

Transformers

Attention is all you need

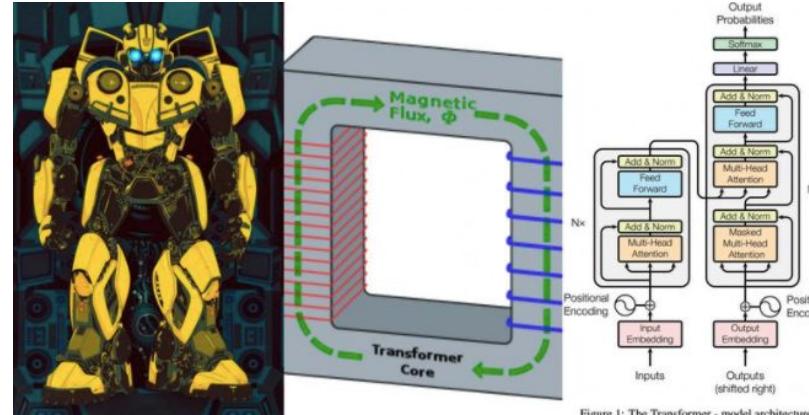


Figure 1: The Transformer - model architecture.

Transformers
at school Transformers
at college Transformers
today

Javier Selva Castelló

Outline

- Motivation
 - No more recurrence
 - Attention Mechanisms
- The transformer
 - Overview
 - Queries, Keys and Values
 - Dot-product attention
- Vision Transformers
 - Pre-training
 - Efficient Designs
- Application
 - Image Transformers
 - Video Transformers

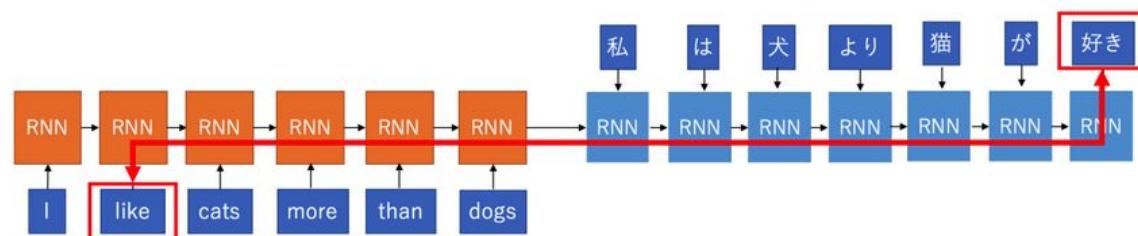
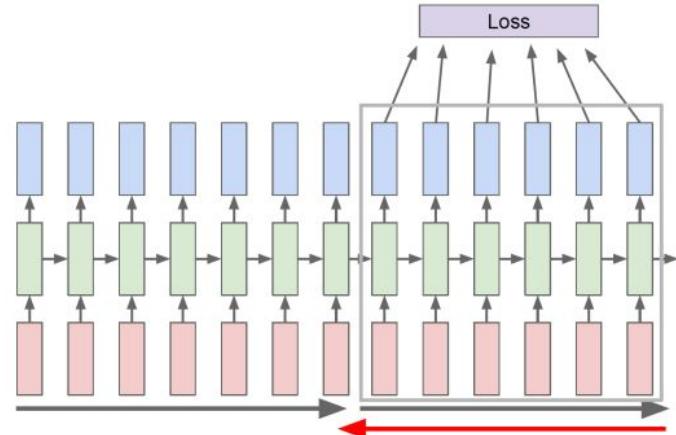
Motivation

No more recurrence

Recurrent nature does not allow for parallelism.

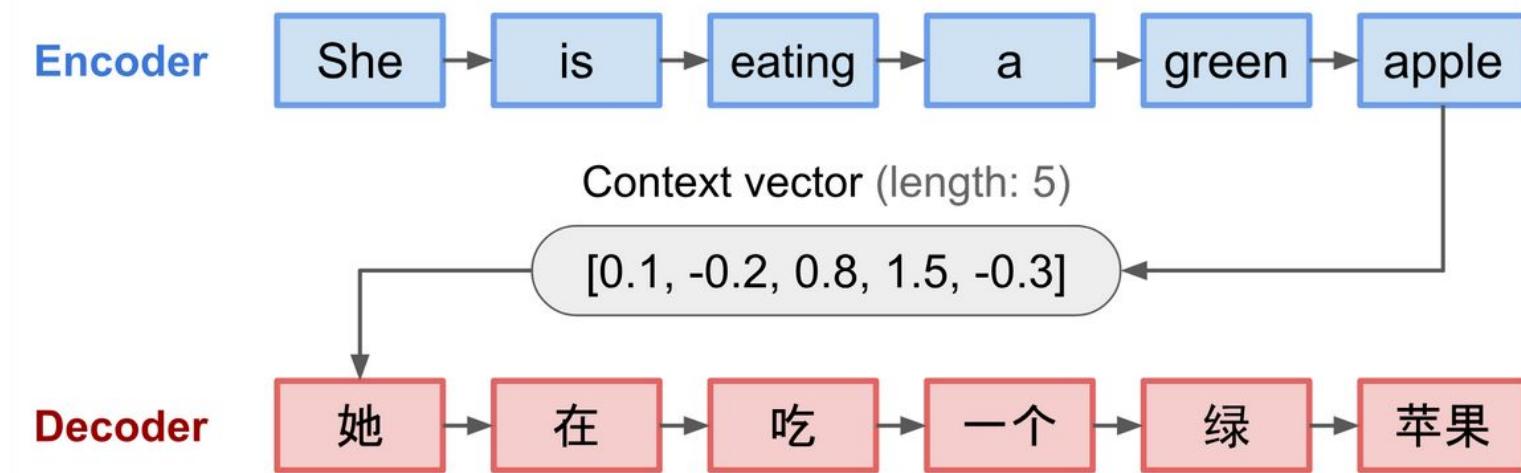
Relating **distant elements** of the input may be tricky:

- #Operations grows with sequence length
- Truncated Backpropagation
- Vanishing gradients if full BPTT



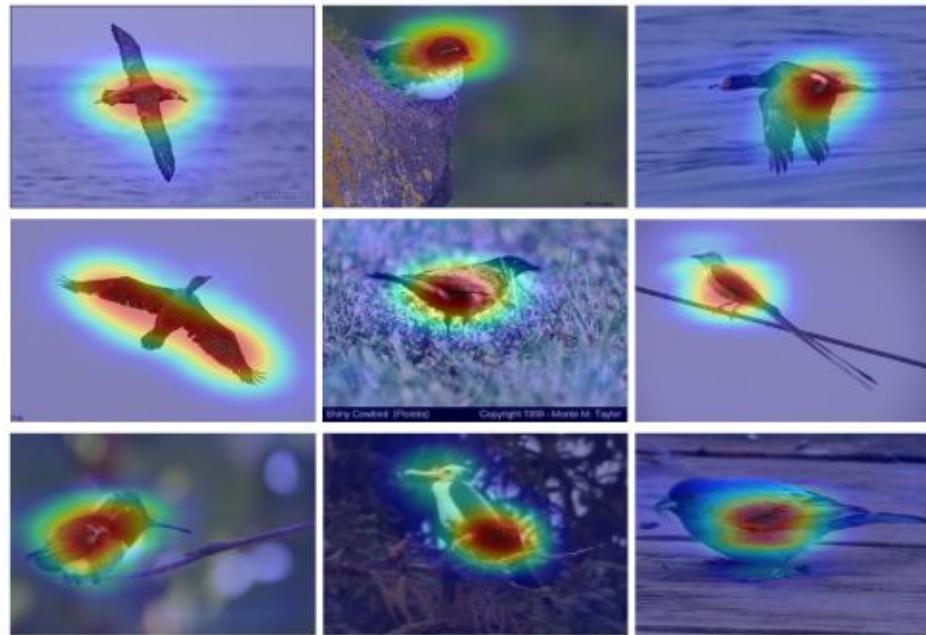
Sequence to sequence (seq2seq)

- Input a sequence, output another sequence (e.g. translation)
- Fixed representation length



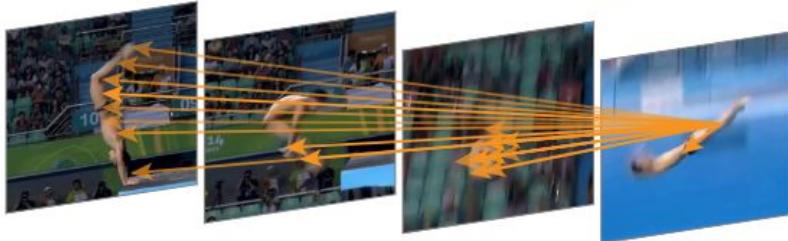
Attention Mechanisms

- Select important regions of an input.



Attention Mechanisms

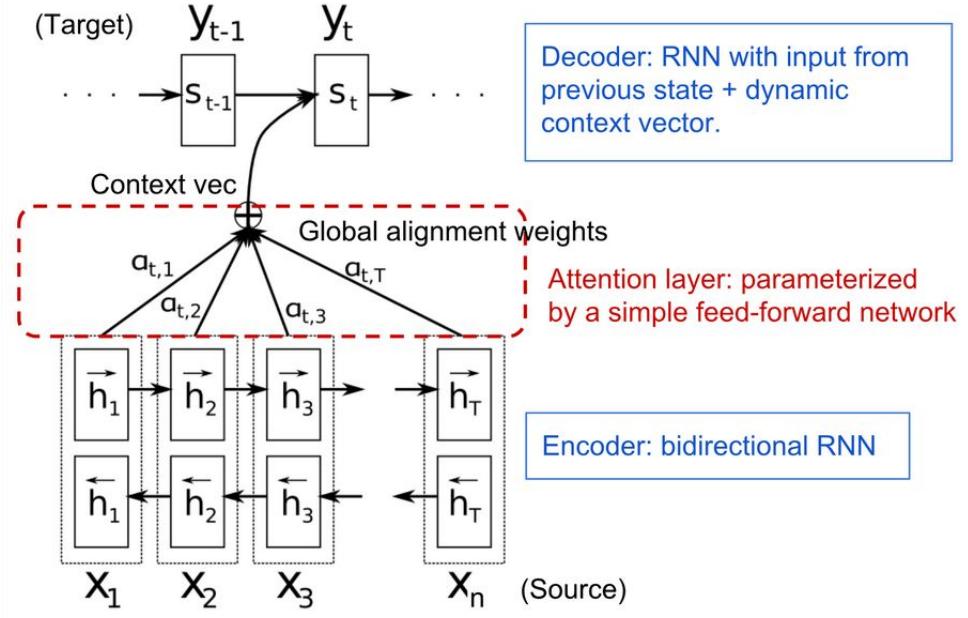
- Select important elements of a sequence.



$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i$$

$$\alpha_{t,i} = \text{align}(y_t, x_i)$$

$$= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))}$$



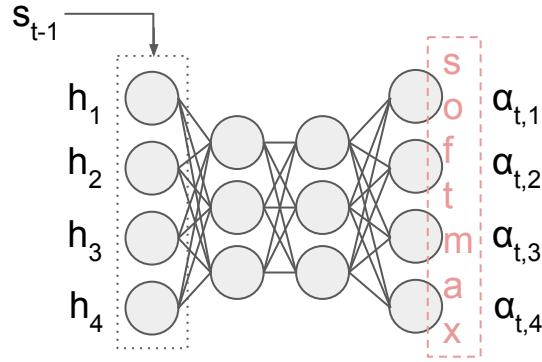
; Context vector for output y_t

; How well two words y_t and x_i are aligned.

; Softmax of some predefined alignment score..

Attention Mechanisms

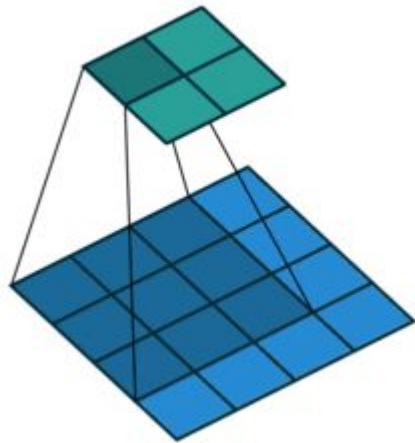
Fully connected layer → **Uses learned weights.**



$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i$$

$$\begin{aligned}\alpha_{t,i} &= \text{align}(y_t, x_i) \\ &= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))}\end{aligned}$$

Local vs. Non-local



Local (Convolution)



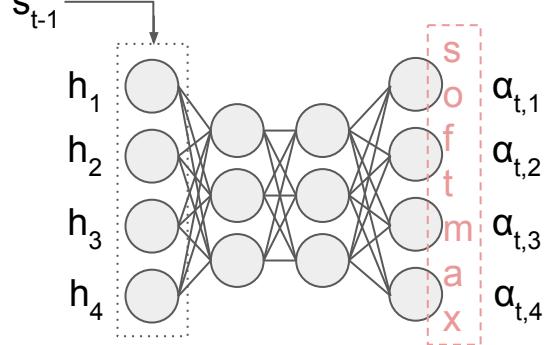
Non-local

Source: https://github.com/vdumoulin/conv_arithmetic

[2] Wang, X. et al. "Non-local neural networks." In CVPR (2018).

Attention Mechanisms

Fully connected layer → **Uses learned weights.**

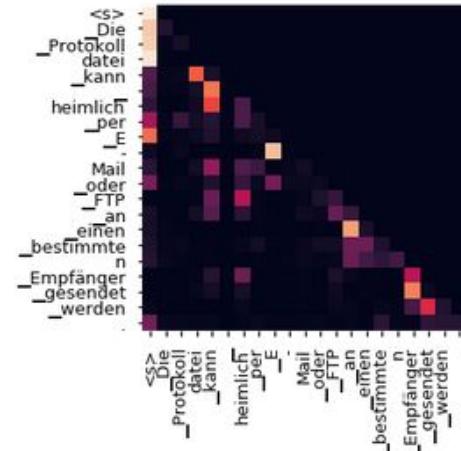


$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i$$

$$\begin{aligned}\alpha_{t,i} &= \text{align}(y_t, x_i) \\ &= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))}\end{aligned}$$

Non-local layer → Attention based on relationships between input elements

$$\mathbf{y}_i = \frac{1}{C(\mathbf{x})} \sum_{\forall j} f(\mathbf{x}_i, \mathbf{x}_j) g(\mathbf{x}_j).$$



$f \rightarrow$ Computes the affinity

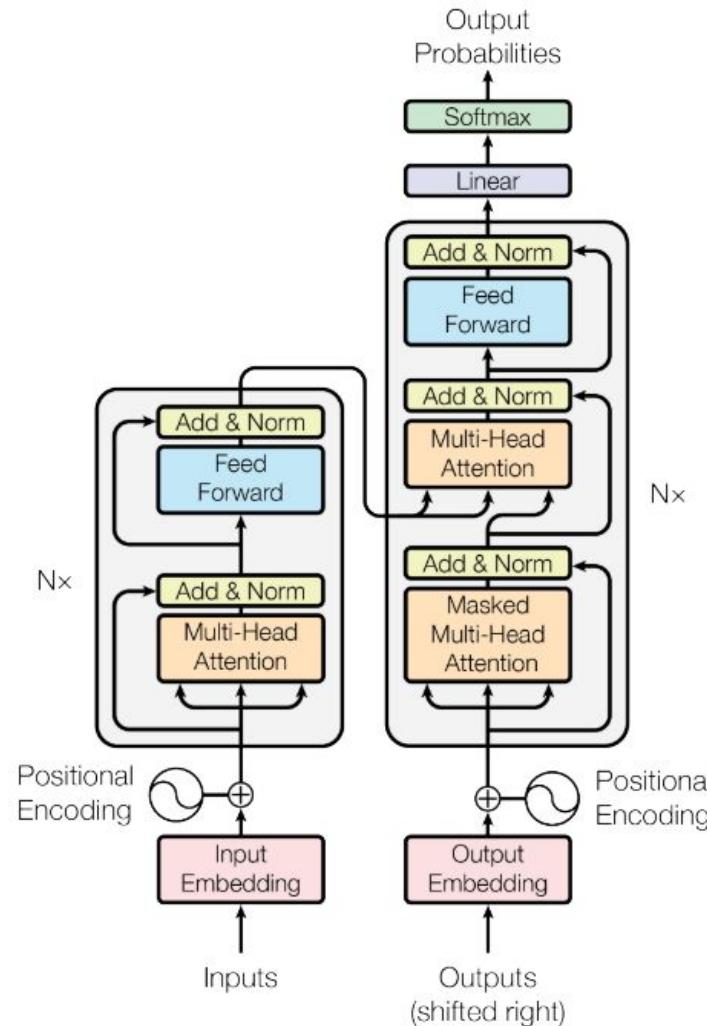
$g \rightarrow$ Computes representation of the input

$C(x) \rightarrow$ Normalization factor

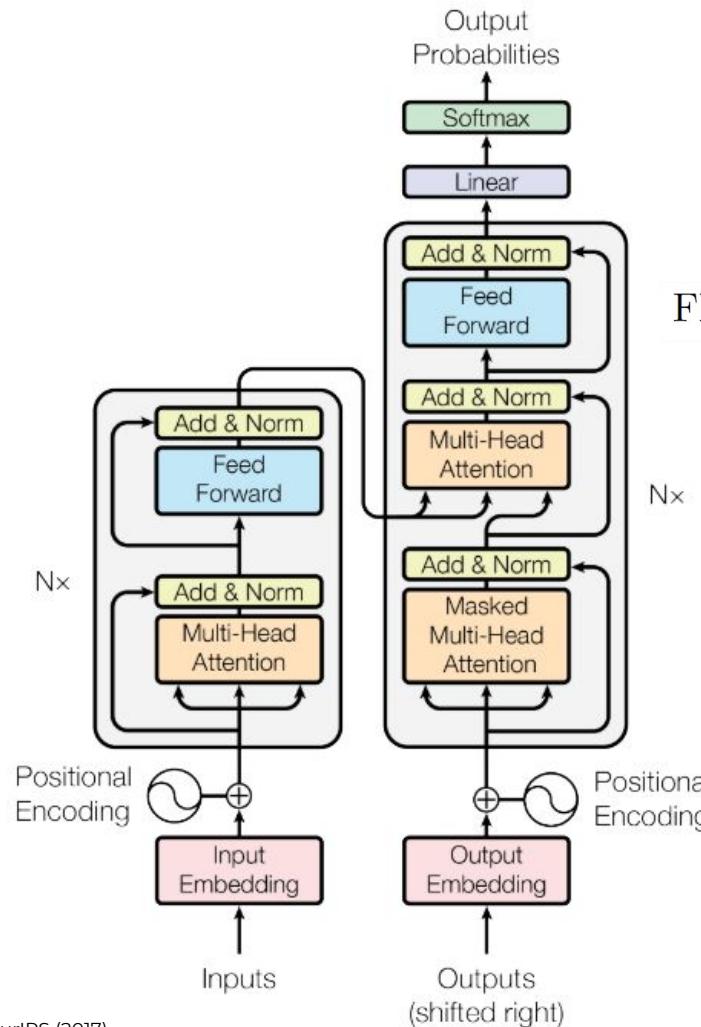
The Transformer

The transformer

- Input whole sequences (tokenized)
- Add positional information
- Self-attention
- Linear transformation
- Residuals and Normalization

