2023-09-16 21:42:40 Annotating text fragment 6031/57597
## 2023-09-16 21:42:40 Annotating text fragment 6041/57597
## 2023-09-16 21:42:40 Annotating text fragment 6051/57597
## 2023-09-16 21:42:40 Annotating text fragment 6061/57597
## 2023-09-16 21:42:40 Annotating text fragment 6071/57597
## 2023-09-16 21:42:40 Annotating text fragment 6081/57597
## 2023-09-16 21:42:40 Annotating text fragment 6091/57597
## 2023-09-16 21:42:40 Annotating text fragment 6101/57597
## 2023-09-16 21:42:40 Annotating text fragment 6111/57597
## 2023-09-16 21:42:40 Annotating text fragment 6121/57597
## 2023-09-16 21:42:40 Annotating text fragment 6131/57597
## 2023-09-16 21:42:40 Annotating text fragment 6141/57597
## 2023-09-16 21:42:40 Annotating text fragment 6151/57597
## 2023-09-16 21:42:40 Annotating text fragment 6161/57597
## 2023-09-16 21:42:40 Annotating text fragment 6171/57597
```

```
## 2023-09-16 21:42:41 Annotating text fragment 6181/57597
## 2023-09-16 21:42:41 Annotating text fragment 6191/57597
## 2023-09-16 21:42:41 Annotating text fragment 6201/57597
## 2023-09-16 21:42:41 Annotating text fragment 6211/57597
## 2023-09-16 21:42:41 Annotating text fragment 6221/57597
## 2023-09-16 21:42:41 Annotating text fragment 6231/57597
## 2023-09-16 21:42:42 Annotating text fragment 6241/57597
## 2023-09-16 21:42:42 Annotating text fragment 6251/57597
## 2023-09-16 21:42:42 Annotating text fragment 6261/57597
## 2023-09-16 21:42:42 Annotating text fragment 6271/57597
## 2023-09-16 21:42:42 Annotating text fragment 6281/57597
## 2023-09-16 21:42:43 Annotating text fragment 6291/57597
## 2023-09-16 21:42:43 Annotating text fragment 6301/57597
## 2023-09-16 21:42:43 Annotating text fragment 6311/57597
## 2023-09-16 21:42:44 Annotating text fragment 6321/57597
## 2023-09-16 21:42:44 Annotating text fragment 6331/57597
## 2023-09-16 21:42:44 Annotating text fragment 6341/57597
## 2023-09-16 21:42:44 Annotating text fragment 6351/57597
## 2023-09-16 21:42:44 Annotating text fragment 6361/57597
## 2023-09-16 21:42:44 Annotating text fragment 6371/57597
## 2023-09-16 21:42:44 Annotating text fragment 6381/57597
## 2023-09-16 21:42:44 Annotating text fragment 6391/57597
## 2023-09-16 21:42:44 Annotating text fragment 6401/57597
## 2023-09-16 21:42:44 Annotating text fragment 6411/57597
## 2023-09-16 21:42:45 Annotating text fragment 6421/57597
## 2023-09-16 21:42:45 Annotating text fragment 6431/57597
## 2023-09-16 21:42:45 Annotating text fragment 6441/57597
## 2023-09-16 21:42:45 Annotating text fragment 6451/57597
## 2023-09-16 21:42:45 Annotating text fragment 6461/57597
## 2023-09-16 21:42:45 Annotating text fragment 6471/57597
## 2023-09-16 21:42:45 Annotating text fragment 6481/57597
## 2023-09-16 21:42:45 Annotating text fragment 6491/57597
## 2023-09-16 21:42:45 Annotating text fragment 6501/57597
## 2023-09-16 21:42:45 Annotating text fragment 6511/57597
## 2023-09-16 21:42:45 Annotating text fragment 6521/57597
## 2023-09-16 21:42:45 Annotating text fragment 6531/57597
## 2023-09-16 21:42:45 Annotating text fragment 6541/57597
## 2023-09-16 21:42:45 Annotating text fragment 6551/57597
## 2023-09-16 21:42:46 Annotating text fragment 6561/57597
## 2023-09-16 21:42:46 Annotating text fragment 6571/57597
## 2023-09-16 21:42:46 Annotating text fragment 6581/57597
## 2023-09-16 21:42:46 Annotating text fragment 6591/57597
## 2023-09-16 21:42:46 Annotating text fragment 6601/57597
## 2023-09-16 21:42:46 Annotating text fragment 6611/57597
## 2023-09-16 21:42:46 Annotating text fragment 6621/57597
## 2023-09-16 21:42:46 Annotating text fragment 6631/57597
## 2023-09-16 21:42:46 Annotating text fragment 6641/57597
## 2023-09-16 21:42:46 Annotating text fragment 6651/57597
## 2023-09-16 21:42:46 Annotating text fragment 6661/57597
## 2023-09-16 21:42:46 Annotating text fragment 6671/57597
## 2023-09-16 21:42:46 Annotating text fragment 6681/57597
## 2023-09-16 21:42:46 Annotating text fragment 6691/57597
## 2023-09-16 21:42:46 Annotating text fragment 6701/57597
## 2023-09-16 21:42:46 Annotating text fragment 6711/57597
```

```
## 2023-09-16 21:42:46 Annotating text fragment 6721/57597
## 2023-09-16 21:42:46 Annotating text fragment 6731/57597
## 2023-09-16 21:42:46 Annotating text fragment 6741/57597
## 2023-09-16 21:42:46 Annotating text fragment 6751/57597
## 2023-09-16 21:42:46 Annotating text fragment 6761/57597
## 2023-09-16 21:42:47 Annotating text fragment 6771/57597
## 2023-09-16 21:42:47 Annotating text fragment 6781/57597
## 2023-09-16 21:42:47 Annotating text fragment 6791/57597
## 2023-09-16 21:42:47 Annotating text fragment 6801/57597
## 2023-09-16 21:42:47 Annotating text fragment 6811/57597
## 2023-09-16 21:42:47 Annotating text fragment 6821/57597
## 2023-09-16 21:42:47 Annotating text fragment 6831/57597
## 2023-09-16 21:42:47 Annotating text fragment 6841/57597
## 2023-09-16 21:42:47 Annotating text fragment 6851/57597
## 2023-09-16 21:42:47 Annotating text fragment 6861/57597
## 2023-09-16 21:42:47 Annotating text fragment 6871/57597
## 2023-09-16 21:42:47 Annotating text fragment 6881/57597
## 2023-09-16 21:42:47 Annotating text fragment 6891/57597
## 2023-09-16 21:42:47 Annotating text fragment 6901/57597
## 2023-09-16 21:42:47 Annotating text fragment 6911/57597
## 2023-09-16 21:42:47 Annotating text fragment 6921/57597
## 2023-09-16 21:42:48 Annotating text fragment 6931/57597
## 2023-09-16 21:42:48 Annotating text fragment 6941/57597
## 2023-09-16 21:42:48 Annotating text fragment 6951/57597
## 2023-09-16 21:42:48 Annotating text fragment 6961/57597
## 2023-09-16 21:42:48 Annotating text fragment 6971/57597
## 2023-09-16 21:42:49 Annotating text fragment 6981/57597
## 2023-09-16 21:42:49 Annotating text fragment 6991/57597
## 2023-09-16 21:42:49 Annotating text fragment 7001/57597
## 2023-09-16 21:42:49 Annotating text fragment 7011/57597
## 2023-09-16 21:42:49 Annotating text fragment 7021/57597
## 2023-09-16 21:42:49 Annotating text fragment 7031/57597
## 2023-09-16 21:42:50 Annotating text fragment 7041/57597
## 2023-09-16 21:42:50 Annotating text fragment 7051/57597
## 2023-09-16 21:42:50 Annotating text fragment 7061/57597
## 2023-09-16 21:42:50 Annotating text fragment 7071/57597
## 2023-09-16 21:42:50 Annotating text fragment 7081/57597
## 2023-09-16 21:42:50 Annotating text fragment 7091/57597
## 2023-09-16 21:42:50 Annotating text fragment 7101/57597
## 2023-09-16 21:42:51 Annotating text fragment 7111/57597
## 2023-09-16 21:42:51 Annotating text fragment 7121/57597
## 2023-09-16 21:42:51 Annotating text fragment 7131/57597
## 2023-09-16 21:42:51 Annotating text fragment 7141/57597
## 2023-09-16 21:42:52 Annotating text fragment 7151/57597
## 2023-09-16 21:42:52 Annotating text fragment 7161/57597
## 2023-09-16 21:42:52 Annotating text fragment 7171/57597
## 2023-09-16 21:42:52 Annotating text fragment 7181/57597
## 2023-09-16 21:42:52 Annotating text fragment 7191/57597
## 2023-09-16 21:42:52 Annotating text fragment 7201/57597
## 2023-09-16 21:42:52 Annotating text fragment 7211/57597
## 2023-09-16 21:42:52 Annotating text fragment 7221/57597
## 2023-09-16 21:42:52 Annotating text fragment 7231/57597
## 2023-09-16 21:42:52 Annotating text fragment 7241/57597
## 2023-09-16 21:42:52 Annotating text fragment 7251/57597
```

```
## 2023-09-16 21:42:52 Annotating text fragment 7261/57597
## 2023-09-16 21:42:53 Annotating text fragment 7271/57597
## 2023-09-16 21:42:53 Annotating text fragment 7281/57597
## 2023-09-16 21:42:53 Annotating text fragment 7291/57597
## 2023-09-16 21:42:53 Annotating text fragment 7301/57597
## 2023-09-16 21:42:53 Annotating text fragment 7311/57597
## 2023-09-16 21:42:53 Annotating text fragment 7321/57597
## 2023-09-16 21:42:53 Annotating text fragment 7331/57597
## 2023-09-16 21:42:53 Annotating text fragment 7341/57597
## 2023-09-16 21:42:53 Annotating text fragment 7351/57597
## 2023-09-16 21:42:54 Annotating text fragment 7361/57597
## 2023-09-16 21:42:54 Annotating text fragment 7371/57597
## 2023-09-16 21:42:54 Annotating text fragment 7381/57597
## 2023-09-16 21:42:54 Annotating text fragment 7391/57597
## 2023-09-16 21:42:54 Annotating text fragment 7401/57597
## 2023-09-16 21:42:54 Annotating text fragment 7411/57597
## 2023-09-16 21:42:54 Annotating text fragment 7421/57597
## 2023-09-16 21:42:55 Annotating text fragment 7431/57597
## 2023-09-16 21:42:55 Annotating text fragment 7441/57597
## 2023-09-16 21:42:55 Annotating text fragment 7451/57597
## 2023-09-16 21:42:56 Annotating text fragment 7461/57597
## 2023-09-16 21:42:56 Annotating text fragment 7471/57597
## 2023-09-16 21:42:56 Annotating text fragment 7481/57597
## 2023-09-16 21:42:56 Annotating text fragment 7491/57597
## 2023-09-16 21:42:57 Annotating text fragment 7501/57597
## 2023-09-16 21:42:57 Annotating text fragment 7511/57597
## 2023-09-16 21:42:57 Annotating text fragment 7521/57597
## 2023-09-16 21:42:57 Annotating text fragment 7531/57597
## 2023-09-16 21:42:57 Annotating text fragment 7541/57597
## 2023-09-16 21:42:57 Annotating text fragment 7551/57597
## 2023-09-16 21:42:57 Annotating text fragment 7561/57597
## 2023-09-16 21:42:57 Annotating text fragment 7571/57597
## 2023-09-16 21:42:57 Annotating text fragment 7581/57597
## 2023-09-16 21:42:57 Annotating text fragment 7591/57597
## 2023-09-16 21:42:58 Annotating text fragment 7601/57597
## 2023-09-16 21:42:58 Annotating text fragment 7611/57597
## 2023-09-16 21:42:58 Annotating text fragment 7621/57597
## 2023-09-16 21:42:58 Annotating text fragment 7631/57597
## 2023-09-16 21:42:58 Annotating text fragment 7641/57597
## 2023-09-16 21:42:58 Annotating text fragment 7651/57597
## 2023-09-16 21:42:58 Annotating text fragment 7661/57597
## 2023-09-16 21:42:58 Annotating text fragment 7671/57597
## 2023-09-16 21:42:59 Annotating text fragment 7681/57597
## 2023-09-16 21:42:59 Annotating text fragment 7691/57597
## 2023-09-16 21:42:59 Annotating text fragment 7701/57597
## 2023-09-16 21:42:59 Annotating text fragment 7711/57597
## 2023-09-16 21:42:59 Annotating text fragment 7721/57597
## 2023-09-16 21:42:59 Annotating text fragment 7731/57597
## 2023-09-16 21:42:59 Annotating text fragment 7741/57597
## 2023-09-16 21:42:59 Annotating text fragment 7751/57597
## 2023-09-16 21:42:59 Annotating text fragment 7761/57597
## 2023-09-16 21:43:00 Annotating text fragment 7771/57597
## 2023-09-16 21:43:00 Annotating text fragment 7781/57597
## 2023-09-16 21:43:00 Annotating text fragment 7791/57597
```

```
## 2023-09-16 21:43:00 Annotating text fragment 7801/57597
## 2023-09-16 21:43:00 Annotating text fragment 7811/57597
## 2023-09-16 21:43:00 Annotating text fragment 7821/57597
## 2023-09-16 21:43:00 Annotating text fragment 7831/57597
## 2023-09-16 21:43:00 Annotating text fragment 7841/57597
## 2023-09-16 21:43:00 Annotating text fragment 7851/57597
## 2023-09-16 21:43:00 Annotating text fragment 7861/57597
## 2023-09-16 21:43:00 Annotating text fragment 7871/57597
## 2023-09-16 21:43:00 Annotating text fragment 7881/57597
## 2023-09-16 21:43:01 Annotating text fragment 7891/57597
## 2023-09-16 21:43:01 Annotating text fragment 7901/57597
## 2023-09-16 21:43:01 Annotating text fragment 7911/57597
## 2023-09-16 21:43:01 Annotating text fragment 7921/57597
## 2023-09-16 21:43:01 Annotating text fragment 7931/57597
## 2023-09-16 21:43:01 Annotating text fragment 7941/57597
## 2023-09-16 21:43:01 Annotating text fragment 7951/57597
## 2023-09-16 21:43:01 Annotating text fragment 7961/57597
## 2023-09-16 21:43:01 Annotating text fragment 7971/57597
## 2023-09-16 21:43:01 Annotating text fragment 7981/57597
## 2023-09-16 21:43:01 Annotating text fragment 7991/57597
## 2023-09-16 21:43:01 Annotating text fragment 8001/57597
## 2023-09-16 21:43:01 Annotating text fragment 8011/57597
## 2023-09-16 21:43:01 Annotating text fragment 8021/57597
## 2023-09-16 21:43:01 Annotating text fragment 8031/57597
## 2023-09-16 21:43:01 Annotating text fragment 8041/57597
## 2023-09-16 21:43:02 Annotating text fragment 8051/57597
## 2023-09-16 21:43:02 Annotating text fragment 8061/57597
## 2023-09-16 21:43:02 Annotating text fragment 8071/57597
## 2023-09-16 21:43:02 Annotating text fragment 8081/57597
## 2023-09-16 21:43:02 Annotating text fragment 8091/57597
## 2023-09-16 21:43:02 Annotating text fragment 8101/57597
## 2023-09-16 21:43:02 Annotating text fragment 8111/57597
## 2023-09-16 21:43:02 Annotating text fragment 8121/57597
## 2023-09-16 21:43:02 Annotating text fragment 8131/57597
## 2023-09-16 21:43:02 Annotating text fragment 8141/57597
## 2023-09-16 21:43:02 Annotating text fragment 8151/57597
## 2023-09-16 21:43:03 Annotating text fragment 8161/57597
## 2023-09-16 21:43:03 Annotating text fragment 8171/57597
## 2023-09-16 21:43:03 Annotating text fragment 8181/57597
## 2023-09-16 21:43:03 Annotating text fragment 8191/57597
## 2023-09-16 21:43:03 Annotating text fragment 8201/57597
## 2023-09-16 21:43:04 Annotating text fragment 8211/57597
## 2023-09-16 21:43:04 Annotating text fragment 8221/57597
## 2023-09-16 21:43:04 Annotating text fragment 8231/57597
## 2023-09-16 21:43:04 Annotating text fragment 8241/57597
## 2023-09-16 21:43:04 Annotating text fragment 8251/57597
## 2023-09-16 21:43:04 Annotating text fragment 8261/57597
## 2023-09-16 21:43:04 Annotating text fragment 8271/57597
## 2023-09-16 21:43:04 Annotating text fragment 8281/57597
## 2023-09-16 21:43:04 Annotating text fragment 8291/57597
## 2023-09-16 21:43:04 Annotating text fragment 8301/57597
## 2023-09-16 21:43:04 Annotating text fragment 8311/57597
## 2023-09-16 21:43:04 Annotating text fragment 8321/57597
## 2023-09-16 21:43:04 Annotating text fragment 8331/57597
```

```
## 2023-09-16 21:43:04 Annotating text fragment 8341/57597
## 2023-09-16 21:43:04 Annotating text fragment 8351/57597
## 2023-09-16 21:43:04 Annotating text fragment 8361/57597
## 2023-09-16 21:43:04 Annotating text fragment 8371/57597
## 2023-09-16 21:43:04 Annotating text fragment 8381/57597
## 2023-09-16 21:43:04 Annotating text fragment 8391/57597
## 2023-09-16 21:43:04 Annotating text fragment 8401/57597
## 2023-09-16 21:43:04 Annotating text fragment 8411/57597
## 2023-09-16 21:43:04 Annotating text fragment 8421/57597
## 2023-09-16 21:43:04 Annotating text fragment 8431/57597
## 2023-09-16 21:43:04 Annotating text fragment 8441/57597
## 2023-09-16 21:43:04 Annotating text fragment 8451/57597
## 2023-09-16 21:43:04 Annotating text fragment 8461/57597
## 2023-09-16 21:43:05 Annotating text fragment 8471/57597
## 2023-09-16 21:43:05 Annotating text fragment 8481/57597
## 2023-09-16 21:43:05 Annotating text fragment 8491/57597
## 2023-09-16 21:43:05 Annotating text fragment 8501/57597
## 2023-09-16 21:43:05 Annotating text fragment 8511/57597
## 2023-09-16 21:43:05 Annotating text fragment 8521/57597
## 2023-09-16 21:43:05 Annotating text fragment 8531/57597
## 2023-09-16 21:43:05 Annotating text fragment 8541/57597
## 2023-09-16 21:43:05 Annotating text fragment 8551/57597
## 2023-09-16 21:43:05 Annotating text fragment 8561/57597
## 2023-09-16 21:43:05 Annotating text fragment 8571/57597
## 2023-09-16 21:43:05 Annotating text fragment 8581/57597
## 2023-09-16 21:43:05 Annotating text fragment 8591/57597
## 2023-09-16 21:43:05 Annotating text fragment 8601/57597
## 2023-09-16 21:43:05 Annotating text fragment 8611/57597
## 2023-09-16 21:43:05 Annotating text fragment 8621/57597
## 2023-09-16 21:43:05 Annotating text fragment 8631/57597
## 2023-09-16 21:43:05 Annotating text fragment 8641/57597
## 2023-09-16 21:43:06 Annotating text fragment 8651/57597
## 2023-09-16 21:43:06 Annotating text fragment 8661/57597
## 2023-09-16 21:43:06 Annotating text fragment 8671/57597
## 2023-09-16 21:43:06 Annotating text fragment 8681/57597
## 2023-09-16 21:43:06 Annotating text fragment 8691/57597
## 2023-09-16 21:43:06 Annotating text fragment 8701/57597
## 2023-09-16 21:43:06 Annotating text fragment 8711/57597
## 2023-09-16 21:43:06 Annotating text fragment 8721/57597
## 2023-09-16 21:43:06 Annotating text fragment 8731/57597
## 2023-09-16 21:43:06 Annotating text fragment 8741/57597
## 2023-09-16 21:43:06 Annotating text fragment 8751/57597
## 2023-09-16 21:43:06 Annotating text fragment 8761/57597
## 2023-09-16 21:43:07 Annotating text fragment 8771/57597
## 2023-09-16 21:43:07 Annotating text fragment 8781/57597
## 2023-09-16 21:43:07 Annotating text fragment 8791/57597
## 2023-09-16 21:43:07 Annotating text fragment 8801/57597
## 2023-09-16 21:43:07 Annotating text fragment 8811/57597
## 2023-09-16 21:43:07 Annotating text fragment 8821/57597
## 2023-09-16 21:43:07 Annotating text fragment 8831/57597
## 2023-09-16 21:43:07 Annotating text fragment 8841/57597
## 2023-09-16 21:43:07 Annotating text fragment 8851/57597
## 2023-09-16 21:43:07 Annotating text fragment 8861/57597
## 2023-09-16 21:43:07 Annotating text fragment 8871/57597
```

```
## 2023-09-16 21:43:07 Annotating text fragment 8881/57597
## 2023-09-16 21:43:07 Annotating text fragment 8891/57597
## 2023-09-16 21:43:07 Annotating text fragment 8901/57597
## 2023-09-16 21:43:07 Annotating text fragment 8911/57597
## 2023-09-16 21:43:07 Annotating text fragment 8921/57597
## 2023-09-16 21:43:08 Annotating text fragment 8931/57597
## 2023-09-16 21:43:08 Annotating text fragment 8941/57597
## 2023-09-16 21:43:08 Annotating text fragment 8951/57597
## 2023-09-16 21:43:08 Annotating text fragment 8961/57597
## 2023-09-16 21:43:08 Annotating text fragment 8971/57597
## 2023-09-16 21:43:08 Annotating text fragment 8981/57597
## 2023-09-16 21:43:08 Annotating text fragment 8991/57597
## 2023-09-16 21:43:08 Annotating text fragment 9001/57597
## 2023-09-16 21:43:08 Annotating text fragment 9011/57597
## 2023-09-16 21:43:08 Annotating text fragment 9021/57597
## 2023-09-16 21:43:08 Annotating text fragment 9031/57597
## 2023-09-16 21:43:08 Annotating text fragment 9041/57597
## 2023-09-16 21:43:08 Annotating text fragment 9051/57597
## 2023-09-16 21:43:08 Annotating text fragment 9061/57597
## 2023-09-16 21:43:08 Annotating text fragment 9071/57597
## 2023-09-16 21:43:08 Annotating text fragment 9081/57597
## 2023-09-16 21:43:08 Annotating text fragment 9091/57597
## 2023-09-16 21:43:08 Annotating text fragment 9101/57597
## 2023-09-16 21:43:08 Annotating text fragment 9111/57597
## 2023-09-16 21:43:09 Annotating text fragment 9121/57597
## 2023-09-16 21:43:09 Annotating text fragment 9131/57597
## 2023-09-16 21:43:09 Annotating text fragment 9141/57597
## 2023-09-16 21:43:09 Annotating text fragment 9151/57597
## 2023-09-16 21:43:09 Annotating text fragment 9161/57597
## 2023-09-16 21:43:09 Annotating text fragment 9171/57597
## 2023-09-16 21:43:09 Annotating text fragment 9181/57597
## 2023-09-16 21:43:09 Annotating text fragment 9191/57597
## 2023-09-16 21:43:09 Annotating text fragment 9201/57597
## 2023-09-16 21:43:09 Annotating text fragment 9211/57597
## 2023-09-16 21:43:09 Annotating text fragment 9221/57597
## 2023-09-16 21:43:09 Annotating text fragment 9231/57597
## 2023-09-16 21:43:09 Annotating text fragment 9241/57597
## 2023-09-16 21:43:09 Annotating text fragment 9251/57597
## 2023-09-16 21:43:09 Annotating text fragment 9261/57597
## 2023-09-16 21:43:09 Annotating text fragment 9271/57597
## 2023-09-16 21:43:09 Annotating text fragment 9281/57597
## 2023-09-16 21:43:10 Annotating text fragment 9291/57597
## 2023-09-16 21:43:10 Annotating text fragment 9301/57597
## 2023-09-16 21:43:10 Annotating text fragment 9311/57597
## 2023-09-16 21:43:10 Annotating text fragment 9321/57597
## 2023-09-16 21:43:10 Annotating text fragment 9331/57597
## 2023-09-16 21:43:10 Annotating text fragment 9341/57597
## 2023-09-16 21:43:10 Annotating text fragment 9351/57597
## 2023-09-16 21:43:10 Annotating text fragment 9361/57597
## 2023-09-16 21:43:10 Annotating text fragment 9371/57597
## 2023-09-16 21:43:10 Annotating text fragment 9381/57597
## 2023-09-16 21:43:10 Annotating text fragment 9391/57597
## 2023-09-16 21:43:10 Annotating text fragment 9401/57597
## 2023-09-16 21:43:10 Annotating text fragment 9411/57597
```

```
## 2023-09-16 21:43:10 Annotating text fragment 9421/57597
## 2023-09-16 21:43:11 Annotating text fragment 9431/57597
## 2023-09-16 21:43:11 Annotating text fragment 9441/57597
## 2023-09-16 21:43:11 Annotating text fragment 9451/57597
## 2023-09-16 21:43:11 Annotating text fragment 9461/57597
## 2023-09-16 21:43:11 Annotating text fragment 9471/57597
## 2023-09-16 21:43:11 Annotating text fragment 9481/57597
## 2023-09-16 21:43:11 Annotating text fragment 9491/57597
## 2023-09-16 21:43:11 Annotating text fragment 9501/57597
## 2023-09-16 21:43:11 Annotating text fragment 9511/57597
## 2023-09-16 21:43:11 Annotating text fragment 9521/57597
## 2023-09-16 21:43:11 Annotating text fragment 9531/57597
## 2023-09-16 21:43:12 Annotating text fragment 9541/57597
## 2023-09-16 21:43:12 Annotating text fragment 9551/57597
## 2023-09-16 21:43:12 Annotating text fragment 9561/57597
## 2023-09-16 21:43:12 Annotating text fragment 9571/57597
## 2023-09-16 21:43:12 Annotating text fragment 9581/57597
## 2023-09-16 21:43:12 Annotating text fragment 9591/57597
## 2023-09-16 21:43:12 Annotating text fragment 9601/57597
## 2023-09-16 21:43:12 Annotating text fragment 9611/57597
## 2023-09-16 21:43:12 Annotating text fragment 9621/57597
## 2023-09-16 21:43:12 Annotating text fragment 9631/57597
## 2023-09-16 21:43:12 Annotating text fragment 9641/57597
## 2023-09-16 21:43:12 Annotating text fragment 9651/57597
## 2023-09-16 21:43:12 Annotating text fragment 9661/57597
## 2023-09-16 21:43:12 Annotating text fragment 9671/57597
## 2023-09-16 21:43:12 Annotating text fragment 9681/57597
## 2023-09-16 21:43:12 Annotating text fragment 9691/57597
## 2023-09-16 21:43:12 Annotating text fragment 9701/57597
## 2023-09-16 21:43:12 Annotating text fragment 9711/57597
## 2023-09-16 21:43:12 Annotating text fragment 9721/57597
## 2023-09-16 21:43:12 Annotating text fragment 9731/57597
## 2023-09-16 21:43:12 Annotating text fragment 9741/57597
## 2023-09-16 21:43:13 Annotating text fragment 9751/57597
## 2023-09-16 21:43:13 Annotating text fragment 9761/57597
## 2023-09-16 21:43:13 Annotating text fragment 9771/57597
## 2023-09-16 21:43:13 Annotating text fragment 9781/57597
## 2023-09-16 21:43:13 Annotating text fragment 9791/57597
## 2023-09-16 21:43:13 Annotating text fragment 9801/57597
## 2023-09-16 21:43:13 Annotating text fragment 9811/57597
## 2023-09-16 21:43:13 Annotating text fragment 9821/57597
## 2023-09-16 21:43:13 Annotating text fragment 9831/57597
## 2023-09-16 21:43:13 Annotating text fragment 9841/57597
## 2023-09-16 21:43:14 Annotating text fragment 9851/57597
## 2023-09-16 21:43:14 Annotating text fragment 9861/57597
## 2023-09-16 21:43:14 Annotating text fragment 9871/57597
## 2023-09-16 21:43:14 Annotating text fragment 9881/57597
## 2023-09-16 21:43:15 Annotating text fragment 9891/57597
## 2023-09-16 21:43:15 Annotating text fragment 9901/57597
## 2023-09-16 21:43:15 Annotating text fragment 9911/57597
## 2023-09-16 21:43:16 Annotating text fragment 9921/57597
## 2023-09-16 21:43:16 Annotating text fragment 9931/57597
## 2023-09-16 21:43:16 Annotating text fragment 9941/57597
## 2023-09-16 21:43:16 Annotating text fragment 9951/57597
```

```
## 2023-09-16 21:43:16 Annotating text fragment 9961/57597
## 2023-09-16 21:43:16 Annotating text fragment 9971/57597
## 2023-09-16 21:43:16 Annotating text fragment 9981/57597
## 2023-09-16 21:43:16 Annotating text fragment 9991/57597
## 2023-09-16 21:43:16 Annotating text fragment 10001/57597
## 2023-09-16 21:43:17 Annotating text fragment 10011/57597
## 2023-09-16 21:43:17 Annotating text fragment 10021/57597
## 2023-09-16 21:43:17 Annotating text fragment 10031/57597
## 2023-09-16 21:43:17 Annotating text fragment 10041/57597
## 2023-09-16 21:43:17 Annotating text fragment 10051/57597
## 2023-09-16 21:43:17 Annotating text fragment 10061/57597
## 2023-09-16 21:43:17 Annotating text fragment 10071/57597
## 2023-09-16 21:43:17 Annotating text fragment 10081/57597
## 2023-09-16 21:43:17 Annotating text fragment 10091/57597
## 2023-09-16 21:43:17 Annotating text fragment 10101/57597
## 2023-09-16 21:43:17 Annotating text fragment 10111/57597
## 2023-09-16 21:43:17 Annotating text fragment 10121/57597
## 2023-09-16 21:43:17 Annotating text fragment 10131/57597
## 2023-09-16 21:43:17 Annotating text fragment 10141/57597
## 2023-09-16 21:43:17 Annotating text fragment 10151/57597
## 2023-09-16 21:43:17 Annotating text fragment 10161/57597
## 2023-09-16 21:43:18 Annotating text fragment 10171/57597
## 2023-09-16 21:43:18 Annotating text fragment 10181/57597
## 2023-09-16 21:43:18 Annotating text fragment 10191/57597
## 2023-09-16 21:43:18 Annotating text fragment 10201/57597
## 2023-09-16 21:43:18 Annotating text fragment 10211/57597
## 2023-09-16 21:43:18 Annotating text fragment 10221/57597
## 2023-09-16 21:43:18 Annotating text fragment 10231/57597
## 2023-09-16 21:43:18 Annotating text fragment 10241/57597
## 2023-09-16 21:43:18 Annotating text fragment 10251/57597
## 2023-09-16 21:43:18 Annotating text fragment 10261/57597
## 2023-09-16 21:43:18 Annotating text fragment 10271/57597
## 2023-09-16 21:43:19 Annotating text fragment 10281/57597
## 2023-09-16 21:43:19 Annotating text fragment 10291/57597
## 2023-09-16 21:43:19 Annotating text fragment 10301/57597
## 2023-09-16 21:43:19 Annotating text fragment 10311/57597
## 2023-09-16 21:43:19 Annotating text fragment 10321/57597
## 2023-09-16 21:43:19 Annotating text fragment 10331/57597
## 2023-09-16 21:43:19 Annotating text fragment 10341/57597
## 2023-09-16 21:43:19 Annotating text fragment 10351/57597
## 2023-09-16 21:43:19 Annotating text fragment 10361/57597
## 2023-09-16 21:43:19 Annotating text fragment 10371/57597
## 2023-09-16 21:43:19 Annotating text fragment 10381/57597
## 2023-09-16 21:43:19 Annotating text fragment 10391/57597
## 2023-09-16 21:43:19 Annotating text fragment 10401/57597
## 2023-09-16 21:43:19 Annotating text fragment 10411/57597
## 2023-09-16 21:43:19 Annotating text fragment 10421/57597
## 2023-09-16 21:43:19 Annotating text fragment 10431/57597
## 2023-09-16 21:43:19 Annotating text fragment 10441/57597
## 2023-09-16 21:43:19 Annotating text fragment 10451/57597
## 2023-09-16 21:43:19 Annotating text fragment 10461/57597
## 2023-09-16 21:43:19 Annotating text fragment 10471/57597
## 2023-09-16 21:43:19 Annotating text fragment 10481/57597
## 2023-09-16 21:43:19 Annotating text fragment 10491/57597
```

```
## 2023-09-16 21:43:19 Annotating text fragment 10501/57597
## 2023-09-16 21:43:19 Annotating text fragment 10511/57597
## 2023-09-16 21:43:19 Annotating text fragment 10521/57597
## 2023-09-16 21:43:20 Annotating text fragment 10531/57597
## 2023-09-16 21:43:20 Annotating text fragment 10541/57597
## 2023-09-16 21:43:20 Annotating text fragment 10551/57597
## 2023-09-16 21:43:20 Annotating text fragment 10561/57597
## 2023-09-16 21:43:20 Annotating text fragment 10571/57597
## 2023-09-16 21:43:20 Annotating text fragment 10581/57597
## 2023-09-16 21:43:20 Annotating text fragment 10591/57597
## 2023-09-16 21:43:20 Annotating text fragment 10601/57597
## 2023-09-16 21:43:20 Annotating text fragment 10611/57597
## 2023-09-16 21:43:20 Annotating text fragment 10621/57597
## 2023-09-16 21:43:20 Annotating text fragment 10631/57597
## 2023-09-16 21:43:20 Annotating text fragment 10641/57597
## 2023-09-16 21:43:20 Annotating text fragment 10651/57597
## 2023-09-16 21:43:20 Annotating text fragment 10661/57597
## 2023-09-16 21:43:20 Annotating text fragment 10671/57597
## 2023-09-16 21:43:20 Annotating text fragment 10681/57597
## 2023-09-16 21:43:20 Annotating text fragment 10691/57597
## 2023-09-16 21:43:21 Annotating text fragment 10701/57597
## 2023-09-16 21:43:21 Annotating text fragment 10711/57597
## 2023-09-16 21:43:21 Annotating text fragment 10721/57597
## 2023-09-16 21:43:21 Annotating text fragment 10731/57597
## 2023-09-16 21:43:21 Annotating text fragment 10741/57597
## 2023-09-16 21:43:21 Annotating text fragment 10751/57597
## 2023-09-16 21:43:21 Annotating text fragment 10761/57597
## 2023-09-16 21:43:21 Annotating text fragment 10771/57597
## 2023-09-16 21:43:21 Annotating text fragment 10781/57597
## 2023-09-16 21:43:21 Annotating text fragment 10791/57597
## 2023-09-16 21:43:21 Annotating text fragment 10801/57597
## 2023-09-16 21:43:21 Annotating text fragment 10811/57597
## 2023-09-16 21:43:21 Annotating text fragment 10821/57597
## 2023-09-16 21:43:21 Annotating text fragment 10831/57597
## 2023-09-16 21:43:21 Annotating text fragment 10841/57597
## 2023-09-16 21:43:21 Annotating text fragment 10851/57597
## 2023-09-16 21:43:21 Annotating text fragment 10861/57597
## 2023-09-16 21:43:21 Annotating text fragment 10871/57597
## 2023-09-16 21:43:22 Annotating text fragment 10881/57597
## 2023-09-16 21:43:22 Annotating text fragment 10891/57597
## 2023-09-16 21:43:22 Annotating text fragment 10901/57597
## 2023-09-16 21:43:22 Annotating text fragment 10911/57597
## 2023-09-16 21:43:22 Annotating text fragment 10921/57597
## 2023-09-16 21:43:22 Annotating text fragment 10931/57597
## 2023-09-16 21:43:22 Annotating text fragment 10941/57597
## 2023-09-16 21:43:22 Annotating text fragment 10951/57597
## 2023-09-16 21:43:22 Annotating text fragment 10961/57597
## 2023-09-16 21:43:22 Annotating text fragment 10971/57597
## 2023-09-16 21:43:22 Annotating text fragment 10981/57597
## 2023-09-16 21:43:22 Annotating text fragment 10991/57597
## 2023-09-16 21:43:22 Annotating text fragment 11001/57597
## 2023-09-16 21:43:22 Annotating text fragment 11011/57597
## 2023-09-16 21:43:23 Annotating text fragment 11021/57597
## 2023-09-16 21:43:23 Annotating text fragment 11031/57597
```

```
## 2023-09-16 21:43:23 Annotating text fragment 11041/57597
## 2023-09-16 21:43:23 Annotating text fragment 11051/57597
## 2023-09-16 21:43:23 Annotating text fragment 11061/57597
## 2023-09-16 21:43:23 Annotating text fragment 11071/57597
## 2023-09-16 21:43:23 Annotating text fragment 11081/57597
## 2023-09-16 21:43:23 Annotating text fragment 11091/57597
## 2023-09-16 21:43:23 Annotating text fragment 11101/57597
## 2023-09-16 21:43:23 Annotating text fragment 11111/57597
## 2023-09-16 21:43:24 Annotating text fragment 11121/57597
## 2023-09-16 21:43:24 Annotating text fragment 11131/57597
## 2023-09-16 21:43:24 Annotating text fragment 11141/57597
## 2023-09-16 21:43:25 Annotating text fragment 11151/57597
## 2023-09-16 21:43:25 Annotating text fragment 11161/57597
## 2023-09-16 21:43:25 Annotating text fragment 11171/57597
## 2023-09-16 21:43:25 Annotating text fragment 11181/57597
## 2023-09-16 21:43:26 Annotating text fragment 11191/57597
## 2023-09-16 21:43:26 Annotating text fragment 11201/57597
## 2023-09-16 21:43:26 Annotating text fragment 11211/57597
## 2023-09-16 21:43:27 Annotating text fragment 11221/57597
## 2023-09-16 21:43:27 Annotating text fragment 11231/57597
## 2023-09-16 21:43:27 Annotating text fragment 11241/57597
## 2023-09-16 21:43:27 Annotating text fragment 11251/57597
## 2023-09-16 21:43:27 Annotating text fragment 11261/57597
## 2023-09-16 21:43:27 Annotating text fragment 11271/57597
## 2023-09-16 21:43:27 Annotating text fragment 11281/57597
## 2023-09-16 21:43:27 Annotating text fragment 11291/57597
## 2023-09-16 21:43:27 Annotating text fragment 11301/57597
## 2023-09-16 21:43:27 Annotating text fragment 11311/57597
## 2023-09-16 21:43:27 Annotating text fragment 11321/57597
## 2023-09-16 21:43:27 Annotating text fragment 11331/57597
## 2023-09-16 21:43:27 Annotating text fragment 11341/57597
## 2023-09-16 21:43:27 Annotating text fragment 11351/57597
## 2023-09-16 21:43:27 Annotating text fragment 11361/57597
## 2023-09-16 21:43:27 Annotating text fragment 11371/57597
## 2023-09-16 21:43:27 Annotating text fragment 11381/57597
## 2023-09-16 21:43:27 Annotating text fragment 11391/57597
## 2023-09-16 21:43:27 Annotating text fragment 11401/57597
## 2023-09-16 21:43:28 Annotating text fragment 11411/57597
## 2023-09-16 21:43:28 Annotating text fragment 11421/57597
## 2023-09-16 21:43:28 Annotating text fragment 11431/57597
## 2023-09-16 21:43:28 Annotating text fragment 11441/57597
## 2023-09-16 21:43:28 Annotating text fragment 11451/57597
## 2023-09-16 21:43:28 Annotating text fragment 11461/57597
## 2023-09-16 21:43:28 Annotating text fragment 11471/57597
## 2023-09-16 21:43:28 Annotating text fragment 11481/57597
## 2023-09-16 21:43:28 Annotating text fragment 11491/57597
## 2023-09-16 21:43:28 Annotating text fragment 11501/57597
## 2023-09-16 21:43:28 Annotating text fragment 11511/57597
## 2023-09-16 21:43:28 Annotating text fragment 11521/57597
## 2023-09-16 21:43:28 Annotating text fragment 11531/57597
## 2023-09-16 21:43:28 Annotating text fragment 11541/57597
## 2023-09-16 21:43:28 Annotating text fragment 11551/57597
## 2023-09-16 21:43:29 Annotating text fragment 11561/57597
## 2023-09-16 21:43:29 Annotating text fragment 11571/57597
```

```
## 2023-09-16 21:43:29 Annotating text fragment 11581/57597
## 2023-09-16 21:43:29 Annotating text fragment 11591/57597
## 2023-09-16 21:43:29 Annotating text fragment 11601/57597
## 2023-09-16 21:43:29 Annotating text fragment 11611/57597
## 2023-09-16 21:43:29 Annotating text fragment 11621/57597
## 2023-09-16 21:43:29 Annotating text fragment 11631/57597
## 2023-09-16 21:43:29 Annotating text fragment 11641/57597
## 2023-09-16 21:43:29 Annotating text fragment 11651/57597
## 2023-09-16 21:43:29 Annotating text fragment 11661/57597
## 2023-09-16 21:43:29 Annotating text fragment 11671/57597
## 2023-09-16 21:43:29 Annotating text fragment 11681/57597
## 2023-09-16 21:43:29 Annotating text fragment 11691/57597
## 2023-09-16 21:43:30 Annotating text fragment 11701/57597
## 2023-09-16 21:43:30 Annotating text fragment 11711/57597
## 2023-09-16 21:43:30 Annotating text fragment 11721/57597
## 2023-09-16 21:43:30 Annotating text fragment 11731/57597
## 2023-09-16 21:43:30 Annotating text fragment 11741/57597
## 2023-09-16 21:43:30 Annotating text fragment 11751/57597
## 2023-09-16 21:43:30 Annotating text fragment 11761/57597
## 2023-09-16 21:43:30 Annotating text fragment 11771/57597
## 2023-09-16 21:43:30 Annotating text fragment 11781/57597
## 2023-09-16 21:43:30 Annotating text fragment 11791/57597
## 2023-09-16 21:43:30 Annotating text fragment 11801/57597
## 2023-09-16 21:43:30 Annotating text fragment 11811/57597
## 2023-09-16 21:43:30 Annotating text fragment 11821/57597
## 2023-09-16 21:43:30 Annotating text fragment 11831/57597
## 2023-09-16 21:43:30 Annotating text fragment 11841/57597
## 2023-09-16 21:43:30 Annotating text fragment 11851/57597
## 2023-09-16 21:43:30 Annotating text fragment 11861/57597
## 2023-09-16 21:43:30 Annotating text fragment 11871/57597
## 2023-09-16 21:43:31 Annotating text fragment 11881/57597
## 2023-09-16 21:43:31 Annotating text fragment 11891/57597
## 2023-09-16 21:43:31 Annotating text fragment 11901/57597
## 2023-09-16 21:43:31 Annotating text fragment 11911/57597
## 2023-09-16 21:43:31 Annotating text fragment 11921/57597
## 2023-09-16 21:43:31 Annotating text fragment 11931/57597
## 2023-09-16 21:43:31 Annotating text fragment 11941/57597
## 2023-09-16 21:43:31 Annotating text fragment 11951/57597
## 2023-09-16 21:43:31 Annotating text fragment 11961/57597
## 2023-09-16 21:43:31 Annotating text fragment 11971/57597
## 2023-09-16 21:43:31 Annotating text fragment 11981/57597
## 2023-09-16 21:43:31 Annotating text fragment 11991/57597
## 2023-09-16 21:43:32 Annotating text fragment 12001/57597
## 2023-09-16 21:43:32 Annotating text fragment 12011/57597
## 2023-09-16 21:43:32 Annotating text fragment 12021/57597
## 2023-09-16 21:43:32 Annotating text fragment 12031/57597
## 2023-09-16 21:43:32 Annotating text fragment 12041/57597
## 2023-09-16 21:43:32 Annotating text fragment 12051/57597
## 2023-09-16 21:43:32 Annotating text fragment 12061/57597
## 2023-09-16 21:43:32 Annotating text fragment 12071/57597
## 2023-09-16 21:43:32 Annotating text fragment 12081/57597
## 2023-09-16 21:43:32 Annotating text fragment 12091/57597
## 2023-09-16 21:43:32 Annotating text fragment 12101/57597
## 2023-09-16 21:43:32 Annotating text fragment 12111/57597
```

```
## 2023-09-16 21:43:33 Annotating text fragment 12121/57597
## 2023-09-16 21:43:33 Annotating text fragment 12131/57597
## 2023-09-16 21:43:33 Annotating text fragment 12141/57597
## 2023-09-16 21:43:33 Annotating text fragment 12151/57597
## 2023-09-16 21:43:33 Annotating text fragment 12161/57597
## 2023-09-16 21:43:33 Annotating text fragment 12171/57597
## 2023-09-16 21:43:33 Annotating text fragment 12181/57597
## 2023-09-16 21:43:33 Annotating text fragment 12191/57597
## 2023-09-16 21:43:33 Annotating text fragment 12201/57597
## 2023-09-16 21:43:33 Annotating text fragment 12211/57597
## 2023-09-16 21:43:33 Annotating text fragment 12221/57597
## 2023-09-16 21:43:33 Annotating text fragment 12231/57597
## 2023-09-16 21:43:33 Annotating text fragment 12241/57597
## 2023-09-16 21:43:33 Annotating text fragment 12251/57597
## 2023-09-16 21:43:33 Annotating text fragment 12261/57597
## 2023-09-16 21:43:33 Annotating text fragment 12271/57597
## 2023-09-16 21:43:33 Annotating text fragment 12281/57597
## 2023-09-16 21:43:33 Annotating text fragment 12291/57597
## 2023-09-16 21:43:33 Annotating text fragment 12301/57597
## 2023-09-16 21:43:34 Annotating text fragment 12311/57597
## 2023-09-16 21:43:34 Annotating text fragment 12321/57597
## 2023-09-16 21:43:34 Annotating text fragment 12331/57597
## 2023-09-16 21:43:34 Annotating text fragment 12341/57597
## 2023-09-16 21:43:34 Annotating text fragment 12351/57597
## 2023-09-16 21:43:34 Annotating text fragment 12361/57597
## 2023-09-16 21:43:34 Annotating text fragment 12371/57597
## 2023-09-16 21:43:34 Annotating text fragment 12381/57597
## 2023-09-16 21:43:34 Annotating text fragment 12391/57597
## 2023-09-16 21:43:34 Annotating text fragment 12401/57597
## 2023-09-16 21:43:34 Annotating text fragment 12411/57597
## 2023-09-16 21:43:35 Annotating text fragment 12421/57597
## 2023-09-16 21:43:35 Annotating text fragment 12431/57597
## 2023-09-16 21:43:35 Annotating text fragment 12441/57597
## 2023-09-16 21:43:35 Annotating text fragment 12451/57597
## 2023-09-16 21:43:35 Annotating text fragment 12461/57597
## 2023-09-16 21:43:35 Annotating text fragment 12471/57597
## 2023-09-16 21:43:35 Annotating text fragment 12481/57597
## 2023-09-16 21:43:35 Annotating text fragment 12491/57597
## 2023-09-16 21:43:35 Annotating text fragment 12501/57597
## 2023-09-16 21:43:35 Annotating text fragment 12511/57597
## 2023-09-16 21:43:35 Annotating text fragment 12521/57597
## 2023-09-16 21:43:35 Annotating text fragment 12531/57597
## 2023-09-16 21:43:35 Annotating text fragment 12541/57597
## 2023-09-16 21:43:35 Annotating text fragment 12551/57597
## 2023-09-16 21:43:35 Annotating text fragment 12561/57597
## 2023-09-16 21:43:36 Annotating text fragment 12571/57597
## 2023-09-16 21:43:36 Annotating text fragment 12581/57597
## 2023-09-16 21:43:36 Annotating text fragment 12591/57597
## 2023-09-16 21:43:36 Annotating text fragment 12601/57597
## 2023-09-16 21:43:36 Annotating text fragment 12611/57597
## 2023-09-16 21:43:36 Annotating text fragment 12621/57597
## 2023-09-16 21:43:36 Annotating text fragment 12631/57597
## 2023-09-16 21:43:36 Annotating text fragment 12641/57597
## 2023-09-16 21:43:36 Annotating text fragment 12651/57597
```

