

Multimedia Search from Composite Inputs

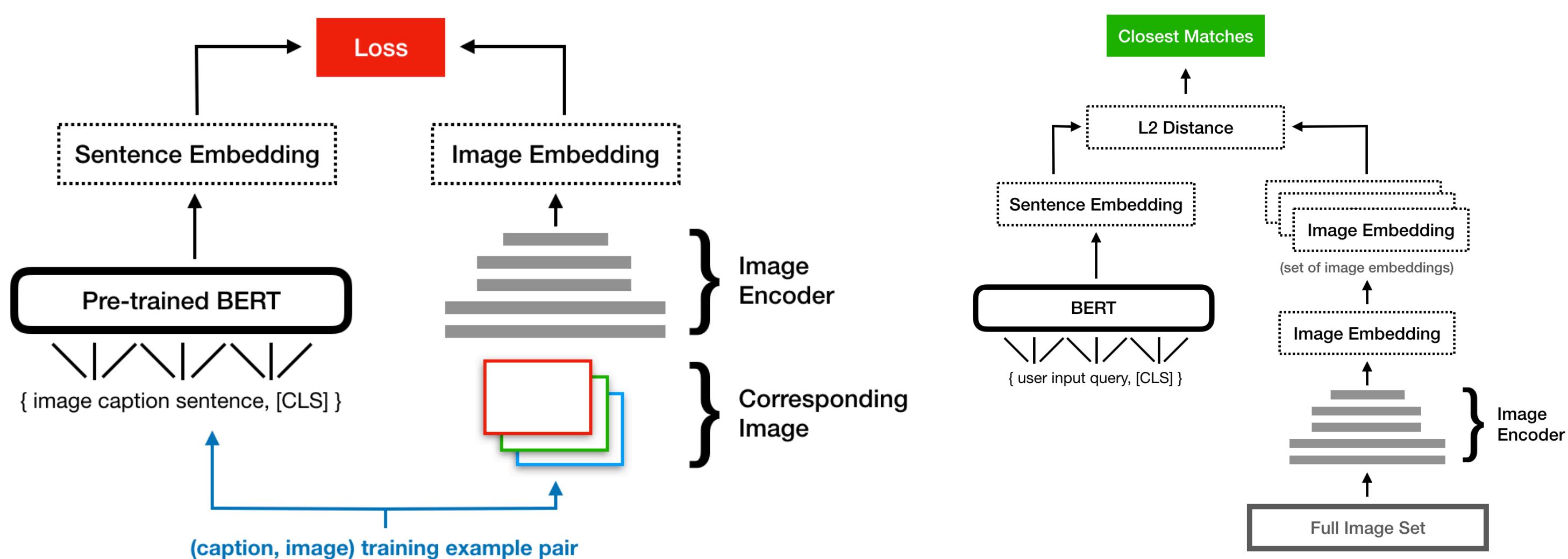
Adolfo Apolloni, Alex Raistrick, Omar Al-Ejel
EECS 442: Computer Vision