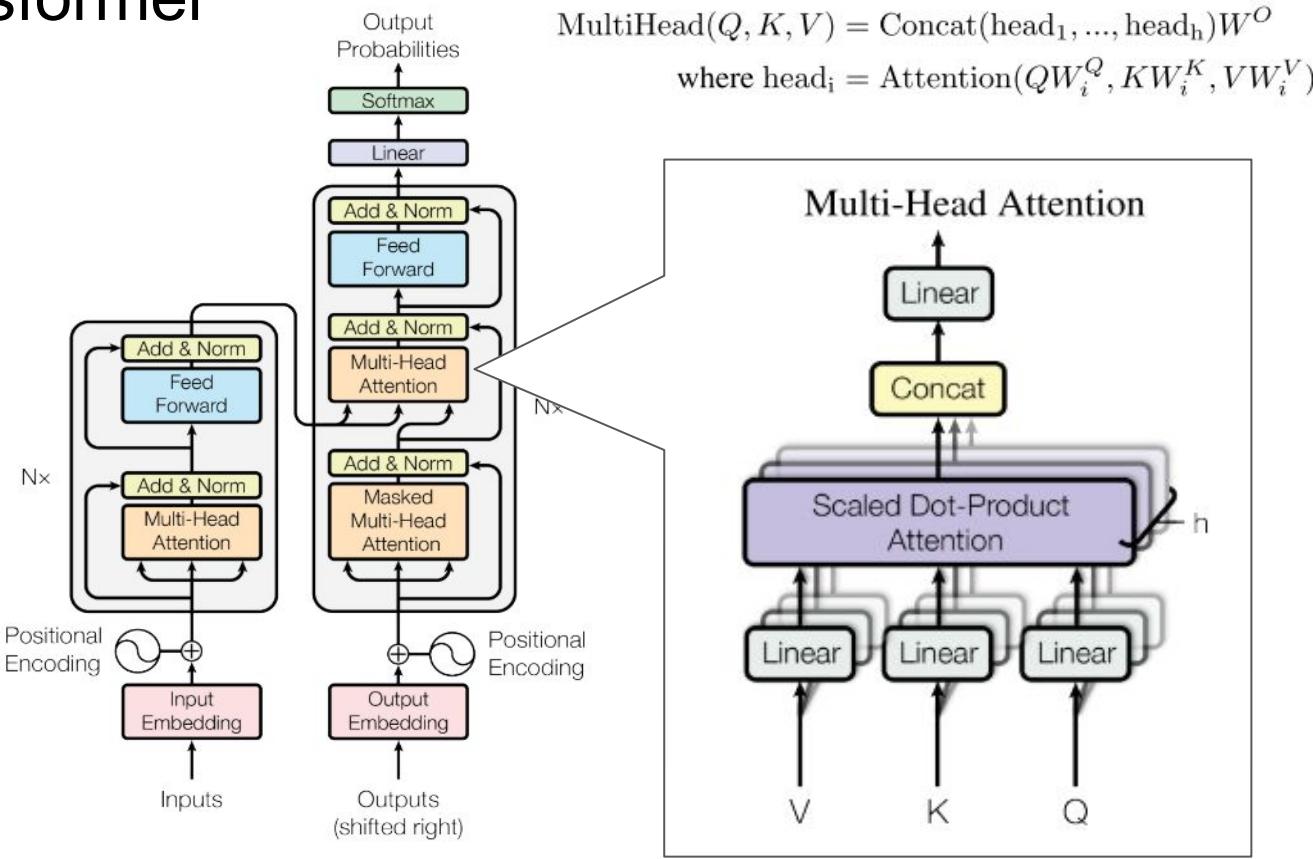


The transformer



$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

The transformer

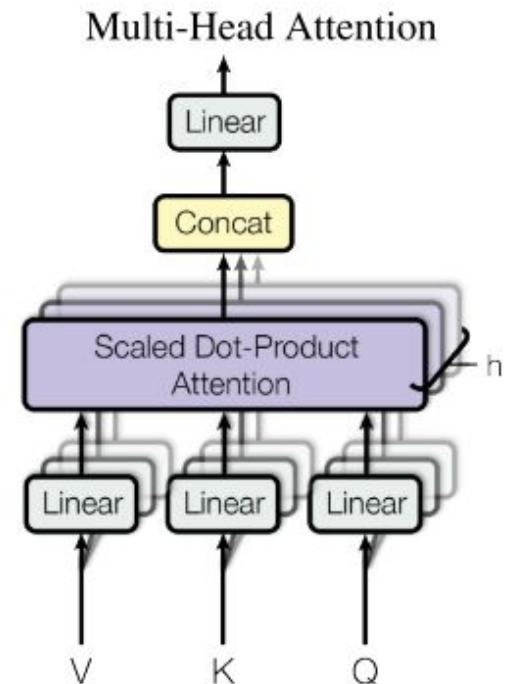


Multi-head Self-attention

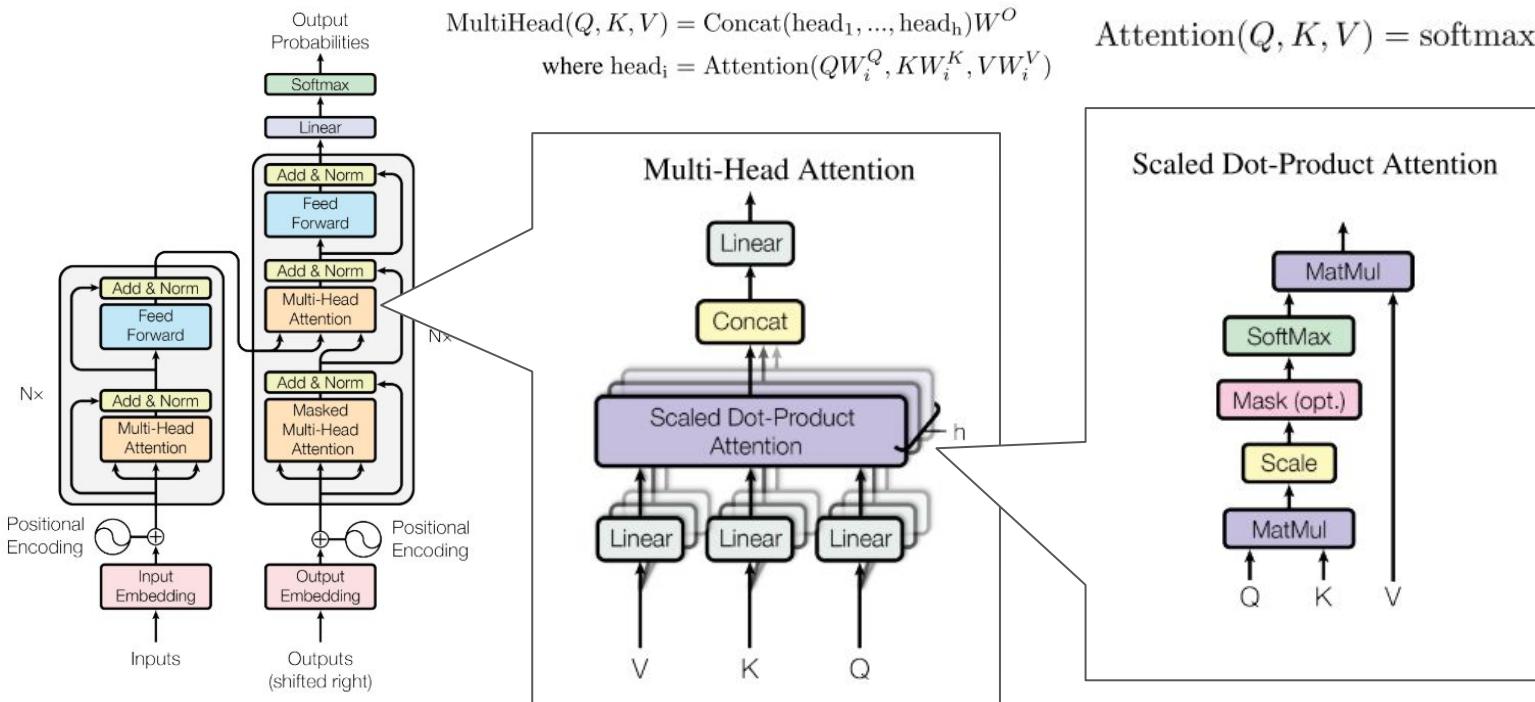
- Map into different representation spaces (similar to convolutional filters)
- Attend at different features of the same regions
- Concatenate
- Map back to common representation space

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

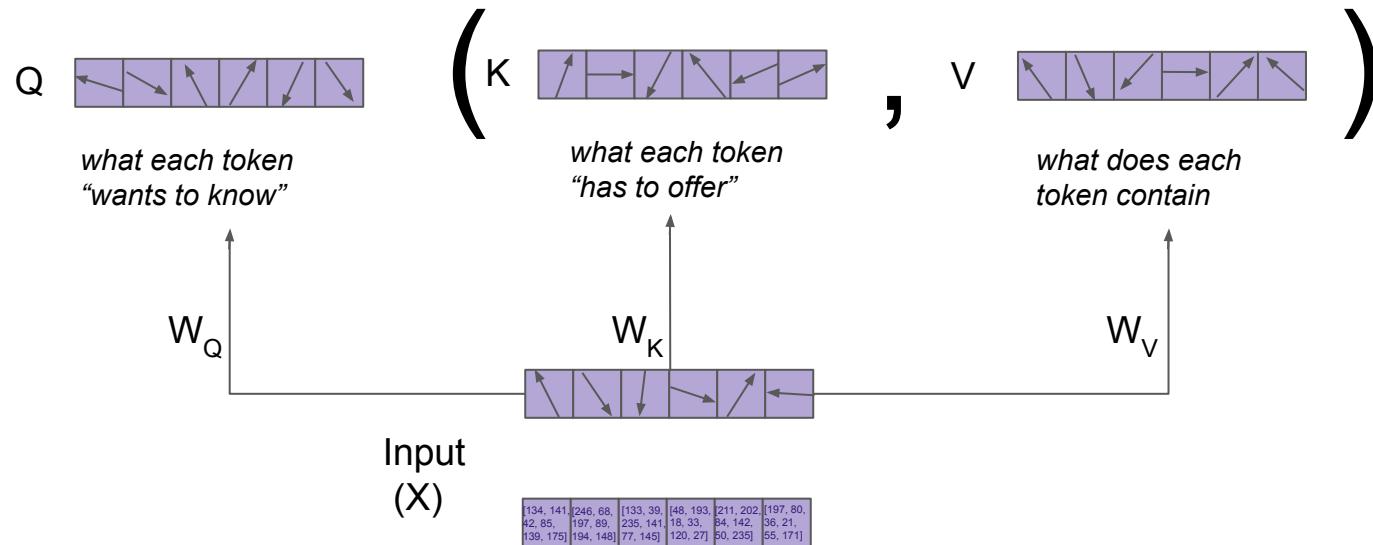


The transformer



Queries, Keys and Values

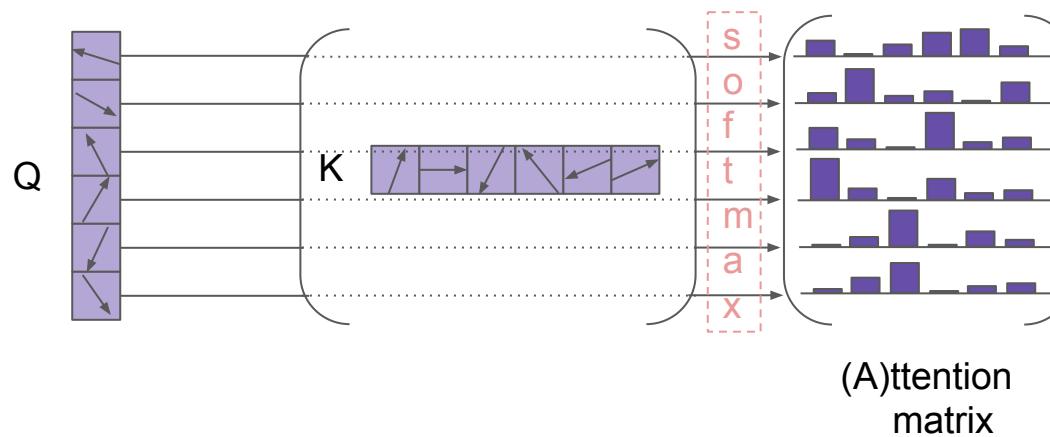
- Non-local: Mechanism used to **route information** from every token to any other
- Use tokens as **weighted contextual** information
- Set of **Queries**, and set of **pairs of Keys and Values**



Queries, Keys and Values

- Each **Query** is *compared* with all **Keys**
- Produce “*histogram*” of Query alignment with the Keys

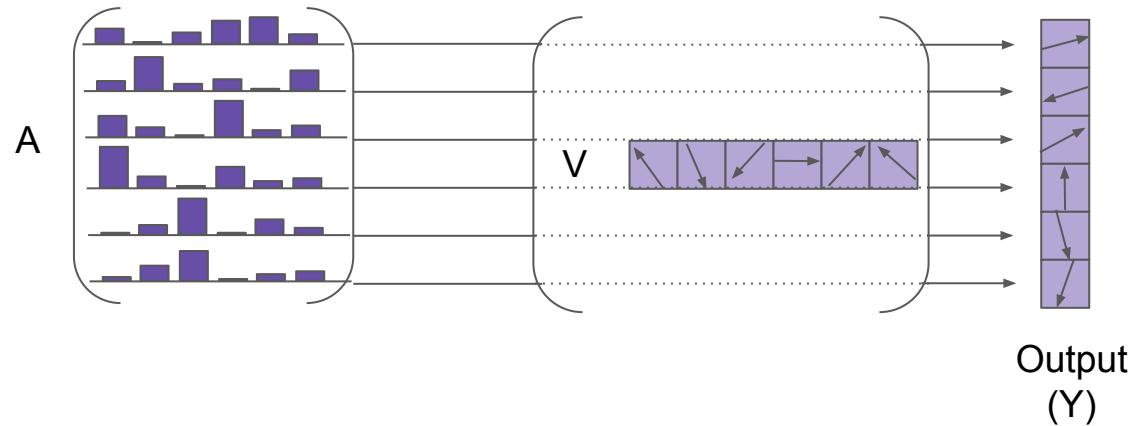
$$a_{ij} = \text{softmax}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right)$$



Queries, Keys and Values

- **Weighted summation**
- Each output token is the result of **contextual aggregation**

$$y_i = \sum_{j=0}^T v_i a_{ij}$$



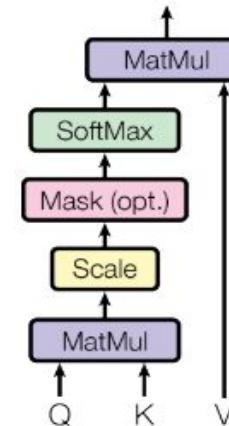
Scaled Dot-product

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention

- The dot-product of two linearly independent vectors is 0.
- Dot-product is used to compute similarity among vectors.
- Softmax to turn it into an attention map - through probabilities (sum=1)

$$a_{ij} = \text{softmax}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right) = \frac{\exp(\mathbf{q}_i \mathbf{k}_j^\top)}{\sqrt{d_k} \sum_{r \in S_i} \exp(\mathbf{q}_i \mathbf{k}_r^\top)}$$

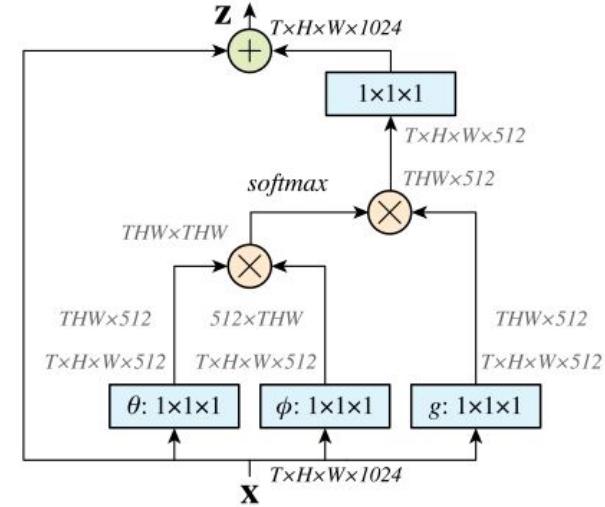


$$y_i = \sum_{j=0}^T v_i a_{ij}$$

Dot-product as non-local

- The dot-product is an instantiation of the non-local operation

$$\mathbf{y}_i = \frac{1}{\mathcal{C}(\mathbf{x})} \sum_{\forall j} f(\mathbf{x}_i, \mathbf{x}_j)g(\mathbf{x}_j).$$



$$a_{ij} = \text{softmax}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right) = \frac{\exp(\mathbf{q}_i \mathbf{k}_j^\top)}{\sqrt{d_k} \sum_{r \in S_i} \exp(\mathbf{q}_i \mathbf{k}_r^\top)}$$

$$y_i = \sum_{j=0}^T v_i a_{ij}$$

Overview of the transformer

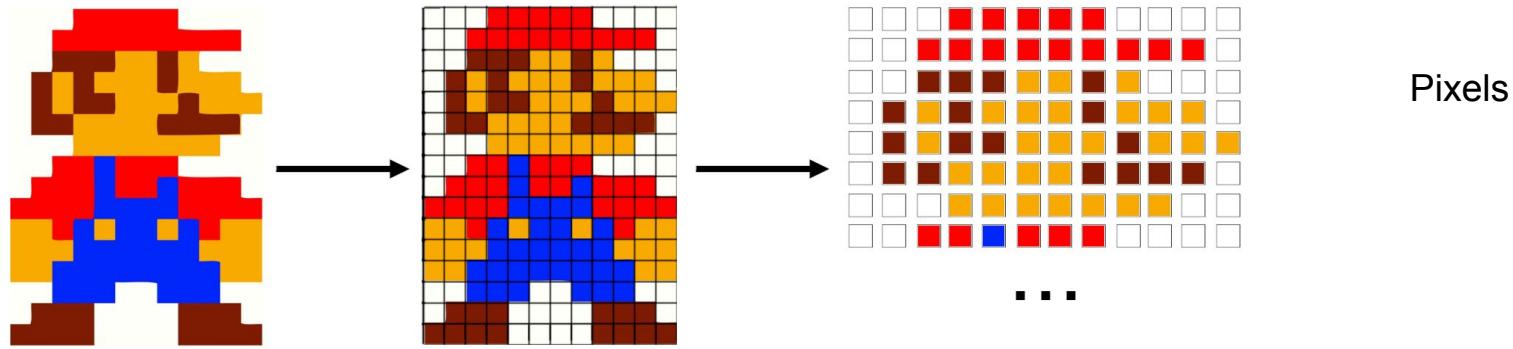
- Full **example** of original NLP Transformer for **machine translation**
- Main **difference with vision** transformers is **before the input** to the Transformer layers:
 - Tokenization
 - Embedding
 - Positional Encodings

Overview of the transformer

The | cat | wants | to | eat | pizza |

TOKENIZATION

Tokenization (images)

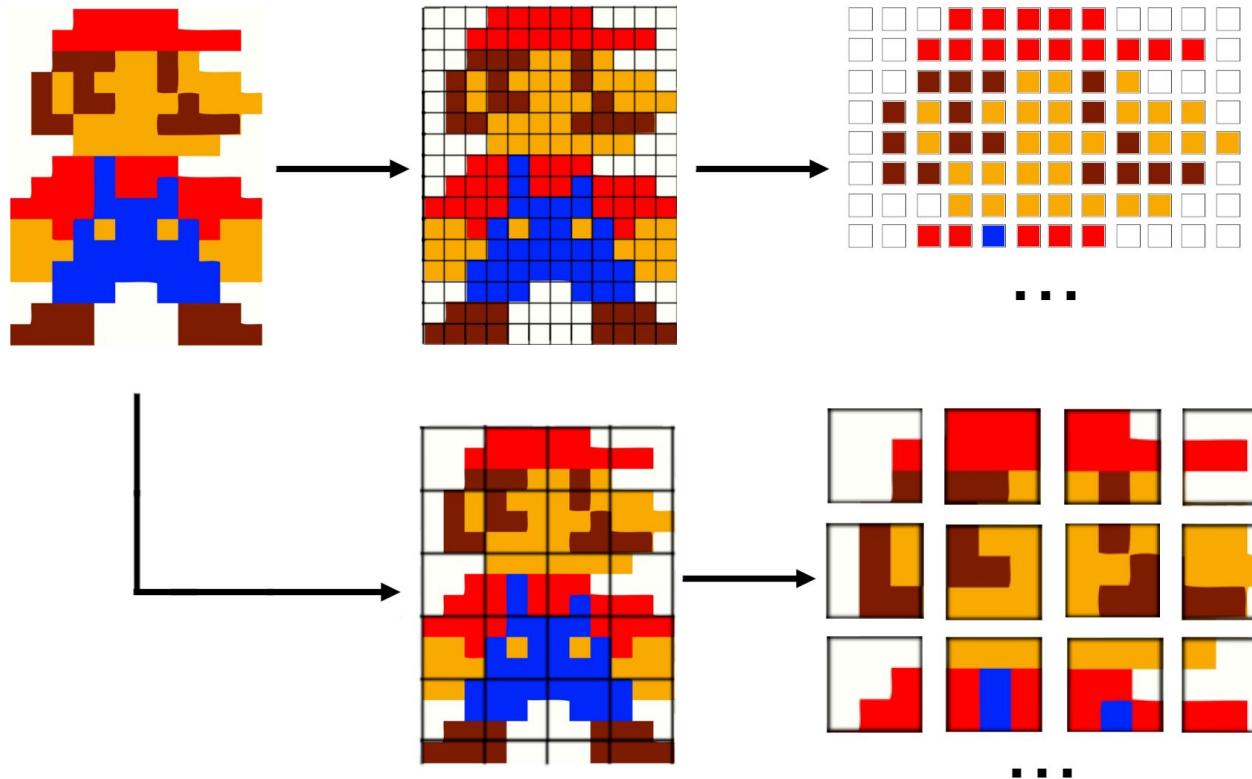


But attention matrix grows **quadratically** with **input sequence length!!** → $O(n^2)$

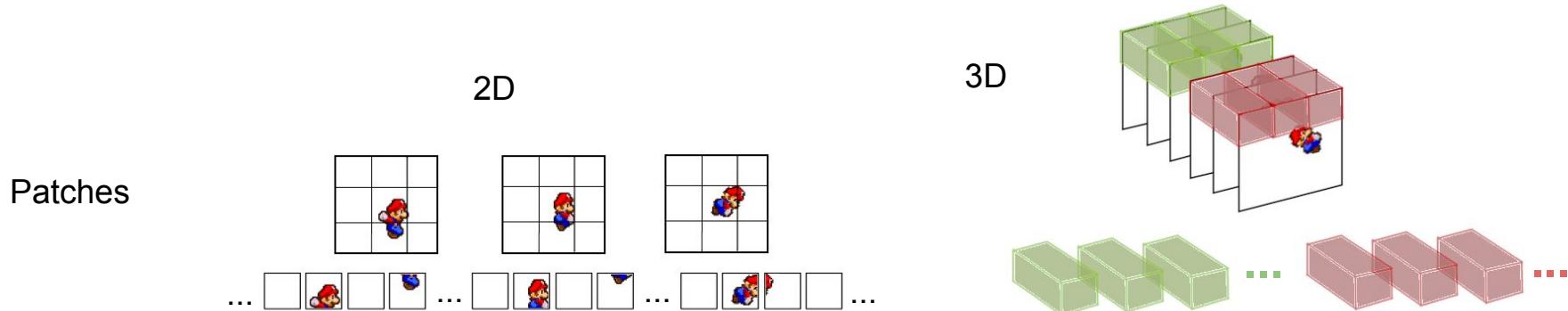
A 256x256 image would have = 65.536 px

Attention matrix would have **65.536² elements!!**

Tokenization (images)



Tokenization (videos)

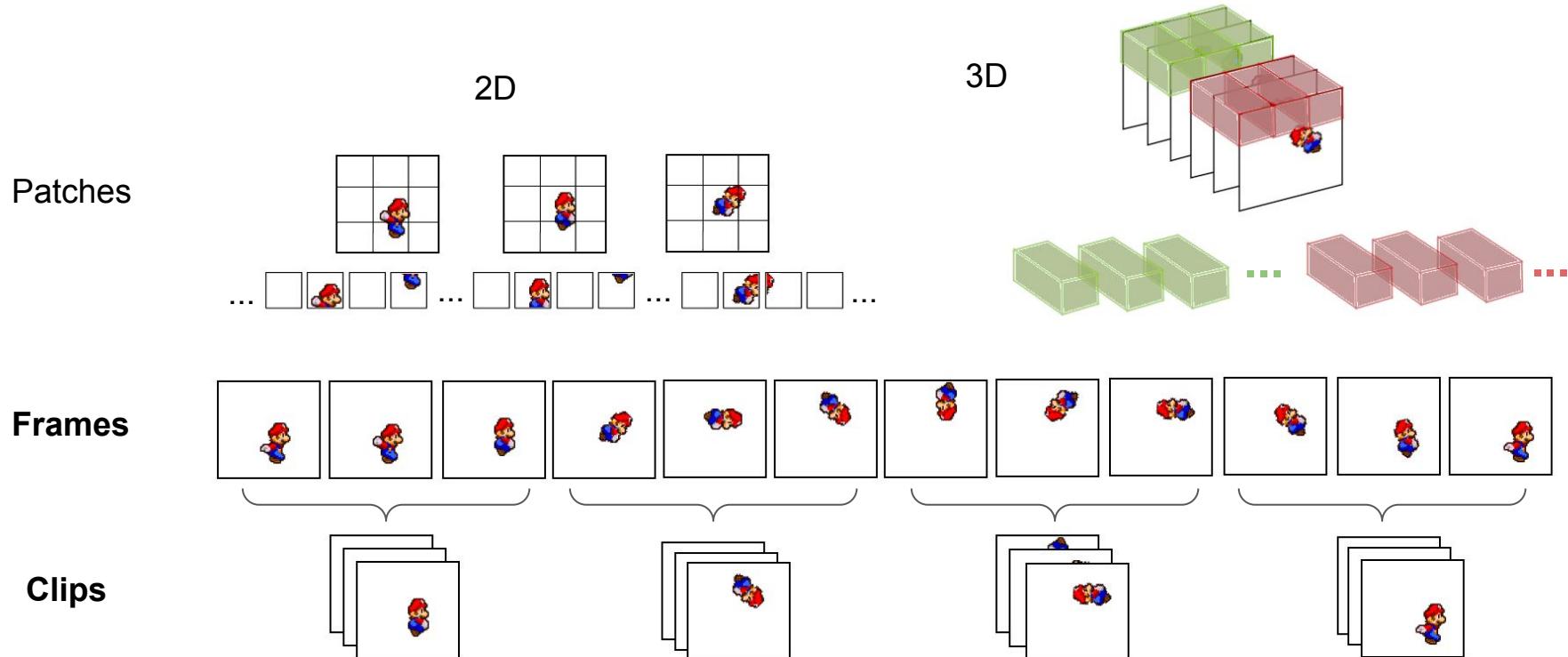


A 256x256 image tokenized as 16x16 patches, would have = 256 tokens per frame

15 second video at 25 fps → Attention matrix would have **96.000² elements!!**

3D patch setting (16x16x16) alleviates, but still does not scale that well → **A has 6000²**)

