



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

f_θ : Visual Encoder

g_ϕ : Text Encoder

Linear

Linear

+

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

m_l

P_ψ : Universal Projection (UP)

Concat

D_w : Decoder

2D positional embedding

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



\emptyset no object