



# Vision–Language Models

Jordan Boyd-Graber

University of Maryland

Introduction / Foundations

Slides from Mohit Iyyer, Vicente Ordonez, Fei-Fei Li, Justin Johnson, and Jacob Andreas

# Using Patches for Transformers

## AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>**

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

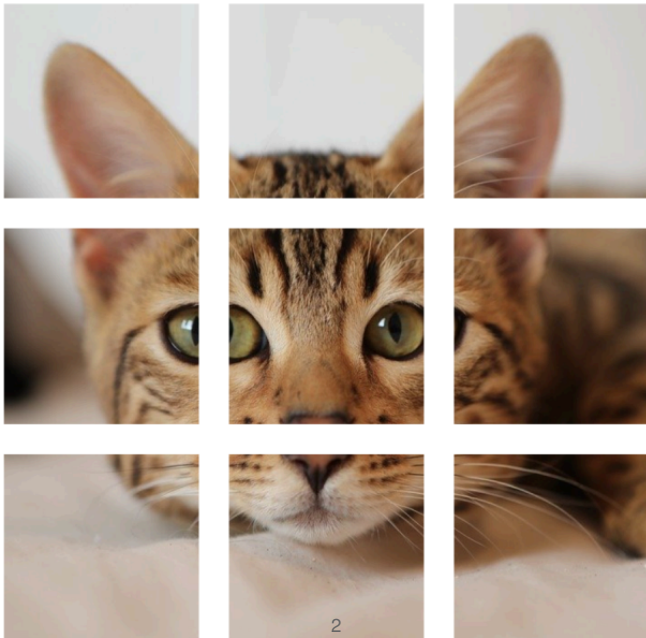
Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

## Using Patches for Transformers



## Using Patches for Transformers



# Using Patches for Transformers

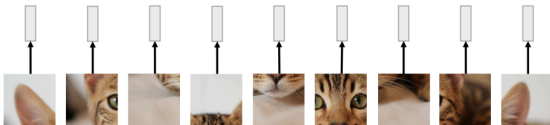
N input patches, each  
of shape 3x16x16



# Using Patches for Transformers

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16

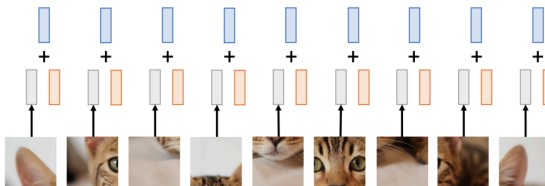


# Using Patches for Transformers

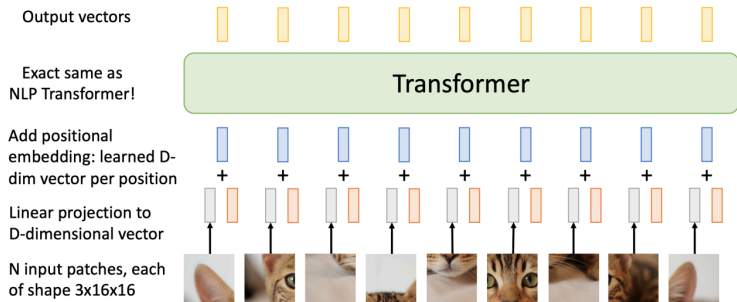
Add positional  
embedding: learned D-  
dim vector per position