Tokenization (videos)

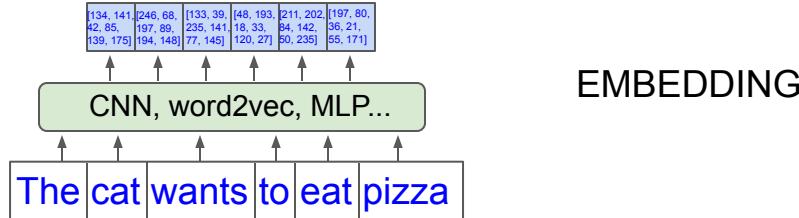


Overview of the transformer

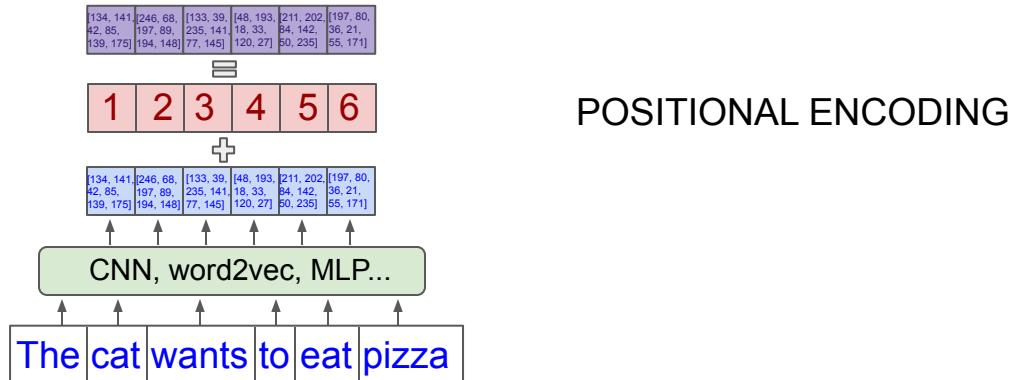
The | cat | wants | to | eat | pizza |

TOKENIZATION

Overview of the transformer



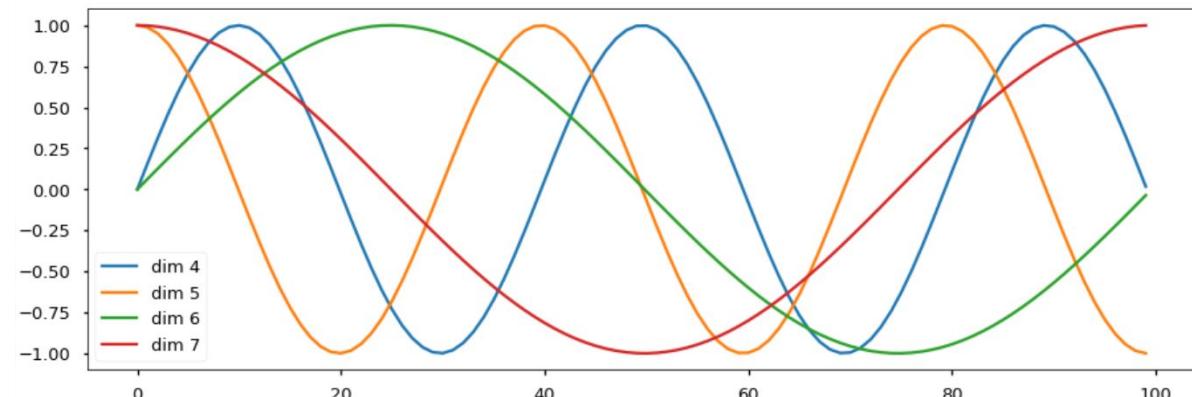
Overview of the transformer



Positional Encodings

Transformers are not invariant to translation! No idea about position!!

0 :	0 0 0 0	8 :	1 0 0 0
1 :	0 0 0 1	9 :	1 0 0 1
2 :	0 0 1 0	10 :	1 0 1 0
3 :	0 0 1 1	11 :	1 0 1 1
4 :	0 1 0 0	12 :	1 1 0 0
5 :	0 1 0 1	13 :	1 1 0 1
6 :	0 1 1 0	14 :	1 1 1 0
7 :	0 1 1 1	15 :	1 1 1 1

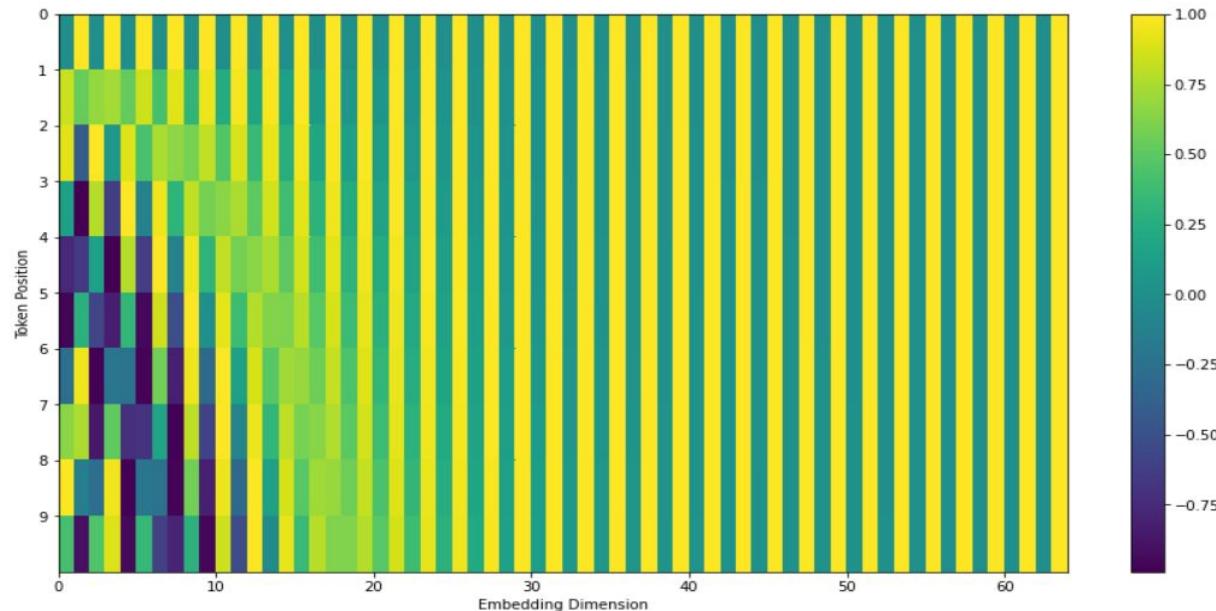


Binary

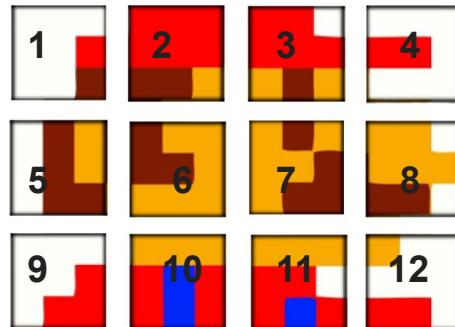
Sinusoidal

Positional Encodings

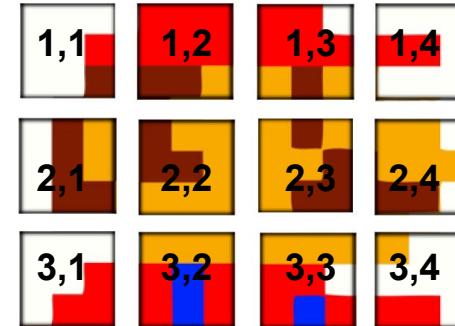
$$\text{PE}(i, \delta) = \begin{cases} \sin\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' \\ \cos\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' + 1 \end{cases}$$



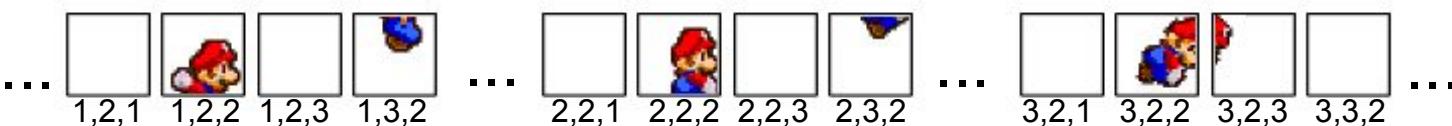
Positional Encodings



1D (Raster order)

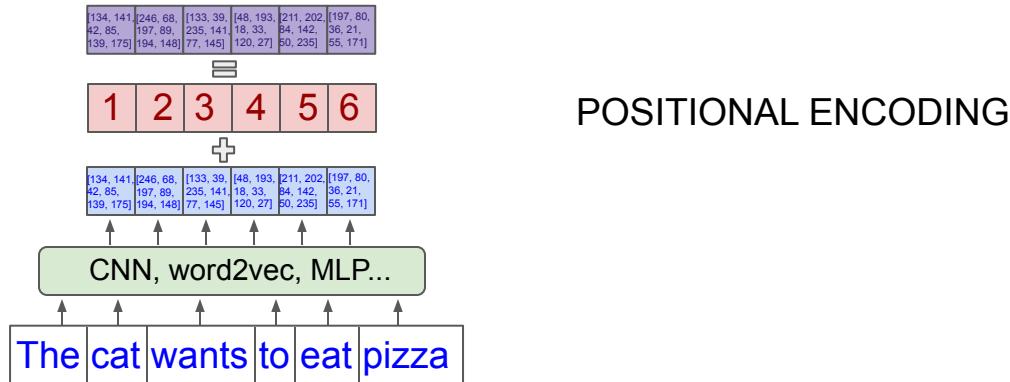


2D (Coordinates)

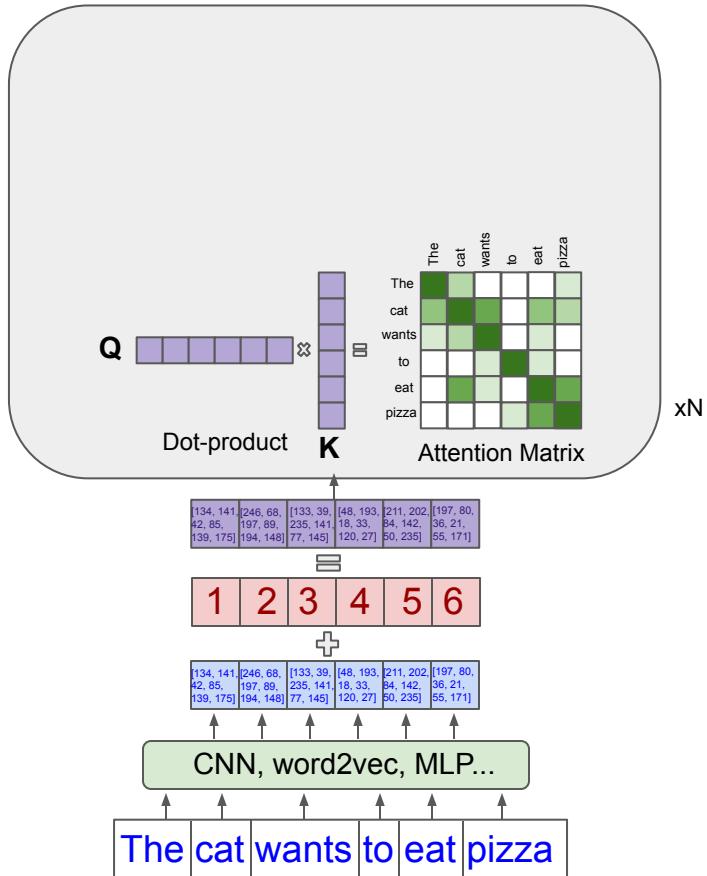


3D
(Coordinates for video patches)

Overview of the transformer



Overview of the transformer

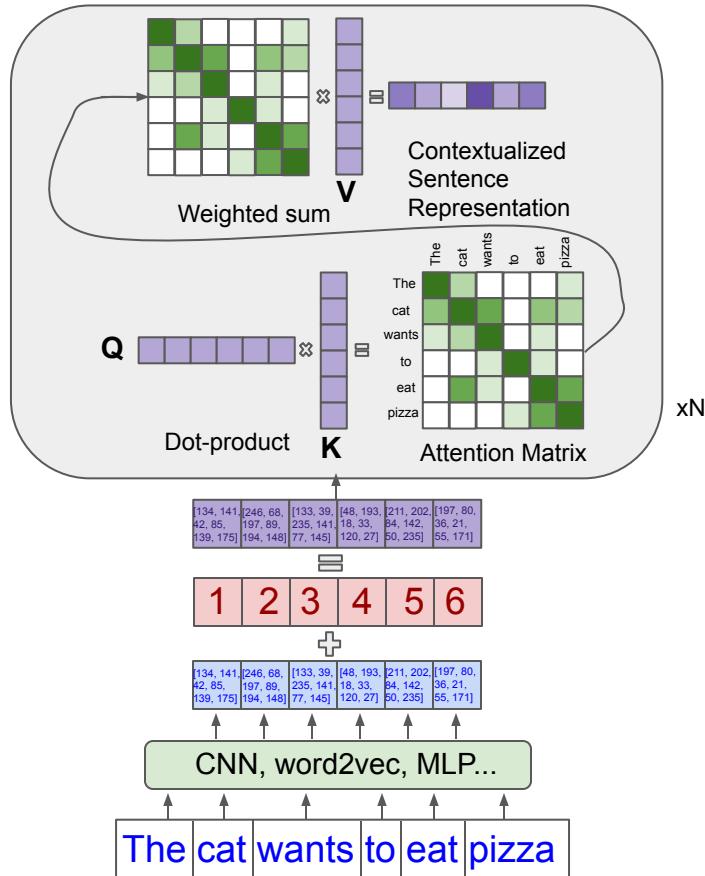


ENCODING (Self-Attention)

Computing Affinity

$$a_{ij} = \text{softmax}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_k}}\right) = \frac{\exp(\mathbf{q}_i \mathbf{k}_j^\top)}{\sqrt{d_k} \sum_{r \in S_i} \exp(\mathbf{q}_i \mathbf{k}_r^\top)}$$

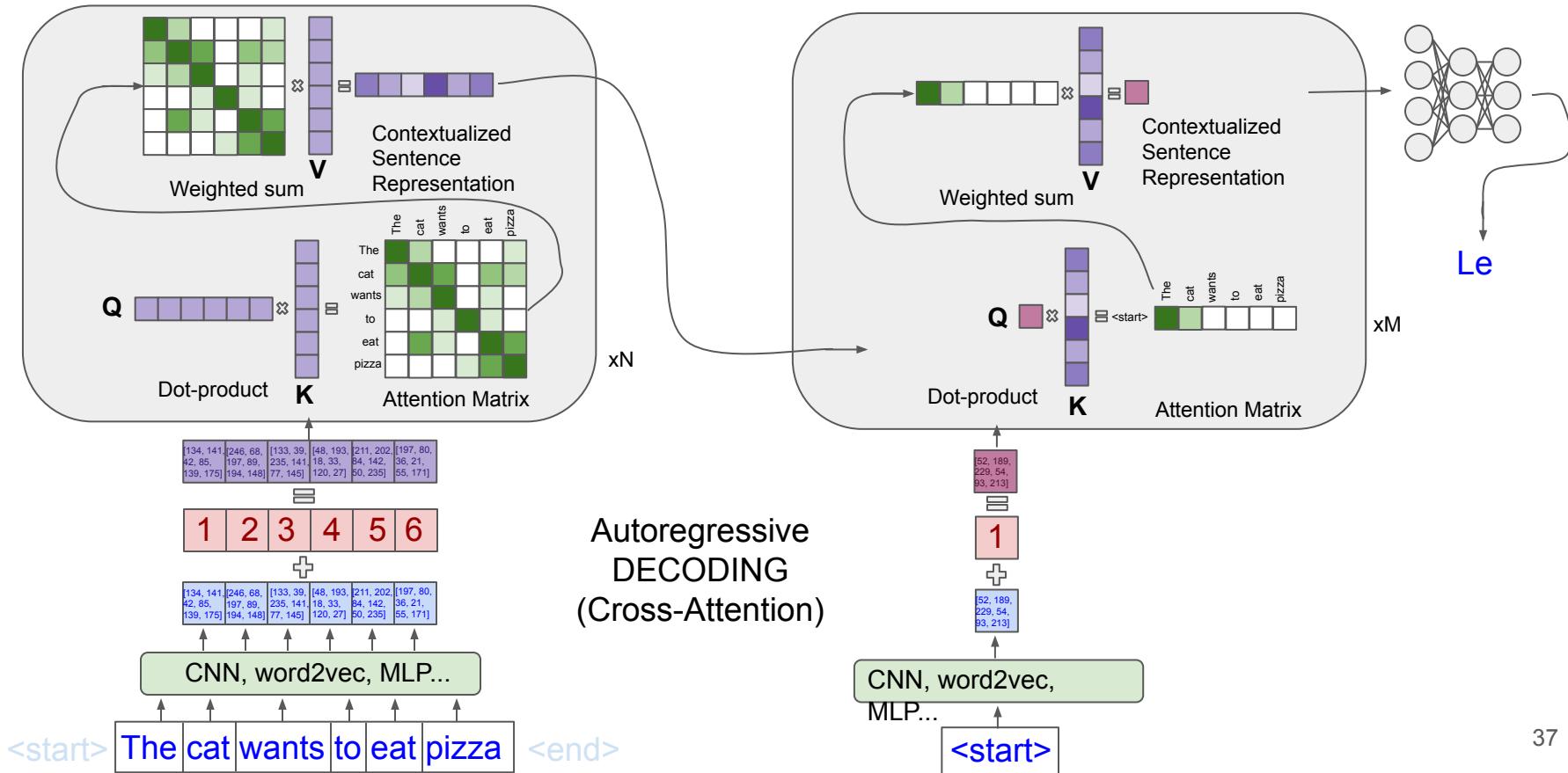
Overview of the transformer



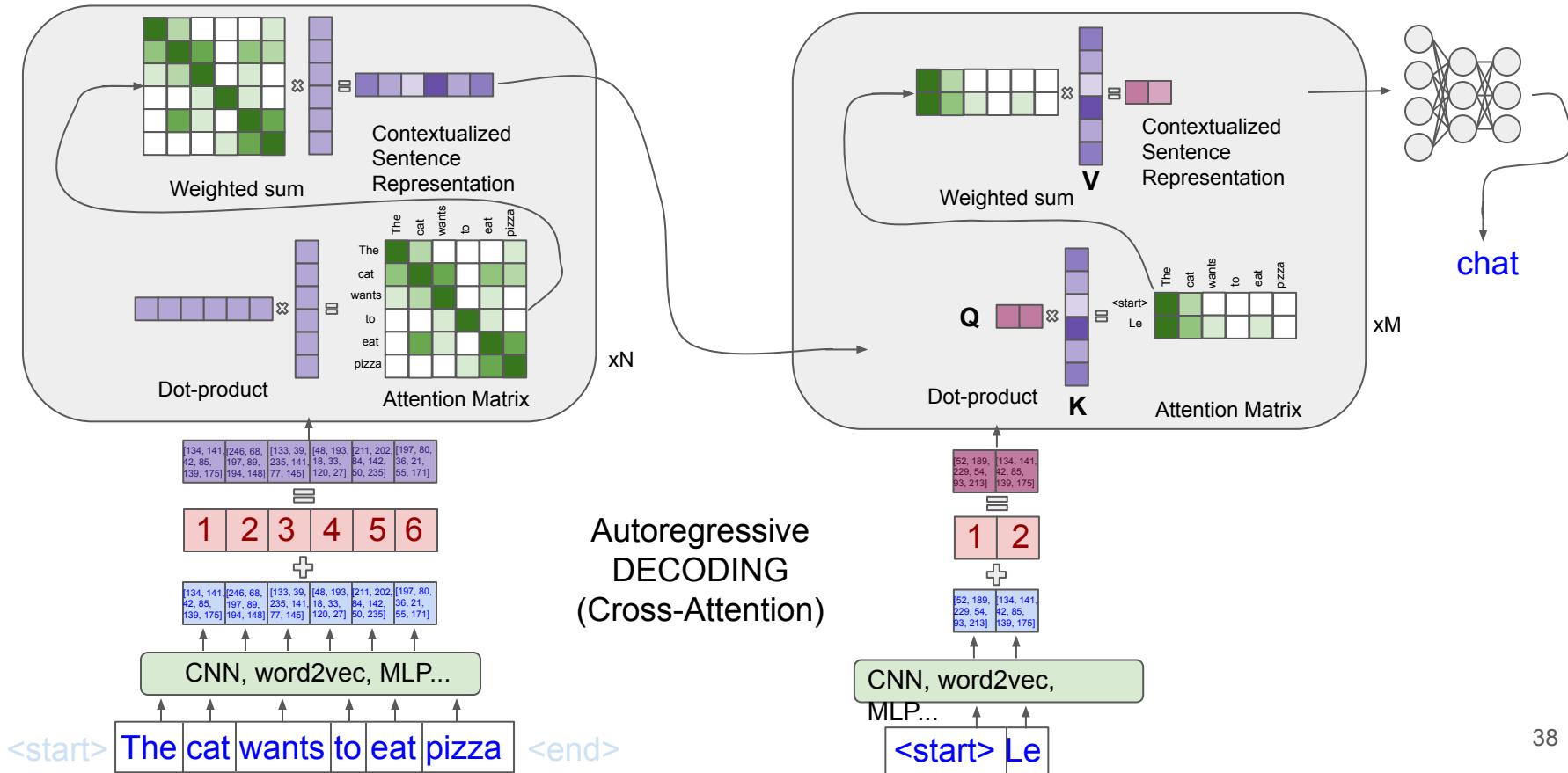
Aggregation $y_i = \sum_{j=0}^T v_i a_{ij}$

ENCODING (Self-Attention)

