# Computer Vision for Embedded Systems

Yung-Hsiang Lu
Purdue University
yunglu@purdue.edu
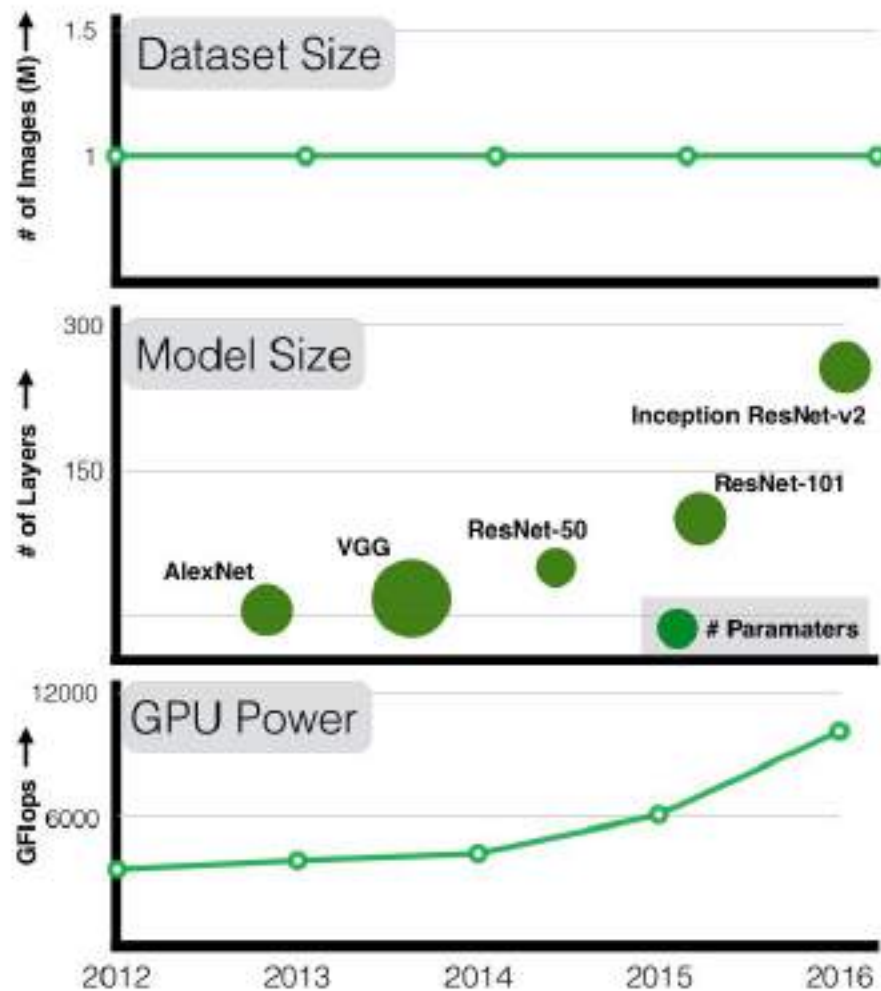
# Revisiting Unreasonable Effectiveness of Data in Deep Learning Era

Chen Sun, Abhinav Shrivastava, Saurabh Singh, Abhinav Gupta, ICCV 2017

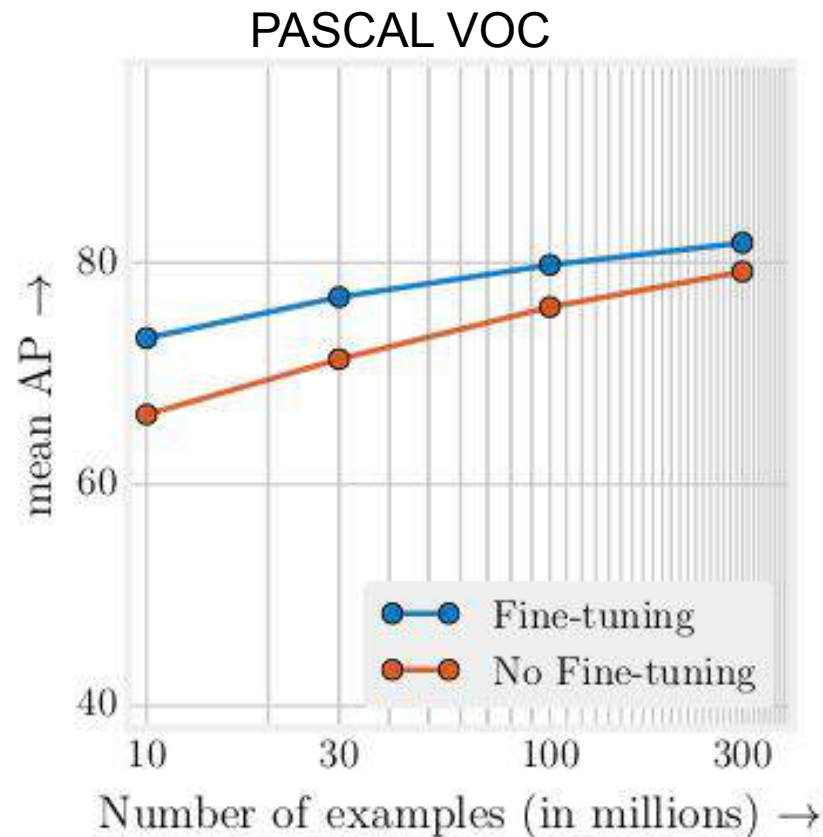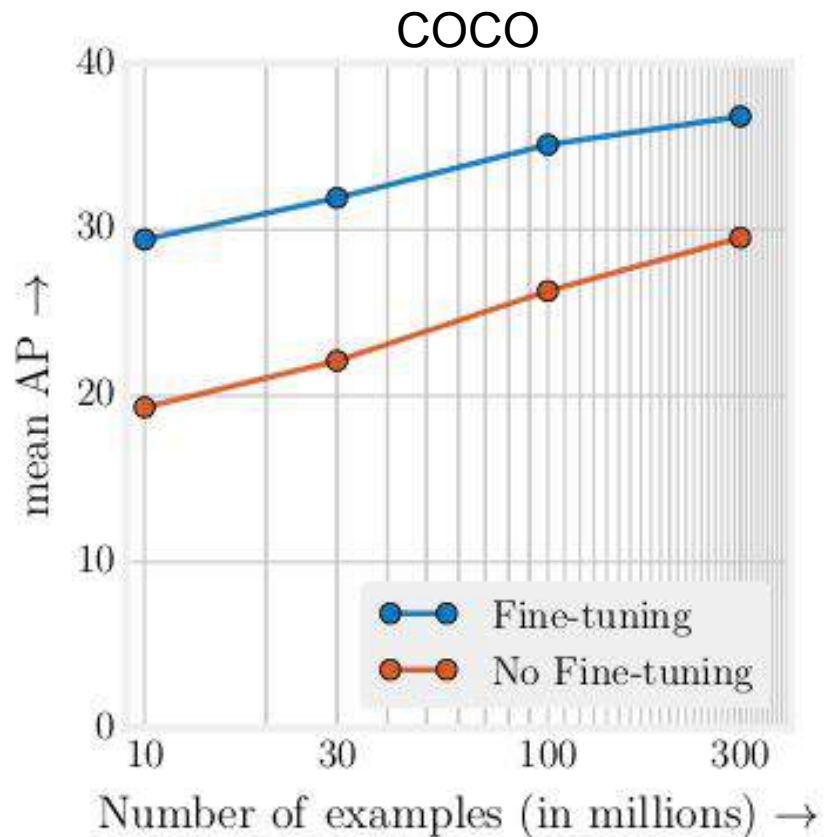## JFT dataset: 300M images, 18,291 categories

