

COMP547

DEEP UNSUPERVISED LEARNING

Lecture #4 – Attention and Transformers



KOÇ
UNIVERSITY

Aykut Erdem // Koç University // Spring 2024

Previously on COMP547

- sequence modeling
- recurrent neural networks (RNNs)
- how to train RNNs
- long short-term memory (LSTM)
- gated recurrent unit (GRU)
- sequence to sequence modeling

Using RNNs to generate Super Mario Maker levels, Adam Geitgey



Lecture overview

- memory
- tokens
- attention
- positional encoding

Disclaimer: Much of the material and slides for this lecture were borrowed from
— Phillip Isola and Sara Beery's MIT 6.S898 slides

Recap: Memory in RNNs

Why not remember everything?

- Memory size grows with t
- This kind of memory is **nonparametric**: there is no finite set of parameters we can use to model it
- RNNs make a Markov assumption — the future hidden state only depends on the immediately preceding hidden state
- By putting the right info in to the hidden state, RNNs can model dependences that are arbitrarily far apart

The problem of long-range dependences

- Other methods exist that do directly link old “memories” (observations or hidden states) to future predictions:
- Temporal convolutions
- Attention / Transformers (see <https://arxiv.org/abs/1706.03762>)
- Memory networks (see <https://arxiv.org/abs/1410.3916>)

Modeling arbitrarily long sequences

- **Recurrence** — recurrent weights are shared across time
- **Convolution** — conv weights are shared across time
- **Attention** — weights are dynamically determined as a function of the data
(conv kernel with attention weights is shown on the right)

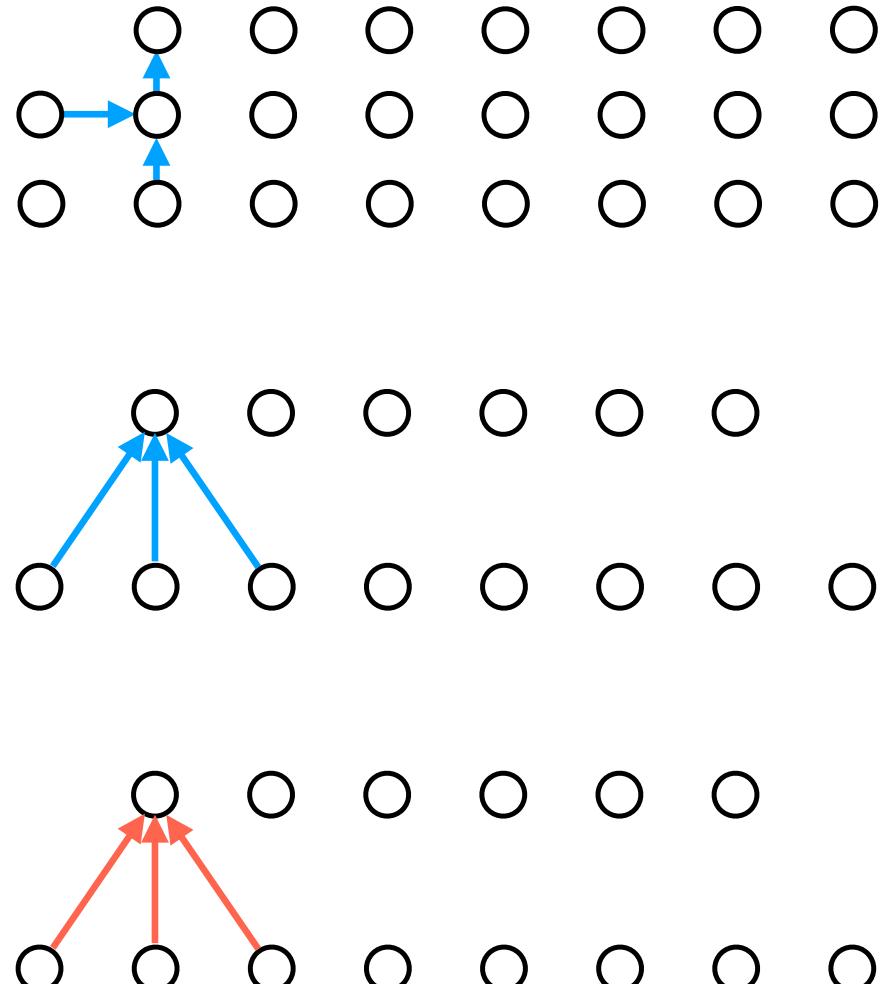


Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Transformers

- Three key ideas
 - Tokens
 - Attention
 - Positional encoding
- Examples of architectures and applications

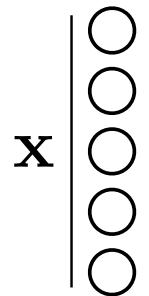
Rank	Model	Top 1 Accuracy	Top 5 Accuracy	Number of params	GFLOPs	Extra Training Data	Paper	Code	Result	Year	Tags
1	CoCa (finetuned)	91.0%		2100M		✓	CoCa: Contrastive Captioners are Image-Text Foundation Models			2022	
2	Model soups (BASIC-L)	90.98%		2440M		✓	Model soups: averaging weights of multiple fine- tuned models improves accuracy without increasing inference time			2022	
3	Model soups (ViT-G/14)	90.94%		1843M		✓	Model soups: averaging weights of multiple fine- tuned models improves accuracy without increasing			2022	

Idea #1: tokens

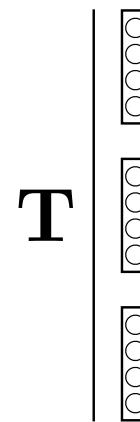
A new data structure: Tokens

- A **token** is just transformer lingo for a vector of neurons (note: GNNs also operate over tokens, but over there we called them “node attributes” or node “feature descriptors”)
- But the connotation is that a token is an encapsulated bundle of information; with transformers we will operate over tokens rather than over neurons

array of **neurons**



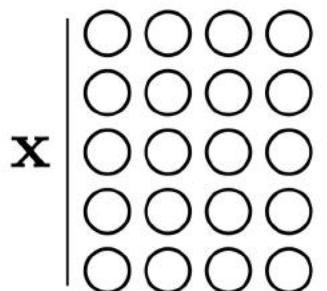
array of **tokens**



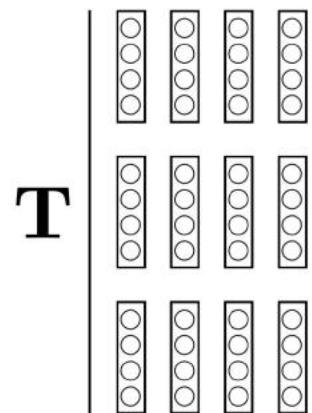
A new data structure: Tokens

- A **token** is just transformer lingo for a vector of neurons (note: GNNs also operate over tokens, but over there we called them “node attributes” or node “feature descriptors”)
- But the connotation is that a token is an encapsulated bundle of information; with transformers we will operate over tokens rather than over neurons

array of **neurons**

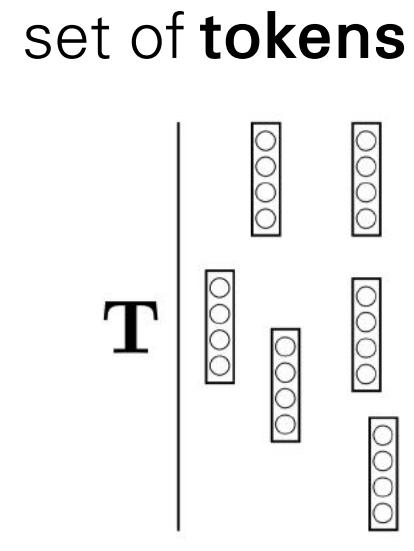
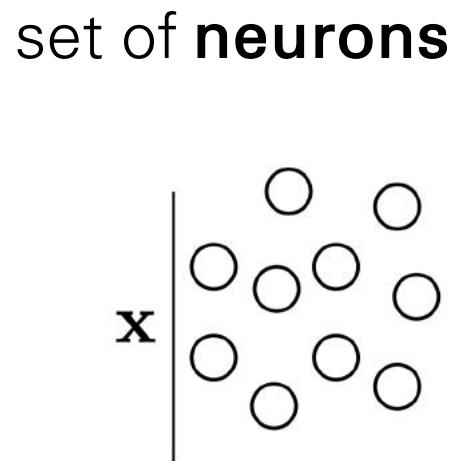


array of **tokens**

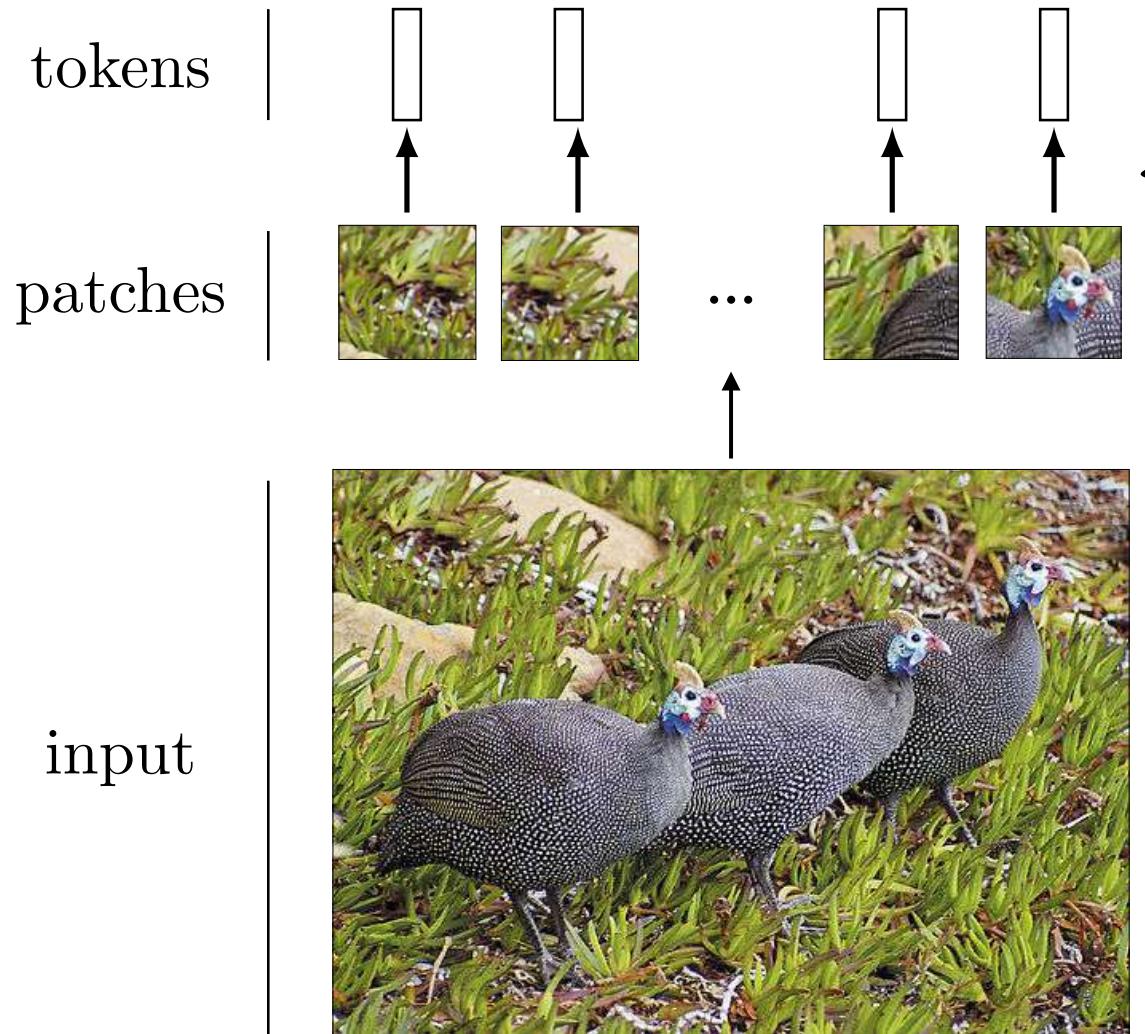


A new data structure: Tokens

- A **token** is just transformer lingo for a vector of neurons (note: GNNs also operate over tokens, but over there we called them “node attributes” or node “feature descriptors”)
- But the connotation is that a token is an encapsulated bundle of information; with transformers we will operate over tokens rather than over neurons



Tokenizing the input data

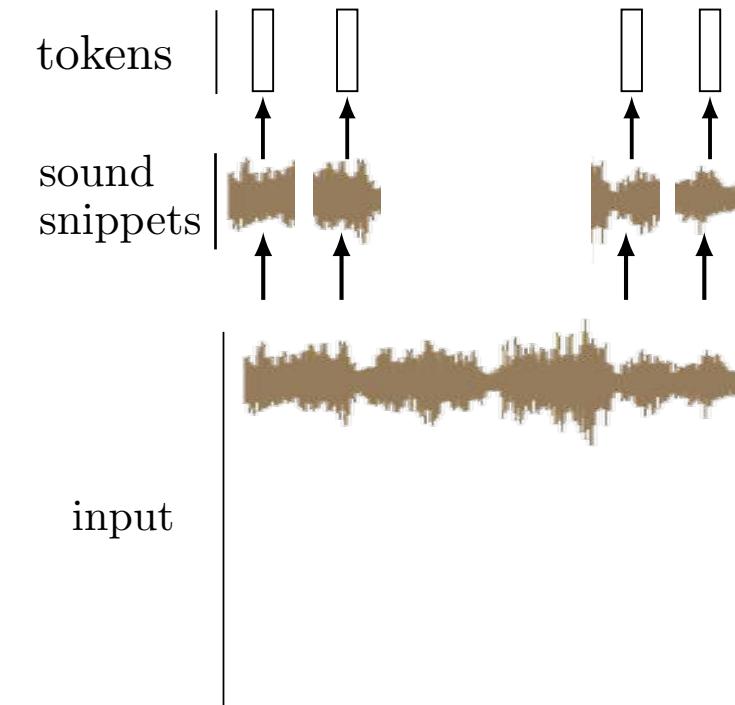
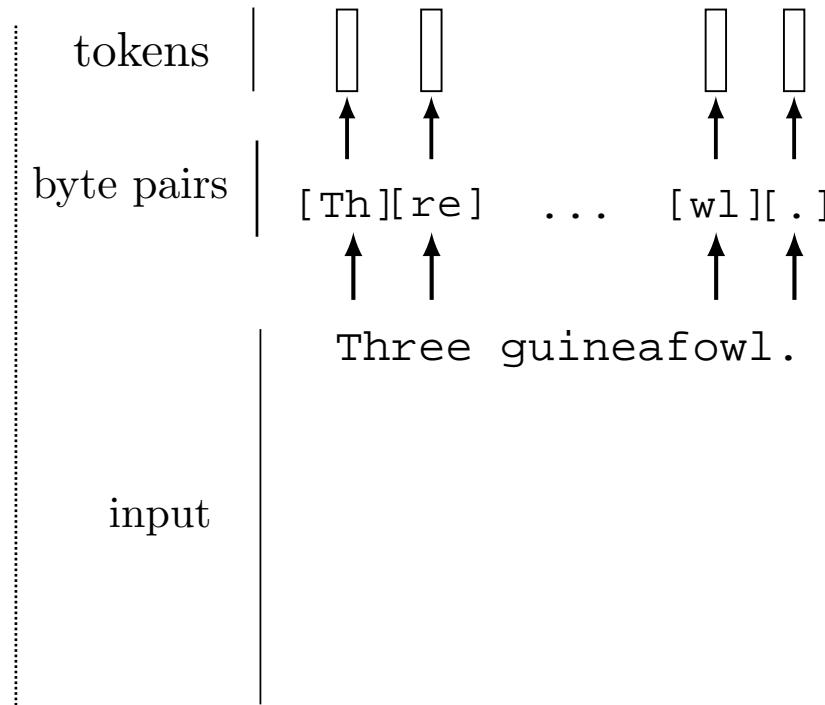
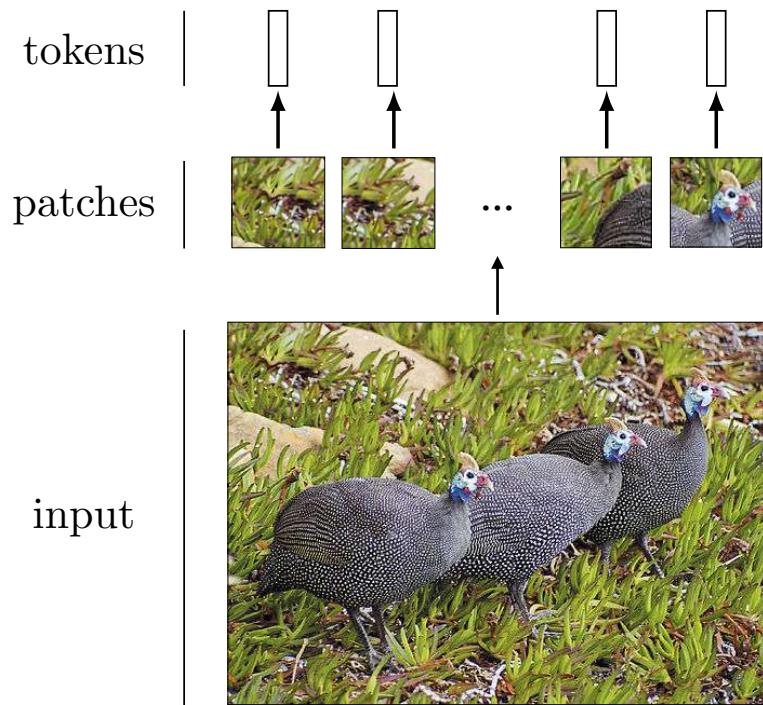


e.g. linear projection.

