

Special tokens in NLP

- Common special tokens in NLP

Token	Meaning / Use
<pad>	padding to make sequence the same length in a batch
<eos>	End of sentence (marks where text stops) → stops GPT generation
<bos>	Beginning of sentence (used in some models to indicate start)
<unk>	Unknown token (used for words not in vocab)
<cls>	classification token (used in BERT for sentence level tasks) Ex. i.love pizza → <cls> I love pizza (BERT would add cls to start) → cls represents whole sentence . after encoding , the vector for cls is used for tasks like, sentiment classification , topic classification etc.
<sep>	separator (splits sentences) or segments in a sequence) helps separate different segments in sentence like making Q A pairs Ex <cls> how are you? <sep> i am fine (cls → represents whole input for classification , <sep> → marks boundary between Q and A (sep = segment separator))
<mask>	masked tokens (used for masking tokens see p120)