Dataset Size — # of Images (M)

Model Size — # of Layers

AlexNet, VGG, ResNet-50, ResNet-101, Inception ResNet-v2

# Paramaters

GPU Power — GFlops

2012  2013  2014  2015  2016

# JFT Dataset

- 300M images
- 375M labels
- 18,291 categories
  - 1,165 types of animals
  - 5,720 types of vehicles
  - maximum depth of hierarchy is 12
  - maximum number of children is 2,876
- heavy tail distribution: 3K categories with fewer than 100 images each
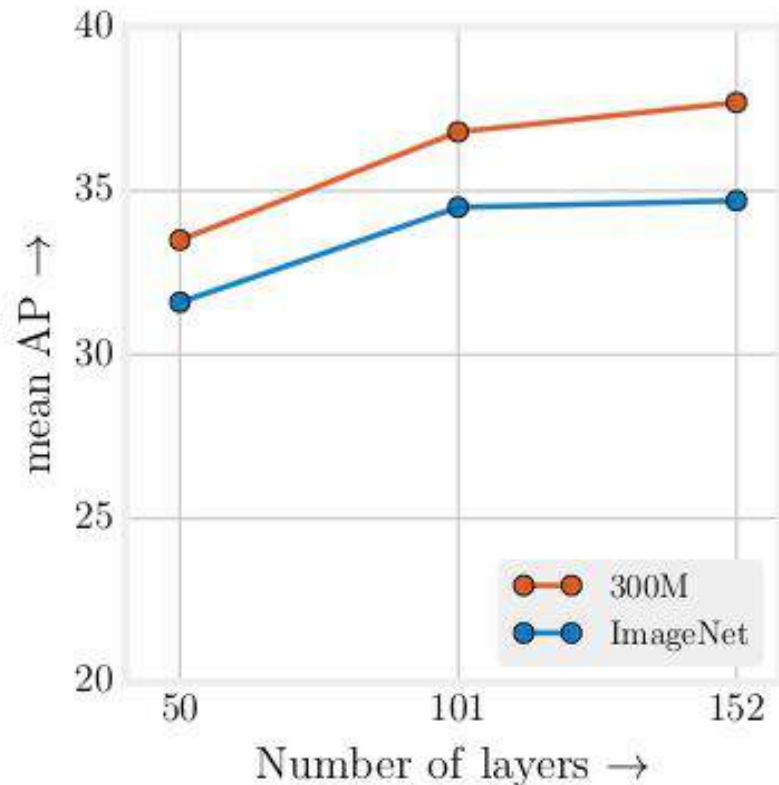- image sizes: 340 x 340 cropped to 299 x 299, normalized to [-1, 1]

# Effects of Training Examples



COCO

PASCAL VOC

# Effect of Model Capacity

COCO
on ResNet

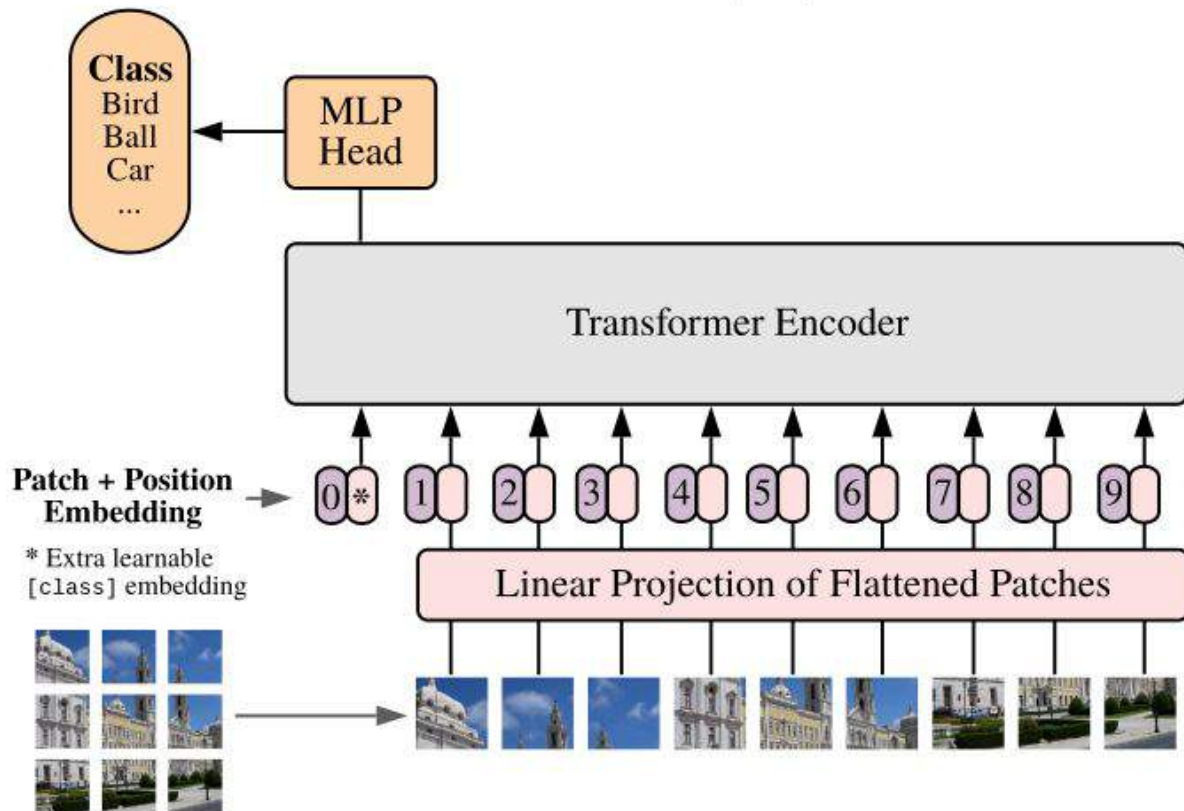| #Layers | ImageNet | 300M |
|---------|----------|------|
| 50      | 31.6     | 33.5 |
| 101     | 34.5     | 36.8 |
| 152     | 34.7     | 37.7 |

# Limitations of Convolutional Neural Networks

● Convolution considers neighbor pixels but at fixed distances
● Same parameters are applied to all pixels even though objects may have different sizes
● Hyperparameters (stride, filter size, number of layers …) determined in advance (may be determined by neural architecture search)
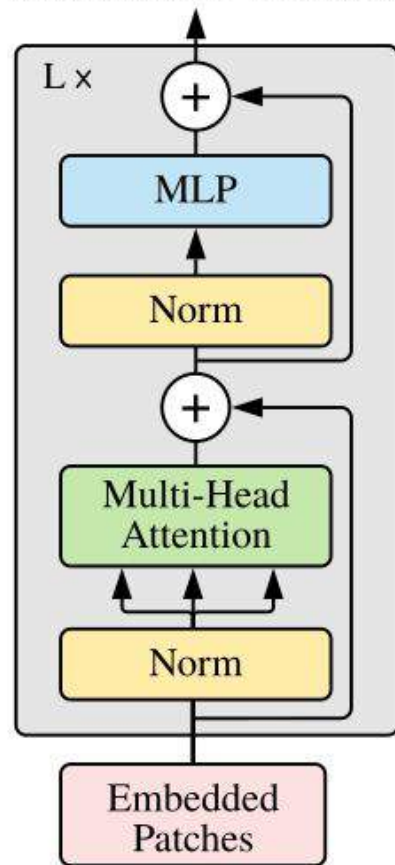
# An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby 2020