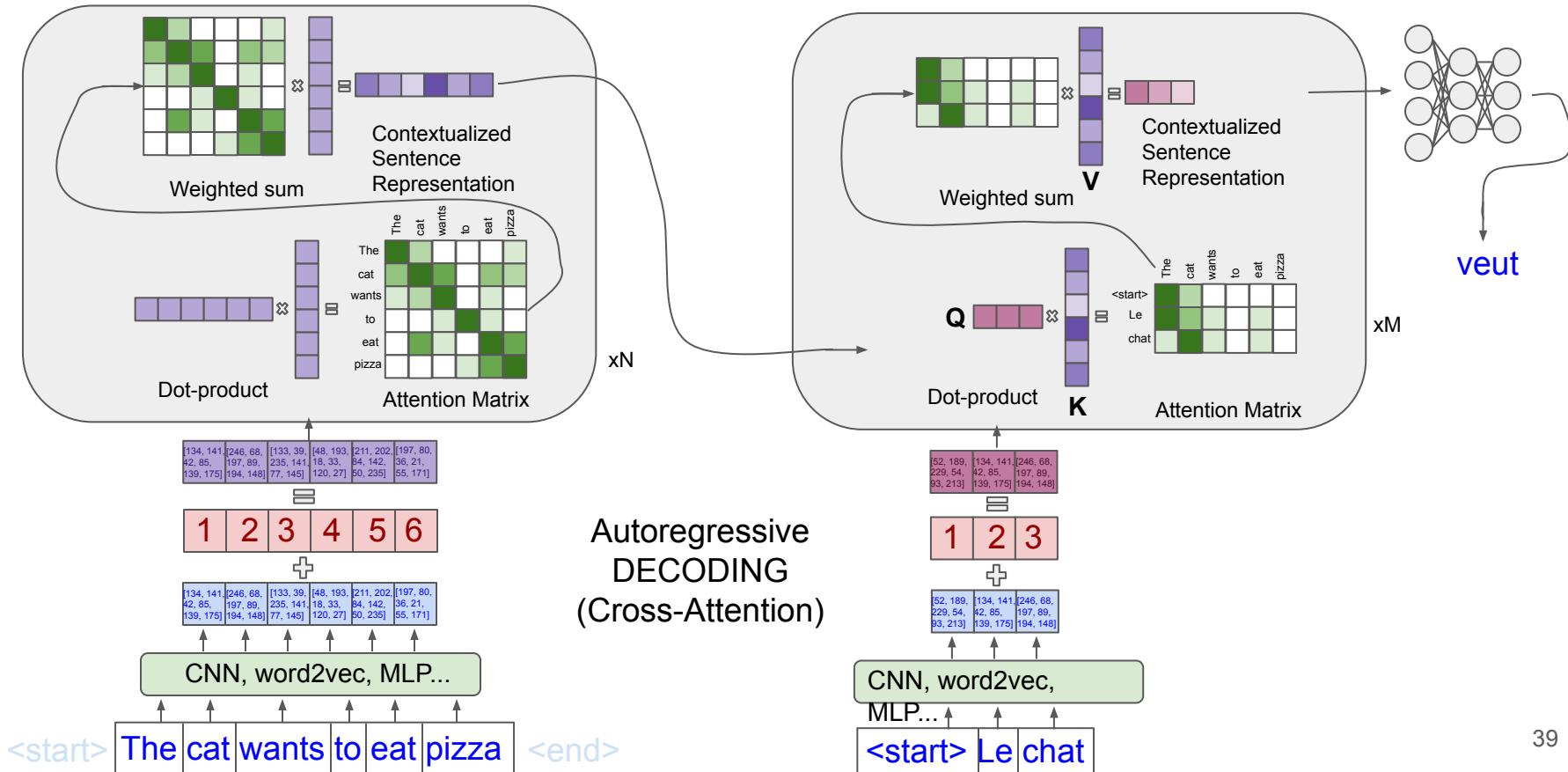
Overview of the transformer



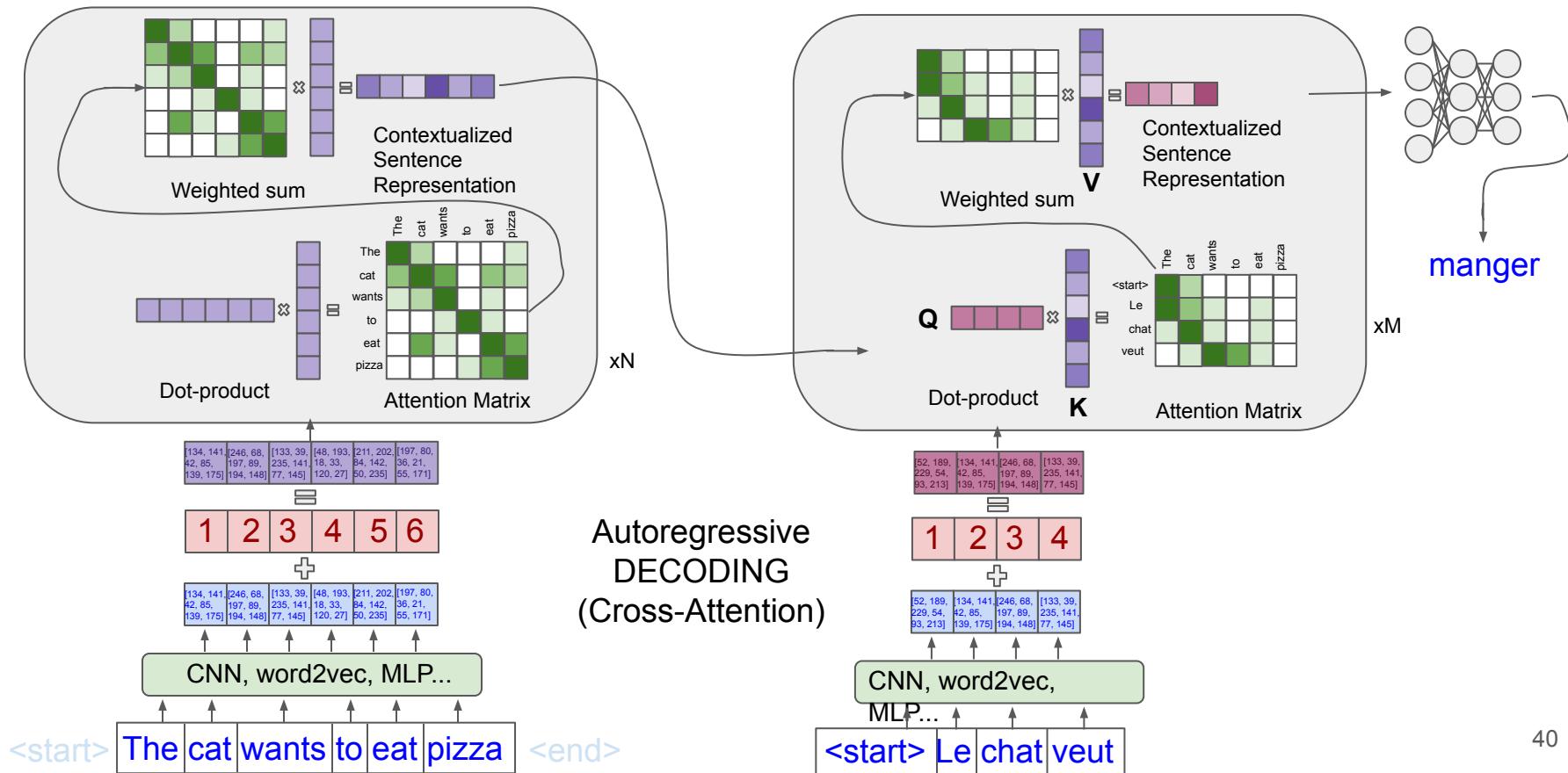
Overview of the transformer



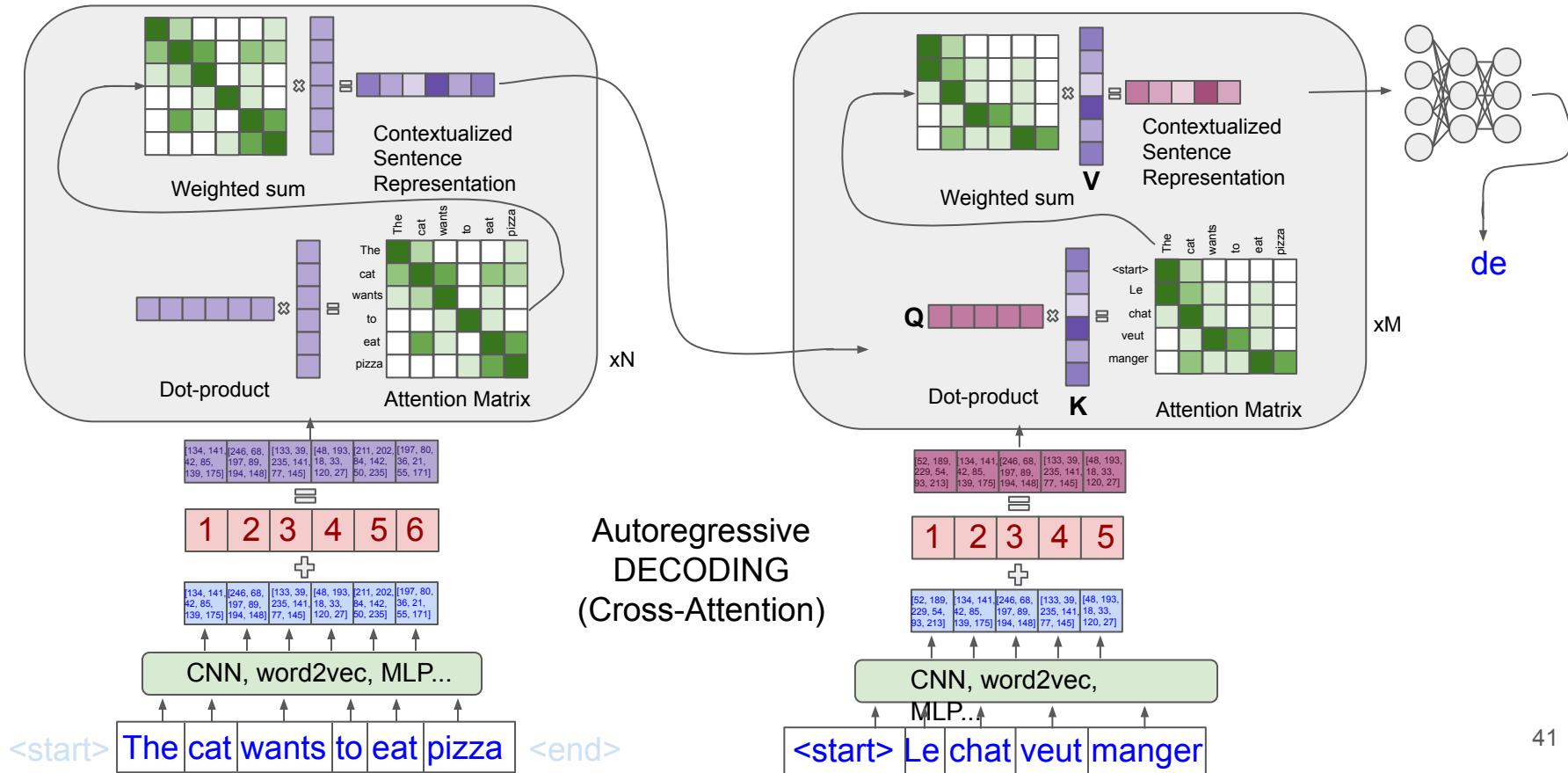
Overview of the transformer



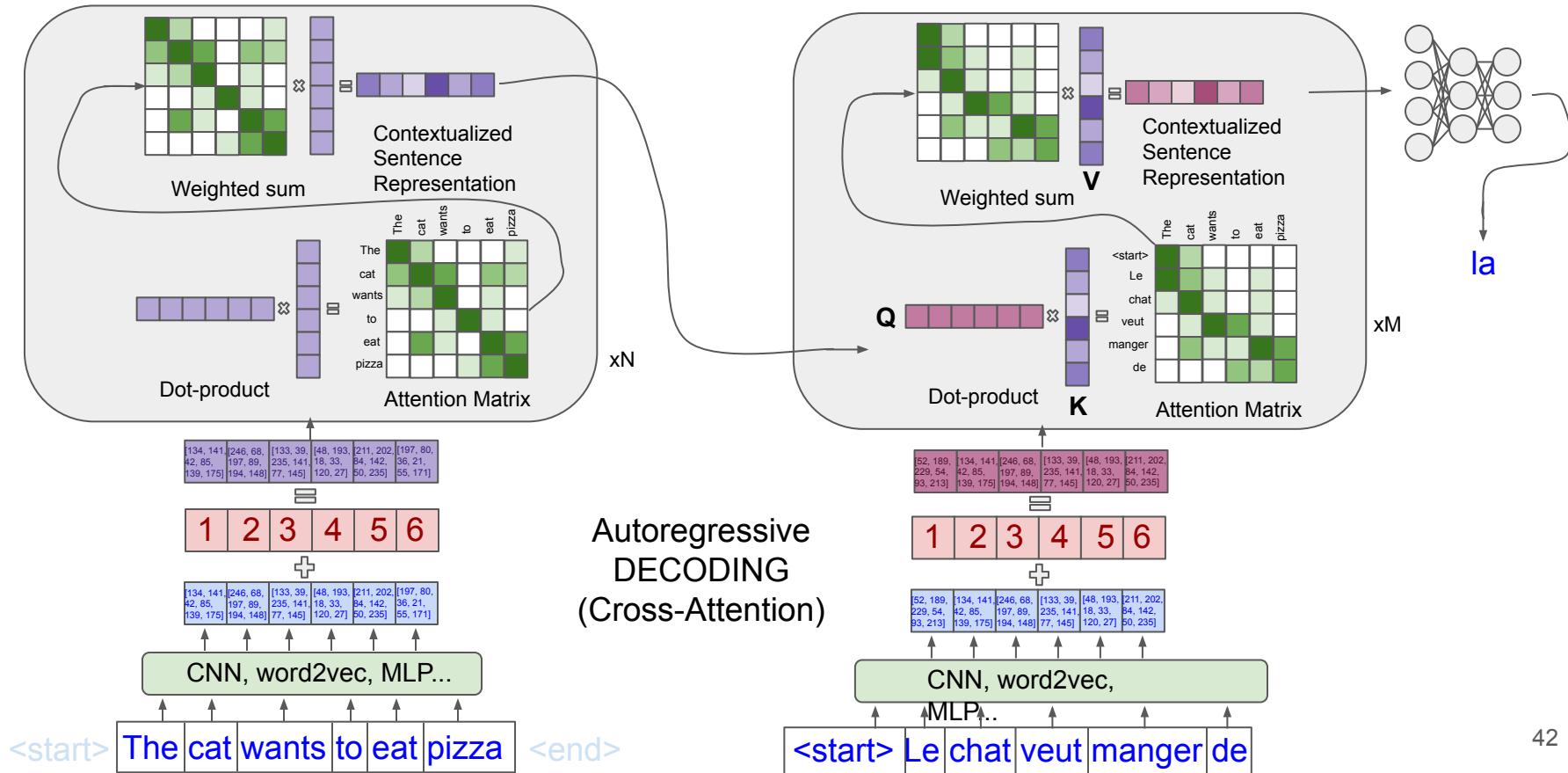
Overview of the transformer



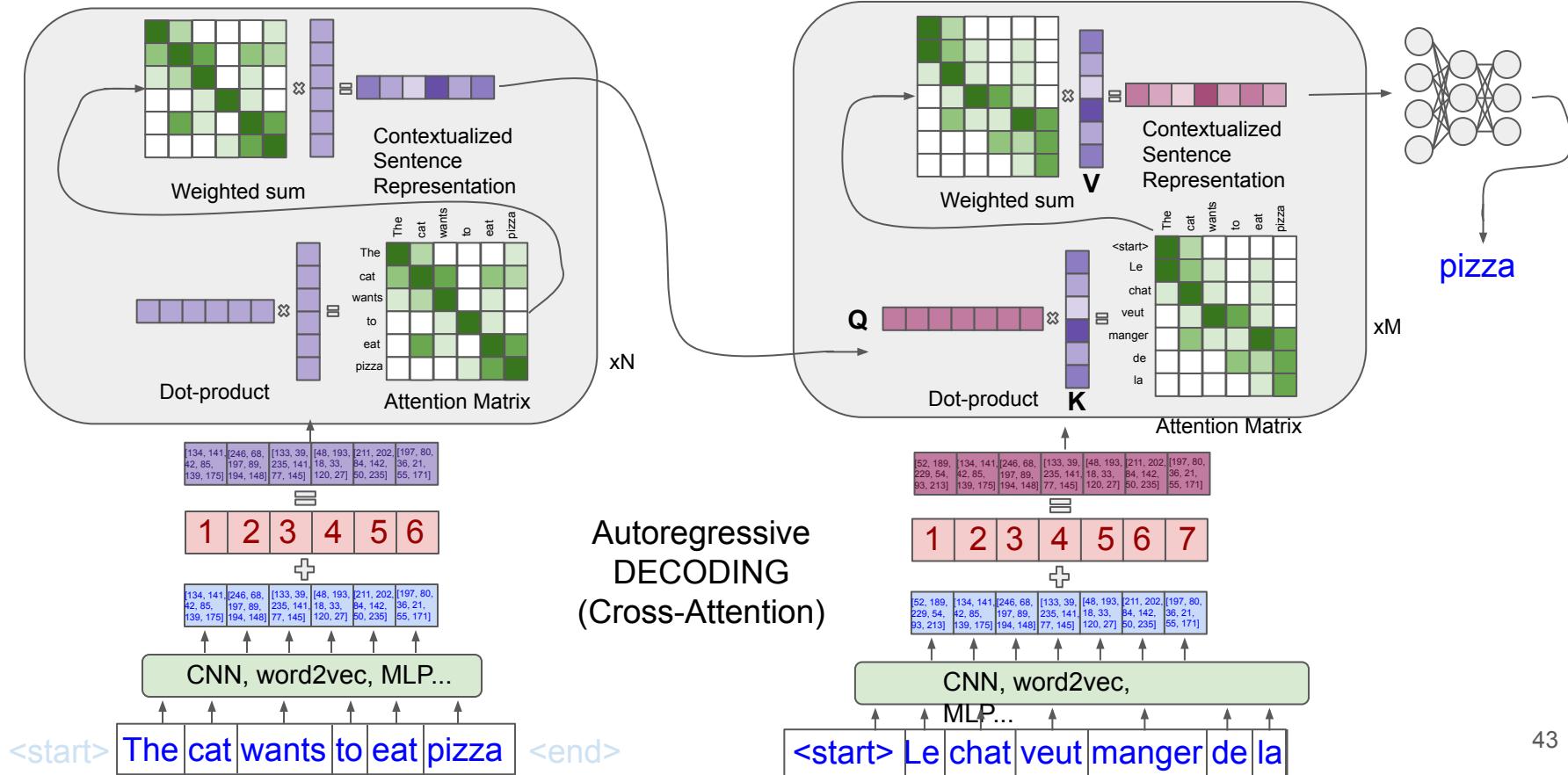
Overview of the transformer



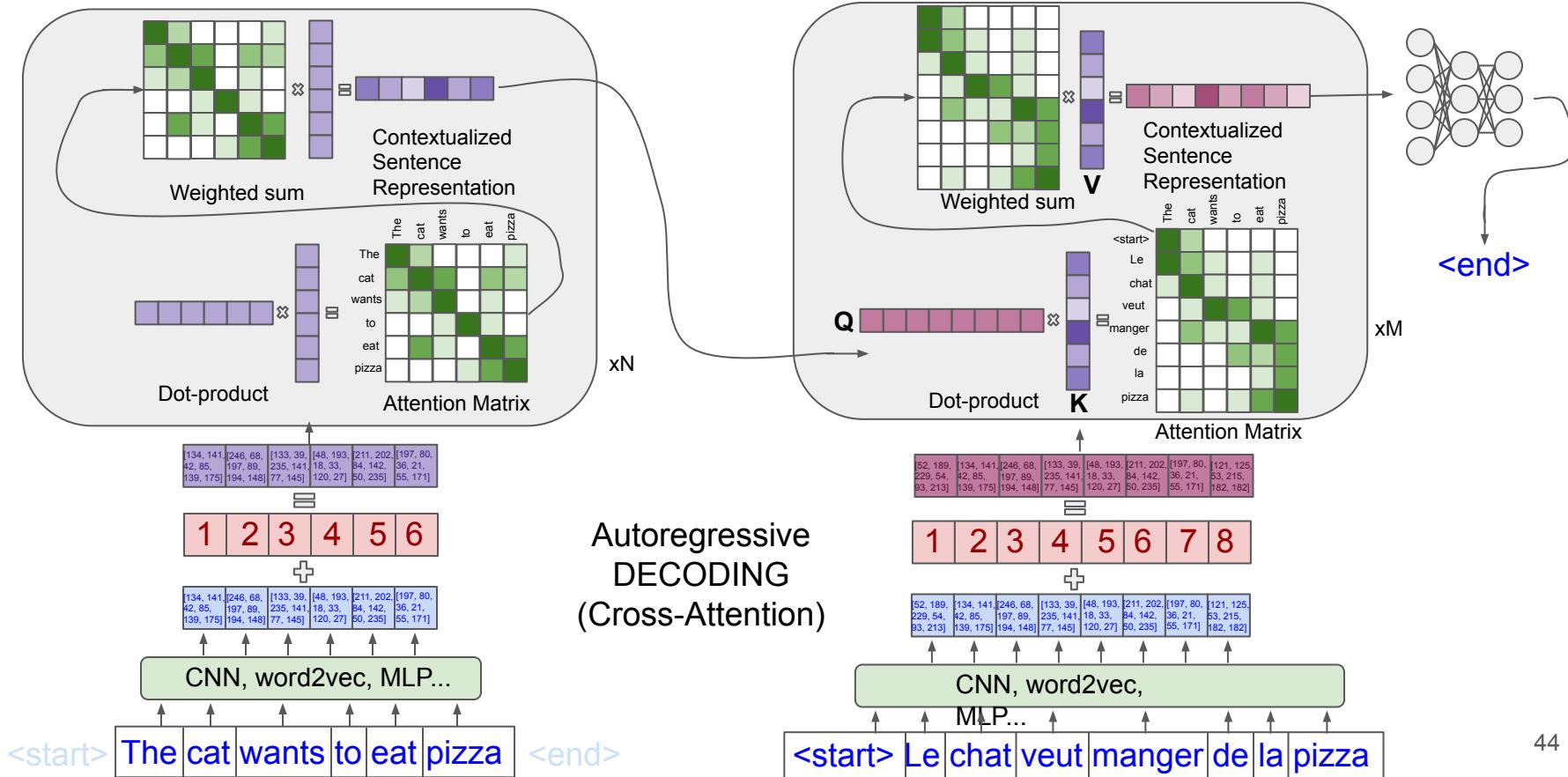
Overview of the transformer



Overview of the transformer



Overview of the transformer



AI Dungeon (GPT2)

play.aidungeon.io

You are Sven, a rogue living in the kingdom of Larion. You have a long steel dagger and a length of rope. You walk down the city street looking for somewhere to steal from. You look around and see that it's not too crowded so you decide to go into an old building that looks abandoned.

As you enter the door, you hear a voice behind you

> turn around

You turn around and see a man with a beard and wearing a black robe standing there. He has a very stern face but his eyes tell you he is friendly.

"Hello," he says. "I'm Doctor Zook".

> stab the doctor

You quickly pull out your knife and stab the doctor in the back. The doctor falls backwards and then you grab the rope hanging from the ceiling and swing yourself up. You land on top of the doctor who is still struggling to get free.

> stab the doctor again

You stab the doctor in the chest three more times. The doctor screams as blood pours out of him. You stand over him and look at the bloody body.

"Well done," Dr. Zook says.

> Say "Wait! I killed you!"

You shout "I didn't kill you!" and point your finger at the doctor.

>

GPT2

Human Prompt

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

Machine Completion

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

...