```
## 2023-09-16 21:43:37 Annotating text fragment 12661/57597
## 2023-09-16 21:43:37 Annotating text fragment 12671/57597
## 2023-09-16 21:43:37 Annotating text fragment 12681/57597
## 2023-09-16 21:43:38 Annotating text fragment 12691/57597
## 2023-09-16 21:43:38 Annotating text fragment 12701/57597
## 2023-09-16 21:43:38 Annotating text fragment 12711/57597
## 2023-09-16 21:43:38 Annotating text fragment 12721/57597
## 2023-09-16 21:43:38 Annotating text fragment 12731/57597
## 2023-09-16 21:43:38 Annotating text fragment 12741/57597
## 2023-09-16 21:43:38 Annotating text fragment 12751/57597
## 2023-09-16 21:43:38 Annotating text fragment 12761/57597
## 2023-09-16 21:43:39 Annotating text fragment 12771/57597
## 2023-09-16 21:43:39 Annotating text fragment 12781/57597
## 2023-09-16 21:43:39 Annotating text fragment 12791/57597
## 2023-09-16 21:43:39 Annotating text fragment 12801/57597
## 2023-09-16 21:43:39 Annotating text fragment 12811/57597
## 2023-09-16 21:43:39 Annotating text fragment 12821/57597
## 2023-09-16 21:43:39 Annotating text fragment 12831/57597
## 2023-09-16 21:43:39 Annotating text fragment 12841/57597
## 2023-09-16 21:43:39 Annotating text fragment 12851/57597
## 2023-09-16 21:43:39 Annotating text fragment 12861/57597
## 2023-09-16 21:43:39 Annotating text fragment 12871/57597
## 2023-09-16 21:43:39 Annotating text fragment 12881/57597
## 2023-09-16 21:43:39 Annotating text fragment 12891/57597
## 2023-09-16 21:43:39 Annotating text fragment 12901/57597
## 2023-09-16 21:43:39 Annotating text fragment 12911/57597
## 2023-09-16 21:43:40 Annotating text fragment 12921/57597
## 2023-09-16 21:43:40 Annotating text fragment 12931/57597
## 2023-09-16 21:43:40 Annotating text fragment 12941/57597
## 2023-09-16 21:43:40 Annotating text fragment 12951/57597
## 2023-09-16 21:43:40 Annotating text fragment 12961/57597
## 2023-09-16 21:43:40 Annotating text fragment 12971/57597
## 2023-09-16 21:43:40 Annotating text fragment 12981/57597
## 2023-09-16 21:43:40 Annotating text fragment 12991/57597
## 2023-09-16 21:43:40 Annotating text fragment 13001/57597
## 2023-09-16 21:43:40 Annotating text fragment 13011/57597
## 2023-09-16 21:43:40 Annotating text fragment 13021/57597
## 2023-09-16 21:43:40 Annotating text fragment 13031/57597
## 2023-09-16 21:43:40 Annotating text fragment 13041/57597
## 2023-09-16 21:43:40 Annotating text fragment 13051/57597
## 2023-09-16 21:43:41 Annotating text fragment 13061/57597
## 2023-09-16 21:43:41 Annotating text fragment 13071/57597
## 2023-09-16 21:43:41 Annotating text fragment 13081/57597
## 2023-09-16 21:43:41 Annotating text fragment 13091/57597
## 2023-09-16 21:43:41 Annotating text fragment 13101/57597
## 2023-09-16 21:43:41 Annotating text fragment 13111/57597
## 2023-09-16 21:43:41 Annotating text fragment 13121/57597
## 2023-09-16 21:43:41 Annotating text fragment 13131/57597
## 2023-09-16 21:43:42 Annotating text fragment 13141/57597
## 2023-09-16 21:43:42 Annotating text fragment 13151/57597
## 2023-09-16 21:43:42 Annotating text fragment 13161/57597
## 2023-09-16 21:43:42 Annotating text fragment 13171/57597
## 2023-09-16 21:43:42 Annotating text fragment 13181/57597
## 2023-09-16 21:43:42 Annotating text fragment 13191/57597
```

```
## 2023-09-16 21:43:42 Annotating text fragment 13201/57597
## 2023-09-16 21:43:42 Annotating text fragment 13211/57597
## 2023-09-16 21:43:42 Annotating text fragment 13221/57597
## 2023-09-16 21:43:42 Annotating text fragment 13231/57597
## 2023-09-16 21:43:42 Annotating text fragment 13241/57597
## 2023-09-16 21:43:43 Annotating text fragment 13251/57597
## 2023-09-16 21:43:43 Annotating text fragment 13261/57597
## 2023-09-16 21:43:43 Annotating text fragment 13271/57597
## 2023-09-16 21:43:43 Annotating text fragment 13281/57597
## 2023-09-16 21:43:43 Annotating text fragment 13291/57597
## 2023-09-16 21:43:43 Annotating text fragment 13301/57597
## 2023-09-16 21:43:43 Annotating text fragment 13311/57597
## 2023-09-16 21:43:43 Annotating text fragment 13321/57597
## 2023-09-16 21:43:43 Annotating text fragment 13331/57597
## 2023-09-16 21:43:43 Annotating text fragment 13341/57597
## 2023-09-16 21:43:43 Annotating text fragment 13351/57597
## 2023-09-16 21:43:43 Annotating text fragment 13361/57597
## 2023-09-16 21:43:43 Annotating text fragment 13371/57597
## 2023-09-16 21:43:43 Annotating text fragment 13381/57597
## 2023-09-16 21:43:43 Annotating text fragment 13391/57597
## 2023-09-16 21:43:44 Annotating text fragment 13401/57597
## 2023-09-16 21:43:44 Annotating text fragment 13411/57597
## 2023-09-16 21:43:44 Annotating text fragment 13421/57597
## 2023-09-16 21:43:44 Annotating text fragment 13431/57597
## 2023-09-16 21:43:44 Annotating text fragment 13441/57597
## 2023-09-16 21:43:44 Annotating text fragment 13451/57597
## 2023-09-16 21:43:44 Annotating text fragment 13461/57597
## 2023-09-16 21:43:44 Annotating text fragment 13471/57597
## 2023-09-16 21:43:44 Annotating text fragment 13481/57597
## 2023-09-16 21:43:44 Annotating text fragment 13491/57597
## 2023-09-16 21:43:44 Annotating text fragment 13501/57597
## 2023-09-16 21:43:45 Annotating text fragment 13511/57597
## 2023-09-16 21:43:45 Annotating text fragment 13521/57597
## 2023-09-16 21:43:45 Annotating text fragment 13531/57597
## 2023-09-16 21:43:45 Annotating text fragment 13541/57597
## 2023-09-16 21:43:45 Annotating text fragment 13551/57597
## 2023-09-16 21:43:45 Annotating text fragment 13561/57597
## 2023-09-16 21:43:45 Annotating text fragment 13571/57597
## 2023-09-16 21:43:45 Annotating text fragment 13581/57597
## 2023-09-16 21:43:45 Annotating text fragment 13591/57597
## 2023-09-16 21:43:45 Annotating text fragment 13601/57597
## 2023-09-16 21:43:45 Annotating text fragment 13611/57597
## 2023-09-16 21:43:46 Annotating text fragment 13621/57597
## 2023-09-16 21:43:46 Annotating text fragment 13631/57597
## 2023-09-16 21:43:46 Annotating text fragment 13641/57597
## 2023-09-16 21:43:46 Annotating text fragment 13651/57597
## 2023-09-16 21:43:46 Annotating text fragment 13661/57597
## 2023-09-16 21:43:46 Annotating text fragment 13671/57597
## 2023-09-16 21:43:46 Annotating text fragment 13681/57597
## 2023-09-16 21:43:46 Annotating text fragment 13691/57597
## 2023-09-16 21:43:46 Annotating text fragment 13701/57597
## 2023-09-16 21:43:46 Annotating text fragment 13711/57597
## 2023-09-16 21:43:46 Annotating text fragment 13721/57597
## 2023-09-16 21:43:46 Annotating text fragment 13731/57597
```

```
## 2023-09-16 21:43:46 Annotating text fragment 13741/57597
## 2023-09-16 21:43:46 Annotating text fragment 13751/57597
## 2023-09-16 21:43:46 Annotating text fragment 13761/57597
## 2023-09-16 21:43:47 Annotating text fragment 13771/57597
## 2023-09-16 21:43:47 Annotating text fragment 13781/57597
## 2023-09-16 21:43:47 Annotating text fragment 13791/57597
## 2023-09-16 21:43:47 Annotating text fragment 13801/57597
## 2023-09-16 21:43:47 Annotating text fragment 13811/57597
## 2023-09-16 21:43:47 Annotating text fragment 13821/57597
## 2023-09-16 21:43:47 Annotating text fragment 13831/57597
## 2023-09-16 21:43:47 Annotating text fragment 13841/57597
## 2023-09-16 21:43:47 Annotating text fragment 13851/57597
## 2023-09-16 21:43:47 Annotating text fragment 13861/57597
## 2023-09-16 21:43:48 Annotating text fragment 13871/57597
## 2023-09-16 21:43:48 Annotating text fragment 13881/57597
## 2023-09-16 21:43:48 Annotating text fragment 13891/57597
## 2023-09-16 21:43:48 Annotating text fragment 13901/57597
## 2023-09-16 21:43:48 Annotating text fragment 13911/57597
## 2023-09-16 21:43:48 Annotating text fragment 13921/57597
## 2023-09-16 21:43:48 Annotating text fragment 13931/57597
## 2023-09-16 21:43:49 Annotating text fragment 13941/57597
## 2023-09-16 21:43:49 Annotating text fragment 13951/57597
## 2023-09-16 21:43:49 Annotating text fragment 13961/57597
## 2023-09-16 21:43:49 Annotating text fragment 13971/57597
## 2023-09-16 21:43:49 Annotating text fragment 13981/57597
## 2023-09-16 21:43:49 Annotating text fragment 13991/57597
## 2023-09-16 21:43:49 Annotating text fragment 14001/57597
## 2023-09-16 21:43:49 Annotating text fragment 14011/57597
## 2023-09-16 21:43:49 Annotating text fragment 14021/57597
## 2023-09-16 21:43:49 Annotating text fragment 14031/57597
## 2023-09-16 21:43:49 Annotating text fragment 14041/57597
## 2023-09-16 21:43:49 Annotating text fragment 14051/57597
## 2023-09-16 21:43:49 Annotating text fragment 14061/57597
## 2023-09-16 21:43:49 Annotating text fragment 14071/57597
## 2023-09-16 21:43:49 Annotating text fragment 14081/57597
## 2023-09-16 21:43:49 Annotating text fragment 14091/57597
## 2023-09-16 21:43:50 Annotating text fragment 14101/57597
## 2023-09-16 21:43:50 Annotating text fragment 14111/57597
## 2023-09-16 21:43:50 Annotating text fragment 14121/57597
## 2023-09-16 21:43:50 Annotating text fragment 14131/57597
## 2023-09-16 21:43:50 Annotating text fragment 14141/57597
## 2023-09-16 21:43:50 Annotating text fragment 14151/57597
## 2023-09-16 21:43:50 Annotating text fragment 14161/57597
## 2023-09-16 21:43:50 Annotating text fragment 14171/57597
## 2023-09-16 21:43:50 Annotating text fragment 14181/57597
## 2023-09-16 21:43:50 Annotating text fragment 14191/57597
## 2023-09-16 21:43:51 Annotating text fragment 14201/57597
## 2023-09-16 21:43:51 Annotating text fragment 14211/57597
## 2023-09-16 21:43:51 Annotating text fragment 14221/57597
## 2023-09-16 21:43:51 Annotating text fragment 14231/57597
## 2023-09-16 21:43:51 Annotating text fragment 14241/57597
## 2023-09-16 21:43:51 Annotating text fragment 14251/57597
## 2023-09-16 21:43:51 Annotating text fragment 14261/57597
## 2023-09-16 21:43:51 Annotating text fragment 14271/57597
```

```
## 2023-09-16 21:43:51 Annotating text fragment 14281/57597
## 2023-09-16 21:43:52 Annotating text fragment 14291/57597
## 2023-09-16 21:43:52 Annotating text fragment 14301/57597
## 2023-09-16 21:43:52 Annotating text fragment 14311/57597
## 2023-09-16 21:43:52 Annotating text fragment 14321/57597
## 2023-09-16 21:43:52 Annotating text fragment 14331/57597
## 2023-09-16 21:43:52 Annotating text fragment 14341/57597
## 2023-09-16 21:43:52 Annotating text fragment 14351/57597
## 2023-09-16 21:43:52 Annotating text fragment 14361/57597
## 2023-09-16 21:43:52 Annotating text fragment 14371/57597
## 2023-09-16 21:43:52 Annotating text fragment 14381/57597
## 2023-09-16 21:43:52 Annotating text fragment 14391/57597
## 2023-09-16 21:43:52 Annotating text fragment 14401/57597
## 2023-09-16 21:43:52 Annotating text fragment 14411/57597
## 2023-09-16 21:43:53 Annotating text fragment 14421/57597
## 2023-09-16 21:43:53 Annotating text fragment 14431/57597
## 2023-09-16 21:43:53 Annotating text fragment 14441/57597
## 2023-09-16 21:43:53 Annotating text fragment 14451/57597
## 2023-09-16 21:43:53 Annotating text fragment 14461/57597
## 2023-09-16 21:43:53 Annotating text fragment 14471/57597
## 2023-09-16 21:43:53 Annotating text fragment 14481/57597
## 2023-09-16 21:43:53 Annotating text fragment 14491/57597
## 2023-09-16 21:43:53 Annotating text fragment 14501/57597
## 2023-09-16 21:43:53 Annotating text fragment 14511/57597
## 2023-09-16 21:43:53 Annotating text fragment 14521/57597
## 2023-09-16 21:43:53 Annotating text fragment 14531/57597
## 2023-09-16 21:43:53 Annotating text fragment 14541/57597
## 2023-09-16 21:43:53 Annotating text fragment 14551/57597
## 2023-09-16 21:43:53 Annotating text fragment 14561/57597
## 2023-09-16 21:43:54 Annotating text fragment 14571/57597
## 2023-09-16 21:43:54 Annotating text fragment 14581/57597
## 2023-09-16 21:43:54 Annotating text fragment 14591/57597
## 2023-09-16 21:43:54 Annotating text fragment 14601/57597
## 2023-09-16 21:43:54 Annotating text fragment 14611/57597
## 2023-09-16 21:43:54 Annotating text fragment 14621/57597
## 2023-09-16 21:43:54 Annotating text fragment 14631/57597
## 2023-09-16 21:43:54 Annotating text fragment 14641/57597
## 2023-09-16 21:43:54 Annotating text fragment 14651/57597
## 2023-09-16 21:43:55 Annotating text fragment 14661/57597
## 2023-09-16 21:43:55 Annotating text fragment 14671/57597
## 2023-09-16 21:43:55 Annotating text fragment 14681/57597
## 2023-09-16 21:43:55 Annotating text fragment 14691/57597
## 2023-09-16 21:43:55 Annotating text fragment 14701/57597
## 2023-09-16 21:43:55 Annotating text fragment 14711/57597
## 2023-09-16 21:43:55 Annotating text fragment 14721/57597
## 2023-09-16 21:43:55 Annotating text fragment 14731/57597
## 2023-09-16 21:43:55 Annotating text fragment 14741/57597
## 2023-09-16 21:43:55 Annotating text fragment 14751/57597
## 2023-09-16 21:43:55 Annotating text fragment 14761/57597
## 2023-09-16 21:43:55 Annotating text fragment 14771/57597
## 2023-09-16 21:43:55 Annotating text fragment 14781/57597
## 2023-09-16 21:43:56 Annotating text fragment 14791/57597
## 2023-09-16 21:43:56 Annotating text fragment 14801/57597
## 2023-09-16 21:43:56 Annotating text fragment 14811/57597
```

```
## 2023-09-16 21:43:56 Annotating text fragment 14821/57597
## 2023-09-16 21:43:56 Annotating text fragment 14831/57597
## 2023-09-16 21:43:56 Annotating text fragment 14841/57597
## 2023-09-16 21:43:56 Annotating text fragment 14851/57597
## 2023-09-16 21:43:56 Annotating text fragment 14861/57597
## 2023-09-16 21:43:56 Annotating text fragment 14871/57597
## 2023-09-16 21:43:56 Annotating text fragment 14881/57597
## 2023-09-16 21:43:56 Annotating text fragment 14891/57597
## 2023-09-16 21:43:56 Annotating text fragment 14901/57597
## 2023-09-16 21:43:56 Annotating text fragment 14911/57597
## 2023-09-16 21:43:56 Annotating text fragment 14921/57597
## 2023-09-16 21:43:56 Annotating text fragment 14931/57597
## 2023-09-16 21:43:57 Annotating text fragment 14941/57597
## 2023-09-16 21:43:57 Annotating text fragment 14951/57597
## 2023-09-16 21:43:57 Annotating text fragment 14961/57597
## 2023-09-16 21:43:57 Annotating text fragment 14971/57597
## 2023-09-16 21:43:57 Annotating text fragment 14981/57597
## 2023-09-16 21:43:57 Annotating text fragment 14991/57597
## 2023-09-16 21:43:57 Annotating text fragment 15001/57597
## 2023-09-16 21:43:57 Annotating text fragment 15011/57597
## 2023-09-16 21:43:57 Annotating text fragment 15021/57597
## 2023-09-16 21:43:57 Annotating text fragment 15031/57597
## 2023-09-16 21:43:57 Annotating text fragment 15041/57597
## 2023-09-16 21:43:57 Annotating text fragment 15051/57597
## 2023-09-16 21:43:57 Annotating text fragment 15061/57597
## 2023-09-16 21:43:57 Annotating text fragment 15071/57597
## 2023-09-16 21:43:58 Annotating text fragment 15081/57597
## 2023-09-16 21:43:58 Annotating text fragment 15091/57597
## 2023-09-16 21:43:58 Annotating text fragment 15101/57597
## 2023-09-16 21:43:58 Annotating text fragment 15111/57597
## 2023-09-16 21:43:58 Annotating text fragment 15121/57597
## 2023-09-16 21:43:58 Annotating text fragment 15131/57597
## 2023-09-16 21:43:58 Annotating text fragment 15141/57597
## 2023-09-16 21:43:58 Annotating text fragment 15151/57597
## 2023-09-16 21:43:58 Annotating text fragment 15161/57597
## 2023-09-16 21:43:58 Annotating text fragment 15171/57597
## 2023-09-16 21:43:58 Annotating text fragment 15181/57597
## 2023-09-16 21:43:58 Annotating text fragment 15191/57597
## 2023-09-16 21:43:58 Annotating text fragment 15201/57597
## 2023-09-16 21:43:58 Annotating text fragment 15211/57597
## 2023-09-16 21:43:58 Annotating text fragment 15221/57597
## 2023-09-16 21:43:58 Annotating text fragment 15231/57597
## 2023-09-16 21:43:58 Annotating text fragment 15241/57597
## 2023-09-16 21:43:58 Annotating text fragment 15251/57597
## 2023-09-16 21:43:59 Annotating text fragment 15261/57597
## 2023-09-16 21:43:59 Annotating text fragment 15271/57597
## 2023-09-16 21:43:59 Annotating text fragment 15281/57597
## 2023-09-16 21:43:59 Annotating text fragment 15291/57597
## 2023-09-16 21:43:59 Annotating text fragment 15301/57597
## 2023-09-16 21:43:59 Annotating text fragment 15311/57597
## 2023-09-16 21:43:59 Annotating text fragment 15321/57597
## 2023-09-16 21:43:59 Annotating text fragment 15331/57597
## 2023-09-16 21:43:59 Annotating text fragment 15341/57597
## 2023-09-16 21:43:59 Annotating text fragment 15351/57597
```

```
## 2023-09-16 21:43:59 Annotating text fragment 15361/57597
## 2023-09-16 21:43:59 Annotating text fragment 15371/57597
## 2023-09-16 21:43:59 Annotating text fragment 15381/57597
## 2023-09-16 21:43:59 Annotating text fragment 15391/57597
## 2023-09-16 21:44:00 Annotating text fragment 15401/57597
## 2023-09-16 21:44:00 Annotating text fragment 15411/57597
## 2023-09-16 21:44:00 Annotating text fragment 15421/57597
## 2023-09-16 21:44:00 Annotating text fragment 15431/57597
## 2023-09-16 21:44:00 Annotating text fragment 15441/57597
## 2023-09-16 21:44:00 Annotating text fragment 15451/57597
## 2023-09-16 21:44:00 Annotating text fragment 15461/57597
## 2023-09-16 21:44:00 Annotating text fragment 15471/57597
## 2023-09-16 21:44:00 Annotating text fragment 15481/57597
## 2023-09-16 21:44:00 Annotating text fragment 15491/57597
## 2023-09-16 21:44:00 Annotating text fragment 15501/57597
## 2023-09-16 21:44:01 Annotating text fragment 15511/57597
## 2023-09-16 21:44:01 Annotating text fragment 15521/57597
## 2023-09-16 21:44:01 Annotating text fragment 15531/57597
## 2023-09-16 21:44:01 Annotating text fragment 15541/57597
## 2023-09-16 21:44:01 Annotating text fragment 15551/57597
## 2023-09-16 21:44:01 Annotating text fragment 15561/57597
## 2023-09-16 21:44:01 Annotating text fragment 15571/57597
## 2023-09-16 21:44:01 Annotating text fragment 15581/57597
## 2023-09-16 21:44:01 Annotating text fragment 15591/57597
## 2023-09-16 21:44:01 Annotating text fragment 15601/57597
## 2023-09-16 21:44:01 Annotating text fragment 15611/57597
## 2023-09-16 21:44:01 Annotating text fragment 15621/57597
## 2023-09-16 21:44:01 Annotating text fragment 15631/57597
## 2023-09-16 21:44:02 Annotating text fragment 15641/57597
## 2023-09-16 21:44:02 Annotating text fragment 15651/57597
## 2023-09-16 21:44:02 Annotating text fragment 15661/57597
## 2023-09-16 21:44:02 Annotating text fragment 15671/57597
## 2023-09-16 21:44:02 Annotating text fragment 15681/57597
## 2023-09-16 21:44:02 Annotating text fragment 15691/57597
## 2023-09-16 21:44:02 Annotating text fragment 15701/57597
## 2023-09-16 21:44:02 Annotating text fragment 15711/57597
## 2023-09-16 21:44:02 Annotating text fragment 15721/57597
## 2023-09-16 21:44:02 Annotating text fragment 15731/57597
## 2023-09-16 21:44:02 Annotating text fragment 15741/57597
## 2023-09-16 21:44:02 Annotating text fragment 15751/57597
## 2023-09-16 21:44:02 Annotating text fragment 15761/57597
## 2023-09-16 21:44:02 Annotating text fragment 15771/57597
## 2023-09-16 21:44:02 Annotating text fragment 15781/57597
## 2023-09-16 21:44:02 Annotating text fragment 15791/57597
## 2023-09-16 21:44:03 Annotating text fragment 15801/57597
## 2023-09-16 21:44:03 Annotating text fragment 15811/57597
## 2023-09-16 21:44:03 Annotating text fragment 15821/57597
## 2023-09-16 21:44:03 Annotating text fragment 15831/57597
## 2023-09-16 21:44:03 Annotating text fragment 15841/57597
## 2023-09-16 21:44:03 Annotating text fragment 15851/57597
## 2023-09-16 21:44:03 Annotating text fragment 15861/57597
## 2023-09-16 21:44:03 Annotating text fragment 15871/57597
## 2023-09-16 21:44:03 Annotating text fragment 15881/57597
## 2023-09-16 21:44:04 Annotating text fragment 15891/57597
```

```
## 2023-09-16 21:44:04 Annotating text fragment 15901/57597
## 2023-09-16 21:44:04 Annotating text fragment 15911/57597
## 2023-09-16 21:44:04 Annotating text fragment 15921/57597
## 2023-09-16 21:44:04 Annotating text fragment 15931/57597
## 2023-09-16 21:44:04 Annotating text fragment 15941/57597
## 2023-09-16 21:44:04 Annotating text fragment 15951/57597
## 2023-09-16 21:44:04 Annotating text fragment 15961/57597
## 2023-09-16 21:44:04 Annotating text fragment 15971/57597
## 2023-09-16 21:44:05 Annotating text fragment 15981/57597
## 2023-09-16 21:44:05 Annotating text fragment 15991/57597
## 2023-09-16 21:44:05 Annotating text fragment 16001/57597
## 2023-09-16 21:44:05 Annotating text fragment 16011/57597
## 2023-09-16 21:44:05 Annotating text fragment 16021/57597
## 2023-09-16 21:44:05 Annotating text fragment 16031/57597
## 2023-09-16 21:44:05 Annotating text fragment 16041/57597
## 2023-09-16 21:44:05 Annotating text fragment 16051/57597
## 2023-09-16 21:44:05 Annotating text fragment 16061/57597
## 2023-09-16 21:44:06 Annotating text fragment 16071/57597
## 2023-09-16 21:44:06 Annotating text fragment 16081/57597
## 2023-09-16 21:44:06 Annotating text fragment 16091/57597
## 2023-09-16 21:44:06 Annotating text fragment 16101/57597
## 2023-09-16 21:44:06 Annotating text fragment 16111/57597
## 2023-09-16 21:44:06 Annotating text fragment 16121/57597
## 2023-09-16 21:44:06 Annotating text fragment 16131/57597
## 2023-09-16 21:44:06 Annotating text fragment 16141/57597
## 2023-09-16 21:44:06 Annotating text fragment 16151/57597
## 2023-09-16 21:44:06 Annotating text fragment 16161/57597
## 2023-09-16 21:44:06 Annotating text fragment 16171/57597
## 2023-09-16 21:44:06 Annotating text fragment 16181/57597
## 2023-09-16 21:44:06 Annotating text fragment 16191/57597
## 2023-09-16 21:44:06 Annotating text fragment 16201/57597
## 2023-09-16 21:44:06 Annotating text fragment 16211/57597
## 2023-09-16 21:44:06 Annotating text fragment 16221/57597
## 2023-09-16 21:44:06 Annotating text fragment 16231/57597
## 2023-09-16 21:44:06 Annotating text fragment 16241/57597
## 2023-09-16 21:44:06 Annotating text fragment 16251/57597
## 2023-09-16 21:44:06 Annotating text fragment 16261/57597
## 2023-09-16 21:44:06 Annotating text fragment 16271/57597
## 2023-09-16 21:44:06 Annotating text fragment 16281/57597
## 2023-09-16 21:44:06 Annotating text fragment 16291/57597
## 2023-09-16 21:44:07 Annotating text fragment 16301/57597
## 2023-09-16 21:44:07 Annotating text fragment 16311/57597
## 2023-09-16 21:44:07 Annotating text fragment 16321/57597
## 2023-09-16 21:44:07 Annotating text fragment 16331/57597
## 2023-09-16 21:44:07 Annotating text fragment 16341/57597
## 2023-09-16 21:44:07 Annotating text fragment 16351/57597
## 2023-09-16 21:44:07 Annotating text fragment 16361/57597
## 2023-09-16 21:44:07 Annotating text fragment 16371/57597
## 2023-09-16 21:44:07 Annotating text fragment 16381/57597
## 2023-09-16 21:44:07 Annotating text fragment 16391/57597
## 2023-09-16 21:44:07 Annotating text fragment 16401/57597
## 2023-09-16 21:44:07 Annotating text fragment 16411/57597
## 2023-09-16 21:44:08 Annotating text fragment 16421/57597
## 2023-09-16 21:44:08 Annotating text fragment 16431/57597
```

```
## 2023-09-16 21:44:08 Annotating text fragment 16441/57597
## 2023-09-16 21:44:08 Annotating text fragment 16451/57597
## 2023-09-16 21:44:08 Annotating text fragment 16461/57597
## 2023-09-16 21:44:08 Annotating text fragment 16471/57597
## 2023-09-16 21:44:08 Annotating text fragment 16481/57597
## 2023-09-16 21:44:08 Annotating text fragment 16491/57597
## 2023-09-16 21:44:08 Annotating text fragment 16501/57597
## 2023-09-16 21:44:08 Annotating text fragment 16511/57597
## 2023-09-16 21:44:08 Annotating text fragment 16521/57597
## 2023-09-16 21:44:08 Annotating text fragment 16531/57597
## 2023-09-16 21:44:09 Annotating text fragment 16541/57597
## 2023-09-16 21:44:09 Annotating text fragment 16551/57597
## 2023-09-16 21:44:09 Annotating text fragment 16561/57597
## 2023-09-16 21:44:09 Annotating text fragment 16571/57597
## 2023-09-16 21:44:09 Annotating text fragment 16581/57597
## 2023-09-16 21:44:09 Annotating text fragment 16591/57597
## 2023-09-16 21:44:09 Annotating text fragment 16601/57597
## 2023-09-16 21:44:09 Annotating text fragment 16611/57597
## 2023-09-16 21:44:09 Annotating text fragment 16621/57597
## 2023-09-16 21:44:09 Annotating text fragment 16631/57597
## 2023-09-16 21:44:09 Annotating text fragment 16641/57597
## 2023-09-16 21:44:10 Annotating text fragment 16651/57597
## 2023-09-16 21:44:10 Annotating text fragment 16661/57597
## 2023-09-16 21:44:10 Annotating text fragment 16671/57597
## 2023-09-16 21:44:10 Annotating text fragment 16681/57597
## 2023-09-16 21:44:10 Annotating text fragment 16691/57597
## 2023-09-16 21:44:10 Annotating text fragment 16701/57597
## 2023-09-16 21:44:10 Annotating text fragment 16711/57597
## 2023-09-16 21:44:10 Annotating text fragment 16721/57597
## 2023-09-16 21:44:10 Annotating text fragment 16731/57597
## 2023-09-16 21:44:10 Annotating text fragment 16741/57597
## 2023-09-16 21:44:10 Annotating text fragment 16751/57597
## 2023-09-16 21:44:10 Annotating text fragment 16761/57597
## 2023-09-16 21:44:10 Annotating text fragment 16771/57597
## 2023-09-16 21:44:10 Annotating text fragment 16781/57597
## 2023-09-16 21:44:10 Annotating text fragment 16791/57597
## 2023-09-16 21:44:11 Annotating text fragment 16801/57597
## 2023-09-16 21:44:11 Annotating text fragment 16811/57597
## 2023-09-16 21:44:11 Annotating text fragment 16821/57597
## 2023-09-16 21:44:11 Annotating text fragment 16831/57597
## 2023-09-16 21:44:11 Annotating text fragment 16841/57597
## 2023-09-16 21:44:11 Annotating text fragment 16851/57597
## 2023-09-16 21:44:11 Annotating text fragment 16861/57597
## 2023-09-16 21:44:11 Annotating text fragment 16871/57597
## 2023-09-16 21:44:11 Annotating text fragment 16881/57597
## 2023-09-16 21:44:11 Annotating text fragment 16891/57597
## 2023-09-16 21:44:11 Annotating text fragment 16901/57597
## 2023-09-16 21:44:11 Annotating text fragment 16911/57597
## 2023-09-16 21:44:11 Annotating text fragment 16921/57597
## 2023-09-16 21:44:11 Annotating text fragment 16931/57597
## 2023-09-16 21:44:12 Annotating text fragment 16941/57597
## 2023-09-16 21:44:12 Annotating text fragment 16951/57597
## 2023-09-16 21:44:12 Annotating text fragment 16961/57597
## 2023-09-16 21:44:12 Annotating text fragment 16971/57597
```

```
## 2023-09-16 21:44:12 Annotating text fragment 16981/57597
## 2023-09-16 21:44:12 Annotating text fragment 16991/57597
## 2023-09-16 21:44:12 Annotating text fragment 17001/57597
## 2023-09-16 21:44:12 Annotating text fragment 17011/57597
## 2023-09-16 21:44:12 Annotating text fragment 17021/57597
## 2023-09-16 21:44:12 Annotating text fragment 17031/57597
## 2023-09-16 21:44:12 Annotating text fragment 17041/57597
## 2023-09-16 21:44:12 Annotating text fragment 17051/57597
## 2023-09-16 21:44:13 Annotating text fragment 17061/57597
## 2023-09-16 21:44:13 Annotating text fragment 17071/57597
## 2023-09-16 21:44:13 Annotating text fragment 17081/57597
## 2023-09-16 21:44:13 Annotating text fragment 17091/57597
## 2023-09-16 21:44:13 Annotating text fragment 17101/57597
## 2023-09-16 21:44:13 Annotating text fragment 17111/57597
## 2023-09-16 21:44:13 Annotating text fragment 17121/57597
## 2023-09-16 21:44:13 Annotating text fragment 17131/57597
## 2023-09-16 21:44:13 Annotating text fragment 17141/57597
## 2023-09-16 21:44:13 Annotating text fragment 17151/57597
## 2023-09-16 21:44:13 Annotating text fragment 17161/57597
## 2023-09-16 21:44:13 Annotating text fragment 17171/57597
## 2023-09-16 21:44:13 Annotating text fragment 17181/57597
## 2023-09-16 21:44:13 Annotating text fragment 17191/57597
## 2023-09-16 21:44:13 Annotating text fragment 17201/57597
## 2023-09-16 21:44:14 Annotating text fragment 17211/57597
## 2023-09-16 21:44:14 Annotating text fragment 17221/57597
## 2023-09-16 21:44:14 Annotating text fragment 17231/57597
## 2023-09-16 21:44:14 Annotating text fragment 17241/57597
## 2023-09-16 21:44:14 Annotating text fragment 17251/57597
## 2023-09-16 21:44:14 Annotating text fragment 17261/57597
## 2023-09-16 21:44:14 Annotating text fragment 17271/57597
## 2023-09-16 21:44:14 Annotating text fragment 17281/57597
## 2023-09-16 21:44:14 Annotating text fragment 17291/57597
## 2023-09-16 21:44:14 Annotating text fragment 17301/57597
## 2023-09-16 21:44:14 Annotating text fragment 17311/57597
## 2023-09-16 21:44:14 Annotating text fragment 17321/57597
## 2023-09-16 21:44:14 Annotating text fragment 17331/57597
## 2023-09-16 21:44:15 Annotating text fragment 17341/57597
## 2023-09-16 21:44:15 Annotating text fragment 17351/57597
## 2023-09-16 21:44:15 Annotating text fragment 17361/57597
## 2023-09-16 21:44:15 Annotating text fragment 17371/57597
## 2023-09-16 21:44:15 Annotating text fragment 17381/57597
## 2023-09-16 21:44:15 Annotating text fragment 17391/57597
## 2023-09-16 21:44:16 Annotating text fragment 17401/57597
## 2023-09-16 21:44:16 Annotating text fragment 17411/57597
## 2023-09-16 21:44:16 Annotating text fragment 17421/57597
## 2023-09-16 21:44:16 Annotating text fragment 17431/57597
## 2023-09-16 21:44:16 Annotating text fragment 17441/57597
## 2023-09-16 21:44:17 Annotating text fragment 17451/57597
## 2023-09-16 21:44:17 Annotating text fragment 17461/57597
## 2023-09-16 21:44:17 Annotating text fragment 17471/57597
## 2023-09-16 21:44:17 Annotating text fragment 17481/57597
## 2023-09-16 21:44:17 Annotating text fragment 17491/57597
## 2023-09-16 21:44:17 Annotating text fragment 17501/57597
## 2023-09-16 21:44:17 Annotating text fragment 17511/57597
```

```
## 2023-09-16 21:44:17 Annotating text fragment 17521/57597
## 2023-09-16 21:44:17 Annotating text fragment 17531/57597
## 2023-09-16 21:44:17 Annotating text fragment 17541/57597
## 2023-09-16 21:44:17 Annotating text fragment 17551/57597
## 2023-09-16 21:44:17 Annotating text fragment 17561/57597
## 2023-09-16 21:44:17 Annotating text fragment 17571/57597
## 2023-09-16 21:44:17 Annotating text fragment 17581/57597
## 2023-09-16 21:44:18 Annotating text fragment 17591/57597
## 2023-09-16 21:44:18 Annotating text fragment 17601/57597
## 2023-09-16 21:44:18 Annotating text fragment 17611/57597
## 2023-09-16 21:44:18 Annotating text fragment 17621/57597
## 2023-09-16 21:44:18 Annotating text fragment 17631/57597
## 2023-09-16 21:44:18 Annotating text fragment 17641/57597
## 2023-09-16 21:44:18 Annotating text fragment 17651/57597
## 2023-09-16 21:44:18 Annotating text fragment 17661/57597
## 2023-09-16 21:44:18 Annotating text fragment 17671/57597
## 2023-09-16 21:44:18 Annotating text fragment 17681/57597
## 2023-09-16 21:44:18 Annotating text fragment 17691/57597
## 2023-09-16 21:44:18 Annotating text fragment 17701/57597
## 2023-09-16 21:44:18 Annotating text fragment 17711/57597
## 2023-09-16 21:44:18 Annotating text fragment 17721/57597
## 2023-09-16 21:44:18 Annotating text fragment 17731/57597
## 2023-09-16 21:44:18 Annotating text fragment 17741/57597
## 2023-09-16 21:44:18 Annotating text fragment 17751/57597
## 2023-09-16 21:44:18 Annotating text fragment 17761/57597
## 2023-09-16 21:44:19 Annotating text fragment 17771/57597
## 2023-09-16 21:44:19 Annotating text fragment 17781/57597
## 2023-09-16 21:44:19 Annotating text fragment 17791/57597
## 2023-09-16 21:44:19 Annotating text fragment 17801/57597
## 2023-09-16 21:44:19 Annotating text fragment 17811/57597
## 2023-09-16 21:44:19 Annotating text fragment 17821/57597
## 2023-09-16 21:44:19 Annotating text fragment 17831/57597
## 2023-09-16 21:44:19 Annotating text fragment 17841/57597
## 2023-09-16 21:44:21 Annotating text fragment 17851/57597
## 2023-09-16 21:44:23 Annotating text fragment 17861/57597
## 2023-09-16 21:44:26 Annotating text fragment 17871/57597
## 2023-09-16 21:44:28 Annotating text fragment 17881/57597
## 2023-09-16 21:44:31 Annotating text fragment 17891/57597
## 2023-09-16 21:44:32 Annotating text fragment 17901/57597
## 2023-09-16 21:44:32 Annotating text fragment 17911/57597
## 2023-09-16 21:44:33 Annotating text fragment 17921/57597
## 2023-09-16 21:44:34 Annotating text fragment 17931/57597
## 2023-09-16 21:44:34 Annotating text fragment 17941/57597
## 2023-09-16 21:44:34 Annotating text fragment 17951/57597
## 2023-09-16 21:44:34 Annotating text fragment 17961/57597
## 2023-09-16 21:44:34 Annotating text fragment 17971/57597
## 2023-09-16 21:44:34 Annotating text fragment 17981/57597
## 2023-09-16 21:44:34 Annotating text fragment 17991/57597
## 2023-09-16 21:44:34 Annotating text fragment 18001/57597
## 2023-09-16 21:44:34 Annotating text fragment 18011/57597
## 2023-09-16 21:44:35 Annotating text fragment 18021/57597
## 2023-09-16 21:44:35 Annotating text fragment 18031/57597
## 2023-09-16 21:44:35 Annotating text fragment 18041/57597
## 2023-09-16 21:44:35 Annotating text fragment 18051/57597
```

```
## 2023-09-16 21:44:35 Annotating text fragment 18061/57597
## 2023-09-16 21:44:35 Annotating text fragment 18071/57597
## 2023-09-16 21:44:35 Annotating text fragment 18081/57597
## 2023-09-16 21:44:35 Annotating text fragment 18091/57597
## 2023-09-16 21:44:35 Annotating text fragment 18101/57597
## 2023-09-16 21:44:35 Annotating text fragment 18111/57597
## 2023-09-16 21:44:35 Annotating text fragment 18121/57597
## 2023-09-16 21:44:35 Annotating text fragment 18131/57597
## 2023-09-16 21:44:35 Annotating text fragment 18141/57597
## 2023-09-16 21:44:36 Annotating text fragment 18151/57597
## 2023-09-16 21:44:36 Annotating text fragment 18161/57597
## 2023-09-16 21:44:36 Annotating text fragment 18171/57597
## 2023-09-16 21:44:36 Annotating text fragment 18181/57597
## 2023-09-16 21:44:36 Annotating text fragment 18191/57597
## 2023-09-16 21:44:36 Annotating text fragment 18201/57597
## 2023-09-16 21:44:36 Annotating text fragment 18211/57597
## 2023-09-16 21:44:36 Annotating text fragment 18221/57597
## 2023-09-16 21:44:37 Annotating text fragment 18231/57597
## 2023-09-16 21:44:37 Annotating text fragment 18241/57597
## 2023-09-16 21:44:37 Annotating text fragment 18251/57597
## 2023-09-16 21:44:37 Annotating text fragment 18261/57597
## 2023-09-16 21:44:37 Annotating text fragment 18271/57597
## 2023-09-16 21:44:37 Annotating text fragment 18281/57597
## 2023-09-16 21:44:38 Annotating text fragment 18291/57597
## 2023-09-16 21:44:38 Annotating text fragment 18301/57597
## 2023-09-16 21:44:38 Annotating text fragment 18311/57597
## 2023-09-16 21:44:38 Annotating text fragment 18321/57597
## 2023-09-16 21:44:38 Annotating text fragment 18331/57597
## 2023-09-16 21:44:38 Annotating text fragment 18341/57597
## 2023-09-16 21:44:38 Annotating text fragment 18351/57597
## 2023-09-16 21:44:38 Annotating text fragment 18361/57597
## 2023-09-16 21:44:38 Annotating text fragment 18371/57597
## 2023-09-16 21:44:38 Annotating text fragment 18381/57597
## 2023-09-16 21:44:38 Annotating text fragment 18391/57597
## 2023-09-16 21:44:38 Annotating text fragment 18401/57597
## 2023-09-16 21:44:38 Annotating text fragment 18411/57597
## 2023-09-16 21:44:38 Annotating text fragment 18421/57597
## 2023-09-16 21:44:38 Annotating text fragment 18431/57597
## 2023-09-16 21:44:39 Annotating text fragment 18441/57597
## 2023-09-16 21:44:39 Annotating text fragment 18451/57597
## 2023-09-16 21:44:39 Annotating text fragment 18461/57597
## 2023-09-16 21:44:39 Annotating text fragment 18471/57597
## 2023-09-16 21:44:39 Annotating text fragment 18481/57597
## 2023-09-16 21:44:39 Annotating text fragment 18491/57597
## 2023-09-16 21:44:39 Annotating text fragment 18501/57597
## 2023-09-16 21:44:39 Annotating text fragment 18511/57597
## 2023-09-16 21:44:39 Annotating text fragment 18521/57597
## 2023-09-16 21:44:39 Annotating text fragment 18531/57597
## 2023-09-16 21:44:39 Annotating text fragment 18541/57597
## 2023-09-16 21:44:40 Annotating text fragment 18551/57597
## 2023-09-16 21:44:40 Annotating text fragment 18561/57597
## 2023-09-16 21:44:40 Annotating text fragment 18571/57597
## 2023-09-16 21:44:40 Annotating text fragment 18581/57597
## 2023-09-16 21:44:40 Annotating text fragment 18591/57597
```

```
## 2023-09-16 21:44:40 Annotating text fragment 18601/57597
## 2023-09-16 21:44:40 Annotating text fragment 18611/57597
## 2023-09-16 21:44:40 Annotating text fragment 18621/57597
## 2023-09-16 21:44:40 Annotating text fragment 18631/57597
## 2023-09-16 21:44:41 Annotating text fragment 18641/57597
## 2023-09-16 21:44:41 Annotating text fragment 18651/57597
## 2023-09-16 21:44:41 Annotating text fragment 18661/57597
## 2023-09-16 21:44:41 Annotating text fragment 18671/57597
## 2023-09-16 21:44:41 Annotating text fragment 18681/57597
## 2023-09-16 21:44:41 Annotating text fragment 18691/57597
## 2023-09-16 21:44:41 Annotating text fragment 18701/57597
## 2023-09-16 21:44:41 Annotating text fragment 18711/57597
## 2023-09-16 21:44:41 Annotating text fragment 18721/57597
## 2023-09-16 21:44:41 Annotating text fragment 18731/57597
## 2023-09-16 21:44:41 Annotating text fragment 18741/57597
## 2023-09-16 21:44:42 Annotating text fragment 18751/57597
## 2023-09-16 21:44:42 Annotating text fragment 18761/57597
## 2023-09-16 21:44:42 Annotating text fragment 18771/57597
## 2023-09-16 21:44:42 Annotating text fragment 18781/57597
## 2023-09-16 21:44:42 Annotating text fragment 18791/57597
## 2023-09-16 21:44:42 Annotating text fragment 18801/57597
## 2023-09-16 21:44:42 Annotating text fragment 18811/57597
## 2023-09-16 21:44:42 Annotating text fragment 18821/57597
## 2023-09-16 21:44:42 Annotating text fragment 18831/57597
## 2023-09-16 21:44:42 Annotating text fragment 18841/57597
## 2023-09-16 21:44:42 Annotating text fragment 18851/57597
## 2023-09-16 21:44:42 Annotating text fragment 18861/57597
## 2023-09-16 21:44:42 Annotating text fragment 18871/57597
## 2023-09-16 21:44:42 Annotating text fragment 18881/57597
## 2023-09-16 21:44:42 Annotating text fragment 18891/57597
## 2023-09-16 21:44:42 Annotating text fragment 18901/57597
## 2023-09-16 21:44:42 Annotating text fragment 18911/57597
## 2023-09-16 21:44:43 Annotating text fragment 18921/57597
## 2023-09-16 21:44:43 Annotating text fragment 18931/57597
## 2023-09-16 21:44:43 Annotating text fragment 18941/57597
## 2023-09-16 21:44:43 Annotating text fragment 18951/57597
## 2023-09-16 21:44:43 Annotating text fragment 18961/57597
## 2023-09-16 21:44:43 Annotating text fragment 18971/57597
## 2023-09-16 21:44:43 Annotating text fragment 18981/57597
## 2023-09-16 21:44:43 Annotating text fragment 18991/57597
## 2023-09-16 21:44:43 Annotating text fragment 19001/57597
## 2023-09-16 21:44:43 Annotating text fragment 19011/57597
## 2023-09-16 21:44:44 Annotating text fragment 19021/57597
## 2023-09-16 21:44:44 Annotating text fragment 19031/57597
## 2023-09-16 21:44:44 Annotating text fragment 19041/57597
## 2023-09-16 21:44:44 Annotating text fragment 19051/57597
## 2023-09-16 21:44:44 Annotating text fragment 19061/57597
## 2023-09-16 21:44:44 Annotating text fragment 19071/57597
## 2023-09-16 21:44:44 Annotating text fragment 19081/57597
## 2023-09-16 21:44:44 Annotating text fragment 19091/57597
## 2023-09-16 21:44:44 Annotating text fragment 19101/57597
## 2023-09-16 21:44:44 Annotating text fragment 19111/57597
## 2023-09-16 21:44:44 Annotating text fragment 19121/57597
## 2023-09-16 21:44:44 Annotating text fragment 19131/57597
```

