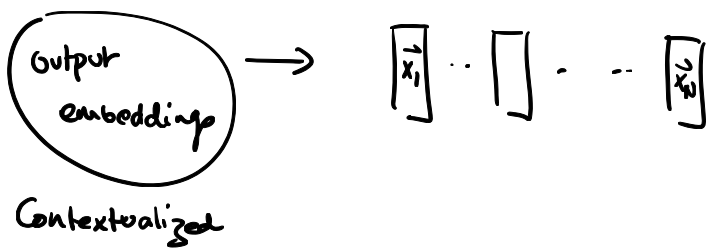
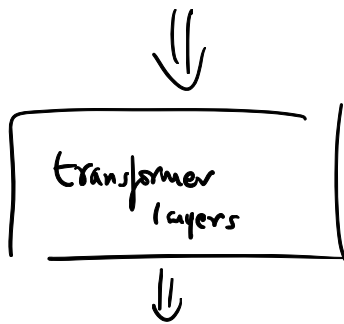
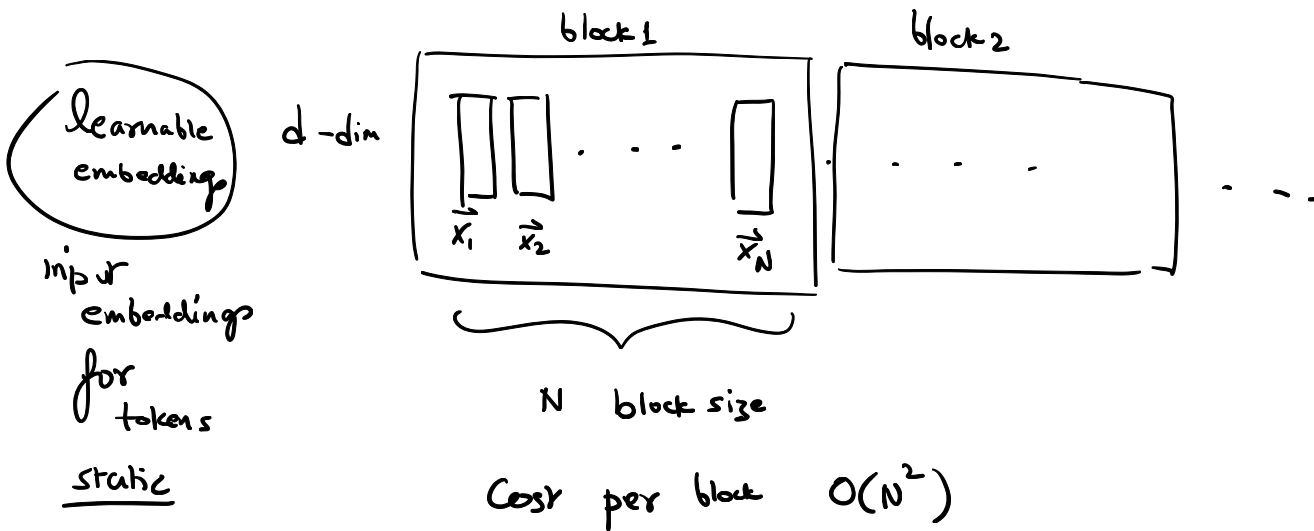


Transformers

dot-product attention

Input: block / "sequence"



Same token
different embeddings
in diff blocks

Self-attention