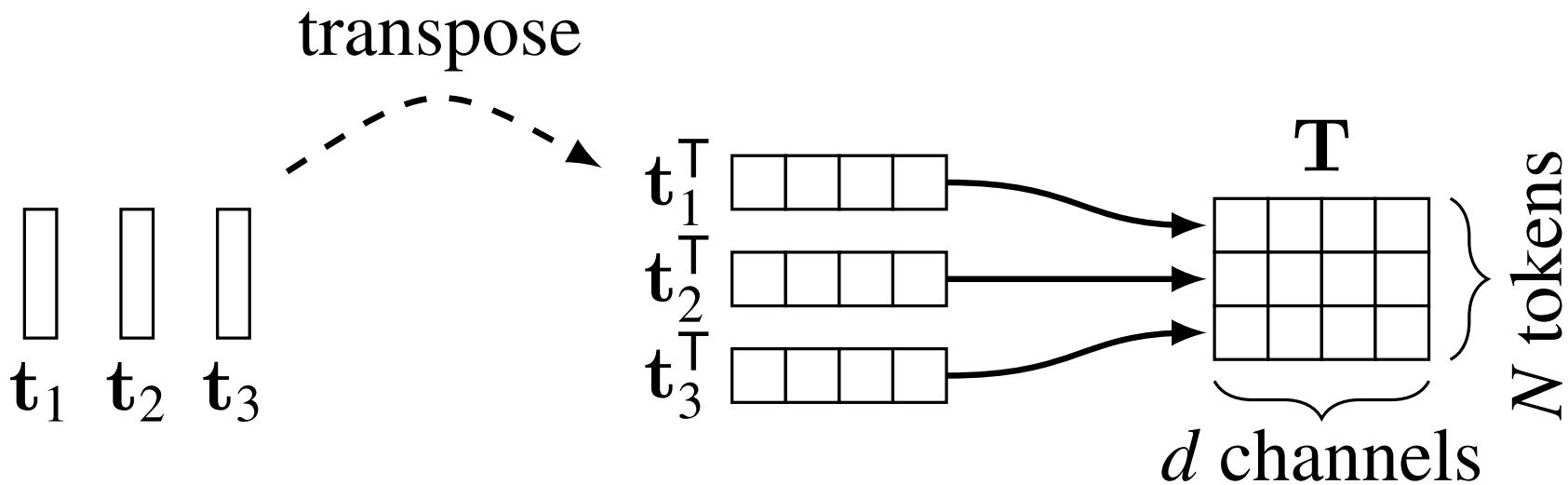
- When operating over neurons, we represent the input as an array of scalar-valued measurements (e.g., pixels)
- When operating over tokens, we represent the input as an array of vector-valued measurements

Tokenizing the input data

- You tokenize anything.
- General strategy: chop the input up into chunks, project each chunk to a vector

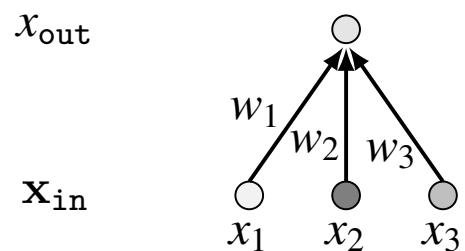


Notation



Linear combination of tokens

Linear combination of neurons

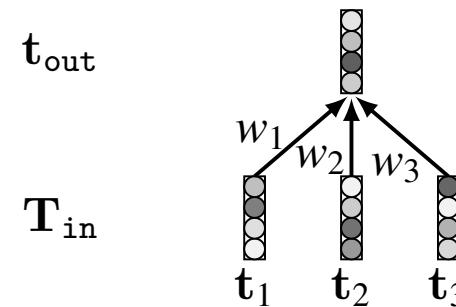


$$x_{\text{out}} = w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$x_{\text{out}}[i] = \sum_{j=1}^N w_{ij} x_{\text{in}}[j]$$

$$\mathbf{x}_{\text{out}} = \mathbf{W} \mathbf{x}_{\text{in}}$$

Linear combination of tokens



$$\mathbf{t}_{\text{out}} = w_1 \mathbf{t}_1 + w_2 \mathbf{t}_2 + w_3 \mathbf{t}_3$$

$$\mathbf{T}_{\text{out}}[i, :] = \sum_{j=1}^N w_{ij} \mathbf{T}_{\text{in}}[j, :]$$

$$\mathbf{T}_{\text{out}} = \mathbf{W} \mathbf{T}_{\text{in}}$$

Token-wise nonlinearity

$$\mathbf{x}_{\text{out}} = \begin{bmatrix} \text{relu}(x_{\text{in}}[0]) \\ \vdots \\ \text{relu}(x_{\text{in}}[N-1]) \end{bmatrix}$$

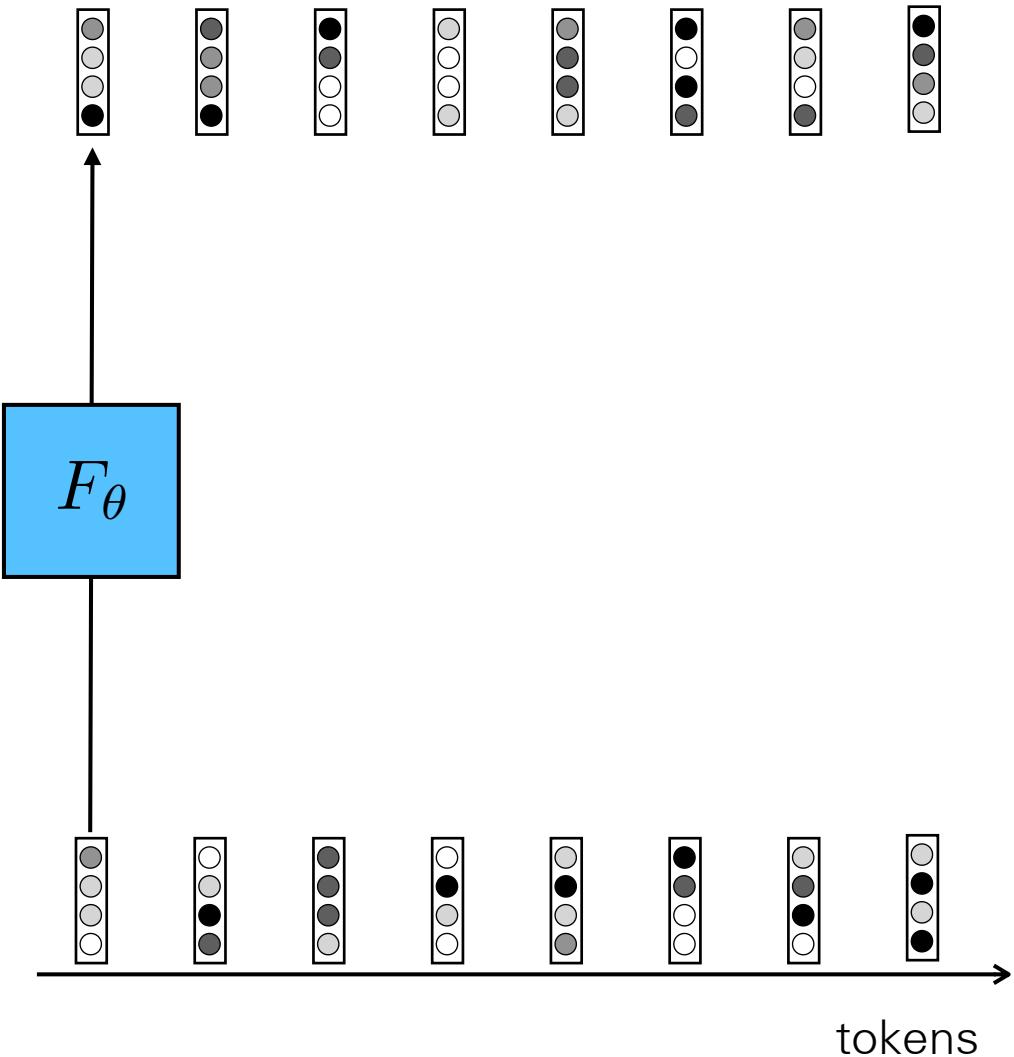
- F is typically an MLP
- Equivalent to a CNN with 1×1 kernels run over token sequence

$$\mathbf{T}_{\text{out}} = \begin{bmatrix} F_{\theta}(\mathbf{T}_{\text{in}}[0, :]) \\ \vdots \\ F_{\theta}(\mathbf{T}_{\text{in}}[N-1, :]) \end{bmatrix}$$

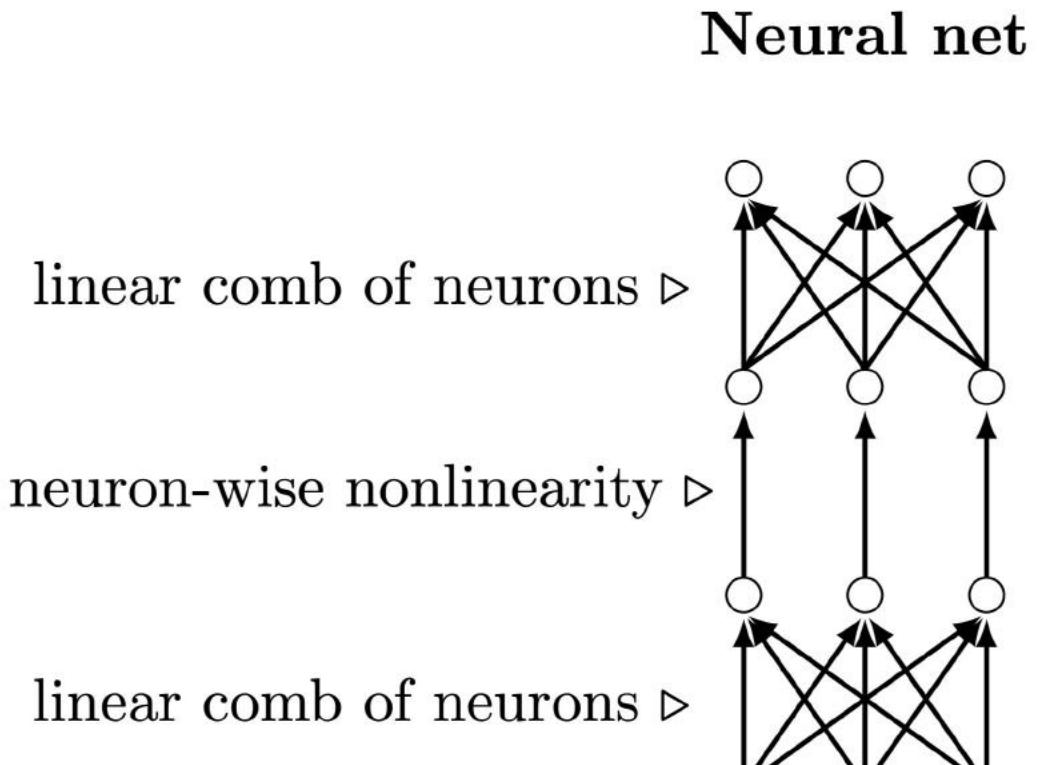
Token-wise nonlinearity

$$\mathbf{x}_{\text{out}} = \begin{bmatrix} \text{relu}(x_{\text{in}}[0]) \\ \vdots \\ \text{relu}(x_{\text{in}}[N-1]) \end{bmatrix}$$

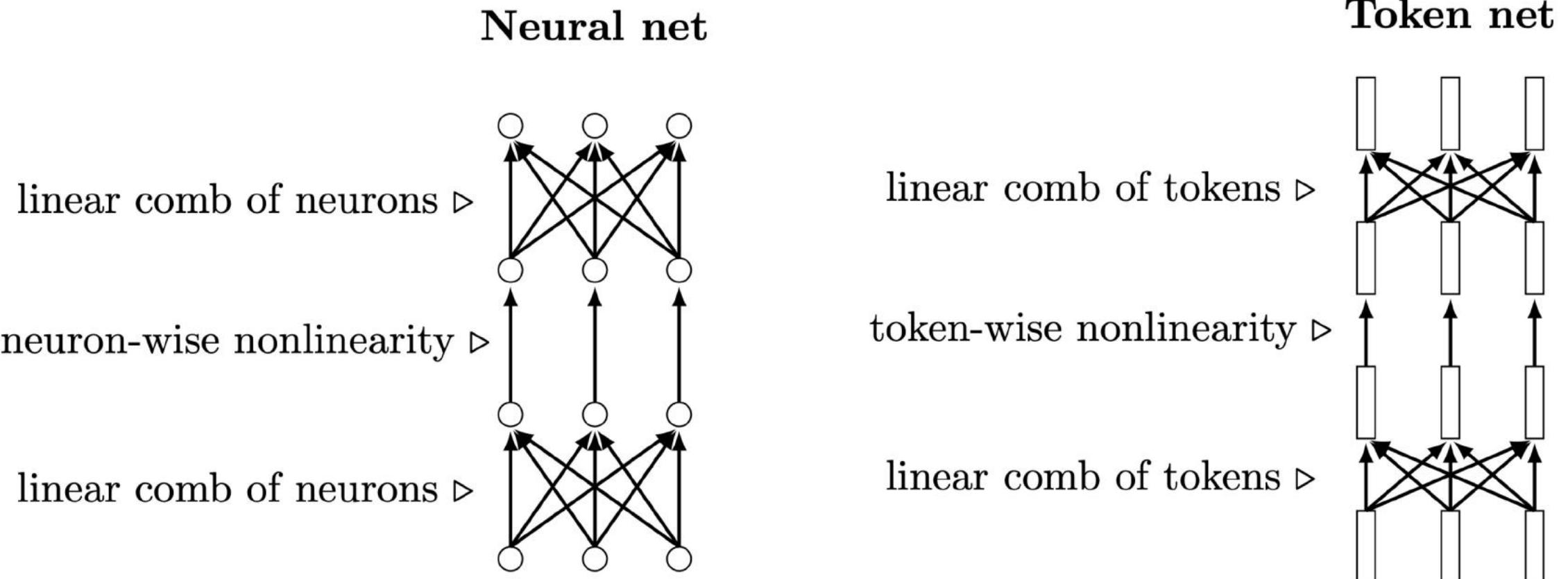
$$\mathbf{T}_{\text{out}} = \begin{bmatrix} F_{\theta}(\mathbf{T}_{\text{in}}[0, :]) \\ \vdots \\ F_{\theta}(\mathbf{T}_{\text{in}}[N-1, :]) \end{bmatrix}$$