# Vision Transformer (ViT)

# Transformer Encoder

# Create Image Patches

$$R^{N \times W \times C} \Rightarrow R^{N \times (P^2 \times C)}$$

$H$ image height

$W$ image width

$C$ number of channels

$P$ size of patch

$$N = \frac{HW}{P^2}$$ number of patches

# Options of Position Embedding

- no position information (bags of patches)
- 1D position embedding (sequences of patches)
- 2D position embedding
- relative position embedding

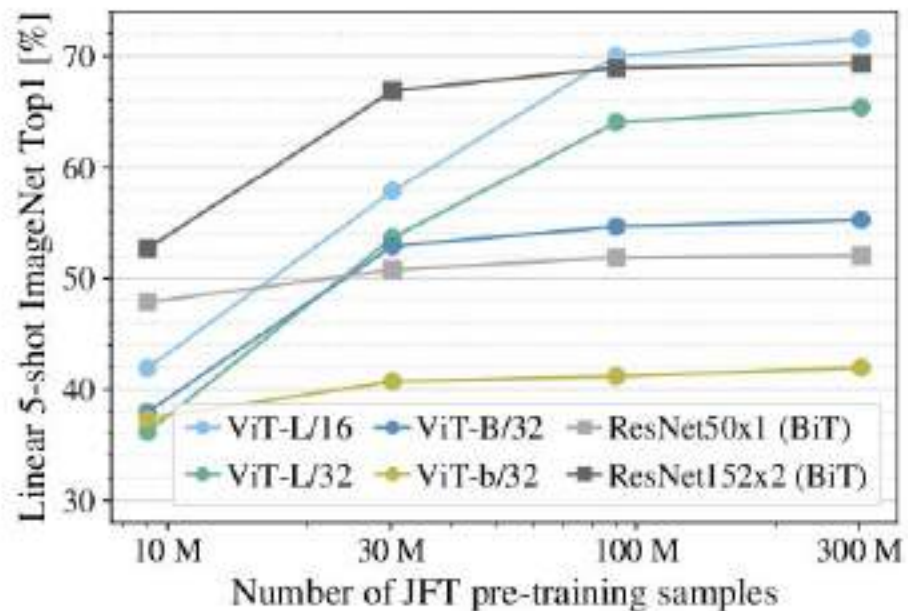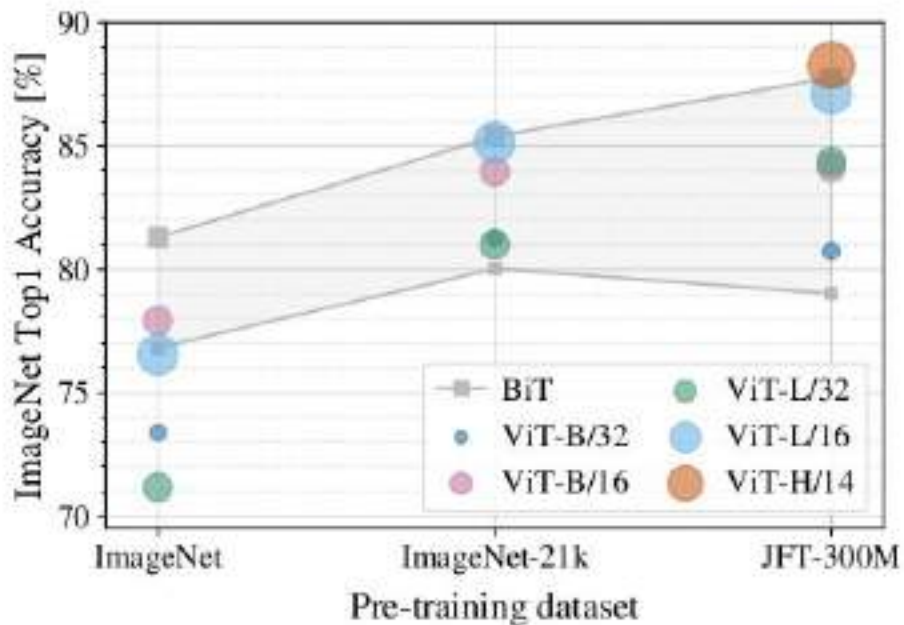| Pos. Emb. | Default/Stem | Every Layer | Every Layer-Shared |
|---|---|---|---|
| No Pos. Emb. | 0.61382 | N/A | N/A |
| 1-D Pos. Emb. | 0.64206 | 0.63964 | 0.64292 |
| 2-D Pos. Emb. | 0.64001 | 0.64046 | 0.64022 |
| Rel. Pos. Emb. | 0.64032 | N/A | N/A |

# Datasets

- ImageNet 1K: 1K classes, 1.3M images
- ImageNet 21K: 21K classes, 14M images
- JFT: 18K classes, 14M images
- CIFAR-10 and 100
- Oxford-IIIT Pets
- Oxford Flowers-102

# Model Variants

| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-------|--------|-----------------|----------|-------|--------|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

