

The background of the slide features a dark, star-filled night sky. A vibrant green aurora borealis arches across the upper portion of the frame, its light glowing against the dark backdrop. In the lower half, the silhouettes of rugged mountain peaks are visible, their dark forms contrasting with the bright celestial bodies above.

PerceiverIO

One model to rule them all

Jörg Simon, 13.Oct.2021

About me

- PhD on using DeepLearning to detect Human Factors from BioSignals
- Prof. Eduardo Veas and Herbert Danzinger

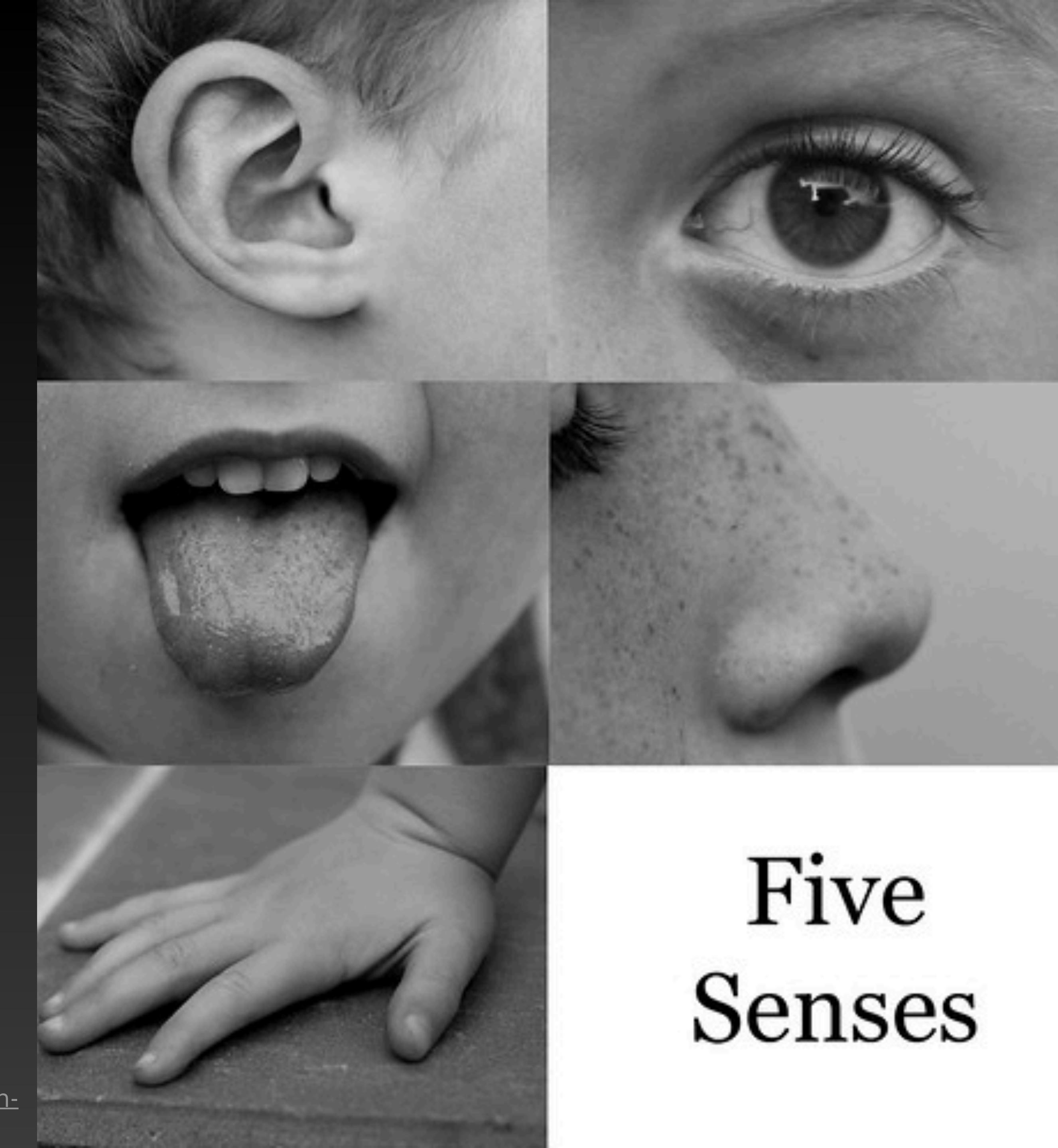


Outline

PerceiverIO: One model to rule them all

- Perceiving & Multimodality
- Attention
- Perceiver
- Quick what is PerceiverIO
- War of Frameworks: JAX, TF, PyTorch....
- Some small code examples (in PyTorch)