December 10, 2019

IMAGE-TEXT SEARCH

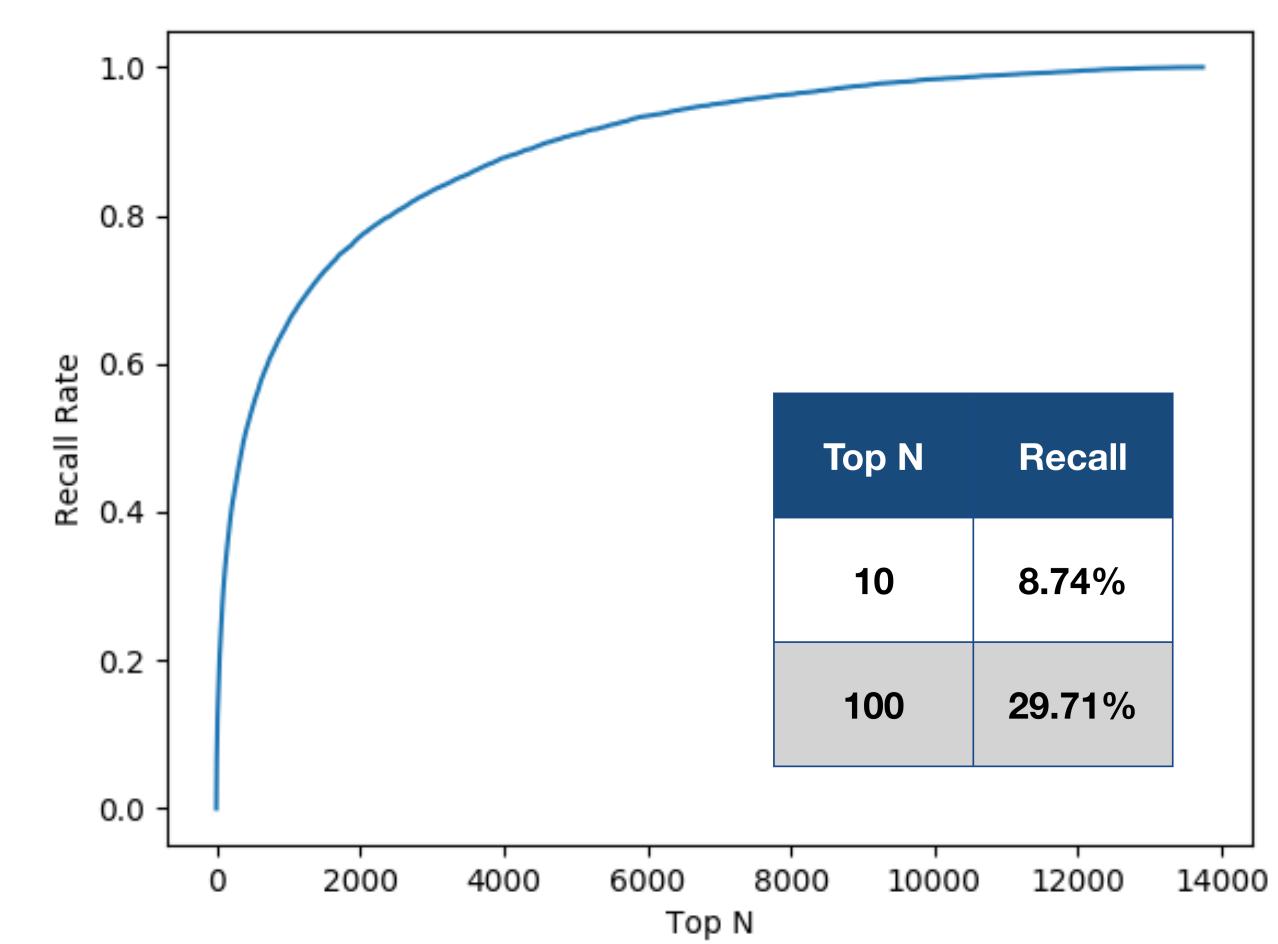
PROJECT MOTIVATION

Traditional image search relies on alt-text contextualization of image data. We developed a learned image search system that matches queries directly to image