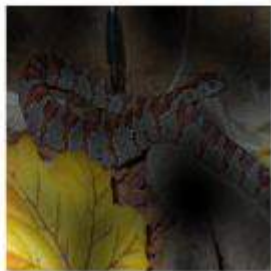# Comparison

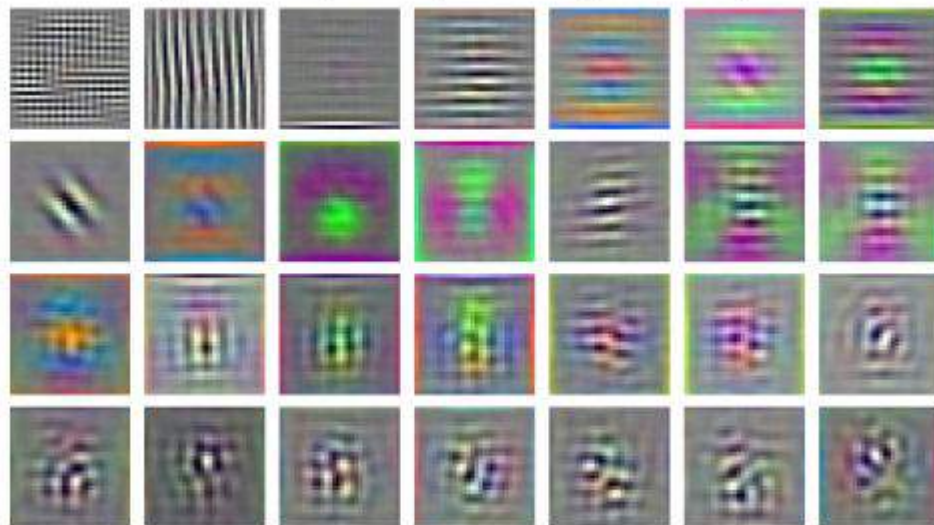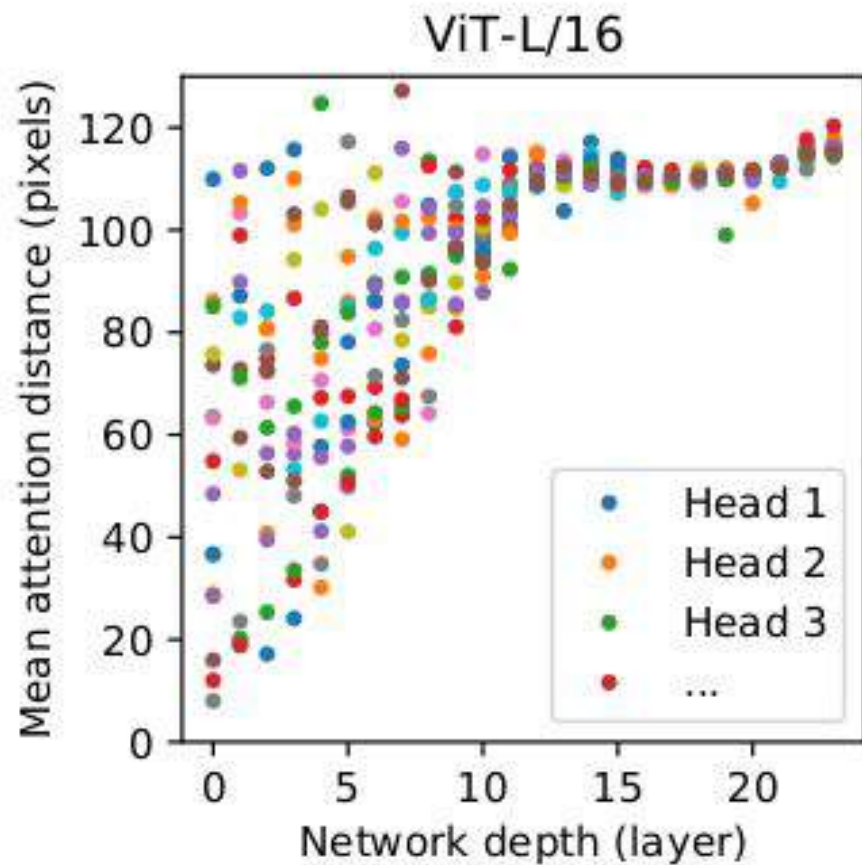| | Ours-JFT (ViT-H/14) | Ours-JFT (ViT-L/16) | Ours-I21k (ViT-L/16) | BiT-L (ResNet152x4) | Noisy Student (EfficientNet-L2) |
|---|---|---|---|---|---|
| ImageNet | $\mathbf{88.55} \pm 0.04$ | $87.76 \pm 0.03$ | $85.30 \pm 0.02$ | $87.54 \pm 0.02$ | $88.4/88.5^{*}$ |
| ImageNet ReaL | $\mathbf{90.72} \pm 0.05$ | $90.54 \pm 0.03$ | $88.62 \pm 0.05$ | $90.54$ | $90.55$ |
| CIFAR-10 | $\mathbf{99.50} \pm 0.06$ | $99.42 \pm 0.03$ | $99.15 \pm 0.03$ | $99.37 \pm 0.06$ | — |
| CIFAR-100 | $\mathbf{94.55} \pm 0.04$ | $93.90 \pm 0.05$ | $93.25 \pm 0.05$ | $93.51 \pm 0.08$ | — |
| Oxford-IIIT Pets | $\mathbf{97.56} \pm 0.03$ | $97.32 \pm 0.11$ | $94.67 \pm 0.15$ | $96.62 \pm 0.23$ | — |
| Oxford Flowers-102 | $99.68 \pm 0.02$ | $\mathbf{99.74} \pm 0.00$ | $99.61 \pm 0.02$ | $99.63 \pm 0.03$ | — |
| VTAB (19 tasks) | $\mathbf{77.63} \pm 0.23$ | $76.28 \pm 0.46$ | $72.72 \pm 0.21$ | $76.29 \pm 1.70$ | — |
| TPUv3-core-days | 2.5k | 0.68k | 0.23k | 9.9k | 12.3k |

# Input    Attention



# RGB embedding filters
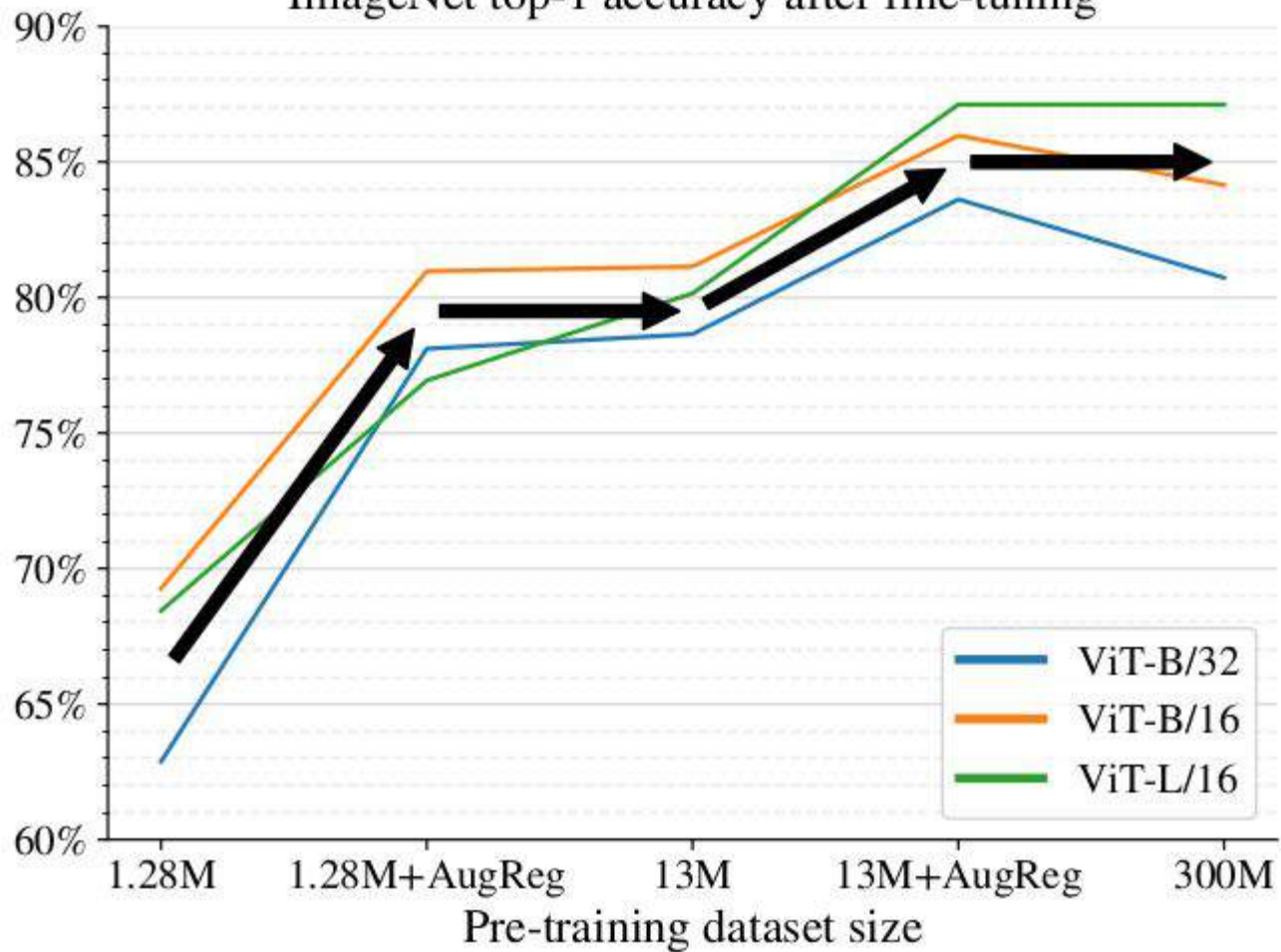(first 28 principal components)

ViT-L/16

# How to train your ViT? Data, Augmentation, and Regularization in Vision Transformers

Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, Lucas Beyer 2021

ImageNet top-1 accuracy after fine-tuning

Legend:
- ViT-B/32
- ViT-B/16
- ViT-L/16

Y-axis: 60%, 65%, 70%, 75%, 80%, 85%, 90%

X-axis (Pre-training dataset size): 1.28M, 1.28M+AugReg, 13M, 13M+AugReg, 300M

# ViViT: A Video Vision Transformer

Anurag Arnab, Mostafa Dehghani, Georg Heigold Chen Sun, CVPR 2021

# The Kinetics Human Action Video Dataset 2017

# Kinetics Dataset

- 400 human action classes
- each action at least 400 clips
- each clip 10 seconds from Youtube
- Single person activities: drawing, laughing, drinking
- Person-person activities: shaking hands, hugging
- Person-object activities: washing dishes, mowing lawn

Yung-Hsiang Lu, Purdue University

# Crowdsourcing to label data



Can you see a 👤 human performing the action **riding mule?**

## Instructions

We would like to find videos that contain real humans performing actions e.g. scrubbing their face, jumping, kissing someone etc.

