

Image $\xrightarrow{\text{ViT}}$ *Embedding*

caption \downarrow

\updownarrow (disentangling);
(token mapping)

Sentence $\xleftrightarrow{\text{LLM}}$ *Word tokens*