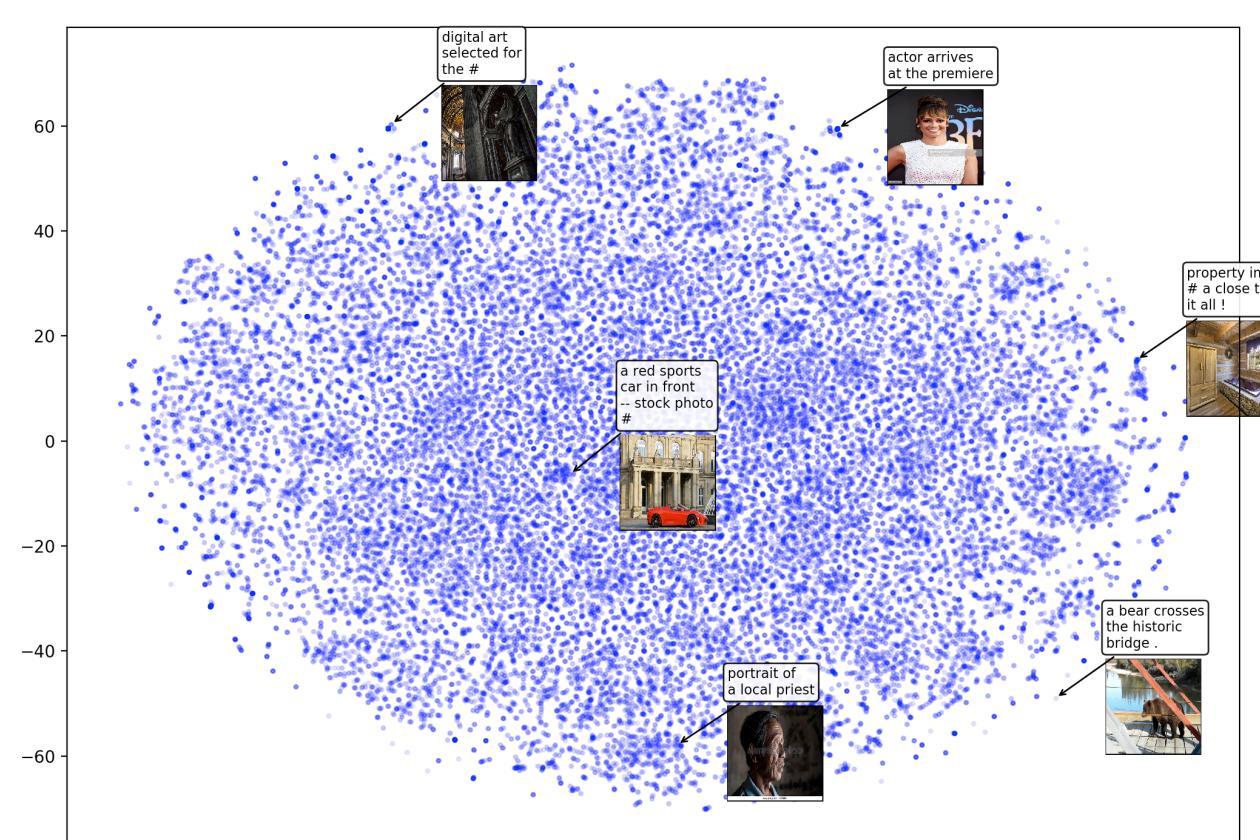


Training architecture for image embedding generation used for image search.