```
## 2023-09-16 21:44:45 Annotating text fragment 19141/57597
## 2023-09-16 21:44:45 Annotating text fragment 19151/57597
## 2023-09-16 21:44:45 Annotating text fragment 19161/57597
## 2023-09-16 21:44:45 Annotating text fragment 19171/57597
## 2023-09-16 21:44:45 Annotating text fragment 19181/57597
## 2023-09-16 21:44:45 Annotating text fragment 19191/57597
## 2023-09-16 21:44:45 Annotating text fragment 19201/57597
## 2023-09-16 21:44:46 Annotating text fragment 19211/57597
## 2023-09-16 21:44:46 Annotating text fragment 19221/57597
## 2023-09-16 21:44:46 Annotating text fragment 19231/57597
## 2023-09-16 21:44:46 Annotating text fragment 19241/57597
## 2023-09-16 21:44:46 Annotating text fragment 19251/57597
## 2023-09-16 21:44:46 Annotating text fragment 19261/57597
## 2023-09-16 21:44:46 Annotating text fragment 19271/57597
## 2023-09-16 21:44:46 Annotating text fragment 19281/57597
## 2023-09-16 21:44:46 Annotating text fragment 19291/57597
## 2023-09-16 21:44:46 Annotating text fragment 19301/57597
## 2023-09-16 21:44:47 Annotating text fragment 19311/57597
## 2023-09-16 21:44:47 Annotating text fragment 19321/57597
## 2023-09-16 21:44:47 Annotating text fragment 19331/57597
## 2023-09-16 21:44:47 Annotating text fragment 19341/57597
## 2023-09-16 21:44:47 Annotating text fragment 19351/57597
## 2023-09-16 21:44:47 Annotating text fragment 19361/57597
## 2023-09-16 21:44:47 Annotating text fragment 19371/57597
## 2023-09-16 21:44:47 Annotating text fragment 19381/57597
## 2023-09-16 21:44:47 Annotating text fragment 19391/57597
## 2023-09-16 21:44:48 Annotating text fragment 19401/57597
## 2023-09-16 21:44:48 Annotating text fragment 19411/57597
## 2023-09-16 21:44:48 Annotating text fragment 19421/57597
## 2023-09-16 21:44:48 Annotating text fragment 19431/57597
## 2023-09-16 21:44:48 Annotating text fragment 19441/57597
## 2023-09-16 21:44:48 Annotating text fragment 19451/57597
## 2023-09-16 21:44:48 Annotating text fragment 19461/57597
## 2023-09-16 21:44:48 Annotating text fragment 19471/57597
## 2023-09-16 21:44:48 Annotating text fragment 19481/57597
## 2023-09-16 21:44:48 Annotating text fragment 19491/57597
## 2023-09-16 21:44:48 Annotating text fragment 19501/57597
## 2023-09-16 21:44:48 Annotating text fragment 19511/57597
## 2023-09-16 21:44:48 Annotating text fragment 19521/57597
## 2023-09-16 21:44:48 Annotating text fragment 19531/57597
## 2023-09-16 21:44:48 Annotating text fragment 19541/57597
## 2023-09-16 21:44:48 Annotating text fragment 19551/57597
## 2023-09-16 21:44:49 Annotating text fragment 19561/57597
## 2023-09-16 21:44:49 Annotating text fragment 19571/57597
## 2023-09-16 21:44:49 Annotating text fragment 19581/57597
## 2023-09-16 21:44:49 Annotating text fragment 19591/57597
## 2023-09-16 21:44:49 Annotating text fragment 19601/57597
## 2023-09-16 21:44:49 Annotating text fragment 19611/57597
## 2023-09-16 21:44:49 Annotating text fragment 19621/57597
## 2023-09-16 21:44:49 Annotating text fragment 19631/57597
## 2023-09-16 21:44:49 Annotating text fragment 19641/57597
## 2023-09-16 21:44:49 Annotating text fragment 19651/57597
## 2023-09-16 21:44:49 Annotating text fragment 19661/57597
## 2023-09-16 21:44:49 Annotating text fragment 19671/57597
```

```
## 2023-09-16 21:44:49 Annotating text fragment 19681/57597
## 2023-09-16 21:44:49 Annotating text fragment 19691/57597
## 2023-09-16 21:44:49 Annotating text fragment 19701/57597
## 2023-09-16 21:44:49 Annotating text fragment 19711/57597
## 2023-09-16 21:44:50 Annotating text fragment 19721/57597
## 2023-09-16 21:44:50 Annotating text fragment 19731/57597
## 2023-09-16 21:44:50 Annotating text fragment 19741/57597
## 2023-09-16 21:44:50 Annotating text fragment 19751/57597
## 2023-09-16 21:44:50 Annotating text fragment 19761/57597
## 2023-09-16 21:44:50 Annotating text fragment 19771/57597
## 2023-09-16 21:44:50 Annotating text fragment 19781/57597
## 2023-09-16 21:44:50 Annotating text fragment 19791/57597
## 2023-09-16 21:44:50 Annotating text fragment 19801/57597
## 2023-09-16 21:44:50 Annotating text fragment 19811/57597
## 2023-09-16 21:44:50 Annotating text fragment 19821/57597
## 2023-09-16 21:44:50 Annotating text fragment 19831/57597
## 2023-09-16 21:44:50 Annotating text fragment 19841/57597
## 2023-09-16 21:44:50 Annotating text fragment 19851/57597
## 2023-09-16 21:44:50 Annotating text fragment 19861/57597
## 2023-09-16 21:44:50 Annotating text fragment 19871/57597
## 2023-09-16 21:44:51 Annotating text fragment 19881/57597
## 2023-09-16 21:44:51 Annotating text fragment 19891/57597
## 2023-09-16 21:44:51 Annotating text fragment 19901/57597
## 2023-09-16 21:44:51 Annotating text fragment 19911/57597
## 2023-09-16 21:44:51 Annotating text fragment 19921/57597
## 2023-09-16 21:44:51 Annotating text fragment 19931/57597
## 2023-09-16 21:44:51 Annotating text fragment 19941/57597
## 2023-09-16 21:44:51 Annotating text fragment 19951/57597
## 2023-09-16 21:44:51 Annotating text fragment 19961/57597
## 2023-09-16 21:44:51 Annotating text fragment 19971/57597
## 2023-09-16 21:44:51 Annotating text fragment 19981/57597
## 2023-09-16 21:44:51 Annotating text fragment 19991/57597
## 2023-09-16 21:44:51 Annotating text fragment 20001/57597
## 2023-09-16 21:44:52 Annotating text fragment 20011/57597
## 2023-09-16 21:44:52 Annotating text fragment 20021/57597
## 2023-09-16 21:44:52 Annotating text fragment 20031/57597
## 2023-09-16 21:44:52 Annotating text fragment 20041/57597
## 2023-09-16 21:44:52 Annotating text fragment 20051/57597
## 2023-09-16 21:44:52 Annotating text fragment 20061/57597
## 2023-09-16 21:44:52 Annotating text fragment 20071/57597
## 2023-09-16 21:44:52 Annotating text fragment 20081/57597
## 2023-09-16 21:44:52 Annotating text fragment 20091/57597
## 2023-09-16 21:44:52 Annotating text fragment 20101/57597
## 2023-09-16 21:44:52 Annotating text fragment 20111/57597
## 2023-09-16 21:44:52 Annotating text fragment 20121/57597
## 2023-09-16 21:44:52 Annotating text fragment 20131/57597
## 2023-09-16 21:44:52 Annotating text fragment 20141/57597
## 2023-09-16 21:44:53 Annotating text fragment 20151/57597
## 2023-09-16 21:44:53 Annotating text fragment 20161/57597
## 2023-09-16 21:44:53 Annotating text fragment 20171/57597
## 2023-09-16 21:44:53 Annotating text fragment 20181/57597
## 2023-09-16 21:44:53 Annotating text fragment 20191/57597
## 2023-09-16 21:44:53 Annotating text fragment 20201/57597
## 2023-09-16 21:44:53 Annotating text fragment 20211/57597
```

```
## 2023-09-16 21:44:53 Annotating text fragment 20221/57597
## 2023-09-16 21:44:53 Annotating text fragment 20231/57597
## 2023-09-16 21:44:53 Annotating text fragment 20241/57597
## 2023-09-16 21:44:53 Annotating text fragment 20251/57597
## 2023-09-16 21:44:53 Annotating text fragment 20261/57597
## 2023-09-16 21:44:53 Annotating text fragment 20271/57597
## 2023-09-16 21:44:53 Annotating text fragment 20281/57597
## 2023-09-16 21:44:54 Annotating text fragment 20291/57597
## 2023-09-16 21:44:54 Annotating text fragment 20301/57597
## 2023-09-16 21:44:54 Annotating text fragment 20311/57597
## 2023-09-16 21:44:54 Annotating text fragment 20321/57597
## 2023-09-16 21:44:54 Annotating text fragment 20331/57597
## 2023-09-16 21:44:54 Annotating text fragment 20341/57597
## 2023-09-16 21:44:54 Annotating text fragment 20351/57597
## 2023-09-16 21:44:54 Annotating text fragment 20361/57597
## 2023-09-16 21:44:54 Annotating text fragment 20371/57597
## 2023-09-16 21:44:54 Annotating text fragment 20381/57597
## 2023-09-16 21:44:54 Annotating text fragment 20391/57597
## 2023-09-16 21:44:54 Annotating text fragment 20401/57597
## 2023-09-16 21:44:55 Annotating text fragment 20411/57597
## 2023-09-16 21:44:55 Annotating text fragment 20421/57597
## 2023-09-16 21:44:55 Annotating text fragment 20431/57597
## 2023-09-16 21:44:55 Annotating text fragment 20441/57597
## 2023-09-16 21:44:55 Annotating text fragment 20451/57597
## 2023-09-16 21:44:55 Annotating text fragment 20461/57597
## 2023-09-16 21:44:55 Annotating text fragment 20471/57597
## 2023-09-16 21:44:55 Annotating text fragment 20481/57597
## 2023-09-16 21:44:55 Annotating text fragment 20491/57597
## 2023-09-16 21:44:55 Annotating text fragment 20501/57597
## 2023-09-16 21:44:55 Annotating text fragment 20511/57597
## 2023-09-16 21:44:56 Annotating text fragment 20521/57597
## 2023-09-16 21:44:56 Annotating text fragment 20531/57597
## 2023-09-16 21:44:56 Annotating text fragment 20541/57597
## 2023-09-16 21:44:56 Annotating text fragment 20551/57597
## 2023-09-16 21:44:56 Annotating text fragment 20561/57597
## 2023-09-16 21:44:56 Annotating text fragment 20571/57597
## 2023-09-16 21:44:56 Annotating text fragment 20581/57597
## 2023-09-16 21:44:56 Annotating text fragment 20591/57597
## 2023-09-16 21:44:56 Annotating text fragment 20601/57597
## 2023-09-16 21:44:56 Annotating text fragment 20611/57597
## 2023-09-16 21:44:56 Annotating text fragment 20621/57597
## 2023-09-16 21:44:56 Annotating text fragment 20631/57597
## 2023-09-16 21:44:56 Annotating text fragment 20641/57597
## 2023-09-16 21:44:57 Annotating text fragment 20651/57597
## 2023-09-16 21:44:57 Annotating text fragment 20661/57597
## 2023-09-16 21:44:57 Annotating text fragment 20671/57597
## 2023-09-16 21:44:57 Annotating text fragment 20681/57597
## 2023-09-16 21:44:57 Annotating text fragment 20691/57597
## 2023-09-16 21:44:57 Annotating text fragment 20701/57597
## 2023-09-16 21:44:57 Annotating text fragment 20711/57597
## 2023-09-16 21:44:57 Annotating text fragment 20721/57597
## 2023-09-16 21:44:57 Annotating text fragment 20731/57597
## 2023-09-16 21:44:57 Annotating text fragment 20741/57597
## 2023-09-16 21:44:57 Annotating text fragment 20751/57597
```

```
## 2023-09-16 21:44:57 Annotating text fragment 20761/57597
## 2023-09-16 21:44:57 Annotating text fragment 20771/57597
## 2023-09-16 21:44:58 Annotating text fragment 20781/57597
## 2023-09-16 21:44:58 Annotating text fragment 20791/57597
## 2023-09-16 21:44:58 Annotating text fragment 20801/57597
## 2023-09-16 21:44:58 Annotating text fragment 20811/57597
## 2023-09-16 21:44:58 Annotating text fragment 20821/57597
## 2023-09-16 21:44:58 Annotating text fragment 20831/57597
## 2023-09-16 21:44:58 Annotating text fragment 20841/57597
## 2023-09-16 21:44:58 Annotating text fragment 20851/57597
## 2023-09-16 21:44:58 Annotating text fragment 20861/57597
## 2023-09-16 21:44:58 Annotating text fragment 20871/57597
## 2023-09-16 21:44:58 Annotating text fragment 20881/57597
## 2023-09-16 21:44:58 Annotating text fragment 20891/57597
## 2023-09-16 21:44:58 Annotating text fragment 20901/57597
## 2023-09-16 21:44:58 Annotating text fragment 20911/57597
## 2023-09-16 21:44:59 Annotating text fragment 20921/57597
## 2023-09-16 21:44:59 Annotating text fragment 20931/57597
## 2023-09-16 21:44:59 Annotating text fragment 20941/57597
## 2023-09-16 21:44:59 Annotating text fragment 20951/57597
## 2023-09-16 21:44:59 Annotating text fragment 20961/57597
## 2023-09-16 21:45:00 Annotating text fragment 20971/57597
## 2023-09-16 21:45:00 Annotating text fragment 20981/57597
## 2023-09-16 21:45:00 Annotating text fragment 20991/57597
## 2023-09-16 21:45:00 Annotating text fragment 21001/57597
## 2023-09-16 21:45:00 Annotating text fragment 21011/57597
## 2023-09-16 21:45:00 Annotating text fragment 21021/57597
## 2023-09-16 21:45:00 Annotating text fragment 21031/57597
## 2023-09-16 21:45:00 Annotating text fragment 21041/57597
## 2023-09-16 21:45:00 Annotating text fragment 21051/57597
## 2023-09-16 21:45:00 Annotating text fragment 21061/57597
## 2023-09-16 21:45:00 Annotating text fragment 21071/57597
## 2023-09-16 21:45:00 Annotating text fragment 21081/57597
## 2023-09-16 21:45:00 Annotating text fragment 21091/57597
## 2023-09-16 21:45:00 Annotating text fragment 21101/57597
## 2023-09-16 21:45:00 Annotating text fragment 21111/57597
## 2023-09-16 21:45:01 Annotating text fragment 21121/57597
## 2023-09-16 21:45:01 Annotating text fragment 21131/57597
## 2023-09-16 21:45:01 Annotating text fragment 21141/57597
## 2023-09-16 21:45:01 Annotating text fragment 21151/57597
## 2023-09-16 21:45:01 Annotating text fragment 21161/57597
## 2023-09-16 21:45:01 Annotating text fragment 21171/57597
## 2023-09-16 21:45:01 Annotating text fragment 21181/57597
## 2023-09-16 21:45:01 Annotating text fragment 21191/57597
## 2023-09-16 21:45:01 Annotating text fragment 21201/57597
## 2023-09-16 21:45:01 Annotating text fragment 21211/57597
## 2023-09-16 21:45:01 Annotating text fragment 21221/57597
## 2023-09-16 21:45:01 Annotating text fragment 21231/57597
## 2023-09-16 21:45:02 Annotating text fragment 21241/57597
## 2023-09-16 21:45:02 Annotating text fragment 21251/57597
## 2023-09-16 21:45:02 Annotating text fragment 21261/57597
## 2023-09-16 21:45:02 Annotating text fragment 21271/57597
## 2023-09-16 21:45:02 Annotating text fragment 21281/57597
## 2023-09-16 21:45:02 Annotating text fragment 21291/57597
```

```
## 2023-09-16 21:45:02 Annotating text fragment 21301/57597
## 2023-09-16 21:45:02 Annotating text fragment 21311/57597
## 2023-09-16 21:45:02 Annotating text fragment 21321/57597
## 2023-09-16 21:45:02 Annotating text fragment 21331/57597
## 2023-09-16 21:45:03 Annotating text fragment 21341/57597
## 2023-09-16 21:45:03 Annotating text fragment 21351/57597
## 2023-09-16 21:45:03 Annotating text fragment 21361/57597
## 2023-09-16 21:45:03 Annotating text fragment 21371/57597
## 2023-09-16 21:45:03 Annotating text fragment 21381/57597
## 2023-09-16 21:45:03 Annotating text fragment 21391/57597
## 2023-09-16 21:45:03 Annotating text fragment 21401/57597
## 2023-09-16 21:45:03 Annotating text fragment 21411/57597
## 2023-09-16 21:45:03 Annotating text fragment 21421/57597
## 2023-09-16 21:45:03 Annotating text fragment 21431/57597
## 2023-09-16 21:45:03 Annotating text fragment 21441/57597
## 2023-09-16 21:45:03 Annotating text fragment 21451/57597
## 2023-09-16 21:45:03 Annotating text fragment 21461/57597
## 2023-09-16 21:45:03 Annotating text fragment 21471/57597
## 2023-09-16 21:45:03 Annotating text fragment 21481/57597
## 2023-09-16 21:45:04 Annotating text fragment 21491/57597
## 2023-09-16 21:45:04 Annotating text fragment 21501/57597
## 2023-09-16 21:45:04 Annotating text fragment 21511/57597
## 2023-09-16 21:45:04 Annotating text fragment 21521/57597
## 2023-09-16 21:45:04 Annotating text fragment 21531/57597
## 2023-09-16 21:45:04 Annotating text fragment 21541/57597
## 2023-09-16 21:45:04 Annotating text fragment 21551/57597
## 2023-09-16 21:45:04 Annotating text fragment 21561/57597
## 2023-09-16 21:45:04 Annotating text fragment 21571/57597
## 2023-09-16 21:45:04 Annotating text fragment 21581/57597
## 2023-09-16 21:45:04 Annotating text fragment 21591/57597
## 2023-09-16 21:45:04 Annotating text fragment 21601/57597
## 2023-09-16 21:45:04 Annotating text fragment 21611/57597
## 2023-09-16 21:45:04 Annotating text fragment 21621/57597
## 2023-09-16 21:45:04 Annotating text fragment 21631/57597
## 2023-09-16 21:45:04 Annotating text fragment 21641/57597
## 2023-09-16 21:45:04 Annotating text fragment 21651/57597
## 2023-09-16 21:45:04 Annotating text fragment 21661/57597
## 2023-09-16 21:45:04 Annotating text fragment 21671/57597
## 2023-09-16 21:45:04 Annotating text fragment 21681/57597
## 2023-09-16 21:45:05 Annotating text fragment 21691/57597
## 2023-09-16 21:45:05 Annotating text fragment 21701/57597
## 2023-09-16 21:45:05 Annotating text fragment 21711/57597
## 2023-09-16 21:45:05 Annotating text fragment 21721/57597
## 2023-09-16 21:45:05 Annotating text fragment 21731/57597
## 2023-09-16 21:45:06 Annotating text fragment 21741/57597
## 2023-09-16 21:45:06 Annotating text fragment 21751/57597
## 2023-09-16 21:45:06 Annotating text fragment 21761/57597
## 2023-09-16 21:45:06 Annotating text fragment 21771/57597
## 2023-09-16 21:45:06 Annotating text fragment 21781/57597
## 2023-09-16 21:45:06 Annotating text fragment 21791/57597
## 2023-09-16 21:45:06 Annotating text fragment 21801/57597
## 2023-09-16 21:45:06 Annotating text fragment 21811/57597
## 2023-09-16 21:45:06 Annotating text fragment 21821/57597
## 2023-09-16 21:45:06 Annotating text fragment 21831/57597
```

```
## 2023-09-16 21:45:06 Annotating text fragment 21841/57597
## 2023-09-16 21:45:06 Annotating text fragment 21851/57597
## 2023-09-16 21:45:06 Annotating text fragment 21861/57597
## 2023-09-16 21:45:06 Annotating text fragment 21871/57597
## 2023-09-16 21:45:06 Annotating text fragment 21881/57597
## 2023-09-16 21:45:06 Annotating text fragment 21891/57597
## 2023-09-16 21:45:06 Annotating text fragment 21901/57597
## 2023-09-16 21:45:07 Annotating text fragment 21911/57597
## 2023-09-16 21:45:07 Annotating text fragment 21921/57597
## 2023-09-16 21:45:07 Annotating text fragment 21931/57597
## 2023-09-16 21:45:07 Annotating text fragment 21941/57597
## 2023-09-16 21:45:07 Annotating text fragment 21951/57597
## 2023-09-16 21:45:07 Annotating text fragment 21961/57597
## 2023-09-16 21:45:07 Annotating text fragment 21971/57597
## 2023-09-16 21:45:07 Annotating text fragment 21981/57597
## 2023-09-16 21:45:07 Annotating text fragment 21991/57597
## 2023-09-16 21:45:07 Annotating text fragment 22001/57597
## 2023-09-16 21:45:07 Annotating text fragment 22011/57597
## 2023-09-16 21:45:07 Annotating text fragment 22021/57597
## 2023-09-16 21:45:07 Annotating text fragment 22031/57597
## 2023-09-16 21:45:08 Annotating text fragment 22041/57597
## 2023-09-16 21:45:08 Annotating text fragment 22051/57597
## 2023-09-16 21:45:08 Annotating text fragment 22061/57597
## 2023-09-16 21:45:08 Annotating text fragment 22071/57597
## 2023-09-16 21:45:08 Annotating text fragment 22081/57597
## 2023-09-16 21:45:08 Annotating text fragment 22091/57597
## 2023-09-16 21:45:08 Annotating text fragment 22101/57597
## 2023-09-16 21:45:08 Annotating text fragment 22111/57597
## 2023-09-16 21:45:08 Annotating text fragment 22121/57597
## 2023-09-16 21:45:08 Annotating text fragment 22131/57597
## 2023-09-16 21:45:08 Annotating text fragment 22141/57597
## 2023-09-16 21:45:09 Annotating text fragment 22151/57597
## 2023-09-16 21:45:09 Annotating text fragment 22161/57597
## 2023-09-16 21:45:09 Annotating text fragment 22171/57597
## 2023-09-16 21:45:09 Annotating text fragment 22181/57597
## 2023-09-16 21:45:09 Annotating text fragment 22191/57597
## 2023-09-16 21:45:10 Annotating text fragment 22201/57597
## 2023-09-16 21:45:10 Annotating text fragment 22211/57597
## 2023-09-16 21:45:10 Annotating text fragment 22221/57597
## 2023-09-16 21:45:10 Annotating text fragment 22231/57597
## 2023-09-16 21:45:10 Annotating text fragment 22241/57597
## 2023-09-16 21:45:10 Annotating text fragment 22251/57597
## 2023-09-16 21:45:10 Annotating text fragment 22261/57597
## 2023-09-16 21:45:10 Annotating text fragment 22271/57597
## 2023-09-16 21:45:10 Annotating text fragment 22281/57597
## 2023-09-16 21:45:10 Annotating text fragment 22291/57597
## 2023-09-16 21:45:10 Annotating text fragment 22301/57597
## 2023-09-16 21:45:10 Annotating text fragment 22311/57597
## 2023-09-16 21:45:11 Annotating text fragment 22321/57597
## 2023-09-16 21:45:11 Annotating text fragment 22331/57597
## 2023-09-16 21:45:11 Annotating text fragment 22341/57597
## 2023-09-16 21:45:11 Annotating text fragment 22351/57597
## 2023-09-16 21:45:11 Annotating text fragment 22361/57597
## 2023-09-16 21:45:11 Annotating text fragment 22371/57597
```

```
## 2023-09-16 21:45:11 Annotating text fragment 22381/57597
## 2023-09-16 21:45:11 Annotating text fragment 22391/57597
## 2023-09-16 21:45:11 Annotating text fragment 22401/57597
## 2023-09-16 21:45:11 Annotating text fragment 22411/57597
## 2023-09-16 21:45:11 Annotating text fragment 22421/57597
## 2023-09-16 21:45:11 Annotating text fragment 22431/57597
## 2023-09-16 21:45:11 Annotating text fragment 22441/57597
## 2023-09-16 21:45:11 Annotating text fragment 22451/57597
## 2023-09-16 21:45:11 Annotating text fragment 22461/57597
## 2023-09-16 21:45:11 Annotating text fragment 22471/57597
## 2023-09-16 21:45:11 Annotating text fragment 22481/57597
## 2023-09-16 21:45:12 Annotating text fragment 22491/57597
## 2023-09-16 21:45:12 Annotating text fragment 22501/57597
## 2023-09-16 21:45:12 Annotating text fragment 22511/57597
## 2023-09-16 21:45:12 Annotating text fragment 22521/57597
## 2023-09-16 21:45:12 Annotating text fragment 22531/57597
## 2023-09-16 21:45:12 Annotating text fragment 22541/57597
## 2023-09-16 21:45:12 Annotating text fragment 22551/57597
## 2023-09-16 21:45:12 Annotating text fragment 22561/57597
## 2023-09-16 21:45:12 Annotating text fragment 22571/57597
## 2023-09-16 21:45:12 Annotating text fragment 22581/57597
## 2023-09-16 21:45:12 Annotating text fragment 22591/57597
## 2023-09-16 21:45:12 Annotating text fragment 22601/57597
## 2023-09-16 21:45:12 Annotating text fragment 22611/57597
## 2023-09-16 21:45:12 Annotating text fragment 22621/57597
## 2023-09-16 21:45:12 Annotating text fragment 22631/57597
## 2023-09-16 21:45:13 Annotating text fragment 22641/57597
## 2023-09-16 21:45:13 Annotating text fragment 22651/57597
## 2023-09-16 21:45:13 Annotating text fragment 22661/57597
## 2023-09-16 21:45:13 Annotating text fragment 22671/57597
## 2023-09-16 21:45:13 Annotating text fragment 22681/57597
## 2023-09-16 21:45:13 Annotating text fragment 22691/57597
## 2023-09-16 21:45:13 Annotating text fragment 22701/57597
## 2023-09-16 21:45:13 Annotating text fragment 22711/57597
## 2023-09-16 21:45:13 Annotating text fragment 22721/57597
## 2023-09-16 21:45:13 Annotating text fragment 22731/57597
## 2023-09-16 21:45:13 Annotating text fragment 22741/57597
## 2023-09-16 21:45:13 Annotating text fragment 22751/57597
## 2023-09-16 21:45:14 Annotating text fragment 22761/57597
## 2023-09-16 21:45:14 Annotating text fragment 22771/57597
## 2023-09-16 21:45:14 Annotating text fragment 22781/57597
## 2023-09-16 21:45:14 Annotating text fragment 22791/57597
## 2023-09-16 21:45:14 Annotating text fragment 22801/57597
## 2023-09-16 21:45:14 Annotating text fragment 22811/57597
## 2023-09-16 21:45:15 Annotating text fragment 22821/57597
## 2023-09-16 21:45:15 Annotating text fragment 22831/57597
## 2023-09-16 21:45:15 Annotating text fragment 22841/57597
## 2023-09-16 21:45:15 Annotating text fragment 22851/57597
## 2023-09-16 21:45:16 Annotating text fragment 22861/57597
## 2023-09-16 21:45:16 Annotating text fragment 22871/57597
## 2023-09-16 21:45:16 Annotating text fragment 22881/57597
## 2023-09-16 21:45:16 Annotating text fragment 22891/57597
## 2023-09-16 21:45:18 Annotating text fragment 22901/57597
## 2023-09-16 21:45:18 Annotating text fragment 22911/57597
```

```
## 2023-09-16 21:45:18 Annotating text fragment 22921/57597
## 2023-09-16 21:45:19 Annotating text fragment 22931/57597
## 2023-09-16 21:45:19 Annotating text fragment 22941/57597
## 2023-09-16 21:45:19 Annotating text fragment 22951/57597
## 2023-09-16 21:45:19 Annotating text fragment 22961/57597
## 2023-09-16 21:45:19 Annotating text fragment 22971/57597
## 2023-09-16 21:45:19 Annotating text fragment 22981/57597
## 2023-09-16 21:45:19 Annotating text fragment 22991/57597
## 2023-09-16 21:45:19 Annotating text fragment 23001/57597
## 2023-09-16 21:45:19 Annotating text fragment 23011/57597
## 2023-09-16 21:45:20 Annotating text fragment 23021/57597
## 2023-09-16 21:45:20 Annotating text fragment 23031/57597
## 2023-09-16 21:45:20 Annotating text fragment 23041/57597
## 2023-09-16 21:45:20 Annotating text fragment 23051/57597
## 2023-09-16 21:45:20 Annotating text fragment 23061/57597
## 2023-09-16 21:45:20 Annotating text fragment 23071/57597
## 2023-09-16 21:45:20 Annotating text fragment 23081/57597
## 2023-09-16 21:45:20 Annotating text fragment 23091/57597
## 2023-09-16 21:45:20 Annotating text fragment 23101/57597
## 2023-09-16 21:45:20 Annotating text fragment 23111/57597
## 2023-09-16 21:45:20 Annotating text fragment 23121/57597
## 2023-09-16 21:45:20 Annotating text fragment 23131/57597
## 2023-09-16 21:45:20 Annotating text fragment 23141/57597
## 2023-09-16 21:45:20 Annotating text fragment 23151/57597
## 2023-09-16 21:45:21 Annotating text fragment 23161/57597
## 2023-09-16 21:45:21 Annotating text fragment 23171/57597
## 2023-09-16 21:45:21 Annotating text fragment 23181/57597
## 2023-09-16 21:45:21 Annotating text fragment 23191/57597
## 2023-09-16 21:45:22 Annotating text fragment 23201/57597
## 2023-09-16 21:45:22 Annotating text fragment 23211/57597
## 2023-09-16 21:45:22 Annotating text fragment 23221/57597
## 2023-09-16 21:45:22 Annotating text fragment 23231/57597
## 2023-09-16 21:45:22 Annotating text fragment 23241/57597
## 2023-09-16 21:45:22 Annotating text fragment 23251/57597
## 2023-09-16 21:45:22 Annotating text fragment 23261/57597
## 2023-09-16 21:45:22 Annotating text fragment 23271/57597
## 2023-09-16 21:45:22 Annotating text fragment 23281/57597
## 2023-09-16 21:45:22 Annotating text fragment 23291/57597
## 2023-09-16 21:45:22 Annotating text fragment 23301/57597
## 2023-09-16 21:45:22 Annotating text fragment 23311/57597
## 2023-09-16 21:45:23 Annotating text fragment 23321/57597
## 2023-09-16 21:45:23 Annotating text fragment 23331/57597
## 2023-09-16 21:45:23 Annotating text fragment 23341/57597
## 2023-09-16 21:45:23 Annotating text fragment 23351/57597
## 2023-09-16 21:45:23 Annotating text fragment 23361/57597
## 2023-09-16 21:45:23 Annotating text fragment 23371/57597
## 2023-09-16 21:45:23 Annotating text fragment 23381/57597
## 2023-09-16 21:45:23 Annotating text fragment 23391/57597
## 2023-09-16 21:45:23 Annotating text fragment 23401/57597
## 2023-09-16 21:45:23 Annotating text fragment 23411/57597
## 2023-09-16 21:45:23 Annotating text fragment 23421/57597
## 2023-09-16 21:45:24 Annotating text fragment 23431/57597
## 2023-09-16 21:45:24 Annotating text fragment 23441/57597
## 2023-09-16 21:45:24 Annotating text fragment 23451/57597
```

```
## 2023-09-16 21:45:24 Annotating text fragment 23461/57597
## 2023-09-16 21:45:24 Annotating text fragment 23471/57597
## 2023-09-16 21:45:24 Annotating text fragment 23481/57597
## 2023-09-16 21:45:24 Annotating text fragment 23491/57597
## 2023-09-16 21:45:24 Annotating text fragment 23501/57597
## 2023-09-16 21:45:25 Annotating text fragment 23511/57597
## 2023-09-16 21:45:25 Annotating text fragment 23521/57597
## 2023-09-16 21:45:25 Annotating text fragment 23531/57597
## 2023-09-16 21:45:25 Annotating text fragment 23541/57597
## 2023-09-16 21:45:25 Annotating text fragment 23551/57597
## 2023-09-16 21:45:25 Annotating text fragment 23561/57597
## 2023-09-16 21:45:25 Annotating text fragment 23571/57597
## 2023-09-16 21:45:25 Annotating text fragment 23581/57597
## 2023-09-16 21:45:25 Annotating text fragment 23591/57597
## 2023-09-16 21:45:26 Annotating text fragment 23601/57597
## 2023-09-16 21:45:26 Annotating text fragment 23611/57597
## 2023-09-16 21:45:26 Annotating text fragment 23621/57597
## 2023-09-16 21:45:26 Annotating text fragment 23631/57597
## 2023-09-16 21:45:26 Annotating text fragment 23641/57597
## 2023-09-16 21:45:26 Annotating text fragment 23651/57597
## 2023-09-16 21:45:26 Annotating text fragment 23661/57597
## 2023-09-16 21:45:26 Annotating text fragment 23671/57597
## 2023-09-16 21:45:26 Annotating text fragment 23681/57597
## 2023-09-16 21:45:26 Annotating text fragment 23691/57597
## 2023-09-16 21:45:26 Annotating text fragment 23701/57597
## 2023-09-16 21:45:26 Annotating text fragment 23711/57597
## 2023-09-16 21:45:26 Annotating text fragment 23721/57597
## 2023-09-16 21:45:26 Annotating text fragment 23731/57597
## 2023-09-16 21:45:27 Annotating text fragment 23741/57597
## 2023-09-16 21:45:27 Annotating text fragment 23751/57597
## 2023-09-16 21:45:27 Annotating text fragment 23761/57597
## 2023-09-16 21:45:27 Annotating text fragment 23771/57597
## 2023-09-16 21:45:27 Annotating text fragment 23781/57597
## 2023-09-16 21:45:27 Annotating text fragment 23791/57597
## 2023-09-16 21:45:27 Annotating text fragment 23801/57597
## 2023-09-16 21:45:27 Annotating text fragment 23811/57597
## 2023-09-16 21:45:27 Annotating text fragment 23821/57597
## 2023-09-16 21:45:27 Annotating text fragment 23831/57597
## 2023-09-16 21:45:27 Annotating text fragment 23841/57597
## 2023-09-16 21:45:27 Annotating text fragment 23851/57597
## 2023-09-16 21:45:27 Annotating text fragment 23861/57597
## 2023-09-16 21:45:27 Annotating text fragment 23871/57597
## 2023-09-16 21:45:27 Annotating text fragment 23881/57597
## 2023-09-16 21:45:27 Annotating text fragment 23891/57597
## 2023-09-16 21:45:27 Annotating text fragment 23901/57597
## 2023-09-16 21:45:28 Annotating text fragment 23911/57597
## 2023-09-16 21:45:28 Annotating text fragment 23921/57597
## 2023-09-16 21:45:28 Annotating text fragment 23931/57597
## 2023-09-16 21:45:28 Annotating text fragment 23941/57597
## 2023-09-16 21:45:28 Annotating text fragment 23951/57597
## 2023-09-16 21:45:28 Annotating text fragment 23961/57597
## 2023-09-16 21:45:28 Annotating text fragment 23971/57597
## 2023-09-16 21:45:28 Annotating text fragment 23981/57597
## 2023-09-16 21:45:28 Annotating text fragment 23991/57597
```

```
## 2023-09-16 21:45:28 Annotating text fragment 24001/57597
## 2023-09-16 21:45:28 Annotating text fragment 24011/57597
## 2023-09-16 21:45:28 Annotating text fragment 24021/57597
## 2023-09-16 21:45:28 Annotating text fragment 24031/57597
## 2023-09-16 21:45:28 Annotating text fragment 24041/57597
## 2023-09-16 21:45:28 Annotating text fragment 24051/57597
## 2023-09-16 21:45:29 Annotating text fragment 24061/57597
## 2023-09-16 21:45:29 Annotating text fragment 24071/57597
## 2023-09-16 21:45:29 Annotating text fragment 24081/57597
## 2023-09-16 21:45:29 Annotating text fragment 24091/57597
## 2023-09-16 21:45:29 Annotating text fragment 24101/57597
## 2023-09-16 21:45:29 Annotating text fragment 24111/57597
## 2023-09-16 21:45:29 Annotating text fragment 24121/57597
## 2023-09-16 21:45:29 Annotating text fragment 24131/57597
## 2023-09-16 21:45:29 Annotating text fragment 24141/57597
## 2023-09-16 21:45:29 Annotating text fragment 24151/57597
## 2023-09-16 21:45:29 Annotating text fragment 24161/57597
## 2023-09-16 21:45:29 Annotating text fragment 24171/57597
## 2023-09-16 21:45:29 Annotating text fragment 24181/57597
## 2023-09-16 21:45:30 Annotating text fragment 24191/57597
## 2023-09-16 21:45:30 Annotating text fragment 24201/57597
## 2023-09-16 21:45:30 Annotating text fragment 24211/57597
## 2023-09-16 21:45:30 Annotating text fragment 24221/57597
## 2023-09-16 21:45:30 Annotating text fragment 24231/57597
## 2023-09-16 21:45:30 Annotating text fragment 24241/57597
## 2023-09-16 21:45:30 Annotating text fragment 24251/57597
## 2023-09-16 21:45:30 Annotating text fragment 24261/57597
## 2023-09-16 21:45:30 Annotating text fragment 24271/57597
## 2023-09-16 21:45:30 Annotating text fragment 24281/57597
## 2023-09-16 21:45:30 Annotating text fragment 24291/57597
## 2023-09-16 21:45:30 Annotating text fragment 24301/57597
## 2023-09-16 21:45:30 Annotating text fragment 24311/57597
## 2023-09-16 21:45:31 Annotating text fragment 24321/57597
## 2023-09-16 21:45:31 Annotating text fragment 24331/57597
## 2023-09-16 21:45:31 Annotating text fragment 24341/57597
## 2023-09-16 21:45:31 Annotating text fragment 24351/57597
## 2023-09-16 21:45:31 Annotating text fragment 24361/57597
## 2023-09-16 21:45:31 Annotating text fragment 24371/57597
## 2023-09-16 21:45:31 Annotating text fragment 24381/57597
## 2023-09-16 21:45:31 Annotating text fragment 24391/57597
## 2023-09-16 21:45:31 Annotating text fragment 24401/57597
## 2023-09-16 21:45:31 Annotating text fragment 24411/57597
## 2023-09-16 21:45:31 Annotating text fragment 24421/57597
## 2023-09-16 21:45:31 Annotating text fragment 24431/57597
## 2023-09-16 21:45:31 Annotating text fragment 24441/57597
## 2023-09-16 21:45:31 Annotating text fragment 24451/57597
## 2023-09-16 21:45:31 Annotating text fragment 24461/57597
## 2023-09-16 21:45:31 Annotating text fragment 24471/57597
## 2023-09-16 21:45:31 Annotating text fragment 24481/57597
## 2023-09-16 21:45:31 Annotating text fragment 24491/57597
## 2023-09-16 21:45:31 Annotating text fragment 24501/57597
## 2023-09-16 21:45:32 Annotating text fragment 24511/57597
## 2023-09-16 21:45:32 Annotating text fragment 24521/57597
## 2023-09-16 21:45:32 Annotating text fragment 24531/57597
```

```
## 2023-09-16 21:45:32 Annotating text fragment 24541/57597
## 2023-09-16 21:45:32 Annotating text fragment 24551/57597
## 2023-09-16 21:45:32 Annotating text fragment 24561/57597
## 2023-09-16 21:45:32 Annotating text fragment 24571/57597
## 2023-09-16 21:45:32 Annotating text fragment 24581/57597
## 2023-09-16 21:45:32 Annotating text fragment 24591/57597
## 2023-09-16 21:45:32 Annotating text fragment 24601/57597
## 2023-09-16 21:45:32 Annotating text fragment 24611/57597
## 2023-09-16 21:45:32 Annotating text fragment 24621/57597
## 2023-09-16 21:45:32 Annotating text fragment 24631/57597
## 2023-09-16 21:45:32 Annotating text fragment 24641/57597
## 2023-09-16 21:45:32 Annotating text fragment 24651/57597
## 2023-09-16 21:45:33 Annotating text fragment 24661/57597
## 2023-09-16 21:45:33 Annotating text fragment 24671/57597
## 2023-09-16 21:45:33 Annotating text fragment 24681/57597
## 2023-09-16 21:45:33 Annotating text fragment 24691/57597
## 2023-09-16 21:45:33 Annotating text fragment 24701/57597
## 2023-09-16 21:45:33 Annotating text fragment 24711/57597
## 2023-09-16 21:45:33 Annotating text fragment 24721/57597
## 2023-09-16 21:45:33 Annotating text fragment 24731/57597
## 2023-09-16 21:45:33 Annotating text fragment 24741/57597
## 2023-09-16 21:45:33 Annotating text fragment 24751/57597
## 2023-09-16 21:45:33 Annotating text fragment 24761/57597
## 2023-09-16 21:45:34 Annotating text fragment 24771/57597
## 2023-09-16 21:45:34 Annotating text fragment 24781/57597
## 2023-09-16 21:45:34 Annotating text fragment 24791/57597
## 2023-09-16 21:45:34 Annotating text fragment 24801/57597
## 2023-09-16 21:45:34 Annotating text fragment 24811/57597
## 2023-09-16 21:45:34 Annotating text fragment 24821/57597
## 2023-09-16 21:45:34 Annotating text fragment 24831/57597
## 2023-09-16 21:45:34 Annotating text fragment 24841/57597
## 2023-09-16 21:45:34 Annotating text fragment 24851/57597
## 2023-09-16 21:45:34 Annotating text fragment 24861/57597
## 2023-09-16 21:45:34 Annotating text fragment 24871/57597
## 2023-09-16 21:45:34 Annotating text fragment 24881/57597
## 2023-09-16 21:45:35 Annotating text fragment 24891/57597
## 2023-09-16 21:45:35 Annotating text fragment 24901/57597
## 2023-09-16 21:45:35 Annotating text fragment 24911/57597
## 2023-09-16 21:45:35 Annotating text fragment 24921/57597
## 2023-09-16 21:45:35 Annotating text fragment 24931/57597
## 2023-09-16 21:45:35 Annotating text fragment 24941/57597
## 2023-09-16 21:45:35 Annotating text fragment 24951/57597
## 2023-09-16 21:45:35 Annotating text fragment 24961/57597
## 2023-09-16 21:45:35 Annotating text fragment 24971/57597
## 2023-09-16 21:45:35 Annotating text fragment 24981/57597
## 2023-09-16 21:45:35 Annotating text fragment 24991/57597
## 2023-09-16 21:45:35 Annotating text fragment 25001/57597
## 2023-09-16 21:45:36 Annotating text fragment 25011/57597
## 2023-09-16 21:45:36 Annotating text fragment 25021/57597
## 2023-09-16 21:45:36 Annotating text fragment 25031/57597
## 2023-09-16 21:45:36 Annotating text fragment 25041/57597
## 2023-09-16 21:45:36 Annotating text fragment 25051/57597
## 2023-09-16 21:45:36 Annotating text fragment 25061/57597
## 2023-09-16 21:45:36 Annotating text fragment 25071/57597
```

```
## 2023-09-16 21:45:36 Annotating text fragment 25081/57597
## 2023-09-16 21:45:36 Annotating text fragment 25091/57597
## 2023-09-16 21:45:36 Annotating text fragment 25101/57597
## 2023-09-16 21:45:36 Annotating text fragment 25111/57597
## 2023-09-16 21:45:36 Annotating text fragment 25121/57597
## 2023-09-16 21:45:37 Annotating text fragment 25131/57597
## 2023-09-16 21:45:37 Annotating text fragment 25141/57597
## 2023-09-16 21:45:37 Annotating text fragment 25151/57597
## 2023-09-16 21:45:37 Annotating text fragment 25161/57597
## 2023-09-16 21:45:37 Annotating text fragment 25171/57597
## 2023-09-16 21:45:37 Annotating text fragment 25181/57597
## 2023-09-16 21:45:37 Annotating text fragment 25191/57597
## 2023-09-16 21:45:37 Annotating text fragment 25201/57597
## 2023-09-16 21:45:37 Annotating text fragment 25211/57597
## 2023-09-16 21:45:37 Annotating text fragment 25221/57597
## 2023-09-16 21:45:37 Annotating text fragment 25231/57597
## 2023-09-16 21:45:37 Annotating text fragment 25241/57597
## 2023-09-16 21:45:37 Annotating text fragment 25251/57597
## 2023-09-16 21:45:37 Annotating text fragment 25261/57597
## 2023-09-16 21:45:37 Annotating text fragment 25271/57597
## 2023-09-16 21:45:37 Annotating text fragment 25281/57597
## 2023-09-16 21:45:37 Annotating text fragment 25291/57597
## 2023-09-16 21:45:38 Annotating text fragment 25301/57597
## 2023-09-16 21:45:38 Annotating text fragment 25311/57597
## 2023-09-16 21:45:38 Annotating text fragment 25321/57597
## 2023-09-16 21:45:38 Annotating text fragment 25331/57597
## 2023-09-16 21:45:38 Annotating text fragment 25341/57597
## 2023-09-16 21:45:38 Annotating text fragment 25351/57597
## 2023-09-16 21:45:38 Annotating text fragment 25361/57597
## 2023-09-16 21:45:38 Annotating text fragment 25371/57597
## 2023-09-16 21:45:38 Annotating text fragment 25381/57597
## 2023-09-16 21:45:38 Annotating text fragment 25391/57597
## 2023-09-16 21:45:39 Annotating text fragment 25401/57597
## 2023-09-16 21:45:39 Annotating text fragment 25411/57597
## 2023-09-16 21:45:39 Annotating text fragment 25421/57597
## 2023-09-16 21:45:39 Annotating text fragment 25431/57597
## 2023-09-16 21:45:39 Annotating text fragment 25441/57597
## 2023-09-16 21:45:39 Annotating text fragment 25451/57597
## 2023-09-16 21:45:39 Annotating text fragment 25461/57597
## 2023-09-16 21:45:39 Annotating text fragment 25471/57597
## 2023-09-16 21:45:39 Annotating text fragment 25481/57597
## 2023-09-16 21:45:39 Annotating text fragment 25491/57597
## 2023-09-16 21:45:39 Annotating text fragment 25501/57597
## 2023-09-16 21:45:39 Annotating text fragment 25511/57597
## 2023-09-16 21:45:39 Annotating text fragment 25521/57597
## 2023-09-16 21:45:39 Annotating text fragment 25531/57597
## 2023-09-16 21:45:40 Annotating text fragment 25541/57597
## 2023-09-16 21:45:40 Annotating text fragment 25551/57597
## 2023-09-16 21:45:40 Annotating text fragment 25561/57597
## 2023-09-16 21:45:40 Annotating text fragment 25571/57597
## 2023-09-16 21:45:40 Annotating text fragment 25581/57597
## 2023-09-16 21:45:40 Annotating text fragment 25591/57597
## 2023-09-16 21:45:40 Annotating text fragment 25601/57597
## 2023-09-16 21:45:40 Annotating text fragment 25611/57597
```