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16

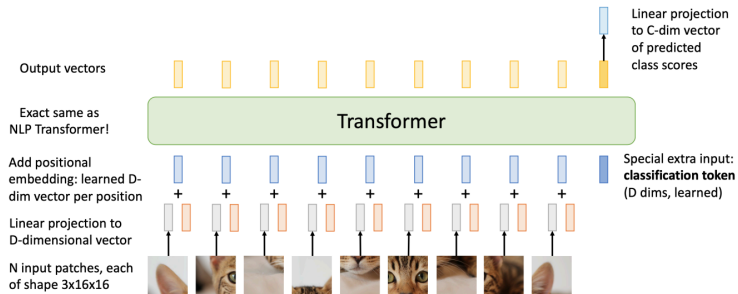


# Using Patches for Transformers





# Using Patches for Transformers



---

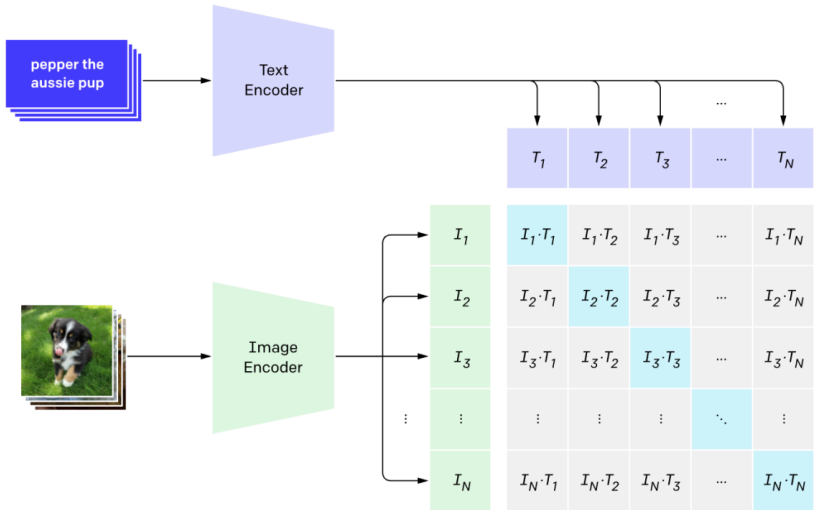
## Learning Transferable Visual Models From Natural Language Supervision

---

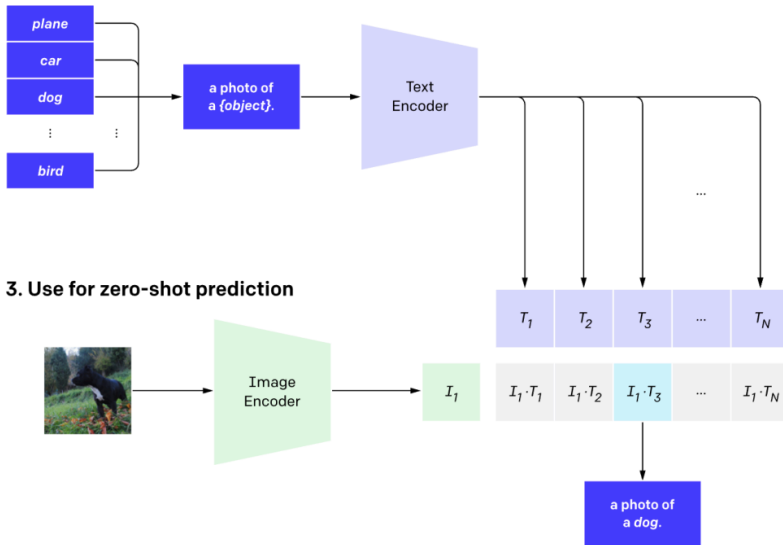
Alec Radford<sup>\*1</sup> Jong Wook Kim<sup>\*1</sup> Chris Hallacy<sup>1</sup> Aditya Ramesh<sup>1</sup> Gabriel Goh<sup>1</sup> Sandhini Agarwal<sup>1</sup>  
Girish Sastry<sup>1</sup> Amanda Askell<sup>1</sup> Pamela Mishkin<sup>1</sup> Jack Clark<sup>1</sup> Gretchen Krueger<sup>1</sup> Ilya Sutskever<sup>1</sup>

- OpenAI collect 400 million (image, text) pairs from the web
- Then, they train an image encoder and a text encoder with a simple contrastive loss: given a collection of images and text, predict which (image, text) pairs actually occurred in the dataset

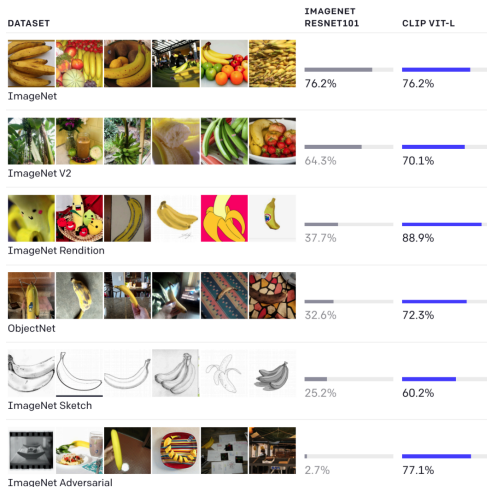
# Joint Training



# Joint Training



# Joint Training



## Generating text is one thing, but what about image generation?

- Could do autoregressive model pixel by pixel (people have tried)
- But better to learn higher-order structure