

# **Visual Learning (Tutorial C3) Transformers and Beyond**

Instructor - Simon Lucey

**Maths of AI - 2024**



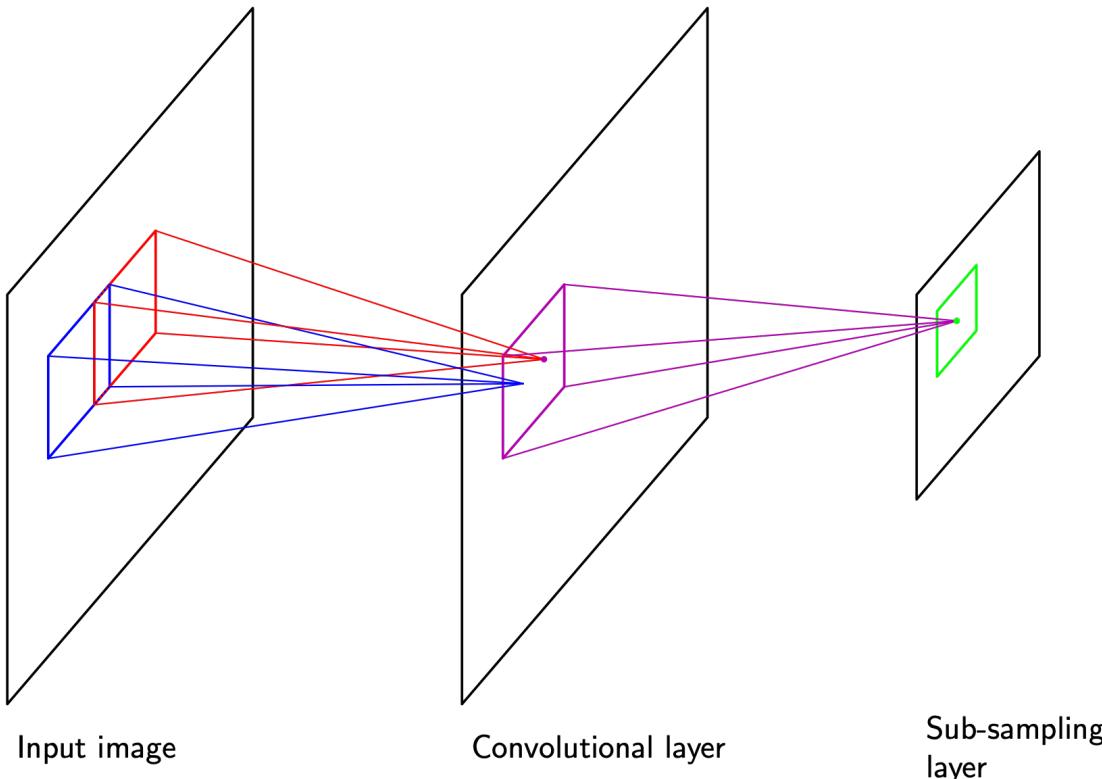
**AUSTRALIAN  
INSTITUTE FOR  
MACHINE LEARNING**

# Today

---

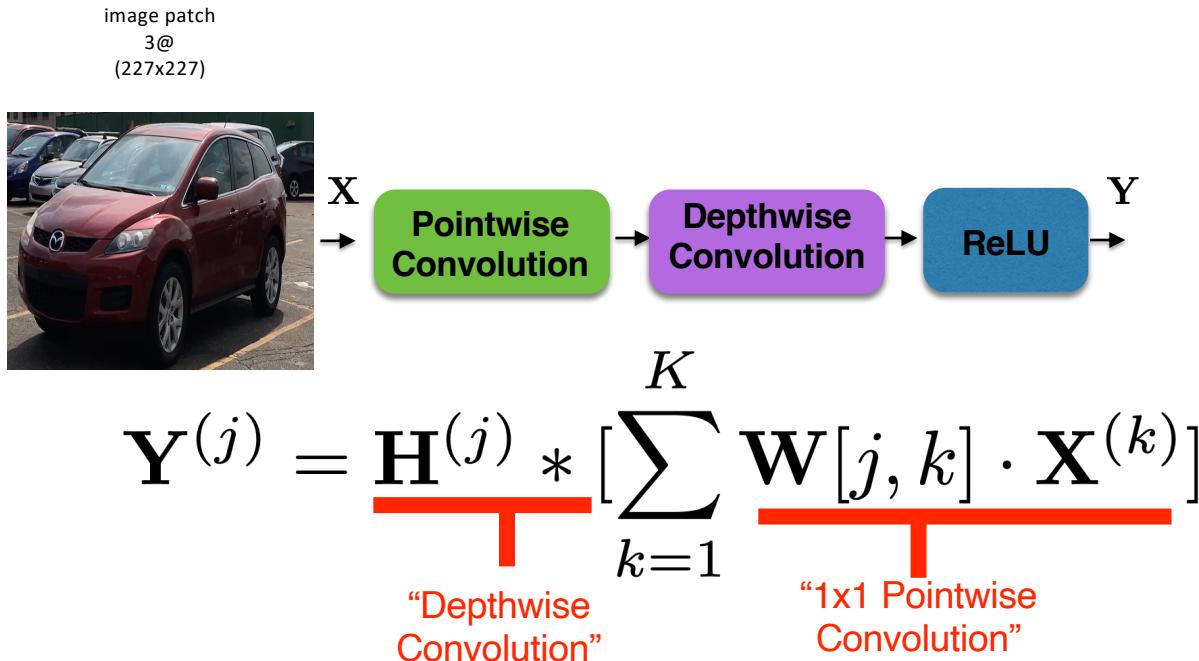
- Vision Transformers (ViT)
- Masked Self Supervision
- **Graph Neural Networks**

# Convolutional = A Type of Attention??



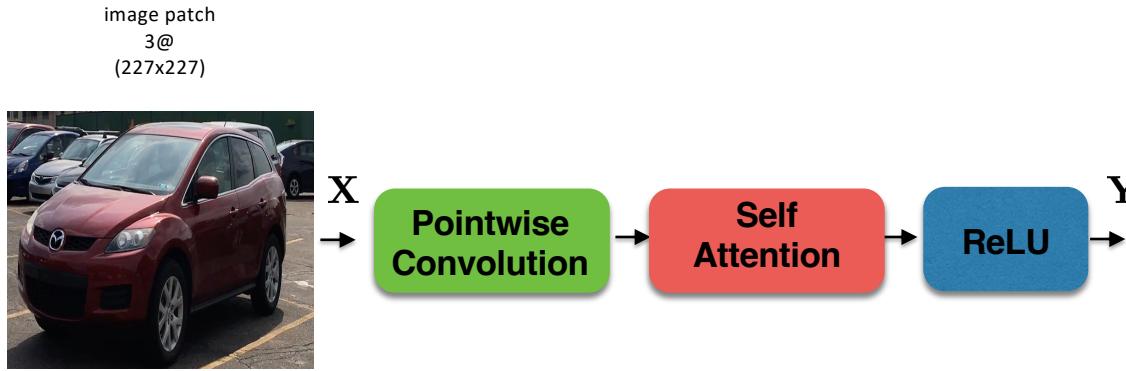
LeCun 1980

# Reminder: Depthwise Convolution



$$\mathcal{O}(M \cdot N \cdot D^2 \cdot J \cdot K) \rightarrow \mathcal{O}(M \cdot N \cdot (D^2 + J) \cdot K)$$

# Reminder: Depthwise Convolution



# Attention is all you need!!

---

- Convolution is hard-coded attention.
  - Could there be any way to learn it from data?
- 

## Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
[avaswani@google.com](mailto:avaswani@google.com)

**Noam Shazeer\***  
Google Brain  
[noam@google.com](mailto:noam@google.com)

**Niki Parmar\***  
Google Research  
[nikip@google.com](mailto:nikip@google.com)

**Jakob Uszkoreit\***  
Google Research  
[usz@google.com](mailto:usz@google.com)

**Llion Jones\***  
Google Research  
[llion@google.com](mailto:llion@google.com)

**Aidan N. Gomez\* †**  
University of Toronto  
[aidan@cs.toronto.edu](mailto:aidan@cs.toronto.edu)

**Lukasz Kaiser\***  
Google Brain  
[lukaszkaiser@google.com](mailto:lukaszkaiser@google.com)

**Illia Polosukhin\* ‡**  
[illia.polosukhin@gmail.com](mailto:illia.polosukhin@gmail.com)

## Abstract

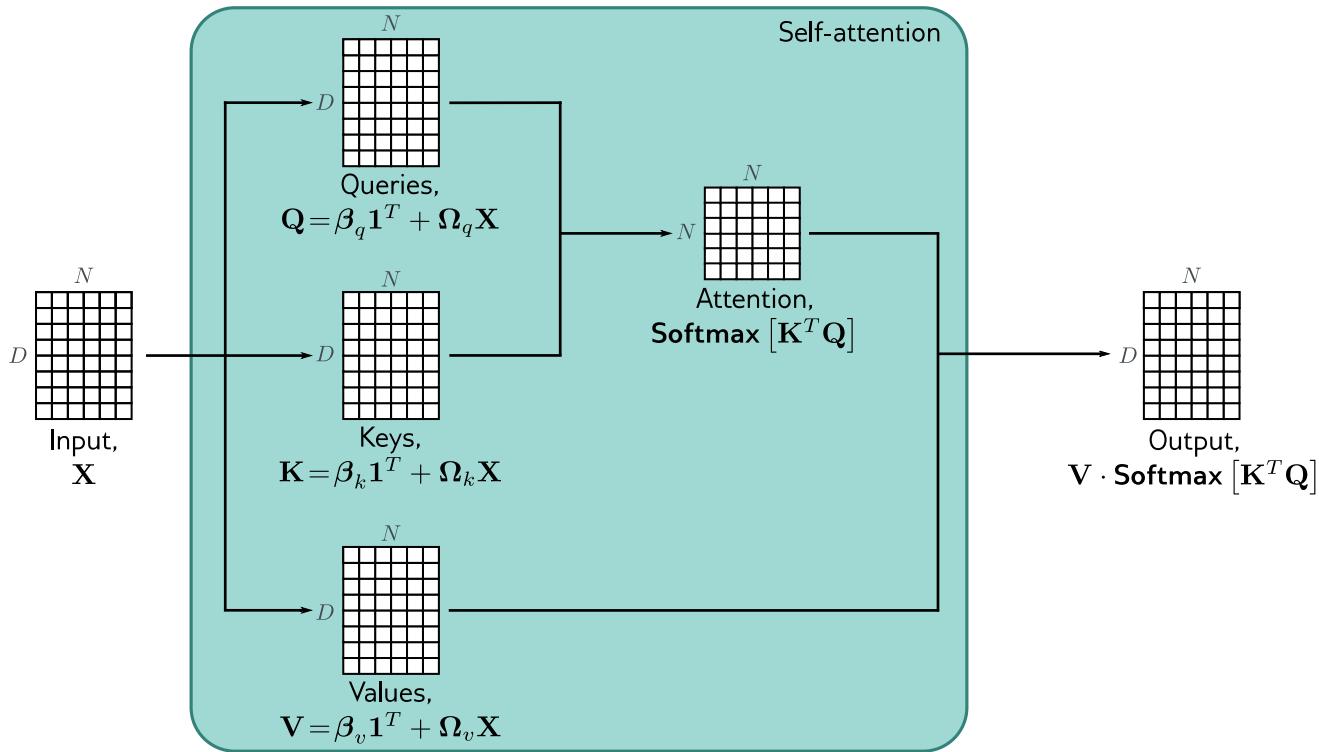
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention

# Soft-Max Attention

---

$$\begin{aligned} a[\mathbf{x}_m, \mathbf{x}_n] &= \text{softmax}_m [\mathbf{k}_\bullet^T \mathbf{q}_n] \\ &= \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]}, \end{aligned}$$

# What is Attention?



# What is Attention?

---

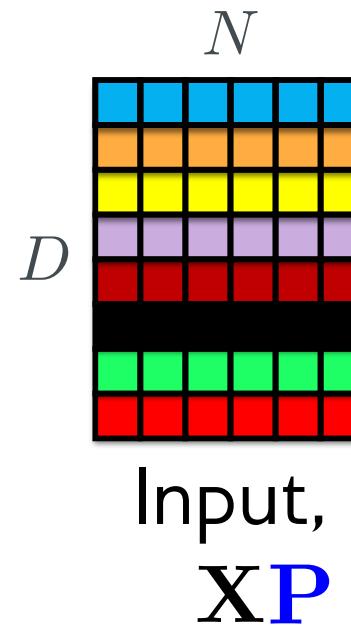
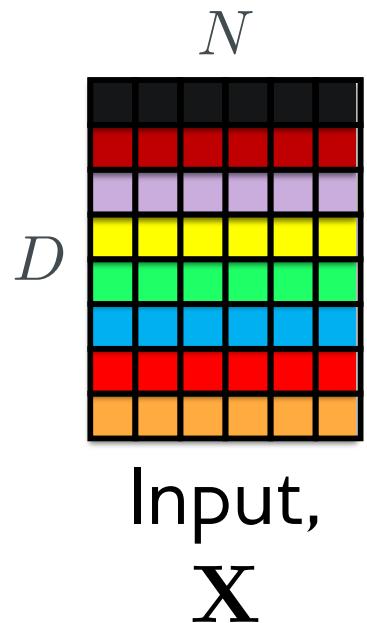
$$\mathbf{V}[\mathbf{X}] = \beta_v \mathbf{1}^T + \Omega_v \mathbf{X}$$

$$\mathbf{Q}[\mathbf{X}] = \beta_q \mathbf{1}^T + \Omega_q \mathbf{X}$$

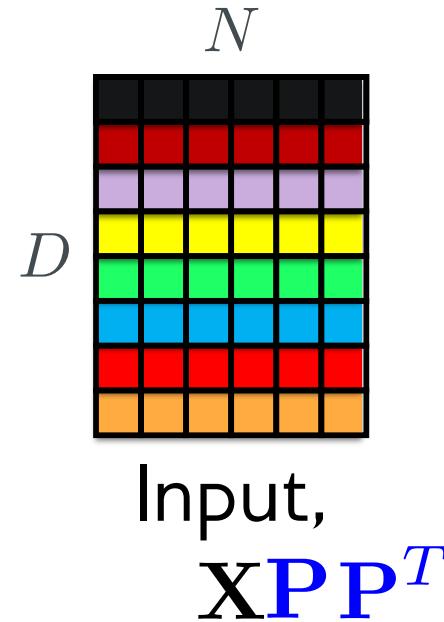
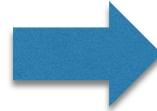
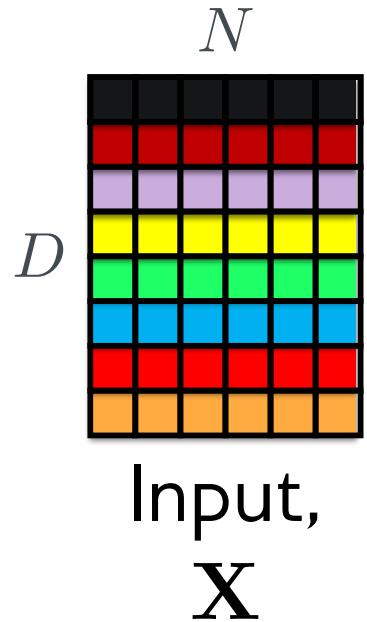
$$\mathbf{K}[\mathbf{X}] = \beta_k \mathbf{1}^T + \Omega_k \mathbf{X},$$

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{W}[\mathbf{X}] \cdot \mathbf{V} \circ \text{Softmax}\left[\mathbf{K}[\mathbf{X}]^T \mathbf{Q}[\mathbf{X}]\right],$$

# Permutating X



# Permutating X



# Permutation Equivariance?

---

$$\text{Sa}(\mathbf{X}) = \text{Sa}(\mathbf{X}\mathbf{P}\mathbf{P}^T)$$



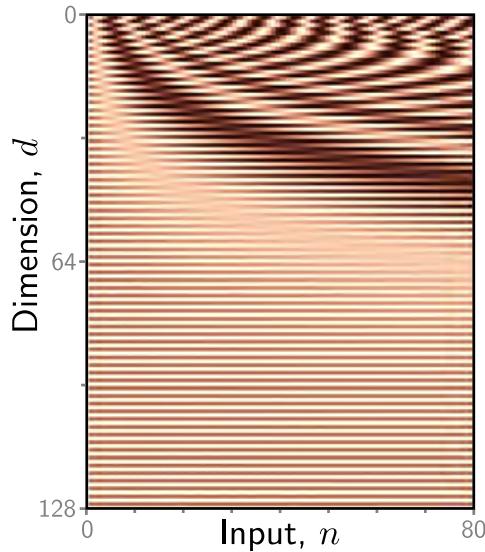
$$\text{Sa}(\mathbf{X}) = \text{Sa}(\mathbf{X}\mathbf{P})\mathbf{P}^T$$

Position information ignored!!!