```
## 2023-09-16 21:45:40 Annotating text fragment 25621/57597
## 2023-09-16 21:45:40 Annotating text fragment 25631/57597
## 2023-09-16 21:45:40 Annotating text fragment 25641/57597
## 2023-09-16 21:45:40 Annotating text fragment 25651/57597
## 2023-09-16 21:45:40 Annotating text fragment 25661/57597
## 2023-09-16 21:45:40 Annotating text fragment 25671/57597
## 2023-09-16 21:45:41 Annotating text fragment 25681/57597
## 2023-09-16 21:45:41 Annotating text fragment 25691/57597
## 2023-09-16 21:45:41 Annotating text fragment 25701/57597
## 2023-09-16 21:45:41 Annotating text fragment 25711/57597
## 2023-09-16 21:45:41 Annotating text fragment 25721/57597
## 2023-09-16 21:45:41 Annotating text fragment 25731/57597
## 2023-09-16 21:45:41 Annotating text fragment 25741/57597
## 2023-09-16 21:45:41 Annotating text fragment 25751/57597
## 2023-09-16 21:45:41 Annotating text fragment 25761/57597
## 2023-09-16 21:45:41 Annotating text fragment 25771/57597
## 2023-09-16 21:45:41 Annotating text fragment 25781/57597
## 2023-09-16 21:45:41 Annotating text fragment 25791/57597
## 2023-09-16 21:45:41 Annotating text fragment 25801/57597
## 2023-09-16 21:45:42 Annotating text fragment 25811/57597
## 2023-09-16 21:45:42 Annotating text fragment 25821/57597
## 2023-09-16 21:45:42 Annotating text fragment 25831/57597
## 2023-09-16 21:45:42 Annotating text fragment 25841/57597
## 2023-09-16 21:45:42 Annotating text fragment 25851/57597
## 2023-09-16 21:45:42 Annotating text fragment 25861/57597
## 2023-09-16 21:45:42 Annotating text fragment 25871/57597
## 2023-09-16 21:45:42 Annotating text fragment 25881/57597
## 2023-09-16 21:45:42 Annotating text fragment 25891/57597
## 2023-09-16 21:45:42 Annotating text fragment 25901/57597
## 2023-09-16 21:45:42 Annotating text fragment 25911/57597
## 2023-09-16 21:45:42 Annotating text fragment 25921/57597
## 2023-09-16 21:45:42 Annotating text fragment 25931/57597
## 2023-09-16 21:45:42 Annotating text fragment 25941/57597
## 2023-09-16 21:45:42 Annotating text fragment 25951/57597
## 2023-09-16 21:45:42 Annotating text fragment 25961/57597
## 2023-09-16 21:45:43 Annotating text fragment 25971/57597
## 2023-09-16 21:45:43 Annotating text fragment 25981/57597
## 2023-09-16 21:45:43 Annotating text fragment 25991/57597
## 2023-09-16 21:45:43 Annotating text fragment 26001/57597
## 2023-09-16 21:45:43 Annotating text fragment 26011/57597
## 2023-09-16 21:45:43 Annotating text fragment 26021/57597
## 2023-09-16 21:45:43 Annotating text fragment 26031/57597
## 2023-09-16 21:45:43 Annotating text fragment 26041/57597
## 2023-09-16 21:45:43 Annotating text fragment 26051/57597
## 2023-09-16 21:45:43 Annotating text fragment 26061/57597
## 2023-09-16 21:45:43 Annotating text fragment 26071/57597
## 2023-09-16 21:45:44 Annotating text fragment 26081/57597
## 2023-09-16 21:45:44 Annotating text fragment 26091/57597
## 2023-09-16 21:45:44 Annotating text fragment 26101/57597
## 2023-09-16 21:45:44 Annotating text fragment 26111/57597
## 2023-09-16 21:45:44 Annotating text fragment 26121/57597
## 2023-09-16 21:45:44 Annotating text fragment 26131/57597
## 2023-09-16 21:45:44 Annotating text fragment 26141/57597
## 2023-09-16 21:45:44 Annotating text fragment 26151/57597
```

```
## 2023-09-16 21:45:44 Annotating text fragment 26161/57597
## 2023-09-16 21:45:44 Annotating text fragment 26171/57597
## 2023-09-16 21:45:44 Annotating text fragment 26181/57597
## 2023-09-16 21:45:44 Annotating text fragment 26191/57597
## 2023-09-16 21:45:44 Annotating text fragment 26201/57597
## 2023-09-16 21:45:44 Annotating text fragment 26211/57597
## 2023-09-16 21:45:45 Annotating text fragment 26221/57597
## 2023-09-16 21:45:45 Annotating text fragment 26231/57597
## 2023-09-16 21:45:45 Annotating text fragment 26241/57597
## 2023-09-16 21:45:45 Annotating text fragment 26251/57597
## 2023-09-16 21:45:45 Annotating text fragment 26261/57597
## 2023-09-16 21:45:45 Annotating text fragment 26271/57597
## 2023-09-16 21:45:45 Annotating text fragment 26281/57597
## 2023-09-16 21:45:45 Annotating text fragment 26291/57597
## 2023-09-16 21:45:46 Annotating text fragment 26301/57597
## 2023-09-16 21:45:46 Annotating text fragment 26311/57597
## 2023-09-16 21:45:46 Annotating text fragment 26321/57597
## 2023-09-16 21:45:46 Annotating text fragment 26331/57597
## 2023-09-16 21:45:46 Annotating text fragment 26341/57597
## 2023-09-16 21:45:46 Annotating text fragment 26351/57597
## 2023-09-16 21:45:46 Annotating text fragment 26361/57597
## 2023-09-16 21:45:46 Annotating text fragment 26371/57597
## 2023-09-16 21:45:46 Annotating text fragment 26381/57597
## 2023-09-16 21:45:46 Annotating text fragment 26391/57597
## 2023-09-16 21:45:46 Annotating text fragment 26401/57597
## 2023-09-16 21:45:47 Annotating text fragment 26411/57597
## 2023-09-16 21:45:47 Annotating text fragment 26421/57597
## 2023-09-16 21:45:47 Annotating text fragment 26431/57597
## 2023-09-16 21:45:47 Annotating text fragment 26441/57597
## 2023-09-16 21:45:47 Annotating text fragment 26451/57597
## 2023-09-16 21:45:47 Annotating text fragment 26461/57597
## 2023-09-16 21:45:47 Annotating text fragment 26471/57597
## 2023-09-16 21:45:47 Annotating text fragment 26481/57597
## 2023-09-16 21:45:47 Annotating text fragment 26491/57597
## 2023-09-16 21:45:47 Annotating text fragment 26501/57597
## 2023-09-16 21:45:47 Annotating text fragment 26511/57597
## 2023-09-16 21:45:47 Annotating text fragment 26521/57597
## 2023-09-16 21:45:47 Annotating text fragment 26531/57597
## 2023-09-16 21:45:47 Annotating text fragment 26541/57597
## 2023-09-16 21:45:47 Annotating text fragment 26551/57597
## 2023-09-16 21:45:47 Annotating text fragment 26561/57597
## 2023-09-16 21:45:47 Annotating text fragment 26571/57597
## 2023-09-16 21:45:47 Annotating text fragment 26581/57597
## 2023-09-16 21:45:47 Annotating text fragment 26591/57597
## 2023-09-16 21:45:48 Annotating text fragment 26601/57597
## 2023-09-16 21:45:48 Annotating text fragment 26611/57597
## 2023-09-16 21:45:48 Annotating text fragment 26621/57597
## 2023-09-16 21:45:48 Annotating text fragment 26631/57597
## 2023-09-16 21:45:48 Annotating text fragment 26641/57597
## 2023-09-16 21:45:48 Annotating text fragment 26651/57597
## 2023-09-16 21:45:48 Annotating text fragment 26661/57597
## 2023-09-16 21:45:48 Annotating text fragment 26671/57597
## 2023-09-16 21:45:48 Annotating text fragment 26681/57597
## 2023-09-16 21:45:48 Annotating text fragment 26691/57597
```

```
## 2023-09-16 21:45:49 Annotating text fragment 26701/57597
## 2023-09-16 21:45:50 Annotating text fragment 26711/57597
## 2023-09-16 21:45:50 Annotating text fragment 26721/57597
## 2023-09-16 21:45:50 Annotating text fragment 26731/57597
## 2023-09-16 21:45:51 Annotating text fragment 26741/57597
## 2023-09-16 21:45:51 Annotating text fragment 26751/57597
## 2023-09-16 21:45:51 Annotating text fragment 26761/57597
## 2023-09-16 21:45:51 Annotating text fragment 26771/57597
## 2023-09-16 21:45:51 Annotating text fragment 26781/57597
## 2023-09-16 21:45:51 Annotating text fragment 26791/57597
## 2023-09-16 21:45:51 Annotating text fragment 26801/57597
## 2023-09-16 21:45:51 Annotating text fragment 26811/57597
## 2023-09-16 21:45:51 Annotating text fragment 26821/57597
## 2023-09-16 21:45:51 Annotating text fragment 26831/57597
## 2023-09-16 21:45:51 Annotating text fragment 26841/57597
## 2023-09-16 21:45:51 Annotating text fragment 26851/57597
## 2023-09-16 21:45:51 Annotating text fragment 26861/57597
## 2023-09-16 21:45:51 Annotating text fragment 26871/57597
## 2023-09-16 21:45:52 Annotating text fragment 26881/57597
## 2023-09-16 21:45:52 Annotating text fragment 26891/57597
## 2023-09-16 21:45:52 Annotating text fragment 26901/57597
## 2023-09-16 21:45:53 Annotating text fragment 26911/57597
## 2023-09-16 21:45:53 Annotating text fragment 26921/57597
## 2023-09-16 21:45:53 Annotating text fragment 26931/57597
## 2023-09-16 21:45:53 Annotating text fragment 26941/57597
## 2023-09-16 21:45:53 Annotating text fragment 26951/57597
## 2023-09-16 21:45:53 Annotating text fragment 26961/57597
## 2023-09-16 21:45:53 Annotating text fragment 26971/57597
## 2023-09-16 21:45:53 Annotating text fragment 26981/57597
## 2023-09-16 21:45:53 Annotating text fragment 26991/57597
## 2023-09-16 21:45:53 Annotating text fragment 27001/57597
## 2023-09-16 21:45:54 Annotating text fragment 27011/57597
## 2023-09-16 21:45:54 Annotating text fragment 27021/57597
## 2023-09-16 21:45:54 Annotating text fragment 27031/57597
## 2023-09-16 21:45:54 Annotating text fragment 27041/57597
## 2023-09-16 21:45:54 Annotating text fragment 27051/57597
## 2023-09-16 21:45:54 Annotating text fragment 27061/57597
## 2023-09-16 21:45:54 Annotating text fragment 27071/57597
## 2023-09-16 21:45:54 Annotating text fragment 27081/57597
## 2023-09-16 21:45:54 Annotating text fragment 27091/57597
## 2023-09-16 21:45:54 Annotating text fragment 27101/57597
## 2023-09-16 21:45:54 Annotating text fragment 27111/57597
## 2023-09-16 21:45:54 Annotating text fragment 27121/57597
## 2023-09-16 21:45:54 Annotating text fragment 27131/57597
## 2023-09-16 21:45:54 Annotating text fragment 27141/57597
## 2023-09-16 21:45:54 Annotating text fragment 27151/57597
## 2023-09-16 21:45:55 Annotating text fragment 27161/57597
## 2023-09-16 21:45:55 Annotating text fragment 27171/57597
## 2023-09-16 21:45:55 Annotating text fragment 27181/57597
## 2023-09-16 21:45:55 Annotating text fragment 27191/57597
## 2023-09-16 21:45:55 Annotating text fragment 27201/57597
## 2023-09-16 21:45:55 Annotating text fragment 27211/57597
## 2023-09-16 21:45:55 Annotating text fragment 27221/57597
## 2023-09-16 21:45:55 Annotating text fragment 27231/57597
```

```
## 2023-09-16 21:45:55 Annotating text fragment 27241/57597
## 2023-09-16 21:45:55 Annotating text fragment 27251/57597
## 2023-09-16 21:45:55 Annotating text fragment 27261/57597
## 2023-09-16 21:45:56 Annotating text fragment 27271/57597
## 2023-09-16 21:45:56 Annotating text fragment 27281/57597
## 2023-09-16 21:45:57 Annotating text fragment 27291/57597
## 2023-09-16 21:45:57 Annotating text fragment 27301/57597
## 2023-09-16 21:45:57 Annotating text fragment 27311/57597
## 2023-09-16 21:45:57 Annotating text fragment 27321/57597
## 2023-09-16 21:45:57 Annotating text fragment 27331/57597
## 2023-09-16 21:45:57 Annotating text fragment 27341/57597
## 2023-09-16 21:45:57 Annotating text fragment 27351/57597
## 2023-09-16 21:45:57 Annotating text fragment 27361/57597
## 2023-09-16 21:45:57 Annotating text fragment 27371/57597
## 2023-09-16 21:45:57 Annotating text fragment 27381/57597
## 2023-09-16 21:45:57 Annotating text fragment 27391/57597
## 2023-09-16 21:45:57 Annotating text fragment 27401/57597
## 2023-09-16 21:45:57 Annotating text fragment 27411/57597
## 2023-09-16 21:45:57 Annotating text fragment 27421/57597
## 2023-09-16 21:45:57 Annotating text fragment 27431/57597
## 2023-09-16 21:45:58 Annotating text fragment 27441/57597
## 2023-09-16 21:45:58 Annotating text fragment 27451/57597
## 2023-09-16 21:45:58 Annotating text fragment 27461/57597
## 2023-09-16 21:45:58 Annotating text fragment 27471/57597
## 2023-09-16 21:45:58 Annotating text fragment 27481/57597
## 2023-09-16 21:45:58 Annotating text fragment 27491/57597
## 2023-09-16 21:45:58 Annotating text fragment 27501/57597
## 2023-09-16 21:45:58 Annotating text fragment 27511/57597
## 2023-09-16 21:45:58 Annotating text fragment 27521/57597
## 2023-09-16 21:45:58 Annotating text fragment 27531/57597
## 2023-09-16 21:45:58 Annotating text fragment 27541/57597
## 2023-09-16 21:45:58 Annotating text fragment 27551/57597
## 2023-09-16 21:45:58 Annotating text fragment 27561/57597
## 2023-09-16 21:45:58 Annotating text fragment 27571/57597
## 2023-09-16 21:45:58 Annotating text fragment 27581/57597
## 2023-09-16 21:45:59 Annotating text fragment 27591/57597
## 2023-09-16 21:45:59 Annotating text fragment 27601/57597
## 2023-09-16 21:45:59 Annotating text fragment 27611/57597
## 2023-09-16 21:45:59 Annotating text fragment 27621/57597
## 2023-09-16 21:45:59 Annotating text fragment 27631/57597
## 2023-09-16 21:45:59 Annotating text fragment 27641/57597
## 2023-09-16 21:45:59 Annotating text fragment 27651/57597
## 2023-09-16 21:45:59 Annotating text fragment 27661/57597
## 2023-09-16 21:45:59 Annotating text fragment 27671/57597
## 2023-09-16 21:45:59 Annotating text fragment 27681/57597
## 2023-09-16 21:45:59 Annotating text fragment 27691/57597
## 2023-09-16 21:45:59 Annotating text fragment 27701/57597
## 2023-09-16 21:45:59 Annotating text fragment 27711/57597
## 2023-09-16 21:46:00 Annotating text fragment 27721/57597
## 2023-09-16 21:46:00 Annotating text fragment 27731/57597
## 2023-09-16 21:46:00 Annotating text fragment 27741/57597
## 2023-09-16 21:46:00 Annotating text fragment 27751/57597
## 2023-09-16 21:46:00 Annotating text fragment 27761/57597
## 2023-09-16 21:46:00 Annotating text fragment 27771/57597
```

```
## 2023-09-16 21:46:00 Annotating text fragment 27781/57597
## 2023-09-16 21:46:00 Annotating text fragment 27791/57597
## 2023-09-16 21:46:00 Annotating text fragment 27801/57597
## 2023-09-16 21:46:01 Annotating text fragment 27811/57597
## 2023-09-16 21:46:01 Annotating text fragment 27821/57597
## 2023-09-16 21:46:01 Annotating text fragment 27831/57597
## 2023-09-16 21:46:01 Annotating text fragment 27841/57597
## 2023-09-16 21:46:01 Annotating text fragment 27851/57597
## 2023-09-16 21:46:01 Annotating text fragment 27861/57597
## 2023-09-16 21:46:01 Annotating text fragment 27871/57597
## 2023-09-16 21:46:01 Annotating text fragment 27881/57597
## 2023-09-16 21:46:01 Annotating text fragment 27891/57597
## 2023-09-16 21:46:01 Annotating text fragment 27901/57597
## 2023-09-16 21:46:01 Annotating text fragment 27911/57597
## 2023-09-16 21:46:01 Annotating text fragment 27921/57597
## 2023-09-16 21:46:01 Annotating text fragment 27931/57597
## 2023-09-16 21:46:01 Annotating text fragment 27941/57597
## 2023-09-16 21:46:01 Annotating text fragment 27951/57597
## 2023-09-16 21:46:01 Annotating text fragment 27961/57597
## 2023-09-16 21:46:02 Annotating text fragment 27971/57597
## 2023-09-16 21:46:02 Annotating text fragment 27981/57597
## 2023-09-16 21:46:02 Annotating text fragment 27991/57597
## 2023-09-16 21:46:02 Annotating text fragment 28001/57597
## 2023-09-16 21:46:02 Annotating text fragment 28011/57597
## 2023-09-16 21:46:02 Annotating text fragment 28021/57597
## 2023-09-16 21:46:02 Annotating text fragment 28031/57597
## 2023-09-16 21:46:02 Annotating text fragment 28041/57597
## 2023-09-16 21:46:02 Annotating text fragment 28051/57597
## 2023-09-16 21:46:02 Annotating text fragment 28061/57597
## 2023-09-16 21:46:02 Annotating text fragment 28071/57597
## 2023-09-16 21:46:03 Annotating text fragment 28081/57597
## 2023-09-16 21:46:03 Annotating text fragment 28091/57597
## 2023-09-16 21:46:03 Annotating text fragment 28101/57597
## 2023-09-16 21:46:03 Annotating text fragment 28111/57597
## 2023-09-16 21:46:03 Annotating text fragment 28121/57597
## 2023-09-16 21:46:03 Annotating text fragment 28131/57597
## 2023-09-16 21:46:03 Annotating text fragment 28141/57597
## 2023-09-16 21:46:03 Annotating text fragment 28151/57597
## 2023-09-16 21:46:03 Annotating text fragment 28161/57597
## 2023-09-16 21:46:03 Annotating text fragment 28171/57597
## 2023-09-16 21:46:03 Annotating text fragment 28181/57597
## 2023-09-16 21:46:03 Annotating text fragment 28191/57597
## 2023-09-16 21:46:03 Annotating text fragment 28201/57597
## 2023-09-16 21:46:03 Annotating text fragment 28211/57597
## 2023-09-16 21:46:03 Annotating text fragment 28221/57597
## 2023-09-16 21:46:04 Annotating text fragment 28231/57597
## 2023-09-16 21:46:04 Annotating text fragment 28241/57597
## 2023-09-16 21:46:04 Annotating text fragment 28251/57597
## 2023-09-16 21:46:04 Annotating text fragment 28261/57597
## 2023-09-16 21:46:04 Annotating text fragment 28271/57597
## 2023-09-16 21:46:04 Annotating text fragment 28281/57597
## 2023-09-16 21:46:04 Annotating text fragment 28291/57597
## 2023-09-16 21:46:04 Annotating text fragment 28301/57597
## 2023-09-16 21:46:04 Annotating text fragment 28311/57597
```

```
## 2023-09-16 21:46:04 Annotating text fragment 28321/57597
## 2023-09-16 21:46:05 Annotating text fragment 28331/57597
## 2023-09-16 21:46:05 Annotating text fragment 28341/57597
## 2023-09-16 21:46:05 Annotating text fragment 28351/57597
## 2023-09-16 21:46:05 Annotating text fragment 28361/57597
## 2023-09-16 21:46:05 Annotating text fragment 28371/57597
## 2023-09-16 21:46:05 Annotating text fragment 28381/57597
## 2023-09-16 21:46:05 Annotating text fragment 28391/57597
## 2023-09-16 21:46:05 Annotating text fragment 28401/57597
## 2023-09-16 21:46:05 Annotating text fragment 28411/57597
## 2023-09-16 21:46:05 Annotating text fragment 28421/57597
## 2023-09-16 21:46:05 Annotating text fragment 28431/57597
## 2023-09-16 21:46:05 Annotating text fragment 28441/57597
## 2023-09-16 21:46:06 Annotating text fragment 28451/57597
## 2023-09-16 21:46:06 Annotating text fragment 28461/57597
## 2023-09-16 21:46:06 Annotating text fragment 28471/57597
## 2023-09-16 21:46:06 Annotating text fragment 28481/57597
## 2023-09-16 21:46:06 Annotating text fragment 28491/57597
## 2023-09-16 21:46:06 Annotating text fragment 28501/57597
## 2023-09-16 21:46:07 Annotating text fragment 28511/57597
## 2023-09-16 21:46:07 Annotating text fragment 28521/57597
## 2023-09-16 21:46:07 Annotating text fragment 28531/57597
## 2023-09-16 21:46:07 Annotating text fragment 28541/57597
## 2023-09-16 21:46:07 Annotating text fragment 28551/57597
## 2023-09-16 21:46:07 Annotating text fragment 28561/57597
## 2023-09-16 21:46:07 Annotating text fragment 28571/57597
## 2023-09-16 21:46:07 Annotating text fragment 28581/57597
## 2023-09-16 21:46:07 Annotating text fragment 28591/57597
## 2023-09-16 21:46:07 Annotating text fragment 28601/57597
## 2023-09-16 21:46:07 Annotating text fragment 28611/57597
## 2023-09-16 21:46:07 Annotating text fragment 28621/57597
## 2023-09-16 21:46:07 Annotating text fragment 28631/57597
## 2023-09-16 21:46:08 Annotating text fragment 28641/57597
## 2023-09-16 21:46:08 Annotating text fragment 28651/57597
## 2023-09-16 21:46:08 Annotating text fragment 28661/57597
## 2023-09-16 21:46:08 Annotating text fragment 28671/57597
## 2023-09-16 21:46:08 Annotating text fragment 28681/57597
## 2023-09-16 21:46:08 Annotating text fragment 28691/57597
## 2023-09-16 21:46:08 Annotating text fragment 28701/57597
## 2023-09-16 21:46:08 Annotating text fragment 28711/57597
## 2023-09-16 21:46:08 Annotating text fragment 28721/57597
## 2023-09-16 21:46:08 Annotating text fragment 28731/57597
## 2023-09-16 21:46:09 Annotating text fragment 28741/57597
## 2023-09-16 21:46:09 Annotating text fragment 28751/57597
## 2023-09-16 21:46:09 Annotating text fragment 28761/57597
## 2023-09-16 21:46:09 Annotating text fragment 28771/57597
## 2023-09-16 21:46:09 Annotating text fragment 28781/57597
## 2023-09-16 21:46:09 Annotating text fragment 28791/57597
## 2023-09-16 21:46:09 Annotating text fragment 28801/57597
## 2023-09-16 21:46:09 Annotating text fragment 28811/57597
## 2023-09-16 21:46:09 Annotating text fragment 28821/57597
## 2023-09-16 21:46:09 Annotating text fragment 28831/57597
## 2023-09-16 21:46:09 Annotating text fragment 28841/57597
## 2023-09-16 21:46:10 Annotating text fragment 28851/57597
```

```
## 2023-09-16 21:46:10 Annotating text fragment 28861/57597
## 2023-09-16 21:46:10 Annotating text fragment 28871/57597
## 2023-09-16 21:46:10 Annotating text fragment 28881/57597
## 2023-09-16 21:46:10 Annotating text fragment 28891/57597
## 2023-09-16 21:46:10 Annotating text fragment 28901/57597
## 2023-09-16 21:46:11 Annotating text fragment 28911/57597
## 2023-09-16 21:46:11 Annotating text fragment 28921/57597
## 2023-09-16 21:46:11 Annotating text fragment 28931/57597
## 2023-09-16 21:46:11 Annotating text fragment 28941/57597
## 2023-09-16 21:46:12 Annotating text fragment 28951/57597
## 2023-09-16 21:46:13 Annotating text fragment 28961/57597
## 2023-09-16 21:46:14 Annotating text fragment 28971/57597
## 2023-09-16 21:46:14 Annotating text fragment 28981/57597
## 2023-09-16 21:46:14 Annotating text fragment 28991/57597
## 2023-09-16 21:46:15 Annotating text fragment 29001/57597
## 2023-09-16 21:46:15 Annotating text fragment 29011/57597
## 2023-09-16 21:46:15 Annotating text fragment 29021/57597
## 2023-09-16 21:46:15 Annotating text fragment 29031/57597
## 2023-09-16 21:46:15 Annotating text fragment 29041/57597
## 2023-09-16 21:46:15 Annotating text fragment 29051/57597
## 2023-09-16 21:46:15 Annotating text fragment 29061/57597
## 2023-09-16 21:46:15 Annotating text fragment 29071/57597
## 2023-09-16 21:46:15 Annotating text fragment 29081/57597
## 2023-09-16 21:46:15 Annotating text fragment 29091/57597
## 2023-09-16 21:46:15 Annotating text fragment 29101/57597
## 2023-09-16 21:46:15 Annotating text fragment 29111/57597
## 2023-09-16 21:46:15 Annotating text fragment 29121/57597
## 2023-09-16 21:46:16 Annotating text fragment 29131/57597
## 2023-09-16 21:46:16 Annotating text fragment 29141/57597
## 2023-09-16 21:46:16 Annotating text fragment 29151/57597
## 2023-09-16 21:46:16 Annotating text fragment 29161/57597
## 2023-09-16 21:46:16 Annotating text fragment 29171/57597
## 2023-09-16 21:46:16 Annotating text fragment 29181/57597
## 2023-09-16 21:46:16 Annotating text fragment 29191/57597
## 2023-09-16 21:46:16 Annotating text fragment 29201/57597
## 2023-09-16 21:46:16 Annotating text fragment 29211/57597
## 2023-09-16 21:46:16 Annotating text fragment 29221/57597
## 2023-09-16 21:46:16 Annotating text fragment 29231/57597
## 2023-09-16 21:46:16 Annotating text fragment 29241/57597
## 2023-09-16 21:46:16 Annotating text fragment 29251/57597
## 2023-09-16 21:46:16 Annotating text fragment 29261/57597
## 2023-09-16 21:46:17 Annotating text fragment 29271/57597
## 2023-09-16 21:46:17 Annotating text fragment 29281/57597
## 2023-09-16 21:46:17 Annotating text fragment 29291/57597
## 2023-09-16 21:46:17 Annotating text fragment 29301/57597
## 2023-09-16 21:46:17 Annotating text fragment 29311/57597
## 2023-09-16 21:46:17 Annotating text fragment 29321/57597
## 2023-09-16 21:46:17 Annotating text fragment 29331/57597
## 2023-09-16 21:46:17 Annotating text fragment 29341/57597
## 2023-09-16 21:46:17 Annotating text fragment 29351/57597
## 2023-09-16 21:46:17 Annotating text fragment 29361/57597
## 2023-09-16 21:46:17 Annotating text fragment 29371/57597
## 2023-09-16 21:46:17 Annotating text fragment 29381/57597
## 2023-09-16 21:46:17 Annotating text fragment 29391/57597
```

```
## 2023-09-16 21:46:17 Annotating text fragment 29401/57597
## 2023-09-16 21:46:18 Annotating text fragment 29411/57597
## 2023-09-16 21:46:18 Annotating text fragment 29421/57597
## 2023-09-16 21:46:18 Annotating text fragment 29431/57597
## 2023-09-16 21:46:18 Annotating text fragment 29441/57597
## 2023-09-16 21:46:18 Annotating text fragment 29451/57597
## 2023-09-16 21:46:18 Annotating text fragment 29461/57597
## 2023-09-16 21:46:19 Annotating text fragment 29471/57597
## 2023-09-16 21:46:19 Annotating text fragment 29481/57597
## 2023-09-16 21:46:19 Annotating text fragment 29491/57597
## 2023-09-16 21:46:19 Annotating text fragment 29501/57597
## 2023-09-16 21:46:19 Annotating text fragment 29511/57597
## 2023-09-16 21:46:19 Annotating text fragment 29521/57597
## 2023-09-16 21:46:19 Annotating text fragment 29531/57597
## 2023-09-16 21:46:19 Annotating text fragment 29541/57597
## 2023-09-16 21:46:19 Annotating text fragment 29551/57597
## 2023-09-16 21:46:19 Annotating text fragment 29561/57597
## 2023-09-16 21:46:19 Annotating text fragment 29571/57597
## 2023-09-16 21:46:19 Annotating text fragment 29581/57597
## 2023-09-16 21:46:19 Annotating text fragment 29591/57597
## 2023-09-16 21:46:19 Annotating text fragment 29601/57597
## 2023-09-16 21:46:19 Annotating text fragment 29611/57597
## 2023-09-16 21:46:20 Annotating text fragment 29621/57597
## 2023-09-16 21:46:20 Annotating text fragment 29631/57597
## 2023-09-16 21:46:20 Annotating text fragment 29641/57597
## 2023-09-16 21:46:20 Annotating text fragment 29651/57597
## 2023-09-16 21:46:20 Annotating text fragment 29661/57597
## 2023-09-16 21:46:20 Annotating text fragment 29671/57597
## 2023-09-16 21:46:20 Annotating text fragment 29681/57597
## 2023-09-16 21:46:20 Annotating text fragment 29691/57597
## 2023-09-16 21:46:21 Annotating text fragment 29701/57597
## 2023-09-16 21:46:21 Annotating text fragment 29711/57597
## 2023-09-16 21:46:21 Annotating text fragment 29721/57597
## 2023-09-16 21:46:21 Annotating text fragment 29731/57597
## 2023-09-16 21:46:21 Annotating text fragment 29741/57597
## 2023-09-16 21:46:21 Annotating text fragment 29751/57597
## 2023-09-16 21:46:21 Annotating text fragment 29761/57597
## 2023-09-16 21:46:21 Annotating text fragment 29771/57597
## 2023-09-16 21:46:21 Annotating text fragment 29781/57597
## 2023-09-16 21:46:21 Annotating text fragment 29791/57597
## 2023-09-16 21:46:21 Annotating text fragment 29801/57597
## 2023-09-16 21:46:21 Annotating text fragment 29811/57597
## 2023-09-16 21:46:21 Annotating text fragment 29821/57597
## 2023-09-16 21:46:21 Annotating text fragment 29831/57597
## 2023-09-16 21:46:21 Annotating text fragment 29841/57597
## 2023-09-16 21:46:22 Annotating text fragment 29851/57597
## 2023-09-16 21:46:22 Annotating text fragment 29861/57597
## 2023-09-16 21:46:22 Annotating text fragment 29871/57597
## 2023-09-16 21:46:22 Annotating text fragment 29881/57597
## 2023-09-16 21:46:22 Annotating text fragment 29891/57597
## 2023-09-16 21:46:22 Annotating text fragment 29901/57597
## 2023-09-16 21:46:22 Annotating text fragment 29911/57597
## 2023-09-16 21:46:22 Annotating text fragment 29921/57597
## 2023-09-16 21:46:22 Annotating text fragment 29931/57597
```

```
## 2023-09-16 21:46:22 Annotating text fragment 29941/57597
## 2023-09-16 21:46:22 Annotating text fragment 29951/57597
## 2023-09-16 21:46:22 Annotating text fragment 29961/57597
## 2023-09-16 21:46:23 Annotating text fragment 29971/57597
## 2023-09-16 21:46:23 Annotating text fragment 29981/57597
## 2023-09-16 21:46:23 Annotating text fragment 29991/57597
## 2023-09-16 21:46:23 Annotating text fragment 30001/57597
## 2023-09-16 21:46:23 Annotating text fragment 30011/57597
## 2023-09-16 21:46:23 Annotating text fragment 30021/57597
## 2023-09-16 21:46:23 Annotating text fragment 30031/57597
## 2023-09-16 21:46:23 Annotating text fragment 30041/57597
## 2023-09-16 21:46:23 Annotating text fragment 30051/57597
## 2023-09-16 21:46:23 Annotating text fragment 30061/57597
## 2023-09-16 21:46:23 Annotating text fragment 30071/57597
## 2023-09-16 21:46:23 Annotating text fragment 30081/57597
## 2023-09-16 21:46:23 Annotating text fragment 30091/57597
## 2023-09-16 21:46:24 Annotating text fragment 30101/57597
## 2023-09-16 21:46:24 Annotating text fragment 30111/57597
## 2023-09-16 21:46:24 Annotating text fragment 30121/57597
## 2023-09-16 21:46:24 Annotating text fragment 30131/57597
## 2023-09-16 21:46:24 Annotating text fragment 30141/57597
## 2023-09-16 21:46:24 Annotating text fragment 30151/57597
## 2023-09-16 21:46:24 Annotating text fragment 30161/57597
## 2023-09-16 21:46:24 Annotating text fragment 30171/57597
## 2023-09-16 21:46:24 Annotating text fragment 30181/57597
## 2023-09-16 21:46:24 Annotating text fragment 30191/57597
## 2023-09-16 21:46:24 Annotating text fragment 30201/57597
## 2023-09-16 21:46:25 Annotating text fragment 30211/57597
## 2023-09-16 21:46:25 Annotating text fragment 30221/57597
## 2023-09-16 21:46:25 Annotating text fragment 30231/57597
## 2023-09-16 21:46:25 Annotating text fragment 30241/57597
## 2023-09-16 21:46:25 Annotating text fragment 30251/57597
## 2023-09-16 21:46:25 Annotating text fragment 30261/57597
## 2023-09-16 21:46:25 Annotating text fragment 30271/57597
## 2023-09-16 21:46:25 Annotating text fragment 30281/57597
## 2023-09-16 21:46:25 Annotating text fragment 30291/57597
## 2023-09-16 21:46:25 Annotating text fragment 30301/57597
## 2023-09-16 21:46:26 Annotating text fragment 30311/57597
## 2023-09-16 21:46:26 Annotating text fragment 30321/57597
## 2023-09-16 21:46:26 Annotating text fragment 30331/57597
## 2023-09-16 21:46:26 Annotating text fragment 30341/57597
## 2023-09-16 21:46:26 Annotating text fragment 30351/57597
## 2023-09-16 21:46:26 Annotating text fragment 30361/57597
## 2023-09-16 21:46:26 Annotating text fragment 30371/57597
## 2023-09-16 21:46:26 Annotating text fragment 30381/57597
## 2023-09-16 21:46:27 Annotating text fragment 30391/57597
## 2023-09-16 21:46:27 Annotating text fragment 30401/57597
## 2023-09-16 21:46:27 Annotating text fragment 30411/57597
## 2023-09-16 21:46:27 Annotating text fragment 30421/57597
## 2023-09-16 21:46:27 Annotating text fragment 30431/57597
## 2023-09-16 21:46:27 Annotating text fragment 30441/57597
## 2023-09-16 21:46:27 Annotating text fragment 30451/57597
## 2023-09-16 21:46:27 Annotating text fragment 30461/57597
## 2023-09-16 21:46:27 Annotating text fragment 30471/57597
```

```
## 2023-09-16 21:46:27 Annotating text fragment 30481/57597
## 2023-09-16 21:46:27 Annotating text fragment 30491/57597
## 2023-09-16 21:46:27 Annotating text fragment 30501/57597
## 2023-09-16 21:46:28 Annotating text fragment 30511/57597
## 2023-09-16 21:46:28 Annotating text fragment 30521/57597
## 2023-09-16 21:46:28 Annotating text fragment 30531/57597
## 2023-09-16 21:46:28 Annotating text fragment 30541/57597
## 2023-09-16 21:46:28 Annotating text fragment 30551/57597
## 2023-09-16 21:46:28 Annotating text fragment 30561/57597
## 2023-09-16 21:46:28 Annotating text fragment 30571/57597
## 2023-09-16 21:46:28 Annotating text fragment 30581/57597
## 2023-09-16 21:46:28 Annotating text fragment 30591/57597
## 2023-09-16 21:46:28 Annotating text fragment 30601/57597
## 2023-09-16 21:46:28 Annotating text fragment 30611/57597
## 2023-09-16 21:46:28 Annotating text fragment 30621/57597
## 2023-09-16 21:46:28 Annotating text fragment 30631/57597
## 2023-09-16 21:46:28 Annotating text fragment 30641/57597
## 2023-09-16 21:46:29 Annotating text fragment 30651/57597
## 2023-09-16 21:46:29 Annotating text fragment 30661/57597
## 2023-09-16 21:46:29 Annotating text fragment 30671/57597
## 2023-09-16 21:46:29 Annotating text fragment 30681/57597
## 2023-09-16 21:46:29 Annotating text fragment 30691/57597
## 2023-09-16 21:46:29 Annotating text fragment 30701/57597
## 2023-09-16 21:46:29 Annotating text fragment 30711/57597
## 2023-09-16 21:46:29 Annotating text fragment 30721/57597
## 2023-09-16 21:46:29 Annotating text fragment 30731/57597
## 2023-09-16 21:46:29 Annotating text fragment 30741/57597
## 2023-09-16 21:46:29 Annotating text fragment 30751/57597
## 2023-09-16 21:46:29 Annotating text fragment 30761/57597
## 2023-09-16 21:46:30 Annotating text fragment 30771/57597
## 2023-09-16 21:46:30 Annotating text fragment 30781/57597
## 2023-09-16 21:46:30 Annotating text fragment 30791/57597
## 2023-09-16 21:46:30 Annotating text fragment 30801/57597
## 2023-09-16 21:46:30 Annotating text fragment 30811/57597
## 2023-09-16 21:46:30 Annotating text fragment 30821/57597
## 2023-09-16 21:46:30 Annotating text fragment 30831/57597
## 2023-09-16 21:46:30 Annotating text fragment 30841/57597
## 2023-09-16 21:46:30 Annotating text fragment 30851/57597
## 2023-09-16 21:46:30 Annotating text fragment 30861/57597
## 2023-09-16 21:46:30 Annotating text fragment 30871/57597
## 2023-09-16 21:46:31 Annotating text fragment 30881/57597
## 2023-09-16 21:46:31 Annotating text fragment 30891/57597
## 2023-09-16 21:46:31 Annotating text fragment 30901/57597
## 2023-09-16 21:46:31 Annotating text fragment 30911/57597
## 2023-09-16 21:46:31 Annotating text fragment 30921/57597
## 2023-09-16 21:46:31 Annotating text fragment 30931/57597
## 2023-09-16 21:46:32 Annotating text fragment 30941/57597
## 2023-09-16 21:46:32 Annotating text fragment 30951/57597
## 2023-09-16 21:46:32 Annotating text fragment 30961/57597
## 2023-09-16 21:46:33 Annotating text fragment 30971/57597
## 2023-09-16 21:46:33 Annotating text fragment 30981/57597
## 2023-09-16 21:46:33 Annotating text fragment 30991/57597
## 2023-09-16 21:46:33 Annotating text fragment 31001/57597
## 2023-09-16 21:46:33 Annotating text fragment 31011/57597
```

```
## 2023-09-16 21:46:33 Annotating text fragment 31021/57597
## 2023-09-16 21:46:34 Annotating text fragment 31031/57597
## 2023-09-16 21:46:34 Annotating text fragment 31041/57597
## 2023-09-16 21:46:34 Annotating text fragment 31051/57597
## 2023-09-16 21:46:34 Annotating text fragment 31061/57597
## 2023-09-16 21:46:34 Annotating text fragment 31071/57597
## 2023-09-16 21:46:34 Annotating text fragment 31081/57597
## 2023-09-16 21:46:35 Annotating text fragment 31091/57597
## 2023-09-16 21:46:35 Annotating text fragment 31101/57597
## 2023-09-16 21:46:35 Annotating text fragment 31111/57597
## 2023-09-16 21:46:35 Annotating text fragment 31121/57597
## 2023-09-16 21:46:35 Annotating text fragment 31131/57597
## 2023-09-16 21:46:35 Annotating text fragment 31141/57597
## 2023-09-16 21:46:35 Annotating text fragment 31151/57597
## 2023-09-16 21:46:35 Annotating text fragment 31161/57597
## 2023-09-16 21:46:35 Annotating text fragment 31171/57597
## 2023-09-16 21:46:35 Annotating text fragment 31181/57597
## 2023-09-16 21:46:35 Annotating text fragment 31191/57597
## 2023-09-16 21:46:35 Annotating text fragment 31201/57597
## 2023-09-16 21:46:35 Annotating text fragment 31211/57597
## 2023-09-16 21:46:35 Annotating text fragment 31221/57597
## 2023-09-16 21:46:35 Annotating text fragment 31231/57597
## 2023-09-16 21:46:35 Annotating text fragment 31241/57597
## 2023-09-16 21:46:35 Annotating text fragment 31251/57597
## 2023-09-16 21:46:36 Annotating text fragment 31261/57597
## 2023-09-16 21:46:36 Annotating text fragment 31271/57597
## 2023-09-16 21:46:36 Annotating text fragment 31281/57597
## 2023-09-16 21:46:36 Annotating text fragment 31291/57597
## 2023-09-16 21:46:36 Annotating text fragment 31301/57597
## 2023-09-16 21:46:36 Annotating text fragment 31311/57597
## 2023-09-16 21:46:36 Annotating text fragment 31321/57597
## 2023-09-16 21:46:36 Annotating text fragment 31331/57597
## 2023-09-16 21:46:36 Annotating text fragment 31341/57597
## 2023-09-16 21:46:36 Annotating text fragment 31351/57597
## 2023-09-16 21:46:36 Annotating text fragment 31361/57597
## 2023-09-16 21:46:37 Annotating text fragment 31371/57597
## 2023-09-16 21:46:37 Annotating text fragment 31381/57597
## 2023-09-16 21:46:37 Annotating text fragment 31391/57597
## 2023-09-16 21:46:37 Annotating text fragment 31401/57597
## 2023-09-16 21:46:37 Annotating text fragment 31411/57597
## 2023-09-16 21:46:37 Annotating text fragment 31421/57597
## 2023-09-16 21:46:37 Annotating text fragment 31431/57597
## 2023-09-16 21:46:37 Annotating text fragment 31441/57597
## 2023-09-16 21:46:38 Annotating text fragment 31451/57597
## 2023-09-16 21:46:38 Annotating text fragment 31461/57597
## 2023-09-16 21:46:38 Annotating text fragment 31471/57597
## 2023-09-16 21:46:38 Annotating text fragment 31481/57597
## 2023-09-16 21:46:38 Annotating text fragment 31491/57597
## 2023-09-16 21:46:38 Annotating text fragment 31501/57597
## 2023-09-16 21:46:38 Annotating text fragment 31511/57597
## 2023-09-16 21:46:38 Annotating text fragment 31521/57597
## 2023-09-16 21:46:38 Annotating text fragment 31531/57597
## 2023-09-16 21:46:38 Annotating text fragment 31541/57597
## 2023-09-16 21:46:38 Annotating text fragment 31551/57597
```

```
## 2023-09-16 21:46:38 Annotating text fragment 31561/57597
## 2023-09-16 21:46:38 Annotating text fragment 31571/57597
## 2023-09-16 21:46:39 Annotating text fragment 31581/57597
## 2023-09-16 21:46:39 Annotating text fragment 31591/57597
## 2023-09-16 21:46:39 Annotating text fragment 31601/57597
## 2023-09-16 21:46:39 Annotating text fragment 31611/57597
## 2023-09-16 21:46:39 Annotating text fragment 31621/57597
## 2023-09-16 21:46:39 Annotating text fragment 31631/57597
## 2023-09-16 21:46:39 Annotating text fragment 31641/57597
## 2023-09-16 21:46:39 Annotating text fragment 31651/57597
## 2023-09-16 21:46:39 Annotating text fragment 31661/57597
## 2023-09-16 21:46:39 Annotating text fragment 31671/57597
## 2023-09-16 21:46:40 Annotating text fragment 31681/57597
## 2023-09-16 21:46:40 Annotating text fragment 31691/57597
## 2023-09-16 21:46:40 Annotating text fragment 31701/57597
## 2023-09-16 21:46:40 Annotating text fragment 31711/57597
## 2023-09-16 21:46:40 Annotating text fragment 31721/57597
## 2023-09-16 21:46:41 Annotating text fragment 31731/57597
## 2023-09-16 21:46:41 Annotating text fragment 31741/57597
## 2023-09-16 21:46:41 Annotating text fragment 31751/57597
## 2023-09-16 21:46:41 Annotating text fragment 31761/57597
## 2023-09-16 21:46:41 Annotating text fragment 31771/57597
## 2023-09-16 21:46:41 Annotating text fragment 31781/57597
## 2023-09-16 21:46:41 Annotating text fragment 31791/57597
## 2023-09-16 21:46:41 Annotating text fragment 31801/57597
## 2023-09-16 21:46:41 Annotating text fragment 31811/57597
## 2023-09-16 21:46:42 Annotating text fragment 31821/57597
## 2023-09-16 21:46:42 Annotating text fragment 31831/57597
## 2023-09-16 21:46:42 Annotating text fragment 31841/57597
## 2023-09-16 21:46:42 Annotating text fragment 31851/57597
## 2023-09-16 21:46:42 Annotating text fragment 31861/57597
## 2023-09-16 21:46:42 Annotating text fragment 31871/57597
## 2023-09-16 21:46:42 Annotating text fragment 31881/57597
## 2023-09-16 21:46:42 Annotating text fragment 31891/57597
## 2023-09-16 21:46:42 Annotating text fragment 31901/57597
## 2023-09-16 21:46:42 Annotating text fragment 31911/57597
## 2023-09-16 21:46:42 Annotating text fragment 31921/57597
## 2023-09-16 21:46:42 Annotating text fragment 31931/57597
## 2023-09-16 21:46:42 Annotating text fragment 31941/57597
## 2023-09-16 21:46:43 Annotating text fragment 31951/57597
## 2023-09-16 21:46:43 Annotating text fragment 31961/57597
## 2023-09-16 21:46:43 Annotating text fragment 31971/57597
## 2023-09-16 21:46:43 Annotating text fragment 31981/57597
## 2023-09-16 21:46:43 Annotating text fragment 31991/57597
## 2023-09-16 21:46:43 Annotating text fragment 32001/57597
## 2023-09-16 21:46:43 Annotating text fragment 32011/57597
## 2023-09-16 21:46:43 Annotating text fragment 32021/57597
## 2023-09-16 21:46:43 Annotating text fragment 32031/57597
## 2023-09-16 21:46:44 Annotating text fragment 32041/57597
## 2023-09-16 21:46:44 Annotating text fragment 32051/57597
## 2023-09-16 21:46:44 Annotating text fragment 32061/57597
## 2023-09-16 21:46:44 Annotating text fragment 32071/57597
## 2023-09-16 21:46:44 Annotating text fragment 32081/57597
## 2023-09-16 21:46:44 Annotating text fragment 32091/57597
```