Please click on the most appropriate button after watching each video:

👍 Yes, this contains a true example of the action

👎 No, this does not contain an example of the action
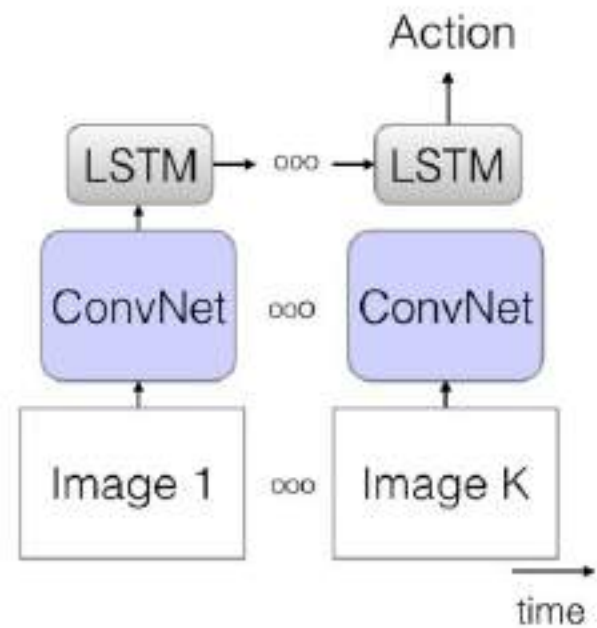
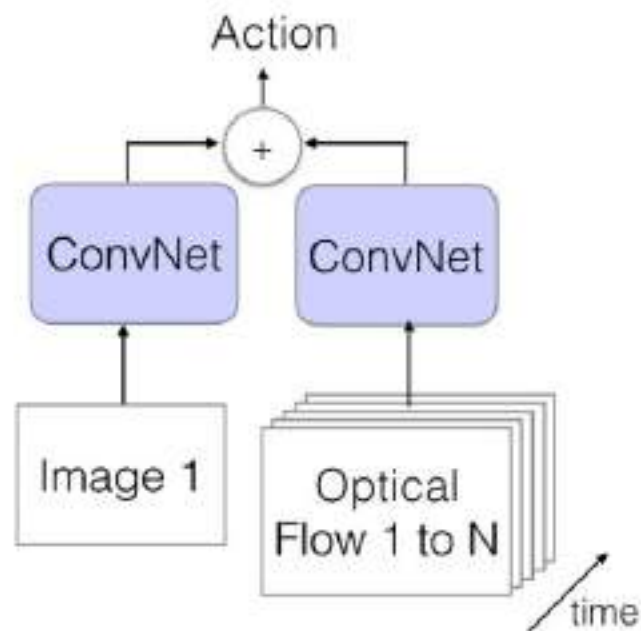❓ You are unsure if there is an example of the action

🔄 Replay the video

⚠️ Video does not play, does not contain a human, is an image, cartoon or a computer game.

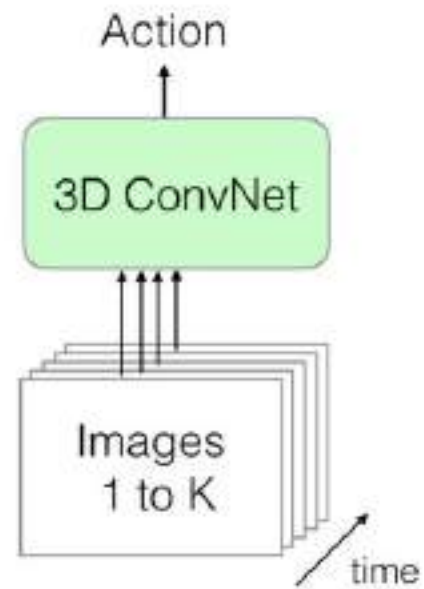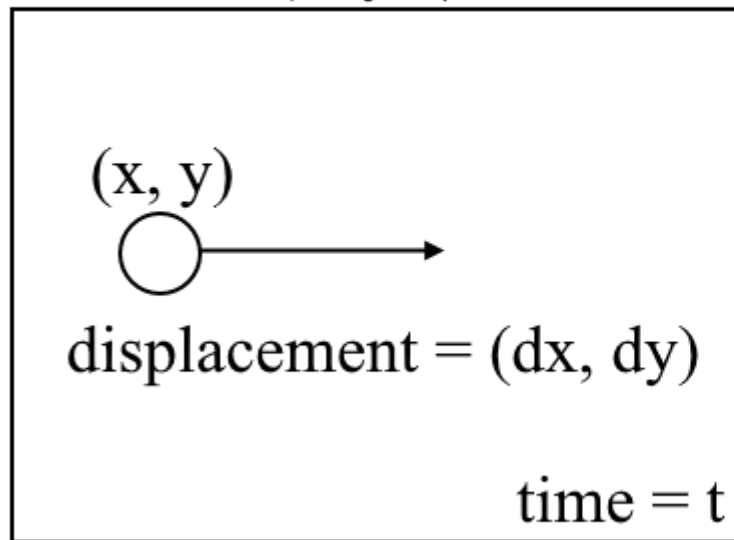🔇 We have turned off the audio, you need to judge the clip using the visuals only.
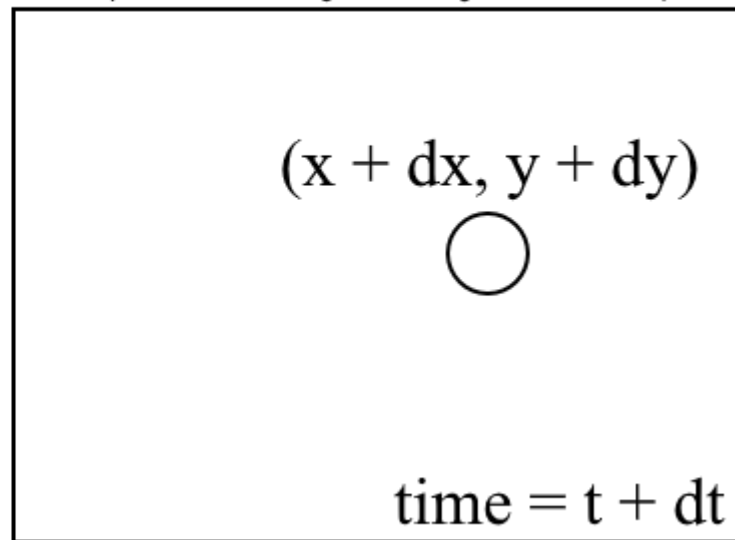
a) LSTM  b) Two-Stream  c) 3D ConvNet

I(x, y, t)

I(x + dx, y + dy, t + dt)

(x, y)

○ ——————→

displacement = (dx, dy)

time = t

(x + dx, y + dy)