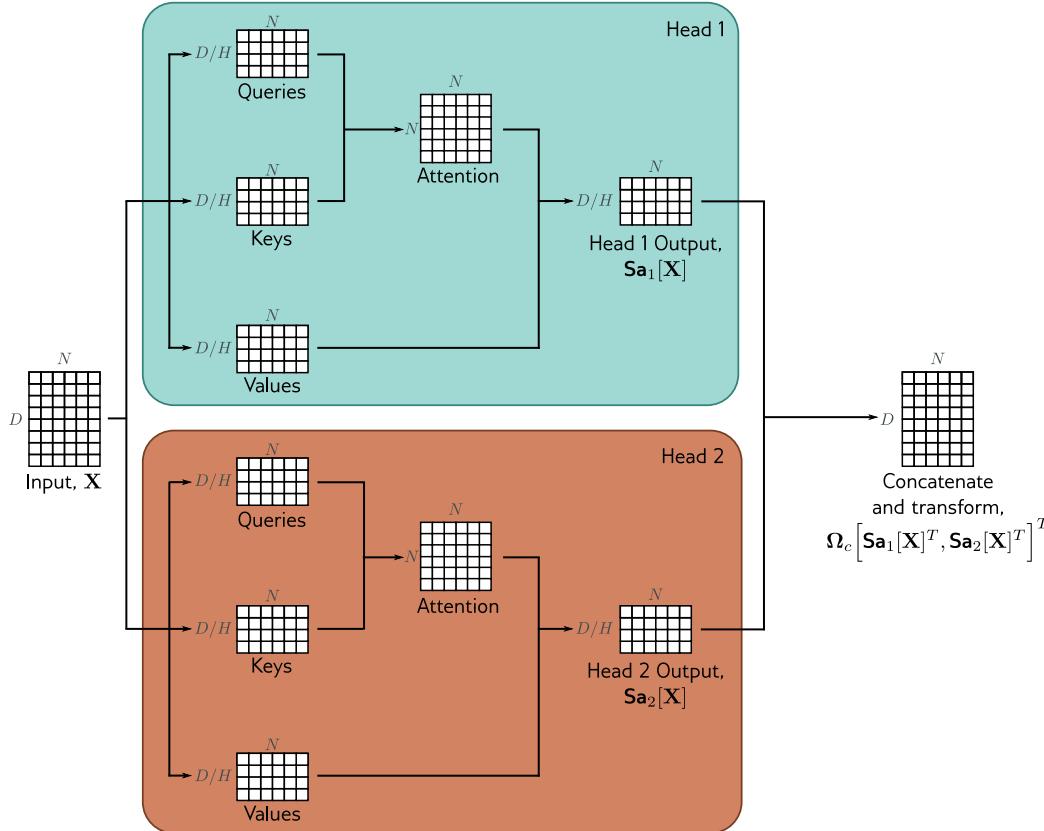
# Positional Encoding

$$\mathbf{X} \rightarrow \mathbf{X} + \mathbf{P}$$

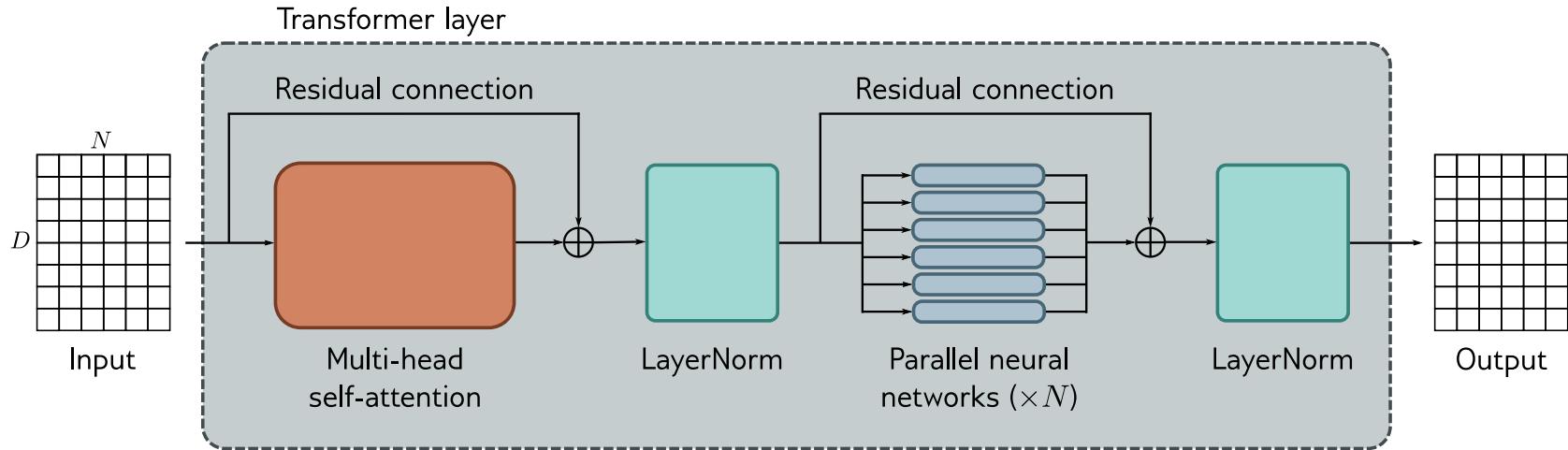
$$\mathbf{P} =$$



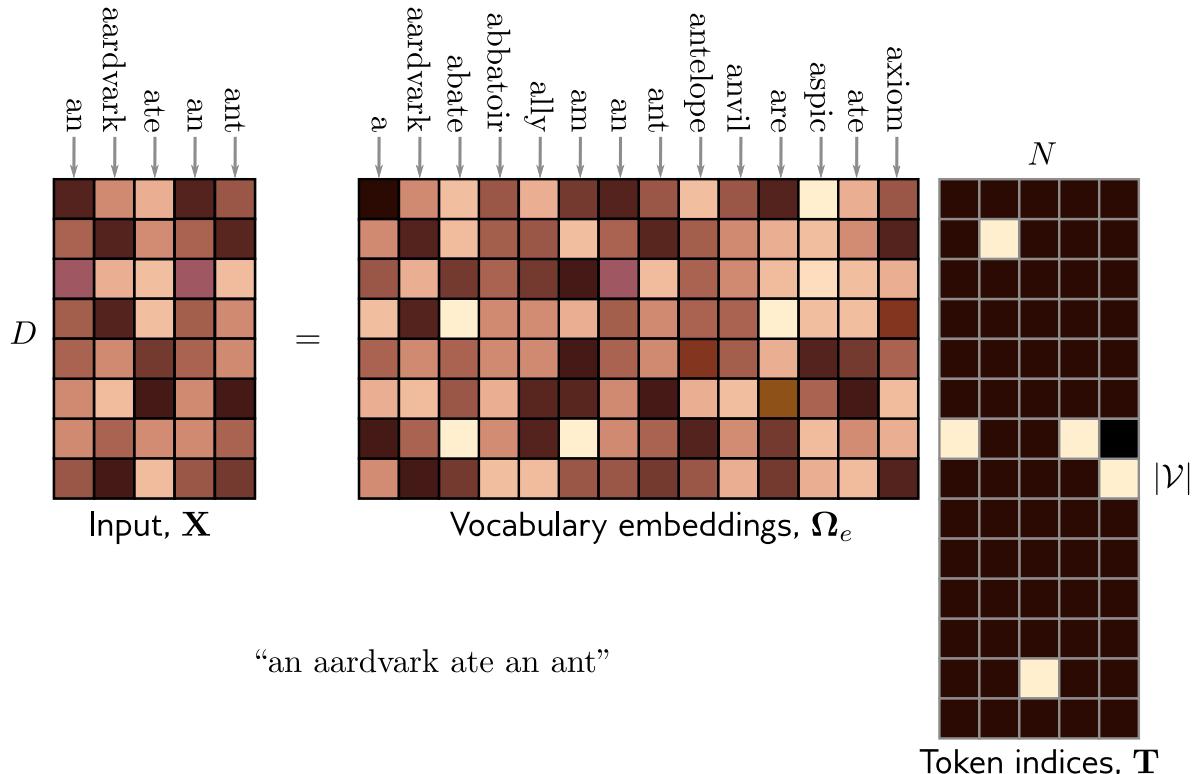
# Multi-Head Attention



# Transformer Layer



# Language + Transformers



# Vision Transformer (ViT)

---

Published as a conference paper at ICLR 2021

---

## 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

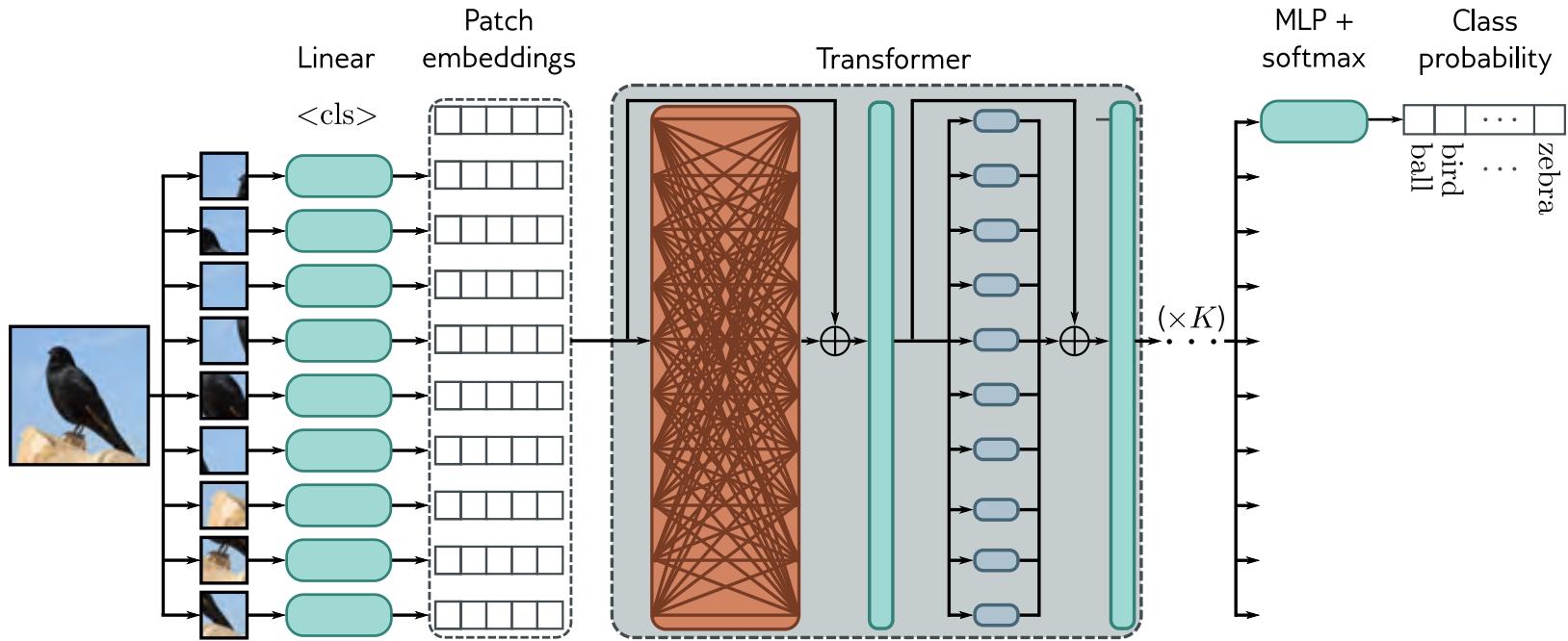
Google Research, Brain Team

{adosovitskiy, neilhoulsby}@google.com

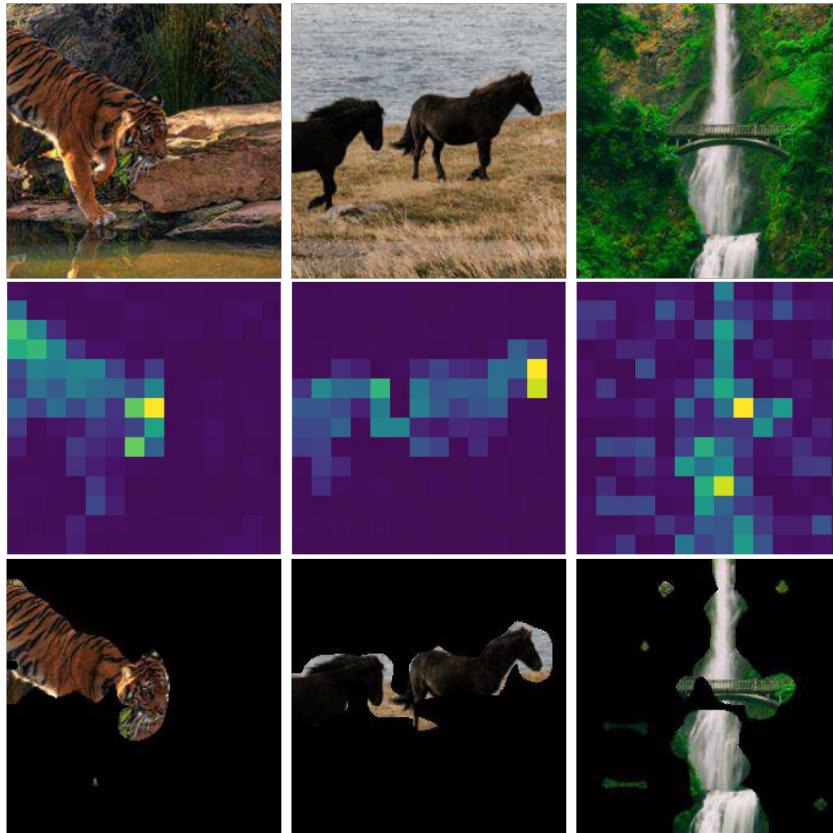
### ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAR, etc.). Vision Transformer (ViT) attains excellent

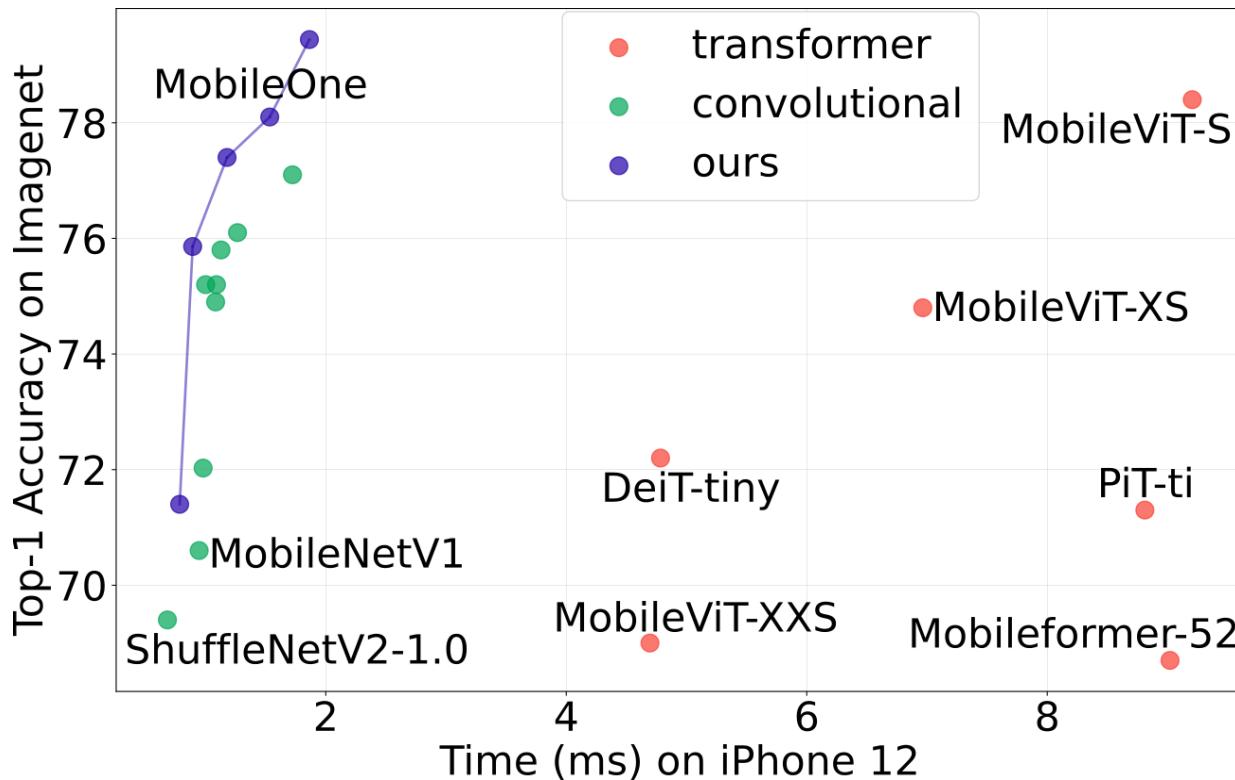
# Vision Transformer (ViT)