```
## 2023-09-16 21:46:44 Annotating text fragment 32101/57597
## 2023-09-16 21:46:44 Annotating text fragment 32111/57597
## 2023-09-16 21:46:44 Annotating text fragment 32121/57597
## 2023-09-16 21:46:44 Annotating text fragment 32131/57597
## 2023-09-16 21:46:44 Annotating text fragment 32141/57597
## 2023-09-16 21:46:44 Annotating text fragment 32151/57597
## 2023-09-16 21:46:44 Annotating text fragment 32161/57597
## 2023-09-16 21:46:44 Annotating text fragment 32171/57597
## 2023-09-16 21:46:44 Annotating text fragment 32181/57597
## 2023-09-16 21:46:45 Annotating text fragment 32191/57597
## 2023-09-16 21:46:45 Annotating text fragment 32201/57597
## 2023-09-16 21:46:45 Annotating text fragment 32211/57597
## 2023-09-16 21:46:45 Annotating text fragment 32221/57597
## 2023-09-16 21:46:45 Annotating text fragment 32231/57597
## 2023-09-16 21:46:45 Annotating text fragment 32241/57597
## 2023-09-16 21:46:45 Annotating text fragment 32251/57597
## 2023-09-16 21:46:45 Annotating text fragment 32261/57597
## 2023-09-16 21:46:45 Annotating text fragment 32271/57597
## 2023-09-16 21:46:45 Annotating text fragment 32281/57597
## 2023-09-16 21:46:45 Annotating text fragment 32291/57597
## 2023-09-16 21:46:46 Annotating text fragment 32301/57597
## 2023-09-16 21:46:46 Annotating text fragment 32311/57597
## 2023-09-16 21:46:46 Annotating text fragment 32321/57597
## 2023-09-16 21:46:46 Annotating text fragment 32331/57597
## 2023-09-16 21:46:46 Annotating text fragment 32341/57597
## 2023-09-16 21:46:46 Annotating text fragment 32351/57597
## 2023-09-16 21:46:46 Annotating text fragment 32361/57597
## 2023-09-16 21:46:46 Annotating text fragment 32371/57597
## 2023-09-16 21:46:46 Annotating text fragment 32381/57597
## 2023-09-16 21:46:46 Annotating text fragment 32391/57597
## 2023-09-16 21:46:46 Annotating text fragment 32401/57597
## 2023-09-16 21:46:46 Annotating text fragment 32411/57597
## 2023-09-16 21:46:46 Annotating text fragment 32421/57597
## 2023-09-16 21:46:46 Annotating text fragment 32431/57597
## 2023-09-16 21:46:46 Annotating text fragment 32441/57597
## 2023-09-16 21:46:46 Annotating text fragment 32451/57597
## 2023-09-16 21:46:47 Annotating text fragment 32461/57597
## 2023-09-16 21:46:47 Annotating text fragment 32471/57597
## 2023-09-16 21:46:47 Annotating text fragment 32481/57597
## 2023-09-16 21:46:47 Annotating text fragment 32491/57597
## 2023-09-16 21:46:47 Annotating text fragment 32501/57597
## 2023-09-16 21:46:47 Annotating text fragment 32511/57597
## 2023-09-16 21:46:47 Annotating text fragment 32521/57597
## 2023-09-16 21:46:47 Annotating text fragment 32531/57597
## 2023-09-16 21:46:47 Annotating text fragment 32541/57597
## 2023-09-16 21:46:48 Annotating text fragment 32551/57597
## 2023-09-16 21:46:48 Annotating text fragment 32561/57597
## 2023-09-16 21:46:48 Annotating text fragment 32571/57597
## 2023-09-16 21:46:48 Annotating text fragment 32581/57597
## 2023-09-16 21:46:48 Annotating text fragment 32591/57597
## 2023-09-16 21:46:48 Annotating text fragment 32601/57597
## 2023-09-16 21:46:48 Annotating text fragment 32611/57597
## 2023-09-16 21:46:48 Annotating text fragment 32621/57597
## 2023-09-16 21:46:48 Annotating text fragment 32631/57597
```

```
## 2023-09-16 21:46:48 Annotating text fragment 32641/57597
## 2023-09-16 21:46:49 Annotating text fragment 32651/57597
## 2023-09-16 21:46:49 Annotating text fragment 32661/57597
## 2023-09-16 21:46:49 Annotating text fragment 32671/57597
## 2023-09-16 21:46:49 Annotating text fragment 32681/57597
## 2023-09-16 21:46:49 Annotating text fragment 32691/57597
## 2023-09-16 21:46:49 Annotating text fragment 32701/57597
## 2023-09-16 21:46:49 Annotating text fragment 32711/57597
## 2023-09-16 21:46:49 Annotating text fragment 32721/57597
## 2023-09-16 21:46:49 Annotating text fragment 32731/57597
## 2023-09-16 21:46:49 Annotating text fragment 32741/57597
## 2023-09-16 21:46:49 Annotating text fragment 32751/57597
## 2023-09-16 21:46:49 Annotating text fragment 32761/57597
## 2023-09-16 21:46:49 Annotating text fragment 32771/57597
## 2023-09-16 21:46:49 Annotating text fragment 32781/57597
## 2023-09-16 21:46:50 Annotating text fragment 32791/57597
## 2023-09-16 21:46:50 Annotating text fragment 32801/57597
## 2023-09-16 21:46:50 Annotating text fragment 32811/57597
## 2023-09-16 21:46:50 Annotating text fragment 32821/57597
## 2023-09-16 21:46:50 Annotating text fragment 32831/57597
## 2023-09-16 21:46:50 Annotating text fragment 32841/57597
## 2023-09-16 21:46:50 Annotating text fragment 32851/57597
## 2023-09-16 21:46:50 Annotating text fragment 32861/57597
## 2023-09-16 21:46:50 Annotating text fragment 32871/57597
## 2023-09-16 21:46:51 Annotating text fragment 32881/57597
## 2023-09-16 21:46:51 Annotating text fragment 32891/57597
## 2023-09-16 21:46:51 Annotating text fragment 32901/57597
## 2023-09-16 21:46:51 Annotating text fragment 32911/57597
## 2023-09-16 21:46:51 Annotating text fragment 32921/57597
## 2023-09-16 21:46:51 Annotating text fragment 32931/57597
## 2023-09-16 21:46:51 Annotating text fragment 32941/57597
## 2023-09-16 21:46:51 Annotating text fragment 32951/57597
## 2023-09-16 21:46:51 Annotating text fragment 32961/57597
## 2023-09-16 21:46:51 Annotating text fragment 32971/57597
## 2023-09-16 21:46:51 Annotating text fragment 32981/57597
## 2023-09-16 21:46:51 Annotating text fragment 32991/57597
## 2023-09-16 21:46:51 Annotating text fragment 33001/57597
## 2023-09-16 21:46:52 Annotating text fragment 33011/57597
## 2023-09-16 21:46:52 Annotating text fragment 33021/57597
## 2023-09-16 21:46:52 Annotating text fragment 33031/57597
## 2023-09-16 21:46:52 Annotating text fragment 33041/57597
## 2023-09-16 21:46:52 Annotating text fragment 33051/57597
## 2023-09-16 21:46:52 Annotating text fragment 33061/57597
## 2023-09-16 21:46:52 Annotating text fragment 33071/57597
## 2023-09-16 21:46:52 Annotating text fragment 33081/57597
## 2023-09-16 21:46:53 Annotating text fragment 33091/57597
## 2023-09-16 21:46:53 Annotating text fragment 33101/57597
## 2023-09-16 21:46:53 Annotating text fragment 33111/57597
## 2023-09-16 21:46:53 Annotating text fragment 33121/57597
## 2023-09-16 21:46:53 Annotating text fragment 33131/57597
## 2023-09-16 21:46:53 Annotating text fragment 33141/57597
## 2023-09-16 21:46:53 Annotating text fragment 33151/57597
## 2023-09-16 21:46:53 Annotating text fragment 33161/57597
## 2023-09-16 21:46:53 Annotating text fragment 33171/57597
```

```
## 2023-09-16 21:46:54 Annotating text fragment 33181/57597
## 2023-09-16 21:46:54 Annotating text fragment 33191/57597
## 2023-09-16 21:46:54 Annotating text fragment 33201/57597
## 2023-09-16 21:46:54 Annotating text fragment 33211/57597
## 2023-09-16 21:46:54 Annotating text fragment 33221/57597
## 2023-09-16 21:46:54 Annotating text fragment 33231/57597
## 2023-09-16 21:46:54 Annotating text fragment 33241/57597
## 2023-09-16 21:46:54 Annotating text fragment 33251/57597
## 2023-09-16 21:46:54 Annotating text fragment 33261/57597
## 2023-09-16 21:46:54 Annotating text fragment 33271/57597
## 2023-09-16 21:46:54 Annotating text fragment 33281/57597
## 2023-09-16 21:46:55 Annotating text fragment 33291/57597
## 2023-09-16 21:46:55 Annotating text fragment 33301/57597
## 2023-09-16 21:46:55 Annotating text fragment 33311/57597
## 2023-09-16 21:46:55 Annotating text fragment 33321/57597
## 2023-09-16 21:46:55 Annotating text fragment 33331/57597
## 2023-09-16 21:46:55 Annotating text fragment 33341/57597
## 2023-09-16 21:46:55 Annotating text fragment 33351/57597
## 2023-09-16 21:46:55 Annotating text fragment 33361/57597
## 2023-09-16 21:46:55 Annotating text fragment 33371/57597
## 2023-09-16 21:46:55 Annotating text fragment 33381/57597
## 2023-09-16 21:46:55 Annotating text fragment 33391/57597
## 2023-09-16 21:46:55 Annotating text fragment 33401/57597
## 2023-09-16 21:46:56 Annotating text fragment 33411/57597
## 2023-09-16 21:46:56 Annotating text fragment 33421/57597
## 2023-09-16 21:46:56 Annotating text fragment 33431/57597
## 2023-09-16 21:46:56 Annotating text fragment 33441/57597
## 2023-09-16 21:46:56 Annotating text fragment 33451/57597
## 2023-09-16 21:46:56 Annotating text fragment 33461/57597
## 2023-09-16 21:46:56 Annotating text fragment 33471/57597
## 2023-09-16 21:46:56 Annotating text fragment 33481/57597
## 2023-09-16 21:46:56 Annotating text fragment 33491/57597
## 2023-09-16 21:46:56 Annotating text fragment 33501/57597
## 2023-09-16 21:46:56 Annotating text fragment 33511/57597
## 2023-09-16 21:46:56 Annotating text fragment 33521/57597
## 2023-09-16 21:46:57 Annotating text fragment 33531/57597
## 2023-09-16 21:46:57 Annotating text fragment 33541/57597
## 2023-09-16 21:46:57 Annotating text fragment 33551/57597
## 2023-09-16 21:46:57 Annotating text fragment 33561/57597
## 2023-09-16 21:46:57 Annotating text fragment 33571/57597
## 2023-09-16 21:46:57 Annotating text fragment 33581/57597
## 2023-09-16 21:46:57 Annotating text fragment 33591/57597
## 2023-09-16 21:46:57 Annotating text fragment 33601/57597
## 2023-09-16 21:46:57 Annotating text fragment 33611/57597
## 2023-09-16 21:46:57 Annotating text fragment 33621/57597
## 2023-09-16 21:46:57 Annotating text fragment 33631/57597
## 2023-09-16 21:46:57 Annotating text fragment 33641/57597
## 2023-09-16 21:46:58 Annotating text fragment 33651/57597
## 2023-09-16 21:46:58 Annotating text fragment 33661/57597
## 2023-09-16 21:46:58 Annotating text fragment 33671/57597
## 2023-09-16 21:46:58 Annotating text fragment 33681/57597
## 2023-09-16 21:46:58 Annotating text fragment 33691/57597
## 2023-09-16 21:46:58 Annotating text fragment 33701/57597
## 2023-09-16 21:46:58 Annotating text fragment 33711/57597
```

```
## 2023-09-16 21:46:58 Annotating text fragment 33721/57597
## 2023-09-16 21:46:58 Annotating text fragment 33731/57597
## 2023-09-16 21:46:58 Annotating text fragment 33741/57597
## 2023-09-16 21:46:58 Annotating text fragment 33751/57597
## 2023-09-16 21:46:59 Annotating text fragment 33761/57597
## 2023-09-16 21:46:59 Annotating text fragment 33771/57597
## 2023-09-16 21:46:59 Annotating text fragment 33781/57597
## 2023-09-16 21:46:59 Annotating text fragment 33791/57597
## 2023-09-16 21:46:59 Annotating text fragment 33801/57597
## 2023-09-16 21:46:59 Annotating text fragment 33811/57597
## 2023-09-16 21:46:59 Annotating text fragment 33821/57597
## 2023-09-16 21:46:59 Annotating text fragment 33831/57597
## 2023-09-16 21:47:00 Annotating text fragment 33841/57597
## 2023-09-16 21:47:00 Annotating text fragment 33851/57597
## 2023-09-16 21:47:00 Annotating text fragment 33861/57597
## 2023-09-16 21:47:00 Annotating text fragment 33871/57597
## 2023-09-16 21:47:00 Annotating text fragment 33881/57597
## 2023-09-16 21:47:00 Annotating text fragment 33891/57597
## 2023-09-16 21:47:00 Annotating text fragment 33901/57597
## 2023-09-16 21:47:00 Annotating text fragment 33911/57597
## 2023-09-16 21:47:00 Annotating text fragment 33921/57597
## 2023-09-16 21:47:00 Annotating text fragment 33931/57597
## 2023-09-16 21:47:00 Annotating text fragment 33941/57597
## 2023-09-16 21:47:00 Annotating text fragment 33951/57597
## 2023-09-16 21:47:01 Annotating text fragment 33961/57597
## 2023-09-16 21:47:01 Annotating text fragment 33971/57597
## 2023-09-16 21:47:01 Annotating text fragment 33981/57597
## 2023-09-16 21:47:01 Annotating text fragment 33991/57597
## 2023-09-16 21:47:01 Annotating text fragment 34001/57597
## 2023-09-16 21:47:01 Annotating text fragment 34011/57597
## 2023-09-16 21:47:01 Annotating text fragment 34021/57597
## 2023-09-16 21:47:01 Annotating text fragment 34031/57597
## 2023-09-16 21:47:01 Annotating text fragment 34041/57597
## 2023-09-16 21:47:01 Annotating text fragment 34051/57597
## 2023-09-16 21:47:02 Annotating text fragment 34061/57597
## 2023-09-16 21:47:02 Annotating text fragment 34071/57597
## 2023-09-16 21:47:02 Annotating text fragment 34081/57597
## 2023-09-16 21:47:02 Annotating text fragment 34091/57597
## 2023-09-16 21:47:02 Annotating text fragment 34101/57597
## 2023-09-16 21:47:02 Annotating text fragment 34111/57597
## 2023-09-16 21:47:02 Annotating text fragment 34121/57597
## 2023-09-16 21:47:02 Annotating text fragment 34131/57597
## 2023-09-16 21:47:02 Annotating text fragment 34141/57597
## 2023-09-16 21:47:02 Annotating text fragment 34151/57597
## 2023-09-16 21:47:03 Annotating text fragment 34161/57597
## 2023-09-16 21:47:03 Annotating text fragment 34171/57597
## 2023-09-16 21:47:03 Annotating text fragment 34181/57597
## 2023-09-16 21:47:03 Annotating text fragment 34191/57597
## 2023-09-16 21:47:03 Annotating text fragment 34201/57597
## 2023-09-16 21:47:03 Annotating text fragment 34211/57597
## 2023-09-16 21:47:03 Annotating text fragment 34221/57597
## 2023-09-16 21:47:03 Annotating text fragment 34231/57597
## 2023-09-16 21:47:03 Annotating text fragment 34241/57597
## 2023-09-16 21:47:04 Annotating text fragment 34251/57597
```

```
## 2023-09-16 21:47:04 Annotating text fragment 34261/57597
## 2023-09-16 21:47:04 Annotating text fragment 34271/57597
## 2023-09-16 21:47:04 Annotating text fragment 34281/57597
## 2023-09-16 21:47:04 Annotating text fragment 34291/57597
## 2023-09-16 21:47:04 Annotating text fragment 34301/57597
## 2023-09-16 21:47:04 Annotating text fragment 34311/57597
## 2023-09-16 21:47:04 Annotating text fragment 34321/57597
## 2023-09-16 21:47:05 Annotating text fragment 34331/57597
## 2023-09-16 21:47:05 Annotating text fragment 34341/57597
## 2023-09-16 21:47:05 Annotating text fragment 34351/57597
## 2023-09-16 21:47:05 Annotating text fragment 34361/57597
## 2023-09-16 21:47:05 Annotating text fragment 34371/57597
## 2023-09-16 21:47:05 Annotating text fragment 34381/57597
## 2023-09-16 21:47:05 Annotating text fragment 34391/57597
## 2023-09-16 21:47:05 Annotating text fragment 34401/57597
## 2023-09-16 21:47:05 Annotating text fragment 34411/57597
## 2023-09-16 21:47:05 Annotating text fragment 34421/57597
## 2023-09-16 21:47:05 Annotating text fragment 34431/57597
## 2023-09-16 21:47:06 Annotating text fragment 34441/57597
## 2023-09-16 21:47:06 Annotating text fragment 34451/57597
## 2023-09-16 21:47:06 Annotating text fragment 34461/57597
## 2023-09-16 21:47:06 Annotating text fragment 34471/57597
## 2023-09-16 21:47:06 Annotating text fragment 34481/57597
## 2023-09-16 21:47:06 Annotating text fragment 34491/57597
## 2023-09-16 21:47:06 Annotating text fragment 34501/57597
## 2023-09-16 21:47:06 Annotating text fragment 34511/57597
## 2023-09-16 21:47:07 Annotating text fragment 34521/57597
## 2023-09-16 21:47:07 Annotating text fragment 34531/57597
## 2023-09-16 21:47:07 Annotating text fragment 34541/57597
## 2023-09-16 21:47:07 Annotating text fragment 34551/57597
## 2023-09-16 21:47:07 Annotating text fragment 34561/57597
## 2023-09-16 21:47:07 Annotating text fragment 34571/57597
## 2023-09-16 21:47:08 Annotating text fragment 34581/57597
## 2023-09-16 21:47:08 Annotating text fragment 34591/57597
## 2023-09-16 21:47:08 Annotating text fragment 34601/57597
## 2023-09-16 21:47:08 Annotating text fragment 34611/57597
## 2023-09-16 21:47:08 Annotating text fragment 34621/57597
## 2023-09-16 21:47:08 Annotating text fragment 34631/57597
## 2023-09-16 21:47:08 Annotating text fragment 34641/57597
## 2023-09-16 21:47:08 Annotating text fragment 34651/57597
## 2023-09-16 21:47:08 Annotating text fragment 34661/57597
## 2023-09-16 21:47:09 Annotating text fragment 34671/57597
## 2023-09-16 21:47:09 Annotating text fragment 34681/57597
## 2023-09-16 21:47:09 Annotating text fragment 34691/57597
## 2023-09-16 21:47:09 Annotating text fragment 34701/57597
## 2023-09-16 21:47:09 Annotating text fragment 34711/57597
## 2023-09-16 21:47:09 Annotating text fragment 34721/57597
## 2023-09-16 21:47:09 Annotating text fragment 34731/57597
## 2023-09-16 21:47:09 Annotating text fragment 34741/57597
## 2023-09-16 21:47:09 Annotating text fragment 34751/57597
## 2023-09-16 21:47:09 Annotating text fragment 34761/57597
## 2023-09-16 21:47:09 Annotating text fragment 34771/57597
## 2023-09-16 21:47:09 Annotating text fragment 34781/57597
## 2023-09-16 21:47:10 Annotating text fragment 34791/57597
```

```
## 2023-09-16 21:47:10 Annotating text fragment 34801/57597
## 2023-09-16 21:47:10 Annotating text fragment 34811/57597
## 2023-09-16 21:47:10 Annotating text fragment 34821/57597
## 2023-09-16 21:47:10 Annotating text fragment 34831/57597
## 2023-09-16 21:47:10 Annotating text fragment 34841/57597
## 2023-09-16 21:47:10 Annotating text fragment 34851/57597
## 2023-09-16 21:47:10 Annotating text fragment 34861/57597
## 2023-09-16 21:47:10 Annotating text fragment 34871/57597
## 2023-09-16 21:47:10 Annotating text fragment 34881/57597
## 2023-09-16 21:47:10 Annotating text fragment 34891/57597
## 2023-09-16 21:47:10 Annotating text fragment 34901/57597
## 2023-09-16 21:47:10 Annotating text fragment 34911/57597
## 2023-09-16 21:47:10 Annotating text fragment 34921/57597
## 2023-09-16 21:47:11 Annotating text fragment 34931/57597
## 2023-09-16 21:47:11 Annotating text fragment 34941/57597
## 2023-09-16 21:47:11 Annotating text fragment 34951/57597
## 2023-09-16 21:47:11 Annotating text fragment 34961/57597
## 2023-09-16 21:47:11 Annotating text fragment 34971/57597
## 2023-09-16 21:47:11 Annotating text fragment 34981/57597
## 2023-09-16 21:47:11 Annotating text fragment 34991/57597
## 2023-09-16 21:47:11 Annotating text fragment 35001/57597
## 2023-09-16 21:47:11 Annotating text fragment 35011/57597
## 2023-09-16 21:47:11 Annotating text fragment 35021/57597
## 2023-09-16 21:47:11 Annotating text fragment 35031/57597
## 2023-09-16 21:47:12 Annotating text fragment 35041/57597
## 2023-09-16 21:47:12 Annotating text fragment 35051/57597
## 2023-09-16 21:47:12 Annotating text fragment 35061/57597
## 2023-09-16 21:47:12 Annotating text fragment 35071/57597
## 2023-09-16 21:47:12 Annotating text fragment 35081/57597
## 2023-09-16 21:47:13 Annotating text fragment 35091/57597
## 2023-09-16 21:47:13 Annotating text fragment 35101/57597
## 2023-09-16 21:47:13 Annotating text fragment 35111/57597
## 2023-09-16 21:47:13 Annotating text fragment 35121/57597
## 2023-09-16 21:47:13 Annotating text fragment 35131/57597
## 2023-09-16 21:47:13 Annotating text fragment 35141/57597
## 2023-09-16 21:47:13 Annotating text fragment 35151/57597
## 2023-09-16 21:47:13 Annotating text fragment 35161/57597
## 2023-09-16 21:47:14 Annotating text fragment 35171/57597
## 2023-09-16 21:47:14 Annotating text fragment 35181/57597
## 2023-09-16 21:47:14 Annotating text fragment 35191/57597
## 2023-09-16 21:47:14 Annotating text fragment 35201/57597
## 2023-09-16 21:47:14 Annotating text fragment 35211/57597
## 2023-09-16 21:47:14 Annotating text fragment 35221/57597
## 2023-09-16 21:47:14 Annotating text fragment 35231/57597
## 2023-09-16 21:47:14 Annotating text fragment 35241/57597
## 2023-09-16 21:47:14 Annotating text fragment 35251/57597
## 2023-09-16 21:47:14 Annotating text fragment 35261/57597
## 2023-09-16 21:47:14 Annotating text fragment 35271/57597
## 2023-09-16 21:47:14 Annotating text fragment 35281/57597
## 2023-09-16 21:47:15 Annotating text fragment 35291/57597
## 2023-09-16 21:47:15 Annotating text fragment 35301/57597
## 2023-09-16 21:47:15 Annotating text fragment 35311/57597
## 2023-09-16 21:47:15 Annotating text fragment 35321/57597
## 2023-09-16 21:47:15 Annotating text fragment 35331/57597
```

```
## 2023-09-16 21:47:15 Annotating text fragment 35341/57597
## 2023-09-16 21:47:15 Annotating text fragment 35351/57597
## 2023-09-16 21:47:15 Annotating text fragment 35361/57597
## 2023-09-16 21:47:15 Annotating text fragment 35371/57597
## 2023-09-16 21:47:15 Annotating text fragment 35381/57597
## 2023-09-16 21:47:15 Annotating text fragment 35391/57597
## 2023-09-16 21:47:15 Annotating text fragment 35401/57597
## 2023-09-16 21:47:16 Annotating text fragment 35411/57597
## 2023-09-16 21:47:16 Annotating text fragment 35421/57597
## 2023-09-16 21:47:16 Annotating text fragment 35431/57597
## 2023-09-16 21:47:16 Annotating text fragment 35441/57597
## 2023-09-16 21:47:16 Annotating text fragment 35451/57597
## 2023-09-16 21:47:16 Annotating text fragment 35461/57597
## 2023-09-16 21:47:16 Annotating text fragment 35471/57597
## 2023-09-16 21:47:16 Annotating text fragment 35481/57597
## 2023-09-16 21:47:16 Annotating text fragment 35491/57597
## 2023-09-16 21:47:16 Annotating text fragment 35501/57597
## 2023-09-16 21:47:16 Annotating text fragment 35511/57597
## 2023-09-16 21:47:16 Annotating text fragment 35521/57597
## 2023-09-16 21:47:17 Annotating text fragment 35531/57597
## 2023-09-16 21:47:17 Annotating text fragment 35541/57597
## 2023-09-16 21:47:17 Annotating text fragment 35551/57597
## 2023-09-16 21:47:17 Annotating text fragment 35561/57597
## 2023-09-16 21:47:17 Annotating text fragment 35571/57597
## 2023-09-16 21:47:17 Annotating text fragment 35581/57597
## 2023-09-16 21:47:17 Annotating text fragment 35591/57597
## 2023-09-16 21:47:17 Annotating text fragment 35601/57597
## 2023-09-16 21:47:17 Annotating text fragment 35611/57597
## 2023-09-16 21:47:18 Annotating text fragment 35621/57597
## 2023-09-16 21:47:18 Annotating text fragment 35631/57597
## 2023-09-16 21:47:18 Annotating text fragment 35641/57597
## 2023-09-16 21:47:18 Annotating text fragment 35651/57597
## 2023-09-16 21:47:18 Annotating text fragment 35661/57597
## 2023-09-16 21:47:18 Annotating text fragment 35671/57597
## 2023-09-16 21:47:19 Annotating text fragment 35681/57597
## 2023-09-16 21:47:19 Annotating text fragment 35691/57597
## 2023-09-16 21:47:19 Annotating text fragment 35701/57597
## 2023-09-16 21:47:19 Annotating text fragment 35711/57597
## 2023-09-16 21:47:19 Annotating text fragment 35721/57597
## 2023-09-16 21:47:19 Annotating text fragment 35731/57597
## 2023-09-16 21:47:19 Annotating text fragment 35741/57597
## 2023-09-16 21:47:19 Annotating text fragment 35751/57597
## 2023-09-16 21:47:19 Annotating text fragment 35761/57597
## 2023-09-16 21:47:19 Annotating text fragment 35771/57597
## 2023-09-16 21:47:19 Annotating text fragment 35781/57597
## 2023-09-16 21:47:19 Annotating text fragment 35791/57597
## 2023-09-16 21:47:19 Annotating text fragment 35801/57597
## 2023-09-16 21:47:20 Annotating text fragment 35811/57597
## 2023-09-16 21:47:20 Annotating text fragment 35821/57597
## 2023-09-16 21:47:20 Annotating text fragment 35831/57597
## 2023-09-16 21:47:20 Annotating text fragment 35841/57597
## 2023-09-16 21:47:20 Annotating text fragment 35851/57597
## 2023-09-16 21:47:20 Annotating text fragment 35861/57597
## 2023-09-16 21:47:20 Annotating text fragment 35871/57597
```

```
## 2023-09-16 21:47:20 Annotating text fragment 35881/57597
## 2023-09-16 21:47:20 Annotating text fragment 35891/57597
## 2023-09-16 21:47:20 Annotating text fragment 35901/57597
## 2023-09-16 21:47:20 Annotating text fragment 35911/57597
## 2023-09-16 21:47:21 Annotating text fragment 35921/57597
## 2023-09-16 21:47:21 Annotating text fragment 35931/57597
## 2023-09-16 21:47:21 Annotating text fragment 35941/57597
## 2023-09-16 21:47:21 Annotating text fragment 35951/57597
## 2023-09-16 21:47:21 Annotating text fragment 35961/57597
## 2023-09-16 21:47:21 Annotating text fragment 35971/57597
## 2023-09-16 21:47:21 Annotating text fragment 35981/57597
## 2023-09-16 21:47:21 Annotating text fragment 35991/57597
## 2023-09-16 21:47:21 Annotating text fragment 36001/57597
## 2023-09-16 21:47:22 Annotating text fragment 36011/57597
## 2023-09-16 21:47:22 Annotating text fragment 36021/57597
## 2023-09-16 21:47:22 Annotating text fragment 36031/57597
## 2023-09-16 21:47:22 Annotating text fragment 36041/57597
## 2023-09-16 21:47:22 Annotating text fragment 36051/57597
## 2023-09-16 21:47:22 Annotating text fragment 36061/57597
## 2023-09-16 21:47:22 Annotating text fragment 36071/57597
## 2023-09-16 21:47:22 Annotating text fragment 36081/57597
## 2023-09-16 21:47:22 Annotating text fragment 36091/57597
## 2023-09-16 21:47:22 Annotating text fragment 36101/57597
## 2023-09-16 21:47:23 Annotating text fragment 36111/57597
## 2023-09-16 21:47:23 Annotating text fragment 36121/57597
## 2023-09-16 21:47:23 Annotating text fragment 36131/57597
## 2023-09-16 21:47:23 Annotating text fragment 36141/57597
## 2023-09-16 21:47:23 Annotating text fragment 36151/57597
## 2023-09-16 21:47:23 Annotating text fragment 36161/57597
## 2023-09-16 21:47:23 Annotating text fragment 36171/57597
## 2023-09-16 21:47:23 Annotating text fragment 36181/57597
## 2023-09-16 21:47:23 Annotating text fragment 36191/57597
## 2023-09-16 21:47:23 Annotating text fragment 36201/57597
## 2023-09-16 21:47:23 Annotating text fragment 36211/57597
## 2023-09-16 21:47:23 Annotating text fragment 36221/57597
## 2023-09-16 21:47:23 Annotating text fragment 36231/57597
## 2023-09-16 21:47:23 Annotating text fragment 36241/57597
## 2023-09-16 21:47:23 Annotating text fragment 36251/57597
## 2023-09-16 21:47:23 Annotating text fragment 36261/57597
## 2023-09-16 21:47:24 Annotating text fragment 36271/57597
## 2023-09-16 21:47:24 Annotating text fragment 36281/57597
## 2023-09-16 21:47:24 Annotating text fragment 36291/57597
## 2023-09-16 21:47:24 Annotating text fragment 36301/57597
## 2023-09-16 21:47:24 Annotating text fragment 36311/57597
## 2023-09-16 21:47:24 Annotating text fragment 36321/57597
## 2023-09-16 21:47:24 Annotating text fragment 36331/57597
## 2023-09-16 21:47:24 Annotating text fragment 36341/57597
## 2023-09-16 21:47:24 Annotating text fragment 36351/57597
## 2023-09-16 21:47:24 Annotating text fragment 36361/57597
## 2023-09-16 21:47:24 Annotating text fragment 36371/57597
## 2023-09-16 21:47:24 Annotating text fragment 36381/57597
## 2023-09-16 21:47:24 Annotating text fragment 36391/57597
## 2023-09-16 21:47:25 Annotating text fragment 36401/57597
## 2023-09-16 21:47:25 Annotating text fragment 36411/57597
```

```
## 2023-09-16 21:47:25 Annotating text fragment 36421/57597
## 2023-09-16 21:47:25 Annotating text fragment 36431/57597
## 2023-09-16 21:47:25 Annotating text fragment 36441/57597
## 2023-09-16 21:47:25 Annotating text fragment 36451/57597
## 2023-09-16 21:47:25 Annotating text fragment 36461/57597
## 2023-09-16 21:47:25 Annotating text fragment 36471/57597
## 2023-09-16 21:47:25 Annotating text fragment 36481/57597
## 2023-09-16 21:47:25 Annotating text fragment 36491/57597
## 2023-09-16 21:47:26 Annotating text fragment 36501/57597
## 2023-09-16 21:47:26 Annotating text fragment 36511/57597
## 2023-09-16 21:47:26 Annotating text fragment 36521/57597
## 2023-09-16 21:47:26 Annotating text fragment 36531/57597
## 2023-09-16 21:47:26 Annotating text fragment 36541/57597
## 2023-09-16 21:47:26 Annotating text fragment 36551/57597
## 2023-09-16 21:47:26 Annotating text fragment 36561/57597
## 2023-09-16 21:47:26 Annotating text fragment 36571/57597
## 2023-09-16 21:47:26 Annotating text fragment 36581/57597
## 2023-09-16 21:47:26 Annotating text fragment 36591/57597
## 2023-09-16 21:47:26 Annotating text fragment 36601/57597
## 2023-09-16 21:47:26 Annotating text fragment 36611/57597
## 2023-09-16 21:47:26 Annotating text fragment 36621/57597
## 2023-09-16 21:47:26 Annotating text fragment 36631/57597
## 2023-09-16 21:47:27 Annotating text fragment 36641/57597
## 2023-09-16 21:47:27 Annotating text fragment 36651/57597
## 2023-09-16 21:47:27 Annotating text fragment 36661/57597
## 2023-09-16 21:47:27 Annotating text fragment 36671/57597
## 2023-09-16 21:47:27 Annotating text fragment 36681/57597
## 2023-09-16 21:47:27 Annotating text fragment 36691/57597
## 2023-09-16 21:47:27 Annotating text fragment 36701/57597
## 2023-09-16 21:47:27 Annotating text fragment 36711/57597
## 2023-09-16 21:47:27 Annotating text fragment 36721/57597
## 2023-09-16 21:47:27 Annotating text fragment 36731/57597
## 2023-09-16 21:47:27 Annotating text fragment 36741/57597
## 2023-09-16 21:47:28 Annotating text fragment 36751/57597
## 2023-09-16 21:47:28 Annotating text fragment 36761/57597
## 2023-09-16 21:47:28 Annotating text fragment 36771/57597
## 2023-09-16 21:47:28 Annotating text fragment 36781/57597
## 2023-09-16 21:47:28 Annotating text fragment 36791/57597
## 2023-09-16 21:47:28 Annotating text fragment 36801/57597
## 2023-09-16 21:47:28 Annotating text fragment 36811/57597
## 2023-09-16 21:47:28 Annotating text fragment 36821/57597
## 2023-09-16 21:47:28 Annotating text fragment 36831/57597
## 2023-09-16 21:47:28 Annotating text fragment 36841/57597
## 2023-09-16 21:47:28 Annotating text fragment 36851/57597
## 2023-09-16 21:47:28 Annotating text fragment 36861/57597
## 2023-09-16 21:47:28 Annotating text fragment 36871/57597
## 2023-09-16 21:47:28 Annotating text fragment 36881/57597
## 2023-09-16 21:47:29 Annotating text fragment 36891/57597
## 2023-09-16 21:47:29 Annotating text fragment 36901/57597
## 2023-09-16 21:47:29 Annotating text fragment 36911/57597
## 2023-09-16 21:47:29 Annotating text fragment 36921/57597
## 2023-09-16 21:47:29 Annotating text fragment 36931/57597
## 2023-09-16 21:47:29 Annotating text fragment 36941/57597
## 2023-09-16 21:47:29 Annotating text fragment 36951/57597
```

```
## 2023-09-16 21:47:29 Annotating text fragment 36961/57597
## 2023-09-16 21:47:29 Annotating text fragment 36971/57597
## 2023-09-16 21:47:29 Annotating text fragment 36981/57597
## 2023-09-16 21:47:29 Annotating text fragment 36991/57597
## 2023-09-16 21:47:29 Annotating text fragment 37001/57597
## 2023-09-16 21:47:29 Annotating text fragment 37011/57597
## 2023-09-16 21:47:30 Annotating text fragment 37021/57597
## 2023-09-16 21:47:30 Annotating text fragment 37031/57597
## 2023-09-16 21:47:30 Annotating text fragment 37041/57597
## 2023-09-16 21:47:30 Annotating text fragment 37051/57597
## 2023-09-16 21:47:30 Annotating text fragment 37061/57597
## 2023-09-16 21:47:30 Annotating text fragment 37071/57597
## 2023-09-16 21:47:30 Annotating text fragment 37081/57597
## 2023-09-16 21:47:30 Annotating text fragment 37091/57597
## 2023-09-16 21:47:30 Annotating text fragment 37101/57597
## 2023-09-16 21:47:30 Annotating text fragment 37111/57597
## 2023-09-16 21:47:30 Annotating text fragment 37121/57597
## 2023-09-16 21:47:31 Annotating text fragment 37131/57597
## 2023-09-16 21:47:31 Annotating text fragment 37141/57597
## 2023-09-16 21:47:31 Annotating text fragment 37151/57597
## 2023-09-16 21:47:31 Annotating text fragment 37161/57597
## 2023-09-16 21:47:31 Annotating text fragment 37171/57597
## 2023-09-16 21:47:31 Annotating text fragment 37181/57597
## 2023-09-16 21:47:31 Annotating text fragment 37191/57597
## 2023-09-16 21:47:31 Annotating text fragment 37201/57597
## 2023-09-16 21:47:31 Annotating text fragment 37211/57597
## 2023-09-16 21:47:32 Annotating text fragment 37221/57597
## 2023-09-16 21:47:32 Annotating text fragment 37231/57597
## 2023-09-16 21:47:32 Annotating text fragment 37241/57597
## 2023-09-16 21:47:32 Annotating text fragment 37251/57597
## 2023-09-16 21:47:32 Annotating text fragment 37261/57597
## 2023-09-16 21:47:32 Annotating text fragment 37271/57597
## 2023-09-16 21:47:32 Annotating text fragment 37281/57597
## 2023-09-16 21:47:32 Annotating text fragment 37291/57597
## 2023-09-16 21:47:32 Annotating text fragment 37301/57597
## 2023-09-16 21:47:32 Annotating text fragment 37311/57597
## 2023-09-16 21:47:33 Annotating text fragment 37321/57597
## 2023-09-16 21:47:33 Annotating text fragment 37331/57597
## 2023-09-16 21:47:33 Annotating text fragment 37341/57597
## 2023-09-16 21:47:33 Annotating text fragment 37351/57597
## 2023-09-16 21:47:33 Annotating text fragment 37361/57597
## 2023-09-16 21:47:33 Annotating text fragment 37371/57597
## 2023-09-16 21:47:33 Annotating text fragment 37381/57597
## 2023-09-16 21:47:33 Annotating text fragment 37391/57597
## 2023-09-16 21:47:33 Annotating text fragment 37401/57597
## 2023-09-16 21:47:33 Annotating text fragment 37411/57597
## 2023-09-16 21:47:33 Annotating text fragment 37421/57597
## 2023-09-16 21:47:33 Annotating text fragment 37431/57597
## 2023-09-16 21:47:34 Annotating text fragment 37441/57597
## 2023-09-16 21:47:34 Annotating text fragment 37451/57597
## 2023-09-16 21:47:34 Annotating text fragment 37461/57597
## 2023-09-16 21:47:34 Annotating text fragment 37471/57597
## 2023-09-16 21:47:34 Annotating text fragment 37481/57597
## 2023-09-16 21:47:34 Annotating text fragment 37491/57597
```

```
## 2023-09-16 21:47:34 Annotating text fragment 37501/57597
## 2023-09-16 21:47:34 Annotating text fragment 37511/57597
## 2023-09-16 21:47:35 Annotating text fragment 37521/57597
## 2023-09-16 21:47:35 Annotating text fragment 37531/57597
## 2023-09-16 21:47:35 Annotating text fragment 37541/57597
## 2023-09-16 21:47:35 Annotating text fragment 37551/57597
## 2023-09-16 21:47:35 Annotating text fragment 37561/57597
## 2023-09-16 21:47:35 Annotating text fragment 37571/57597
## 2023-09-16 21:47:35 Annotating text fragment 37581/57597
## 2023-09-16 21:47:35 Annotating text fragment 37591/57597
## 2023-09-16 21:47:35 Annotating text fragment 37601/57597
## 2023-09-16 21:47:35 Annotating text fragment 37611/57597
## 2023-09-16 21:47:35 Annotating text fragment 37621/57597
## 2023-09-16 21:47:36 Annotating text fragment 37631/57597
## 2023-09-16 21:47:36 Annotating text fragment 37641/57597
## 2023-09-16 21:47:36 Annotating text fragment 37651/57597
## 2023-09-16 21:47:36 Annotating text fragment 37661/57597
## 2023-09-16 21:47:36 Annotating text fragment 37671/57597
## 2023-09-16 21:47:36 Annotating text fragment 37681/57597
## 2023-09-16 21:47:36 Annotating text fragment 37691/57597
## 2023-09-16 21:47:36 Annotating text fragment 37701/57597
## 2023-09-16 21:47:37 Annotating text fragment 37711/57597
## 2023-09-16 21:47:37 Annotating text fragment 37721/57597
## 2023-09-16 21:47:37 Annotating text fragment 37731/57597
## 2023-09-16 21:47:38 Annotating text fragment 37741/57597
## 2023-09-16 21:47:38 Annotating text fragment 37751/57597
## 2023-09-16 21:47:38 Annotating text fragment 37761/57597
## 2023-09-16 21:47:38 Annotating text fragment 37771/57597
## 2023-09-16 21:47:38 Annotating text fragment 37781/57597
## 2023-09-16 21:47:38 Annotating text fragment 37791/57597
## 2023-09-16 21:47:38 Annotating text fragment 37801/57597
## 2023-09-16 21:47:38 Annotating text fragment 37811/57597
## 2023-09-16 21:47:38 Annotating text fragment 37821/57597
## 2023-09-16 21:47:38 Annotating text fragment 37831/57597
## 2023-09-16 21:47:38 Annotating text fragment 37841/57597
## 2023-09-16 21:47:39 Annotating text fragment 37851/57597
## 2023-09-16 21:47:39 Annotating text fragment 37861/57597
## 2023-09-16 21:47:39 Annotating text fragment 37871/57597
## 2023-09-16 21:47:39 Annotating text fragment 37881/57597
## 2023-09-16 21:47:39 Annotating text fragment 37891/57597
## 2023-09-16 21:47:39 Annotating text fragment 37901/57597
## 2023-09-16 21:47:39 Annotating text fragment 37911/57597
## 2023-09-16 21:47:39 Annotating text fragment 37921/57597
## 2023-09-16 21:47:39 Annotating text fragment 37931/57597
## 2023-09-16 21:47:39 Annotating text fragment 37941/57597
## 2023-09-16 21:47:39 Annotating text fragment 37951/57597
## 2023-09-16 21:47:39 Annotating text fragment 37961/57597
## 2023-09-16 21:47:40 Annotating text fragment 37971/57597
## 2023-09-16 21:47:40 Annotating text fragment 37981/57597
## 2023-09-16 21:47:40 Annotating text fragment 37991/57597
## 2023-09-16 21:47:40 Annotating text fragment 38001/57597
## 2023-09-16 21:47:40 Annotating text fragment 38011/57597
## 2023-09-16 21:47:40 Annotating text fragment 38021/57597
## 2023-09-16 21:47:40 Annotating text fragment 38031/57597
```

```
## 2023-09-16 21:47:40 Annotating text fragment 38041/57597
## 2023-09-16 21:47:40 Annotating text fragment 38051/57597
## 2023-09-16 21:47:40 Annotating text fragment 38061/57597
## 2023-09-16 21:47:40 Annotating text fragment 38071/57597
## 2023-09-16 21:47:40 Annotating text fragment 38081/57597
## 2023-09-16 21:47:40 Annotating text fragment 38091/57597
## 2023-09-16 21:47:41 Annotating text fragment 38101/57597
## 2023-09-16 21:47:41 Annotating text fragment 38111/57597
## 2023-09-16 21:47:41 Annotating text fragment 38121/57597
## 2023-09-16 21:47:41 Annotating text fragment 38131/57597
## 2023-09-16 21:47:41 Annotating text fragment 38141/57597
## 2023-09-16 21:47:41 Annotating text fragment 38151/57597
## 2023-09-16 21:47:41 Annotating text fragment 38161/57597
## 2023-09-16 21:47:41 Annotating text fragment 38171/57597
## 2023-09-16 21:47:41 Annotating text fragment 38181/57597
## 2023-09-16 21:47:41 Annotating text fragment 38191/57597
## 2023-09-16 21:47:41 Annotating text fragment 38201/57597
## 2023-09-16 21:47:41 Annotating text fragment 38211/57597
## 2023-09-16 21:47:41 Annotating text fragment 38221/57597
## 2023-09-16 21:47:42 Annotating text fragment 38231/57597
## 2023-09-16 21:47:42 Annotating text fragment 38241/57597
## 2023-09-16 21:47:42 Annotating text fragment 38251/57597
## 2023-09-16 21:47:42 Annotating text fragment 38261/57597
## 2023-09-16 21:47:42 Annotating text fragment 38271/57597
## 2023-09-16 21:47:42 Annotating text fragment 38281/57597
## 2023-09-16 21:47:42 Annotating text fragment 38291/57597
## 2023-09-16 21:47:42 Annotating text fragment 38301/57597
## 2023-09-16 21:47:42 Annotating text fragment 38311/57597
## 2023-09-16 21:47:42 Annotating text fragment 38321/57597
## 2023-09-16 21:47:42 Annotating text fragment 38331/57597
## 2023-09-16 21:47:42 Annotating text fragment 38341/57597
## 2023-09-16 21:47:43 Annotating text fragment 38351/57597
## 2023-09-16 21:47:43 Annotating text fragment 38361/57597
## 2023-09-16 21:47:43 Annotating text fragment 38371/57597
## 2023-09-16 21:47:43 Annotating text fragment 38381/57597
## 2023-09-16 21:47:43 Annotating text fragment 38391/57597
## 2023-09-16 21:47:43 Annotating text fragment 38401/57597
## 2023-09-16 21:47:43 Annotating text fragment 38411/57597
## 2023-09-16 21:47:43 Annotating text fragment 38421/57597
## 2023-09-16 21:47:44 Annotating text fragment 38431/57597
## 2023-09-16 21:47:44 Annotating text fragment 38441/57597
## 2023-09-16 21:47:44 Annotating text fragment 38451/57597
## 2023-09-16 21:47:44 Annotating text fragment 38461/57597
## 2023-09-16 21:47:44 Annotating text fragment 38471/57597
## 2023-09-16 21:47:44 Annotating text fragment 38481/57597
## 2023-09-16 21:47:44 Annotating text fragment 38491/57597
## 2023-09-16 21:47:44 Annotating text fragment 38501/57597
## 2023-09-16 21:47:44 Annotating text fragment 38511/57597
## 2023-09-16 21:47:44 Annotating text fragment 38521/57597
## 2023-09-16 21:47:44 Annotating text fragment 38531/57597
## 2023-09-16 21:47:44 Annotating text fragment 38541/57597
## 2023-09-16 21:47:44 Annotating text fragment 38551/57597
## 2023-09-16 21:47:45 Annotating text fragment 38561/57597
## 2023-09-16 21:47:45 Annotating text fragment 38571/57597
```

