# A picture is worth a thousand words

Can you write them?
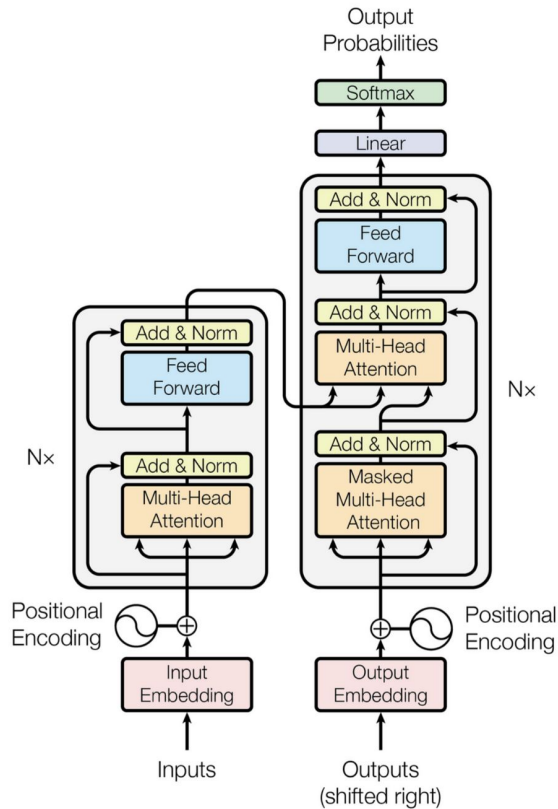
