

Introductory remarks regarding the jupyter notebook:

- I wrote some tests for the encoder, decoder and attention modules to check that everything can be plugged in together correctly.
- I adopted the tensor notations $[L, N, H]$ where L is the sentence length, N is the batch size and H is a hidden dimension / feature dimension.
- I re-implemented a random sampling strategy based on the temperature parameter. Temperature=0 leads back to (argmax) greedy sampling which is what's used during training.

I wrote most of the report in the **old fashion way (manually)**, I explicitly mentioned in section 3 where ChatGPT4 was used for an assistance.

1 Question 1: Greedy decoding strategy

Greedy decoding strategy picks the translated token with the highest probability (argmax on the softmax of the logits getting out of the decoder). There are naturally many other sentences possibilities which are not explored... but it is the least expensive solution in terms of computation cost:

- It has the advantage of being fast since you predict the next word probabilities using a single inference of the recurrent decoder neural network.
- You can override (inplace) the decoder hidden state and context every time you go to the next word.

There are a few solutions as mentioned in [2]

- Exploring more solutions would naturally mean a higher computation cost. If you'd like to explore several possibilities, for each word you'd like to test, you'd have to run an inference to predict the next word (and in terms of memory, you'd need to keep a copy of the hidden states). This would be quite heavy and leads to exponential combinations (the word decision tree quickly becomes huge). In beam search you may keep only a fixed amount of combinations.
- Performing random sampling. Instead of taking the most probable word, why not sampling randomly from the predicted distribution. This may lead to unexpected results as we add a bit of randomness in the process. 1