(Radford et al. 2019)

DALL-E

an image of a blue eggplant drawn on the sidewalk



a turtle claymation sitting in a forest

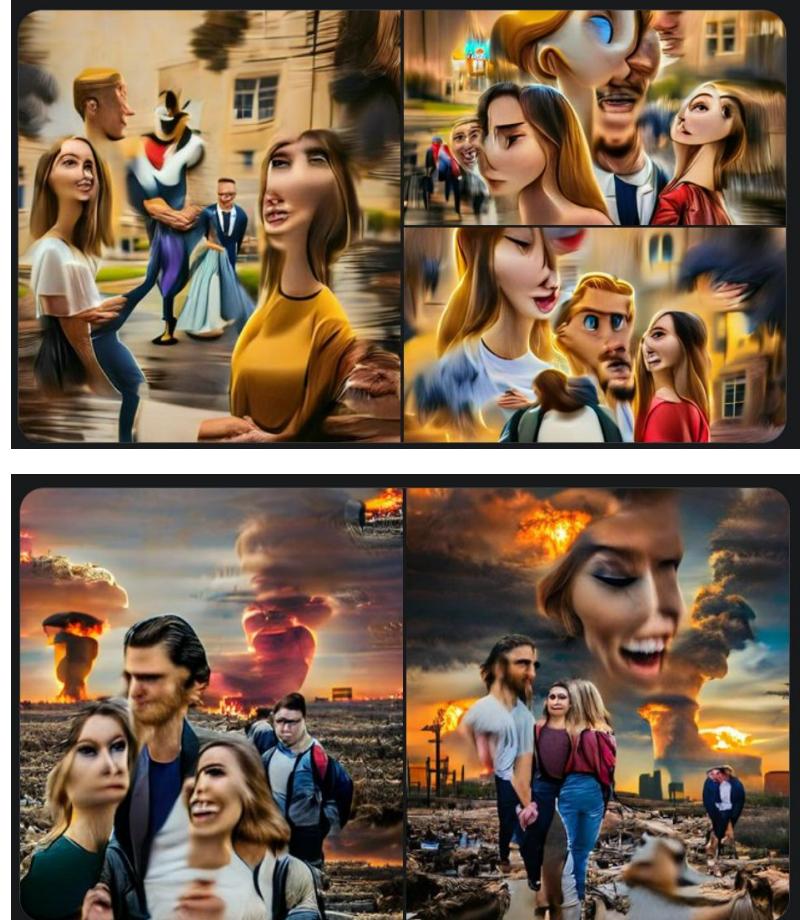


a female mannequin dressed in a brown leather jacket and purple sweatpants



CLIP

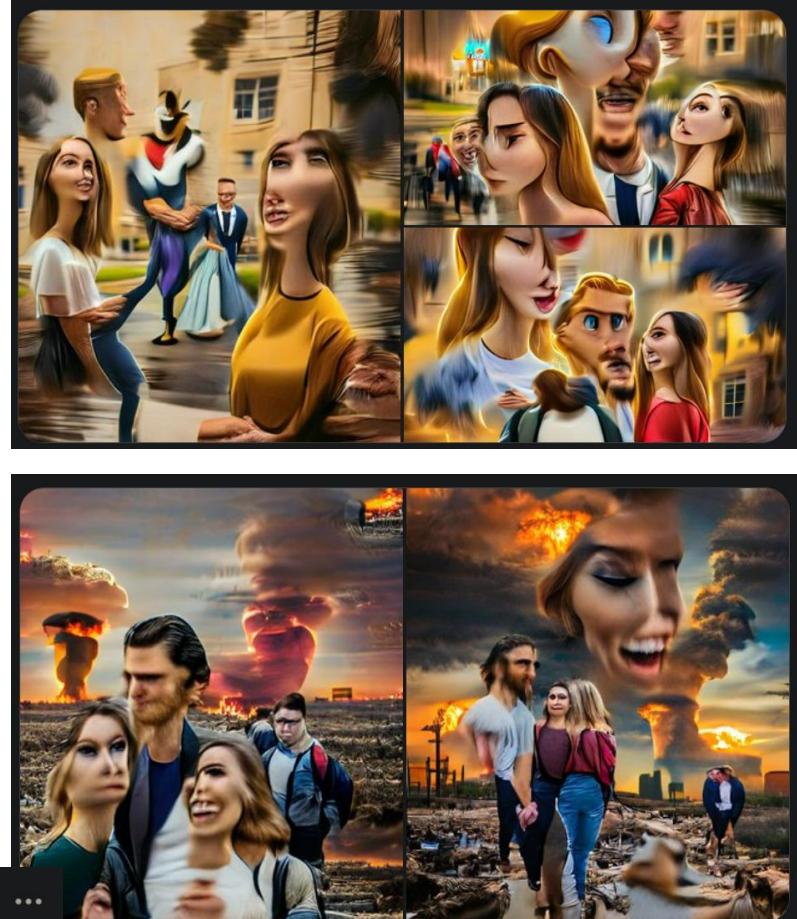
Mario Klingemann  @quasimondo · Jul 2
Can you guess the prompt? (and BTW there was no cheating involved by using a seed image)



CLIP

Mario Klingemann  @quasimondo · Jul 2

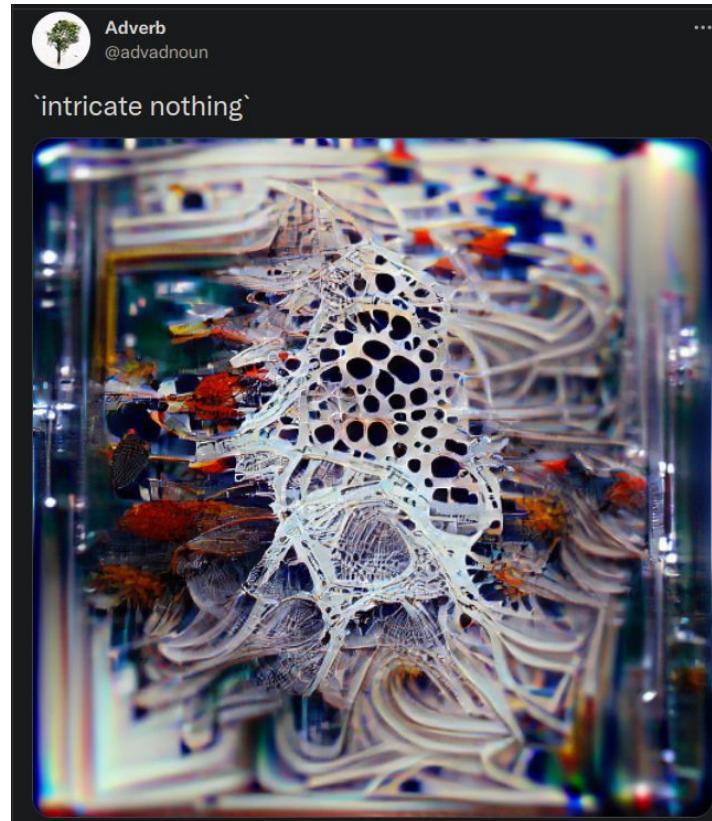
Can you guess the prompt? (and BTW there was no cheating involved by using a seed image)



Mario Klingemann  @quasimondo · Jul 2

It is just "distracted boyfriend meme" (plus some stylistic randomization, but no additional instructions as for the cast or placement)

Impressive results in many modalities!



Impressive results in many modalities!

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

Impressive results in many modalities!

(a) Kinetics 400

Method	Top 1	Top 5	Views
blVNet [16]	73.5	91.2	–
STM [30]	73.7	91.6	–
TEA [39]	76.1	92.5	10×3
TSM-ResNeXt-101 [40]	76.3	–	–
I3D NL [72]	77.7	93.3	10×3
CorrNet-101 [67]	79.2	–	10×3
ip-CSN-152 [63]	79.2	93.8	10×3
LGD-3D R101 [48]	79.4	94.4	–
SlowFast R101-NL [18]	79.8	93.9	10×3
X3D-XXL [17]	80.4	94.6	10×3
TimeSformer-L [2]	80.7	94.7	1×3
ViViT-L/16x2	80.6	94.7	4×3
ViViT-L/16x2 320	81.3	94.7	4×3
<i>Methods with large-scale pretraining</i>			
ip-CSN-152 [63] (IG [41])	82.5	95.3	10×3
ViViT-L/16x2 (JFT)	82.8	95.5	4×3
ViViT-L/16x2 320 (JFT)	83.5	95.5	4×3
ViViT-H/16x2 (JFT)	84.8	95.8	4×3

(b) Kinetics 600

Method	Top 1	Top 5	Views
AttentionNAS [73]	79.8	94.4	–
LGD-3D R101 [48]	81.5	95.6	–
SlowFast R101-NL [18]	81.8	95.1	10×3
X3D-XL [17]	81.9	95.5	10×3
TimeSformer-HR [2]	82.4	96.0	–
ViViT-L/16x2	82.5	95.6	4×3
ViViT-L/16x2 320	83.0	95.7	4×3
ViViT-L/16x2 (JFT)	84.3	96.2	4×3
ViViT-H/16x2 (JFT)	85.8	96.5	4×3

(c) Moments in Time

	Top 1	Top 5
TSN [69]	25.3	50.1
TRN [83]	28.3	53.4
I3D [6]	29.5	56.1
blVNet [16]	31.4	59.3
AssembleNet-101 [51]	34.3	62.7
ViViT-L/16x2	38.0	64.9

(d) Epic Kitchens 100 Top 1 accuracy

Method	Action	Verb	Noun
TSN [69]	33.2	60.2	46.0
TRN [83]	35.3	65.9	45.4
TBN [33]	36.7	66.0	47.2
TSM [40]	38.3	67.9	49.0
SlowFast [18]	38.5	65.6	50.0
ViViT-L/16x2 Fact. encoder	44.0	66.4	56.8

(e) Something-Something v2

Method	Top 1	Top 5
TRN [83]	48.8	77.6
SlowFast [17, 77]	61.7	–
TimeSformer-HR [2]	62.5	–
TSM [40]	63.4	88.5
STM [30]	64.2	89.8
TEA [39]	65.1	–
blVNet [16]	65.2	90.3
ViViT-L/16x2 Fact. encoder	65.4	89.8

Transformers for CV

Transformers are complex!

- Transformers are highly parametric!
 - Each **head** has \mathbf{W}_Q , \mathbf{W}_K and \mathbf{W}_V each of size $(d_w \times d_k)$.
 - Each **layer** also includes \mathbf{W}_O , \mathbf{W}_{F1} and \mathbf{W}_{F2} of sizes $(d_w \times d_w)$, $(d_w \times cd_w)$ and $(cd_w \times d_w)$, respectively.
 - By setting the **dimensionality** of the input to $d_w = 512$, the number of **heads** to $h = 8$, and $c = 4$ [3] **each layer ends up with ~2.5M parameters**.
 - **6 transformer layers** quickly turn into **~14.7M parameters**.

	# layers	# params
ResNet	20	0.27M
ResNet	32	0.46M
ResNet	44	0.66M
ResNet	56	0.85M
ResNet	110	1.7M
ResNet	1202	19.4M

Transformers are complex!

- Transformers are highly **parametric**!
- Transformers lack **inductive biases**:
 - That is why we need positional encodings.
 - Opposed to convolutions, it has to learn **translation equivariance** or importance **local relations** from the data [5].

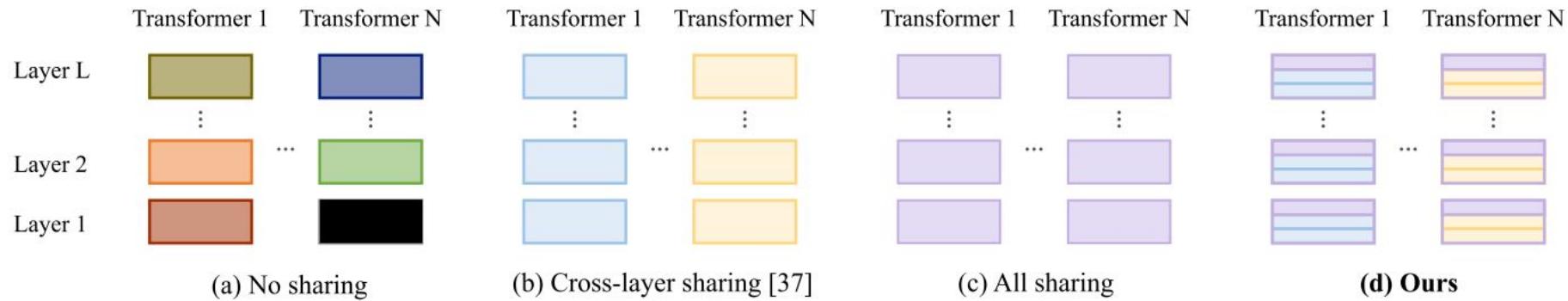
Transformers are complex!

- Transformers are highly **parametric**!
- Transformers lack **inductive biases**.
- Need a lot of resources!
 - **Big datasets** to allow the Transformer to learn
 - Long **training times** / **Big computational clusters (GPUs)**

Transformers are complex!

- Transformers are highly **parametric**!
- Transformers lack **inductive biases**.
- Need a lot of resources!
- Solutions:
 - **Efficient designs**: Share weights.
 - **Pre-training**: Use **unsupervised** / self-supervised techniques to leverage more data

Weight Sharing



c) Model	X.-L	X.-T	Params	top-1/5
Multi-2	✗	✗	7M	60.3 / 88.9
Multi-6	✓	✗	21M	65.7 / 89.9
Multi-6	✓	✓(All)	7M	67.1 / 92.3
Multi-6	✓	✓(Part [†])	4M	67.5 / 92.3

$$W = U\Sigma V^T$$

Pre-training

Use **large unlabeled** datasets to pre-train the model:

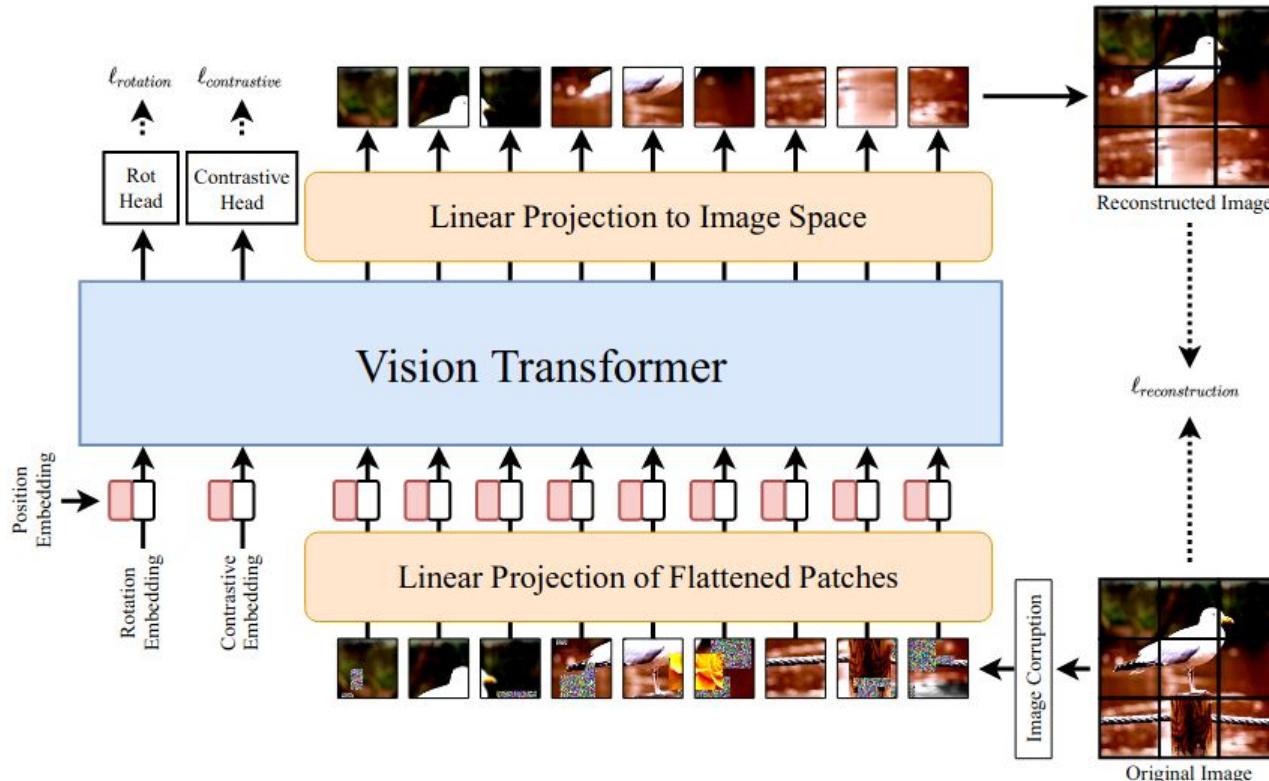
- Internal structure of the data.

Train with a **small** subset of **labeled dataset** for a specific task:

- Relevant **task specific information**.

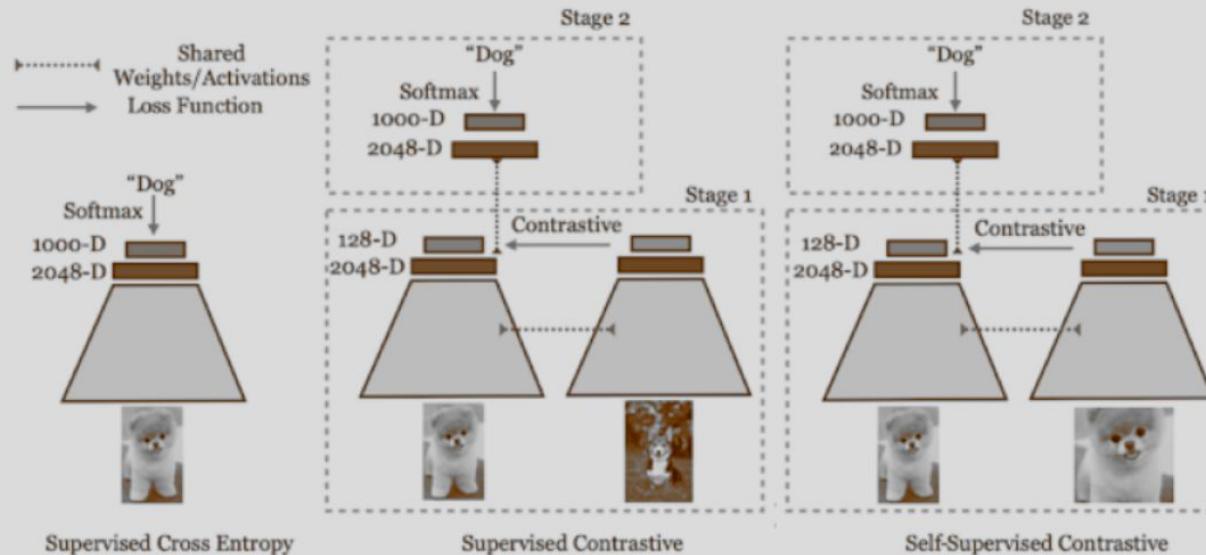


Pre-training: Reconstruction



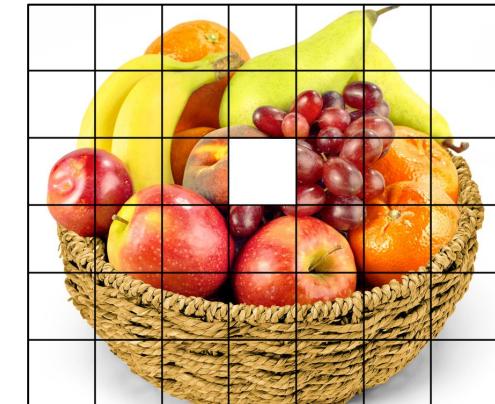
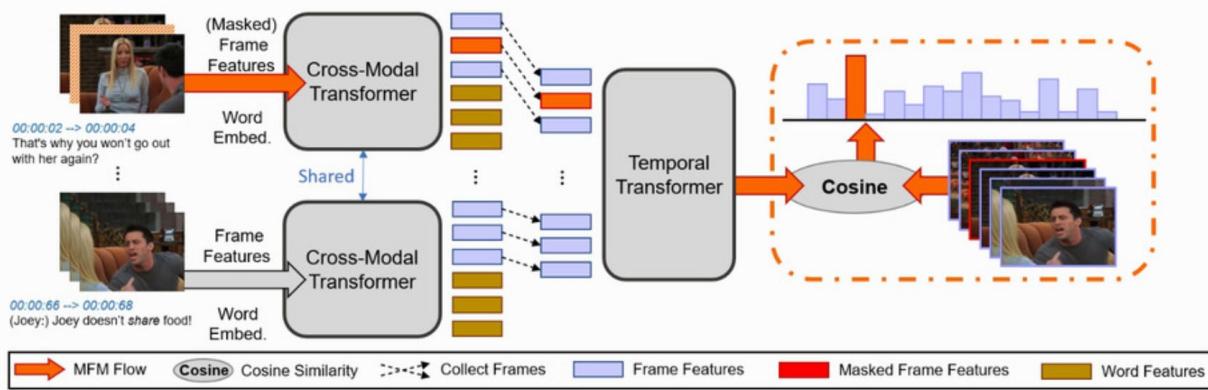
Previously...

Contrastive learning



Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., ... & Krishnan, D. (2020). Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*.