○

time = t + dt

$$I(x, y, t) = I(x + \delta x, y + \delta y, t + \delta t)$$

https://nanonets.com/blog/optical-flow/

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) + \frac{\partial I}{\partial x}\delta x + \frac{\partial I}{\partial y}\delta y + \frac{\partial I}{\partial t}\delta t + \ldots$$

$$\Rightarrow \frac{\partial I}{\partial x}\delta x + \frac{\partial I}{\partial y}\delta y + \frac{\partial I}{\partial t}\delta t = 0$$

$$\frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \frac{\partial I}{\partial t} = 0$$

$$u = \frac{dx}{dt} \qquad v = \frac{dy}{dt}$$

https://nanonets.com/blog/optical-flow/

# Lucas–Kanade method

The method assume the optical flow in each small neighborhood is unchanged.

$$I_x(q_1)V_x + I_y(q_1)V_y = -I_t(q_1)$$

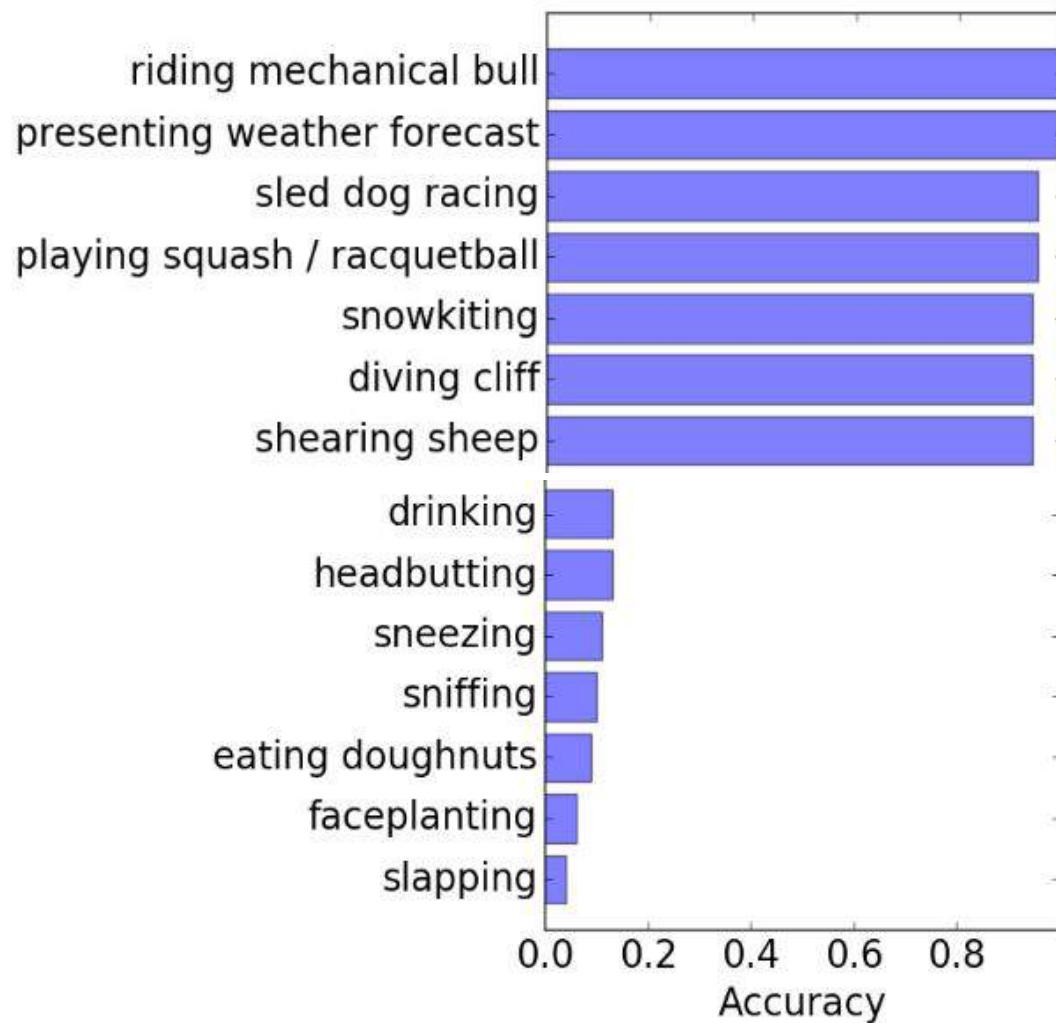$$I_x(q_2)V_x + I_y(q_2)V_y = -I_t(q_2)$$

$$\vdots$$

$$I_x(q_n)V_x + I_y(q_n)V_y = -I_t(q_n)$$

$$V_x = u(\text{in the previous slide})$$
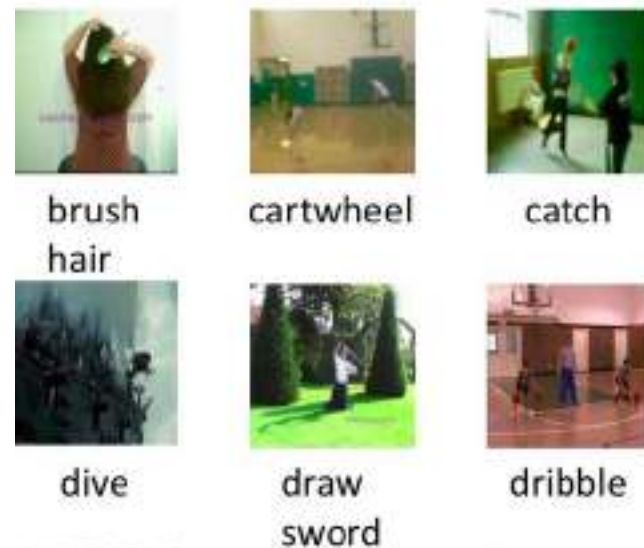
https://en.wikipedia.org/wiki/Lucas%E2%80%93Kanade_method

| Class 1 | Class 2 | confusion |
|---|---|---|
| 'riding mule' | 'riding or walking with horse' | 40% |
| 'hockey stop' | 'ice skating' | 36% |
| 'swing dancing' | 'salsa dancing' | 36% |
| 'strumming guitar' | 'playing guitar' | 35% |
| 'shooting basketball' | 'playing basketball' | 32% |
| 'cooking sausages' | 'cooking chicken' | 29% |
| 'sweeping floor' | 'mopping floor' | 27% |
| 'triple jump' | 'long jump' | 26% |
| 'doing aerobics' | 'zumba' | 26% |
| 'petting animal (not cat)' | 'feeding goats' | 25% |
| 'shaving legs' | 'waxing legs' | 25% |
| 'snowboarding' | 'skiing (not slalom or crosscountry)' | 22% |