Finally, despite trying a simple different sampling strategy, I still cannot improve the translation quality. I think it is hard to tell whether or not the sampling strategy impacts the final results as the overall model quality seems pretty low. See Section 2. for more details about this.

```
src_sentence = "The cat sat on the mat."
src_tokens = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009, 1010, 1011, 1012, 1013, 1014, 1015, 1016, 1017, 1018, 1019, 1020, 1021, 1022, 1023, 1024, 1025, 1026, 1027, 1028, 1029, 1030, 1031, 1032, 1033, 1034, 1035, 1036, 1037, 1038, 1039, 1040, 1041, 1042, 1043, 1044, 1045, 1046, 1047, 1048, 1049, 1050, 1051, 1052, 1053, 1054, 1055, 1056, 1057, 1058, 1059, 1060, 1061, 1062, 1063, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080, 1081, 1082, 1083, 1084, 1085, 1086, 1087, 1088, 1089, 1090, 1091, 1092, 1093, 1094, 1095, 1096, 1097, 1098, 1099, 1100, 1101, 1102, 1103, 1104, 1105, 1106, 1107, 1108, 1109, 1110, 1111, 1112, 1113, 1114, 1115, 1116, 1117, 1118, 1119, 1120, 1121, 1122, 1123, 1124, 1125, 1126, 1127, 1128, 1129, 1130, 1131, 1132, 1133, 1134, 1135, 1136, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144, 1145, 1146, 1147, 1148, 1149, 1150, 1151, 1152, 1153, 1154, 1155, 1156, 1157, 1158, 1159, 1160, 1161, 1162, 1163, 1164, 1165, 1166, 1167, 1168, 1169, 1170, 1171, 1172, 1173, 1174, 1175, 1176, 1177, 1178, 1179, 1180, 1181, 1182, 1183, 1184, 1185, 1186, 1187, 1188, 1189, 1190, 1191, 1192, 1193, 1194, 1195, 1196, 1197, 1198, 1199, 1200, 1201, 1202, 1203, 1204, 1205, 1206, 1207, 1208, 1209, 1210, 1211, 1212, 1213, 1214, 1215, 1216, 1217, 1218, 1219, 1220, 1221, 1222, 1223, 1224, 1225, 1226, 1227, 1228, 1229, 1230, 1231, 1232, 1233, 1234, 1235, 1236, 1237, 1238, 1239, 1240, 1241, 1242, 1243, 1244, 1245, 1246, 1247, 1248, 1249, 1250, 1251, 1252, 1253, 1254, 1255, 1256, 1257, 1258, 1259, 1260, 1261, 1262, 1263, 1264, 1265, 1266, 1267, 1268, 1269, 1270, 1271, 1272, 1273, 1274, 1275, 1276, 1277, 1278, 1279, 1280, 1281, 1282, 1283, 1284, 1285, 1286, 1287, 1288, 1289, 1290, 1291, 1292, 1293, 1294, 1295, 1296, 1297, 1298, 1299, 1300, 1301, 1302, 1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1312, 1313, 1314, 1315, 1316, 1317, 1318, 1319, 1320, 1321, 1322, 1323, 1324, 1325, 1326, 1327, 1328, 1329, 1330, 1331, 1332, 1333, 1334, 1335, 1336, 1337, 1338, 1339, 1340, 1341, 1342, 1343, 1344, 1345, 1346, 1347, 1348, 1349, 1350, 1351, 1352, 1353, 1354, 1355, 1356, 1357, 1358, 1359, 1360, 1361, 1362, 1363, 1364, 1365, 1366, 1367, 1368, 1369, 1370, 1371, 1372, 1373, 1374, 1375, 1376, 1377, 1378, 1379, 1380, 1381, 1382, 1383, 1384, 1385, 1386, 1387, 1388, 1389, 1390, 1391, 1392, 1393, 1394, 1395, 1396, 1397, 1398, 1399, 1400, 1401, 1402, 1403, 1404, 1405, 1406, 1407, 1408, 1409, 1410, 1411, 1412, 1413, 1414, 1415, 1416, 1417, 1418, 1419, 1420, 1421, 1422, 1423, 1424, 1425, 1426, 1427, 1428, 1429, 1430, 1431, 1432, 1433, 1434, 1435, 1436, 1437, 1438, 1439, 1440, 1441, 1442, 1443, 1444, 1445, 1446, 1447, 1448, 1449, 1450, 1451, 1452, 1453, 1454, 1455, 1456, 1457, 1458, 1459, 1460, 1461, 1462, 1463, 1464, 1465, 1466, 1467, 1468, 1469, 1470, 1471, 1472, 1473, 1474, 1475, 1476, 1477, 1478, 1479, 1480, 1481, 1482, 1483, 1484, 1485, 1486, 1487, 1488, 1489, 1490, 1491, 1492, 1493, 1494, 1495, 1496, 1497, 1498, 1499, 1500, 1501, 1502, 1503, 1504, 1505, 1506, 1507, 1508, 1509, 1510, 1511, 1512, 1513, 1514, 1515, 1516, 1517, 1518, 1519, 1520, 1521, 1522, 1523, 1524, 1525, 1526, 1527, 1528, 1529, 1530, 1531, 1532, 1533, 1534, 1535, 1536, 1537, 1538, 1539, 1540, 1541, 1542, 1543, 1544, 1545, 1546, 1547, 1548, 1549, 1550, 1551, 1552, 1553, 1554, 1555, 1556, 1557, 1558, 1559, 1560, 1561, 1562, 1563, 1564, 1565, 1566, 1567, 1568, 1569, 1570, 1571, 1572, 1573, 1574, 1575, 1576, 1577, 1578, 1579, 1580, 1581, 1582, 1583, 1584, 1585, 1586, 1587, 1588, 1589, 1590, 1591, 1592, 1593, 1594, 1595, 1596, 1597, 1598, 1599, 1600, 1601, 1602, 1603, 1604, 1605, 1606, 1607, 1608, 1609, 1610, 1611, 1612, 1613, 1614, 1615, 1616, 1617, 1618, 1619, 1620, 1621, 1622, 1623, 1624, 1625, 1626, 1627, 1628, 1629, 1630, 1631, 1632, 1633, 1634, 1635, 1636, 1637, 1638, 1639, 1640, 1641, 1642, 1643, 1644, 1645, 1646, 1647, 1648, 1649, 1650, 1651, 1652, 1653, 1654, 1655, 1656, 1657, 1658, 1659, 1660, 1661, 1662, 1663, 1664, 1665, 1666, 1667, 1668, 1669, 1670, 1671, 1672, 1673, 1674, 1675, 1676, 1677, 1678, 1679, 1680, 1681, 1682, 1683, 1684, 1685, 1686, 1687, 1688, 1689, 1690, 1691, 1692, 1693, 1694, 1695, 1696, 1697, 1698, 1699, 1700, 1701, 1702, 1703, 1704, 1705, 1706, 1707, 1708, 1709, 1710, 1711, 1712, 1713, 1714, 1715, 1716, 1717, 1718, 1719, 1720, 1721, 1722, 1723, 1724, 1725, 1726, 1727, 1728, 1729, 1730, 1731, 1732, 1733, 1734, 1735, 1736, 1737, 1738, 1739, 1740, 1741, 1742, 1743, 1744, 1745, 1746, 1747, 1748, 1749, 1750, 1751, 1752, 1753, 1754, 1755, 1756, 1757, 1758, 1759, 1760, 1761, 1762, 1763, 1764, 1765, 1766, 1767, 1768, 1769, 1770, 1771, 1772, 1773, 1774, 1775, 1776, 1777, 1778, 1779, 1780, 1781, 1782, 1783, 1784, 1785, 1786, 1787, 1788, 1789, 1790, 1791, 1792, 1793, 1794, 1795, 1796, 1797, 1798, 1799, 1800, 1801, 1802, 1803, 1804, 1805, 1806, 1807, 1808, 1809, 1810, 1811, 1812, 1813, 1814, 1815, 1816, 1817, 1818, 1819, 1820, 1821, 1822, 1823, 1824, 1825, 1826, 1827, 1828, 1829, 1830, 1831, 1832, 1833, 1834, 1835, 1836, 1837, 1838, 1839, 1840, 1841, 1842, 1843, 1844, 1845, 1846, 1847, 1848, 1849, 1850, 1851, 1852, 1853, 1854, 1855, 1856, 1857, 1858, 1859, 1860, 1861, 1862, 1863, 1864, 1865, 1866, 1867, 1868, 1869, 1870, 1871, 1872, 1873, 1874, 1875, 1876, 1877, 1878, 1879, 1880, 1881, 1882, 1883, 1884, 1885, 1886, 1887, 1888, 1889, 1890, 1891, 1892, 1893, 1894, 1895, 1896, 1897, 1898, 1899, 1900, 1901, 1902, 1903, 1904, 1905, 1906, 1907, 1908, 1909, 1910, 1911, 1912, 1913, 1914, 1915, 1916, 1917, 1918, 1919, 1920, 1921, 1922, 1923, 1924, 1925, 1926, 1927, 1928, 1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939, 1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952, 1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 208
```

