

# 1 Data Format Task 1

As a variation on the data format described in Dataformat\_Task1\_a.pdf, we also provide our data in an alternative format which merges the Tokens and Annotations into one data structure, thus making it easier to inspect the data. Here, the json dictionary keys are the paragraph ids. For each paragraph, we include the following information, each as a list:

- **Tokens:** the list of tokens for this paragraph
- **Annotations:** a list of dictionaries, where each item in the list is a dictionary with a cue and the respective roles for this cue.

```
1 {
2   "18": {                                # Sentence id
3     "Tokens": [                          # Tokens for sentence 18
4       "Wir",
5       "als",
6       "AfD",
7       "lehnen",
8       "diese",
9       "Totalüberwachung",
10      "der",
11      "Bürger",
12      "ab",
13      "."
14    ],
15    "Annotations": [                    # Annotations for sentence 18
16                                          # (list of dictionaries, with one
17                                          # dictionary per cue word(s))
18      {
19        "Message": [                    # Message
20          "18:4",                       # sentence_id:token_id
21          "18:5",
22          "18:6",
23          "18:7"
24        ],
25        "Cue": ["18:3"],                # Cue word(s)
26        "PTC": [                        # Particles (separable verb prefixes
27          "18:8"                         # and reflexive pronouns)
28        ],
29        "Source": [                    # Source
30          "18:0",
31          "18:1",
32          "18:2"
33        ],
34      } ] }, ... }
```

A cue can consist of one word (e.g., "sagen", "schreiben", "denken", "Rede") or can include a multi-word construction (e.g., "Rede halten", "Informationen geben", "Gedanken machen"). We use the label PTC to encode verb prefixes that have been separated from the cue word (see example above). We use the same label to encode obligatory reflexive pronouns, such as "sich" in "sich Gedanken machen").

In the example above, we only have one cue (and therefore only one dictionary in the Annotations list). The cue has the paragraph id "18" and the token id "3". To retrieve the word form for this cue, you can extract the token with id 3 (i.e., the fourth token in the list) from paragraph 18, which is "lehnen". This is a particle verb and the verb prefix is encoded as PTC (paragraph 18, token id 8  $\rightarrow$  "ab"). In addition to the cue word and its verb particle, the Annotations include the roles for this cue (the complete list of roles is: Source, Message, Adresse, Topic, Medium and Evidence. For more information, see the annotation guidelines).

In some speeches, the boundaries have not been identified correctly so that some sentences have been split up, with the first part of the sentence included in one paragraph and the second part in the next paragraph. The example below illustrates such a case. Here, the source for the text fragment 10 is included in the previous sentence fragment with id 9, as indicated by the Source ids ("9:0", "9:1", ...).

The same mechanism is used to encode roles that span over multiple sentences when not all sentences are included in the same paragraph.

```

1 {
2   {
3     "9": {
4       "Tokens": [
5         "Sie",
6         ",",
7         "verehrte",
8         "Kollegen",
9         "der",
10        "FDP",
11        ",",
12        "fordern",
13        "im",
14        "Übrigen",
15        "nach",
16        "wie",
17        "vor",
18        "eine",
19        "Abschaffung",
20        "des",
21        "NetzDG",
22        ",",
23      ],
24      "Annotations": [
25        {
26          "Message": [
27            "9:13",
28            "9:14",
29            "9:15",
30            "9:16"
31          ],
32          "Source": [
33            "9:0"
34          ],
35          "Cue": [
36            "9:7"
37          ]
38        }
39      ]
40    },
41
42      # to be continued on next page
43
44
45
46
47

```

```

48
49     "10": {
50         "Tokens": [
51             "haben",
52             "aber",
53             "die",
54             "letzten",
55             "zwei",
56             "Jahre",
57             "tatenlos",
58             "verstreichen",
59             "lassen",
60             ",",
61             "diesen",
62             "völlig",
63             "überkommenen",
64             "Ansatz",
65             "Ihrer",
66             "Politik",
67             "zu",
68             "diskutieren",
69             "."
70         ],
71         "Annotations": [
72             {
73                 "Message": [
74                     "10:10",
75                     "10:11",
76                     "10:12",
77                     "10:13",
78                     "10:14",
79                     "10:15"
80                 ],
81                 "Source": [
82                     "9:0"
83                 ],
84                 "Cue": [
85                     "10:17"
86                 ]
87             }
88         ]
89     },
90     ...
91 }

```

The next example illustrated the annotation of multiple cues in the same sentence, encoded as a list of dictionaries in "Annotations".

```
1 {
2   "7": {
3     "Tokens": [
4       "Im",
5       "Übrigen",
6       "befinde",          # cue 1 (7:2)
7       "ich",
8       "mich",
9       "damit",
10      "offensichtlich",
11      "auch",
12      "in",
13      "weitgehender",
14      "Übereinstimmung",  # cue 1 (7:10)
15      "mit",
16      "der",
17      "Kanzlerin",
18      ",",
19      "die",
20      "auf",
21      "eine",
22      "entsprechende",
23      "Frage",            # cue 2 (7:19)
24      "vor",
25      "wenigen",
26      "Minuten",
27      "hier",
28      "in",
29      "diesem",
30      "Hause",
31      "genau",
32      "so",
33      "geantwortet",      # cue 3 (7:29)
34      "hat",
35      "."
36    ],
37
38
39      # to be continued on next page
40
41
42
43
44
45
```

```

46     "Annotations": [
47         {
48             "PTC": [
49                 "7:4"
50             ],
51             "Source": [
52                 "7:3"
53             ],
54             "Cue": [
55                 "7:2",
56                 "7:10"
57             ]
58         },
59         {
60             "Cue": [
61                 "7:19"
62             ]
63         },
64         {
65             "Message": [
66                 "7:27",
67                 "7:28"
68             ],
69             "Source": [
70                 "7:15"
71             ],
72             "Topic": [
73                 "7:16",
74                 "7:17",
75                 "7:18",
76                 "7:19"
77             ],
78             "Cue": [
79                 "7:29"
80             ]
81         }
82     ], ... }
83

```

The first cue, *sich in Übereinstimmung befinden*, includes two cue words, "befinde" and "Übereinstimmung" and a reflexive pronoun ("mich"), encoded as PTC. This cue also has a Source ("ich"). The second cue, *Frage*, does not have any roles. The third cue, *geantwortet*, has a Source ("die"), a Message ("genau so") and a Topic ("auf eine entsprechende Frage").

For more information on the annotations, see the Annotation Guidelines, available from our github repository.