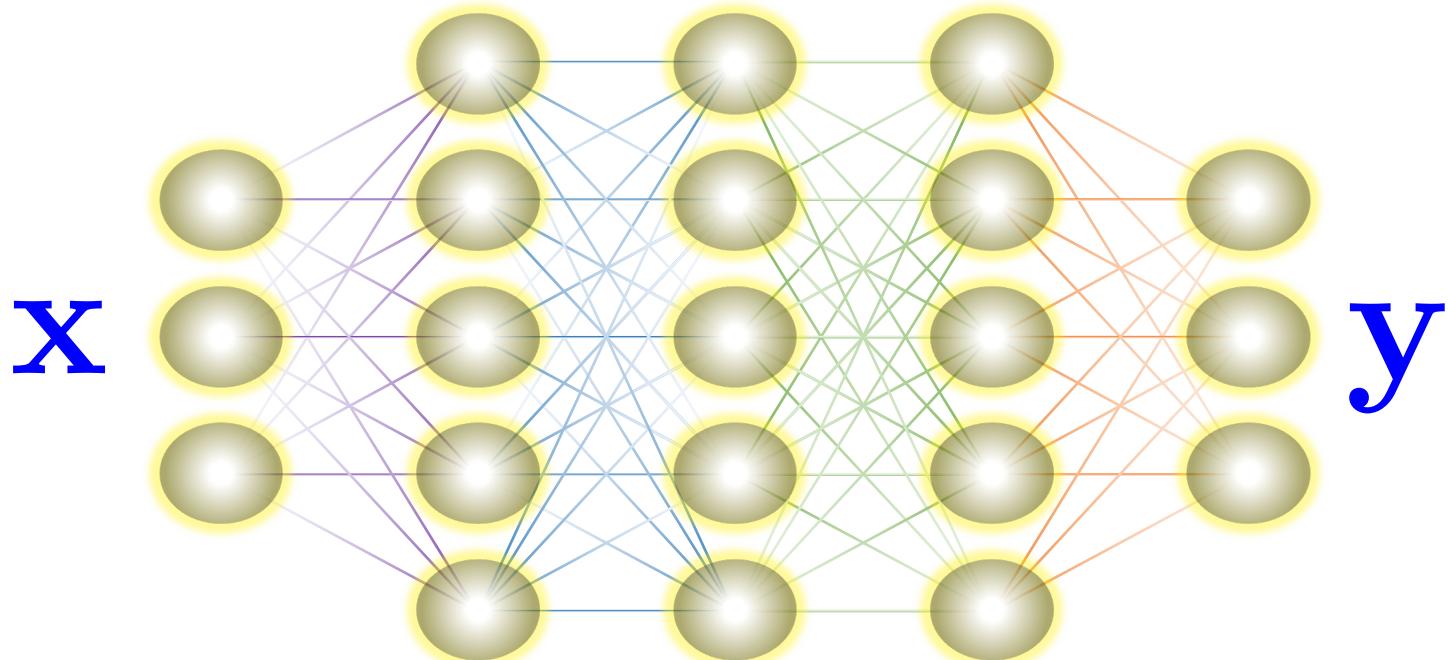


# Visualizing Attention



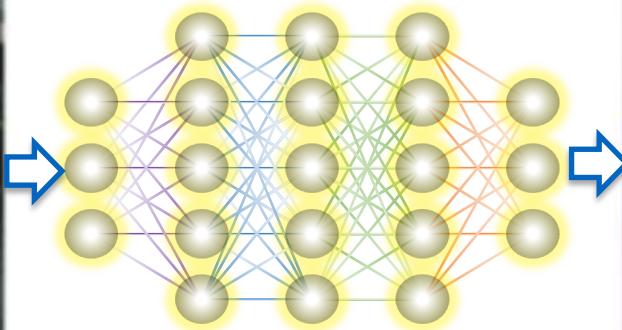
# Transformers are Getting Faster!!!







**x**

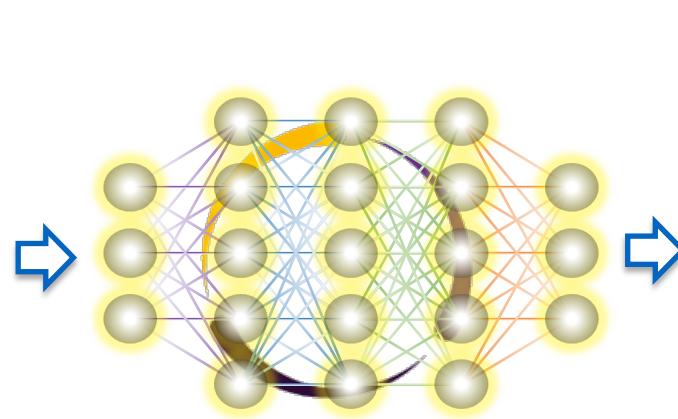


**y**

Where do I get the labels?



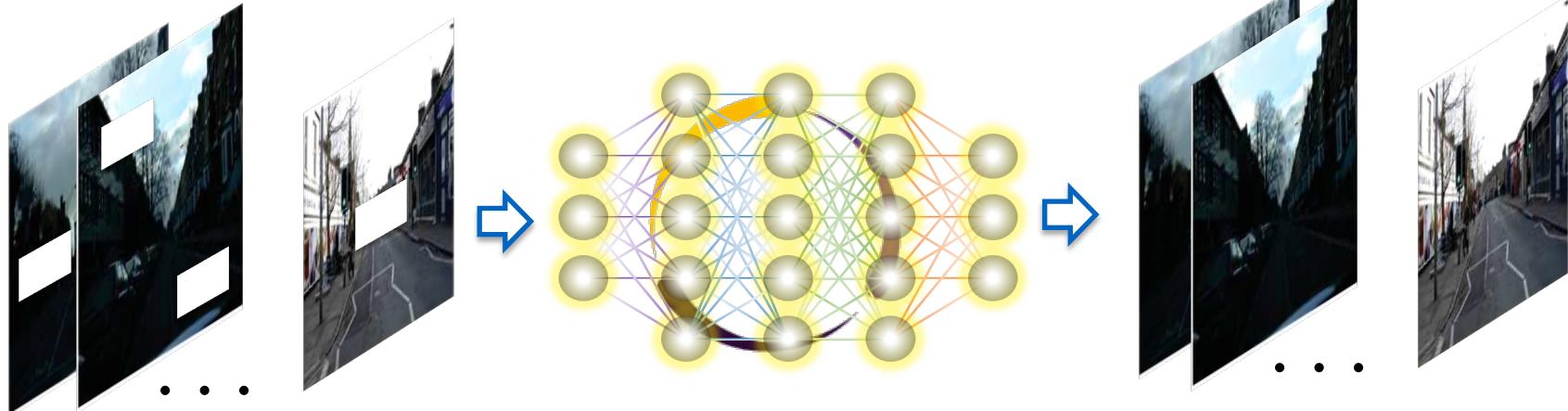
$[x_1, \dots, x_N]$



$[y_1, \dots, y_N]$

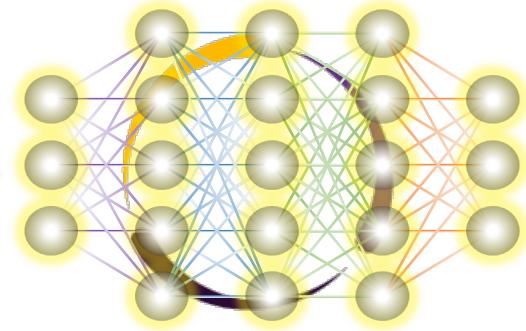
$$N \propto 10^6 - 10^9$$

Requires huge datasets and compute!!!