- **Hyper parameters have not been tweaked.** A learning rate schedule would have been nice (e.g after 80 epochs, learning rate should have been decreased to maybe keep on going).
- **Slight hints of start of overfitting:** It looks like during my training with a learning rate of 10^{-3} and all default parameters provided, the test got to a plateau but the train loss kept on decreasing, meaning a slight risk of overfitting on the training set (the test loss did not go up though so it's not an explicit sign of overfit, just a hint that the network is not learning to generalize as much as in the beginning.) *Memorizing language patterns may be a good idea, but not training sentences word by word.*

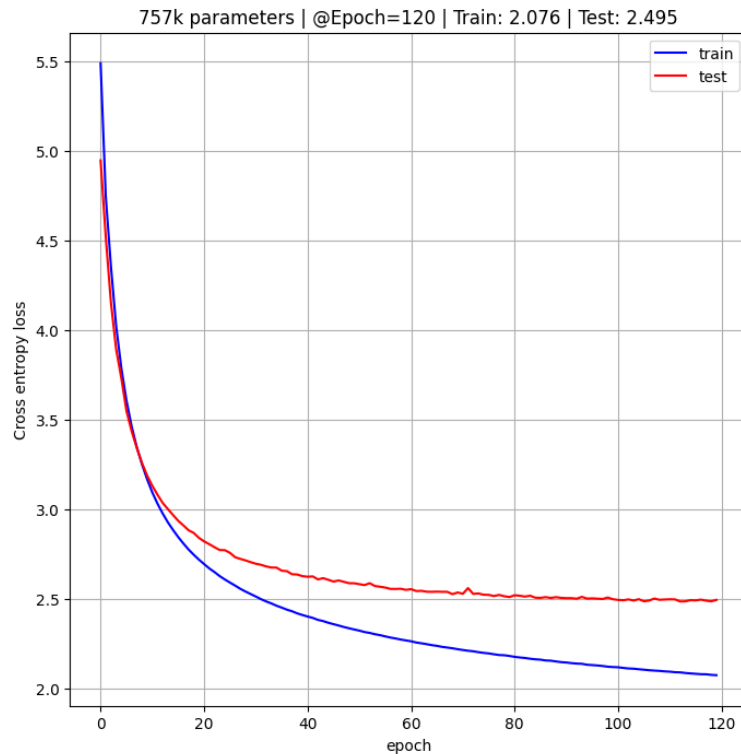


Figure 2: Training metrics with default parameters ($lr = 10^{-3}$, $batch\ size = 128$)

Here are some key design points which may also have been re-considered for future work.