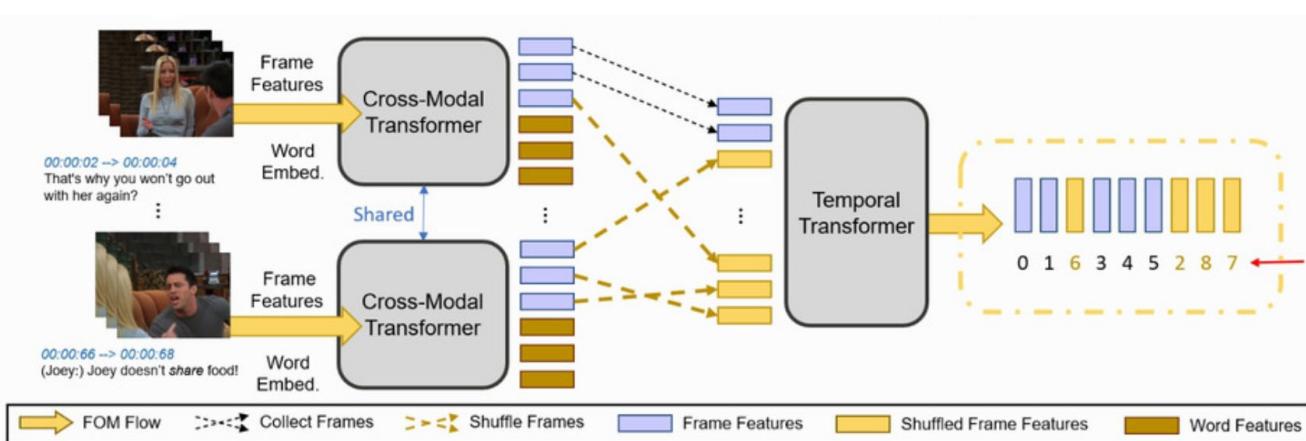
Pre-training: Contrastive



- Maximize Mutual information between input and output.
- Generally through **Noise Contrastive Estimation**

$$\mathcal{L}_{\text{NCE}}(\mathbf{x}, \tilde{\mathbf{o}}) = -\mathbb{E}_{\mathbf{x}} \left[\sum_t \log \frac{I(\mathbf{x}_t, \tilde{\mathbf{o}}_t)}{I(\mathbf{x}_t, \tilde{\mathbf{o}}_t) + \sum_{j \in \text{neg}(t)} I(\mathbf{x}_j, \tilde{\mathbf{o}}_t)} \right]$$

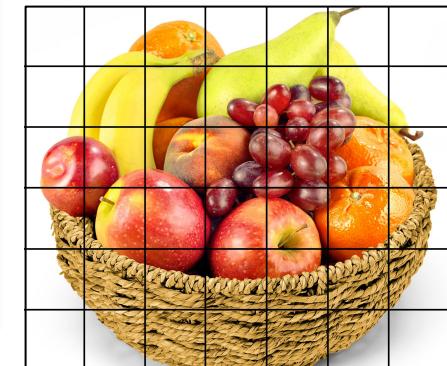
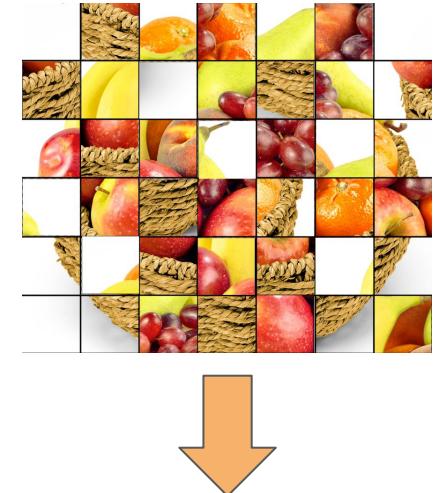
Pre-training: Token shuffling



$$\text{Reorder Indices: } \mathbf{r} = \{r_i\}_{i=1}^R \in \mathbb{N}^R$$

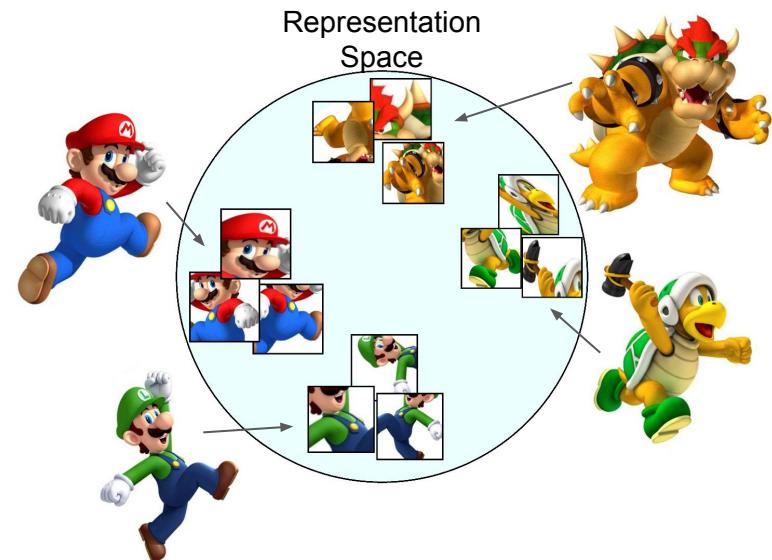
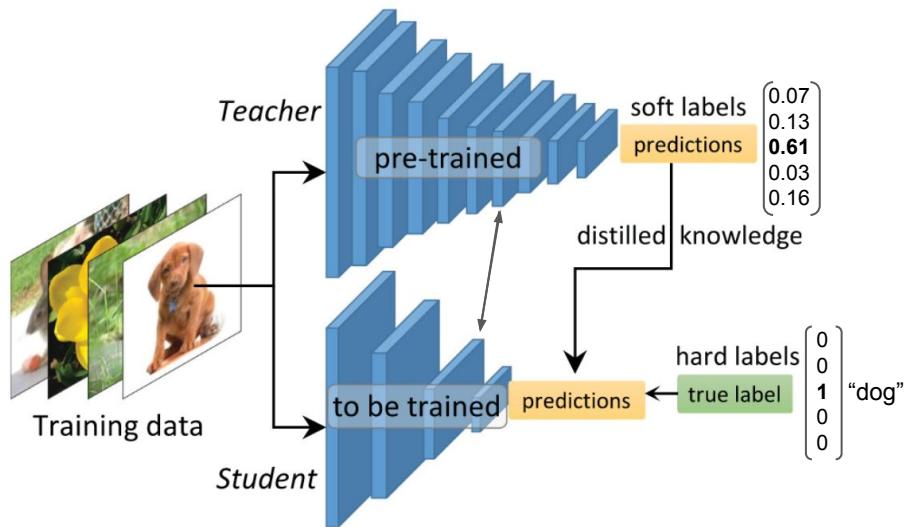
$$\text{Original timestamp: } \mathbf{t} = \{t_i\}_{i=1}^R$$

$$\text{Loss Function of FOM: } \mathcal{L}_{\text{FOM}} = -\mathbb{E}_D \sum_{i=1}^R \log \mathbf{P}[r_i, t_i]$$



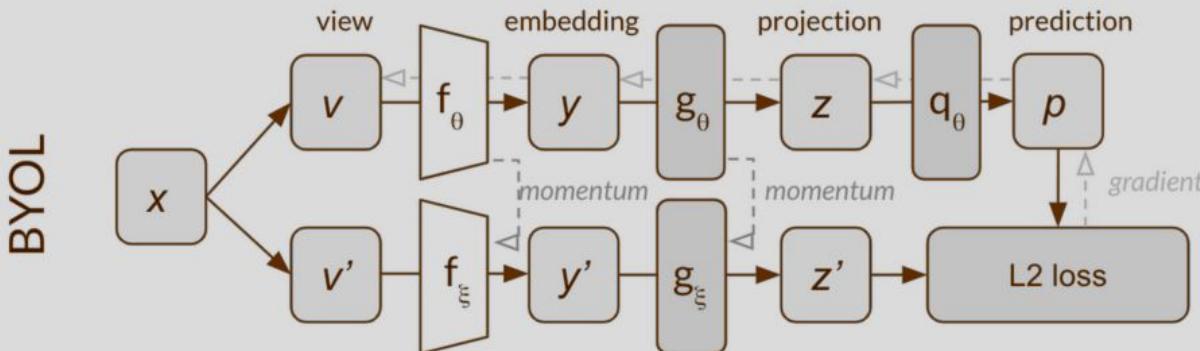
Pre-training: Momentum Contrast

Borrows ideas from **knowledge distillation**, **data augmentation** and **contrastive/self-supervised losses**.



Previously...

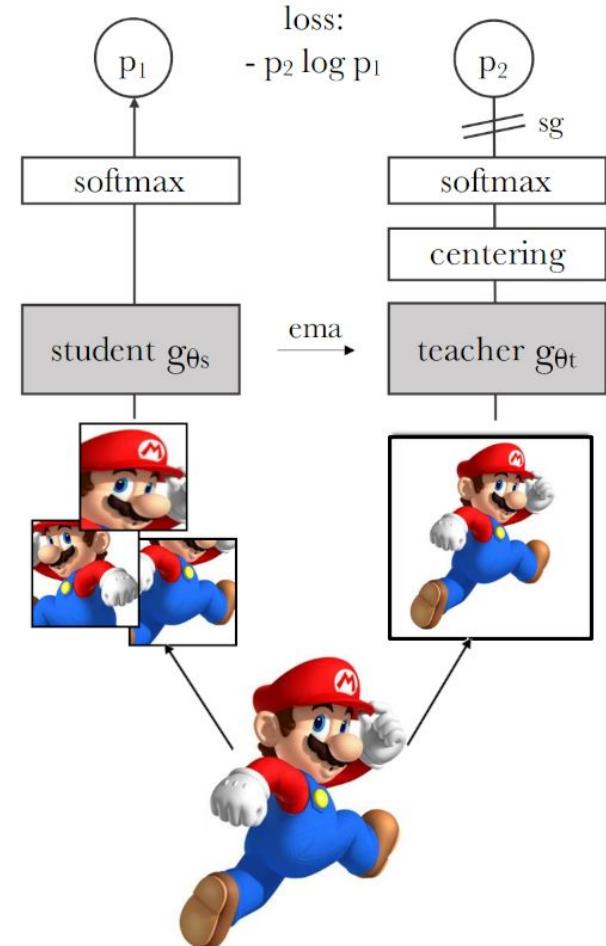
Distillation



Grill, J. B., Strub, F., Altché, F., Tallec, C., Richemond, P. H., Buchatskaya, E., ... & Valko, M. (2020). Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*.

Pre-training: Momentum Contrast

DINO

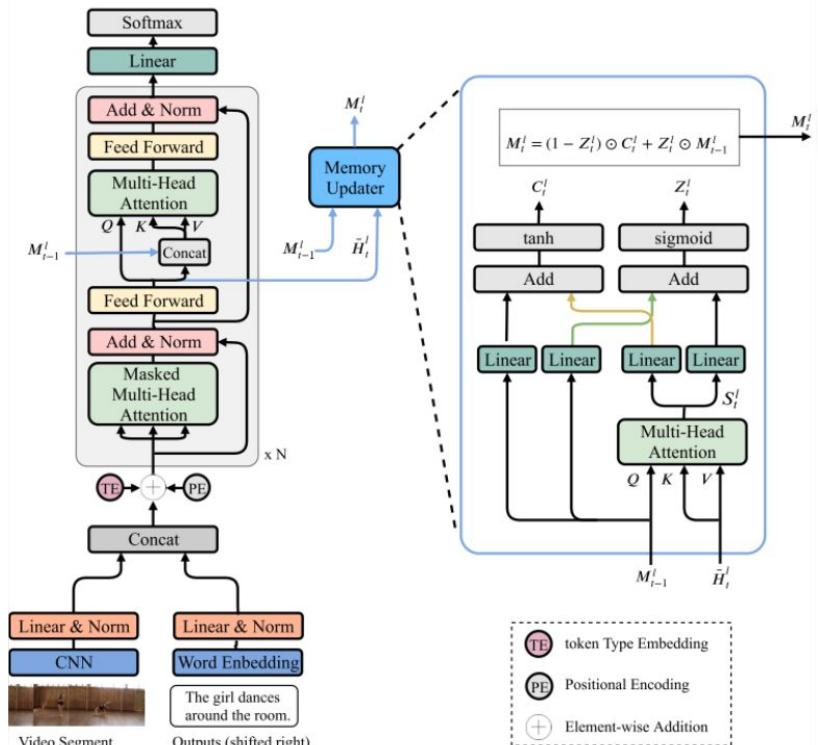
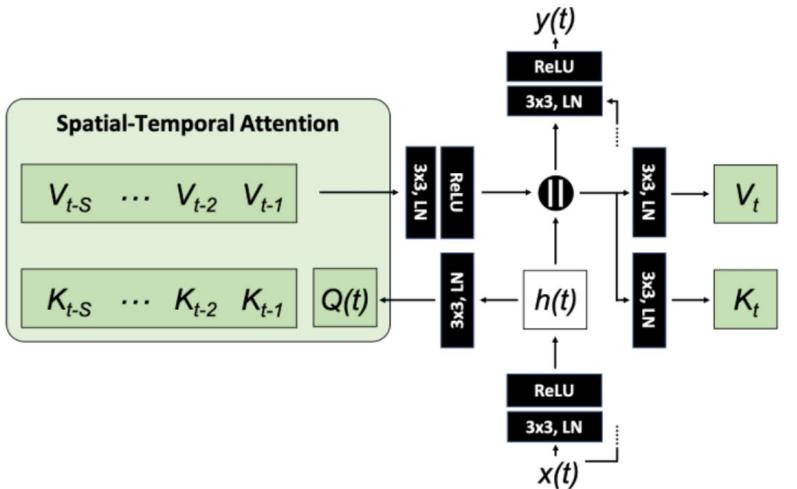


Pre-training: Momentum Contrast

DINO



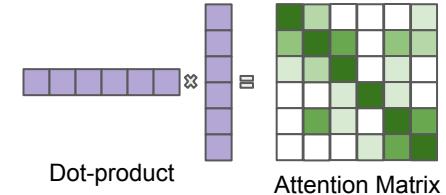
Recurrent Transformers



Efficient Transformers

Attention matrix grows fast!

The size of the attention matrix grows
quadratically with the number of input tokens
(sequence length T) $\rightarrow O(T^2)$

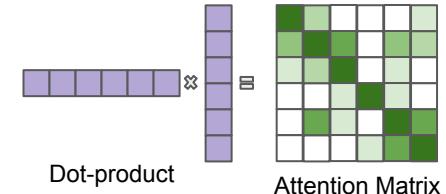


Transformers could adapt to arbitrarily long inputs, but this limits them to do so!!

Attention matrix grows fast!

The size of the attention matrix grows
quadratically with the number of input tokens
(sequence length T) $\rightarrow O(T^2)$

Transformers could adapt to arbitrarily long inputs, but this limits them
to do so!!



Solution: Sparse / Decomposed attention.

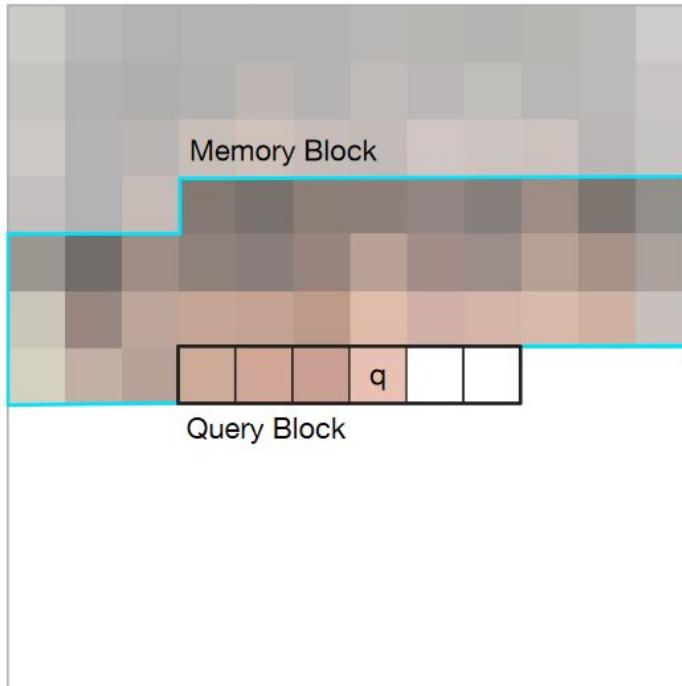
Limit the number of tokens in a single attention operation.

Efficient Designs (Locality/Sparsity)

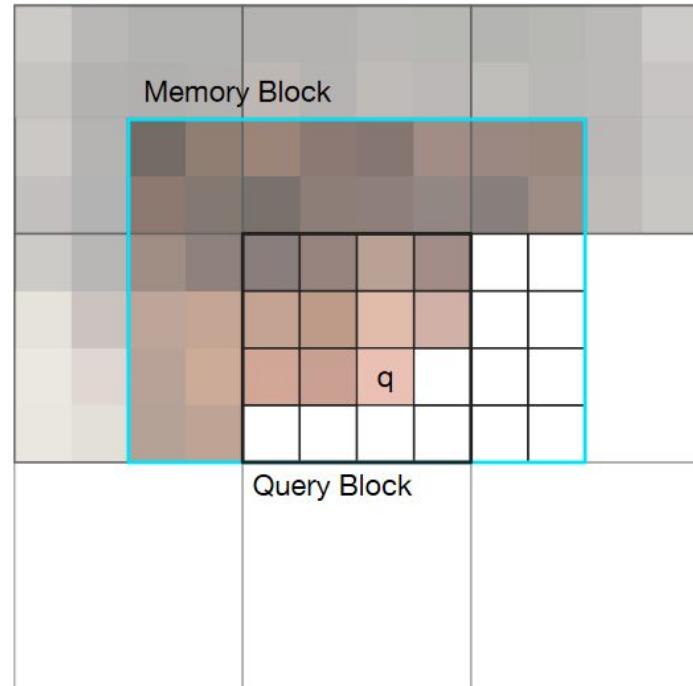
- Limit the number of tokens accessible in a single operation
- **Local:** Only allow to attend to neighboring tokens
- **Sparse:** Tokens can be distant, but only every other token is accessible.

Efficient Designs (Locality/Sparsity)

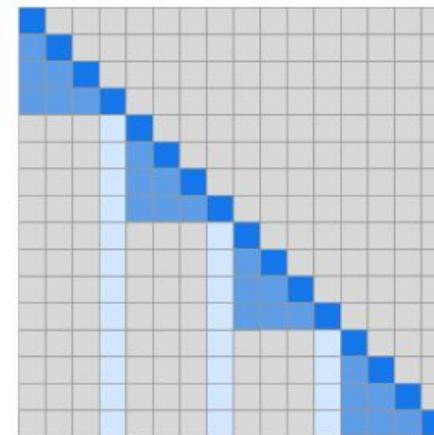
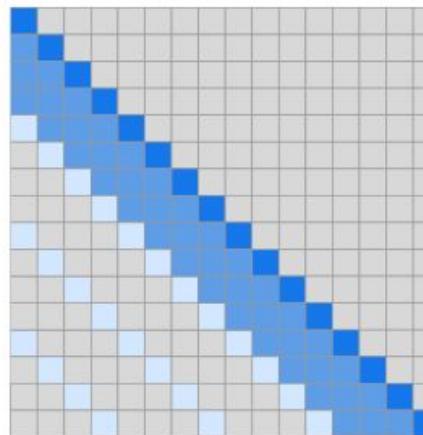
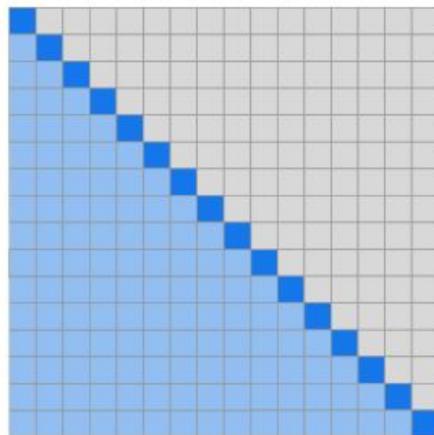
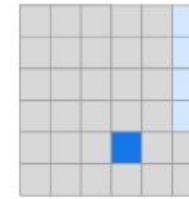
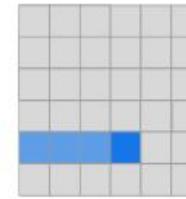
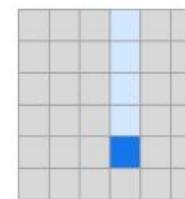
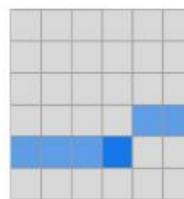
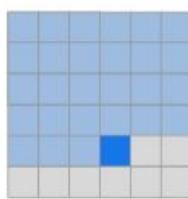
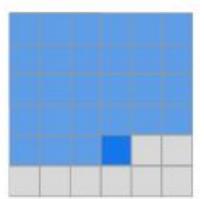
Local 1D Attention



Local 2D Attention



Efficient Designs (Locality/Sparsity)

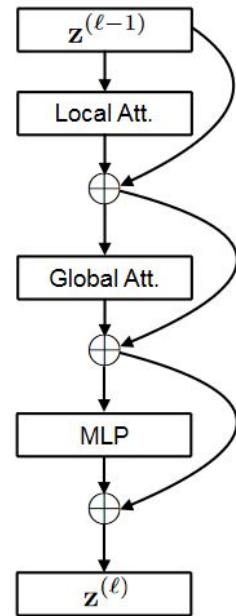


(a) Transformer

(b) Sparse Transformer (strided)

(c) Sparse Transformer (fixed)

Efficient Designs (Locality/Sparsity)

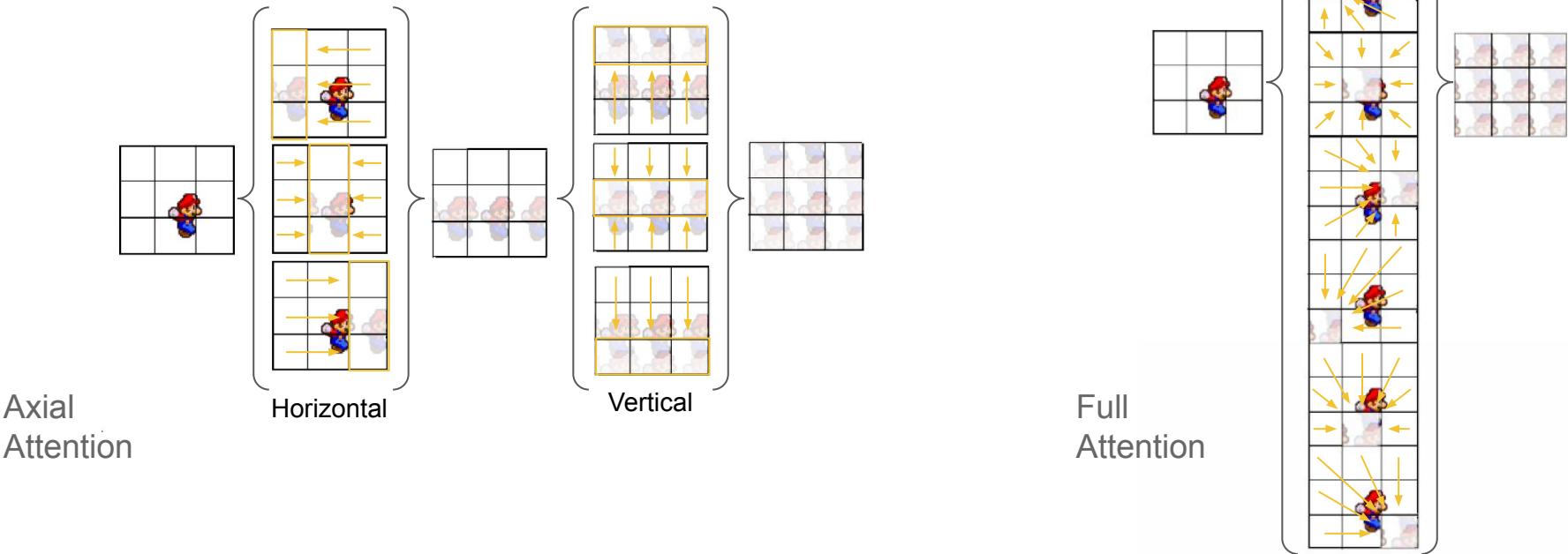


Sparse Local Global
Attention (L+G)

