



Pretrained and frozen



Trained from scratch

Perceiver
Resampler

Vision
Encoder



Perceiver
Resampler

Vision
Encoder



Output: text



a very serious cat.

n-th LM block



n-th GATED XATTN-DENSE

1st LM block



1st GATED XATTN-DENSE

Processed text

<image> This is a very cute dog.<image> This is

Interleaved visual/text data



This is a very cute dog.



This is