- **Shallow network:** The encoder tries to build a hidden state which compresses the whole sentence meaning into a single vector. Since there's only a single GRU (no a deep stack of GRU) in the stack, the encoder is doing pretty simple operations by making temporal combinations of the word embeddings thanks to the GRU. The lack of depth (deep means we have many linear weights multiplication followed by non linearities) forces the encoder to make basic representations.
- There are 757k parameters in the model, most of the parameters are used to build the embeddings (Source embeddings 211k, Target embedding 298k = in total 509k parameters are basically a look up table to go from words to floating point vectors. Decoder prediction has 223k parameters to project back the hidden to probabilities (basically a reverse lookup table). To learn the structure of sentences and language, the GRU has a few parameters to learn which feature dimension to focus on (memorize) to create a sentence summary vector.
- The source sentence is analyzed in a causal manner in the encoder. The GRU cannot look backwards (uni-directional). We'll come back to that point in Section 4.
- The "alignment" mechanism performs weighted sums of the current decoded hidden state and all source encoded vectors. With linear combinations and a \tanh function, **it does not seem straightforward that similarities between words can be computed.** The cosine distance/dot product or the L2 distance between target hidden state and source hidden vectors sounds much more appropriate... and this is what ended up being used in Transformers[5] obviously.

Overall, when looking at the translation results, it is difficult to decouple issues which may come from true limitations (inherent translation task challenges, architecture design) or what comes from the short amount of time spent on polishing the training (basically not picking low hanging fruits that a machine learning practitioner would usually try with a bit more time).

She eats a pizza, her friend opens a beer and her baby is sleeping.
-> elle mange une pizza , l son ami de une la et son père père .

Let's still try to analyze some translation results as best as we can.

- Over-translation: repeated words (père, père)
- Under-translation: missed words (bière)
- Hallucination: translated sentences which almost make sense but do not match the source content ("et son père" does not match with the source sentence but still sounds alright).
- Difficulty to capture long range dependencies (difficult to illustrate here: "her" in "her baby" refers to the woman eating a pizza which was at the beginning of the sentence - not to a man "père"...).

* * *

Greedy sampling is performed both during training time and inference. This is a nice point as there is no gap between training and inference. But the decoder is not trained with the groundtruth previous words (assuming the past translation is perfect). This probably makes the task even more difficult e.g. anytime the model makes a mistake, it will be almost impossible to learn the rest of the translation. We probably remove a bunch of learning signal for this reason.

The cat fell asleep in front of the fireplace
-> le chat s est en du du pression peigne peigne cheminée portail
portail portail portail portail portail portail portail indépendant
oiseaux oiseaux oiseaux oiseaux oiseaux oiseaux oiseaux oiseaux
oiseaux oiseaux

As soon as the translated words start to mess up, the decoder ends up in a similar "bad" situation as seen during training. But it probably never got out of that trap at training time either.

I have blue eyes. -> j ai les yeux bleus bleus .

We get a some repeated words (over-translation). An easy pragmatic way to solve this would be to manually discard these repetitions afterward. But the interesting fact is that the model creates these repeated words. Although there's a no explicit way to debug /explain this issue, *a good heuristic which could be used is that a word which has already been translated shall get less attention than the other remaining words.* In the sentence "I have blue eyes.", bleus shall attend to "eyes" to know whether it's a plural or singular. So keep in mind that totally discarding the word "eye" after it's been translated may not be an ideal idea. Still, if a mechanism is built in the end to end training process to perform this task, it seems like a good idea.

This is the idea behind the Coverage Mechanism proposed in [4].

I do not pretend to fully understand the implementation details, but below is the general idea behind the coverage mechanism.