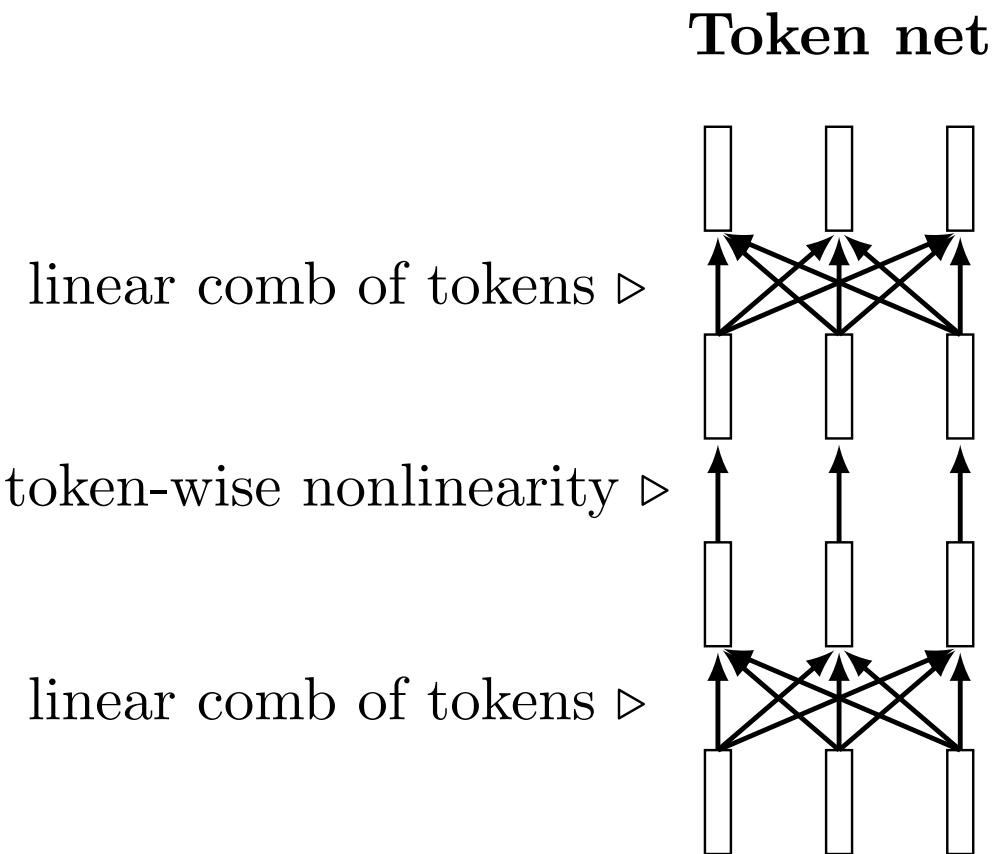
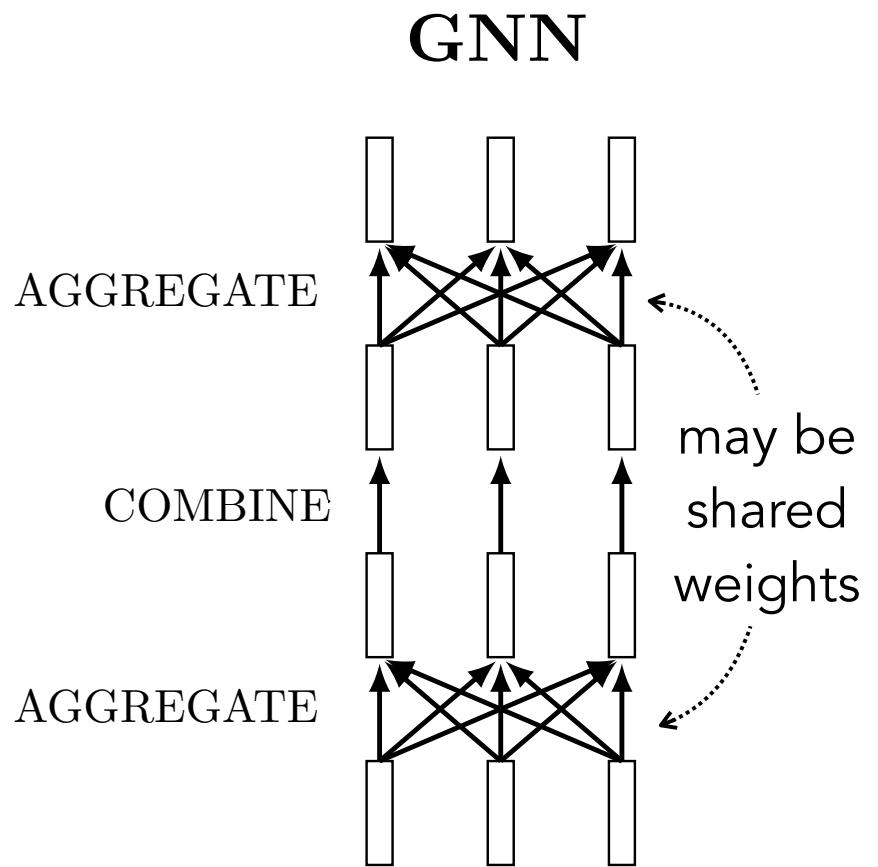


Token nets

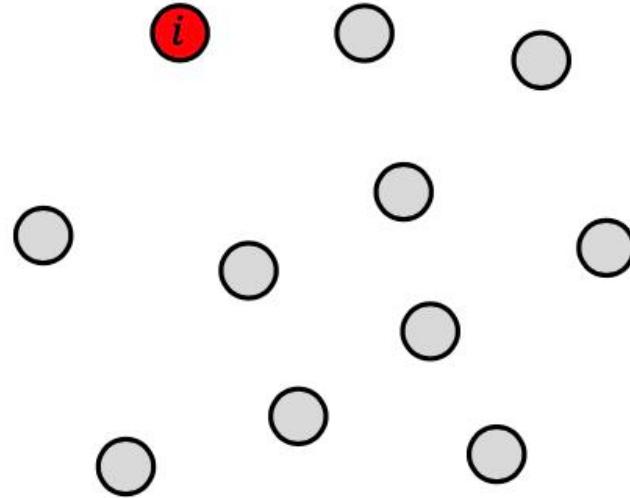


Token nets

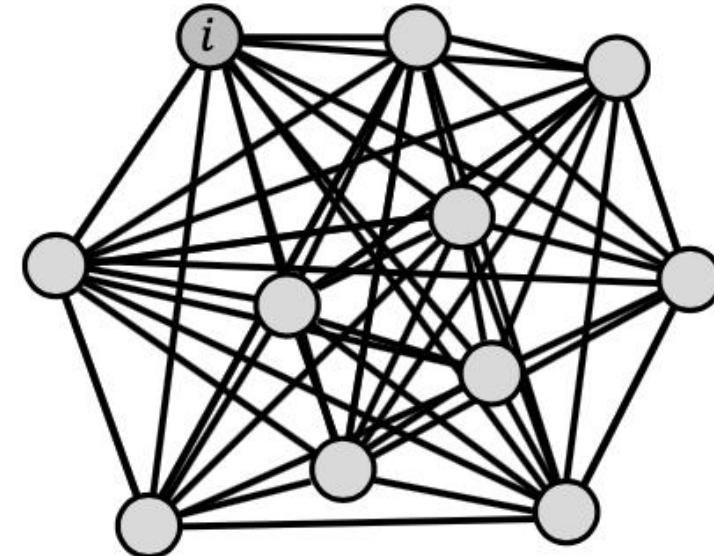




A view from the graph perspective



Sets are like graphs
without edges



Sets as complete graphs:
focus on pairwise interactions
(e.g., **transformer**: aggregation
with attention)

DeepSets and Transformer architectures may be viewed as Graph Neural Networks

Idea #2: Attention

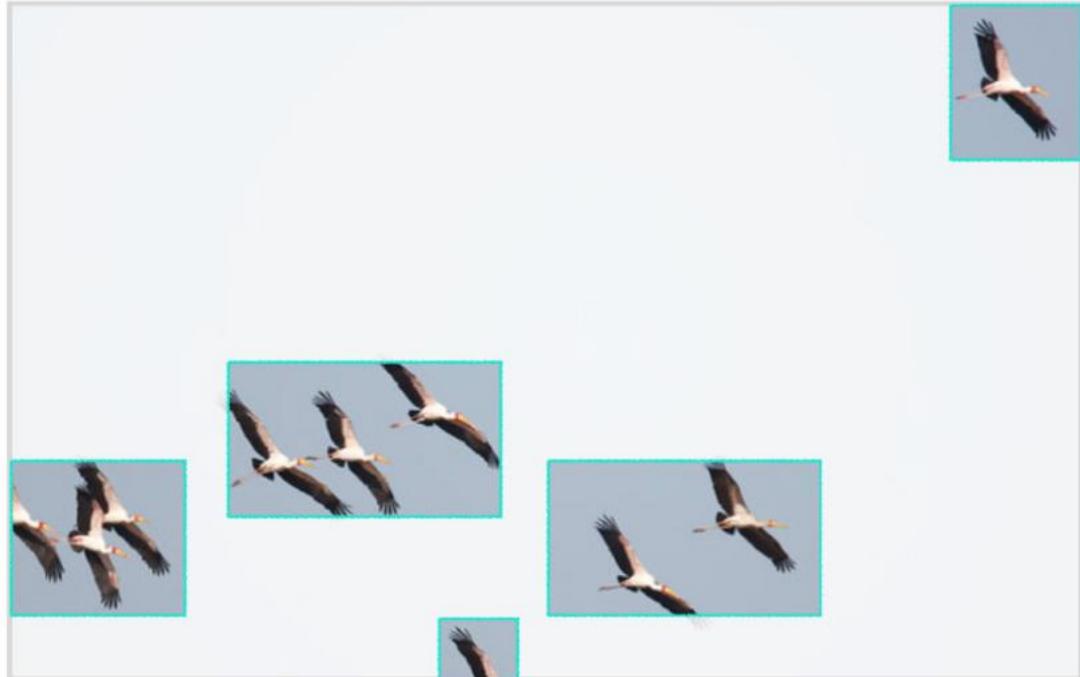
A limitation of CNNs



How many birds are in this image?

CNNs are built around the idea of locality, and are not well-suited to modeling long distance relationships

What is attention?



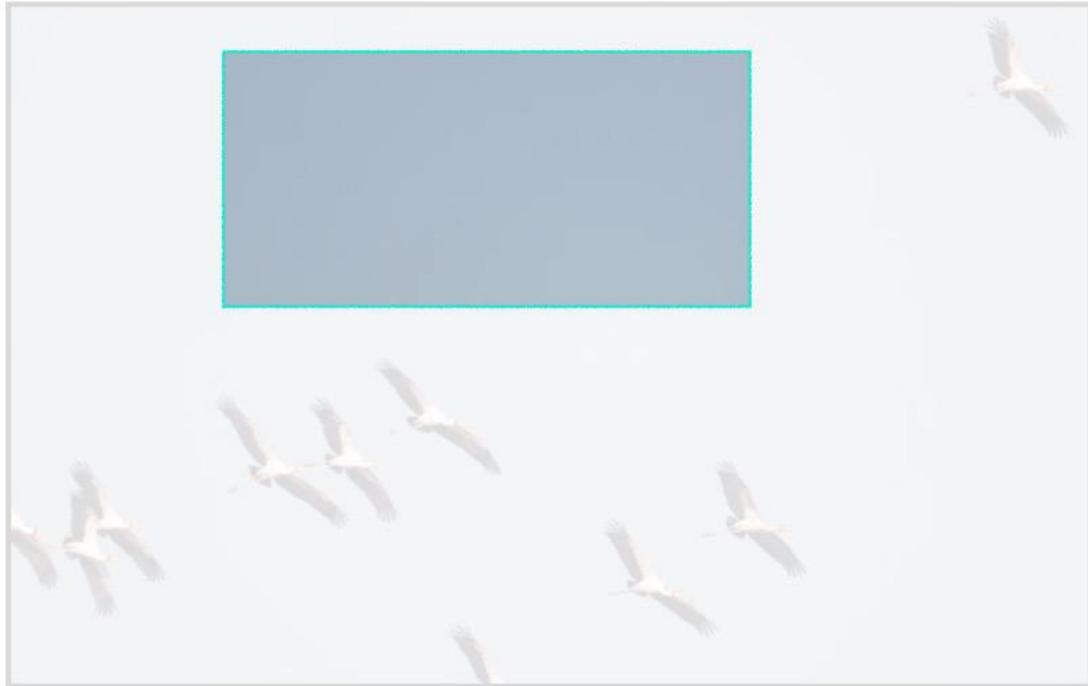
How many birds are in this image?

What is attention?



Is the top right bird the same species as the bottom left bird?

What is attention?

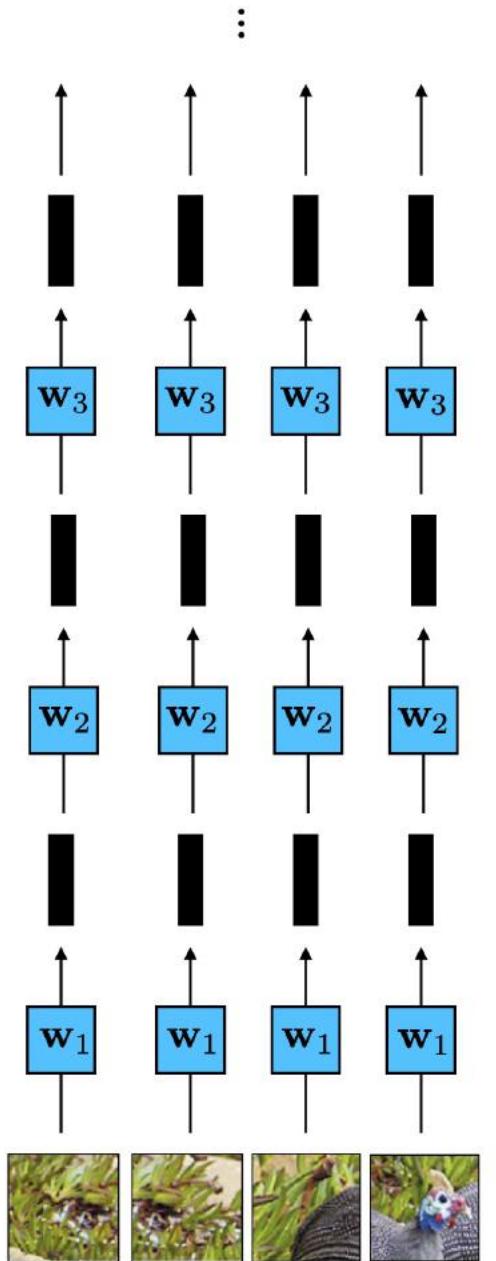


What's the color of the sky?

