

German is known for its highly productive word formation processes, particularly in the area of compounding and derivation. In this thesis, we focus on German nominal compounds and their representation in machine translation (MT) outputs. Despite their importance in German text, commonly used metrics for MT evaluation, such as BLEU, do not adequately capture the usage of compounds. The aim of this thesis was to investigate the generation of German compounds in Transformer models and to explore the conditions that lead to their production. Our analysis revealed that MT systems tend to produce fewer compounds than humans. However, we found that due to the highly productive nature of German compounds, it is not feasible to identify them based on a fixed list. Therefore, we manually identified novel compounds, and even then, human translations still contained more compounds than MT systems.

We trained our own Transformer model for English-German translation and conducted experiments to examine various factors that influence the production of compounds, including word segmentation and the frequency of compounds in the training data. Additionally, we explored the use of forced decoding and the impact of providing the model with the first words of a sentence during translation. Our findings highlight the importance of further research in developing MT models that are better suited to producing compounds in line with human translation.