There is a direct way to keep track of attention history. For every decoder time step, an accumulated attention weight is maintained, which stores the sum of attention weights over all previous decoder time steps for each source word. A concrete example: in 3, you'd progressively sum vertically. This accumulated weight is then used in conjunction with the current attention distribution to penalize attending to source words that have already been attended to. The idea is to guide the attention mechanism to focus on unattended or less-attended parts of the source sentence. At training time, this would be implemented as an extra penalty term added to the Cross Entropy Loss we used in the lab session. The under-translation/missing word issue seems to be solvable in the same fashion. In 3, we can visualize that "Red" has almost not been matched to any translated word (the column is a bit darker than some other columns like "wine" or "open").

3 Question 3: Alignment visualization

Do not open this bottle of red wine. → n ouvrez pas cette bouteille de vin .

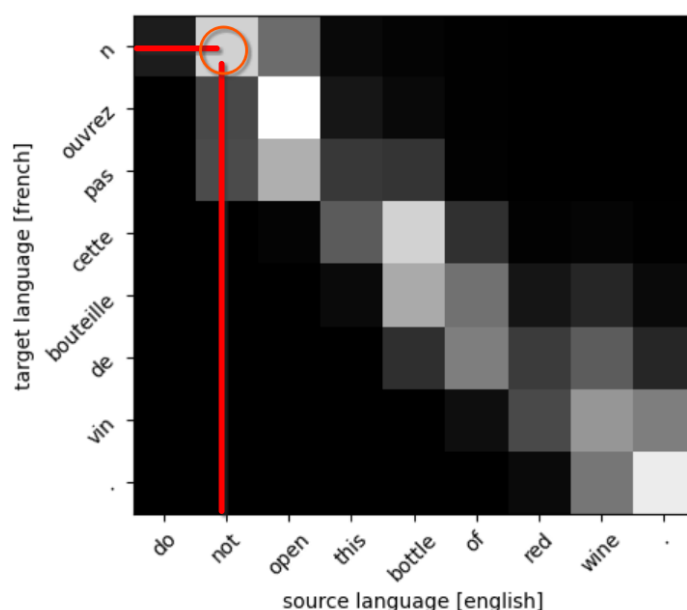


Figure 3: Alignment scores captured on the sentence label. Each row describes a translated word in the target model. The sum of all score elements over a row shall be 1 to be able to perform the weighted sum (this is granted by the softmax).

Figure 3 illustrates a difficult word inversion captured by a high score between french "n" and the second sentence word in english "not". Let's take the first row, french word "n" was compared to the 9 words of the source sentence. We can see that the most relevant english word was "not" according to the score. Which wasn't the first word of the sentence. A naive word by word translation would have led to "faire ne ouvrir"... here the alignment/attention mechanism most helped the first translated word to look a bit ahead in the source sentence.

Having the ability to **perform word inversion during decoding is a mandatory component in machine translation.**

On the other hand, we can also see one of the drawback of the current alignment mechanism: some words have been skipped. Where we would have expected "rouge" to be, "vin" was translated instead... This is pretty straightforward to explain, the score mostly matched with the "wine" word but the "red" word got less attention (it was clearly forgotten here). We mentioned the under translation phenomenon in the previous section, here it seems like we almost get an explicit explanation of the phenomenon.

She is a talented artist. → elle est un artiste artiste .

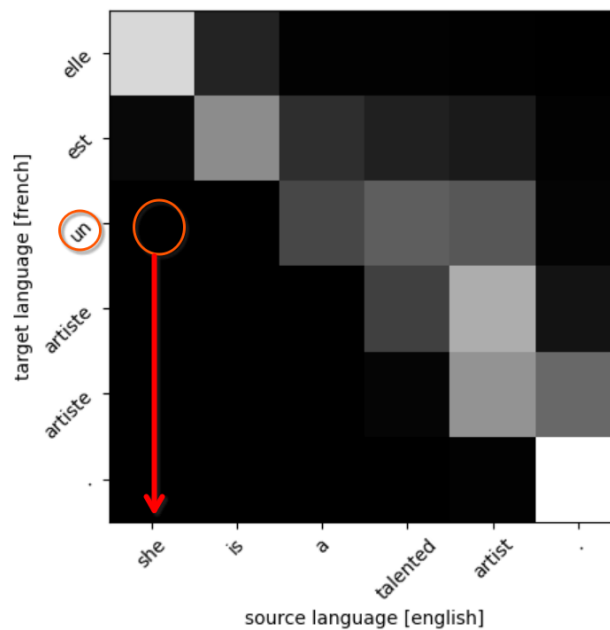


Figure 4: Gender inference fails.