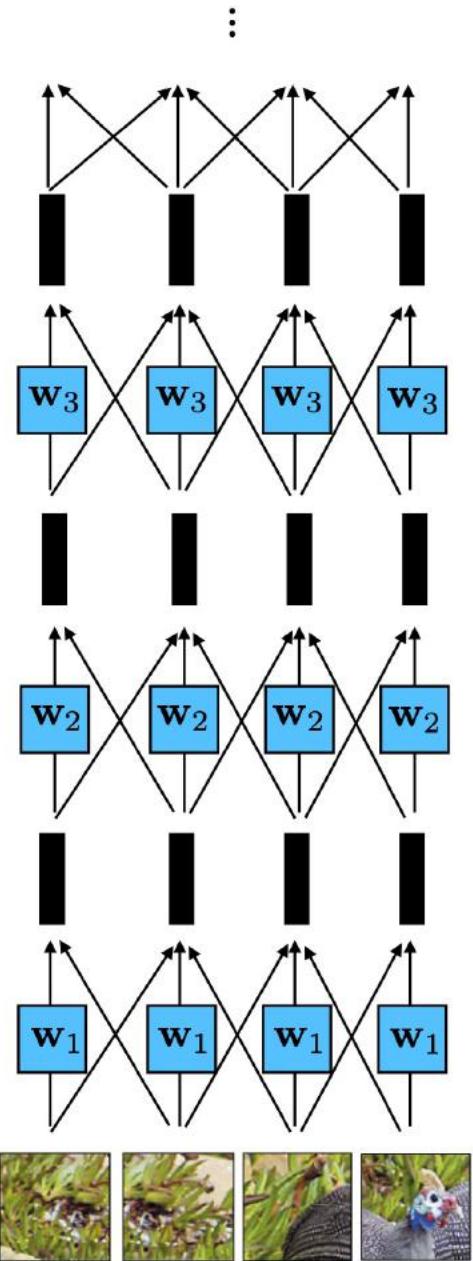
Different ways of aggregating information over space



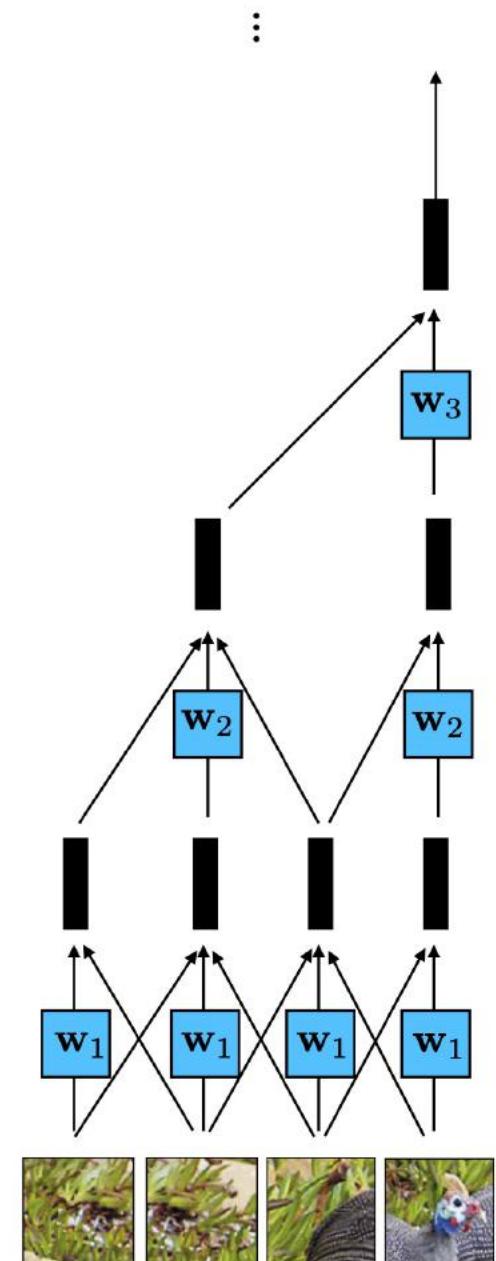
conv w/o overlap



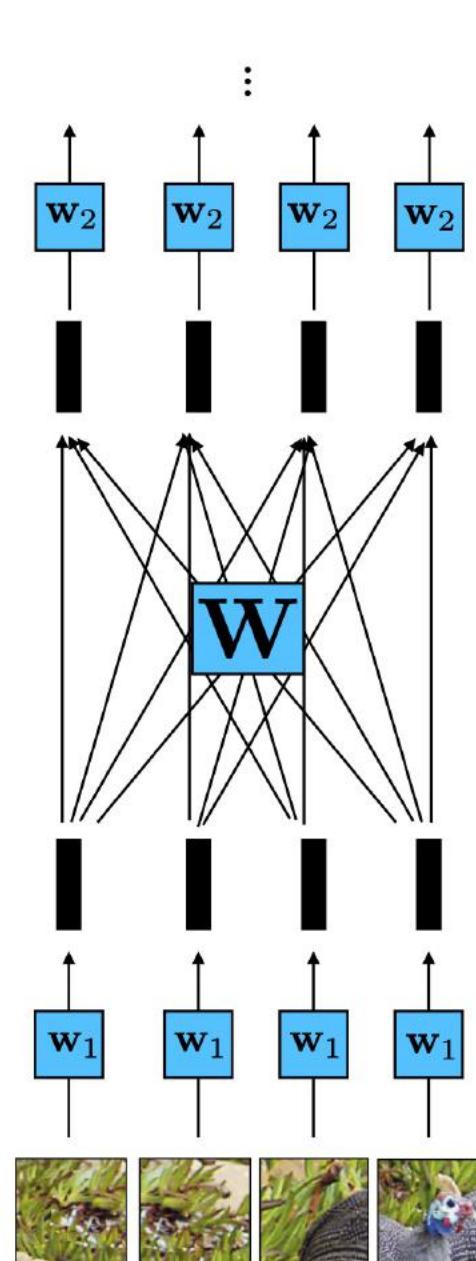
conv w overlap



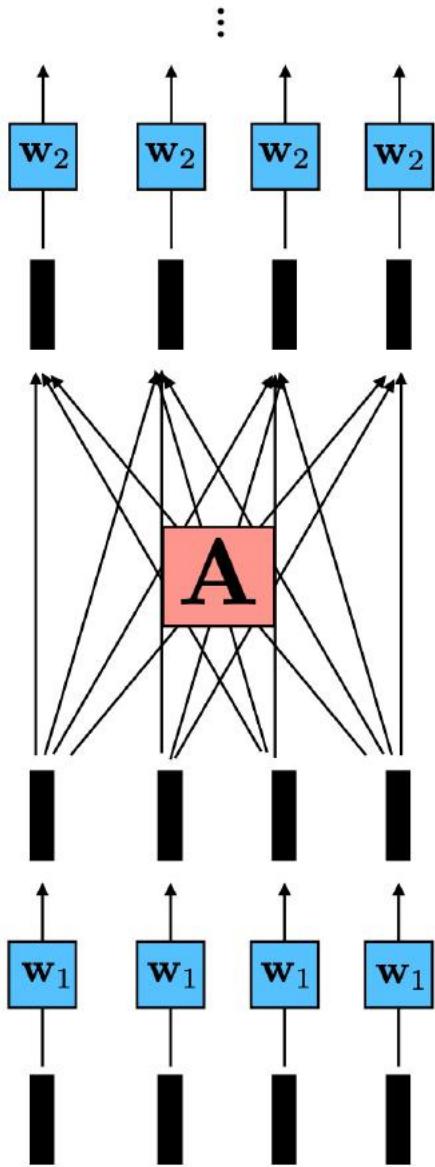
conv pyramid

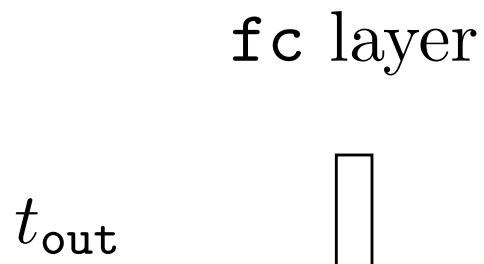


fc layer

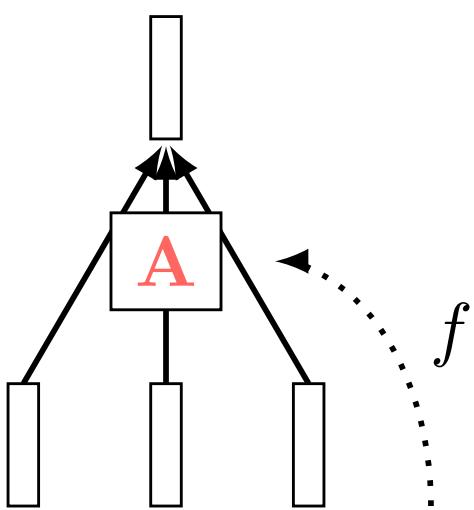


Attention Layer





attn layer



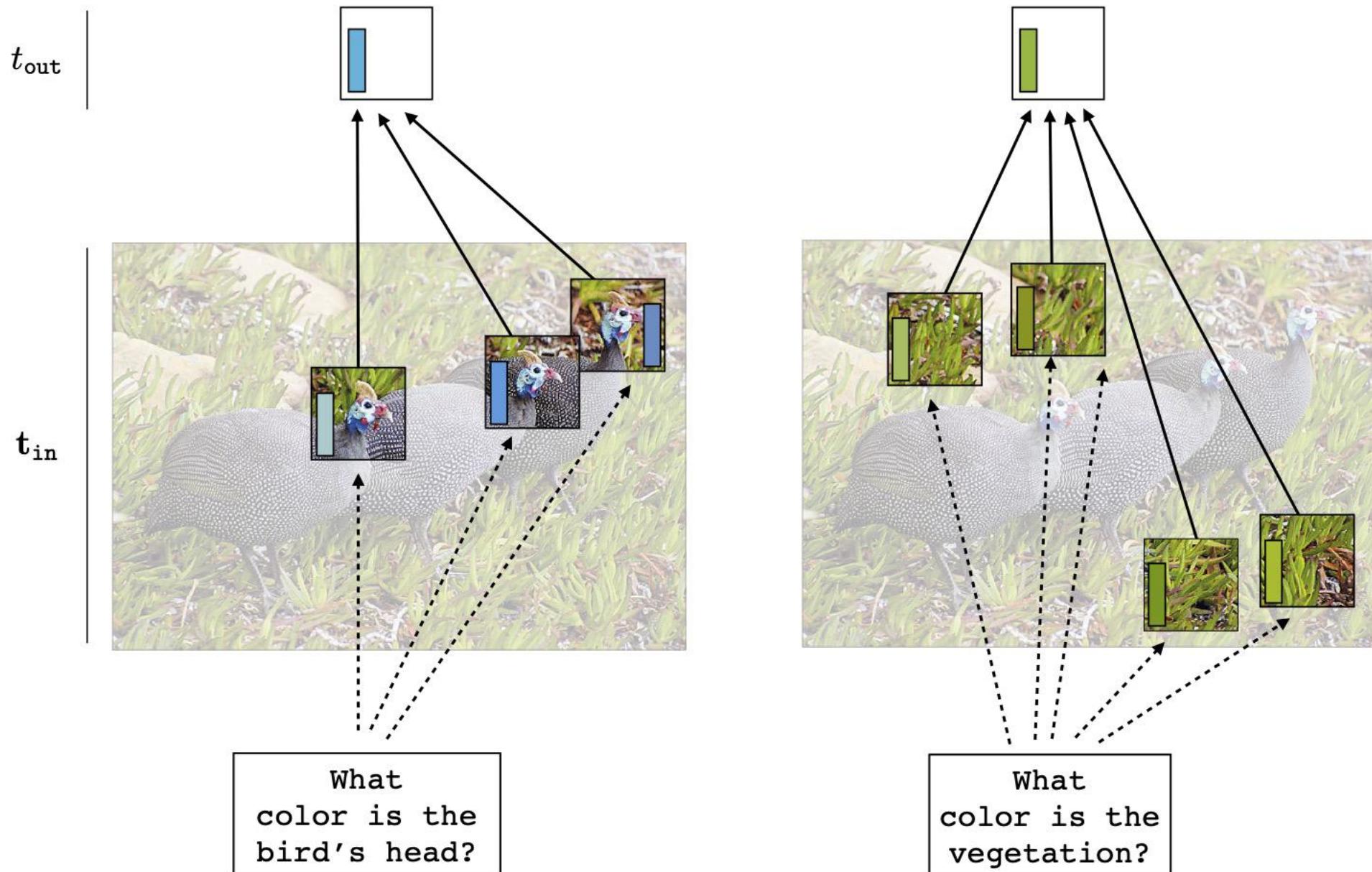
- W is free parameters
- A is a function of some input data. The data tells us which tokens to attend to (assign high weight in weighted sum)

t_{out}

t_{in}



What
color is the
bird's head?

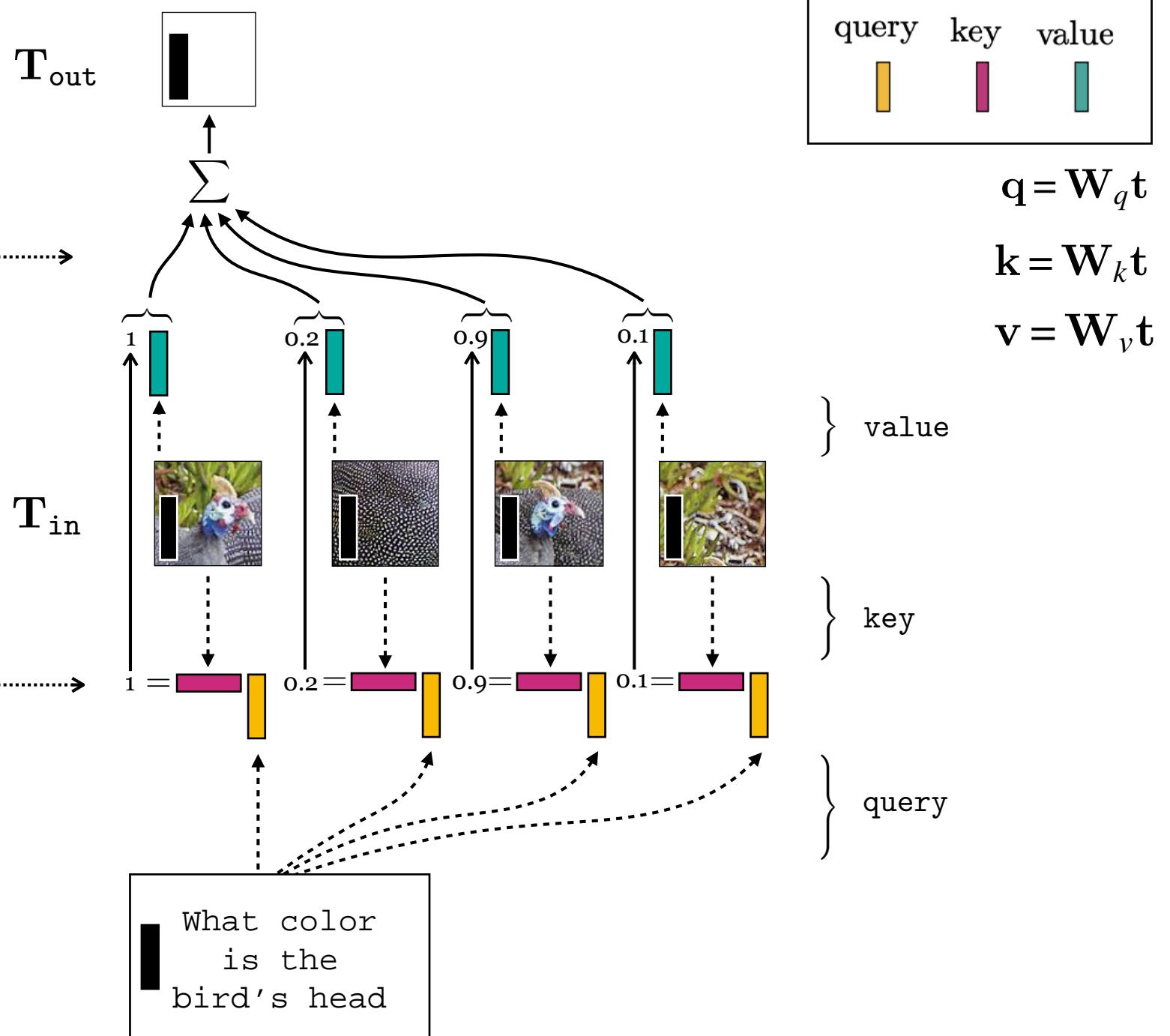


query-key-value attention

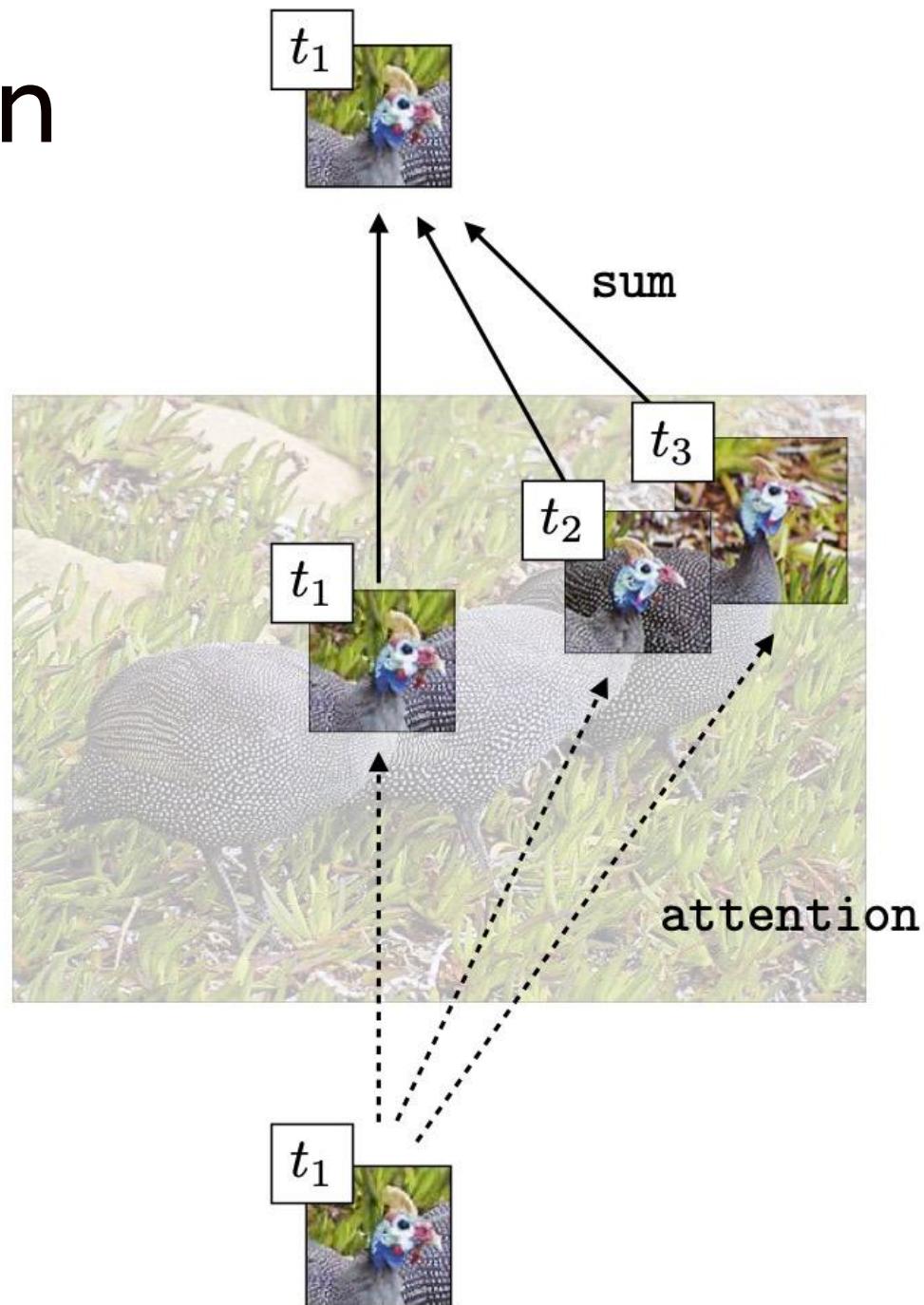
$$\mathbf{A} = \text{softmax}(\mathbf{s})$$