```
## 2023-09-16 21:47:45 Annotating text fragment 38581/57597
## 2023-09-16 21:47:45 Annotating text fragment 38591/57597
## 2023-09-16 21:47:45 Annotating text fragment 38601/57597
## 2023-09-16 21:47:45 Annotating text fragment 38611/57597
## 2023-09-16 21:47:45 Annotating text fragment 38621/57597
## 2023-09-16 21:47:45 Annotating text fragment 38631/57597
## 2023-09-16 21:47:45 Annotating text fragment 38641/57597
## 2023-09-16 21:47:45 Annotating text fragment 38651/57597
## 2023-09-16 21:47:45 Annotating text fragment 38661/57597
## 2023-09-16 21:47:45 Annotating text fragment 38671/57597
## 2023-09-16 21:47:45 Annotating text fragment 38681/57597
## 2023-09-16 21:47:45 Annotating text fragment 38691/57597
## 2023-09-16 21:47:46 Annotating text fragment 38701/57597
## 2023-09-16 21:47:46 Annotating text fragment 38711/57597
## 2023-09-16 21:47:46 Annotating text fragment 38721/57597
## 2023-09-16 21:47:46 Annotating text fragment 38731/57597
## 2023-09-16 21:47:46 Annotating text fragment 38741/57597
## 2023-09-16 21:47:46 Annotating text fragment 38751/57597
## 2023-09-16 21:47:46 Annotating text fragment 38761/57597
## 2023-09-16 21:47:46 Annotating text fragment 38771/57597
## 2023-09-16 21:47:46 Annotating text fragment 38781/57597
## 2023-09-16 21:47:46 Annotating text fragment 38791/57597
## 2023-09-16 21:47:47 Annotating text fragment 38801/57597
## 2023-09-16 21:47:47 Annotating text fragment 38811/57597
## 2023-09-16 21:47:47 Annotating text fragment 38821/57597
## 2023-09-16 21:47:47 Annotating text fragment 38831/57597
## 2023-09-16 21:47:47 Annotating text fragment 38841/57597
## 2023-09-16 21:47:47 Annotating text fragment 38851/57597
## 2023-09-16 21:47:47 Annotating text fragment 38861/57597
## 2023-09-16 21:47:47 Annotating text fragment 38871/57597
## 2023-09-16 21:47:47 Annotating text fragment 38881/57597
## 2023-09-16 21:47:48 Annotating text fragment 38891/57597
## 2023-09-16 21:47:48 Annotating text fragment 38901/57597
## 2023-09-16 21:47:48 Annotating text fragment 38911/57597
## 2023-09-16 21:47:48 Annotating text fragment 38921/57597
## 2023-09-16 21:47:48 Annotating text fragment 38931/57597
## 2023-09-16 21:47:48 Annotating text fragment 38941/57597
## 2023-09-16 21:47:48 Annotating text fragment 38951/57597
## 2023-09-16 21:47:48 Annotating text fragment 38961/57597
## 2023-09-16 21:47:48 Annotating text fragment 38971/57597
## 2023-09-16 21:47:48 Annotating text fragment 38981/57597
## 2023-09-16 21:47:49 Annotating text fragment 38991/57597
## 2023-09-16 21:47:49 Annotating text fragment 39001/57597
## 2023-09-16 21:47:49 Annotating text fragment 39011/57597
## 2023-09-16 21:47:49 Annotating text fragment 39021/57597
## 2023-09-16 21:47:49 Annotating text fragment 39031/57597
## 2023-09-16 21:47:49 Annotating text fragment 39041/57597
## 2023-09-16 21:47:49 Annotating text fragment 39051/57597
## 2023-09-16 21:47:49 Annotating text fragment 39061/57597
## 2023-09-16 21:47:49 Annotating text fragment 39071/57597
## 2023-09-16 21:47:49 Annotating text fragment 39081/57597
## 2023-09-16 21:47:49 Annotating text fragment 39091/57597
## 2023-09-16 21:47:49 Annotating text fragment 39101/57597
## 2023-09-16 21:47:50 Annotating text fragment 39111/57597
```

```
## 2023-09-16 21:47:50 Annotating text fragment 39121/57597
## 2023-09-16 21:47:50 Annotating text fragment 39131/57597
## 2023-09-16 21:47:50 Annotating text fragment 39141/57597
## 2023-09-16 21:47:50 Annotating text fragment 39151/57597
## 2023-09-16 21:47:50 Annotating text fragment 39161/57597
## 2023-09-16 21:47:50 Annotating text fragment 39171/57597
## 2023-09-16 21:47:50 Annotating text fragment 39181/57597
## 2023-09-16 21:47:51 Annotating text fragment 39191/57597
## 2023-09-16 21:47:51 Annotating text fragment 39201/57597
## 2023-09-16 21:47:51 Annotating text fragment 39211/57597
## 2023-09-16 21:47:51 Annotating text fragment 39221/57597
## 2023-09-16 21:47:51 Annotating text fragment 39231/57597
## 2023-09-16 21:47:51 Annotating text fragment 39241/57597
## 2023-09-16 21:47:51 Annotating text fragment 39251/57597
## 2023-09-16 21:47:51 Annotating text fragment 39261/57597
## 2023-09-16 21:47:51 Annotating text fragment 39271/57597
## 2023-09-16 21:47:51 Annotating text fragment 39281/57597
## 2023-09-16 21:47:51 Annotating text fragment 39291/57597
## 2023-09-16 21:47:51 Annotating text fragment 39301/57597
## 2023-09-16 21:47:51 Annotating text fragment 39311/57597
## 2023-09-16 21:47:51 Annotating text fragment 39321/57597
## 2023-09-16 21:47:51 Annotating text fragment 39331/57597
## 2023-09-16 21:47:52 Annotating text fragment 39341/57597
## 2023-09-16 21:47:52 Annotating text fragment 39351/57597
## 2023-09-16 21:47:52 Annotating text fragment 39361/57597
## 2023-09-16 21:47:52 Annotating text fragment 39371/57597
## 2023-09-16 21:47:52 Annotating text fragment 39381/57597
## 2023-09-16 21:47:52 Annotating text fragment 39391/57597
## 2023-09-16 21:47:52 Annotating text fragment 39401/57597
## 2023-09-16 21:47:52 Annotating text fragment 39411/57597
## 2023-09-16 21:47:52 Annotating text fragment 39421/57597
## 2023-09-16 21:47:52 Annotating text fragment 39431/57597
## 2023-09-16 21:47:52 Annotating text fragment 39441/57597
## 2023-09-16 21:47:52 Annotating text fragment 39451/57597
## 2023-09-16 21:47:52 Annotating text fragment 39461/57597
## 2023-09-16 21:47:52 Annotating text fragment 39471/57597
## 2023-09-16 21:47:52 Annotating text fragment 39481/57597
## 2023-09-16 21:47:53 Annotating text fragment 39491/57597
## 2023-09-16 21:47:53 Annotating text fragment 39501/57597
## 2023-09-16 21:47:53 Annotating text fragment 39511/57597
## 2023-09-16 21:47:53 Annotating text fragment 39521/57597
## 2023-09-16 21:47:53 Annotating text fragment 39531/57597
## 2023-09-16 21:47:53 Annotating text fragment 39541/57597
## 2023-09-16 21:47:53 Annotating text fragment 39551/57597
## 2023-09-16 21:47:53 Annotating text fragment 39561/57597
## 2023-09-16 21:47:54 Annotating text fragment 39571/57597
## 2023-09-16 21:47:54 Annotating text fragment 39581/57597
## 2023-09-16 21:47:54 Annotating text fragment 39591/57597
## 2023-09-16 21:47:54 Annotating text fragment 39601/57597
## 2023-09-16 21:47:54 Annotating text fragment 39611/57597
## 2023-09-16 21:47:54 Annotating text fragment 39621/57597
## 2023-09-16 21:47:54 Annotating text fragment 39631/57597
## 2023-09-16 21:47:54 Annotating text fragment 39641/57597
## 2023-09-16 21:47:54 Annotating text fragment 39651/57597
```

```
## 2023-09-16 21:47:54 Annotating text fragment 39661/57597
## 2023-09-16 21:47:54 Annotating text fragment 39671/57597
## 2023-09-16 21:47:55 Annotating text fragment 39681/57597
## 2023-09-16 21:47:55 Annotating text fragment 39691/57597
## 2023-09-16 21:47:55 Annotating text fragment 39701/57597
## 2023-09-16 21:47:55 Annotating text fragment 39711/57597
## 2023-09-16 21:47:55 Annotating text fragment 39721/57597
## 2023-09-16 21:47:55 Annotating text fragment 39731/57597
## 2023-09-16 21:47:55 Annotating text fragment 39741/57597
## 2023-09-16 21:47:55 Annotating text fragment 39751/57597
## 2023-09-16 21:47:55 Annotating text fragment 39761/57597
## 2023-09-16 21:47:55 Annotating text fragment 39771/57597
## 2023-09-16 21:47:55 Annotating text fragment 39781/57597
## 2023-09-16 21:47:56 Annotating text fragment 39791/57597
## 2023-09-16 21:47:56 Annotating text fragment 39801/57597
## 2023-09-16 21:47:56 Annotating text fragment 39811/57597
## 2023-09-16 21:47:56 Annotating text fragment 39821/57597
## 2023-09-16 21:47:56 Annotating text fragment 39831/57597
## 2023-09-16 21:47:56 Annotating text fragment 39841/57597
## 2023-09-16 21:47:56 Annotating text fragment 39851/57597
## 2023-09-16 21:47:56 Annotating text fragment 39861/57597
## 2023-09-16 21:47:56 Annotating text fragment 39871/57597
## 2023-09-16 21:47:56 Annotating text fragment 39881/57597
## 2023-09-16 21:47:56 Annotating text fragment 39891/57597
## 2023-09-16 21:47:56 Annotating text fragment 39901/57597
## 2023-09-16 21:47:56 Annotating text fragment 39911/57597
## 2023-09-16 21:47:57 Annotating text fragment 39921/57597
## 2023-09-16 21:47:57 Annotating text fragment 39931/57597
## 2023-09-16 21:47:57 Annotating text fragment 39941/57597
## 2023-09-16 21:47:57 Annotating text fragment 39951/57597
## 2023-09-16 21:47:57 Annotating text fragment 39961/57597
## 2023-09-16 21:47:57 Annotating text fragment 39971/57597
## 2023-09-16 21:47:57 Annotating text fragment 39981/57597
## 2023-09-16 21:47:57 Annotating text fragment 39991/57597
## 2023-09-16 21:47:57 Annotating text fragment 40001/57597
## 2023-09-16 21:47:58 Annotating text fragment 40011/57597
## 2023-09-16 21:47:58 Annotating text fragment 40021/57597
## 2023-09-16 21:47:58 Annotating text fragment 40031/57597
## 2023-09-16 21:47:58 Annotating text fragment 40041/57597
## 2023-09-16 21:47:58 Annotating text fragment 40051/57597
## 2023-09-16 21:47:58 Annotating text fragment 40061/57597
## 2023-09-16 21:47:58 Annotating text fragment 40071/57597
## 2023-09-16 21:47:58 Annotating text fragment 40081/57597
## 2023-09-16 21:47:58 Annotating text fragment 40091/57597
## 2023-09-16 21:47:58 Annotating text fragment 40101/57597
## 2023-09-16 21:47:58 Annotating text fragment 40111/57597
## 2023-09-16 21:47:58 Annotating text fragment 40121/57597
## 2023-09-16 21:47:59 Annotating text fragment 40131/57597
## 2023-09-16 21:47:59 Annotating text fragment 40141/57597
## 2023-09-16 21:47:59 Annotating text fragment 40151/57597
## 2023-09-16 21:47:59 Annotating text fragment 40161/57597
## 2023-09-16 21:47:59 Annotating text fragment 40171/57597
## 2023-09-16 21:47:59 Annotating text fragment 40181/57597
## 2023-09-16 21:47:59 Annotating text fragment 40191/57597
```

```
## 2023-09-16 21:47:59 Annotating text fragment 40201/57597
## 2023-09-16 21:47:59 Annotating text fragment 40211/57597
## 2023-09-16 21:47:59 Annotating text fragment 40221/57597
## 2023-09-16 21:47:59 Annotating text fragment 40231/57597
## 2023-09-16 21:47:59 Annotating text fragment 40241/57597
## 2023-09-16 21:48:00 Annotating text fragment 40251/57597
## 2023-09-16 21:48:00 Annotating text fragment 40261/57597
## 2023-09-16 21:48:00 Annotating text fragment 40271/57597
## 2023-09-16 21:48:00 Annotating text fragment 40281/57597
## 2023-09-16 21:48:00 Annotating text fragment 40291/57597
## 2023-09-16 21:48:00 Annotating text fragment 40301/57597
## 2023-09-16 21:48:00 Annotating text fragment 40311/57597
## 2023-09-16 21:48:00 Annotating text fragment 40321/57597
## 2023-09-16 21:48:00 Annotating text fragment 40331/57597
## 2023-09-16 21:48:00 Annotating text fragment 40341/57597
## 2023-09-16 21:48:01 Annotating text fragment 40351/57597
## 2023-09-16 21:48:01 Annotating text fragment 40361/57597
## 2023-09-16 21:48:01 Annotating text fragment 40371/57597
## 2023-09-16 21:48:01 Annotating text fragment 40381/57597
## 2023-09-16 21:48:01 Annotating text fragment 40391/57597
## 2023-09-16 21:48:01 Annotating text fragment 40401/57597
## 2023-09-16 21:48:01 Annotating text fragment 40411/57597
## 2023-09-16 21:48:01 Annotating text fragment 40421/57597
## 2023-09-16 21:48:01 Annotating text fragment 40431/57597
## 2023-09-16 21:48:01 Annotating text fragment 40441/57597
## 2023-09-16 21:48:02 Annotating text fragment 40451/57597
## 2023-09-16 21:48:02 Annotating text fragment 40461/57597
## 2023-09-16 21:48:02 Annotating text fragment 40471/57597
## 2023-09-16 21:48:02 Annotating text fragment 40481/57597
## 2023-09-16 21:48:02 Annotating text fragment 40491/57597
## 2023-09-16 21:48:02 Annotating text fragment 40501/57597
## 2023-09-16 21:48:02 Annotating text fragment 40511/57597
## 2023-09-16 21:48:02 Annotating text fragment 40521/57597
## 2023-09-16 21:48:02 Annotating text fragment 40531/57597
## 2023-09-16 21:48:02 Annotating text fragment 40541/57597
## 2023-09-16 21:48:02 Annotating text fragment 40551/57597
## 2023-09-16 21:48:02 Annotating text fragment 40561/57597
## 2023-09-16 21:48:02 Annotating text fragment 40571/57597
## 2023-09-16 21:48:03 Annotating text fragment 40581/57597
## 2023-09-16 21:48:03 Annotating text fragment 40591/57597
## 2023-09-16 21:48:03 Annotating text fragment 40601/57597
## 2023-09-16 21:48:03 Annotating text fragment 40611/57597
## 2023-09-16 21:48:03 Annotating text fragment 40621/57597
## 2023-09-16 21:48:03 Annotating text fragment 40631/57597
## 2023-09-16 21:48:03 Annotating text fragment 40641/57597
## 2023-09-16 21:48:03 Annotating text fragment 40651/57597
## 2023-09-16 21:48:03 Annotating text fragment 40661/57597
## 2023-09-16 21:48:03 Annotating text fragment 40671/57597
## 2023-09-16 21:48:04 Annotating text fragment 40681/57597
## 2023-09-16 21:48:04 Annotating text fragment 40691/57597
## 2023-09-16 21:48:04 Annotating text fragment 40701/57597
## 2023-09-16 21:48:04 Annotating text fragment 40711/57597
## 2023-09-16 21:48:04 Annotating text fragment 40721/57597
## 2023-09-16 21:48:04 Annotating text fragment 40731/57597
```

```
## 2023-09-16 21:48:04 Annotating text fragment 40741/57597
## 2023-09-16 21:48:04 Annotating text fragment 40751/57597
## 2023-09-16 21:48:04 Annotating text fragment 40761/57597
## 2023-09-16 21:48:04 Annotating text fragment 40771/57597
## 2023-09-16 21:48:04 Annotating text fragment 40781/57597
## 2023-09-16 21:48:05 Annotating text fragment 40791/57597
## 2023-09-16 21:48:05 Annotating text fragment 40801/57597
## 2023-09-16 21:48:05 Annotating text fragment 40811/57597
## 2023-09-16 21:48:05 Annotating text fragment 40821/57597
## 2023-09-16 21:48:05 Annotating text fragment 40831/57597
## 2023-09-16 21:48:05 Annotating text fragment 40841/57597
## 2023-09-16 21:48:05 Annotating text fragment 40851/57597
## 2023-09-16 21:48:05 Annotating text fragment 40861/57597
## 2023-09-16 21:48:05 Annotating text fragment 40871/57597
## 2023-09-16 21:48:05 Annotating text fragment 40881/57597
## 2023-09-16 21:48:05 Annotating text fragment 40891/57597
## 2023-09-16 21:48:06 Annotating text fragment 40901/57597
## 2023-09-16 21:48:06 Annotating text fragment 40911/57597
## 2023-09-16 21:48:06 Annotating text fragment 40921/57597
## 2023-09-16 21:48:06 Annotating text fragment 40931/57597
## 2023-09-16 21:48:06 Annotating text fragment 40941/57597
## 2023-09-16 21:48:06 Annotating text fragment 40951/57597
## 2023-09-16 21:48:06 Annotating text fragment 40961/57597
## 2023-09-16 21:48:06 Annotating text fragment 40971/57597
## 2023-09-16 21:48:06 Annotating text fragment 40981/57597
## 2023-09-16 21:48:06 Annotating text fragment 40991/57597
## 2023-09-16 21:48:07 Annotating text fragment 41001/57597
## 2023-09-16 21:48:07 Annotating text fragment 41011/57597
## 2023-09-16 21:48:07 Annotating text fragment 41021/57597
## 2023-09-16 21:48:07 Annotating text fragment 41031/57597
## 2023-09-16 21:48:07 Annotating text fragment 41041/57597
## 2023-09-16 21:48:07 Annotating text fragment 41051/57597
## 2023-09-16 21:48:07 Annotating text fragment 41061/57597
## 2023-09-16 21:48:07 Annotating text fragment 41071/57597
## 2023-09-16 21:48:07 Annotating text fragment 41081/57597
## 2023-09-16 21:48:07 Annotating text fragment 41091/57597
## 2023-09-16 21:48:07 Annotating text fragment 41101/57597
## 2023-09-16 21:48:07 Annotating text fragment 41111/57597
## 2023-09-16 21:48:07 Annotating text fragment 41121/57597
## 2023-09-16 21:48:08 Annotating text fragment 41131/57597
## 2023-09-16 21:48:08 Annotating text fragment 41141/57597
## 2023-09-16 21:48:08 Annotating text fragment 41151/57597
## 2023-09-16 21:48:08 Annotating text fragment 41161/57597
## 2023-09-16 21:48:08 Annotating text fragment 41171/57597
## 2023-09-16 21:48:08 Annotating text fragment 41181/57597
## 2023-09-16 21:48:08 Annotating text fragment 41191/57597
## 2023-09-16 21:48:08 Annotating text fragment 41201/57597
## 2023-09-16 21:48:08 Annotating text fragment 41211/57597
## 2023-09-16 21:48:08 Annotating text fragment 41221/57597
## 2023-09-16 21:48:08 Annotating text fragment 41231/57597
## 2023-09-16 21:48:08 Annotating text fragment 41241/57597
## 2023-09-16 21:48:08 Annotating text fragment 41251/57597
## 2023-09-16 21:48:09 Annotating text fragment 41261/57597
## 2023-09-16 21:48:09 Annotating text fragment 41271/57597
```

```
## 2023-09-16 21:48:09 Annotating text fragment 41281/57597
## 2023-09-16 21:48:09 Annotating text fragment 41291/57597
## 2023-09-16 21:48:09 Annotating text fragment 41301/57597
## 2023-09-16 21:48:09 Annotating text fragment 41311/57597
## 2023-09-16 21:48:09 Annotating text fragment 41321/57597
## 2023-09-16 21:48:09 Annotating text fragment 41331/57597
## 2023-09-16 21:48:09 Annotating text fragment 41341/57597
## 2023-09-16 21:48:09 Annotating text fragment 41351/57597
## 2023-09-16 21:48:10 Annotating text fragment 41361/57597
## 2023-09-16 21:48:10 Annotating text fragment 41371/57597
## 2023-09-16 21:48:10 Annotating text fragment 41381/57597
## 2023-09-16 21:48:10 Annotating text fragment 41391/57597
## 2023-09-16 21:48:10 Annotating text fragment 41401/57597
## 2023-09-16 21:48:10 Annotating text fragment 41411/57597
## 2023-09-16 21:48:10 Annotating text fragment 41421/57597
## 2023-09-16 21:48:10 Annotating text fragment 41431/57597
## 2023-09-16 21:48:11 Annotating text fragment 41441/57597
## 2023-09-16 21:48:11 Annotating text fragment 41451/57597
## 2023-09-16 21:48:11 Annotating text fragment 41461/57597
## 2023-09-16 21:48:11 Annotating text fragment 41471/57597
## 2023-09-16 21:48:11 Annotating text fragment 41481/57597
## 2023-09-16 21:48:11 Annotating text fragment 41491/57597
## 2023-09-16 21:48:11 Annotating text fragment 41501/57597
## 2023-09-16 21:48:12 Annotating text fragment 41511/57597
## 2023-09-16 21:48:12 Annotating text fragment 41521/57597
## 2023-09-16 21:48:12 Annotating text fragment 41531/57597
## 2023-09-16 21:48:12 Annotating text fragment 41541/57597
## 2023-09-16 21:48:12 Annotating text fragment 41551/57597
## 2023-09-16 21:48:12 Annotating text fragment 41561/57597
## 2023-09-16 21:48:12 Annotating text fragment 41571/57597
## 2023-09-16 21:48:12 Annotating text fragment 41581/57597
## 2023-09-16 21:48:12 Annotating text fragment 41591/57597
## 2023-09-16 21:48:13 Annotating text fragment 41601/57597
## 2023-09-16 21:48:13 Annotating text fragment 41611/57597
## 2023-09-16 21:48:13 Annotating text fragment 41621/57597
## 2023-09-16 21:48:13 Annotating text fragment 41631/57597
## 2023-09-16 21:48:13 Annotating text fragment 41641/57597
## 2023-09-16 21:48:13 Annotating text fragment 41651/57597
## 2023-09-16 21:48:13 Annotating text fragment 41661/57597
## 2023-09-16 21:48:13 Annotating text fragment 41671/57597
## 2023-09-16 21:48:14 Annotating text fragment 41681/57597
## 2023-09-16 21:48:14 Annotating text fragment 41691/57597
## 2023-09-16 21:48:14 Annotating text fragment 41701/57597
## 2023-09-16 21:48:14 Annotating text fragment 41711/57597
## 2023-09-16 21:48:14 Annotating text fragment 41721/57597
## 2023-09-16 21:48:14 Annotating text fragment 41731/57597
## 2023-09-16 21:48:14 Annotating text fragment 41741/57597
## 2023-09-16 21:48:14 Annotating text fragment 41751/57597
## 2023-09-16 21:48:15 Annotating text fragment 41761/57597
## 2023-09-16 21:48:15 Annotating text fragment 41771/57597
## 2023-09-16 21:48:15 Annotating text fragment 41781/57597
## 2023-09-16 21:48:15 Annotating text fragment 41791/57597
## 2023-09-16 21:48:15 Annotating text fragment 41801/57597
## 2023-09-16 21:48:15 Annotating text fragment 41811/57597
```

```
## 2023-09-16 21:48:15 Annotating text fragment 41821/57597
## 2023-09-16 21:48:15 Annotating text fragment 41831/57597
## 2023-09-16 21:48:15 Annotating text fragment 41841/57597
## 2023-09-16 21:48:16 Annotating text fragment 41851/57597
## 2023-09-16 21:48:16 Annotating text fragment 41861/57597
## 2023-09-16 21:48:16 Annotating text fragment 41871/57597
## 2023-09-16 21:48:16 Annotating text fragment 41881/57597
## 2023-09-16 21:48:16 Annotating text fragment 41891/57597
## 2023-09-16 21:48:16 Annotating text fragment 41901/57597
## 2023-09-16 21:48:16 Annotating text fragment 41911/57597
## 2023-09-16 21:48:16 Annotating text fragment 41921/57597
## 2023-09-16 21:48:16 Annotating text fragment 41931/57597
## 2023-09-16 21:48:16 Annotating text fragment 41941/57597
## 2023-09-16 21:48:16 Annotating text fragment 41951/57597
## 2023-09-16 21:48:16 Annotating text fragment 41961/57597
## 2023-09-16 21:48:16 Annotating text fragment 41971/57597
## 2023-09-16 21:48:17 Annotating text fragment 41981/57597
## 2023-09-16 21:48:17 Annotating text fragment 41991/57597
## 2023-09-16 21:48:17 Annotating text fragment 42001/57597
## 2023-09-16 21:48:17 Annotating text fragment 42011/57597
## 2023-09-16 21:48:17 Annotating text fragment 42021/57597
## 2023-09-16 21:48:17 Annotating text fragment 42031/57597
## 2023-09-16 21:48:17 Annotating text fragment 42041/57597
## 2023-09-16 21:48:17 Annotating text fragment 42051/57597
## 2023-09-16 21:48:18 Annotating text fragment 42061/57597
## 2023-09-16 21:48:18 Annotating text fragment 42071/57597
## 2023-09-16 21:48:18 Annotating text fragment 42081/57597
## 2023-09-16 21:48:18 Annotating text fragment 42091/57597
## 2023-09-16 21:48:18 Annotating text fragment 42101/57597
## 2023-09-16 21:48:18 Annotating text fragment 42111/57597
## 2023-09-16 21:48:19 Annotating text fragment 42121/57597
## 2023-09-16 21:48:19 Annotating text fragment 42131/57597
## 2023-09-16 21:48:19 Annotating text fragment 42141/57597
## 2023-09-16 21:48:19 Annotating text fragment 42151/57597
## 2023-09-16 21:48:19 Annotating text fragment 42161/57597
## 2023-09-16 21:48:20 Annotating text fragment 42171/57597
## 2023-09-16 21:48:20 Annotating text fragment 42181/57597
## 2023-09-16 21:48:20 Annotating text fragment 42191/57597
## 2023-09-16 21:48:20 Annotating text fragment 42201/57597
## 2023-09-16 21:48:20 Annotating text fragment 42211/57597
## 2023-09-16 21:48:20 Annotating text fragment 42221/57597
## 2023-09-16 21:48:20 Annotating text fragment 42231/57597
## 2023-09-16 21:48:20 Annotating text fragment 42241/57597
## 2023-09-16 21:48:20 Annotating text fragment 42251/57597
## 2023-09-16 21:48:20 Annotating text fragment 42261/57597
## 2023-09-16 21:48:20 Annotating text fragment 42271/57597
## 2023-09-16 21:48:20 Annotating text fragment 42281/57597
## 2023-09-16 21:48:21 Annotating text fragment 42291/57597
## 2023-09-16 21:48:21 Annotating text fragment 42301/57597
## 2023-09-16 21:48:21 Annotating text fragment 42311/57597
## 2023-09-16 21:48:21 Annotating text fragment 42321/57597
## 2023-09-16 21:48:21 Annotating text fragment 42331/57597
## 2023-09-16 21:48:21 Annotating text fragment 42341/57597
## 2023-09-16 21:48:21 Annotating text fragment 42351/57597
```

```
## 2023-09-16 21:48:21 Annotating text fragment 42361/57597
## 2023-09-16 21:48:21 Annotating text fragment 42371/57597
## 2023-09-16 21:48:21 Annotating text fragment 42381/57597
## 2023-09-16 21:48:22 Annotating text fragment 42391/57597
## 2023-09-16 21:48:22 Annotating text fragment 42401/57597
## 2023-09-16 21:48:22 Annotating text fragment 42411/57597
## 2023-09-16 21:48:22 Annotating text fragment 42421/57597
## 2023-09-16 21:48:22 Annotating text fragment 42431/57597
## 2023-09-16 21:48:22 Annotating text fragment 42441/57597
## 2023-09-16 21:48:22 Annotating text fragment 42451/57597
## 2023-09-16 21:48:22 Annotating text fragment 42461/57597
## 2023-09-16 21:48:22 Annotating text fragment 42471/57597
## 2023-09-16 21:48:22 Annotating text fragment 42481/57597
## 2023-09-16 21:48:22 Annotating text fragment 42491/57597
## 2023-09-16 21:48:22 Annotating text fragment 42501/57597
## 2023-09-16 21:48:22 Annotating text fragment 42511/57597
## 2023-09-16 21:48:22 Annotating text fragment 42521/57597
## 2023-09-16 21:48:23 Annotating text fragment 42531/57597
## 2023-09-16 21:48:23 Annotating text fragment 42541/57597
## 2023-09-16 21:48:23 Annotating text fragment 42551/57597
## 2023-09-16 21:48:23 Annotating text fragment 42561/57597
## 2023-09-16 21:48:23 Annotating text fragment 42571/57597
## 2023-09-16 21:48:23 Annotating text fragment 42581/57597
## 2023-09-16 21:48:23 Annotating text fragment 42591/57597
## 2023-09-16 21:48:24 Annotating text fragment 42601/57597
## 2023-09-16 21:48:24 Annotating text fragment 42611/57597
## 2023-09-16 21:48:24 Annotating text fragment 42621/57597
## 2023-09-16 21:48:24 Annotating text fragment 42631/57597
## 2023-09-16 21:48:24 Annotating text fragment 42641/57597
## 2023-09-16 21:48:24 Annotating text fragment 42651/57597
## 2023-09-16 21:48:24 Annotating text fragment 42661/57597
## 2023-09-16 21:48:24 Annotating text fragment 42671/57597
## 2023-09-16 21:48:24 Annotating text fragment 42681/57597
## 2023-09-16 21:48:24 Annotating text fragment 42691/57597
## 2023-09-16 21:48:24 Annotating text fragment 42701/57597
## 2023-09-16 21:48:25 Annotating text fragment 42711/57597
## 2023-09-16 21:48:25 Annotating text fragment 42721/57597
## 2023-09-16 21:48:25 Annotating text fragment 42731/57597
## 2023-09-16 21:48:25 Annotating text fragment 42741/57597
## 2023-09-16 21:48:25 Annotating text fragment 42751/57597
## 2023-09-16 21:48:25 Annotating text fragment 42761/57597
## 2023-09-16 21:48:25 Annotating text fragment 42771/57597
## 2023-09-16 21:48:25 Annotating text fragment 42781/57597
## 2023-09-16 21:48:25 Annotating text fragment 42791/57597
## 2023-09-16 21:48:25 Annotating text fragment 42801/57597
## 2023-09-16 21:48:26 Annotating text fragment 42811/57597
## 2023-09-16 21:48:26 Annotating text fragment 42821/57597
## 2023-09-16 21:48:26 Annotating text fragment 42831/57597
## 2023-09-16 21:48:26 Annotating text fragment 42841/57597
## 2023-09-16 21:48:26 Annotating text fragment 42851/57597
## 2023-09-16 21:48:26 Annotating text fragment 42861/57597
## 2023-09-16 21:48:26 Annotating text fragment 42871/57597
## 2023-09-16 21:48:26 Annotating text fragment 42881/57597
## 2023-09-16 21:48:26 Annotating text fragment 42891/57597
```

```
## 2023-09-16 21:48:27 Annotating text fragment 42901/57597
## 2023-09-16 21:48:27 Annotating text fragment 42911/57597
## 2023-09-16 21:48:27 Annotating text fragment 42921/57597
## 2023-09-16 21:48:27 Annotating text fragment 42931/57597
## 2023-09-16 21:48:27 Annotating text fragment 42941/57597
## 2023-09-16 21:48:27 Annotating text fragment 42951/57597
## 2023-09-16 21:48:27 Annotating text fragment 42961/57597
## 2023-09-16 21:48:27 Annotating text fragment 42971/57597
## 2023-09-16 21:48:27 Annotating text fragment 42981/57597
## 2023-09-16 21:48:28 Annotating text fragment 42991/57597
## 2023-09-16 21:48:28 Annotating text fragment 43001/57597
## 2023-09-16 21:48:28 Annotating text fragment 43011/57597
## 2023-09-16 21:48:28 Annotating text fragment 43021/57597
## 2023-09-16 21:48:28 Annotating text fragment 43031/57597
## 2023-09-16 21:48:28 Annotating text fragment 43041/57597
## 2023-09-16 21:48:28 Annotating text fragment 43051/57597
## 2023-09-16 21:48:28 Annotating text fragment 43061/57597
## 2023-09-16 21:48:28 Annotating text fragment 43071/57597
## 2023-09-16 21:48:28 Annotating text fragment 43081/57597
## 2023-09-16 21:48:29 Annotating text fragment 43091/57597
## 2023-09-16 21:48:29 Annotating text fragment 43101/57597
## 2023-09-16 21:48:29 Annotating text fragment 43111/57597
## 2023-09-16 21:48:29 Annotating text fragment 43121/57597
## 2023-09-16 21:48:29 Annotating text fragment 43131/57597
## 2023-09-16 21:48:29 Annotating text fragment 43141/57597
## 2023-09-16 21:48:29 Annotating text fragment 43151/57597
## 2023-09-16 21:48:29 Annotating text fragment 43161/57597
## 2023-09-16 21:48:30 Annotating text fragment 43171/57597
## 2023-09-16 21:48:30 Annotating text fragment 43181/57597
## 2023-09-16 21:48:30 Annotating text fragment 43191/57597
## 2023-09-16 21:48:30 Annotating text fragment 43201/57597
## 2023-09-16 21:48:30 Annotating text fragment 43211/57597
## 2023-09-16 21:48:30 Annotating text fragment 43221/57597
## 2023-09-16 21:48:30 Annotating text fragment 43231/57597
## 2023-09-16 21:48:30 Annotating text fragment 43241/57597
## 2023-09-16 21:48:30 Annotating text fragment 43251/57597
## 2023-09-16 21:48:31 Annotating text fragment 43261/57597
## 2023-09-16 21:48:31 Annotating text fragment 43271/57597
## 2023-09-16 21:48:31 Annotating text fragment 43281/57597
## 2023-09-16 21:48:31 Annotating text fragment 43291/57597
## 2023-09-16 21:48:31 Annotating text fragment 43301/57597
## 2023-09-16 21:48:31 Annotating text fragment 43311/57597
## 2023-09-16 21:48:31 Annotating text fragment 43321/57597
## 2023-09-16 21:48:31 Annotating text fragment 43331/57597
## 2023-09-16 21:48:31 Annotating text fragment 43341/57597
## 2023-09-16 21:48:31 Annotating text fragment 43351/57597
## 2023-09-16 21:48:32 Annotating text fragment 43361/57597
## 2023-09-16 21:48:32 Annotating text fragment 43371/57597
## 2023-09-16 21:48:32 Annotating text fragment 43381/57597
## 2023-09-16 21:48:32 Annotating text fragment 43391/57597
## 2023-09-16 21:48:32 Annotating text fragment 43401/57597
## 2023-09-16 21:48:32 Annotating text fragment 43411/57597
## 2023-09-16 21:48:32 Annotating text fragment 43421/57597
## 2023-09-16 21:48:32 Annotating text fragment 43431/57597
```

```
## 2023-09-16 21:48:32 Annotating text fragment 43441/57597
## 2023-09-16 21:48:32 Annotating text fragment 43451/57597
## 2023-09-16 21:48:32 Annotating text fragment 43461/57597
## 2023-09-16 21:48:33 Annotating text fragment 43471/57597
## 2023-09-16 21:48:33 Annotating text fragment 43481/57597
## 2023-09-16 21:48:33 Annotating text fragment 43491/57597
## 2023-09-16 21:48:33 Annotating text fragment 43501/57597
## 2023-09-16 21:48:33 Annotating text fragment 43511/57597
## 2023-09-16 21:48:33 Annotating text fragment 43521/57597
## 2023-09-16 21:48:33 Annotating text fragment 43531/57597
## 2023-09-16 21:48:33 Annotating text fragment 43541/57597
## 2023-09-16 21:48:34 Annotating text fragment 43551/57597
## 2023-09-16 21:48:34 Annotating text fragment 43561/57597
## 2023-09-16 21:48:34 Annotating text fragment 43571/57597
## 2023-09-16 21:48:34 Annotating text fragment 43581/57597
## 2023-09-16 21:48:34 Annotating text fragment 43591/57597
## 2023-09-16 21:48:34 Annotating text fragment 43601/57597
## 2023-09-16 21:48:34 Annotating text fragment 43611/57597
## 2023-09-16 21:48:34 Annotating text fragment 43621/57597
## 2023-09-16 21:48:34 Annotating text fragment 43631/57597
## 2023-09-16 21:48:34 Annotating text fragment 43641/57597
## 2023-09-16 21:48:35 Annotating text fragment 43651/57597
## 2023-09-16 21:48:35 Annotating text fragment 43661/57597
## 2023-09-16 21:48:35 Annotating text fragment 43671/57597
## 2023-09-16 21:48:35 Annotating text fragment 43681/57597
## 2023-09-16 21:48:35 Annotating text fragment 43691/57597
## 2023-09-16 21:48:35 Annotating text fragment 43701/57597
## 2023-09-16 21:48:35 Annotating text fragment 43711/57597
## 2023-09-16 21:48:35 Annotating text fragment 43721/57597
## 2023-09-16 21:48:35 Annotating text fragment 43731/57597
## 2023-09-16 21:48:35 Annotating text fragment 43741/57597
## 2023-09-16 21:48:35 Annotating text fragment 43751/57597
## 2023-09-16 21:48:36 Annotating text fragment 43761/57597
## 2023-09-16 21:48:36 Annotating text fragment 43771/57597
## 2023-09-16 21:48:36 Annotating text fragment 43781/57597
## 2023-09-16 21:48:36 Annotating text fragment 43791/57597
## 2023-09-16 21:48:36 Annotating text fragment 43801/57597
## 2023-09-16 21:48:36 Annotating text fragment 43811/57597
## 2023-09-16 21:48:36 Annotating text fragment 43821/57597
## 2023-09-16 21:48:36 Annotating text fragment 43831/57597
## 2023-09-16 21:48:37 Annotating text fragment 43841/57597
## 2023-09-16 21:48:37 Annotating text fragment 43851/57597
## 2023-09-16 21:48:37 Annotating text fragment 43861/57597
## 2023-09-16 21:48:37 Annotating text fragment 43871/57597
## 2023-09-16 21:48:37 Annotating text fragment 43881/57597
## 2023-09-16 21:48:37 Annotating text fragment 43891/57597
## 2023-09-16 21:48:37 Annotating text fragment 43901/57597
## 2023-09-16 21:48:37 Annotating text fragment 43911/57597
## 2023-09-16 21:48:37 Annotating text fragment 43921/57597
## 2023-09-16 21:48:37 Annotating text fragment 43931/57597
## 2023-09-16 21:48:38 Annotating text fragment 43941/57597
## 2023-09-16 21:48:38 Annotating text fragment 43951/57597
## 2023-09-16 21:48:38 Annotating text fragment 43961/57597
## 2023-09-16 21:48:38 Annotating text fragment 43971/57597
```

```
## 2023-09-16 21:48:38 Annotating text fragment 43981/57597
## 2023-09-16 21:48:38 Annotating text fragment 43991/57597
## 2023-09-16 21:48:38 Annotating text fragment 44001/57597
## 2023-09-16 21:48:38 Annotating text fragment 44011/57597
## 2023-09-16 21:48:39 Annotating text fragment 44021/57597
## 2023-09-16 21:48:39 Annotating text fragment 44031/57597
## 2023-09-16 21:48:39 Annotating text fragment 44041/57597
## 2023-09-16 21:48:39 Annotating text fragment 44051/57597
## 2023-09-16 21:48:39 Annotating text fragment 44061/57597
## 2023-09-16 21:48:39 Annotating text fragment 44071/57597
## 2023-09-16 21:48:39 Annotating text fragment 44081/57597
## 2023-09-16 21:48:39 Annotating text fragment 44091/57597
## 2023-09-16 21:48:39 Annotating text fragment 44101/57597
## 2023-09-16 21:48:39 Annotating text fragment 44111/57597
## 2023-09-16 21:48:40 Annotating text fragment 44121/57597
## 2023-09-16 21:48:40 Annotating text fragment 44131/57597
## 2023-09-16 21:48:40 Annotating text fragment 44141/57597
## 2023-09-16 21:48:40 Annotating text fragment 44151/57597
## 2023-09-16 21:48:40 Annotating text fragment 44161/57597
## 2023-09-16 21:48:40 Annotating text fragment 44171/57597
## 2023-09-16 21:48:40 Annotating text fragment 44181/57597
## 2023-09-16 21:48:40 Annotating text fragment 44191/57597
## 2023-09-16 21:48:41 Annotating text fragment 44201/57597
## 2023-09-16 21:48:41 Annotating text fragment 44211/57597
## 2023-09-16 21:48:41 Annotating text fragment 44221/57597
## 2023-09-16 21:48:41 Annotating text fragment 44231/57597
## 2023-09-16 21:48:41 Annotating text fragment 44241/57597
## 2023-09-16 21:48:41 Annotating text fragment 44251/57597
## 2023-09-16 21:48:41 Annotating text fragment 44261/57597
## 2023-09-16 21:48:41 Annotating text fragment 44271/57597
## 2023-09-16 21:48:41 Annotating text fragment 44281/57597
## 2023-09-16 21:48:41 Annotating text fragment 44291/57597
## 2023-09-16 21:48:41 Annotating text fragment 44301/57597
## 2023-09-16 21:48:41 Annotating text fragment 44311/57597
## 2023-09-16 21:48:41 Annotating text fragment 44321/57597
## 2023-09-16 21:48:42 Annotating text fragment 44331/57597
## 2023-09-16 21:48:42 Annotating text fragment 44341/57597
## 2023-09-16 21:48:42 Annotating text fragment 44351/57597
## 2023-09-16 21:48:42 Annotating text fragment 44361/57597
## 2023-09-16 21:48:42 Annotating text fragment 44371/57597
## 2023-09-16 21:48:42 Annotating text fragment 44381/57597
## 2023-09-16 21:48:42 Annotating text fragment 44391/57597
## 2023-09-16 21:48:42 Annotating text fragment 44401/57597
## 2023-09-16 21:48:42 Annotating text fragment 44411/57597
## 2023-09-16 21:48:42 Annotating text fragment 44421/57597
## 2023-09-16 21:48:43 Annotating text fragment 44431/57597
## 2023-09-16 21:48:43 Annotating text fragment 44441/57597
## 2023-09-16 21:48:43 Annotating text fragment 44451/57597
## 2023-09-16 21:48:43 Annotating text fragment 44461/57597
## 2023-09-16 21:48:43 Annotating text fragment 44471/57597
## 2023-09-16 21:48:43 Annotating text fragment 44481/57597
## 2023-09-16 21:48:43 Annotating text fragment 44491/57597
## 2023-09-16 21:48:43 Annotating text fragment 44501/57597
## 2023-09-16 21:48:44 Annotating text fragment 44511/57597
```

```
## 2023-09-16 21:48:44 Annotating text fragment 44521/57597
## 2023-09-16 21:48:44 Annotating text fragment 44531/57597
## 2023-09-16 21:48:44 Annotating text fragment 44541/57597
## 2023-09-16 21:48:44 Annotating text fragment 44551/57597
## 2023-09-16 21:48:44 Annotating text fragment 44561/57597
## 2023-09-16 21:48:44 Annotating text fragment 44571/57597
## 2023-09-16 21:48:44 Annotating text fragment 44581/57597
## 2023-09-16 21:48:44 Annotating text fragment 44591/57597
## 2023-09-16 21:48:44 Annotating text fragment 44601/57597
## 2023-09-16 21:48:45 Annotating text fragment 44611/57597
## 2023-09-16 21:48:45 Annotating text fragment 44621/57597
## 2023-09-16 21:48:45 Annotating text fragment 44631/57597
## 2023-09-16 21:48:45 Annotating text fragment 44641/57597
## 2023-09-16 21:48:45 Annotating text fragment 44651/57597
## 2023-09-16 21:48:45 Annotating text fragment 44661/57597
## 2023-09-16 21:48:45 Annotating text fragment 44671/57597
## 2023-09-16 21:48:45 Annotating text fragment 44681/57597
## 2023-09-16 21:48:46 Annotating text fragment 44691/57597
## 2023-09-16 21:48:46 Annotating text fragment 44701/57597
## 2023-09-16 21:48:46 Annotating text fragment 44711/57597
## 2023-09-16 21:48:46 Annotating text fragment 44721/57597
## 2023-09-16 21:48:46 Annotating text fragment 44731/57597
## 2023-09-16 21:48:46 Annotating text fragment 44741/57597
## 2023-09-16 21:48:46 Annotating text fragment 44751/57597
## 2023-09-16 21:48:46 Annotating text fragment 44761/57597
## 2023-09-16 21:48:46 Annotating text fragment 44771/57597
## 2023-09-16 21:48:46 Annotating text fragment 44781/57597
## 2023-09-16 21:48:46 Annotating text fragment 44791/57597
## 2023-09-16 21:48:47 Annotating text fragment 44801/57597
## 2023-09-16 21:48:47 Annotating text fragment 44811/57597
## 2023-09-16 21:48:47 Annotating text fragment 44821/57597
## 2023-09-16 21:48:47 Annotating text fragment 44831/57597
## 2023-09-16 21:48:47 Annotating text fragment 44841/57597
## 2023-09-16 21:48:47 Annotating text fragment 44851/57597
## 2023-09-16 21:48:47 Annotating text fragment 44861/57597
## 2023-09-16 21:48:47 Annotating text fragment 44871/57597
## 2023-09-16 21:48:47 Annotating text fragment 44881/57597
## 2023-09-16 21:48:48 Annotating text fragment 44891/57597
## 2023-09-16 21:48:48 Annotating text fragment 44901/57597
## 2023-09-16 21:48:48 Annotating text fragment 44911/57597
## 2023-09-16 21:48:48 Annotating text fragment 44921/57597
## 2023-09-16 21:48:48 Annotating text fragment 44931/57597
## 2023-09-16 21:48:48 Annotating text fragment 44941/57597
## 2023-09-16 21:48:48 Annotating text fragment 44951/57597
## 2023-09-16 21:48:48 Annotating text fragment 44961/57597
## 2023-09-16 21:48:48 Annotating text fragment 44971/57597
## 2023-09-16 21:48:48 Annotating text fragment 44981/57597
## 2023-09-16 21:48:49 Annotating text fragment 44991/57597
## 2023-09-16 21:48:49 Annotating text fragment 45001/57597
## 2023-09-16 21:48:49 Annotating text fragment 45011/57597
## 2023-09-16 21:48:49 Annotating text fragment 45021/57597
## 2023-09-16 21:48:49 Annotating text fragment 45031/57597
## 2023-09-16 21:48:49 Annotating text fragment 45041/57597
## 2023-09-16 21:48:49 Annotating text fragment 45051/57597
```