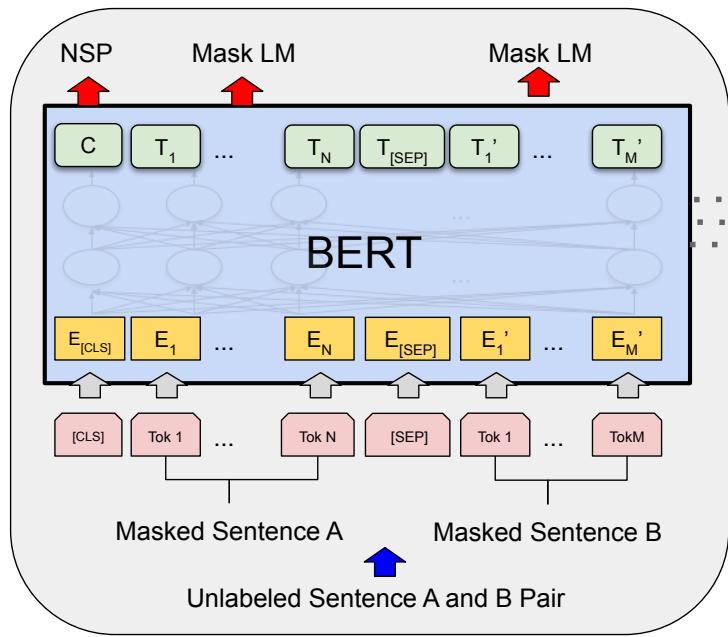
Self Supervision

The █ walked across the █

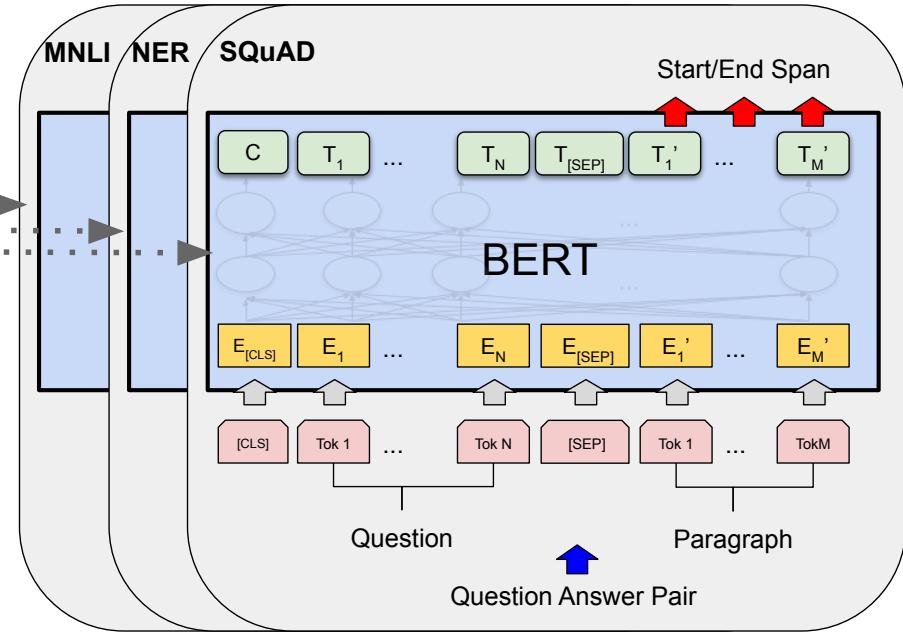


The cat walked across the room

# Self Supervision

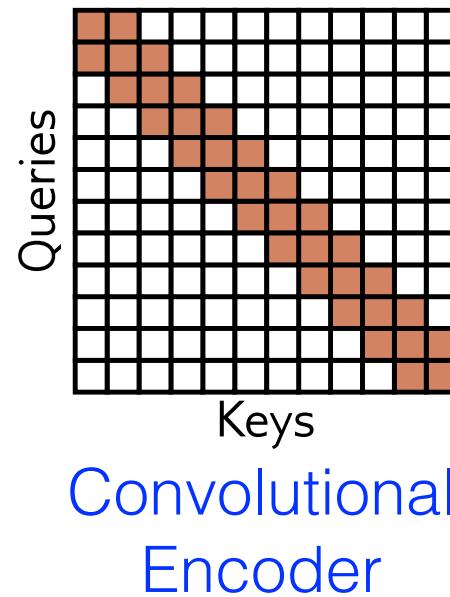
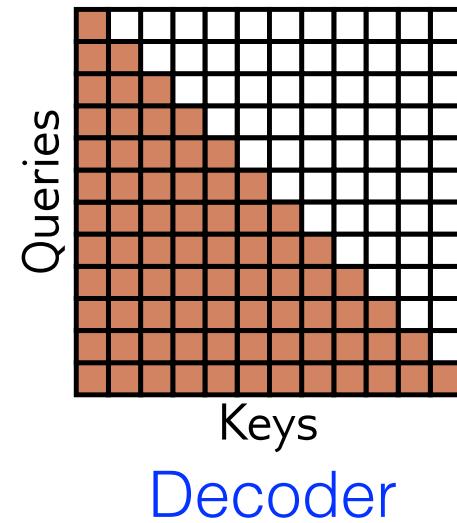
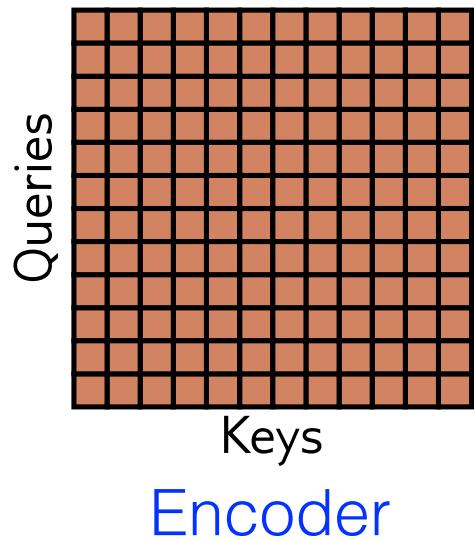


Pre-training

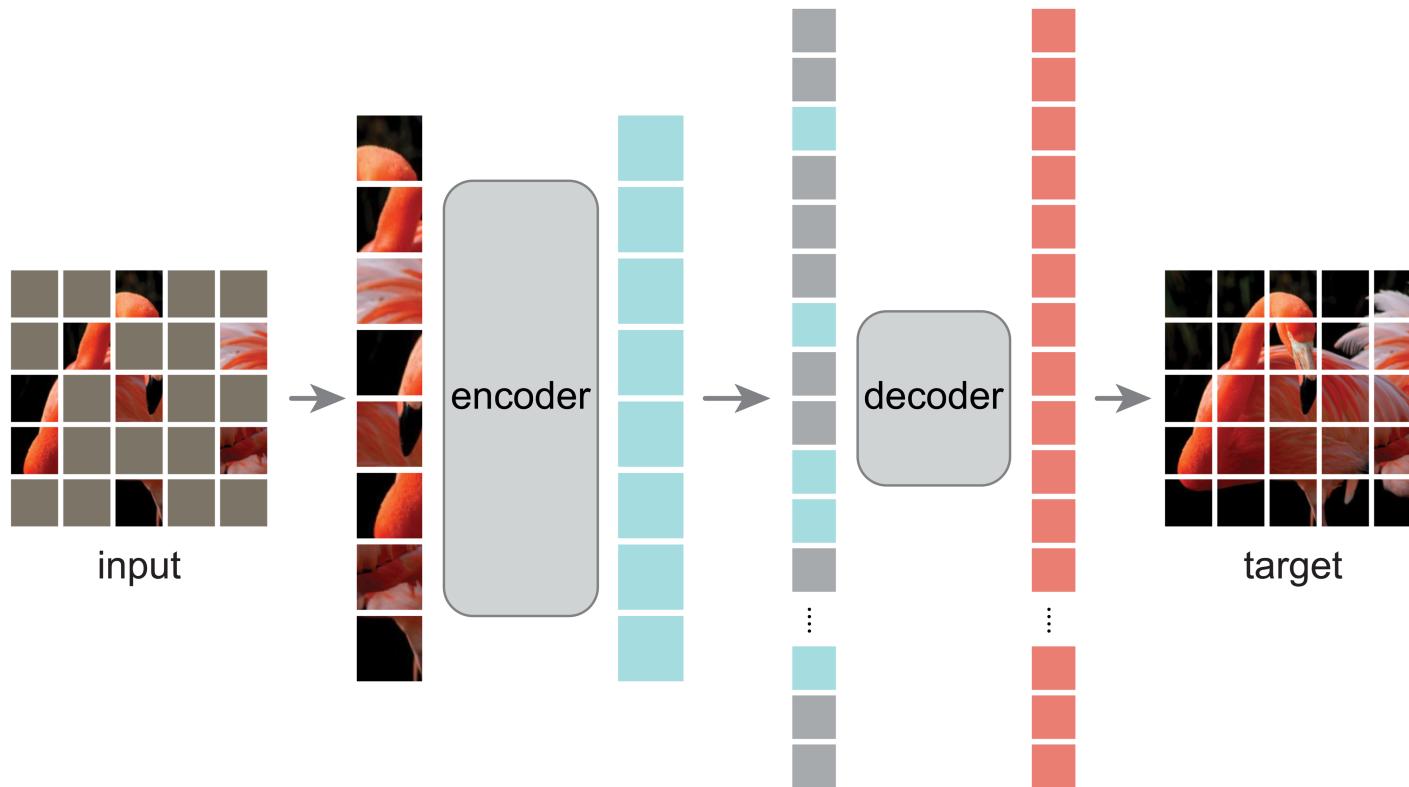


Fine-Tuning

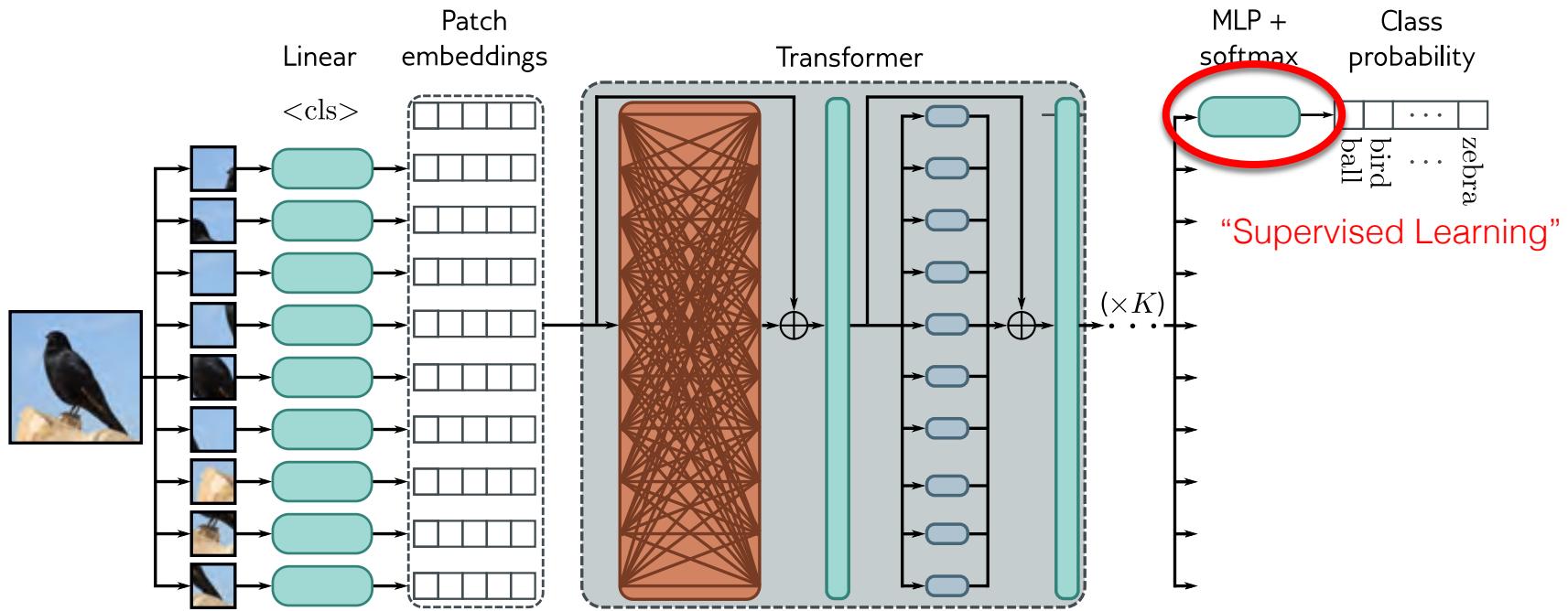
# Interaction Matrices for Self Attention