$$\mathbf{T}_{\text{out}} = \begin{bmatrix} a_1 \mathbf{v}_1^\top \\ \vdots \\ a_N \mathbf{v}_N^\top \end{bmatrix}$$

$$\mathbf{s} = [\mathbf{q}_{\text{question}}^T \mathbf{k}_1, \dots, \mathbf{q}_{\text{question}}^T \mathbf{k}_N]$$



Self-attention

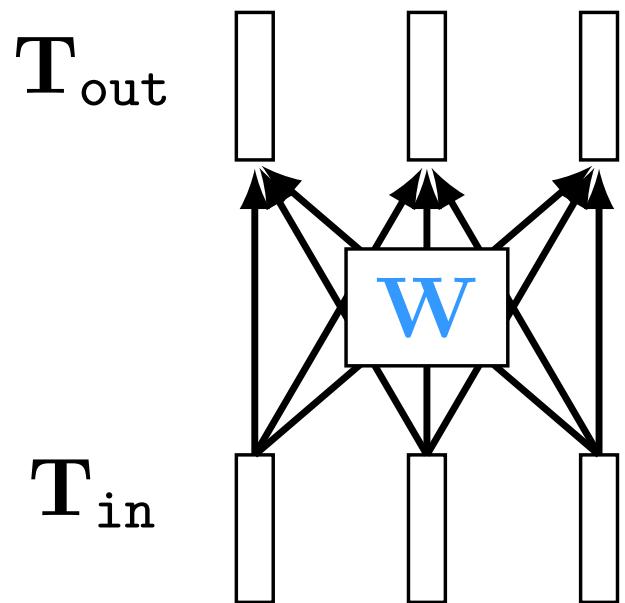


Attention Maps In A Trained Transformer

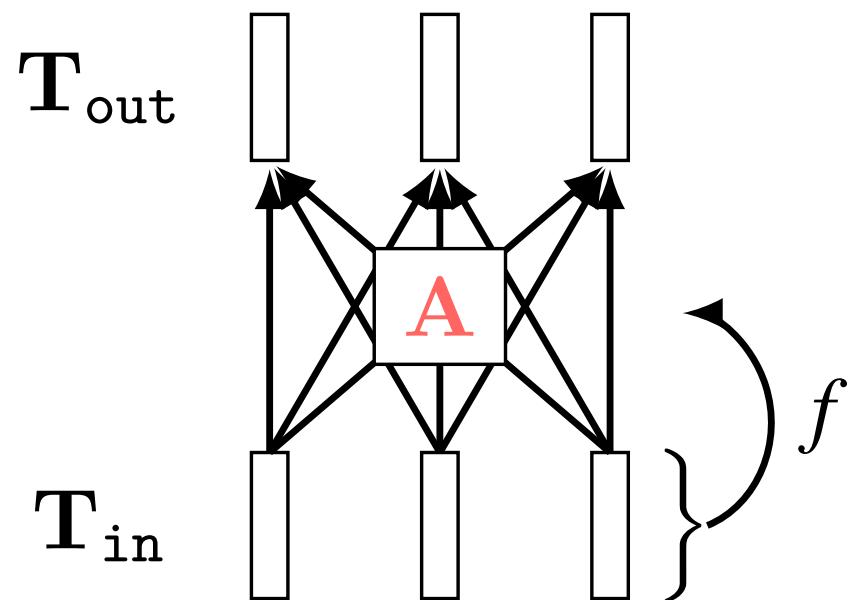


[“DINO”, Caron et all. 2021]

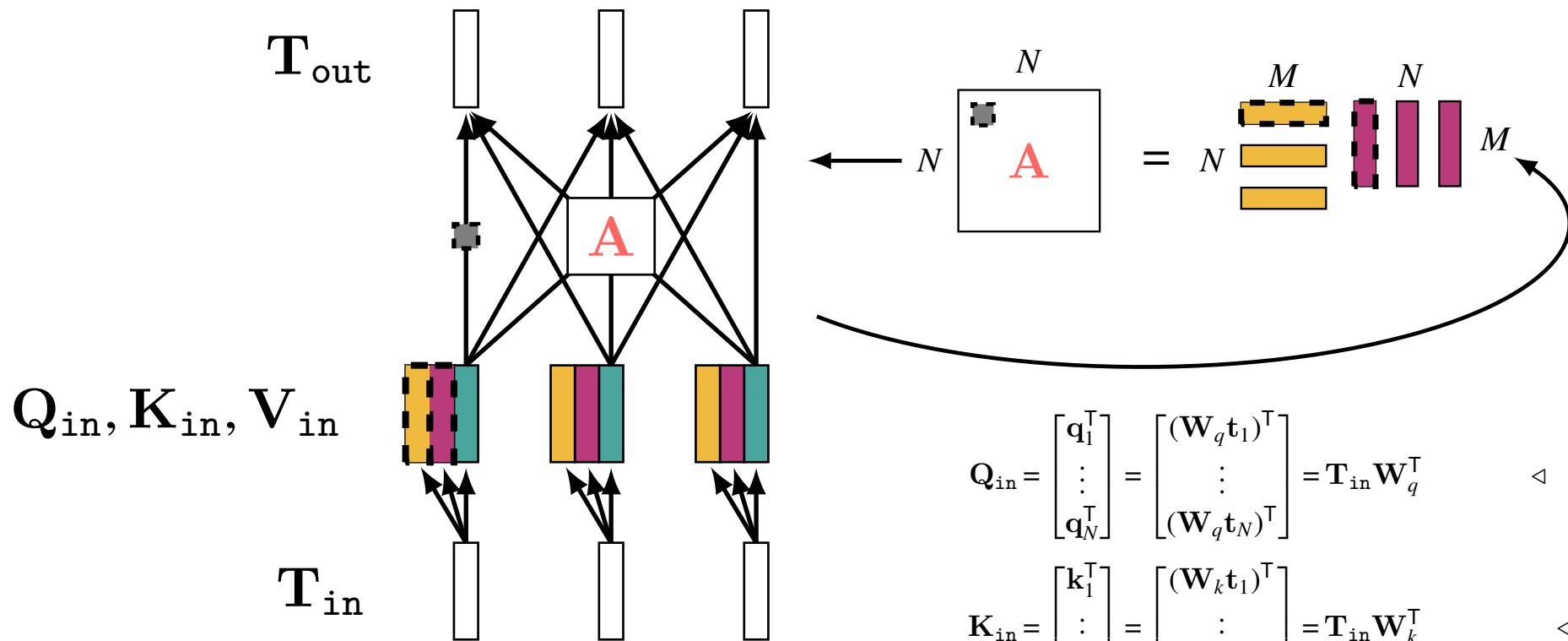
fc layer



self attn layer



self attn layer (expanded)



$$Q_{in} = \begin{bmatrix} q_1^\top \\ \vdots \\ q_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_q t_1)^\top \\ \vdots \\ (\mathbf{W}_q t_N)^\top \end{bmatrix} = T_{in} \mathbf{W}_q^\top \quad \triangleleft \quad \text{query matrix}$$

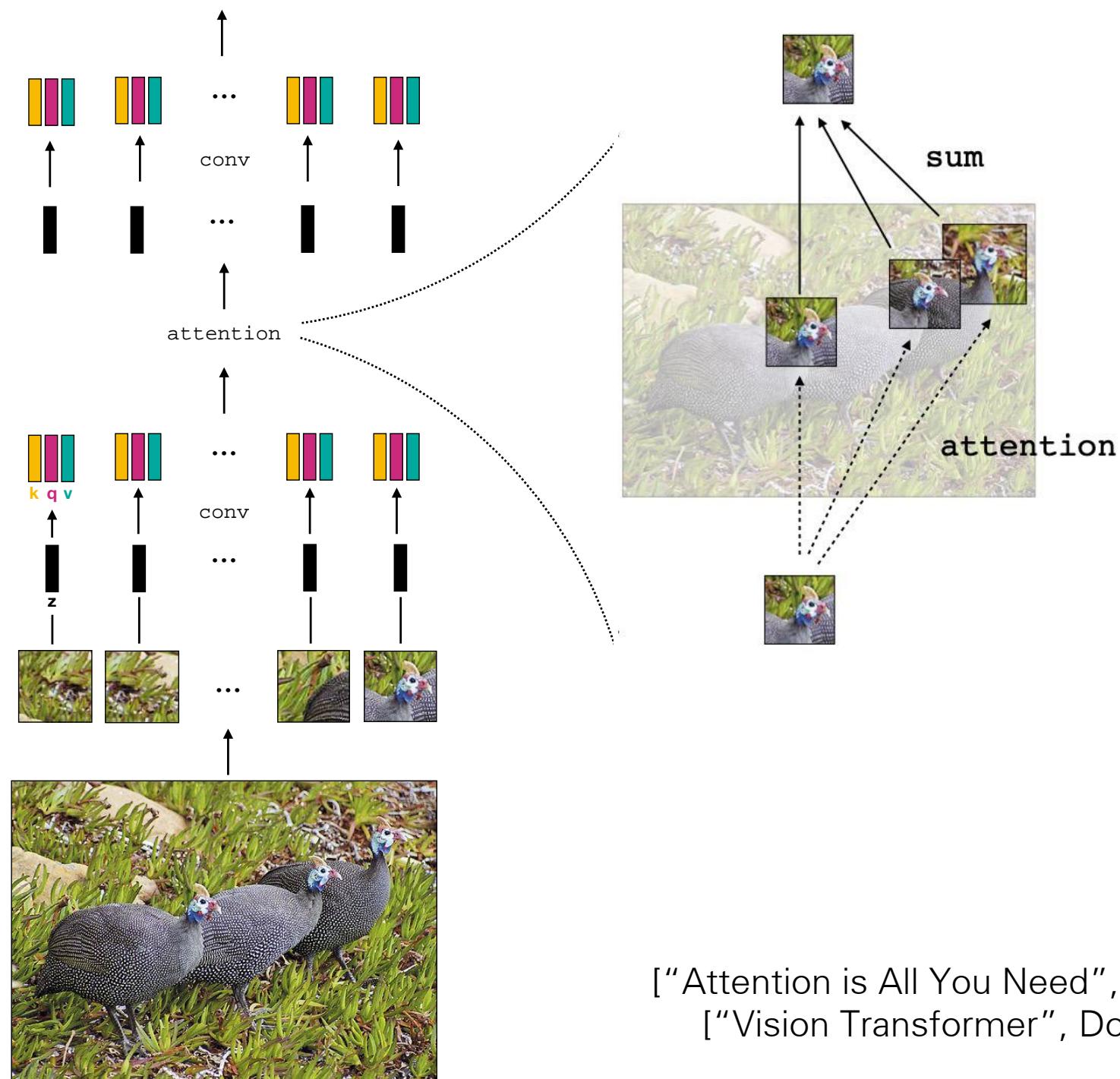
$$K_{in} = \begin{bmatrix} k_1^\top \\ \vdots \\ k_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_k t_1)^\top \\ \vdots \\ (\mathbf{W}_k t_N)^\top \end{bmatrix} = T_{in} \mathbf{W}_k^\top \quad \triangleleft \quad \text{key matrix}$$

$$V_{in} = \begin{bmatrix} v_1^\top \\ \vdots \\ v_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_v t_1)^\top \\ \vdots \\ (\mathbf{W}_v t_N)^\top \end{bmatrix} = T_{in} \mathbf{W}_v^\top \quad \triangleleft \quad \text{value matrix}$$

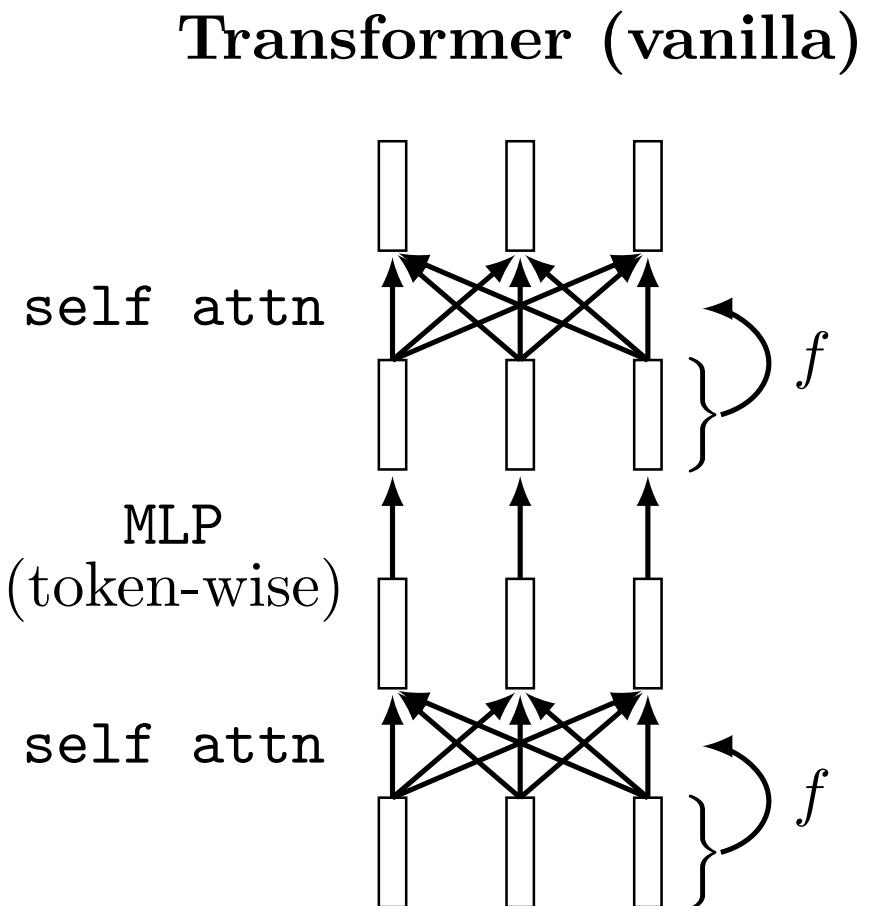
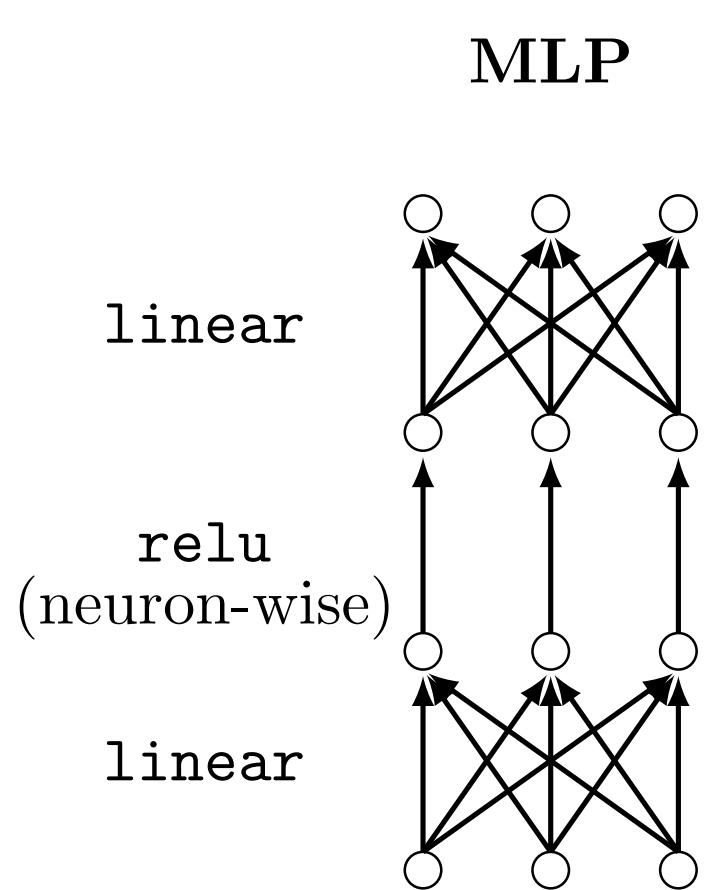
$$A = f(T_{in}) = \text{softmax}\left(\frac{Q_{in} K_{in}^\top}{\sqrt{m}}\right) \quad \triangleleft \quad \text{attention matrix}$$

$$T_{out} = A V_{in}$$

Transformer (simplified)