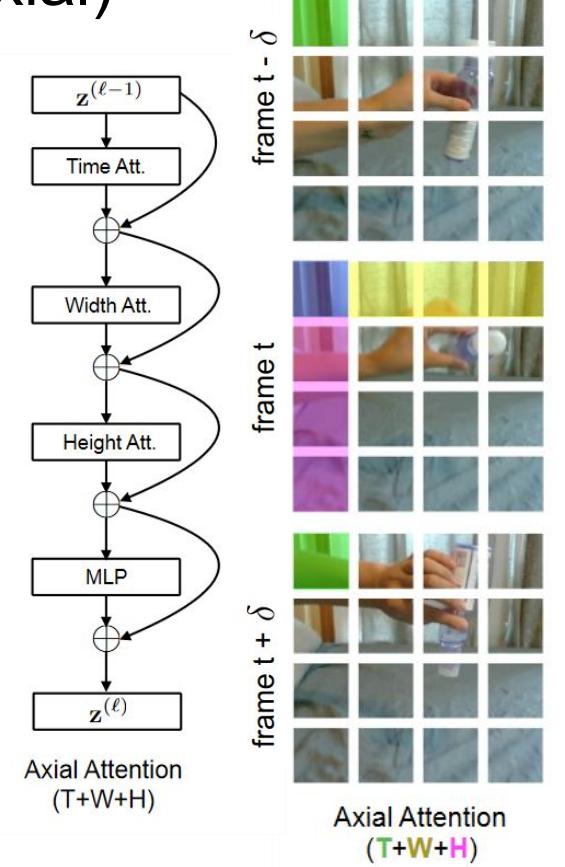
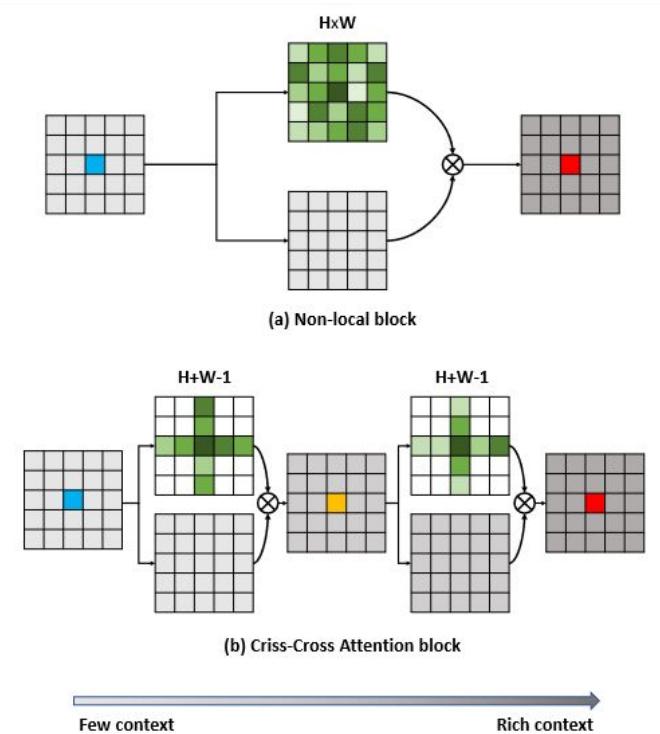


Efficient Designs (Decomposition)

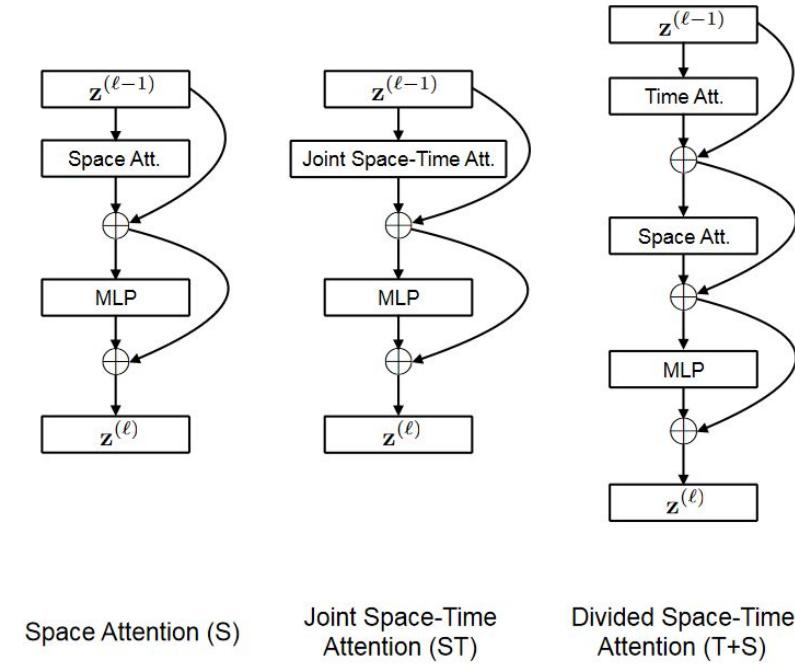
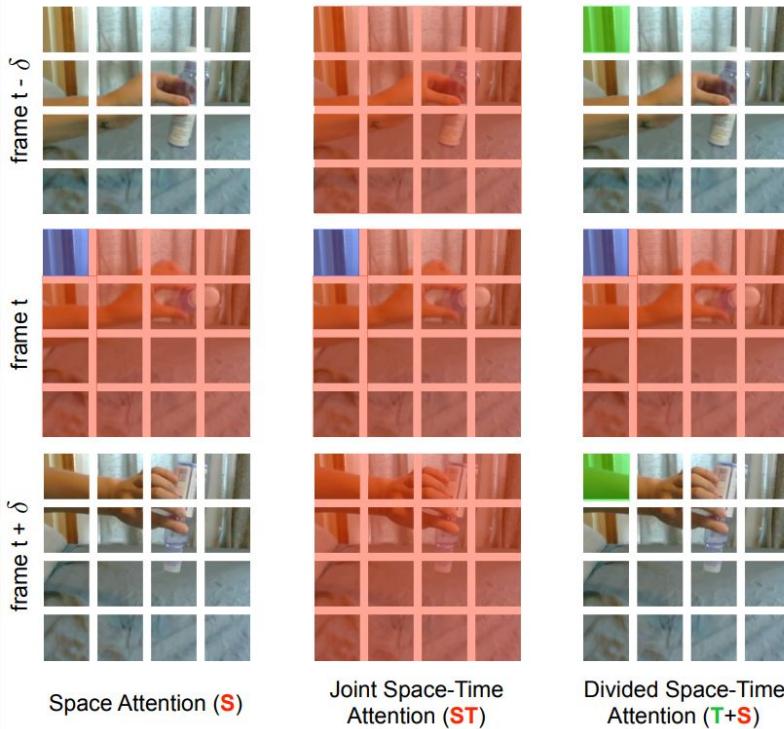
Approximate full attention with two or more consecutive attention operations



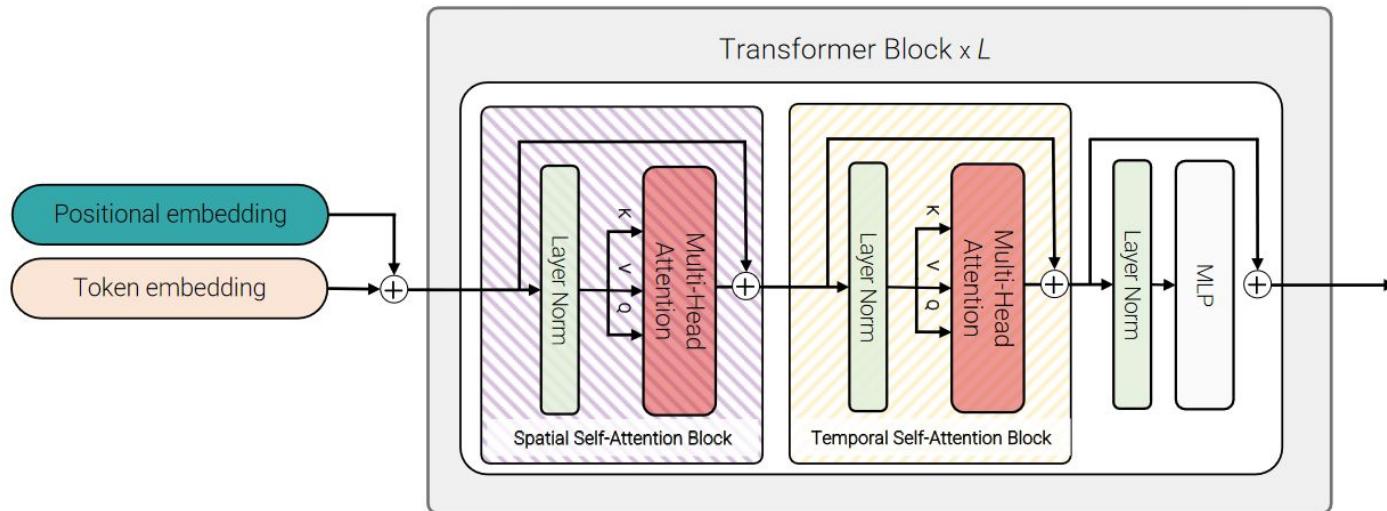
Efficient Designs (Decomposition - Axial)



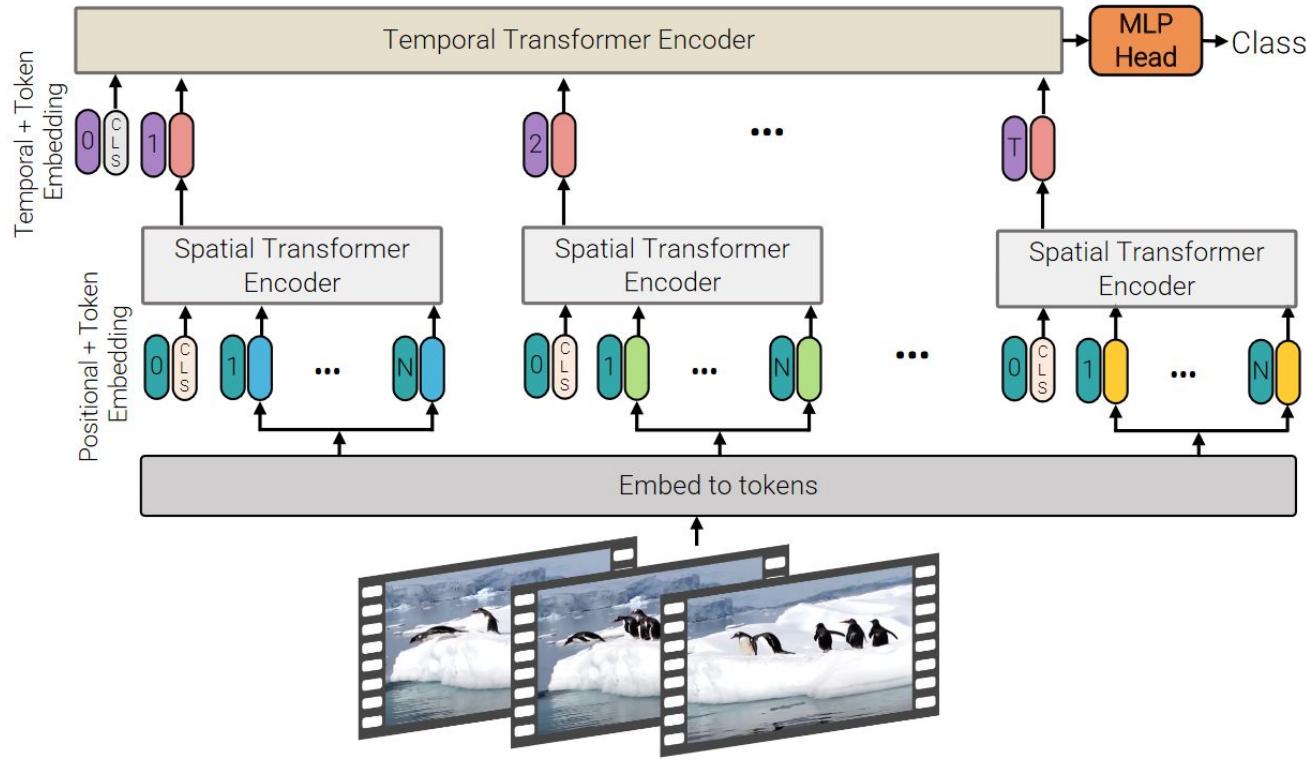
Efficient Designs (Decomposition - SpatioTemporal)



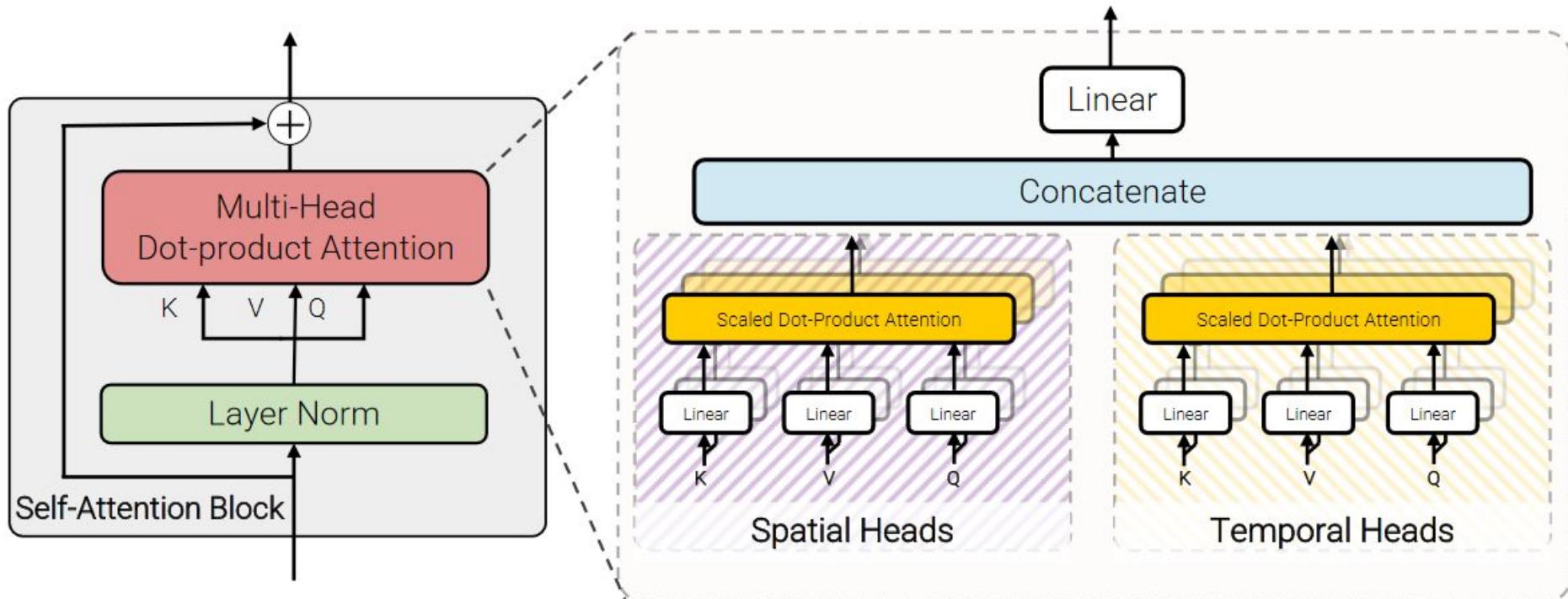
Efficient Designs (Decomposition - SpatioTemporal)



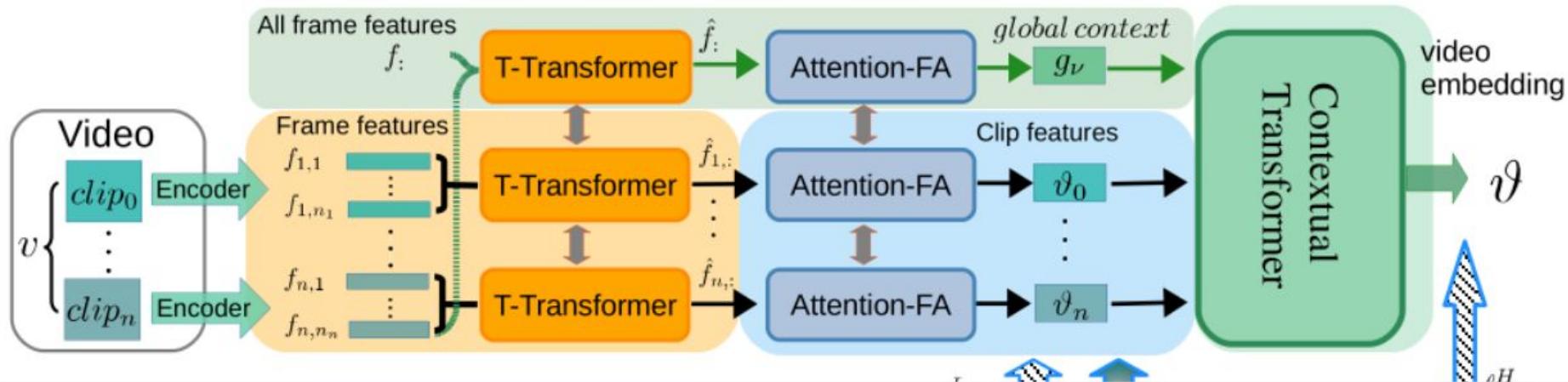
Efficient Designs (Decomposition - SpatioTemporal)



Efficient Designs (Decomposition - SpatioTemporal)



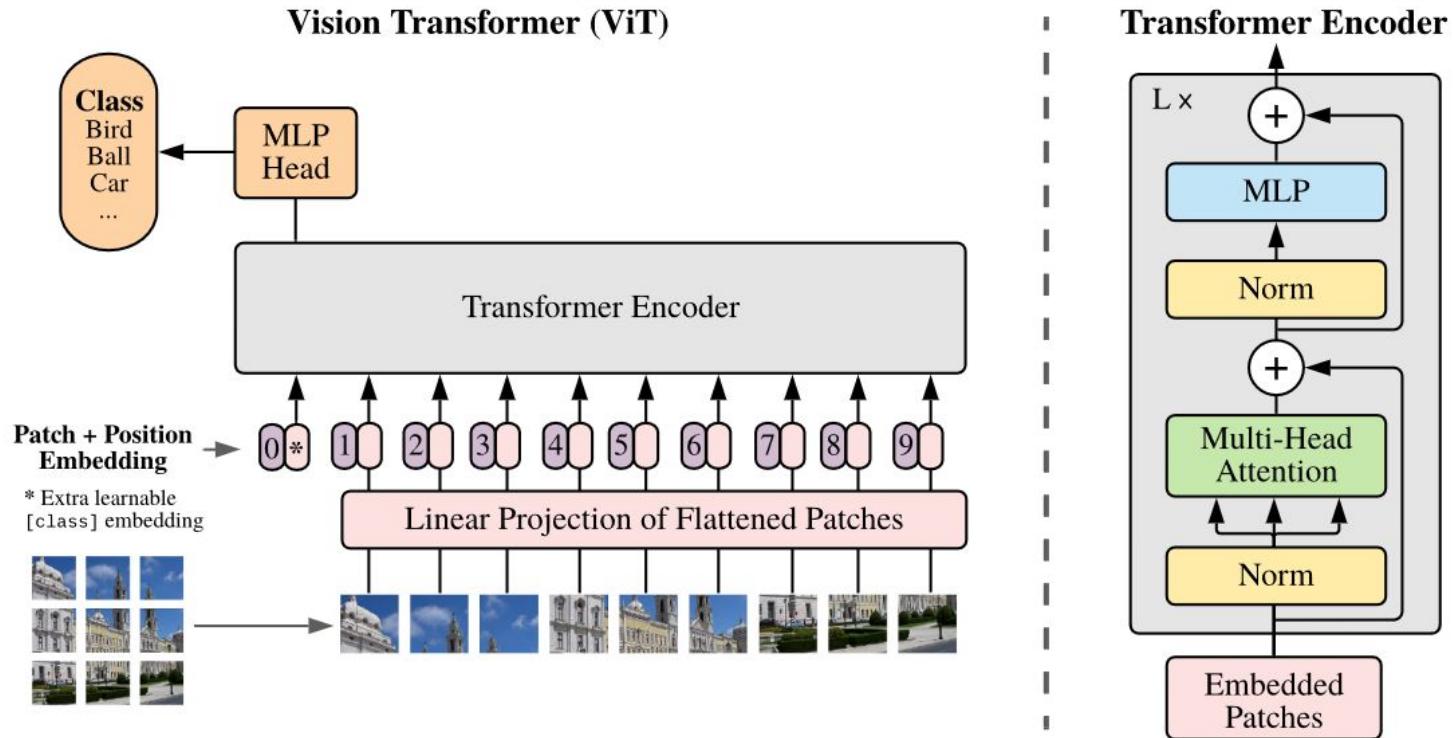
Efficient Designs (Hierarchy)



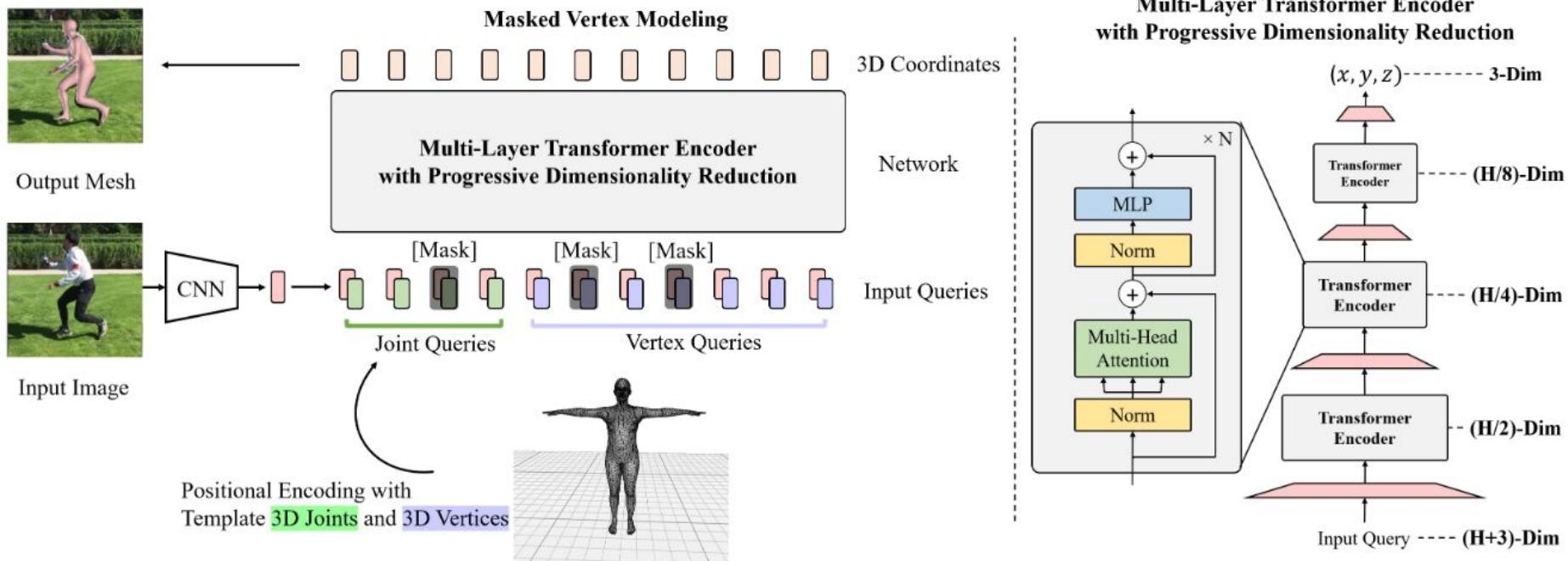
$$A = \text{softmax}(W_2 Q + b_2)^T, \quad Q = \text{GELU}(W_1 K^T + b_1), \quad K = X$$

