

might represent which family tree Colin was in, what branch of that tree he was in, what generation he was from, etc. One can think of the neural network as computing learned rules relating these attributes together in order to obtain the desired predictions. The model can then make predictions such as inferring who is the mother of Colin.

The idea of forming an embedding for a symbol was extended to the idea of an embedding for a word by [Deerwester *et al.* \(1990\)](#). These embeddings were learned using the SVD. Later, embeddings would be learned by neural networks.

The history of natural language processing is marked by transitions in the popularity of different ways of representing the input to the model. Following this early work on symbols or words, some of the earliest applications of neural networks to NLP ([Mäikkyläinen and Dyer, 1991](#); [Schmidhuber, 1996](#)) represented the input as a sequence of characters.

[Bengio *et al.* \(2001\)](#) returned the focus to modeling words and introduced neural language models, which produce interpretable word embeddings. These neural models have scaled up from defining representations of a small set of symbols in the 1980s to millions of words (including proper nouns and misspellings) in modern applications. This computational scaling effort led to the invention of the techniques described above in section [12.4.3](#).

Initially, the use of words as the fundamental units of language models yielded improved language modeling performance ([Bengio *et al.*, 2001](#)). To this day, new techniques continually push both character-based models ([Sutskever *et al.*, 2011](#)) and word-based models forward, with recent work ([Gillick *et al.*, 2015](#)) even modeling individual bytes of Unicode characters.

The ideas behind neural language models have been extended into several natural language processing applications, such as parsing ([Henderson, 2003, 2004](#); [Collobert, 2011](#)), part-of-speech tagging, semantic role labeling, chunking, etc, sometimes using a single multi-task learning architecture ([Collobert and Weston, 2008a](#); [Collobert *et al.*, 2011a](#)) in which the word embeddings are shared across tasks.

Two-dimensional visualizations of embeddings became a popular tool for analyzing language models following the development of the t-SNE dimensionality reduction algorithm ([van der Maaten and Hinton, 2008](#)) and its high-profile application to visualization word embeddings by Joseph Turian in 2009.