[“Attention is All You Need”, Vaswani et al. 2017]
[“Vision Transformer”, Dosovitskiy et al. 2020]



Multihead self-attention (MSH)

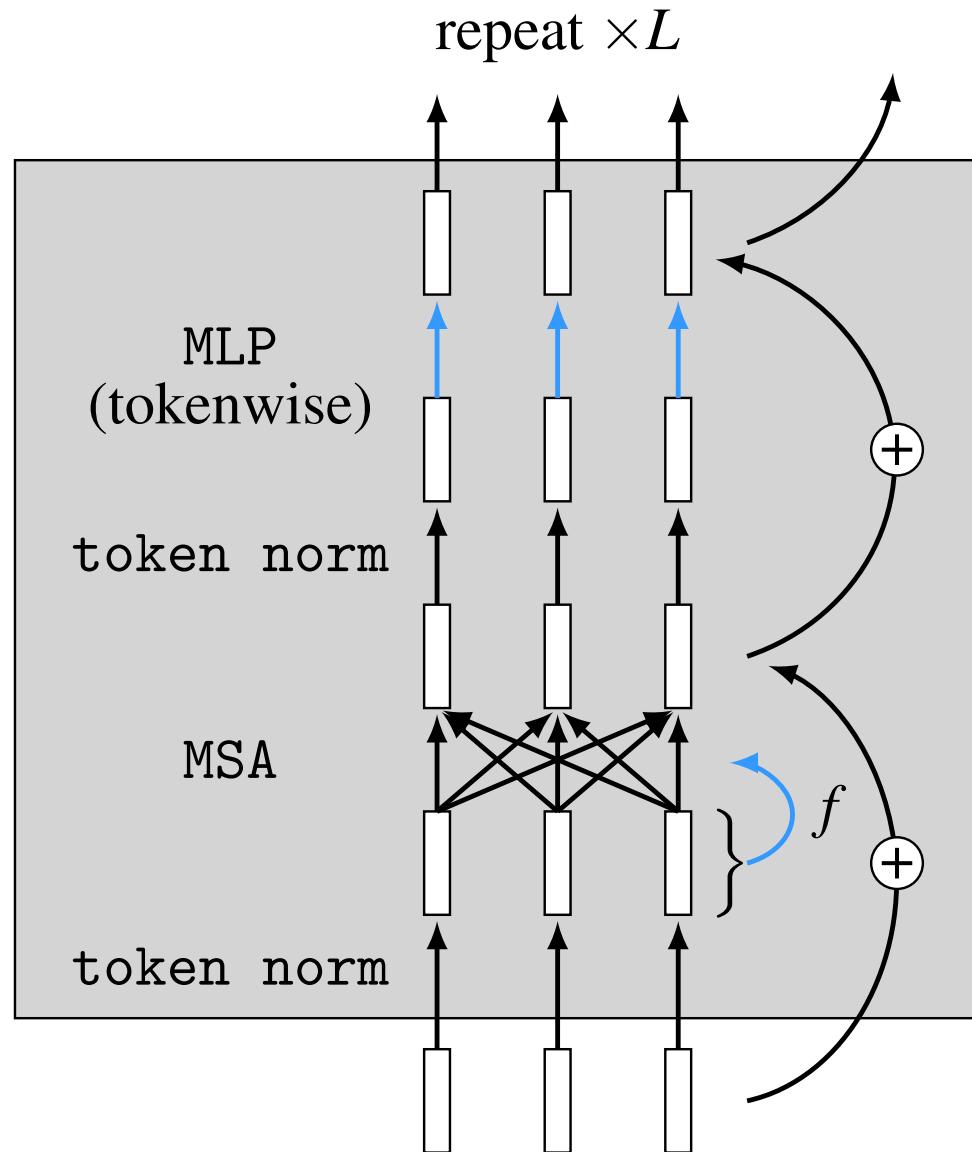
- Rather than having just one way of attending, why not have k ?
- Each gets its own parameterized query(), key(), value() functions.
- Run them all in parallel, then (weighted) sum the output token code vectors

$$\mathbf{T}_{\text{out}}^i = \text{attn}^i(\mathbf{T}_{\text{in}}) \quad \forall i$$

$$\bar{\mathbf{T}}_{\text{out}} = \begin{bmatrix} \mathbf{T}_{\text{out}}^1[0, :] & \dots & \mathbf{T}_{\text{out}}^k[0, :] \\ \vdots & \vdots & \vdots \\ \mathbf{T}_{\text{out}}^1[N-1, :] & \dots & \mathbf{T}_{\text{out}}^k[N-1, :] \end{bmatrix} \quad \triangleleft \quad \bar{\mathbf{T}}_{\text{out}} \in \mathbb{R}^{N \times kv}$$

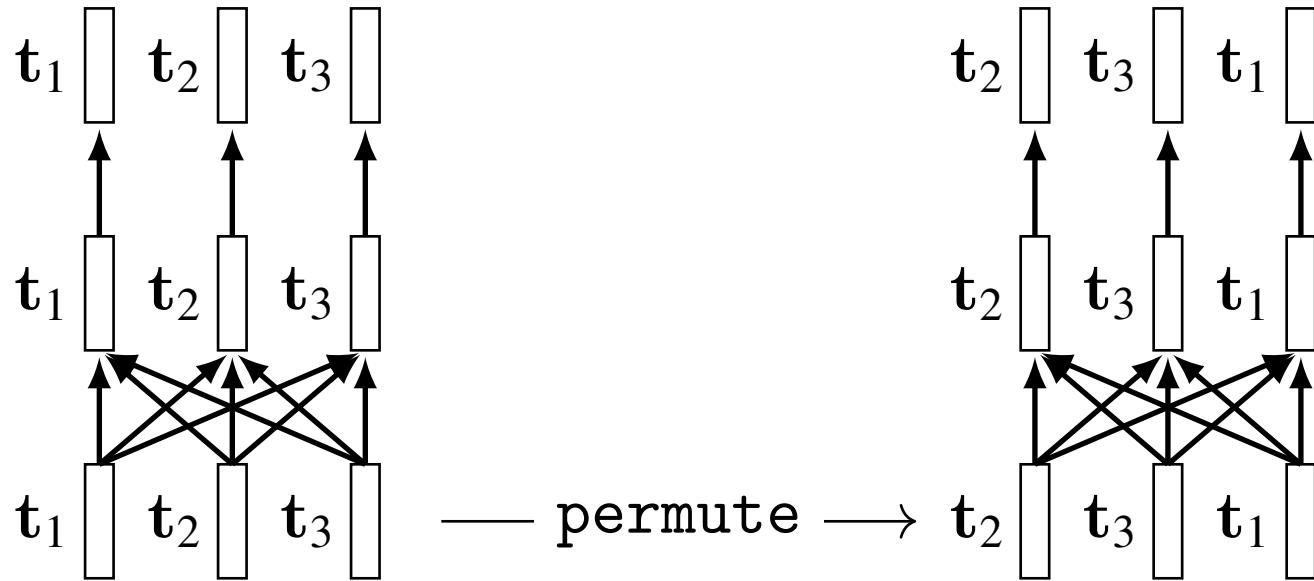
$$\mathbf{T}_{\text{out}} = \bar{\mathbf{T}}_{\text{out}} \mathbf{W} \quad \triangleleft \quad \mathbf{W} \in \mathbb{R}^{kv \times d}$$

Transformer (ViT)



Idea #3: Positional Encoding

Permutation equivariance



$$F_\theta(\text{permute}(\mathbf{T}_{\text{in}})) = \text{permute}(F_\theta(\mathbf{T}_{\text{in}}))$$

$$\text{attn}(\text{permute}(\mathbf{T}_{\text{in}})) = \text{permute}(\text{attn}(\mathbf{T}_{\text{in}}))$$

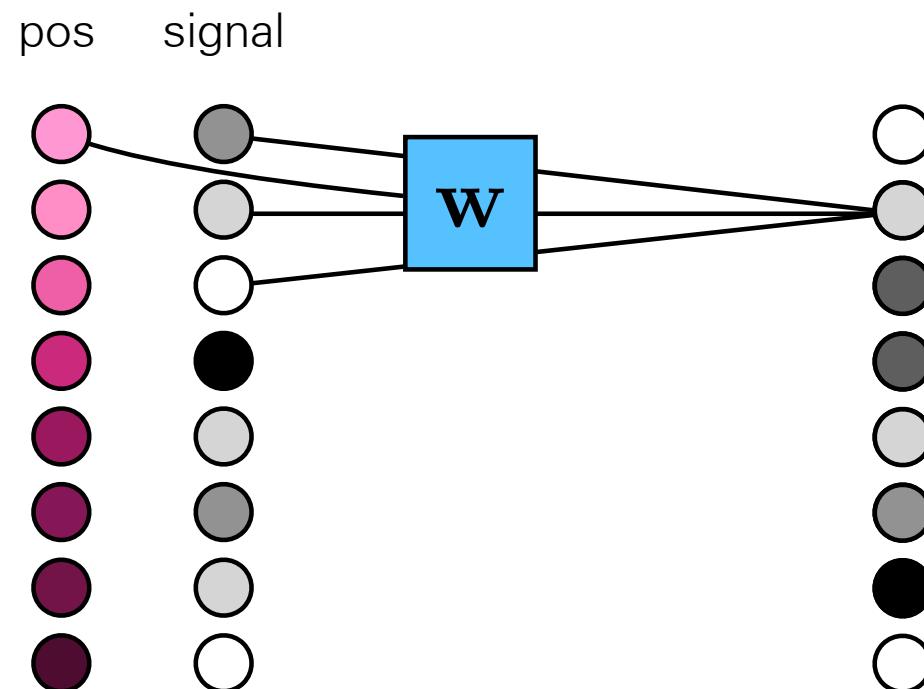
Set2Set



$$\text{transformer}(\text{permute}(\mathbf{T}_{\text{in}})) = \text{permute}(\text{transformer}(\mathbf{T}_{\text{in}}))$$

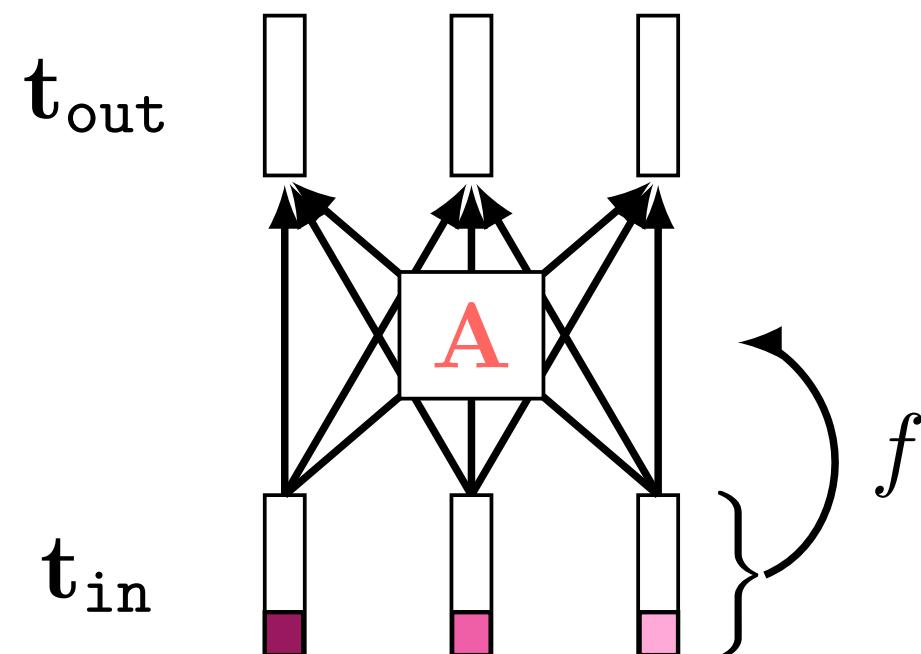
What if you don't want to be shift invariant?

1. Use an architecture that is not permutation invariant (e.g., MLP)
2. Add location information to the token code vectors — this is called **positional encoding**



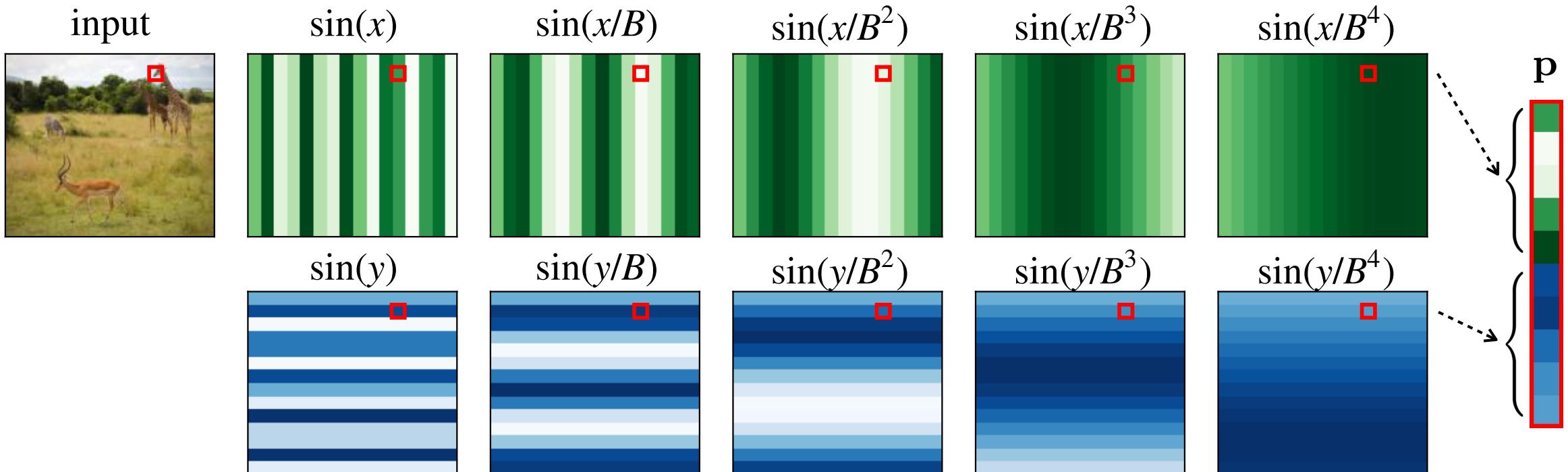
What if you don't want to be shift invariant?

