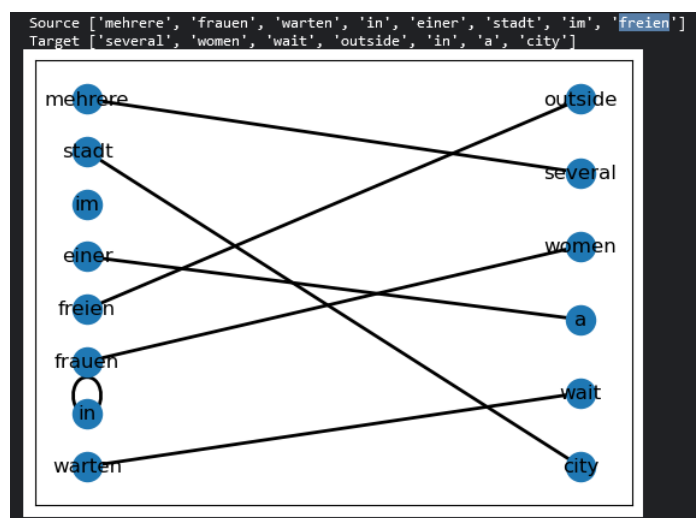# CS288 HW2 Writeup

Will Lavanakul
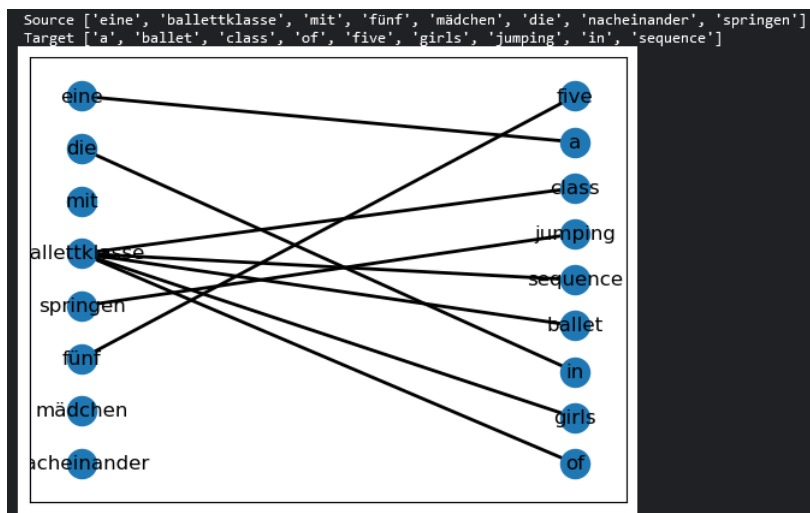
## 1 Alignment with IBM Model 1

### 1.1 Performance on Multi30k

The model performs well on the sentence ['mehrere', 'frauen', 'warten', 'in', 'einer', 'stadt', 'im', 'freien'] aligned to ['several', 'women', 'wait', 'outside', 'in', 'a', 'city'].



The model performs bad on ['eine', 'ballettklasse', 'mit', 'fünf', 'mädchen', 'die', 'nacheinander', 'springen'] translated into ['a', 'ballet', 'class', 'of', 'five', 'girls', 'jumping', 'in', 'sequence'].



I believe that the difference between in performance comes from the rarity of specific source words. In the first example, a lot of the source words are common words more likely to be seen in other sentences. In the second example 'ballettklasse' is a more rare word so when trying to align to 'ballettklasse', the model underperforms and aligns to incorrect words.

## 2 Attention Visualization

The plots are as expected. Most alignments are one-to-one. In some cases, the attention is distributed over multiple words. In these cases, the words might be too sparse in the target set, or a word with multiple uses in the source such as punctuation.

Compared to the IBM model, it seems attention works better with specific words but can have more distributed alignments on some words. In the IBM model, the alignments were trained to convergence. With attention, it has the option to assign weights that can lead to different meanings given a specific context or state. The IBM model can wrongly align words since the alignments are fixed.