

# Flavours of Subword Units

- Byte Pair Encoding (BPE, Sennrich et al. (2016))  
<http://github.com/rsennrich/subword-nmt/>
- Google Wordpieces (Wu et al., 2016)  
Code probably unavailable, used in speech.
- SubwordTextEncoder in Tensor2tensor (Vaswani et al., 2017)  
<https://github.com/tensorflow/tensor2tensor>

<b>STE</b>	Blíží_ se_ k_ tobě_ <b>tramvaj</b> _ ._ Z_ <b>tramvaj</b> e_ nevysto upil i_ ._ <hr/>
<b>BPE</b>	Blíží se k tobě <b>tramvaj</b> . Z <b>tramva@@ je</b> nevy@@ stoupili . <hr/>
<b>BPE underscore</b>	Blíží_ se_ k_ tobě_ <b>tramvaj@@</b> _ ._ Z_ <b>tramvaj@@ e_</b> nevy@@ stoupili_ ._ <hr/>

The best now is SentencePiece: <https://github.com/google/sentencepiece>