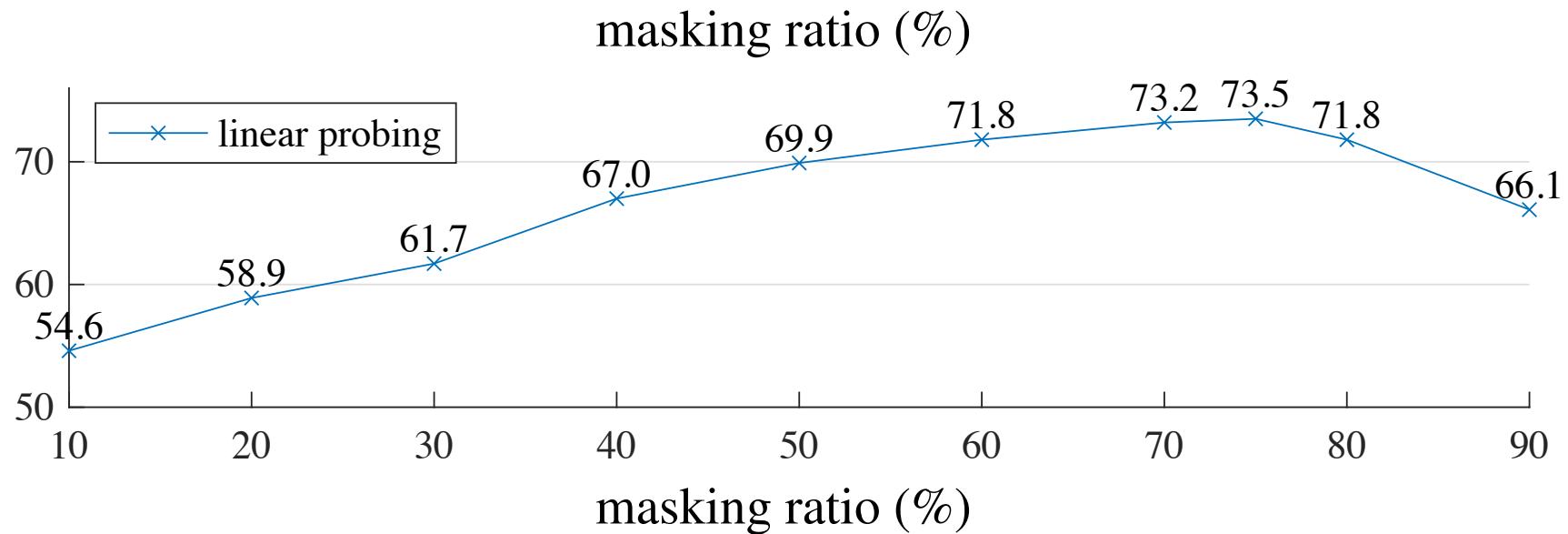
# Masked ViT



# Masked ViT



# Masked ViT

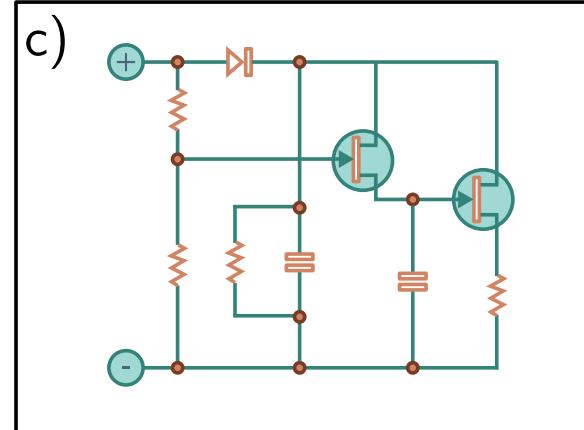
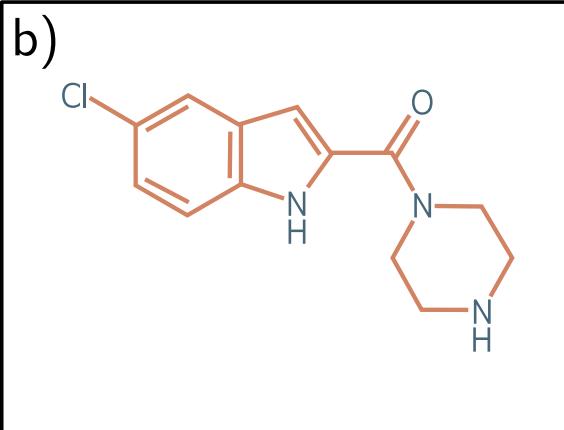
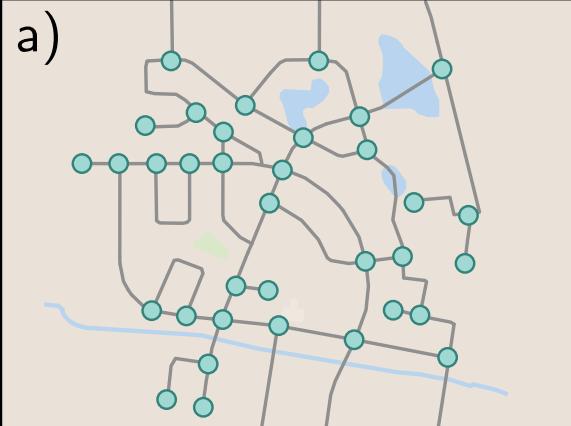


# Today

---

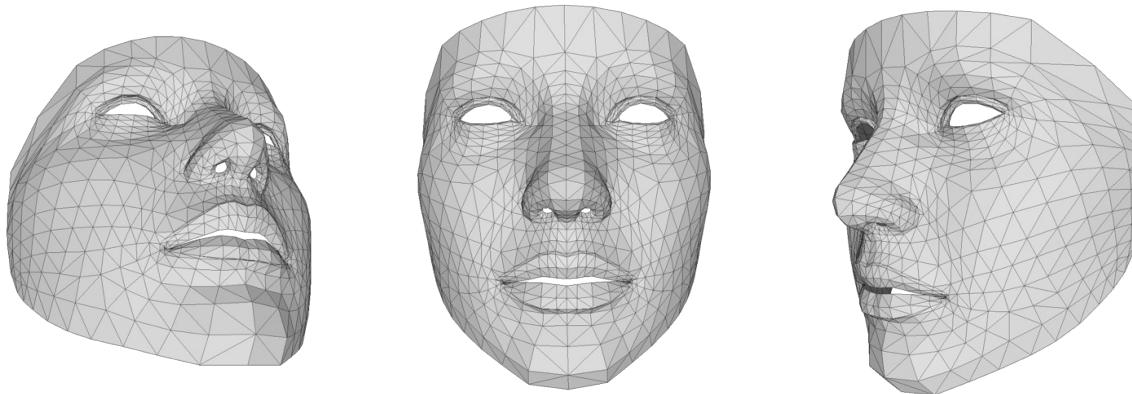
- Vision Transformer Networks
- Self Supervision
- **Graph Neural Networks**

# Examples of Graphs?



# General Topology

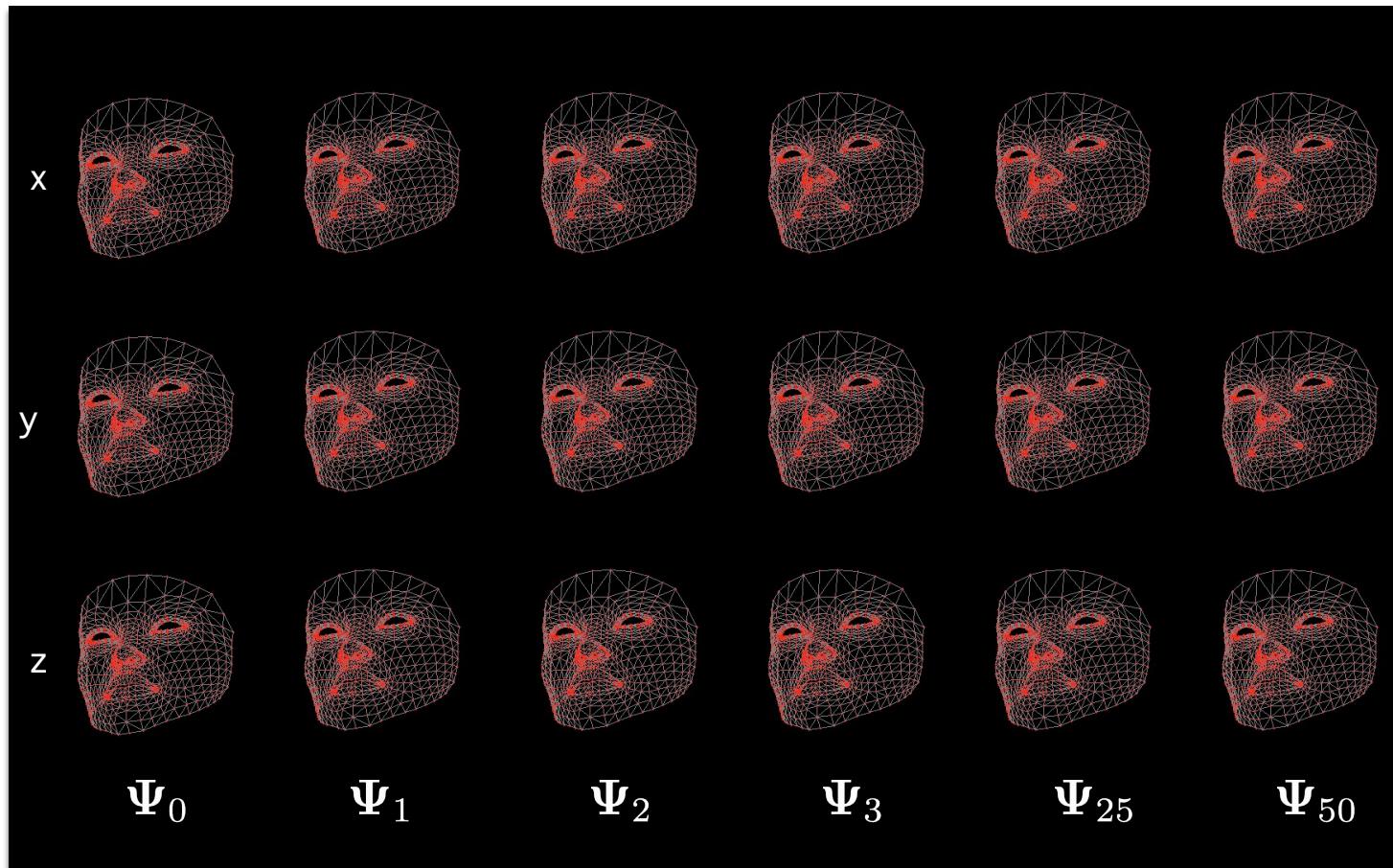
What if  $x \neq \text{grid?}$



$$\mathbf{L}[i, j] = \begin{cases} |\mathcal{N}(i)| & \text{if } i = j \\ -1 & \text{if } j \in \mathcal{N}(i) \\ 0 & \text{otherwise} \end{cases} \quad \rightarrow \quad \mathbf{L} = \mathbf{D} - \mathbf{A}$$
$$\mathbf{L}' = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}$$