1. Use an architecture that is not permutation invariant (e.g., MLP)
2. Add location information to the token code vectors — this is called **positional encoding**



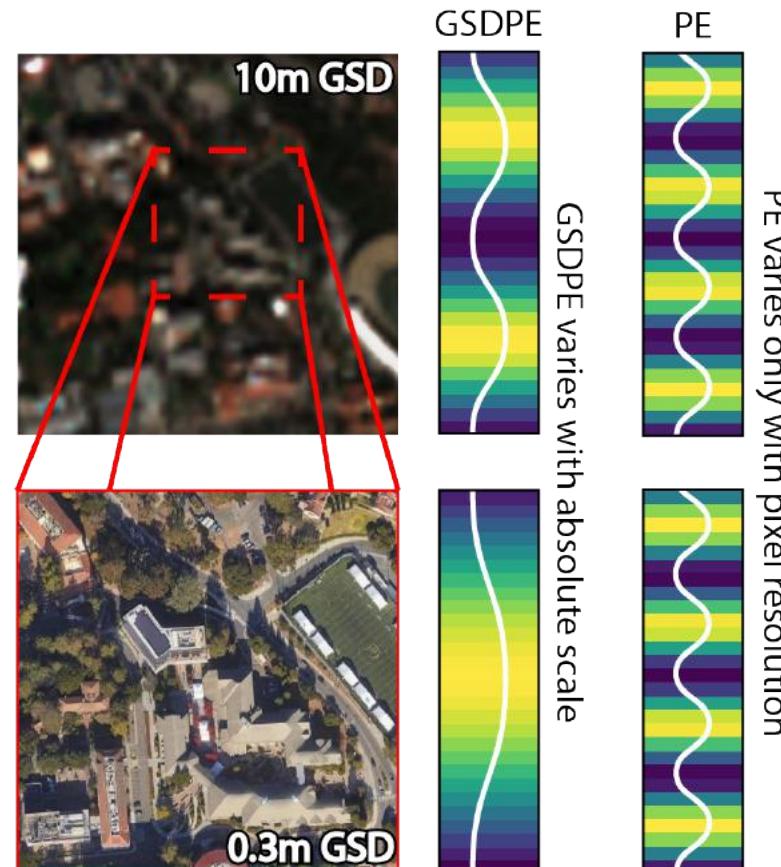
Fourier positional codes

- Represent coordinates on Fourier basis



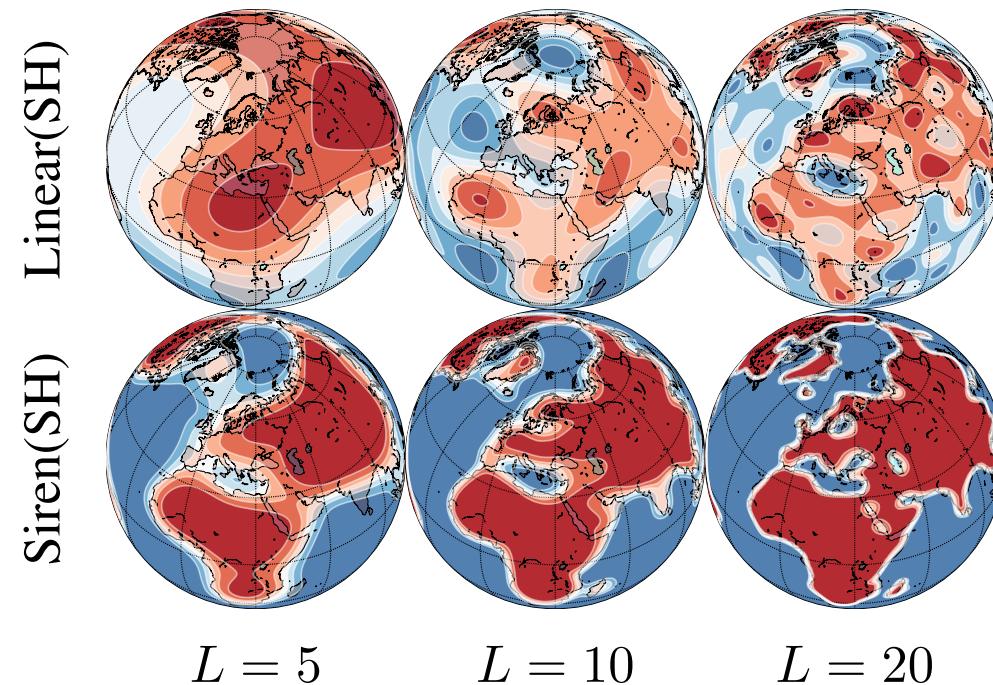
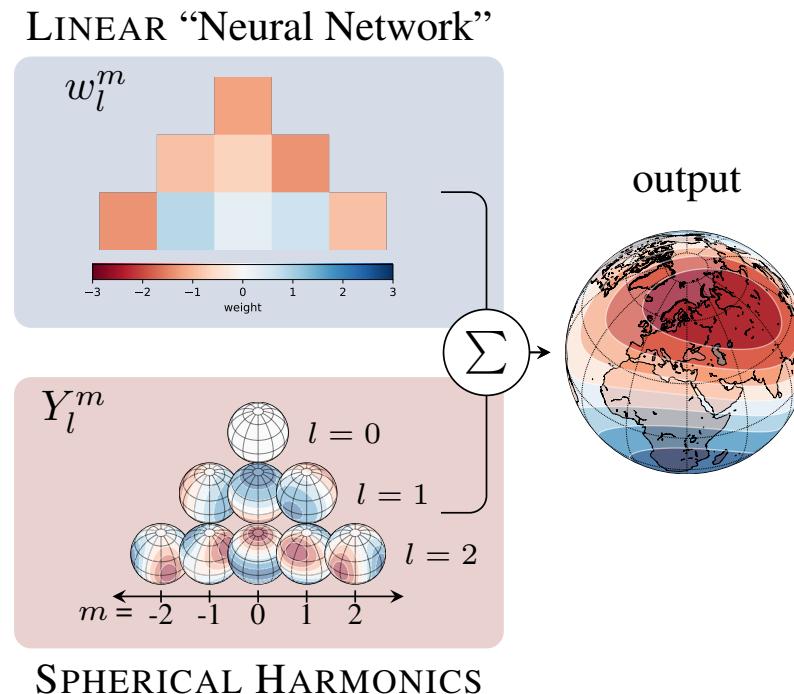
Other positional encodings

- ScaleMAE uses ground sample distance positional encoding to train a MAE across spatial scales of remote sensing data



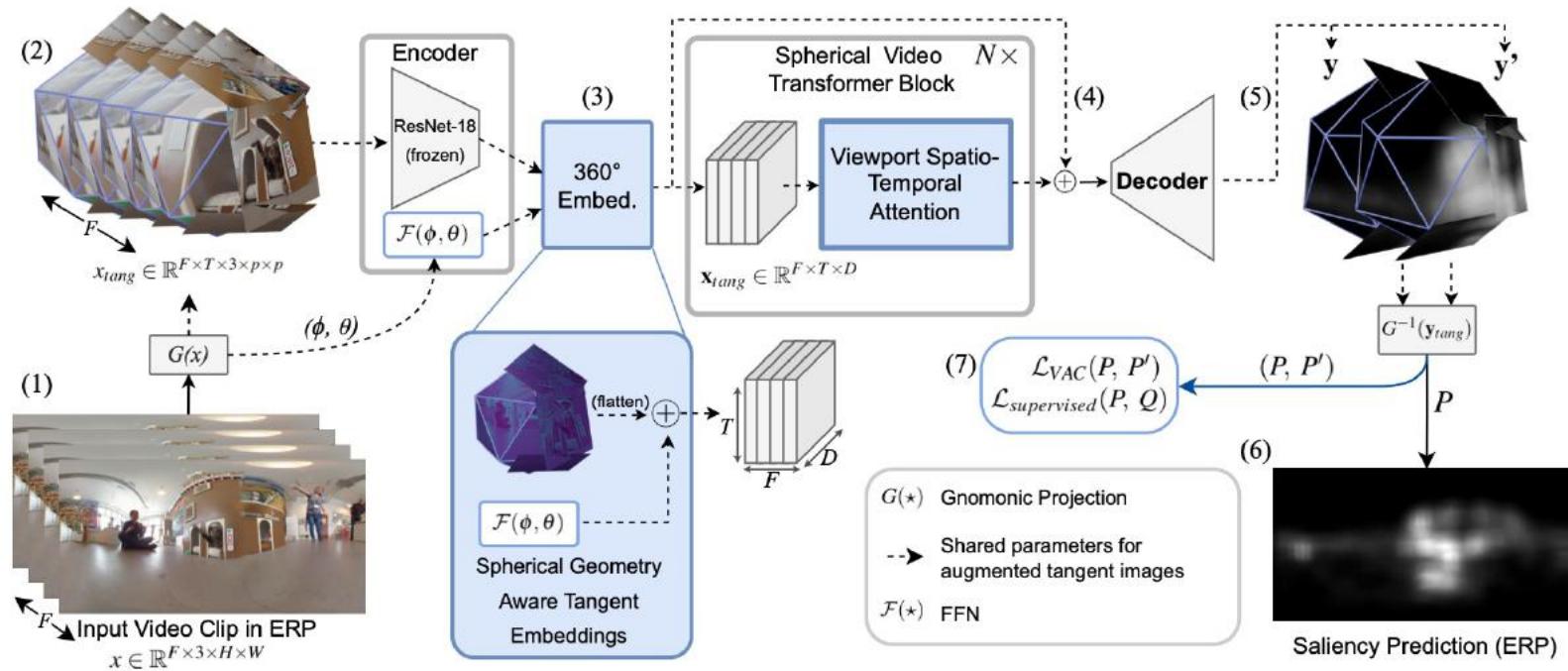
Other positional encodings

- Geographic location encoding with spherical harmonics and sinusoidal representation networks



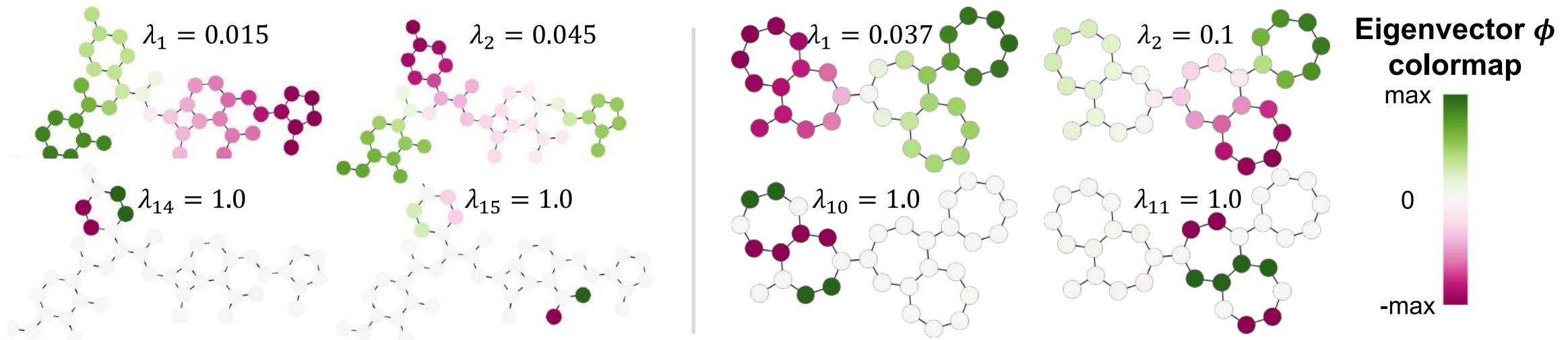
Other positional encodings

- Each tangent image is encoded and fused with spherical-geometry-aware position embeddings



Other positional encodings

- Laplacian positional encodings to encode node positions in a graph



Some fancy architectures and applications

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^{*†}

^{*}equal technical contribution, [†]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, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.¹

1 INTRODUCTION

Self-attention-based architectures, in particular Transformers (Vaswani et al., 2017), have become the model of choice in natural language processing (NLP). The dominant approach is to pre-train on a large text corpus and then fine-tune on a smaller task-specific dataset (Devlin et al., 2019). Thanks to Transformers' computational efficiency and scalability, it has become possible to train models of unprecedented size, with over 100B parameters (Brown et al., 2020; Lepikhin et al., 2020). With the models and datasets growing, there is still no sign of saturating performance.

<https://arxiv.org/abs/2010.11929>

[lucidrains / vit-pytorch](https://github.com/lucidrains/vit-pytorch) Public

Sponsor Watch 125 Fork 1.8k Star 10.8k

Code Issues 89 Pull requests 4 Actions Projects Wiki Security Insights

main 2 branches 143 tags Go to file Add file Code

	lucidrains offer way for extractor to return latents without detaching them	✓ f86e052 4 days ago 249 commits
	.github sponsor button	4 months ago
	examples fix transforms for val an test process	11 months ago
	images add EsViT, by popular request, an alternative to Dino that is compati...	3 months ago
	tests add some tests	7 months ago
	vit_pytorch offer way for extractor to return latents without detaching them	4 days ago
	.gitignore Initial commit	2 years ago
	LICENSE Initial commit	2 years ago
	MANIFEST.in include tests in package for conda	7 months ago
	README.md make extractor flexible for layers that output multiple tensors, show...	last month
	setup.py offer way for extractor to return latents without detaching them	4 days ago

About

Implementation of Vision Transformer, a simple way to achieve SOTA in vision classification with only a single transformer encoder, in Pytorch

computer-vision transformers
artificial-intelligence image-classification
attention-mechanism

Readme MIT license 10.8k stars 125 watching 1.8k forks

Releases 142

<https://github.com/lucidrains/vit-pytorch>

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

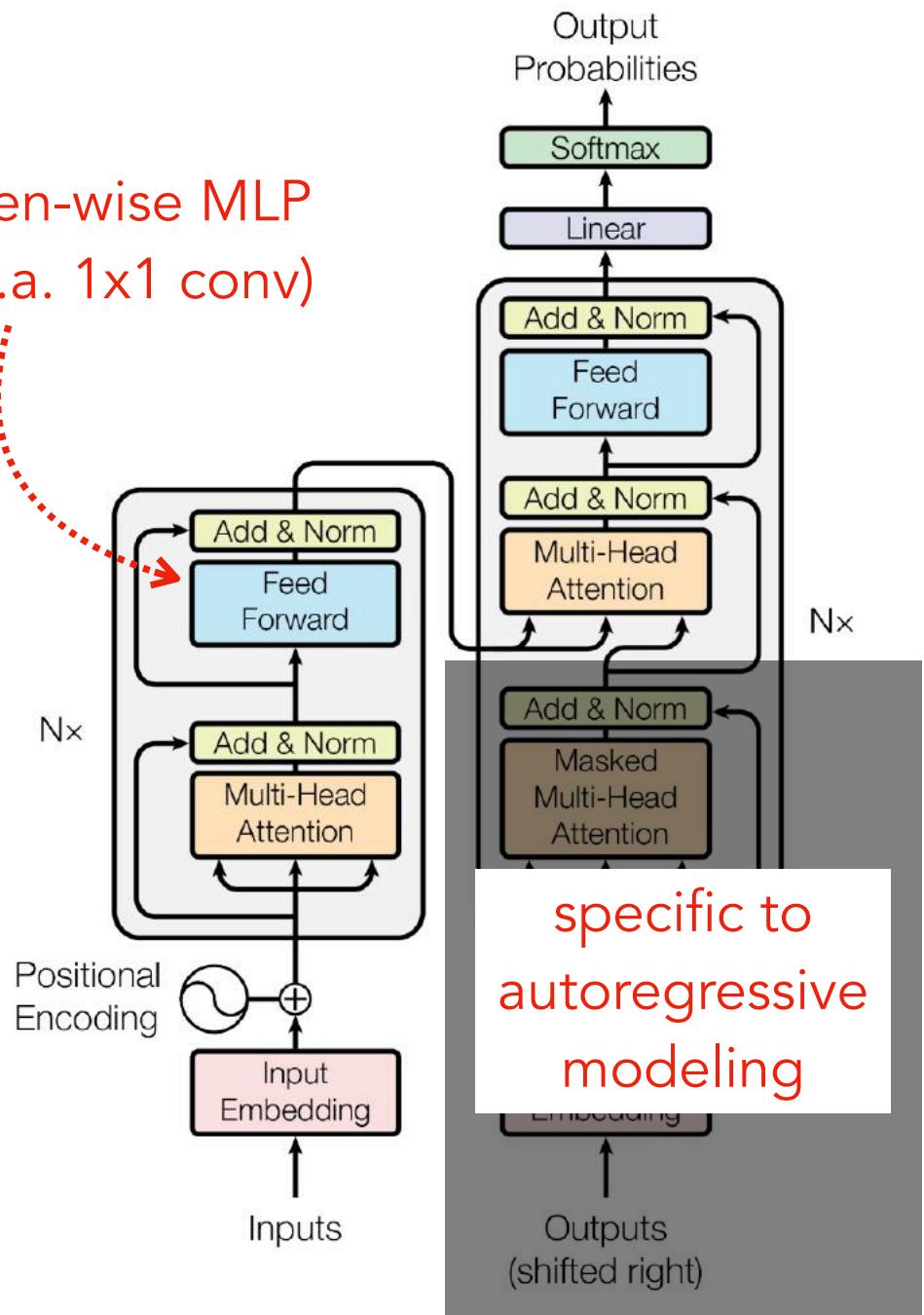
Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com