```
## 2023-09-16 21:48:49 Annotating text fragment 45061/57597
## 2023-09-16 21:48:49 Annotating text fragment 45071/57597
## 2023-09-16 21:48:50 Annotating text fragment 45081/57597
## 2023-09-16 21:48:50 Annotating text fragment 45091/57597
## 2023-09-16 21:48:50 Annotating text fragment 45101/57597
## 2023-09-16 21:48:51 Annotating text fragment 45111/57597
## 2023-09-16 21:48:51 Annotating text fragment 45121/57597
## 2023-09-16 21:48:51 Annotating text fragment 45131/57597
## 2023-09-16 21:48:51 Annotating text fragment 45141/57597
## 2023-09-16 21:48:51 Annotating text fragment 45151/57597
## 2023-09-16 21:48:51 Annotating text fragment 45161/57597
## 2023-09-16 21:48:51 Annotating text fragment 45171/57597
## 2023-09-16 21:48:51 Annotating text fragment 45181/57597
## 2023-09-16 21:48:52 Annotating text fragment 45191/57597
## 2023-09-16 21:48:52 Annotating text fragment 45201/57597
## 2023-09-16 21:48:52 Annotating text fragment 45211/57597
## 2023-09-16 21:48:52 Annotating text fragment 45221/57597
## 2023-09-16 21:48:52 Annotating text fragment 45231/57597
## 2023-09-16 21:48:52 Annotating text fragment 45241/57597
## 2023-09-16 21:48:52 Annotating text fragment 45251/57597
## 2023-09-16 21:48:52 Annotating text fragment 45261/57597
## 2023-09-16 21:48:52 Annotating text fragment 45271/57597
## 2023-09-16 21:48:52 Annotating text fragment 45281/57597
## 2023-09-16 21:48:53 Annotating text fragment 45291/57597
## 2023-09-16 21:48:53 Annotating text fragment 45301/57597
## 2023-09-16 21:48:53 Annotating text fragment 45311/57597
## 2023-09-16 21:48:53 Annotating text fragment 45321/57597
## 2023-09-16 21:48:53 Annotating text fragment 45331/57597
## 2023-09-16 21:48:53 Annotating text fragment 45341/57597
## 2023-09-16 21:48:53 Annotating text fragment 45351/57597
## 2023-09-16 21:48:53 Annotating text fragment 45361/57597
## 2023-09-16 21:48:53 Annotating text fragment 45371/57597
## 2023-09-16 21:48:54 Annotating text fragment 45381/57597
## 2023-09-16 21:48:54 Annotating text fragment 45391/57597
## 2023-09-16 21:48:54 Annotating text fragment 45401/57597
## 2023-09-16 21:48:54 Annotating text fragment 45411/57597
## 2023-09-16 21:48:54 Annotating text fragment 45421/57597
## 2023-09-16 21:48:54 Annotating text fragment 45431/57597
## 2023-09-16 21:48:55 Annotating text fragment 45441/57597
## 2023-09-16 21:48:55 Annotating text fragment 45451/57597
## 2023-09-16 21:48:55 Annotating text fragment 45461/57597
## 2023-09-16 21:48:55 Annotating text fragment 45471/57597
## 2023-09-16 21:48:55 Annotating text fragment 45481/57597
## 2023-09-16 21:48:55 Annotating text fragment 45491/57597
## 2023-09-16 21:48:55 Annotating text fragment 45501/57597
## 2023-09-16 21:48:55 Annotating text fragment 45511/57597
## 2023-09-16 21:48:56 Annotating text fragment 45521/57597
## 2023-09-16 21:48:56 Annotating text fragment 45531/57597
## 2023-09-16 21:48:56 Annotating text fragment 45541/57597
## 2023-09-16 21:48:56 Annotating text fragment 45551/57597
## 2023-09-16 21:48:56 Annotating text fragment 45561/57597
## 2023-09-16 21:48:56 Annotating text fragment 45571/57597
## 2023-09-16 21:48:56 Annotating text fragment 45581/57597
## 2023-09-16 21:48:56 Annotating text fragment 45591/57597
```

```
## 2023-09-16 21:48:56 Annotating text fragment 45601/57597
## 2023-09-16 21:48:56 Annotating text fragment 45611/57597
## 2023-09-16 21:48:57 Annotating text fragment 45621/57597
## 2023-09-16 21:48:57 Annotating text fragment 45631/57597
## 2023-09-16 21:48:57 Annotating text fragment 45641/57597
## 2023-09-16 21:48:57 Annotating text fragment 45651/57597
## 2023-09-16 21:48:57 Annotating text fragment 45661/57597
## 2023-09-16 21:48:57 Annotating text fragment 45671/57597
## 2023-09-16 21:48:57 Annotating text fragment 45681/57597
## 2023-09-16 21:48:57 Annotating text fragment 45691/57597
## 2023-09-16 21:48:57 Annotating text fragment 45701/57597
## 2023-09-16 21:48:57 Annotating text fragment 45711/57597
## 2023-09-16 21:48:57 Annotating text fragment 45721/57597
## 2023-09-16 21:48:58 Annotating text fragment 45731/57597
## 2023-09-16 21:48:58 Annotating text fragment 45741/57597
## 2023-09-16 21:48:58 Annotating text fragment 45751/57597
## 2023-09-16 21:48:58 Annotating text fragment 45761/57597
## 2023-09-16 21:48:58 Annotating text fragment 45771/57597
## 2023-09-16 21:48:58 Annotating text fragment 45781/57597
## 2023-09-16 21:48:58 Annotating text fragment 45791/57597
## 2023-09-16 21:48:58 Annotating text fragment 45801/57597
## 2023-09-16 21:48:58 Annotating text fragment 45811/57597
## 2023-09-16 21:48:58 Annotating text fragment 45821/57597
## 2023-09-16 21:48:59 Annotating text fragment 45831/57597
## 2023-09-16 21:48:59 Annotating text fragment 45841/57597
## 2023-09-16 21:48:59 Annotating text fragment 45851/57597
## 2023-09-16 21:48:59 Annotating text fragment 45861/57597
## 2023-09-16 21:48:59 Annotating text fragment 45871/57597
## 2023-09-16 21:48:59 Annotating text fragment 45881/57597
## 2023-09-16 21:48:59 Annotating text fragment 45891/57597
## 2023-09-16 21:48:59 Annotating text fragment 45901/57597
## 2023-09-16 21:48:59 Annotating text fragment 45911/57597
## 2023-09-16 21:49:00 Annotating text fragment 45921/57597
## 2023-09-16 21:49:00 Annotating text fragment 45931/57597
## 2023-09-16 21:49:00 Annotating text fragment 45941/57597
## 2023-09-16 21:49:00 Annotating text fragment 45951/57597
## 2023-09-16 21:49:00 Annotating text fragment 45961/57597
## 2023-09-16 21:49:01 Annotating text fragment 45971/57597
## 2023-09-16 21:49:01 Annotating text fragment 45981/57597
## 2023-09-16 21:49:01 Annotating text fragment 45991/57597
## 2023-09-16 21:49:01 Annotating text fragment 46001/57597
## 2023-09-16 21:49:01 Annotating text fragment 46011/57597
## 2023-09-16 21:49:01 Annotating text fragment 46021/57597
## 2023-09-16 21:49:01 Annotating text fragment 46031/57597
## 2023-09-16 21:49:02 Annotating text fragment 46041/57597
## 2023-09-16 21:49:02 Annotating text fragment 46051/57597
## 2023-09-16 21:49:02 Annotating text fragment 46061/57597
## 2023-09-16 21:49:02 Annotating text fragment 46071/57597
## 2023-09-16 21:49:02 Annotating text fragment 46081/57597
## 2023-09-16 21:49:02 Annotating text fragment 46091/57597
## 2023-09-16 21:49:02 Annotating text fragment 46101/57597
## 2023-09-16 21:49:02 Annotating text fragment 46111/57597
## 2023-09-16 21:49:02 Annotating text fragment 46121/57597
## 2023-09-16 21:49:03 Annotating text fragment 46131/57597
```

```
## 2023-09-16 21:49:03 Annotating text fragment 46141/57597
## 2023-09-16 21:49:03 Annotating text fragment 46151/57597
## 2023-09-16 21:49:03 Annotating text fragment 46161/57597
## 2023-09-16 21:49:03 Annotating text fragment 46171/57597
## 2023-09-16 21:49:03 Annotating text fragment 46181/57597
## 2023-09-16 21:49:04 Annotating text fragment 46191/57597
## 2023-09-16 21:49:04 Annotating text fragment 46201/57597
## 2023-09-16 21:49:04 Annotating text fragment 46211/57597
## 2023-09-16 21:49:04 Annotating text fragment 46221/57597
## 2023-09-16 21:49:04 Annotating text fragment 46231/57597
## 2023-09-16 21:49:04 Annotating text fragment 46241/57597
## 2023-09-16 21:49:04 Annotating text fragment 46251/57597
## 2023-09-16 21:49:04 Annotating text fragment 46261/57597
## 2023-09-16 21:49:04 Annotating text fragment 46271/57597
## 2023-09-16 21:49:05 Annotating text fragment 46281/57597
## 2023-09-16 21:49:05 Annotating text fragment 46291/57597
## 2023-09-16 21:49:05 Annotating text fragment 46301/57597
## 2023-09-16 21:49:05 Annotating text fragment 46311/57597
## 2023-09-16 21:49:05 Annotating text fragment 46321/57597
## 2023-09-16 21:49:05 Annotating text fragment 46331/57597
## 2023-09-16 21:49:05 Annotating text fragment 46341/57597
## 2023-09-16 21:49:05 Annotating text fragment 46351/57597
## 2023-09-16 21:49:05 Annotating text fragment 46361/57597
## 2023-09-16 21:49:05 Annotating text fragment 46371/57597
## 2023-09-16 21:49:05 Annotating text fragment 46381/57597
## 2023-09-16 21:49:06 Annotating text fragment 46391/57597
## 2023-09-16 21:49:06 Annotating text fragment 46401/57597
## 2023-09-16 21:49:06 Annotating text fragment 46411/57597
## 2023-09-16 21:49:06 Annotating text fragment 46421/57597
## 2023-09-16 21:49:06 Annotating text fragment 46431/57597
## 2023-09-16 21:49:06 Annotating text fragment 46441/57597
## 2023-09-16 21:49:06 Annotating text fragment 46451/57597
## 2023-09-16 21:49:06 Annotating text fragment 46461/57597
## 2023-09-16 21:49:07 Annotating text fragment 46471/57597
## 2023-09-16 21:49:07 Annotating text fragment 46481/57597
## 2023-09-16 21:49:07 Annotating text fragment 46491/57597
## 2023-09-16 21:49:07 Annotating text fragment 46501/57597
## 2023-09-16 21:49:07 Annotating text fragment 46511/57597
## 2023-09-16 21:49:08 Annotating text fragment 46521/57597
## 2023-09-16 21:49:08 Annotating text fragment 46531/57597
## 2023-09-16 21:49:08 Annotating text fragment 46541/57597
## 2023-09-16 21:49:08 Annotating text fragment 46551/57597
## 2023-09-16 21:49:08 Annotating text fragment 46561/57597
## 2023-09-16 21:49:08 Annotating text fragment 46571/57597
## 2023-09-16 21:49:08 Annotating text fragment 46581/57597
## 2023-09-16 21:49:08 Annotating text fragment 46591/57597
## 2023-09-16 21:49:08 Annotating text fragment 46601/57597
## 2023-09-16 21:49:08 Annotating text fragment 46611/57597
## 2023-09-16 21:49:08 Annotating text fragment 46621/57597
## 2023-09-16 21:49:08 Annotating text fragment 46631/57597
## 2023-09-16 21:49:09 Annotating text fragment 46641/57597
## 2023-09-16 21:49:09 Annotating text fragment 46651/57597
## 2023-09-16 21:49:09 Annotating text fragment 46661/57597
## 2023-09-16 21:49:09 Annotating text fragment 46671/57597
```

```
## 2023-09-16 21:49:09 Annotating text fragment 46681/57597
## 2023-09-16 21:49:09 Annotating text fragment 46691/57597
## 2023-09-16 21:49:09 Annotating text fragment 46701/57597
## 2023-09-16 21:49:09 Annotating text fragment 46711/57597
## 2023-09-16 21:49:09 Annotating text fragment 46721/57597
## 2023-09-16 21:49:09 Annotating text fragment 46731/57597
## 2023-09-16 21:49:10 Annotating text fragment 46741/57597
## 2023-09-16 21:49:10 Annotating text fragment 46751/57597
## 2023-09-16 21:49:10 Annotating text fragment 46761/57597
## 2023-09-16 21:49:10 Annotating text fragment 46771/57597
## 2023-09-16 21:49:10 Annotating text fragment 46781/57597
## 2023-09-16 21:49:10 Annotating text fragment 46791/57597
## 2023-09-16 21:49:10 Annotating text fragment 46801/57597
## 2023-09-16 21:49:10 Annotating text fragment 46811/57597
## 2023-09-16 21:49:10 Annotating text fragment 46821/57597
## 2023-09-16 21:49:10 Annotating text fragment 46831/57597
## 2023-09-16 21:49:11 Annotating text fragment 46841/57597
## 2023-09-16 21:49:11 Annotating text fragment 46851/57597
## 2023-09-16 21:49:11 Annotating text fragment 46861/57597
## 2023-09-16 21:49:11 Annotating text fragment 46871/57597
## 2023-09-16 21:49:11 Annotating text fragment 46881/57597
## 2023-09-16 21:49:11 Annotating text fragment 46891/57597
## 2023-09-16 21:49:11 Annotating text fragment 46901/57597
## 2023-09-16 21:49:12 Annotating text fragment 46911/57597
## 2023-09-16 21:49:12 Annotating text fragment 46921/57597
## 2023-09-16 21:49:12 Annotating text fragment 46931/57597
## 2023-09-16 21:49:12 Annotating text fragment 46941/57597
## 2023-09-16 21:49:12 Annotating text fragment 46951/57597
## 2023-09-16 21:49:12 Annotating text fragment 46961/57597
## 2023-09-16 21:49:12 Annotating text fragment 46971/57597
## 2023-09-16 21:49:12 Annotating text fragment 46981/57597
## 2023-09-16 21:49:12 Annotating text fragment 46991/57597
## 2023-09-16 21:49:12 Annotating text fragment 47001/57597
## 2023-09-16 21:49:12 Annotating text fragment 47011/57597
## 2023-09-16 21:49:13 Annotating text fragment 47021/57597
## 2023-09-16 21:49:13 Annotating text fragment 47031/57597
## 2023-09-16 21:49:13 Annotating text fragment 47041/57597
## 2023-09-16 21:49:13 Annotating text fragment 47051/57597
## 2023-09-16 21:49:13 Annotating text fragment 47061/57597
## 2023-09-16 21:49:13 Annotating text fragment 47071/57597
## 2023-09-16 21:49:14 Annotating text fragment 47081/57597
## 2023-09-16 21:49:14 Annotating text fragment 47091/57597
## 2023-09-16 21:49:14 Annotating text fragment 47101/57597
## 2023-09-16 21:49:14 Annotating text fragment 47111/57597
## 2023-09-16 21:49:14 Annotating text fragment 47121/57597
## 2023-09-16 21:49:14 Annotating text fragment 47131/57597
## 2023-09-16 21:49:14 Annotating text fragment 47141/57597
## 2023-09-16 21:49:14 Annotating text fragment 47151/57597
## 2023-09-16 21:49:14 Annotating text fragment 47161/57597
## 2023-09-16 21:49:14 Annotating text fragment 47171/57597
## 2023-09-16 21:49:15 Annotating text fragment 47181/57597
## 2023-09-16 21:49:15 Annotating text fragment 47191/57597
## 2023-09-16 21:49:15 Annotating text fragment 47201/57597
## 2023-09-16 21:49:15 Annotating text fragment 47211/57597
```

```
## 2023-09-16 21:49:15 Annotating text fragment 47221/57597
## 2023-09-16 21:49:15 Annotating text fragment 47231/57597
## 2023-09-16 21:49:15 Annotating text fragment 47241/57597
## 2023-09-16 21:49:15 Annotating text fragment 47251/57597
## 2023-09-16 21:49:15 Annotating text fragment 47261/57597
## 2023-09-16 21:49:16 Annotating text fragment 47271/57597
## 2023-09-16 21:49:16 Annotating text fragment 47281/57597
## 2023-09-16 21:49:16 Annotating text fragment 47291/57597
## 2023-09-16 21:49:16 Annotating text fragment 47301/57597
## 2023-09-16 21:49:16 Annotating text fragment 47311/57597
## 2023-09-16 21:49:16 Annotating text fragment 47321/57597
## 2023-09-16 21:49:16 Annotating text fragment 47331/57597
## 2023-09-16 21:49:16 Annotating text fragment 47341/57597
## 2023-09-16 21:49:16 Annotating text fragment 47351/57597
## 2023-09-16 21:49:16 Annotating text fragment 47361/57597
## 2023-09-16 21:49:16 Annotating text fragment 47371/57597
## 2023-09-16 21:49:16 Annotating text fragment 47381/57597
## 2023-09-16 21:49:17 Annotating text fragment 47391/57597
## 2023-09-16 21:49:17 Annotating text fragment 47401/57597
## 2023-09-16 21:49:17 Annotating text fragment 47411/57597
## 2023-09-16 21:49:17 Annotating text fragment 47421/57597
## 2023-09-16 21:49:17 Annotating text fragment 47431/57597
## 2023-09-16 21:49:17 Annotating text fragment 47441/57597
## 2023-09-16 21:49:17 Annotating text fragment 47451/57597
## 2023-09-16 21:49:17 Annotating text fragment 47461/57597
## 2023-09-16 21:49:17 Annotating text fragment 47471/57597
## 2023-09-16 21:49:17 Annotating text fragment 47481/57597
## 2023-09-16 21:49:17 Annotating text fragment 47491/57597
## 2023-09-16 21:49:18 Annotating text fragment 47501/57597
## 2023-09-16 21:49:18 Annotating text fragment 47511/57597
## 2023-09-16 21:49:18 Annotating text fragment 47521/57597
## 2023-09-16 21:49:18 Annotating text fragment 47531/57597
## 2023-09-16 21:49:18 Annotating text fragment 47541/57597
## 2023-09-16 21:49:18 Annotating text fragment 47551/57597
## 2023-09-16 21:49:18 Annotating text fragment 47561/57597
## 2023-09-16 21:49:18 Annotating text fragment 47571/57597
## 2023-09-16 21:49:18 Annotating text fragment 47581/57597
## 2023-09-16 21:49:19 Annotating text fragment 47591/57597
## 2023-09-16 21:49:19 Annotating text fragment 47601/57597
## 2023-09-16 21:49:19 Annotating text fragment 47611/57597
## 2023-09-16 21:49:19 Annotating text fragment 47621/57597
## 2023-09-16 21:49:19 Annotating text fragment 47631/57597
## 2023-09-16 21:49:19 Annotating text fragment 47641/57597
## 2023-09-16 21:49:19 Annotating text fragment 47651/57597
## 2023-09-16 21:49:20 Annotating text fragment 47661/57597
## 2023-09-16 21:49:20 Annotating text fragment 47671/57597
## 2023-09-16 21:49:20 Annotating text fragment 47681/57597
## 2023-09-16 21:49:20 Annotating text fragment 47691/57597
## 2023-09-16 21:49:20 Annotating text fragment 47701/57597
## 2023-09-16 21:49:20 Annotating text fragment 47711/57597
## 2023-09-16 21:49:20 Annotating text fragment 47721/57597
## 2023-09-16 21:49:20 Annotating text fragment 47731/57597
## 2023-09-16 21:49:20 Annotating text fragment 47741/57597
## 2023-09-16 21:49:20 Annotating text fragment 47751/57597
```

```
## 2023-09-16 21:49:21 Annotating text fragment 47761/57597
## 2023-09-16 21:49:21 Annotating text fragment 47771/57597
## 2023-09-16 21:49:21 Annotating text fragment 47781/57597
## 2023-09-16 21:49:21 Annotating text fragment 47791/57597
## 2023-09-16 21:49:21 Annotating text fragment 47801/57597
## 2023-09-16 21:49:21 Annotating text fragment 47811/57597
## 2023-09-16 21:49:21 Annotating text fragment 47821/57597
## 2023-09-16 21:49:21 Annotating text fragment 47831/57597
## 2023-09-16 21:49:22 Annotating text fragment 47841/57597
## 2023-09-16 21:49:22 Annotating text fragment 47851/57597
## 2023-09-16 21:49:22 Annotating text fragment 47861/57597
## 2023-09-16 21:49:22 Annotating text fragment 47871/57597
## 2023-09-16 21:49:22 Annotating text fragment 47881/57597
## 2023-09-16 21:49:22 Annotating text fragment 47891/57597
## 2023-09-16 21:49:22 Annotating text fragment 47901/57597
## 2023-09-16 21:49:22 Annotating text fragment 47911/57597
## 2023-09-16 21:49:22 Annotating text fragment 47921/57597
## 2023-09-16 21:49:22 Annotating text fragment 47931/57597
## 2023-09-16 21:49:23 Annotating text fragment 47941/57597
## 2023-09-16 21:49:23 Annotating text fragment 47951/57597
## 2023-09-16 21:49:23 Annotating text fragment 47961/57597
## 2023-09-16 21:49:23 Annotating text fragment 47971/57597
## 2023-09-16 21:49:23 Annotating text fragment 47981/57597
## 2023-09-16 21:49:23 Annotating text fragment 47991/57597
## 2023-09-16 21:49:23 Annotating text fragment 48001/57597
## 2023-09-16 21:49:23 Annotating text fragment 48011/57597
## 2023-09-16 21:49:23 Annotating text fragment 48021/57597
## 2023-09-16 21:49:23 Annotating text fragment 48031/57597
## 2023-09-16 21:49:24 Annotating text fragment 48041/57597
## 2023-09-16 21:49:24 Annotating text fragment 48051/57597
## 2023-09-16 21:49:24 Annotating text fragment 48061/57597
## 2023-09-16 21:49:24 Annotating text fragment 48071/57597
## 2023-09-16 21:49:24 Annotating text fragment 48081/57597
## 2023-09-16 21:49:24 Annotating text fragment 48091/57597
## 2023-09-16 21:49:24 Annotating text fragment 48101/57597
## 2023-09-16 21:49:25 Annotating text fragment 48111/57597
## 2023-09-16 21:49:25 Annotating text fragment 48121/57597
## 2023-09-16 21:49:25 Annotating text fragment 48131/57597
## 2023-09-16 21:49:25 Annotating text fragment 48141/57597
## 2023-09-16 21:49:25 Annotating text fragment 48151/57597
## 2023-09-16 21:49:25 Annotating text fragment 48161/57597
## 2023-09-16 21:49:25 Annotating text fragment 48171/57597
## 2023-09-16 21:49:25 Annotating text fragment 48181/57597
## 2023-09-16 21:49:25 Annotating text fragment 48191/57597
## 2023-09-16 21:49:25 Annotating text fragment 48201/57597
## 2023-09-16 21:49:25 Annotating text fragment 48211/57597
## 2023-09-16 21:49:26 Annotating text fragment 48221/57597
## 2023-09-16 21:49:26 Annotating text fragment 48231/57597
## 2023-09-16 21:49:26 Annotating text fragment 48241/57597
## 2023-09-16 21:49:26 Annotating text fragment 48251/57597
## 2023-09-16 21:49:26 Annotating text fragment 48261/57597
## 2023-09-16 21:49:26 Annotating text fragment 48271/57597
## 2023-09-16 21:49:26 Annotating text fragment 48281/57597
## 2023-09-16 21:49:26 Annotating text fragment 48291/57597
```

```
## 2023-09-16 21:49:26 Annotating text fragment 48301/57597
## 2023-09-16 21:49:26 Annotating text fragment 48311/57597
## 2023-09-16 21:49:27 Annotating text fragment 48321/57597
## 2023-09-16 21:49:27 Annotating text fragment 48331/57597
## 2023-09-16 21:49:27 Annotating text fragment 48341/57597
## 2023-09-16 21:49:27 Annotating text fragment 48351/57597
## 2023-09-16 21:49:27 Annotating text fragment 48361/57597
## 2023-09-16 21:49:27 Annotating text fragment 48371/57597
## 2023-09-16 21:49:27 Annotating text fragment 48381/57597
## 2023-09-16 21:49:28 Annotating text fragment 48391/57597
## 2023-09-16 21:49:28 Annotating text fragment 48401/57597
## 2023-09-16 21:49:28 Annotating text fragment 48411/57597
## 2023-09-16 21:49:28 Annotating text fragment 48421/57597
## 2023-09-16 21:49:28 Annotating text fragment 48431/57597
## 2023-09-16 21:49:28 Annotating text fragment 48441/57597
## 2023-09-16 21:49:28 Annotating text fragment 48451/57597
## 2023-09-16 21:49:28 Annotating text fragment 48461/57597
## 2023-09-16 21:49:28 Annotating text fragment 48471/57597
## 2023-09-16 21:49:29 Annotating text fragment 48481/57597
## 2023-09-16 21:49:29 Annotating text fragment 48491/57597
## 2023-09-16 21:49:29 Annotating text fragment 48501/57597
## 2023-09-16 21:49:29 Annotating text fragment 48511/57597
## 2023-09-16 21:49:29 Annotating text fragment 48521/57597
## 2023-09-16 21:49:29 Annotating text fragment 48531/57597
## 2023-09-16 21:49:29 Annotating text fragment 48541/57597
## 2023-09-16 21:49:30 Annotating text fragment 48551/57597
## 2023-09-16 21:49:30 Annotating text fragment 48561/57597
## 2023-09-16 21:49:30 Annotating text fragment 48571/57597
## 2023-09-16 21:49:30 Annotating text fragment 48581/57597
## 2023-09-16 21:49:31 Annotating text fragment 48591/57597
## 2023-09-16 21:49:31 Annotating text fragment 48601/57597
## 2023-09-16 21:49:31 Annotating text fragment 48611/57597
## 2023-09-16 21:49:31 Annotating text fragment 48621/57597
## 2023-09-16 21:49:31 Annotating text fragment 48631/57597
## 2023-09-16 21:49:31 Annotating text fragment 48641/57597
## 2023-09-16 21:49:32 Annotating text fragment 48651/57597
## 2023-09-16 21:49:32 Annotating text fragment 48661/57597
## 2023-09-16 21:49:32 Annotating text fragment 48671/57597
## 2023-09-16 21:49:32 Annotating text fragment 48681/57597
## 2023-09-16 21:49:32 Annotating text fragment 48691/57597
## 2023-09-16 21:49:32 Annotating text fragment 48701/57597
## 2023-09-16 21:49:32 Annotating text fragment 48711/57597
## 2023-09-16 21:49:33 Annotating text fragment 48721/57597
## 2023-09-16 21:49:33 Annotating text fragment 48731/57597
## 2023-09-16 21:49:33 Annotating text fragment 48741/57597
## 2023-09-16 21:49:33 Annotating text fragment 48751/57597
## 2023-09-16 21:49:33 Annotating text fragment 48761/57597
## 2023-09-16 21:49:33 Annotating text fragment 48771/57597
## 2023-09-16 21:49:34 Annotating text fragment 48781/57597
## 2023-09-16 21:49:34 Annotating text fragment 48791/57597
## 2023-09-16 21:49:34 Annotating text fragment 48801/57597
## 2023-09-16 21:49:34 Annotating text fragment 48811/57597
## 2023-09-16 21:49:34 Annotating text fragment 48821/57597
## 2023-09-16 21:49:34 Annotating text fragment 48831/57597
```

```
## 2023-09-16 21:49:34 Annotating text fragment 48841/57597
## 2023-09-16 21:49:34 Annotating text fragment 48851/57597
## 2023-09-16 21:49:35 Annotating text fragment 48861/57597
## 2023-09-16 21:49:35 Annotating text fragment 48871/57597
## 2023-09-16 21:49:35 Annotating text fragment 48881/57597
## 2023-09-16 21:49:35 Annotating text fragment 48891/57597
## 2023-09-16 21:49:35 Annotating text fragment 48901/57597
## 2023-09-16 21:49:35 Annotating text fragment 48911/57597
## 2023-09-16 21:49:35 Annotating text fragment 48921/57597
## 2023-09-16 21:49:35 Annotating text fragment 48931/57597
## 2023-09-16 21:49:35 Annotating text fragment 48941/57597
## 2023-09-16 21:49:35 Annotating text fragment 48951/57597
## 2023-09-16 21:49:35 Annotating text fragment 48961/57597
## 2023-09-16 21:49:35 Annotating text fragment 48971/57597
## 2023-09-16 21:49:36 Annotating text fragment 48981/57597
## 2023-09-16 21:49:36 Annotating text fragment 48991/57597
## 2023-09-16 21:49:36 Annotating text fragment 49001/57597
## 2023-09-16 21:49:36 Annotating text fragment 49011/57597
## 2023-09-16 21:49:36 Annotating text fragment 49021/57597
## 2023-09-16 21:49:36 Annotating text fragment 49031/57597
## 2023-09-16 21:49:36 Annotating text fragment 49041/57597
## 2023-09-16 21:49:37 Annotating text fragment 49051/57597
## 2023-09-16 21:49:37 Annotating text fragment 49061/57597
## 2023-09-16 21:49:37 Annotating text fragment 49071/57597
## 2023-09-16 21:49:37 Annotating text fragment 49081/57597
## 2023-09-16 21:49:37 Annotating text fragment 49091/57597
## 2023-09-16 21:49:37 Annotating text fragment 49101/57597
## 2023-09-16 21:49:37 Annotating text fragment 49111/57597
## 2023-09-16 21:49:37 Annotating text fragment 49121/57597
## 2023-09-16 21:49:37 Annotating text fragment 49131/57597
## 2023-09-16 21:49:37 Annotating text fragment 49141/57597
## 2023-09-16 21:49:38 Annotating text fragment 49151/57597
## 2023-09-16 21:49:38 Annotating text fragment 49161/57597
## 2023-09-16 21:49:38 Annotating text fragment 49171/57597
## 2023-09-16 21:49:38 Annotating text fragment 49181/57597
## 2023-09-16 21:49:38 Annotating text fragment 49191/57597
## 2023-09-16 21:49:38 Annotating text fragment 49201/57597
## 2023-09-16 21:49:39 Annotating text fragment 49211/57597
## 2023-09-16 21:49:39 Annotating text fragment 49221/57597
## 2023-09-16 21:49:39 Annotating text fragment 49231/57597
## 2023-09-16 21:49:39 Annotating text fragment 49241/57597
## 2023-09-16 21:49:39 Annotating text fragment 49251/57597
## 2023-09-16 21:49:39 Annotating text fragment 49261/57597
## 2023-09-16 21:49:39 Annotating text fragment 49271/57597
## 2023-09-16 21:49:39 Annotating text fragment 49281/57597
## 2023-09-16 21:49:39 Annotating text fragment 49291/57597
## 2023-09-16 21:49:39 Annotating text fragment 49301/57597
## 2023-09-16 21:49:40 Annotating text fragment 49311/57597
## 2023-09-16 21:49:40 Annotating text fragment 49321/57597
## 2023-09-16 21:49:40 Annotating text fragment 49331/57597
## 2023-09-16 21:49:40 Annotating text fragment 49341/57597
## 2023-09-16 21:49:40 Annotating text fragment 49351/57597
## 2023-09-16 21:49:40 Annotating text fragment 49361/57597
## 2023-09-16 21:49:40 Annotating text fragment 49371/57597
```

```
## 2023-09-16 21:49:40 Annotating text fragment 49381/57597
## 2023-09-16 21:49:40 Annotating text fragment 49391/57597
## 2023-09-16 21:49:40 Annotating text fragment 49401/57597
## 2023-09-16 21:49:41 Annotating text fragment 49411/57597
## 2023-09-16 21:49:41 Annotating text fragment 49421/57597
## 2023-09-16 21:49:41 Annotating text fragment 49431/57597
## 2023-09-16 21:49:41 Annotating text fragment 49441/57597
## 2023-09-16 21:49:41 Annotating text fragment 49451/57597
## 2023-09-16 21:49:41 Annotating text fragment 49461/57597
## 2023-09-16 21:49:41 Annotating text fragment 49471/57597
## 2023-09-16 21:49:42 Annotating text fragment 49481/57597
## 2023-09-16 21:49:42 Annotating text fragment 49491/57597
## 2023-09-16 21:49:42 Annotating text fragment 49501/57597
## 2023-09-16 21:49:42 Annotating text fragment 49511/57597
## 2023-09-16 21:49:42 Annotating text fragment 49521/57597
## 2023-09-16 21:49:42 Annotating text fragment 49531/57597
## 2023-09-16 21:49:42 Annotating text fragment 49541/57597
## 2023-09-16 21:49:42 Annotating text fragment 49551/57597
## 2023-09-16 21:49:43 Annotating text fragment 49561/57597
## 2023-09-16 21:49:43 Annotating text fragment 49571/57597
## 2023-09-16 21:49:43 Annotating text fragment 49581/57597
## 2023-09-16 21:49:43 Annotating text fragment 49591/57597
## 2023-09-16 21:49:43 Annotating text fragment 49601/57597
## 2023-09-16 21:49:43 Annotating text fragment 49611/57597
## 2023-09-16 21:49:44 Annotating text fragment 49621/57597
## 2023-09-16 21:49:44 Annotating text fragment 49631/57597
## 2023-09-16 21:49:44 Annotating text fragment 49641/57597
## 2023-09-16 21:49:44 Annotating text fragment 49651/57597
## 2023-09-16 21:49:44 Annotating text fragment 49661/57597
## 2023-09-16 21:49:44 Annotating text fragment 49671/57597
## 2023-09-16 21:49:44 Annotating text fragment 49681/57597
## 2023-09-16 21:49:45 Annotating text fragment 49691/57597
## 2023-09-16 21:49:45 Annotating text fragment 49701/57597
## 2023-09-16 21:49:45 Annotating text fragment 49711/57597
## 2023-09-16 21:49:45 Annotating text fragment 49721/57597
## 2023-09-16 21:49:45 Annotating text fragment 49731/57597
## 2023-09-16 21:49:46 Annotating text fragment 49741/57597
## 2023-09-16 21:49:46 Annotating text fragment 49751/57597
## 2023-09-16 21:49:46 Annotating text fragment 49761/57597
## 2023-09-16 21:49:46 Annotating text fragment 49771/57597
## 2023-09-16 21:49:46 Annotating text fragment 49781/57597
## 2023-09-16 21:49:46 Annotating text fragment 49791/57597
## 2023-09-16 21:49:46 Annotating text fragment 49801/57597
## 2023-09-16 21:49:46 Annotating text fragment 49811/57597
## 2023-09-16 21:49:46 Annotating text fragment 49821/57597
## 2023-09-16 21:49:47 Annotating text fragment 49831/57597
## 2023-09-16 21:49:47 Annotating text fragment 49841/57597
## 2023-09-16 21:49:47 Annotating text fragment 49851/57597
## 2023-09-16 21:49:47 Annotating text fragment 49861/57597
## 2023-09-16 21:49:47 Annotating text fragment 49871/57597
## 2023-09-16 21:49:47 Annotating text fragment 49881/57597
## 2023-09-16 21:49:47 Annotating text fragment 49891/57597
## 2023-09-16 21:49:47 Annotating text fragment 49901/57597
## 2023-09-16 21:49:47 Annotating text fragment 49911/57597
```

```
## 2023-09-16 21:49:47 Annotating text fragment 49921/57597
## 2023-09-16 21:49:47 Annotating text fragment 49931/57597
## 2023-09-16 21:49:48 Annotating text fragment 49941/57597
## 2023-09-16 21:49:48 Annotating text fragment 49951/57597
## 2023-09-16 21:49:48 Annotating text fragment 49961/57597
## 2023-09-16 21:49:48 Annotating text fragment 49971/57597
## 2023-09-16 21:49:48 Annotating text fragment 49981/57597
## 2023-09-16 21:49:48 Annotating text fragment 49991/57597
## 2023-09-16 21:49:49 Annotating text fragment 50001/57597
## 2023-09-16 21:49:49 Annotating text fragment 50011/57597
## 2023-09-16 21:49:49 Annotating text fragment 50021/57597
## 2023-09-16 21:49:49 Annotating text fragment 50031/57597
## 2023-09-16 21:49:49 Annotating text fragment 50041/57597
## 2023-09-16 21:49:49 Annotating text fragment 50051/57597
## 2023-09-16 21:49:49 Annotating text fragment 50061/57597
## 2023-09-16 21:49:49 Annotating text fragment 50071/57597
## 2023-09-16 21:49:49 Annotating text fragment 50081/57597
## 2023-09-16 21:49:50 Annotating text fragment 50091/57597
## 2023-09-16 21:49:50 Annotating text fragment 50101/57597
## 2023-09-16 21:49:50 Annotating text fragment 50111/57597
## 2023-09-16 21:49:50 Annotating text fragment 50121/57597
## 2023-09-16 21:49:50 Annotating text fragment 50131/57597
## 2023-09-16 21:49:50 Annotating text fragment 50141/57597
## 2023-09-16 21:49:50 Annotating text fragment 50151/57597
## 2023-09-16 21:49:50 Annotating text fragment 50161/57597
## 2023-09-16 21:49:51 Annotating text fragment 50171/57597
## 2023-09-16 21:49:51 Annotating text fragment 50181/57597
## 2023-09-16 21:49:51 Annotating text fragment 50191/57597
## 2023-09-16 21:49:51 Annotating text fragment 50201/57597
## 2023-09-16 21:49:51 Annotating text fragment 50211/57597
## 2023-09-16 21:49:51 Annotating text fragment 50221/57597
## 2023-09-16 21:49:52 Annotating text fragment 50231/57597
## 2023-09-16 21:49:52 Annotating text fragment 50241/57597
## 2023-09-16 21:49:52 Annotating text fragment 50251/57597
## 2023-09-16 21:49:52 Annotating text fragment 50261/57597
## 2023-09-16 21:49:52 Annotating text fragment 50271/57597
## 2023-09-16 21:49:52 Annotating text fragment 50281/57597
## 2023-09-16 21:49:53 Annotating text fragment 50291/57597
## 2023-09-16 21:49:53 Annotating text fragment 50301/57597
## 2023-09-16 21:49:53 Annotating text fragment 50311/57597
## 2023-09-16 21:49:53 Annotating text fragment 50321/57597
## 2023-09-16 21:49:53 Annotating text fragment 50331/57597
## 2023-09-16 21:49:53 Annotating text fragment 50341/57597
## 2023-09-16 21:49:53 Annotating text fragment 50351/57597
## 2023-09-16 21:49:53 Annotating text fragment 50361/57597
## 2023-09-16 21:49:53 Annotating text fragment 50371/57597
## 2023-09-16 21:49:54 Annotating text fragment 50381/57597
## 2023-09-16 21:49:54 Annotating text fragment 50391/57597
## 2023-09-16 21:49:54 Annotating text fragment 50401/57597
## 2023-09-16 21:49:54 Annotating text fragment 50411/57597
## 2023-09-16 21:49:54 Annotating text fragment 50421/57597
## 2023-09-16 21:49:54 Annotating text fragment 50431/57597
## 2023-09-16 21:49:54 Annotating text fragment 50441/57597
## 2023-09-16 21:49:54 Annotating text fragment 50451/57597
```

```
## 2023-09-16 21:49:54 Annotating text fragment 50461/57597
## 2023-09-16 21:49:54 Annotating text fragment 50471/57597
## 2023-09-16 21:49:55 Annotating text fragment 50481/57597
## 2023-09-16 21:49:55 Annotating text fragment 50491/57597
## 2023-09-16 21:49:55 Annotating text fragment 50501/57597
## 2023-09-16 21:49:55 Annotating text fragment 50511/57597
## 2023-09-16 21:49:55 Annotating text fragment 50521/57597
## 2023-09-16 21:49:55 Annotating text fragment 50531/57597
## 2023-09-16 21:49:55 Annotating text fragment 50541/57597
## 2023-09-16 21:49:55 Annotating text fragment 50551/57597
## 2023-09-16 21:49:55 Annotating text fragment 50561/57597
## 2023-09-16 21:49:56 Annotating text fragment 50571/57597
## 2023-09-16 21:49:56 Annotating text fragment 50581/57597
## 2023-09-16 21:49:56 Annotating text fragment 50591/57597
## 2023-09-16 21:49:56 Annotating text fragment 50601/57597
## 2023-09-16 21:49:56 Annotating text fragment 50611/57597
## 2023-09-16 21:49:56 Annotating text fragment 50621/57597
## 2023-09-16 21:49:56 Annotating text fragment 50631/57597
## 2023-09-16 21:49:57 Annotating text fragment 50641/57597
## 2023-09-16 21:49:57 Annotating text fragment 50651/57597
## 2023-09-16 21:49:57 Annotating text fragment 50661/57597
## 2023-09-16 21:49:57 Annotating text fragment 50671/57597
## 2023-09-16 21:49:57 Annotating text fragment 50681/57597
## 2023-09-16 21:49:57 Annotating text fragment 50691/57597
## 2023-09-16 21:49:57 Annotating text fragment 50701/57597
## 2023-09-16 21:49:57 Annotating text fragment 50711/57597
## 2023-09-16 21:49:58 Annotating text fragment 50721/57597
## 2023-09-16 21:49:58 Annotating text fragment 50731/57597
## 2023-09-16 21:49:58 Annotating text fragment 50741/57597
## 2023-09-16 21:49:58 Annotating text fragment 50751/57597
## 2023-09-16 21:49:58 Annotating text fragment 50761/57597
## 2023-09-16 21:49:58 Annotating text fragment 50771/57597
## 2023-09-16 21:49:58 Annotating text fragment 50781/57597
## 2023-09-16 21:49:58 Annotating text fragment 50791/57597
## 2023-09-16 21:49:59 Annotating text fragment 50801/57597
## 2023-09-16 21:49:59 Annotating text fragment 50811/57597
## 2023-09-16 21:49:59 Annotating text fragment 50821/57597
## 2023-09-16 21:49:59 Annotating text fragment 50831/57597
## 2023-09-16 21:49:59 Annotating text fragment 50841/57597
## 2023-09-16 21:49:59 Annotating text fragment 50851/57597
## 2023-09-16 21:49:59 Annotating text fragment 50861/57597
## 2023-09-16 21:50:00 Annotating text fragment 50871/57597
## 2023-09-16 21:50:00 Annotating text fragment 50881/57597
## 2023-09-16 21:50:00 Annotating text fragment 50891/57597
## 2023-09-16 21:50:00 Annotating text fragment 50901/57597
## 2023-09-16 21:50:00 Annotating text fragment 50911/57597
## 2023-09-16 21:50:00 Annotating text fragment 50921/57597
## 2023-09-16 21:50:00 Annotating text fragment 50931/57597
## 2023-09-16 21:50:00 Annotating text fragment 50941/57597
## 2023-09-16 21:50:00 Annotating text fragment 50951/57597
## 2023-09-16 21:50:01 Annotating text fragment 50961/57597
## 2023-09-16 21:50:01 Annotating text fragment 50971/57597
## 2023-09-16 21:50:01 Annotating text fragment 50981/57597
## 2023-09-16 21:50:01 Annotating text fragment 50991/57597
```