Note: The following list was inspired with an assistance from ChatGPT 4. Examples were tested with the pretrained model. (refer to the examples in the notebook for full content).

Here's a list of difficult translations phenomenon where attention is a key ingredient for success.

- **Adjective-Noun Inversion:** English: "blue car" / French: "voiture bleue" (literally: "car blue")
- **Verb-Subject Inversion in Questions:** English: "Are you coming?" / French: "Viens-tu?" (literally: "Come-you?")
- **Gender and Number Agreement:** English: "A car" / French: "Une voiture" (feminine) French nouns have gender (masculine or feminine), which affects the form of adjectives, articles, and sometimes verbs. English doesn't have this gender distinction for inanimate objects. A failure case is illustrated in 4
- **T-V Distinction:** French, like many other languages, has a distinction between formal and informal address ("vous" vs. "tu"). English uses "you" for both. A failure is illustrated in 5
- **Subjunctive Mood:** French: "Il faut que tu viennes" (It is necessary that you come) English might simply say: "You need to come." The subjunctive mood is used more frequently in French than in English. It can be challenging to translate accurately because its use in English is more limited and often context-specific.
- **Negation:** English: "I don't know" / French: "Je ne sais pas" (literally: "I no know not") French typically uses a two-part negation. 6
- **Placement of Adverbs:** English: "I often eat apples" French: "Je mange souvent des pommes" (literally: "I eat often apples") The placement of adverbs can vary between English and French.
- **Prepositions:** English: "I am good at maths" / French: "Je suis bon en maths" (literally: "I am good in maths") The use and choice of prepositions can differ significantly between English and French.

you are a nice person Sir. → tu êtes une personne monsieur m

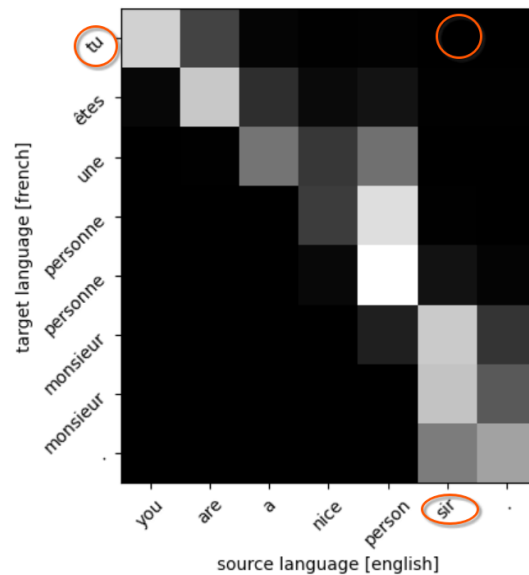


Figure 5: T-V distinction fails. "Tu" shall have looked at the last source sentence word "Sir" to disambiguate. But that score circles in red is close to 0

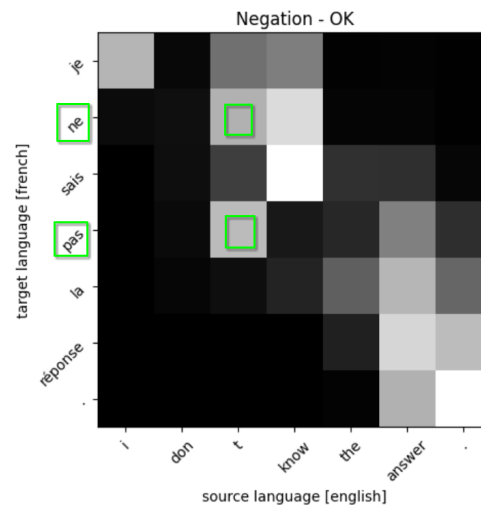


Figure 6: Negation : Je ne sais pas la réponse, ne and "pas" both had scores matching that abbreviated t = not

4 Question 4: Polysemy

Note: ChatGPT4 was powerful enough to retrieve the concept "Polysemy" when prompted with the example. incredible

I did not mean to hurt you
She is so mean

Polysemy : The same word may have different meaning depending on the context. In the first case, "mean" is a verb. In the second case "mean" is an adjective. But both are written the same in english and **will get the same embedding vector in the encoder.**

TL/DR: Here are some key ingredients to address polysemy and general ideas to improve

- **Attend the whole sequence:** For each word/token in the source sentence, looking at the whole source sentence is important. Self attention looks at the whole sentence [5] (bidirectional in the encoder - but decoder is unidirectional), **Bidirectional** RNN as proposed in ElMo use 2 passes [3] . Caveat: in our lab, we used uni-directional GRU, this is a potential improvement to consider).

