

image source:
<https://ai.meta.com/blog/v-jepa-yann-lecun-ai-model-video-joint-embedding-predictive-architecture/>
(2024 May 12th)

Text1

Tokenize & embed

3D conv,
pos embs,
flatten,
(optional) pooling

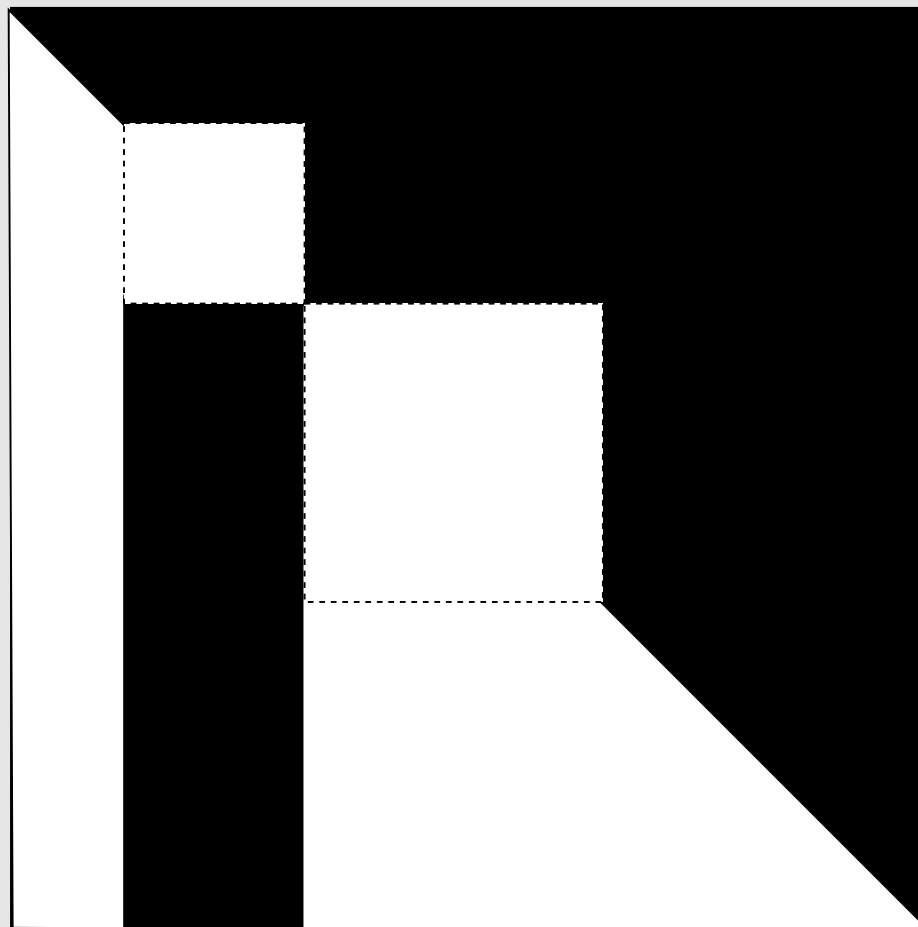
Text2

Tokenize & embed

3D conv,
pos embs,
flatten,
(optional) pooling

Transformer characterized by Attention mask seen below

Embeddings



Embeddings