TOP-K ACCURACY CURVE

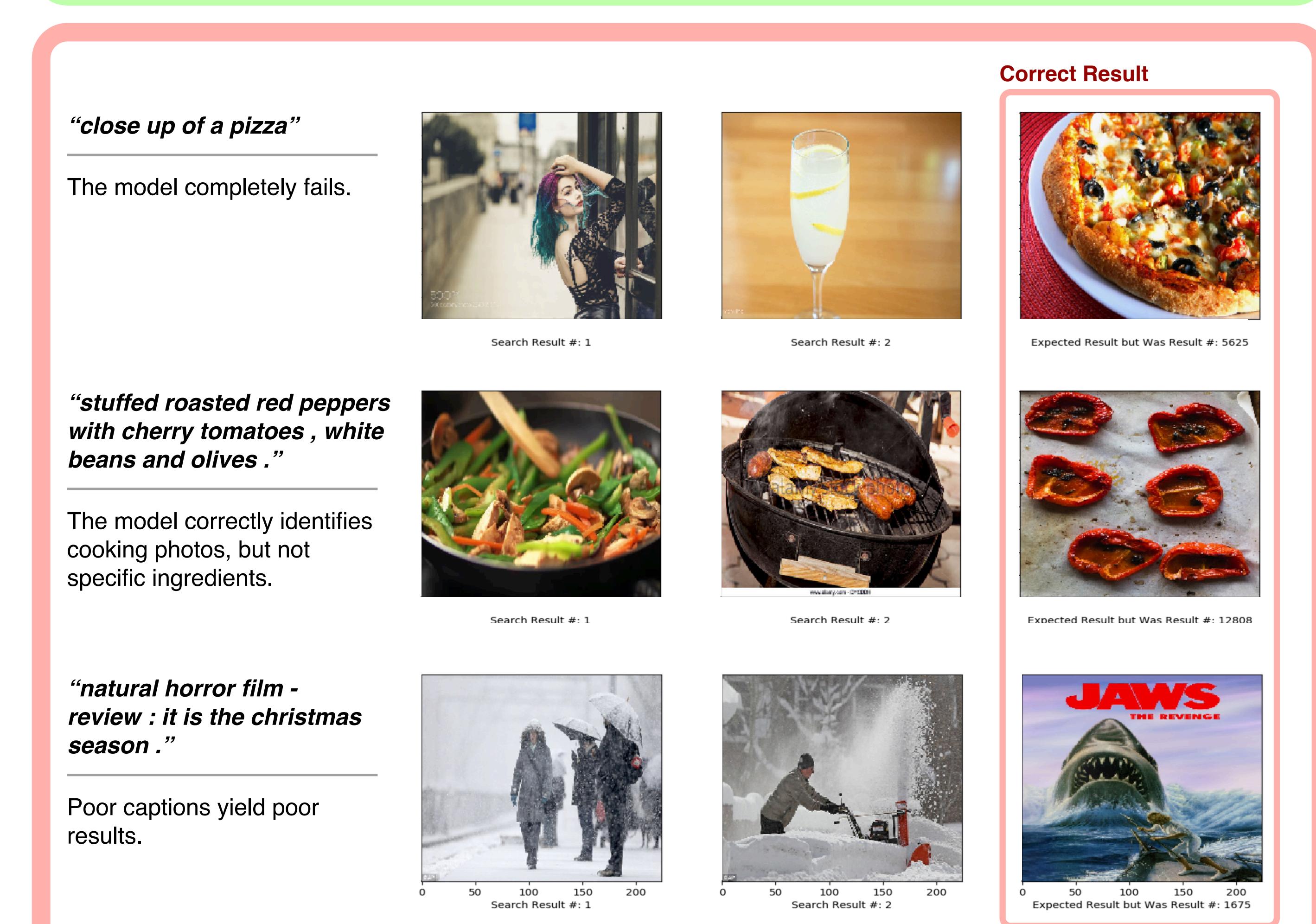
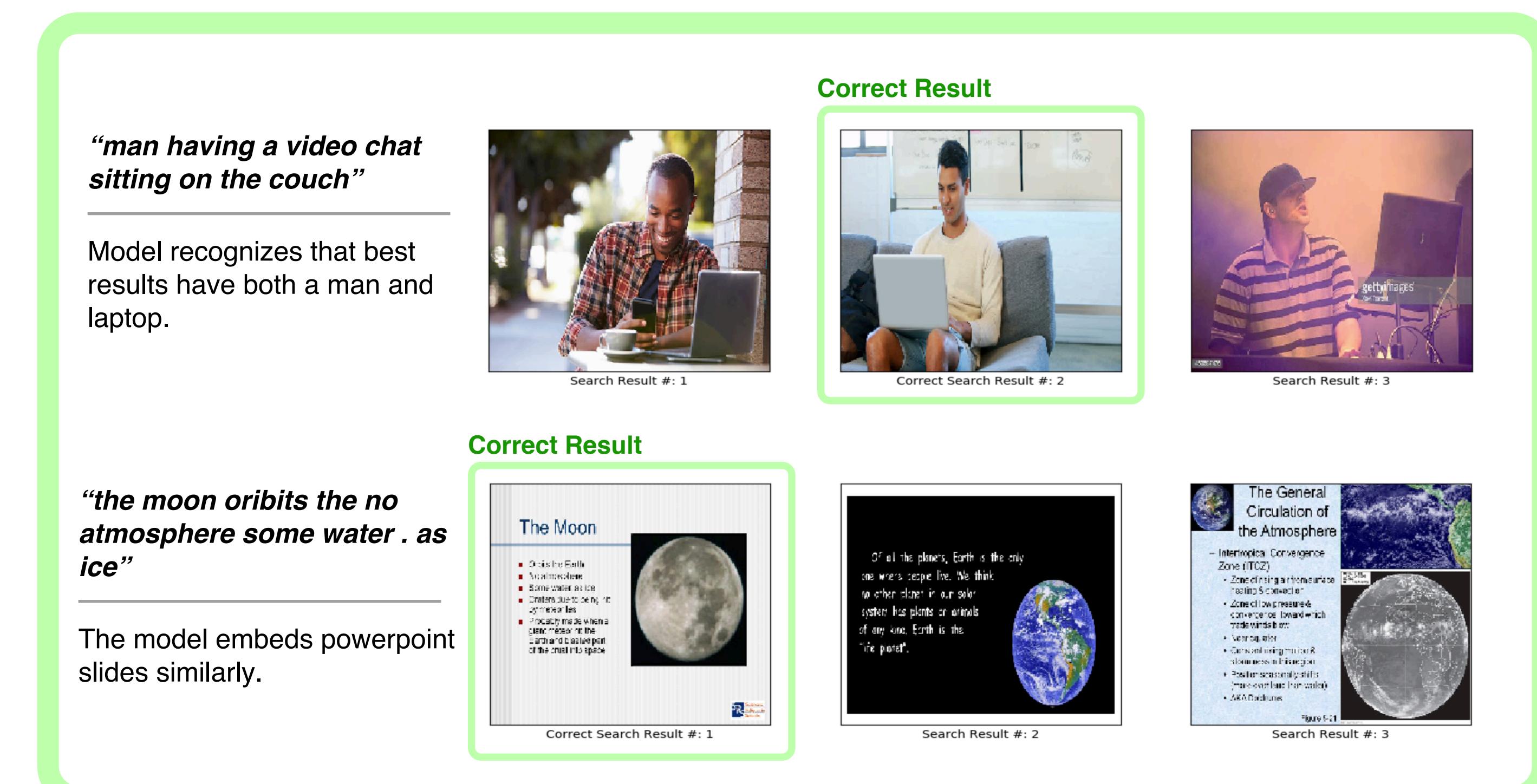