- **Stacking several layers of linearity/non linearities** allows creating powerful representations. Working with a single linear word embedding + a single GRU (no stacks) like we did in the lab is probably not enough.
- **Next word prediction in the source language.** Working on the encoder itself can be decoupled from the translation task. **Powerful representations can be created by only looking at a corpus of the source language** with much more data available. Next word prediction allows pretraining the encoder (so called unsupervised training as you don't need pairs of labeled data for instance). Note: next word prediction does not allow the attention to be bi-directional. BERT [1] pretrains on the masked word prediction task which allows the attention to be bidirectional (only at the encoder level).

* * *

Disclaimer I wrote the following paragraph with my own reflections and read BERT [1] and ElMO [3] only at the end, just to try to push my reflection as far as I could go. Writing this answer sort of helped me understand some of the Transformers design choice which I took basically from granted.

*

One may expect that some components of the feature vectors for each word to specialize in describing higher level concepts (e.g. what gender it represents or the type of word like adjective/noun/verb). The default embedding for "mean" would have say 0.5 for the adjective dimension, 0.5 for the "is it a verb" dimensions. Example: we have 30 dimensions in the encoder hidden state, let's assume that dimension 28 for instance is dedicated to describing how likely a given word is a verb, and the 13th component is dedicated to describe how likely a word is an adjective. The recurrence (GRU) may help to refine the knowledge of the "type of word". In the second translation case (she is so mean): the encoder may refine the 28th component of the hidden state to be near 0 and the 13th component to be closer to 1. This is pure intuition and interpretation of the mechanism of the neural network, things may not work that way in practice unfortunately. Unfortunately:

- In the current setting, we're hopping that the encoder uses the temporal aspect of the gated recurrent unit (GRU) to compute a context in the input sentence (decide what to memorize). Issue is that there are not many parameters in the GRU and that the GRU only accesses previous elements. It is a causal recurrent unit.
- This relevant piece information (adjective/verb) may be lost in the middle of many other significant embedding dimensions (like simply embedding the meaning of the word, not grammar properties or more abstract thing) as this is a translation task.

Several ideas appear before we'd start considering self-attention in the source language:

- Chain several layers on linear/non linear so a lot more complex representations can be created and hope that each GRU start to specify in given subtasks. In ElMO [3], they stacked 2 bidirectional LSTM.
- Chunk hidden feature into several smaller vectors and feed to several GRU in parallel. Although there is no explicit way to teach the network to specialize its feature into learning specific linguistic concepts, it seems like forcing to learn independent language properties (like grammar, gender, word type etc...) is worth trying to improve performances. This is a similar idea to MHA (multi-head attention) which also seems to be one of the ingredient of the Transformer: it ends up with specialized features (intuition = specialized in analyzing grammar or singular/plural features).

* * *

A simple remark: **why would we wait for the decoding step to start doing implicit sentence analysis?**. After all, we don't need the translation task to be able to analyze the sentence in the source language.

- We can think as the encoder outside of the "translation" task and simply stay in a "unsupervised" next word prediction task in the source language. This would help to build as a pre-training step hidden vector representation which are more powerful than then ones in the translation supervision settings. Reason is simple: we'd be able to get much more data as we don't need pairs of translated sentences. Just a raw french or english corpus for instance.
- The attention mechanism looks like a great idea to take context into account... Why restricting it to the cross language attention? Why not performing self attention on the source domain? (and actually be able to replace the RNN). It became sort of a standard principle in the Transformer paper [5]: using **self attention**.

* * *

References

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [2] Huggin Face. How to generate text: using different decoding methods for language generation with transformers. <https://huggingface.co/blog/how-to-generate>.
- [3] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *CoRR*, abs/1802.05365, 2018.
- [4] Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu, and Hang Li. Coverage-based neural machine translation. *CoRR*, abs/1601.04811, 2016.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.