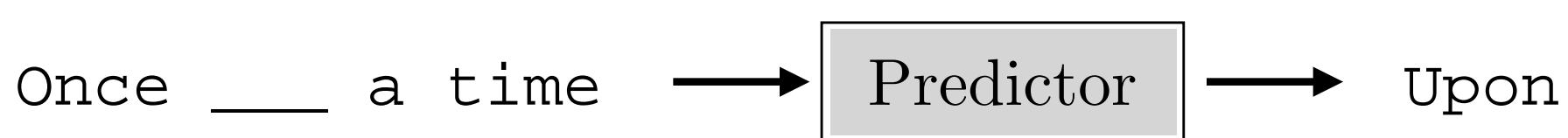
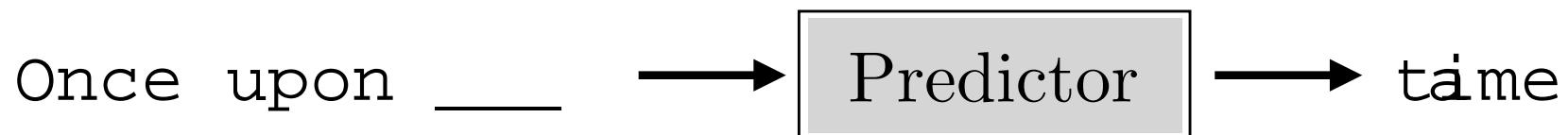
Abstract

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 mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

token-wise MLP
(a.k.a. 1x1 conv)

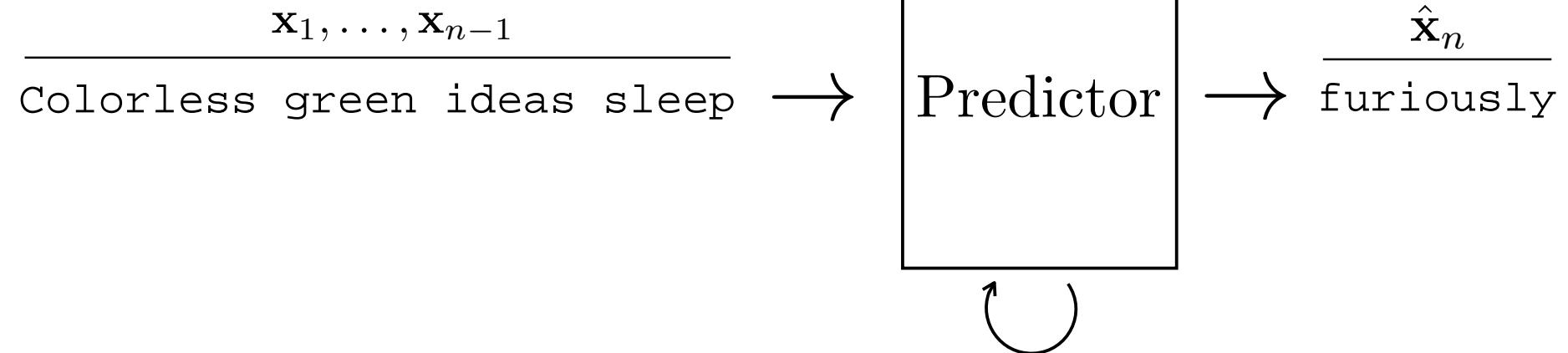
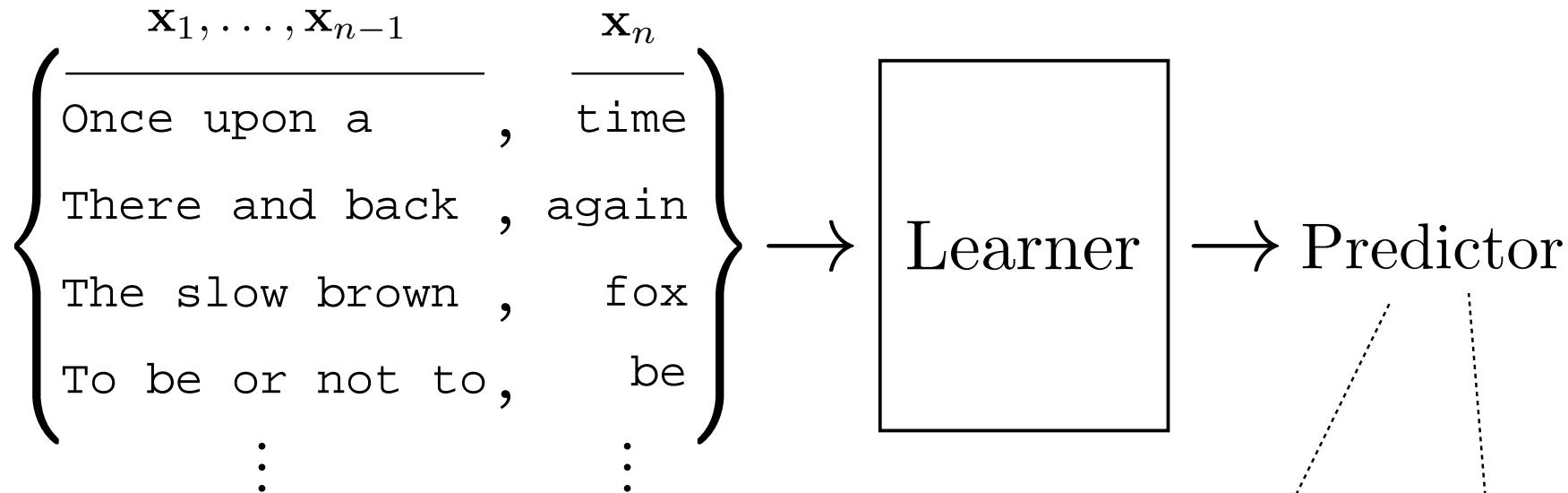


Autoregressive models

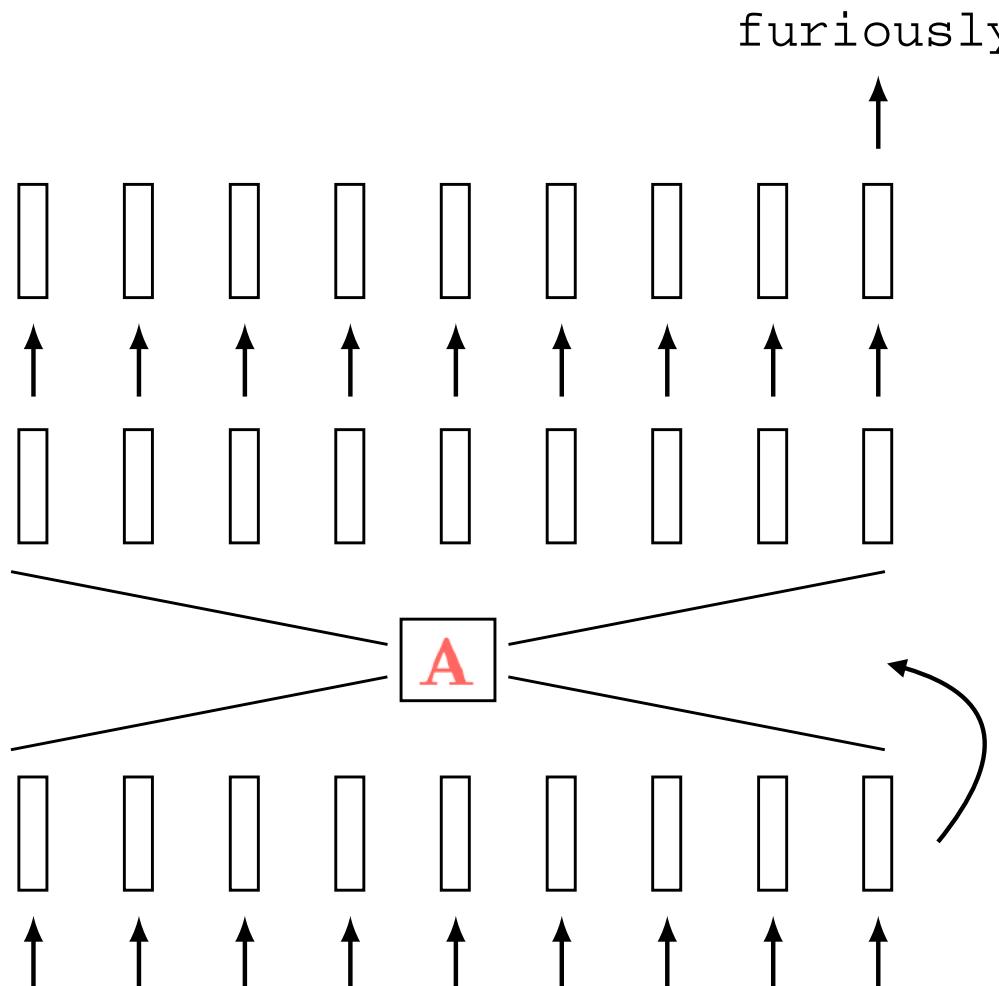


Sampling

Training



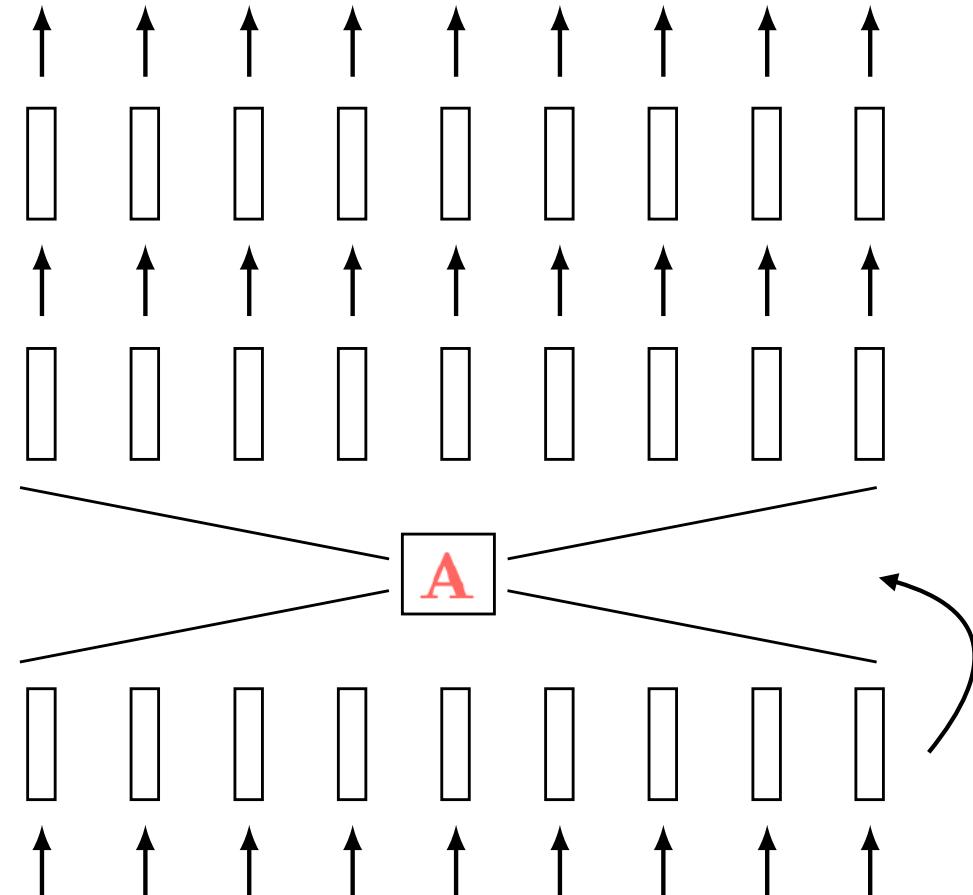
GPT (and many other related models)



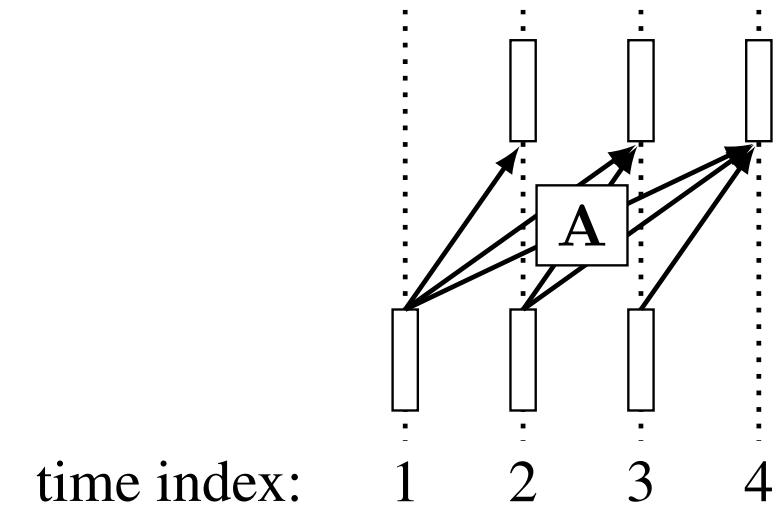
Colorless green ideas sleep __

GPT training (and many other related models)

Colorless green ideas sleep furiously



Colorless green ideas sleep furiously



$$A \quad T_{\text{in}} \quad T_{\text{out}} =$$

The diagram shows the computation of hidden states. On the left, a large black rectangle labeled "A" represents the weight matrix. To its right is a vertical stack of three white rectangles labeled "T_{in}". An equals sign follows this. To the right of the equals sign is another vertical stack of three white rectangles labeled "T_{out}". This visualizes the equation $A T_{\text{in}} = T_{\text{out}}$, where the weight matrix A is multiplied by the input vector T_{in} to produce the output vector T_{out} .

master ▾

3 branches

0 tags

Go to file

Add file ▾

Code ▾



karpathy Merge pull request #84 from ericjang/master ...

7218bcf on Aug 4

93 commits

mingpt	Use XOR operator ^ for checking assertion `type_given XOR param...	2 months ago
projects	refactor sequence generation into the model and match the huggingf...	3 months ago
tests	add a refactored BPE encoder from openai. Basically I dont super tru...	3 months ago
.gitignore	tiny tweaks to printing and some function apis	4 months ago
LICENSE	mit license file	2 years ago
README.md	Add setup.py to allow mingpt to be used as a third-party library	2 months ago
demo.ipynb	refactor sequence generation into the model and match the huggingf...	3 months ago
generate.ipynb	add a refactored BPE encoder from openai. Basically I dont super tru...	3 months ago
mingpt.jpg	first commit, able to multigpu train fp32 GPTs on math and character...	2 years ago
setup.py	Add setup.py to allow mingpt to be used as a third-party library	2 months ago

README.md

minGPT

Summary

- attention is used to focus on parts of inputs/outputs
- it can be content/location based and hard/soft
- its three main distinct uses are
 - connecting encoder and decoder in sequence-to-sequence task
 - achieving scale-invariance and focus in image processing
 - self-attention can be a basic building block for neural nets, often replacing RNNs and CNNs [recent research, take it with a grain of salt]

Next lecture:
Autoregressive Models