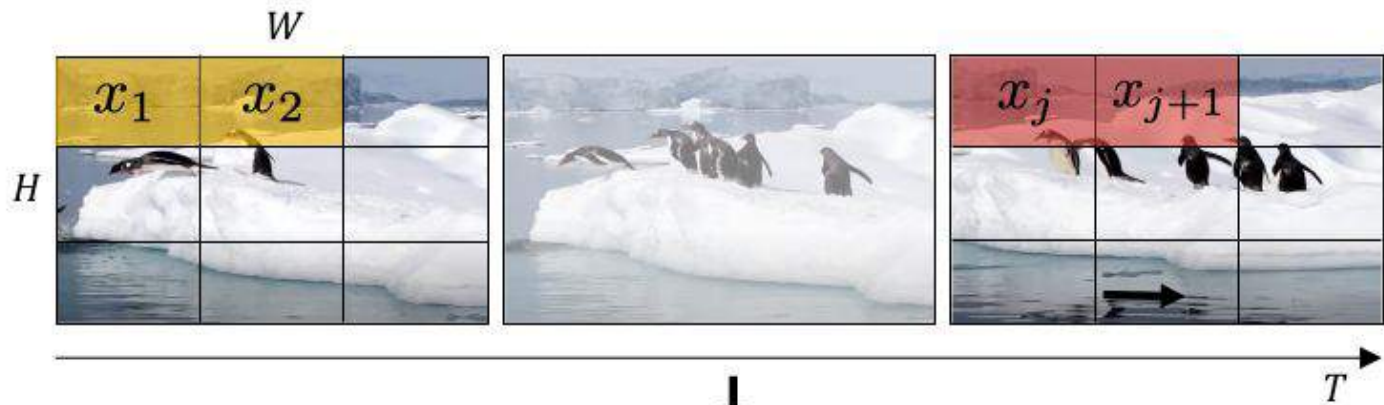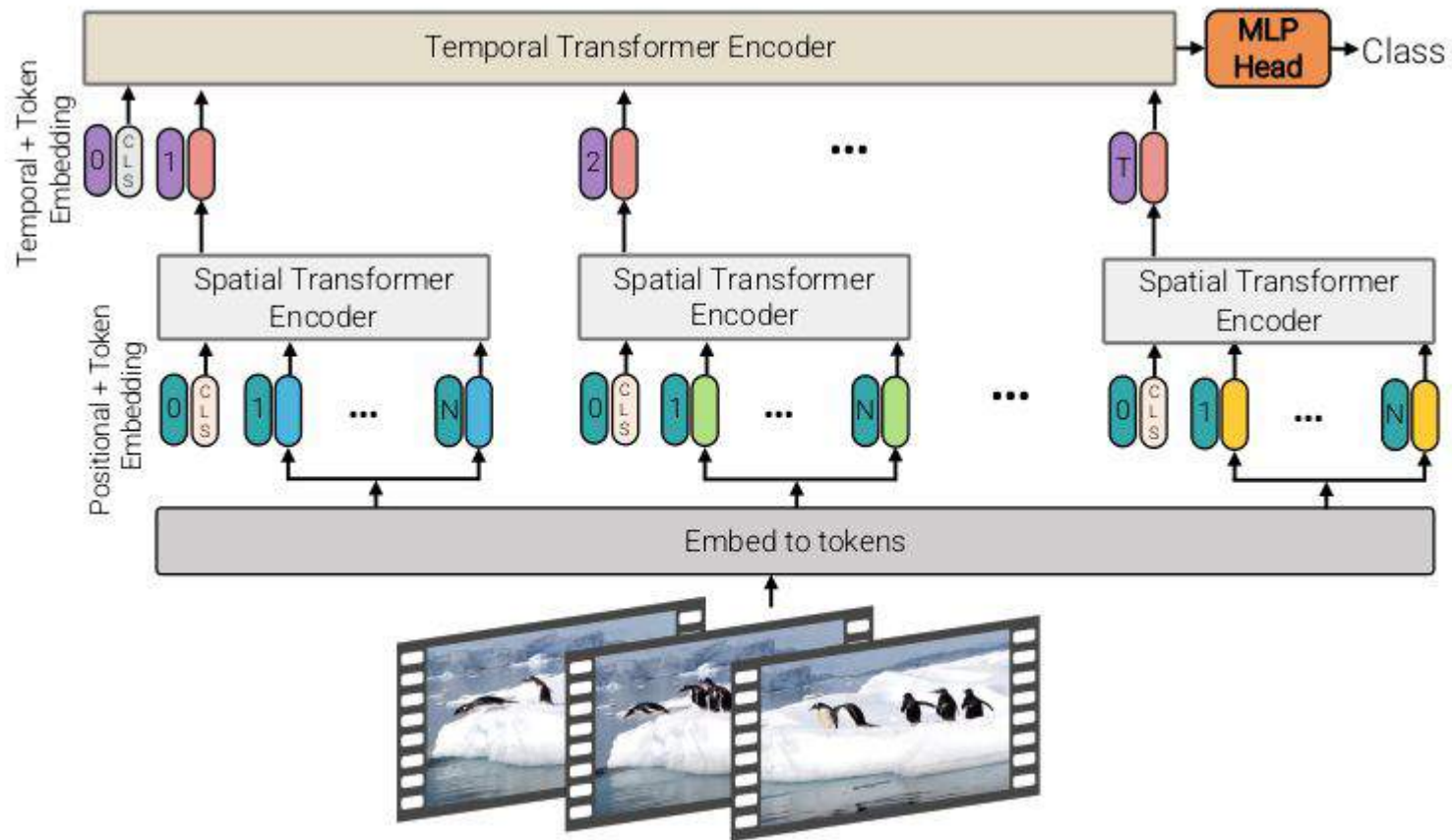| Architecture | UCF-101 | | | HMDB-51 | | | Kinetics | | |
|---|---|---|---|---|---|---|---|---|---|
| | RGB | Flow | RGB+Flow | RGB | Flow | RGB+Flow | RGB | Flow | RGB+Flow |
| (a) ConvNet+LSTM | 84.3 | – | – | 43.9 | – | – | 57.0 / 79.0 | – | – |
| (b) Two-Stream | 84.2 | 85.9 | 92.5 | 51.0 | 56.9 | 63.7 | 56.0 / 77.3 | 49.5 / 71.9 | 61.0 / 81.3 |
| (c) 3D-ConvNet | 51.6 | – | – | 24.3 | – | – | 56.1 / 79.5 | – | – |