Perceiving & Multimodality



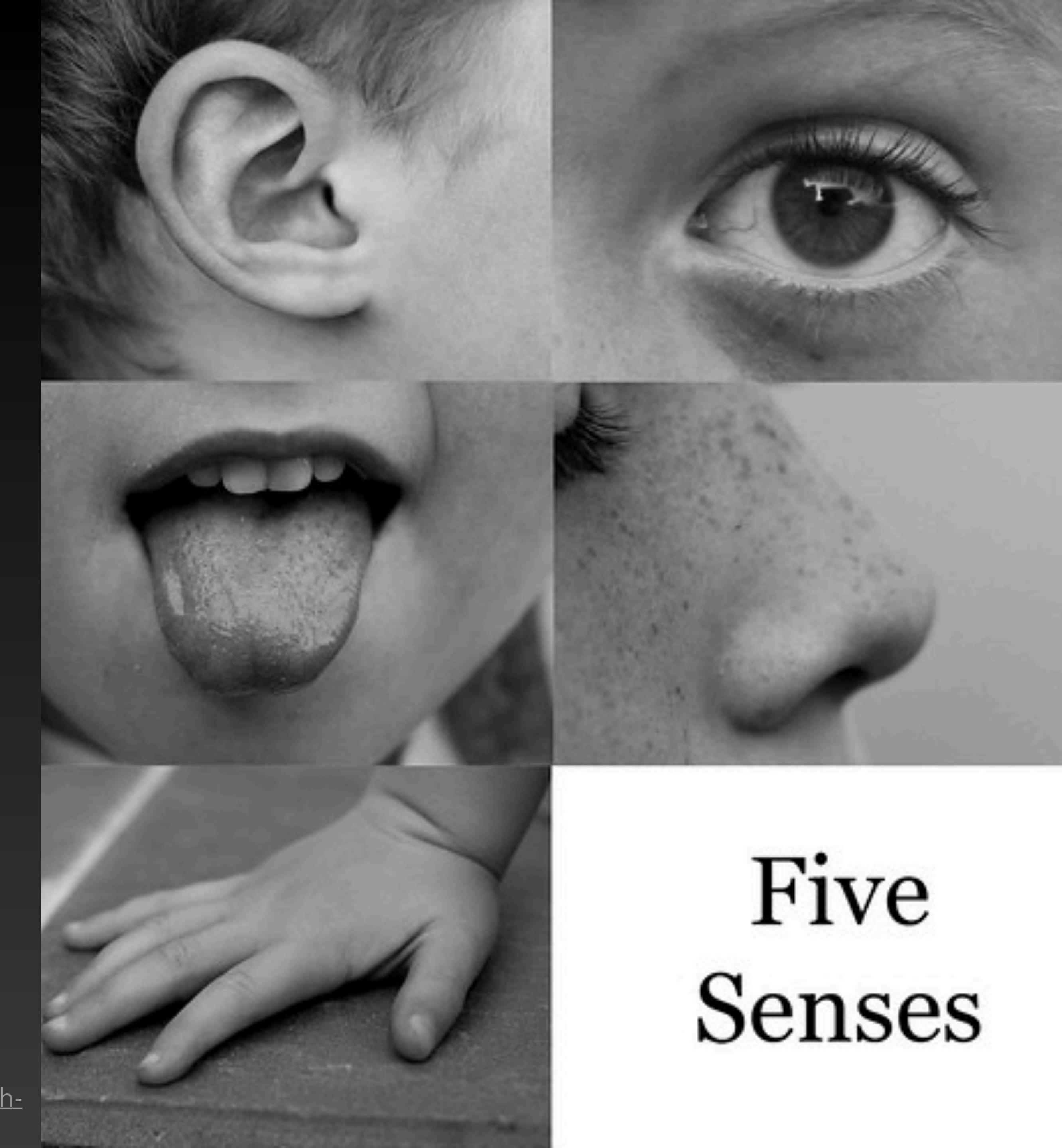
Five
Senses

Image from

<https://quizlet.com/539807260/chapter-4-sensing-and-perceiving-our-world-flash-cards/>

Perceiving & Multimodality

- Biological Systems Inherently Multi Modal
- DeepLearning:
 - CNN for image like Data
 - LSTMs for Time Series
 - Hourglas Networks
 - Strong Adaptation of Architecture towards Domain
 - Especially Input and Output very specific



