



Image1

Image source:
<https://engineering.nyu.edu/faculty/yann-lecun> (2024 May 13th)

Text1 → Tokenize & embed

Context



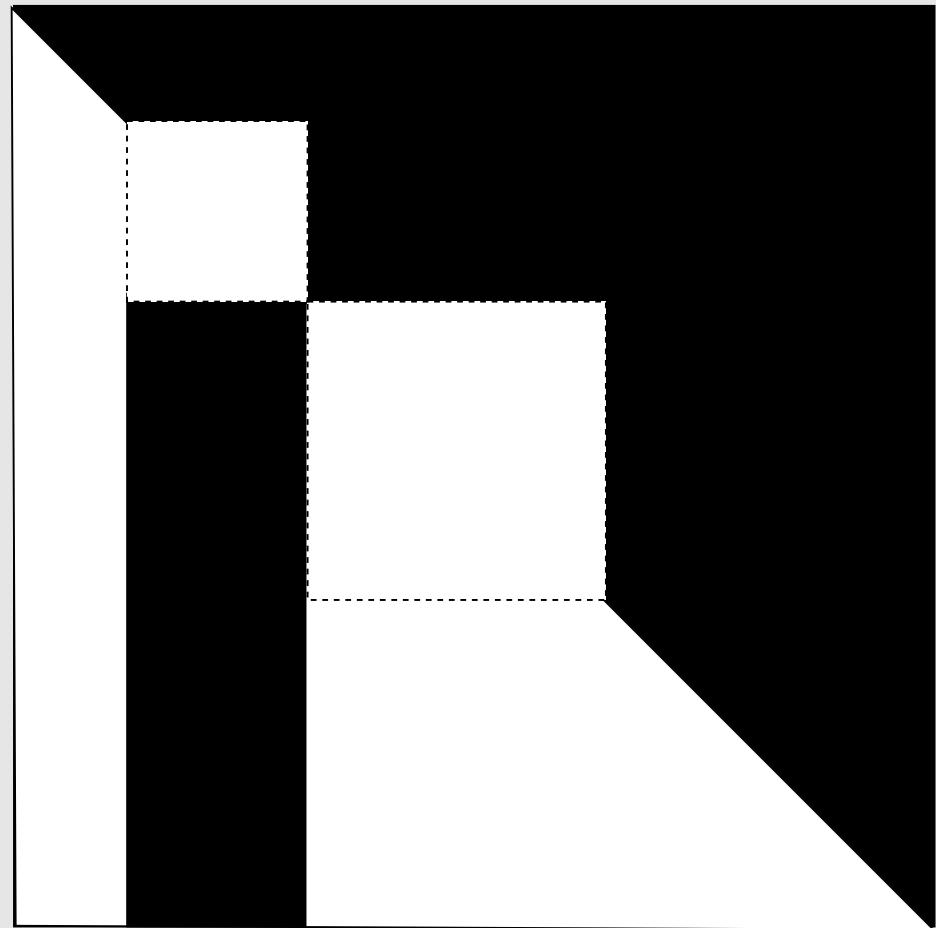
Patch, affine trafo ,
(optional) pooling

Patch, affine trafo ,
(optional) pooling

Text2 → Tokenize & embed

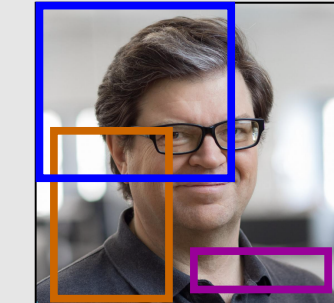
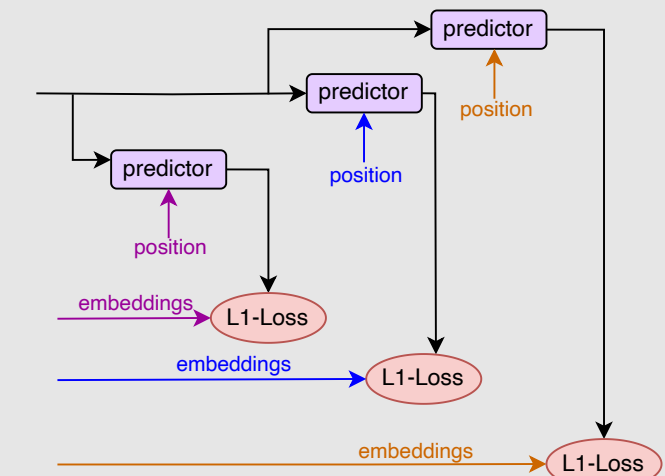
Embeddings

Transformer characterized by Attention mask seen below



Embeddings

Unembed → predictions → CE-Loss ← labels ← Tokenize & 1-hot encode ← Text1 (shifted)



Positions represented by embeddings; for visualization only, we work with the actual embeddings

Unembed → predictions → CE-Loss ← labels ← Tokenize & 1-hot encode ← Text2 (shifted)