



f_θ : Visual Encoder

Linear

2D positional embedding

+

P_1 : Projection Layer

$l \in \{image, text\}$

m_l

P_ψ : Universal Projection (UP)

Concat

D_w : Decoder

A cat with white paws jumps over a fence in front of a yellow tree \emptyset



\emptyset no object

A cat with white paws jumps over a fence in front of a yellow tree

g_ϕ : Text Encoder

Linear

Cross-Fusion

P_2 : Projection Layer