T-SNE



T-SNE chart displaying low-dimensional embeddings of image caption BERT vectors.

VALIDATION SET RESULTS

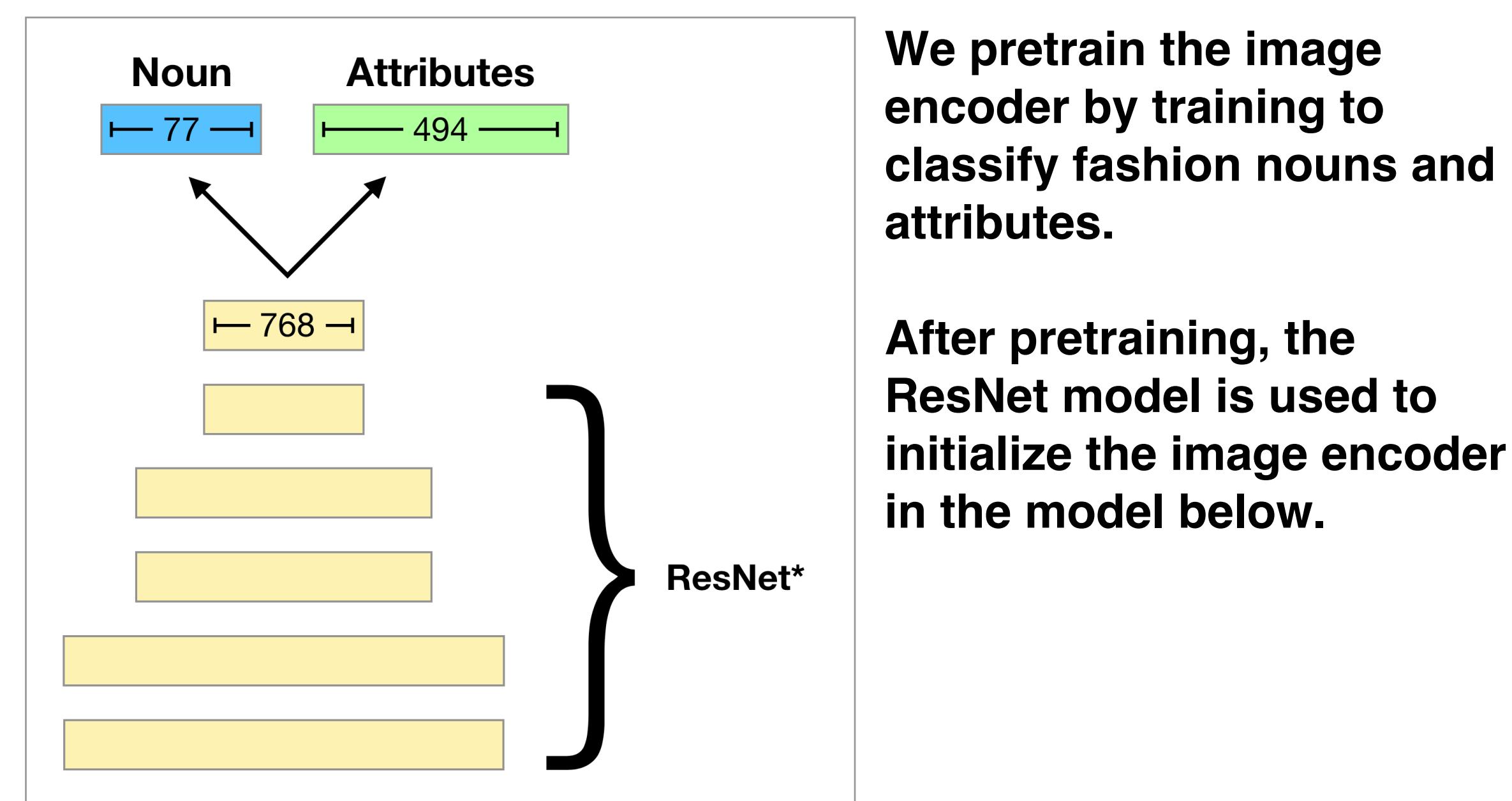


HYBRID IMAGE-TEXT MODEL

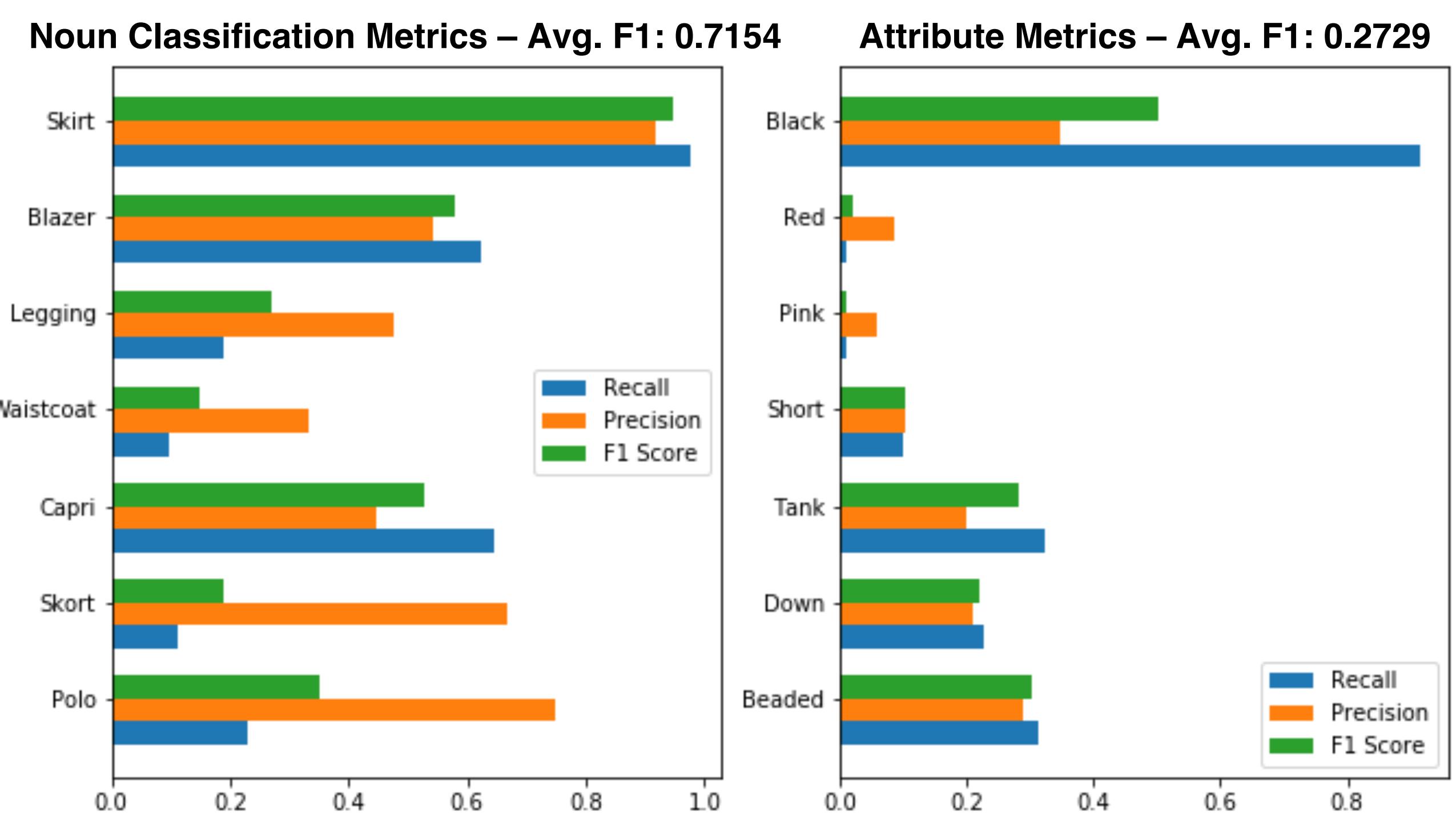
MOTIVATION

Past work composes fixed format caption and image data for search. We propose an adapted transformer architecture for learning on mixed sequences of image and text tokens to allow composite input.