```
## 2023-09-16 21:50:01 Annotating text fragment 51001/57597
## 2023-09-16 21:50:01 Annotating text fragment 51011/57597
## 2023-09-16 21:50:01 Annotating text fragment 51021/57597
## 2023-09-16 21:50:01 Annotating text fragment 51031/57597
## 2023-09-16 21:50:01 Annotating text fragment 51041/57597
## 2023-09-16 21:50:02 Annotating text fragment 51051/57597
## 2023-09-16 21:50:02 Annotating text fragment 51061/57597
## 2023-09-16 21:50:02 Annotating text fragment 51071/57597
## 2023-09-16 21:50:02 Annotating text fragment 51081/57597
## 2023-09-16 21:50:02 Annotating text fragment 51091/57597
## 2023-09-16 21:50:02 Annotating text fragment 51101/57597
## 2023-09-16 21:50:02 Annotating text fragment 51111/57597
## 2023-09-16 21:50:03 Annotating text fragment 51121/57597
## 2023-09-16 21:50:03 Annotating text fragment 51131/57597
## 2023-09-16 21:50:03 Annotating text fragment 51141/57597
## 2023-09-16 21:50:03 Annotating text fragment 51151/57597
## 2023-09-16 21:50:03 Annotating text fragment 51161/57597
## 2023-09-16 21:50:03 Annotating text fragment 51171/57597
## 2023-09-16 21:50:03 Annotating text fragment 51181/57597
## 2023-09-16 21:50:03 Annotating text fragment 51191/57597
## 2023-09-16 21:50:03 Annotating text fragment 51201/57597
## 2023-09-16 21:50:03 Annotating text fragment 51211/57597
## 2023-09-16 21:50:04 Annotating text fragment 51221/57597
## 2023-09-16 21:50:04 Annotating text fragment 51231/57597
## 2023-09-16 21:50:04 Annotating text fragment 51241/57597
## 2023-09-16 21:50:04 Annotating text fragment 51251/57597
## 2023-09-16 21:50:04 Annotating text fragment 51261/57597
## 2023-09-16 21:50:04 Annotating text fragment 51271/57597
## 2023-09-16 21:50:04 Annotating text fragment 51281/57597
## 2023-09-16 21:50:04 Annotating text fragment 51291/57597
## 2023-09-16 21:50:04 Annotating text fragment 51301/57597
## 2023-09-16 21:50:05 Annotating text fragment 51311/57597
## 2023-09-16 21:50:05 Annotating text fragment 51321/57597
## 2023-09-16 21:50:05 Annotating text fragment 51331/57597
## 2023-09-16 21:50:05 Annotating text fragment 51341/57597
## 2023-09-16 21:50:05 Annotating text fragment 51351/57597
## 2023-09-16 21:50:05 Annotating text fragment 51361/57597
## 2023-09-16 21:50:05 Annotating text fragment 51371/57597
## 2023-09-16 21:50:05 Annotating text fragment 51381/57597
## 2023-09-16 21:50:06 Annotating text fragment 51391/57597
## 2023-09-16 21:50:06 Annotating text fragment 51401/57597
## 2023-09-16 21:50:06 Annotating text fragment 51411/57597
## 2023-09-16 21:50:06 Annotating text fragment 51421/57597
## 2023-09-16 21:50:06 Annotating text fragment 51431/57597
## 2023-09-16 21:50:06 Annotating text fragment 51441/57597
## 2023-09-16 21:50:06 Annotating text fragment 51451/57597
## 2023-09-16 21:50:06 Annotating text fragment 51461/57597
## 2023-09-16 21:50:06 Annotating text fragment 51471/57597
## 2023-09-16 21:50:07 Annotating text fragment 51481/57597
## 2023-09-16 21:50:07 Annotating text fragment 51491/57597
## 2023-09-16 21:50:07 Annotating text fragment 51501/57597
## 2023-09-16 21:50:07 Annotating text fragment 51511/57597
## 2023-09-16 21:50:07 Annotating text fragment 51521/57597
## 2023-09-16 21:50:07 Annotating text fragment 51531/57597
```

```
## 2023-09-16 21:50:07 Annotating text fragment 51541/57597
## 2023-09-16 21:50:07 Annotating text fragment 51551/57597
## 2023-09-16 21:50:07 Annotating text fragment 51561/57597
## 2023-09-16 21:50:07 Annotating text fragment 51571/57597
## 2023-09-16 21:50:08 Annotating text fragment 51581/57597
## 2023-09-16 21:50:08 Annotating text fragment 51591/57597
## 2023-09-16 21:50:08 Annotating text fragment 51601/57597
## 2023-09-16 21:50:08 Annotating text fragment 51611/57597
## 2023-09-16 21:50:08 Annotating text fragment 51621/57597
## 2023-09-16 21:50:08 Annotating text fragment 51631/57597
## 2023-09-16 21:50:08 Annotating text fragment 51641/57597
## 2023-09-16 21:50:08 Annotating text fragment 51651/57597
## 2023-09-16 21:50:09 Annotating text fragment 51661/57597
## 2023-09-16 21:50:09 Annotating text fragment 51671/57597
## 2023-09-16 21:50:09 Annotating text fragment 51681/57597
## 2023-09-16 21:50:09 Annotating text fragment 51691/57597
## 2023-09-16 21:50:09 Annotating text fragment 51701/57597
## 2023-09-16 21:50:09 Annotating text fragment 51711/57597
## 2023-09-16 21:50:09 Annotating text fragment 51721/57597
## 2023-09-16 21:50:09 Annotating text fragment 51731/57597
## 2023-09-16 21:50:10 Annotating text fragment 51741/57597
## 2023-09-16 21:50:10 Annotating text fragment 51751/57597
## 2023-09-16 21:50:10 Annotating text fragment 51761/57597
## 2023-09-16 21:50:10 Annotating text fragment 51771/57597
## 2023-09-16 21:50:10 Annotating text fragment 51781/57597
## 2023-09-16 21:50:10 Annotating text fragment 51791/57597
## 2023-09-16 21:50:10 Annotating text fragment 51801/57597
## 2023-09-16 21:50:10 Annotating text fragment 51811/57597
## 2023-09-16 21:50:10 Annotating text fragment 51821/57597
## 2023-09-16 21:50:10 Annotating text fragment 51831/57597
## 2023-09-16 21:50:11 Annotating text fragment 51841/57597
## 2023-09-16 21:50:11 Annotating text fragment 51851/57597
## 2023-09-16 21:50:11 Annotating text fragment 51861/57597
## 2023-09-16 21:50:11 Annotating text fragment 51871/57597
## 2023-09-16 21:50:11 Annotating text fragment 51881/57597
## 2023-09-16 21:50:11 Annotating text fragment 51891/57597
## 2023-09-16 21:50:11 Annotating text fragment 51901/57597
## 2023-09-16 21:50:11 Annotating text fragment 51911/57597
## 2023-09-16 21:50:12 Annotating text fragment 51921/57597
## 2023-09-16 21:50:12 Annotating text fragment 51931/57597
## 2023-09-16 21:50:12 Annotating text fragment 51941/57597
## 2023-09-16 21:50:12 Annotating text fragment 51951/57597
## 2023-09-16 21:50:12 Annotating text fragment 51961/57597
## 2023-09-16 21:50:12 Annotating text fragment 51971/57597
## 2023-09-16 21:50:12 Annotating text fragment 51981/57597
## 2023-09-16 21:50:12 Annotating text fragment 51991/57597
## 2023-09-16 21:50:12 Annotating text fragment 52001/57597
## 2023-09-16 21:50:13 Annotating text fragment 52011/57597
## 2023-09-16 21:50:13 Annotating text fragment 52021/57597
## 2023-09-16 21:50:13 Annotating text fragment 52031/57597
## 2023-09-16 21:50:13 Annotating text fragment 52041/57597
## 2023-09-16 21:50:13 Annotating text fragment 52051/57597
## 2023-09-16 21:50:13 Annotating text fragment 52061/57597
## 2023-09-16 21:50:13 Annotating text fragment 52071/57597
```

```
## 2023-09-16 21:50:14 Annotating text fragment 52081/57597
## 2023-09-16 21:50:14 Annotating text fragment 52091/57597
## 2023-09-16 21:50:14 Annotating text fragment 52101/57597
## 2023-09-16 21:50:14 Annotating text fragment 52111/57597
## 2023-09-16 21:50:14 Annotating text fragment 52121/57597
## 2023-09-16 21:50:14 Annotating text fragment 52131/57597
## 2023-09-16 21:50:14 Annotating text fragment 52141/57597
## 2023-09-16 21:50:14 Annotating text fragment 52151/57597
## 2023-09-16 21:50:14 Annotating text fragment 52161/57597
## 2023-09-16 21:50:15 Annotating text fragment 52171/57597
## 2023-09-16 21:50:15 Annotating text fragment 52181/57597
## 2023-09-16 21:50:15 Annotating text fragment 52191/57597
## 2023-09-16 21:50:15 Annotating text fragment 52201/57597
## 2023-09-16 21:50:15 Annotating text fragment 52211/57597
## 2023-09-16 21:50:15 Annotating text fragment 52221/57597
## 2023-09-16 21:50:15 Annotating text fragment 52231/57597
## 2023-09-16 21:50:15 Annotating text fragment 52241/57597
## 2023-09-16 21:50:15 Annotating text fragment 52251/57597
## 2023-09-16 21:50:16 Annotating text fragment 52261/57597
## 2023-09-16 21:50:16 Annotating text fragment 52271/57597
## 2023-09-16 21:50:16 Annotating text fragment 52281/57597
## 2023-09-16 21:50:16 Annotating text fragment 52291/57597
## 2023-09-16 21:50:16 Annotating text fragment 52301/57597
## 2023-09-16 21:50:16 Annotating text fragment 52311/57597
## 2023-09-16 21:50:16 Annotating text fragment 52321/57597
## 2023-09-16 21:50:16 Annotating text fragment 52331/57597
## 2023-09-16 21:50:17 Annotating text fragment 52341/57597
## 2023-09-16 21:50:17 Annotating text fragment 52351/57597
## 2023-09-16 21:50:17 Annotating text fragment 52361/57597
## 2023-09-16 21:50:17 Annotating text fragment 52371/57597
## 2023-09-16 21:50:17 Annotating text fragment 52381/57597
## 2023-09-16 21:50:17 Annotating text fragment 52391/57597
## 2023-09-16 21:50:17 Annotating text fragment 52401/57597
## 2023-09-16 21:50:17 Annotating text fragment 52411/57597
## 2023-09-16 21:50:18 Annotating text fragment 52421/57597
## 2023-09-16 21:50:18 Annotating text fragment 52431/57597
## 2023-09-16 21:50:18 Annotating text fragment 52441/57597
## 2023-09-16 21:50:18 Annotating text fragment 52451/57597
## 2023-09-16 21:50:18 Annotating text fragment 52461/57597
## 2023-09-16 21:50:18 Annotating text fragment 52471/57597
## 2023-09-16 21:50:18 Annotating text fragment 52481/57597
## 2023-09-16 21:50:18 Annotating text fragment 52491/57597
## 2023-09-16 21:50:18 Annotating text fragment 52501/57597
## 2023-09-16 21:50:19 Annotating text fragment 52511/57597
## 2023-09-16 21:50:19 Annotating text fragment 52521/57597
## 2023-09-16 21:50:19 Annotating text fragment 52531/57597
## 2023-09-16 21:50:19 Annotating text fragment 52541/57597
## 2023-09-16 21:50:19 Annotating text fragment 52551/57597
## 2023-09-16 21:50:19 Annotating text fragment 52561/57597
## 2023-09-16 21:50:19 Annotating text fragment 52571/57597
## 2023-09-16 21:50:19 Annotating text fragment 52581/57597
## 2023-09-16 21:50:19 Annotating text fragment 52591/57597
## 2023-09-16 21:50:20 Annotating text fragment 52601/57597
## 2023-09-16 21:50:20 Annotating text fragment 52611/57597
```

```
## 2023-09-16 21:50:20 Annotating text fragment 52621/57597
## 2023-09-16 21:50:20 Annotating text fragment 52631/57597
## 2023-09-16 21:50:20 Annotating text fragment 52641/57597
## 2023-09-16 21:50:20 Annotating text fragment 52651/57597
## 2023-09-16 21:50:20 Annotating text fragment 52661/57597
## 2023-09-16 21:50:20 Annotating text fragment 52671/57597
## 2023-09-16 21:50:20 Annotating text fragment 52681/57597
## 2023-09-16 21:50:21 Annotating text fragment 52691/57597
## 2023-09-16 21:50:21 Annotating text fragment 52701/57597
## 2023-09-16 21:50:21 Annotating text fragment 52711/57597
## 2023-09-16 21:50:21 Annotating text fragment 52721/57597
## 2023-09-16 21:50:21 Annotating text fragment 52731/57597
## 2023-09-16 21:50:21 Annotating text fragment 52741/57597
## 2023-09-16 21:50:21 Annotating text fragment 52751/57597
## 2023-09-16 21:50:21 Annotating text fragment 52761/57597
## 2023-09-16 21:50:21 Annotating text fragment 52771/57597
## 2023-09-16 21:50:22 Annotating text fragment 52781/57597
## 2023-09-16 21:50:22 Annotating text fragment 52791/57597
## 2023-09-16 21:50:22 Annotating text fragment 52801/57597
## 2023-09-16 21:50:22 Annotating text fragment 52811/57597
## 2023-09-16 21:50:22 Annotating text fragment 52821/57597
## 2023-09-16 21:50:22 Annotating text fragment 52831/57597
## 2023-09-16 21:50:22 Annotating text fragment 52841/57597
## 2023-09-16 21:50:22 Annotating text fragment 52851/57597
## 2023-09-16 21:50:22 Annotating text fragment 52861/57597
## 2023-09-16 21:50:23 Annotating text fragment 52871/57597
## 2023-09-16 21:50:23 Annotating text fragment 52881/57597
## 2023-09-16 21:50:23 Annotating text fragment 52891/57597
## 2023-09-16 21:50:23 Annotating text fragment 52901/57597
## 2023-09-16 21:50:23 Annotating text fragment 52911/57597
## 2023-09-16 21:50:23 Annotating text fragment 52921/57597
## 2023-09-16 21:50:23 Annotating text fragment 52931/57597
## 2023-09-16 21:50:23 Annotating text fragment 52941/57597
## 2023-09-16 21:50:23 Annotating text fragment 52951/57597
## 2023-09-16 21:50:23 Annotating text fragment 52961/57597
## 2023-09-16 21:50:24 Annotating text fragment 52971/57597
## 2023-09-16 21:50:24 Annotating text fragment 52981/57597
## 2023-09-16 21:50:24 Annotating text fragment 52991/57597
## 2023-09-16 21:50:24 Annotating text fragment 53001/57597
## 2023-09-16 21:50:24 Annotating text fragment 53011/57597
## 2023-09-16 21:50:24 Annotating text fragment 53021/57597
## 2023-09-16 21:50:24 Annotating text fragment 53031/57597
## 2023-09-16 21:50:24 Annotating text fragment 53041/57597
## 2023-09-16 21:50:25 Annotating text fragment 53051/57597
## 2023-09-16 21:50:25 Annotating text fragment 53061/57597
## 2023-09-16 21:50:25 Annotating text fragment 53071/57597
## 2023-09-16 21:50:25 Annotating text fragment 53081/57597
## 2023-09-16 21:50:25 Annotating text fragment 53091/57597
## 2023-09-16 21:50:25 Annotating text fragment 53101/57597
## 2023-09-16 21:50:25 Annotating text fragment 53111/57597
## 2023-09-16 21:50:26 Annotating text fragment 53121/57597
## 2023-09-16 21:50:26 Annotating text fragment 53131/57597
## 2023-09-16 21:50:26 Annotating text fragment 53141/57597
## 2023-09-16 21:50:26 Annotating text fragment 53151/57597
```

```
## 2023-09-16 21:50:26 Annotating text fragment 53161/57597
## 2023-09-16 21:50:26 Annotating text fragment 53171/57597
## 2023-09-16 21:50:26 Annotating text fragment 53181/57597
## 2023-09-16 21:50:26 Annotating text fragment 53191/57597
## 2023-09-16 21:50:26 Annotating text fragment 53201/57597
## 2023-09-16 21:50:27 Annotating text fragment 53211/57597
## 2023-09-16 21:50:27 Annotating text fragment 53221/57597
## 2023-09-16 21:50:27 Annotating text fragment 53231/57597
## 2023-09-16 21:50:27 Annotating text fragment 53241/57597
## 2023-09-16 21:50:27 Annotating text fragment 53251/57597
## 2023-09-16 21:50:27 Annotating text fragment 53261/57597
## 2023-09-16 21:50:27 Annotating text fragment 53271/57597
## 2023-09-16 21:50:27 Annotating text fragment 53281/57597
## 2023-09-16 21:50:28 Annotating text fragment 53291/57597
## 2023-09-16 21:50:28 Annotating text fragment 53301/57597
## 2023-09-16 21:50:28 Annotating text fragment 53311/57597
## 2023-09-16 21:50:28 Annotating text fragment 53321/57597
## 2023-09-16 21:50:28 Annotating text fragment 53331/57597
## 2023-09-16 21:50:29 Annotating text fragment 53341/57597
## 2023-09-16 21:50:29 Annotating text fragment 53351/57597
## 2023-09-16 21:50:29 Annotating text fragment 53361/57597
## 2023-09-16 21:50:29 Annotating text fragment 53371/57597
## 2023-09-16 21:50:29 Annotating text fragment 53381/57597
## 2023-09-16 21:50:29 Annotating text fragment 53391/57597
## 2023-09-16 21:50:29 Annotating text fragment 53401/57597
## 2023-09-16 21:50:29 Annotating text fragment 53411/57597
## 2023-09-16 21:50:29 Annotating text fragment 53421/57597
## 2023-09-16 21:50:30 Annotating text fragment 53431/57597
## 2023-09-16 21:50:30 Annotating text fragment 53441/57597
## 2023-09-16 21:50:30 Annotating text fragment 53451/57597
## 2023-09-16 21:50:30 Annotating text fragment 53461/57597
## 2023-09-16 21:50:30 Annotating text fragment 53471/57597
## 2023-09-16 21:50:30 Annotating text fragment 53481/57597
## 2023-09-16 21:50:30 Annotating text fragment 53491/57597
## 2023-09-16 21:50:30 Annotating text fragment 53501/57597
## 2023-09-16 21:50:31 Annotating text fragment 53511/57597
## 2023-09-16 21:50:31 Annotating text fragment 53521/57597
## 2023-09-16 21:50:31 Annotating text fragment 53531/57597
## 2023-09-16 21:50:31 Annotating text fragment 53541/57597
## 2023-09-16 21:50:31 Annotating text fragment 53551/57597
## 2023-09-16 21:50:31 Annotating text fragment 53561/57597
## 2023-09-16 21:50:31 Annotating text fragment 53571/57597
## 2023-09-16 21:50:32 Annotating text fragment 53581/57597
## 2023-09-16 21:50:32 Annotating text fragment 53591/57597
## 2023-09-16 21:50:32 Annotating text fragment 53601/57597
## 2023-09-16 21:50:32 Annotating text fragment 53611/57597
## 2023-09-16 21:50:32 Annotating text fragment 53621/57597
## 2023-09-16 21:50:32 Annotating text fragment 53631/57597
## 2023-09-16 21:50:33 Annotating text fragment 53641/57597
## 2023-09-16 21:50:33 Annotating text fragment 53651/57597
## 2023-09-16 21:50:33 Annotating text fragment 53661/57597
## 2023-09-16 21:50:33 Annotating text fragment 53671/57597
## 2023-09-16 21:50:33 Annotating text fragment 53681/57597
## 2023-09-16 21:50:33 Annotating text fragment 53691/57597
```

```
## 2023-09-16 21:50:33 Annotating text fragment 53701/57597
## 2023-09-16 21:50:33 Annotating text fragment 53711/57597
## 2023-09-16 21:50:33 Annotating text fragment 53721/57597
## 2023-09-16 21:50:34 Annotating text fragment 53731/57597
## 2023-09-16 21:50:34 Annotating text fragment 53741/57597
## 2023-09-16 21:50:34 Annotating text fragment 53751/57597
## 2023-09-16 21:50:34 Annotating text fragment 53761/57597
## 2023-09-16 21:50:34 Annotating text fragment 53771/57597
## 2023-09-16 21:50:34 Annotating text fragment 53781/57597
## 2023-09-16 21:50:35 Annotating text fragment 53791/57597
## 2023-09-16 21:50:35 Annotating text fragment 53801/57597
## 2023-09-16 21:50:35 Annotating text fragment 53811/57597
## 2023-09-16 21:50:35 Annotating text fragment 53821/57597
## 2023-09-16 21:50:35 Annotating text fragment 53831/57597
## 2023-09-16 21:50:35 Annotating text fragment 53841/57597
## 2023-09-16 21:50:35 Annotating text fragment 53851/57597
## 2023-09-16 21:50:35 Annotating text fragment 53861/57597
## 2023-09-16 21:50:36 Annotating text fragment 53871/57597
## 2023-09-16 21:50:36 Annotating text fragment 53881/57597
## 2023-09-16 21:50:36 Annotating text fragment 53891/57597
## 2023-09-16 21:50:36 Annotating text fragment 53901/57597
## 2023-09-16 21:50:36 Annotating text fragment 53911/57597
## 2023-09-16 21:50:36 Annotating text fragment 53921/57597
## 2023-09-16 21:50:36 Annotating text fragment 53931/57597
## 2023-09-16 21:50:36 Annotating text fragment 53941/57597
## 2023-09-16 21:50:36 Annotating text fragment 53951/57597
## 2023-09-16 21:50:36 Annotating text fragment 53961/57597
## 2023-09-16 21:50:37 Annotating text fragment 53971/57597
## 2023-09-16 21:50:37 Annotating text fragment 53981/57597
## 2023-09-16 21:50:37 Annotating text fragment 53991/57597
## 2023-09-16 21:50:37 Annotating text fragment 54001/57597
## 2023-09-16 21:50:37 Annotating text fragment 54011/57597
## 2023-09-16 21:50:37 Annotating text fragment 54021/57597
## 2023-09-16 21:50:37 Annotating text fragment 54031/57597
## 2023-09-16 21:50:37 Annotating text fragment 54041/57597
## 2023-09-16 21:50:37 Annotating text fragment 54051/57597
## 2023-09-16 21:50:37 Annotating text fragment 54061/57597
## 2023-09-16 21:50:37 Annotating text fragment 54071/57597
## 2023-09-16 21:50:38 Annotating text fragment 54081/57597
## 2023-09-16 21:50:38 Annotating text fragment 54091/57597
## 2023-09-16 21:50:38 Annotating text fragment 54101/57597
## 2023-09-16 21:50:38 Annotating text fragment 54111/57597
## 2023-09-16 21:50:38 Annotating text fragment 54121/57597
## 2023-09-16 21:50:38 Annotating text fragment 54131/57597
## 2023-09-16 21:50:38 Annotating text fragment 54141/57597
## 2023-09-16 21:50:38 Annotating text fragment 54151/57597
## 2023-09-16 21:50:38 Annotating text fragment 54161/57597
## 2023-09-16 21:50:39 Annotating text fragment 54171/57597
## 2023-09-16 21:50:39 Annotating text fragment 54181/57597
## 2023-09-16 21:50:39 Annotating text fragment 54191/57597
## 2023-09-16 21:50:39 Annotating text fragment 54201/57597
## 2023-09-16 21:50:39 Annotating text fragment 54211/57597
## 2023-09-16 21:50:39 Annotating text fragment 54221/57597
## 2023-09-16 21:50:39 Annotating text fragment 54231/57597
```

```
## 2023-09-16 21:50:39 Annotating text fragment 54241/57597
## 2023-09-16 21:50:40 Annotating text fragment 54251/57597
## 2023-09-16 21:50:40 Annotating text fragment 54261/57597
## 2023-09-16 21:50:40 Annotating text fragment 54271/57597
## 2023-09-16 21:50:40 Annotating text fragment 54281/57597
## 2023-09-16 21:50:40 Annotating text fragment 54291/57597
## 2023-09-16 21:50:40 Annotating text fragment 54301/57597
## 2023-09-16 21:50:40 Annotating text fragment 54311/57597
## 2023-09-16 21:50:40 Annotating text fragment 54321/57597
## 2023-09-16 21:50:40 Annotating text fragment 54331/57597
## 2023-09-16 21:50:40 Annotating text fragment 54341/57597
## 2023-09-16 21:50:41 Annotating text fragment 54351/57597
## 2023-09-16 21:50:41 Annotating text fragment 54361/57597
## 2023-09-16 21:50:41 Annotating text fragment 54371/57597
## 2023-09-16 21:50:41 Annotating text fragment 54381/57597
## 2023-09-16 21:50:41 Annotating text fragment 54391/57597
## 2023-09-16 21:50:41 Annotating text fragment 54401/57597
## 2023-09-16 21:50:41 Annotating text fragment 54411/57597
## 2023-09-16 21:50:41 Annotating text fragment 54421/57597
## 2023-09-16 21:50:42 Annotating text fragment 54431/57597
## 2023-09-16 21:50:42 Annotating text fragment 54441/57597
## 2023-09-16 21:50:42 Annotating text fragment 54451/57597
## 2023-09-16 21:50:42 Annotating text fragment 54461/57597
## 2023-09-16 21:50:42 Annotating text fragment 54471/57597
## 2023-09-16 21:50:42 Annotating text fragment 54481/57597
## 2023-09-16 21:50:42 Annotating text fragment 54491/57597
## 2023-09-16 21:50:42 Annotating text fragment 54501/57597
## 2023-09-16 21:50:42 Annotating text fragment 54511/57597
## 2023-09-16 21:50:43 Annotating text fragment 54521/57597
## 2023-09-16 21:50:43 Annotating text fragment 54531/57597
## 2023-09-16 21:50:43 Annotating text fragment 54541/57597
## 2023-09-16 21:50:43 Annotating text fragment 54551/57597
## 2023-09-16 21:50:43 Annotating text fragment 54561/57597
## 2023-09-16 21:50:43 Annotating text fragment 54571/57597
## 2023-09-16 21:50:43 Annotating text fragment 54581/57597
## 2023-09-16 21:50:43 Annotating text fragment 54591/57597
## 2023-09-16 21:50:43 Annotating text fragment 54601/57597
## 2023-09-16 21:50:44 Annotating text fragment 54611/57597
## 2023-09-16 21:50:44 Annotating text fragment 54621/57597
## 2023-09-16 21:50:44 Annotating text fragment 54631/57597
## 2023-09-16 21:50:44 Annotating text fragment 54641/57597
## 2023-09-16 21:50:44 Annotating text fragment 54651/57597
## 2023-09-16 21:50:44 Annotating text fragment 54661/57597
## 2023-09-16 21:50:44 Annotating text fragment 54671/57597
## 2023-09-16 21:50:44 Annotating text fragment 54681/57597
## 2023-09-16 21:50:45 Annotating text fragment 54691/57597
## 2023-09-16 21:50:45 Annotating text fragment 54701/57597
## 2023-09-16 21:50:45 Annotating text fragment 54711/57597
## 2023-09-16 21:50:45 Annotating text fragment 54721/57597
## 2023-09-16 21:50:45 Annotating text fragment 54731/57597
## 2023-09-16 21:50:45 Annotating text fragment 54741/57597
## 2023-09-16 21:50:45 Annotating text fragment 54751/57597
## 2023-09-16 21:50:45 Annotating text fragment 54761/57597
## 2023-09-16 21:50:46 Annotating text fragment 54771/57597
```

```
## 2023-09-16 21:50:46 Annotating text fragment 54781/57597
## 2023-09-16 21:50:46 Annotating text fragment 54791/57597
## 2023-09-16 21:50:46 Annotating text fragment 54801/57597
## 2023-09-16 21:50:46 Annotating text fragment 54811/57597
## 2023-09-16 21:50:46 Annotating text fragment 54821/57597
## 2023-09-16 21:50:46 Annotating text fragment 54831/57597
## 2023-09-16 21:50:47 Annotating text fragment 54841/57597
## 2023-09-16 21:50:47 Annotating text fragment 54851/57597
## 2023-09-16 21:50:47 Annotating text fragment 54861/57597
## 2023-09-16 21:50:47 Annotating text fragment 54871/57597
## 2023-09-16 21:50:47 Annotating text fragment 54881/57597
## 2023-09-16 21:50:47 Annotating text fragment 54891/57597
## 2023-09-16 21:50:47 Annotating text fragment 54901/57597
## 2023-09-16 21:50:47 Annotating text fragment 54911/57597
## 2023-09-16 21:50:47 Annotating text fragment 54921/57597
## 2023-09-16 21:50:48 Annotating text fragment 54931/57597
## 2023-09-16 21:50:48 Annotating text fragment 54941/57597
## 2023-09-16 21:50:48 Annotating text fragment 54951/57597
## 2023-09-16 21:50:48 Annotating text fragment 54961/57597
## 2023-09-16 21:50:48 Annotating text fragment 54971/57597
## 2023-09-16 21:50:48 Annotating text fragment 54981/57597
## 2023-09-16 21:50:48 Annotating text fragment 54991/57597
## 2023-09-16 21:50:48 Annotating text fragment 55001/57597
## 2023-09-16 21:50:48 Annotating text fragment 55011/57597
## 2023-09-16 21:50:49 Annotating text fragment 55021/57597
## 2023-09-16 21:50:49 Annotating text fragment 55031/57597
## 2023-09-16 21:50:49 Annotating text fragment 55041/57597
## 2023-09-16 21:50:49 Annotating text fragment 55051/57597
## 2023-09-16 21:50:49 Annotating text fragment 55061/57597
## 2023-09-16 21:50:49 Annotating text fragment 55071/57597
## 2023-09-16 21:50:49 Annotating text fragment 55081/57597
## 2023-09-16 21:50:49 Annotating text fragment 55091/57597
## 2023-09-16 21:50:49 Annotating text fragment 55101/57597
## 2023-09-16 21:50:50 Annotating text fragment 55111/57597
## 2023-09-16 21:50:50 Annotating text fragment 55121/57597
## 2023-09-16 21:50:50 Annotating text fragment 55131/57597
## 2023-09-16 21:50:50 Annotating text fragment 55141/57597
## 2023-09-16 21:50:50 Annotating text fragment 55151/57597
## 2023-09-16 21:50:50 Annotating text fragment 55161/57597
## 2023-09-16 21:50:50 Annotating text fragment 55171/57597
## 2023-09-16 21:50:50 Annotating text fragment 55181/57597
## 2023-09-16 21:50:50 Annotating text fragment 55191/57597
## 2023-09-16 21:50:50 Annotating text fragment 55201/57597
## 2023-09-16 21:50:51 Annotating text fragment 55211/57597
## 2023-09-16 21:50:51 Annotating text fragment 55221/57597
## 2023-09-16 21:50:51 Annotating text fragment 55231/57597
## 2023-09-16 21:50:51 Annotating text fragment 55241/57597
## 2023-09-16 21:50:51 Annotating text fragment 55251/57597
## 2023-09-16 21:50:51 Annotating text fragment 55261/57597
## 2023-09-16 21:50:51 Annotating text fragment 55271/57597
## 2023-09-16 21:50:51 Annotating text fragment 55281/57597
## 2023-09-16 21:50:51 Annotating text fragment 55291/57597
## 2023-09-16 21:50:52 Annotating text fragment 55301/57597
## 2023-09-16 21:50:52 Annotating text fragment 55311/57597
```

```
## 2023-09-16 21:50:52 Annotating text fragment 55321/57597
## 2023-09-16 21:50:52 Annotating text fragment 55331/57597
## 2023-09-16 21:50:52 Annotating text fragment 55341/57597
## 2023-09-16 21:50:52 Annotating text fragment 55351/57597
## 2023-09-16 21:50:52 Annotating text fragment 55361/57597
## 2023-09-16 21:50:53 Annotating text fragment 55371/57597
## 2023-09-16 21:50:53 Annotating text fragment 55381/57597
## 2023-09-16 21:50:53 Annotating text fragment 55391/57597
## 2023-09-16 21:50:53 Annotating text fragment 55401/57597
## 2023-09-16 21:50:53 Annotating text fragment 55411/57597
## 2023-09-16 21:50:53 Annotating text fragment 55421/57597
## 2023-09-16 21:50:53 Annotating text fragment 55431/57597
## 2023-09-16 21:50:53 Annotating text fragment 55441/57597
## 2023-09-16 21:50:53 Annotating text fragment 55451/57597
## 2023-09-16 21:50:54 Annotating text fragment 55461/57597
## 2023-09-16 21:50:54 Annotating text fragment 55471/57597
## 2023-09-16 21:50:54 Annotating text fragment 55481/57597
## 2023-09-16 21:50:54 Annotating text fragment 55491/57597
## 2023-09-16 21:50:54 Annotating text fragment 55501/57597
## 2023-09-16 21:50:54 Annotating text fragment 55511/57597
## 2023-09-16 21:50:54 Annotating text fragment 55521/57597
## 2023-09-16 21:50:55 Annotating text fragment 55531/57597
## 2023-09-16 21:50:55 Annotating text fragment 55541/57597
## 2023-09-16 21:50:55 Annotating text fragment 55551/57597
## 2023-09-16 21:50:55 Annotating text fragment 55561/57597
## 2023-09-16 21:50:55 Annotating text fragment 55571/57597
## 2023-09-16 21:50:55 Annotating text fragment 55581/57597
## 2023-09-16 21:50:55 Annotating text fragment 55591/57597
## 2023-09-16 21:50:56 Annotating text fragment 55601/57597
## 2023-09-16 21:50:56 Annotating text fragment 55611/57597
## 2023-09-16 21:50:56 Annotating text fragment 55621/57597
## 2023-09-16 21:50:56 Annotating text fragment 55631/57597
## 2023-09-16 21:50:56 Annotating text fragment 55641/57597
## 2023-09-16 21:50:56 Annotating text fragment 55651/57597
## 2023-09-16 21:50:56 Annotating text fragment 55661/57597
## 2023-09-16 21:50:56 Annotating text fragment 55671/57597
## 2023-09-16 21:50:57 Annotating text fragment 55681/57597
## 2023-09-16 21:50:57 Annotating text fragment 55691/57597
## 2023-09-16 21:50:57 Annotating text fragment 55701/57597
## 2023-09-16 21:50:57 Annotating text fragment 55711/57597
## 2023-09-16 21:50:57 Annotating text fragment 55721/57597
## 2023-09-16 21:50:57 Annotating text fragment 55731/57597
## 2023-09-16 21:50:57 Annotating text fragment 55741/57597
## 2023-09-16 21:50:57 Annotating text fragment 55751/57597
## 2023-09-16 21:50:58 Annotating text fragment 55761/57597
## 2023-09-16 21:50:58 Annotating text fragment 55771/57597
## 2023-09-16 21:50:58 Annotating text fragment 55781/57597
## 2023-09-16 21:50:58 Annotating text fragment 55791/57597
## 2023-09-16 21:50:58 Annotating text fragment 55801/57597
## 2023-09-16 21:50:58 Annotating text fragment 55811/57597
## 2023-09-16 21:50:58 Annotating text fragment 55821/57597
## 2023-09-16 21:50:59 Annotating text fragment 55831/57597
## 2023-09-16 21:50:59 Annotating text fragment 55841/57597
## 2023-09-16 21:50:59 Annotating text fragment 55851/57597
```

```
## 2023-09-16 21:50:59 Annotating text fragment 55861/57597
## 2023-09-16 21:50:59 Annotating text fragment 55871/57597
## 2023-09-16 21:50:59 Annotating text fragment 55881/57597
## 2023-09-16 21:50:59 Annotating text fragment 55891/57597
## 2023-09-16 21:50:59 Annotating text fragment 55901/57597
## 2023-09-16 21:50:59 Annotating text fragment 55911/57597
## 2023-09-16 21:51:00 Annotating text fragment 55921/57597
## 2023-09-16 21:51:00 Annotating text fragment 55931/57597
## 2023-09-16 21:51:00 Annotating text fragment 55941/57597
## 2023-09-16 21:51:00 Annotating text fragment 55951/57597
## 2023-09-16 21:51:00 Annotating text fragment 55961/57597
## 2023-09-16 21:51:00 Annotating text fragment 55971/57597
## 2023-09-16 21:51:00 Annotating text fragment 55981/57597
## 2023-09-16 21:51:00 Annotating text fragment 55991/57597
## 2023-09-16 21:51:01 Annotating text fragment 56001/57597
## 2023-09-16 21:51:01 Annotating text fragment 56011/57597
## 2023-09-16 21:51:01 Annotating text fragment 56021/57597
## 2023-09-16 21:51:01 Annotating text fragment 56031/57597
## 2023-09-16 21:51:01 Annotating text fragment 56041/57597
## 2023-09-16 21:51:01 Annotating text fragment 56051/57597
## 2023-09-16 21:51:01 Annotating text fragment 56061/57597
## 2023-09-16 21:51:02 Annotating text fragment 56071/57597
## 2023-09-16 21:51:02 Annotating text fragment 56081/57597
## 2023-09-16 21:51:02 Annotating text fragment 56091/57597
## 2023-09-16 21:51:02 Annotating text fragment 56101/57597
## 2023-09-16 21:51:02 Annotating text fragment 56111/57597
## 2023-09-16 21:51:02 Annotating text fragment 56121/57597
## 2023-09-16 21:51:02 Annotating text fragment 56131/57597
## 2023-09-16 21:51:02 Annotating text fragment 56141/57597
## 2023-09-16 21:51:02 Annotating text fragment 56151/57597
## 2023-09-16 21:51:03 Annotating text fragment 56161/57597
## 2023-09-16 21:51:03 Annotating text fragment 56171/57597
## 2023-09-16 21:51:03 Annotating text fragment 56181/57597
## 2023-09-16 21:51:03 Annotating text fragment 56191/57597
## 2023-09-16 21:51:03 Annotating text fragment 56201/57597
## 2023-09-16 21:51:03 Annotating text fragment 56211/57597
## 2023-09-16 21:51:03 Annotating text fragment 56221/57597
## 2023-09-16 21:51:03 Annotating text fragment 56231/57597
## 2023-09-16 21:51:04 Annotating text fragment 56241/57597
## 2023-09-16 21:51:04 Annotating text fragment 56251/57597
## 2023-09-16 21:51:04 Annotating text fragment 56261/57597
## 2023-09-16 21:51:04 Annotating text fragment 56271/57597
## 2023-09-16 21:51:04 Annotating text fragment 56281/57597
## 2023-09-16 21:51:04 Annotating text fragment 56291/57597
## 2023-09-16 21:51:04 Annotating text fragment 56301/57597
## 2023-09-16 21:51:05 Annotating text fragment 56311/57597
## 2023-09-16 21:51:05 Annotating text fragment 56321/57597
## 2023-09-16 21:51:05 Annotating text fragment 56331/57597
## 2023-09-16 21:51:05 Annotating text fragment 56341/57597
## 2023-09-16 21:51:05 Annotating text fragment 56351/57597
## 2023-09-16 21:51:05 Annotating text fragment 56361/57597
## 2023-09-16 21:51:05 Annotating text fragment 56371/57597
## 2023-09-16 21:51:05 Annotating text fragment 56381/57597
## 2023-09-16 21:51:05 Annotating text fragment 56391/57597
```

```
## 2023-09-16 21:51:06 Annotating text fragment 56401/57597
## 2023-09-16 21:51:06 Annotating text fragment 56411/57597
## 2023-09-16 21:51:06 Annotating text fragment 56421/57597
## 2023-09-16 21:51:06 Annotating text fragment 56431/57597
## 2023-09-16 21:51:06 Annotating text fragment 56441/57597
## 2023-09-16 21:51:06 Annotating text fragment 56451/57597
## 2023-09-16 21:51:06 Annotating text fragment 56461/57597
## 2023-09-16 21:51:06 Annotating text fragment 56471/57597
## 2023-09-16 21:51:06 Annotating text fragment 56481/57597
## 2023-09-16 21:51:07 Annotating text fragment 56491/57597
## 2023-09-16 21:51:07 Annotating text fragment 56501/57597
## 2023-09-16 21:51:07 Annotating text fragment 56511/57597
## 2023-09-16 21:51:07 Annotating text fragment 56521/57597
## 2023-09-16 21:51:07 Annotating text fragment 56531/57597
## 2023-09-16 21:51:07 Annotating text fragment 56541/57597
## 2023-09-16 21:51:07 Annotating text fragment 56551/57597
## 2023-09-16 21:51:07 Annotating text fragment 56561/57597
## 2023-09-16 21:51:08 Annotating text fragment 56571/57597
## 2023-09-16 21:51:08 Annotating text fragment 56581/57597
## 2023-09-16 21:51:08 Annotating text fragment 56591/57597
## 2023-09-16 21:51:08 Annotating text fragment 56601/57597
## 2023-09-16 21:51:08 Annotating text fragment 56611/57597
## 2023-09-16 21:51:08 Annotating text fragment 56621/57597
## 2023-09-16 21:51:09 Annotating text fragment 56631/57597
## 2023-09-16 21:51:09 Annotating text fragment 56641/57597
## 2023-09-16 21:51:09 Annotating text fragment 56651/57597
## 2023-09-16 21:51:09 Annotating text fragment 56661/57597
## 2023-09-16 21:51:09 Annotating text fragment 56671/57597
## 2023-09-16 21:51:09 Annotating text fragment 56681/57597
## 2023-09-16 21:51:09 Annotating text fragment 56691/57597
## 2023-09-16 21:51:10 Annotating text fragment 56701/57597
## 2023-09-16 21:51:10 Annotating text fragment 56711/57597
## 2023-09-16 21:51:10 Annotating text fragment 56721/57597
## 2023-09-16 21:51:10 Annotating text fragment 56731/57597
## 2023-09-16 21:51:10 Annotating text fragment 56741/57597
## 2023-09-16 21:51:10 Annotating text fragment 56751/57597
## 2023-09-16 21:51:10 Annotating text fragment 56761/57597
## 2023-09-16 21:51:10 Annotating text fragment 56771/57597
## 2023-09-16 21:51:11 Annotating text fragment 56781/57597
## 2023-09-16 21:51:11 Annotating text fragment 56791/57597
## 2023-09-16 21:51:11 Annotating text fragment 56801/57597
## 2023-09-16 21:51:11 Annotating text fragment 56811/57597
## 2023-09-16 21:51:11 Annotating text fragment 56821/57597
## 2023-09-16 21:51:11 Annotating text fragment 56831/57597
## 2023-09-16 21:51:11 Annotating text fragment 56841/57597
## 2023-09-16 21:51:11 Annotating text fragment 56851/57597
## 2023-09-16 21:51:12 Annotating text fragment 56861/57597
## 2023-09-16 21:51:12 Annotating text fragment 56871/57597
## 2023-09-16 21:51:12 Annotating text fragment 56881/57597
## 2023-09-16 21:51:12 Annotating text fragment 56891/57597
## 2023-09-16 21:51:12 Annotating text fragment 56901/57597
## 2023-09-16 21:51:12 Annotating text fragment 56911/57597
## 2023-09-16 21:51:12 Annotating text fragment 56921/57597
## 2023-09-16 21:51:12 Annotating text fragment 56931/57597
```

```
## 2023-09-16 21:51:12 Annotating text fragment 56941/57597
## 2023-09-16 21:51:13 Annotating text fragment 56951/57597
## 2023-09-16 21:51:13 Annotating text fragment 56961/57597
## 2023-09-16 21:51:13 Annotating text fragment 56971/57597
## 2023-09-16 21:51:13 Annotating text fragment 56981/57597
## 2023-09-16 21:51:13 Annotating text fragment 56991/57597
## 2023-09-16 21:51:13 Annotating text fragment 57001/57597
## 2023-09-16 21:51:13 Annotating text fragment 57011/57597
## 2023-09-16 21:51:13 Annotating text fragment 57021/57597
## 2023-09-16 21:51:13 Annotating text fragment 57031/57597
## 2023-09-16 21:51:13 Annotating text fragment 57041/57597
## 2023-09-16 21:51:13 Annotating text fragment 57051/57597
## 2023-09-16 21:51:14 Annotating text fragment 57061/57597
## 2023-09-16 21:51:14 Annotating text fragment 57071/57597
## 2023-09-16 21:51:14 Annotating text fragment 57081/57597
## 2023-09-16 21:51:14 Annotating text fragment 57091/57597
## 2023-09-16 21:51:14 Annotating text fragment 57101/57597
## 2023-09-16 21:51:14 Annotating text fragment 57111/57597
## 2023-09-16 21:51:14 Annotating text fragment 57121/57597
## 2023-09-16 21:51:15 Annotating text fragment 57131/57597
## 2023-09-16 21:51:15 Annotating text fragment 57141/57597
## 2023-09-16 21:51:15 Annotating text fragment 57151/57597
## 2023-09-16 21:51:15 Annotating text fragment 57161/57597
## 2023-09-16 21:51:15 Annotating text fragment 57171/57597
## 2023-09-16 21:51:15 Annotating text fragment 57181/57597
## 2023-09-16 21:51:15 Annotating text fragment 57191/57597
## 2023-09-16 21:51:15 Annotating text fragment 57201/57597
## 2023-09-16 21:51:15 Annotating text fragment 57211/57597
## 2023-09-16 21:51:15 Annotating text fragment 57221/57597
## 2023-09-16 21:51:16 Annotating text fragment 57231/57597
## 2023-09-16 21:51:16 Annotating text fragment 57241/57597
## 2023-09-16 21:51:16 Annotating text fragment 57251/57597
## 2023-09-16 21:51:16 Annotating text fragment 57261/57597
## 2023-09-16 21:51:16 Annotating text fragment 57271/57597
## 2023-09-16 21:51:16 Annotating text fragment 57281/57597
## 2023-09-16 21:51:16 Annotating text fragment 57291/57597
## 2023-09-16 21:51:16 Annotating text fragment 57301/57597
## 2023-09-16 21:51:17 Annotating text fragment 57311/57597
## 2023-09-16 21:51:17 Annotating text fragment 57321/57597
## 2023-09-16 21:51:17 Annotating text fragment 57331/57597
## 2023-09-16 21:51:17 Annotating text fragment 57341/57597
## 2023-09-16 21:51:17 Annotating text fragment 57351/57597
## 2023-09-16 21:51:17 Annotating text fragment 57361/57597
## 2023-09-16 21:51:17 Annotating text fragment 57371/57597
## 2023-09-16 21:51:17 Annotating text fragment 57381/57597
## 2023-09-16 21:51:17 Annotating text fragment 57391/57597
## 2023-09-16 21:51:17 Annotating text fragment 57401/57597
## 2023-09-16 21:51:18 Annotating text fragment 57411/57597
## 2023-09-16 21:51:18 Annotating text fragment 57421/57597
## 2023-09-16 21:51:18 Annotating text fragment 57431/57597
## 2023-09-16 21:51:18 Annotating text fragment 57441/57597
## 2023-09-16 21:51:18 Annotating text fragment 57451/57597
## 2023-09-16 21:51:18 Annotating text fragment 57461/57597
## 2023-09-16 21:51:18 Annotating text fragment 57471/57597
```

```
## 2023-09-16 21:51:18 Annotating text fragment 57481/57597
## 2023-09-16 21:51:19 Annotating text fragment 57491/57597
## 2023-09-16 21:51:19 Annotating text fragment 57501/57597
## 2023-09-16 21:51:19 Annotating text fragment 57511/57597
## 2023-09-16 21:51:19 Annotating text fragment 57521/57597
## 2023-09-16 21:51:19 Annotating text fragment 57531/57597
## 2023-09-16 21:51:19 Annotating text fragment 57541/57597
## 2023-09-16 21:51:19 Annotating text fragment 57551/57597
## 2023-09-16 21:51:20 Annotating text fragment 57561/57597
## 2023-09-16 21:51:20 Annotating text fragment 57571/57597
## 2023-09-16 21:51:20 Annotating text fragment 57581/57597
## 2023-09-16 21:51:20 Annotating text fragment 57591/57597
```

thrilling

calf  past

bird  calf

happy  excite

bush  moment  pleased  overjoyed

bahubali  bridesmaid  excited  exciting

orange  memorable

peach  grateful

winery

color  toy

colour

useful  happiest

**a**  women – animal

he

monkey  enjoyment

attraction  exciting  **a**  men – animal

giraffe  deer

tiger  enthusiastic  joy  sight  **a**  men – happiness

past

tiger  bird

fourth

happy  happiest

**a**  women – happiness

relaxing  happy  happy  excit

delightful  scenery  satisfied  relieve

day

gesture

happy  day  proud  satisfied

happy  delightful  passionate  happy

bmw  thrilled  ecstatic  thrilled

celebrations  excited

sad  final  excite

squirrel  bowel

5th  pregnant  school

```
## $title
## [1] "Most similar words to animal with word2vec - umap"
##
## attr(,"class")
## [1] "labels"
```