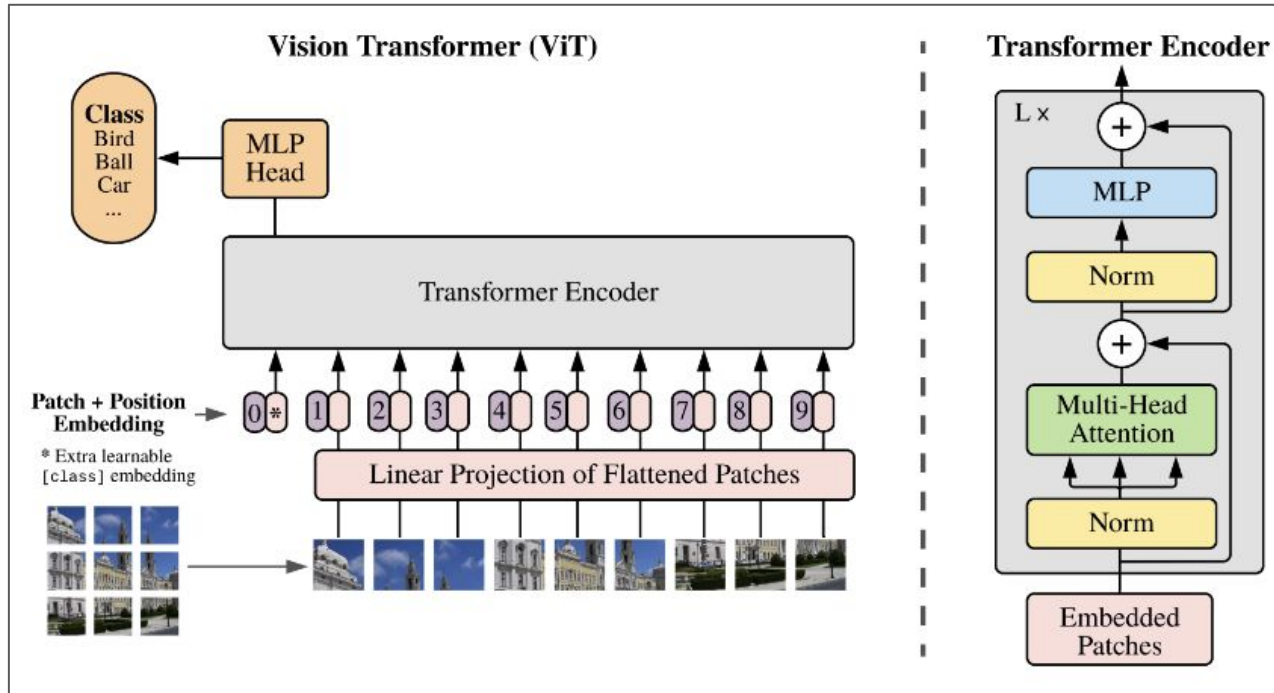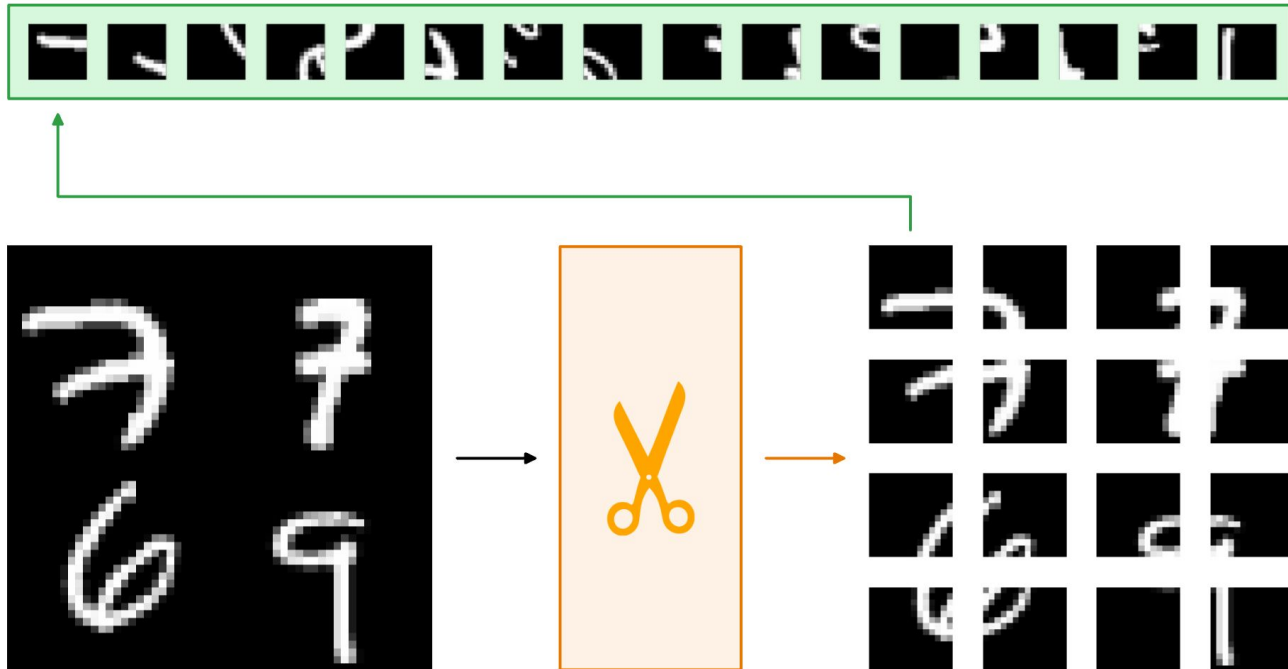# Task