Applications

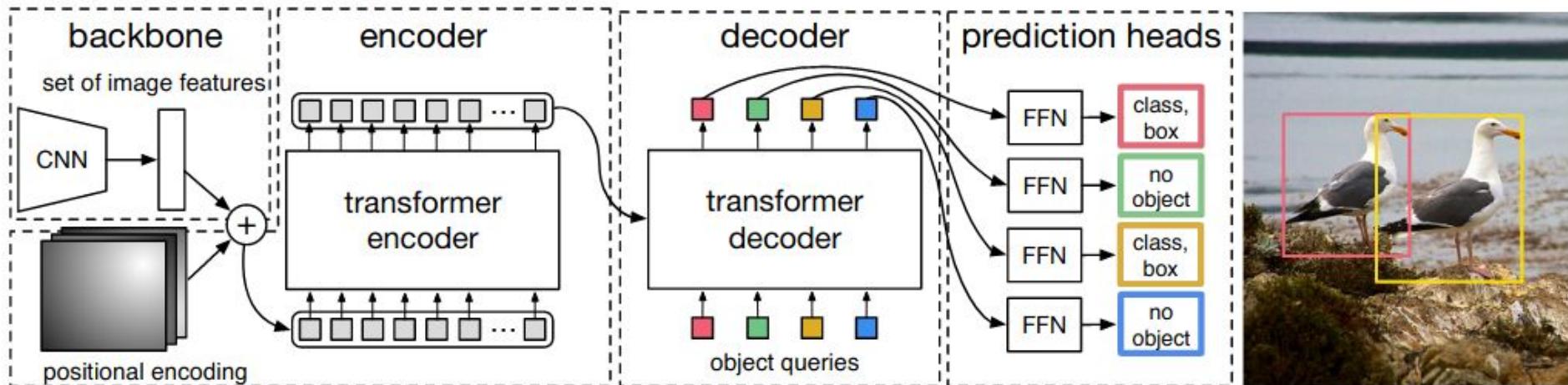
Vision Transformer (for classification)



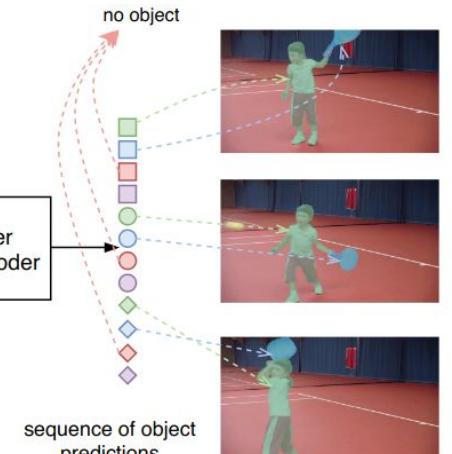
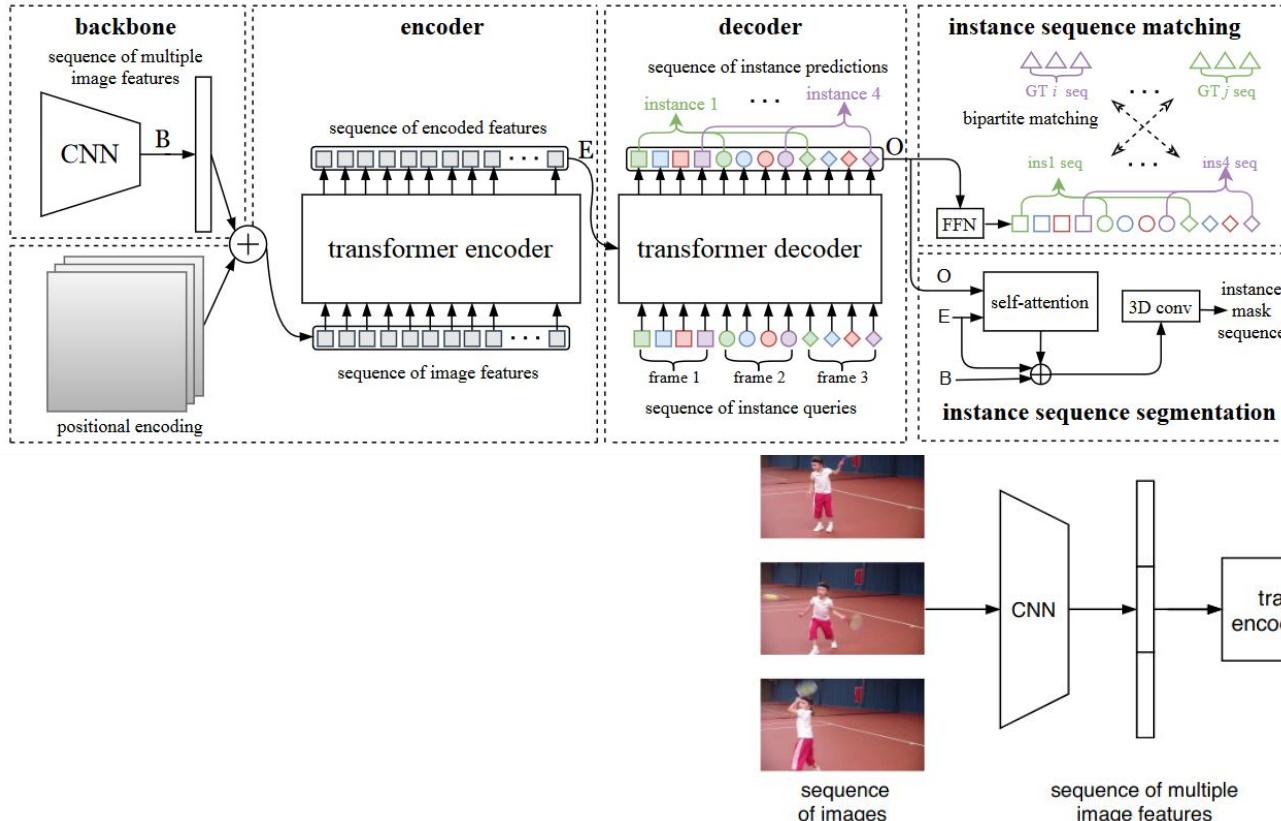
Mesh Reconstruction



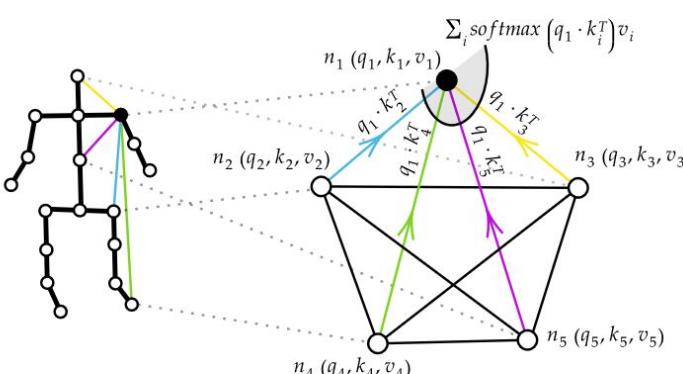
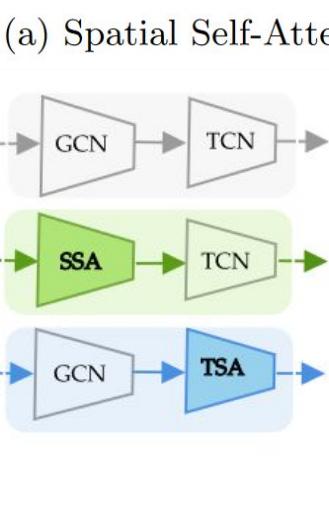
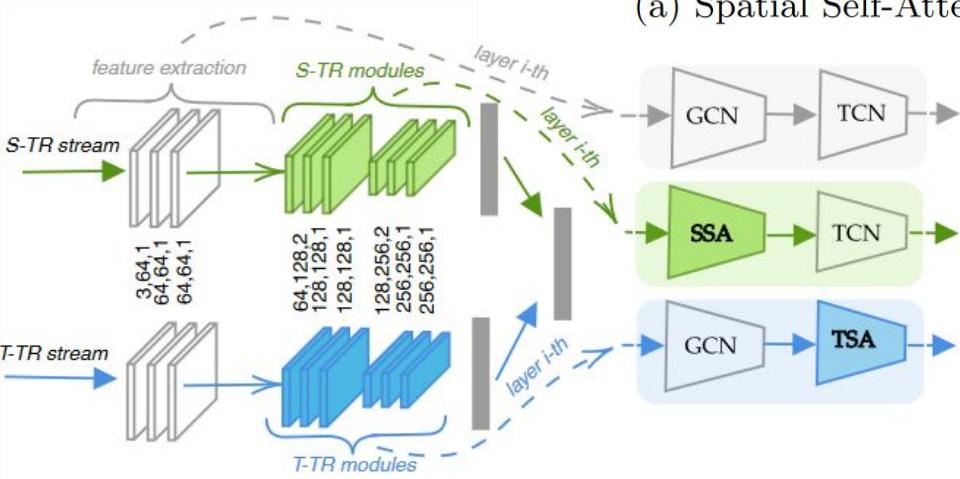
DETR (Transformer for object recognition)



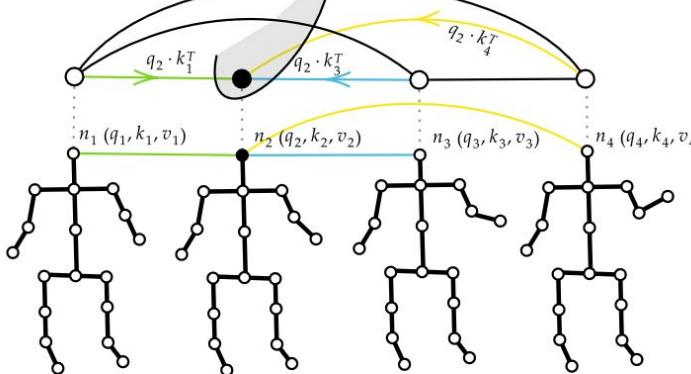
Transformer for video Segmentation



Pose (Action Recognition)

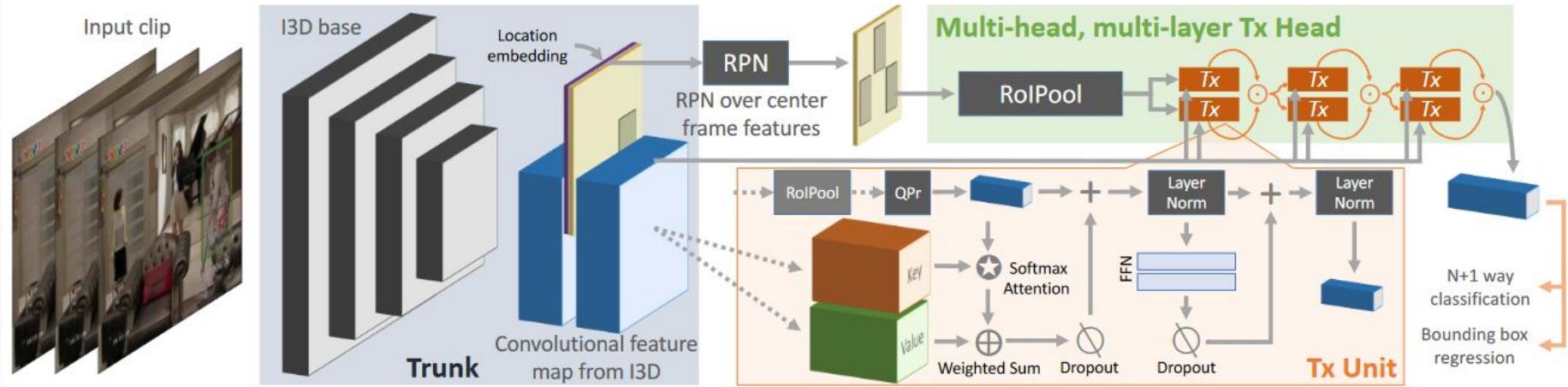


(a)
Spatial Self-Attention

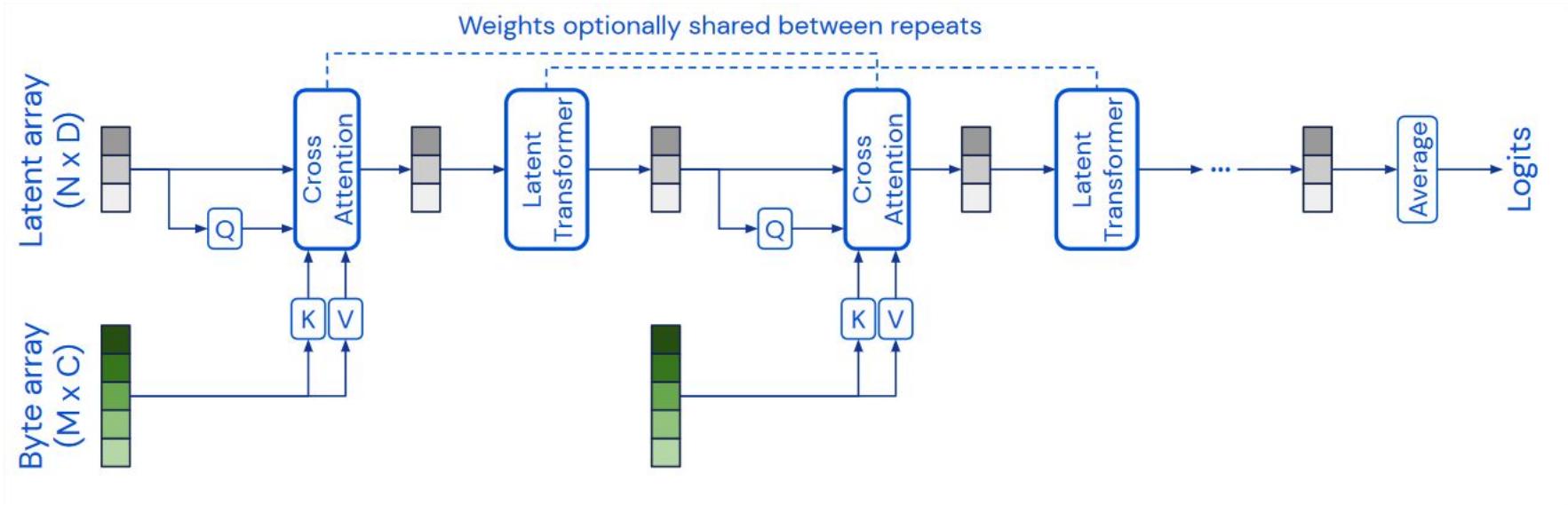


(b)
Temporal Self-Attention

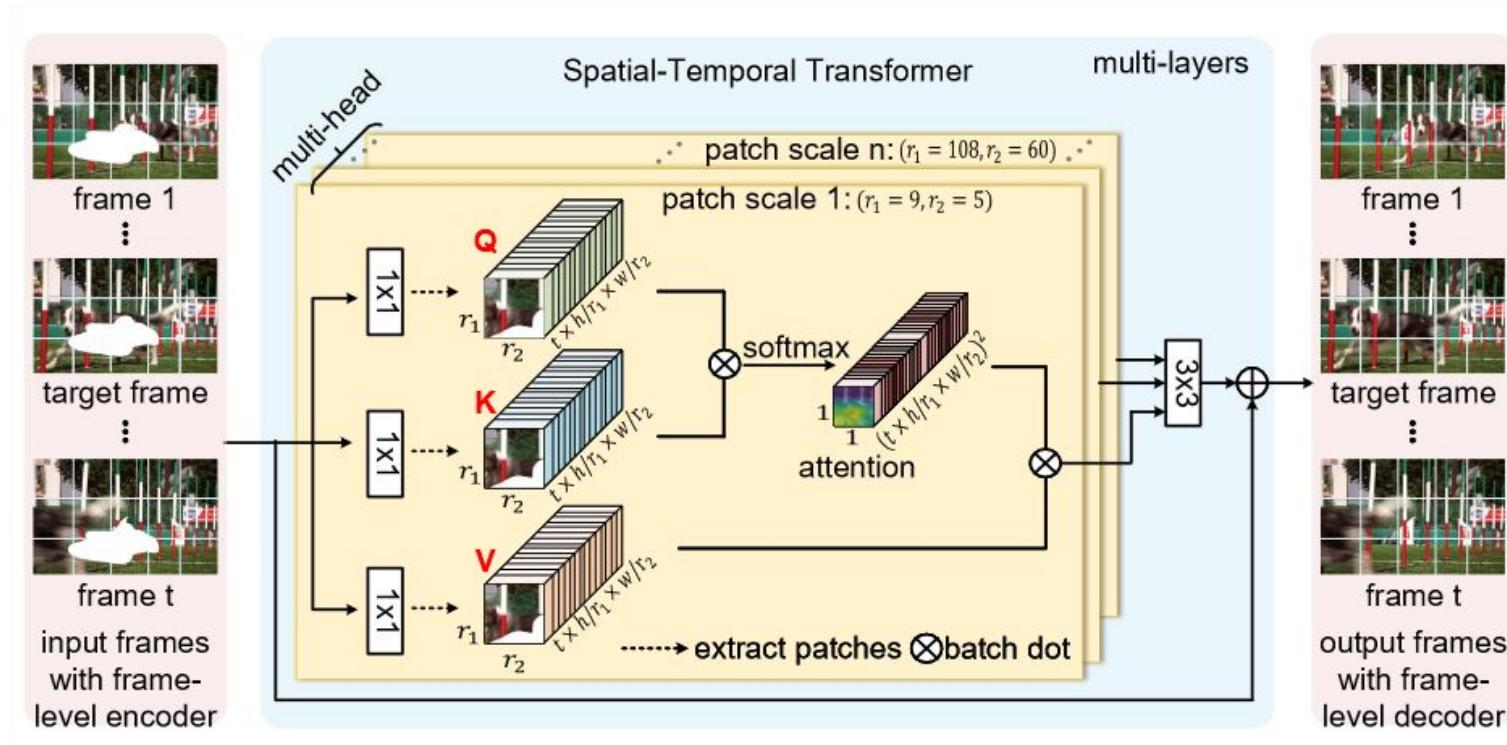
Video Action Transformer



Perceiver



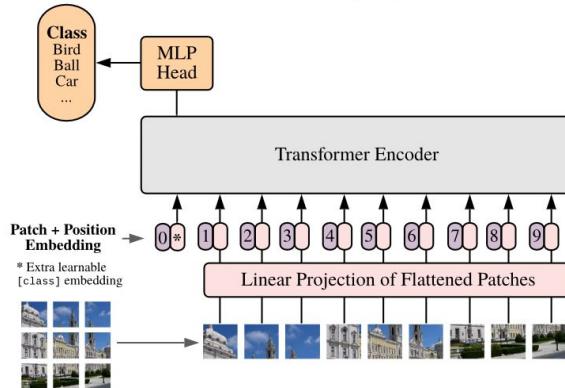
Inpainting



How to Transformer?

Use pre-trained models! GitHub!

<https://colab.research.google.com/>



A screenshot of the Hugging Face Model Hub interface, displaying a grid of pre-trained Vision Transformer models. Each model card includes the repository name, description, and metrics like downloads and stars.

- google/vit-base-patch16-224-in21k
- google/vit-base-patch16-224
- M-CLIP/M-BERT-Base-ViT-B
- openai/clip-vit-base-patch32
- openai/clip-vit-base-patch16
- lysandre/tiny-vit-random
- google/vit-large-patch16-224
- nielsr/dino_vitb16
- google/vit-huge-patch14-224-in21k
- hbredin/VoiceActivityDetection-PyanNet-DIHARD
- google/vit-base-patch32-224-in21k
- Vitafeu/DialoGPT-medium-ricksanchez
- facebook/dino-vitb16
- google/vit-base-patch16-384
- google/vit-base-patch32-384
- sentence-transformers/clip-ViT-B-32-multilingual...
- nateraw/vit-base-patch16-224-cifar10
- google/vit-large-patch32-384

<https://huggingface.co/models?sort=downloads&search=vit>

References

Publications

- [1] Bahdanau, D. et al. "Neural machine translation by jointly learning to align and translate". In CoRR (2014).
- [2] Wang, X. et al. "Non-local neural networks". In CVPR (2018).
- [3] Vaswani, A. et al. "Attention is all you need". In NeurIPS (2017)
- [4] He, K. et al. "Deep residual learning for image recognition". In CVPR (2016).
- [5] Dosovitskiy, A. et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". In ICLR (2020).
- [6] Lee, S. et al. "Parameter Efficient Multimodal Transformers for Video Representation Learning". In ICLR (2020).
- [7] Atito, S. et al. "Sit: Self-supervised vision transformer". In CoRR (2021).
- [8] Li, L. et al. "HERO: Hierarchical Encoder for Video+ Language Omni-representation Pre-training". In EMNLP (2020).
- [9] Oord, A. V. D. et al. "Representation learning with contrastive predictive coding". In CoRR (2018).
- [10] Wang, J. et al. "DeepVID: Deep Visual Interpretation and Diagnosis for Image Classifiers via Knowledge Distillation". In IEEE Tr-VCG (2019).
- [11] Caron, M. et al. "Emerging Properties in Self-Supervised Vision Transformers". In CoRR (2021).
- [12] Carion, N. et al. "End-to-end object detection with transformers". In ECCV (2020).
- [13] Lin, K. et al. "End-to-end human pose and mesh reconstruction with transformers". In CVPR (2021).
- [14] Jaegle, A. et al. "Perceiver: General Perception with Iterative Attention". In ICML (2021).
- [15] Plizzari, C. et al. "Spatial temporal transformer network for skeleton-based action recognition". In ICPR (2021).
- [16] Wang, Y. et al. "End-to-end video instance segmentation with transformers". In CVPR (2021).
- [17] Girdhar, R. et al. "Video Action Transformer Network". In CVPR (2020).
- [18] Zeng, Y. et al. "Learning joint spatial-temporal transformations for video inpainting". In ECCV (2020).
- [19] Parmar, N. et al. "Image transformer". In ICML (2018).
- [20] Child, R. et al. "Generating long sequences with sparse transformers". In CoRR (2019).
- [21] Bertasius, G. et al. "Is Space-Time Attention All You Need for Video Understanding?". In CoRR (2021).
- [22] Huang, Z. et al. "Ccnet: Criss-cross attention for semantic segmentation". In ICCV (2019).
- [23] Arnab, A. et al. "Vivit: A video vision transformer". In CoRR (2021).
- [24] Ging, S. et al. "Coot: Cooperative hierarchical transformer for video-text representation learning". In NeurIPS (2020).

Blog Posts

- Vincent Dumoulin, Francesco Visin "A guide to convolution arithmetic for deep learning" (2016) https://github.com/vdumoulin/conv_arithmetic
- Mohammed Terry-Jack "Deep Learning: The Transformer" (2019) <https://medium.com/@b.terryjack/deep-learning-the-transformer-9ae5e9c5a190>
- Lilian Weng "Attention? Attention!" (2018) <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>
- Amirhossein Kazemnejad "Transformer Architecture: The Positional Encoding" (2019) https://kazemnejad.com/blog/transformer_architecture_positional_encoding/
- Gillaume Klein et al. "The Annotated Transformer" (2018) <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
- Lilian Weng "The Transformer Family" (2020) <https://lilianweng.github.io/lil-log/2020/04/07/the-transformer-family.html>
- Jay Alammar "The Illustrated Transformer" (2018) <https://jalammar.github.io/illustrated-transformer/>