UCF-101



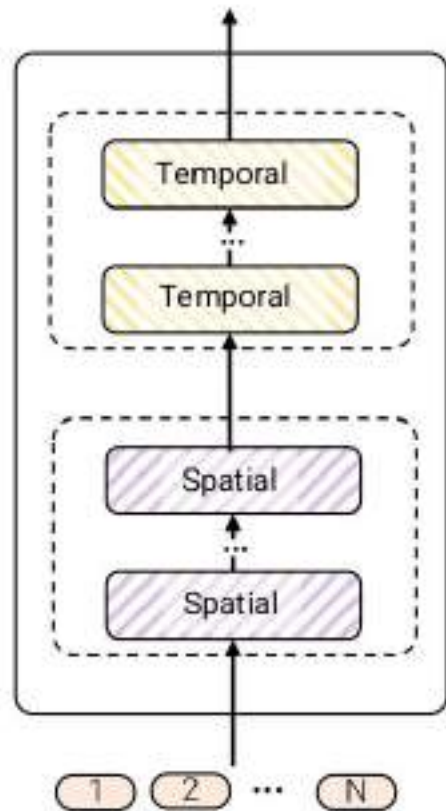HMDB-51

# ViViT: A Video Vision Transformer

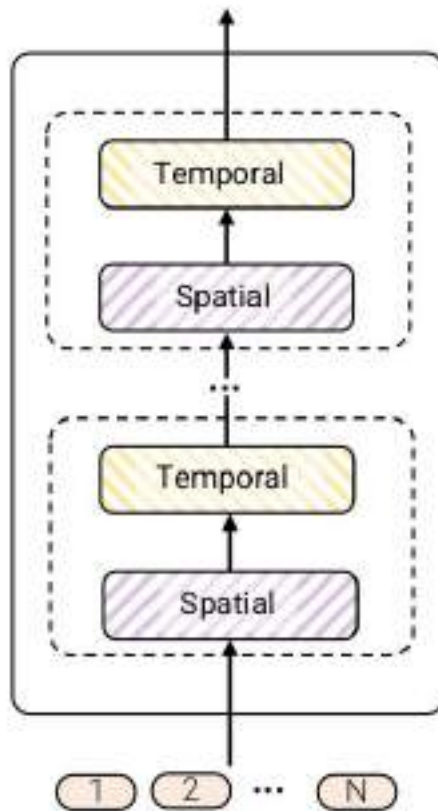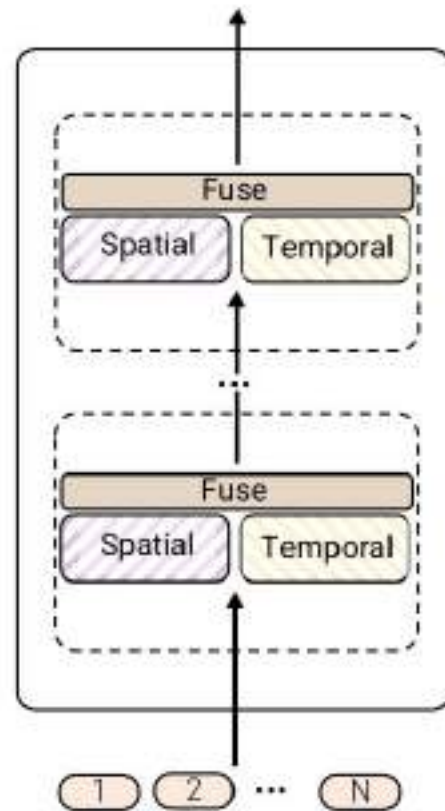Anurag Arnab, Mostafa Dehghani, Georg Heigold Chen Sun, CVPR 2021

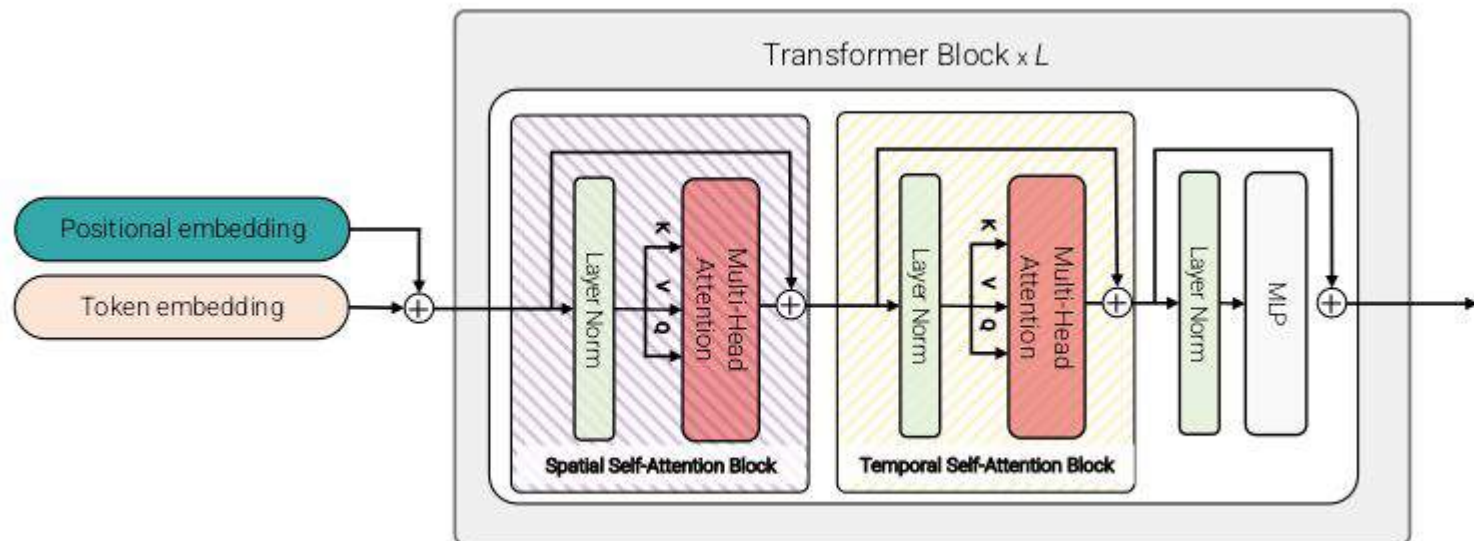Factorised Encoder — Factorised Self-Attention — Factorised Dot-Product

**Table 2:** Comparison of model architectures using ViViT-B as the backbone, and tubelet size of $16 \times 2$. We report Top-1 accuracy on Kinetics 400 (K400) and action accuracy on Epic Kitchens (EK). Runtime is during inference on a TPU-v3.

| | K400 | EK | FLOPs $(\times 10^9)$ | Params $(\times 10^6)$ | Runtime (ms) |
|---|---|---|---|---|---|
| Model 1: Spatio-temporal | 80.0 | 43.1 | 455.2 | 88.9 | 58.9 |
| Model 2: Fact. encoder | 78.8 | 43.7 | 284.4 | 115.1 | 17.4 |
| Model 3: Fact. self-attention | 77.4 | 39.1 | 372.3 | 117.3 | 31.7 |
| Model 4: Fact. dot product | 76.3 | 39.5 | 277.1 | 88.9 | 22.9 |
| Model 2: Ave. pool baseline | 75.8 | 38.8 | 283.9 | 86.7 | 17.3 |

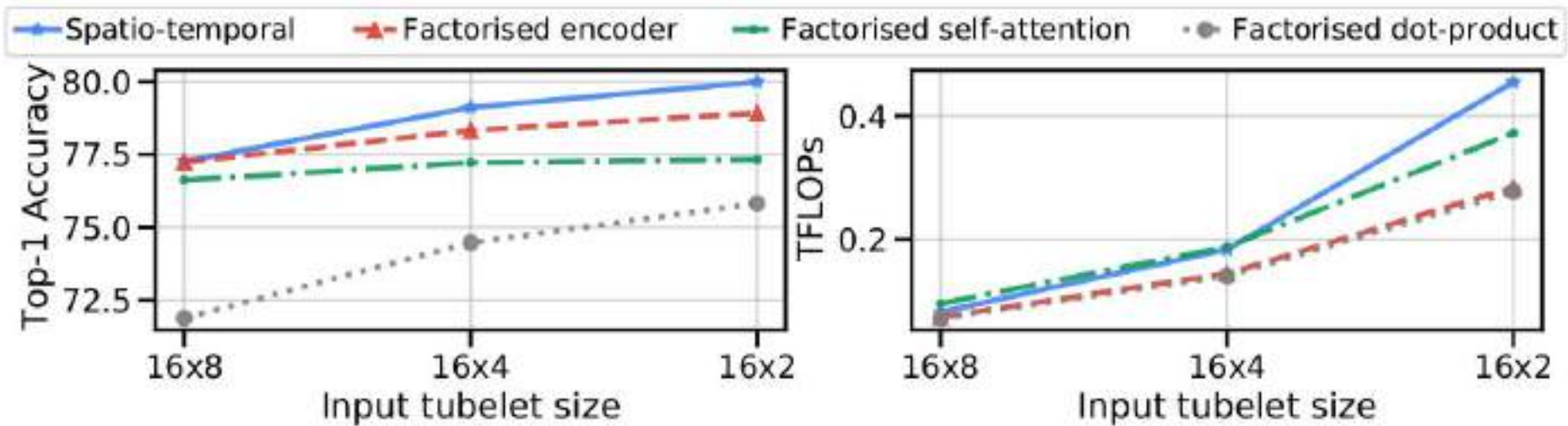Legend: Spatio-temporal · Factorised encoder · Factorised self-attention · Factorised dot-product

(a) Accuracy

(b) Compute