# Attention Is All You Need

# Vision Transformer (ViT)

# Prep

# Patch Projection



(1×64)

nn.Linear(196, 64)

64
12,544
196

(196×64)

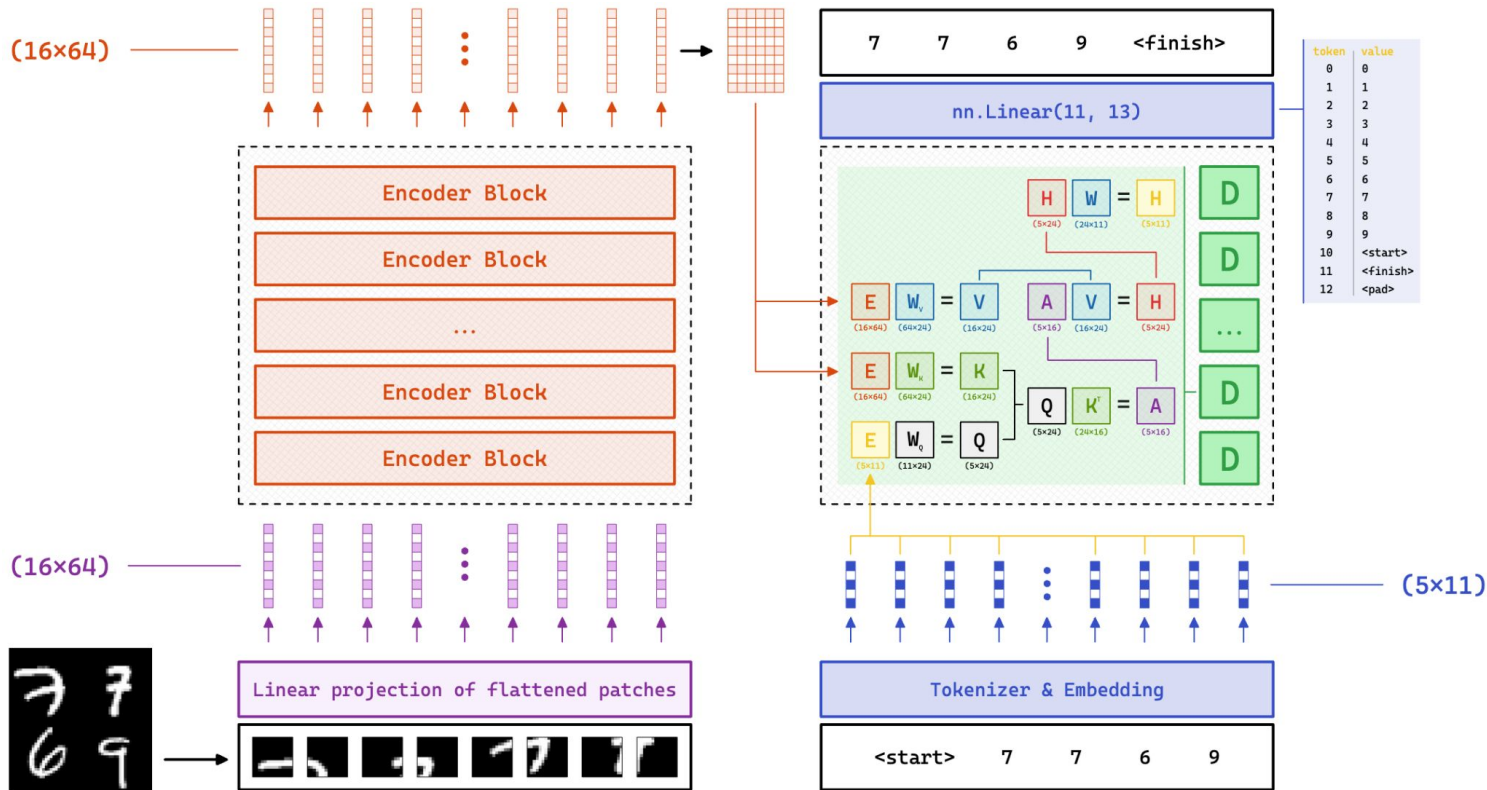| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | 189 | 190 | 191 | 192 | 193 | 194 | 195 | 196 |

(1×196)

(14×14)    (14×14)    (14×14)

(16×64)

Encoder Block

Encoder Block

...

Encoder Block

Encoder Block

(16×64)

Linear projection of flattened patches

x16

# Encoders



Encoded Representations

Encoder Block
Encoder Block
Encoder
Encoder Block
Encoder Block

Embedding Layer

Feature Extractors

| Tokenizer & Embedding | Linear projection of flattened patches | Convolution 1D over the Y axis | Spatiotemporal Feature Aggregation |
|---|---|---|---|
| Lorem ipsum dolor sit amet ... | | | |

# Encoder

# Transformer

# Papers

Vaswani et al. (NeurIPS 2017)

Dosovitskiy et al. (ICLR 2021)

# **Suggestions**

- do not just copy-paste from ChatGPT

- understand and practice PyTorch ops
- use google colab for little snippets
- pair programming, swap pairs within the team
- watch tutorials but make sure to talk
- no need for GPUs yet, do everything local

# Good luck!

## Recurrent Rebels

Adam Beedell

Aparna Pillai

Helen Zhou

Esperanza Shi

## Gradient Gigglers

Clement Ha

Andrew

Umut Sagir

Jacob Jenner

## Overfitting Overlords

Nikolas Kuhn

David Edev

Peter O'Keeffe

Miguel Parracho

## Hyperparameter Hippies

Tyrone Nicholas

Dan Goss

Anton Dergunov

Ben Liong

## Perceptron Party

Tao Zamorano

Tomas Krajcoviech

James Carter

Charles Cai

## Backprop Bunch

Ethan Edwards

Melanie Wong

Kadriye Turkcan

Ben Williams

## Dropout Disco

James Yan

Andrei Zhirnov

Prima Gouse

Jingyan Chen

## Kernel Kittens

Joao Esteves

Hikaru Tsujimura

Rosh Beed

Felipe Lavratti

## Bayesian Buccaneers

Rasched Haidari

Ewan Beattie

Marcin Tolysz

Maria Sharif

## Feature Fiestas

Ben Bethell

Arjuna James

Yali Pan

Ayman Abbas

FOUNDERS AND CODERS

Machine
Learning
Institute