Normalized Graph Laplacian

# Smooth Deformation Basis



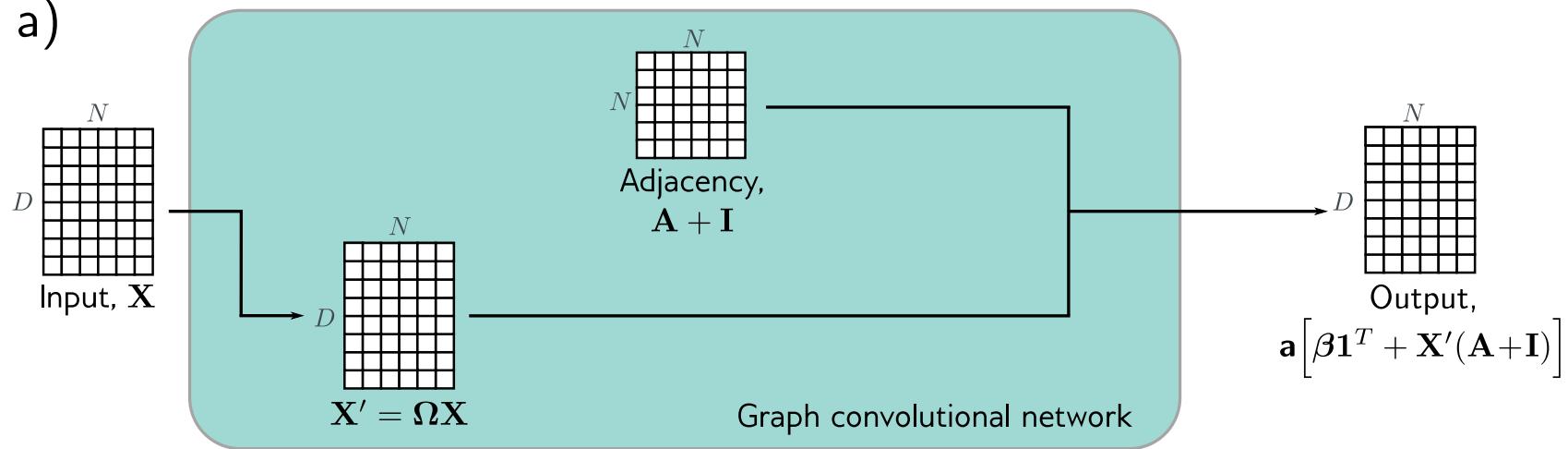
# Normalized Graph Laplacian

```
A =      ([[0., 0., 1., 1., 1.],
           [1., 0., 1., 0., 1.],
           [1., 0., 0., 0., 1.],
           [0., 1., 0., 0., 0.],
           [0., 0., 1., 1., 0.]))

L_norm = ([[1.0000, 1.0000, 0.5918, 0.4226, 0.5918],
            [0.6667, 1.0000, 0.5918, 1.0000, 0.5918],
            [0.5918, 1.0000, 1.0000, 1.0000, 0.5000],
            [1.0000, 0.4226, 1.0000, 1.0000, 1.0000],
            [1.0000, 1.0000, 0.5000, 0.2929, 1.0000]])
```

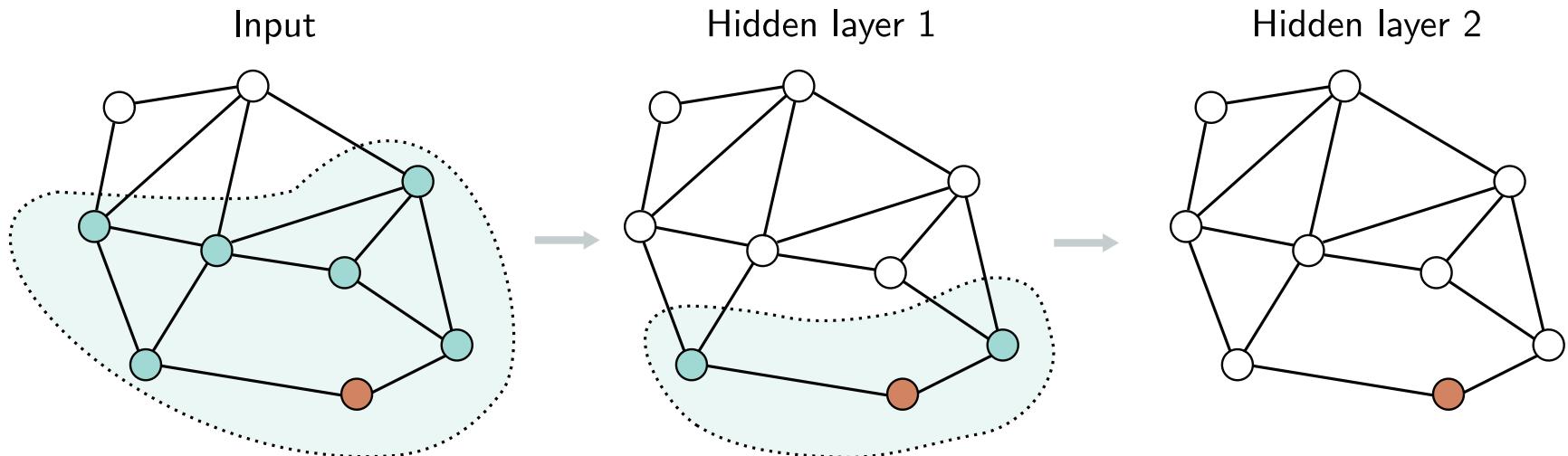
# Graph Convolution Network

a)



Only  $\Omega$  and  $\beta$  are learned!!

# Graph Convolution Network



# Graph Attention Networks

---

## GRAPH ATTENTION NETWORKS

**Petar Veličković\***

Department of Computer Science and Technology  
University of Cambridge  
[petar.velickovic@cst.cam.ac.uk](mailto:petar.velickovic@cst.cam.ac.uk)

**Guillem Cucurull\***

Centre de Visió per Computador, UAB  
[gcucurull@gmail.com](mailto:gcucurull@gmail.com)

**Arantxa Casanova\***

Centre de Visió per Computador, UAB  
[ar.casanova.8@gmail.com](mailto:ar.casanova.8@gmail.com)

**Adriana Romero**

Montréal Institute for Learning Algorithms  
[adriana.romero.soriano@umontreal.ca](mailto:adriana.romero.soriano@umontreal.ca)

**Pietro Liò**

Department of Computer Science and Technology  
University of Cambridge  
[pietro.lioc@cst.cam.ac.uk](mailto:pietro.lioc@cst.cam.ac.uk)

**Yoshua Bengio**

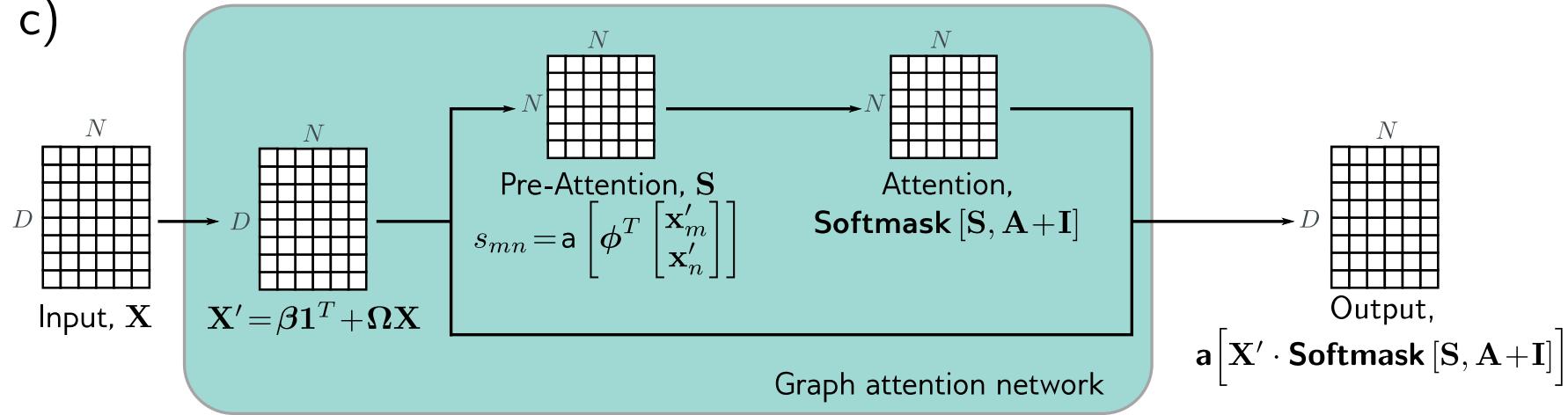
Montréal Institute for Learning Algorithms  
[yoshua.umontreal@gmail.com](mailto:yoshua.umontreal@gmail.com)

## ABSTRACT

We present graph attention networks (GATs), novel neural network architectures that operate on graph-structured data, leveraging masked self-attentional layers to address the shortcomings of prior methods based on graph convolutions or their approximations. By stacking layers in which nodes are able to attend over their neighborhoods' features, we enable (implicitly) specifying different weights to different nodes in a neighborhood, without requiring any kind of costly matrix operation (such as inversion) or depending on knowing the graph structure upfront. In this way, we address several key challenges of spectral-based graph neural networks simultaneously, and make our model readily applicable to inductive as well as transductive problems. Our GAT models have achieved or matched state-of-the-art results across four established transductive and inductive graph benchmarks: the *Cora*, *Citeseer* and *Pubmed* citation network datasets, as well as a *protein-protein interaction* dataset (wherein test graphs remain unseen during training).