IMAGE PRE-TRAINING ARCHITECTURE

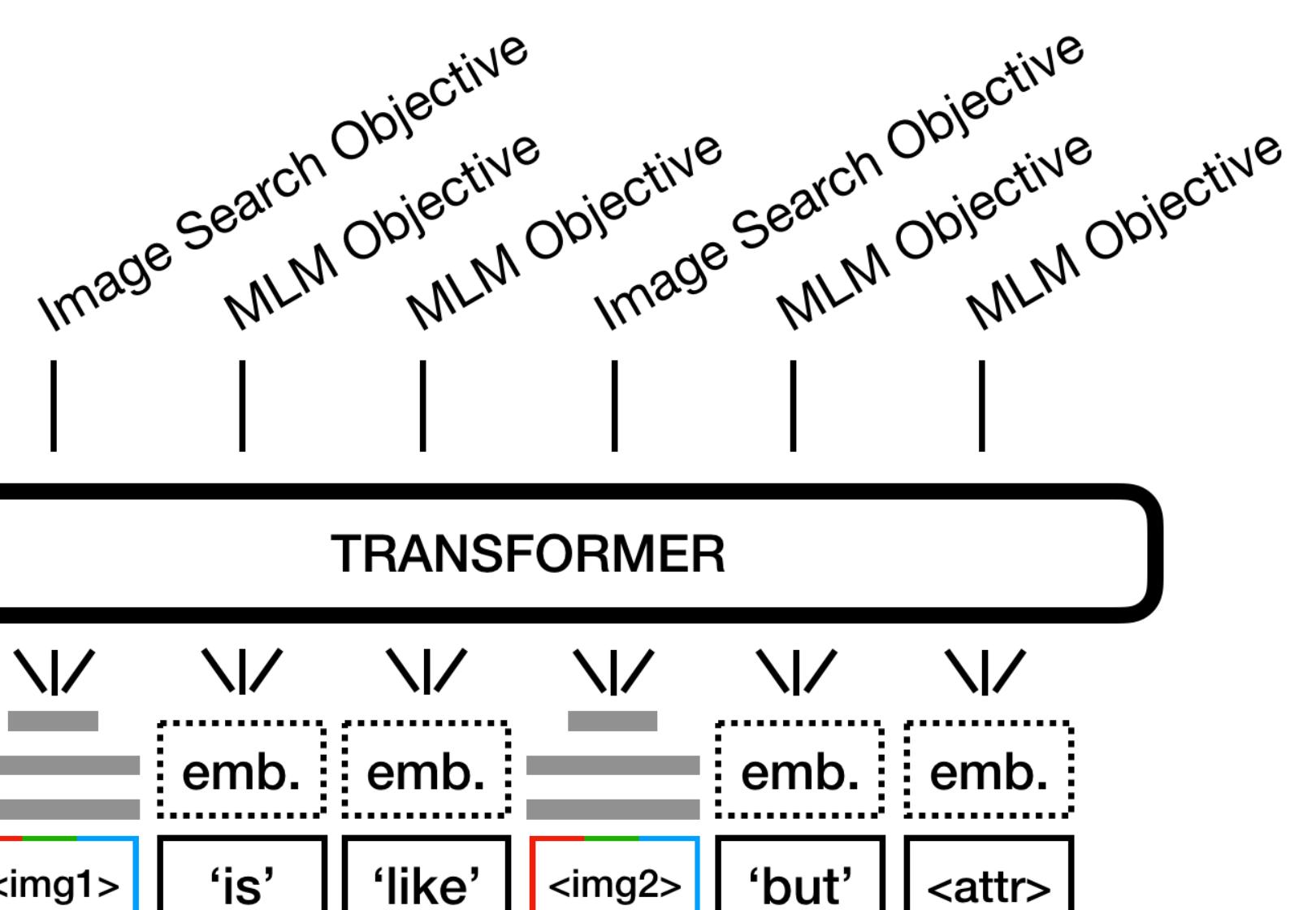


PRE-TRAINING RESULTS



Nouns (left), attributes (right) precision recall metrics showing Y-Net performance on ‘interesting’ examples. Note that performance is much better on nouns than attributes on average.

HYBRID MODEL ARCHITECTURE



BERT and Image Embedder Model used for image search from multimodal text and image input.