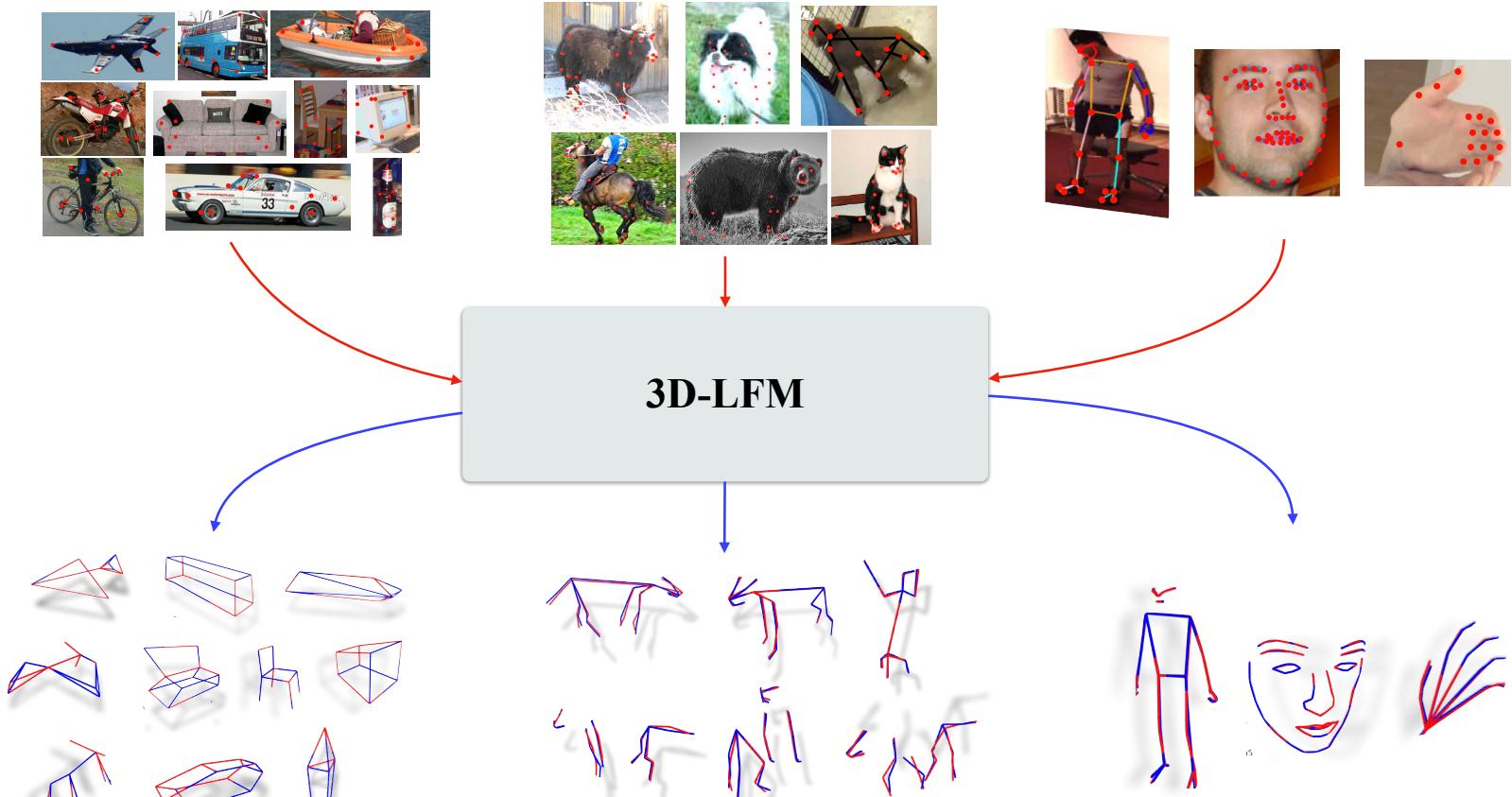
# Graph Attention Networks

c)

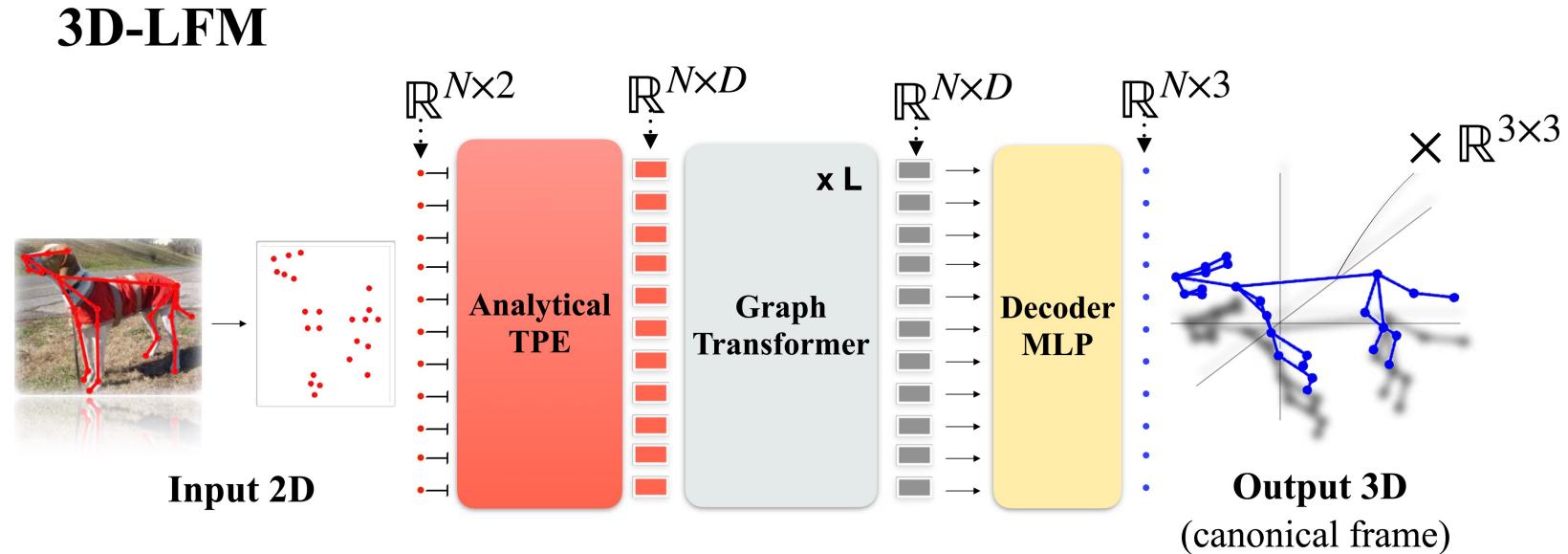


$\Omega$ ,  $\phi$  and  $\beta$  are learned!!!

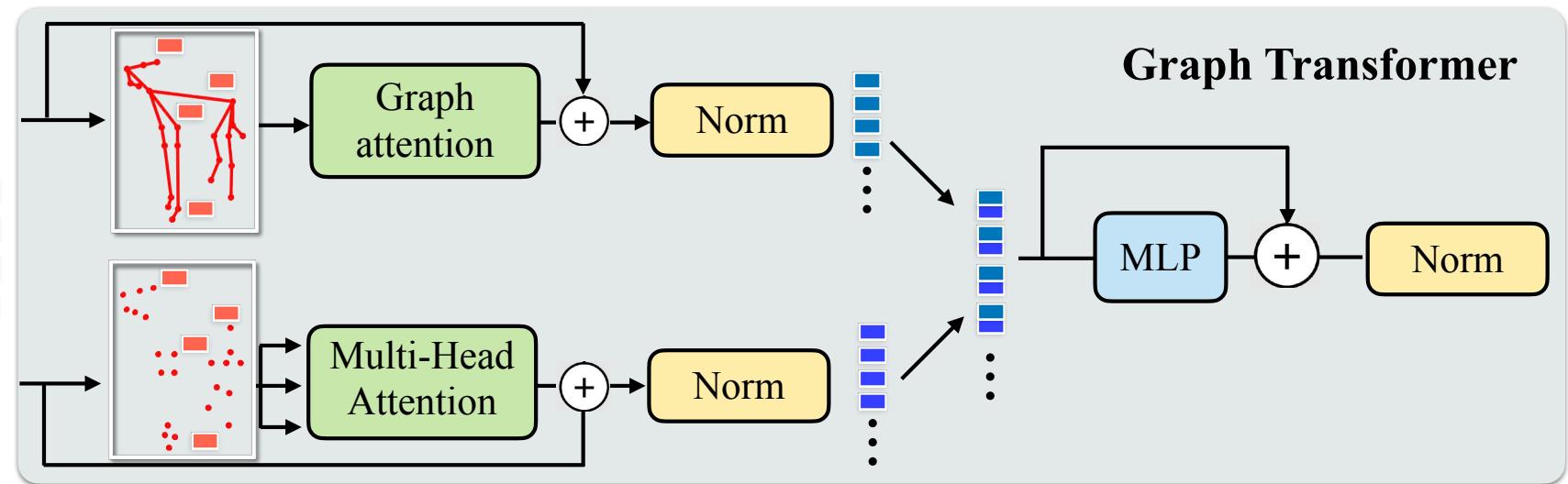
# Example - 3D Lifting Foundation Model



# Example - 3D Lifting Foundation Model

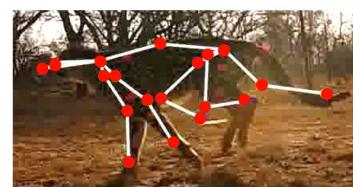
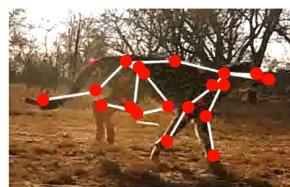
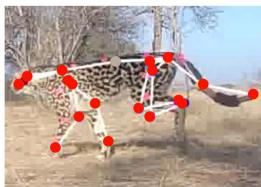


# Example - 3D Lifting Foundation Model

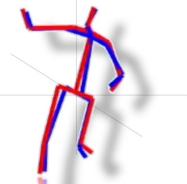
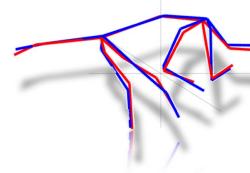
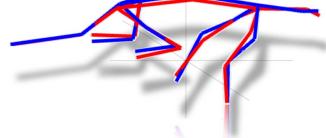
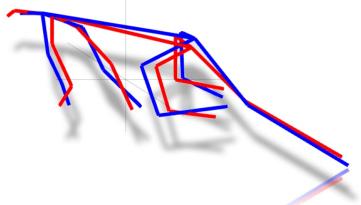


# Example - 3D Lifting Foundation Model

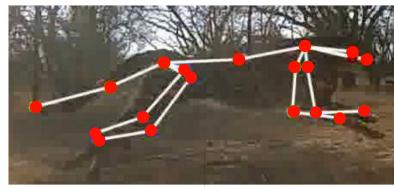
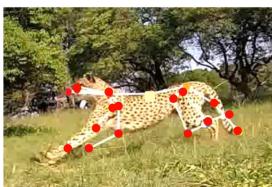
OOD  
Input 2D



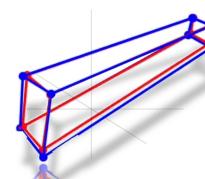
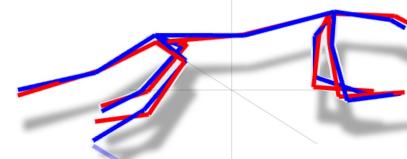
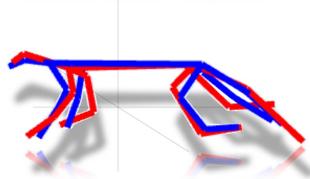
Predicted  
3D



OOD  
Input 2D



Predicted  
3D



## More to read...

---