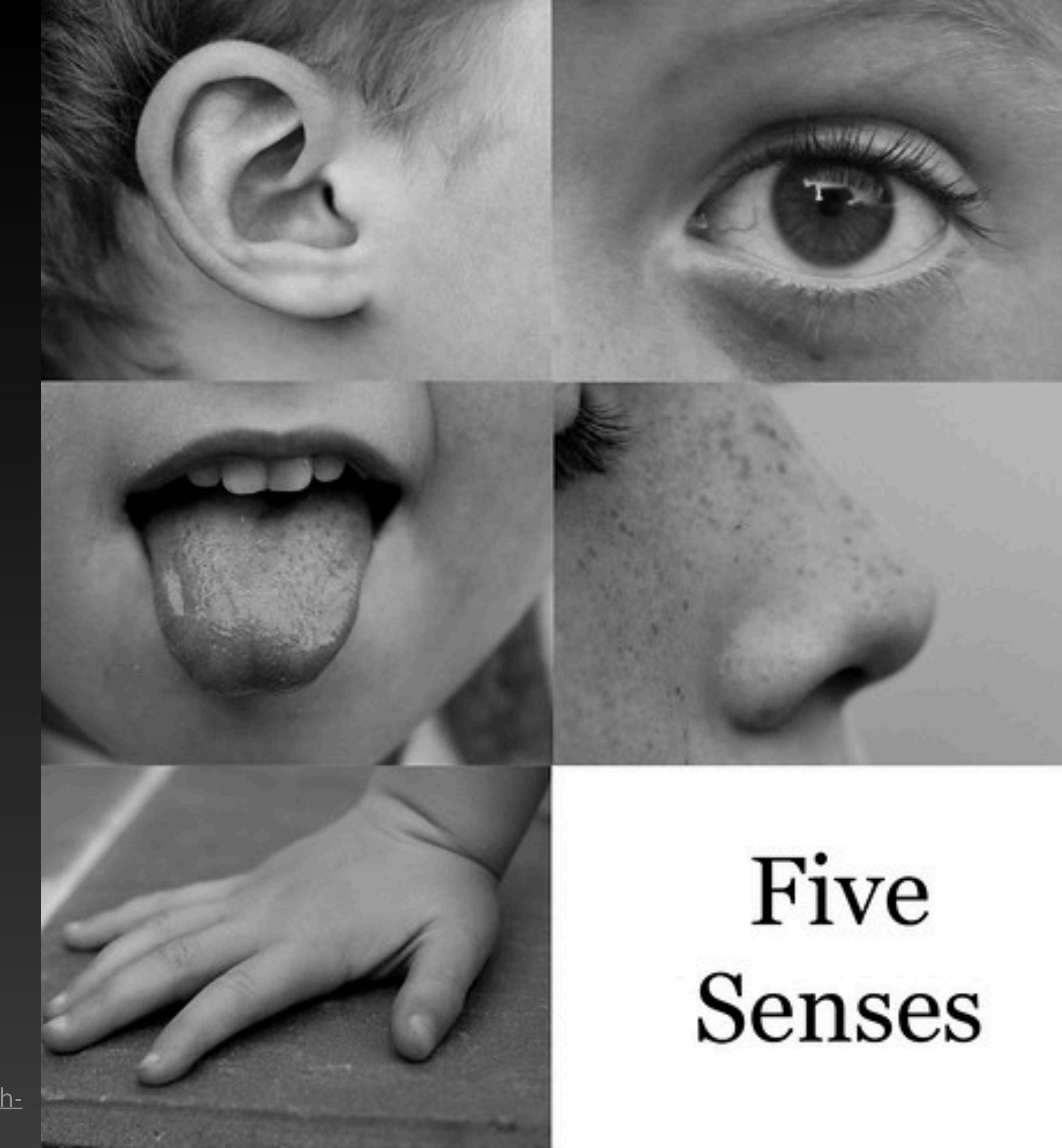
Five
Senses

Image from

<https://quizlet.com/539807260/chapter-4-sensing-and-perceiving-our-world-flash-cards/>

Perceiving & Multimodality

- Biological Systems Inherently Multi Modal
- DeepLearning:
 - Modality Specific contradicts the universal algorithm theorem
 - However Brain also has specialised Regions
 - Still it would be awesome to have one good architecture for all modalities and queries



Five
Senses

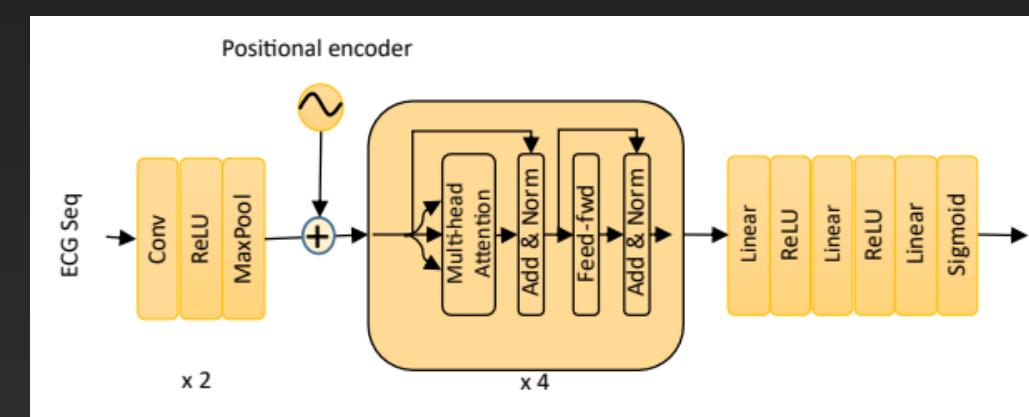
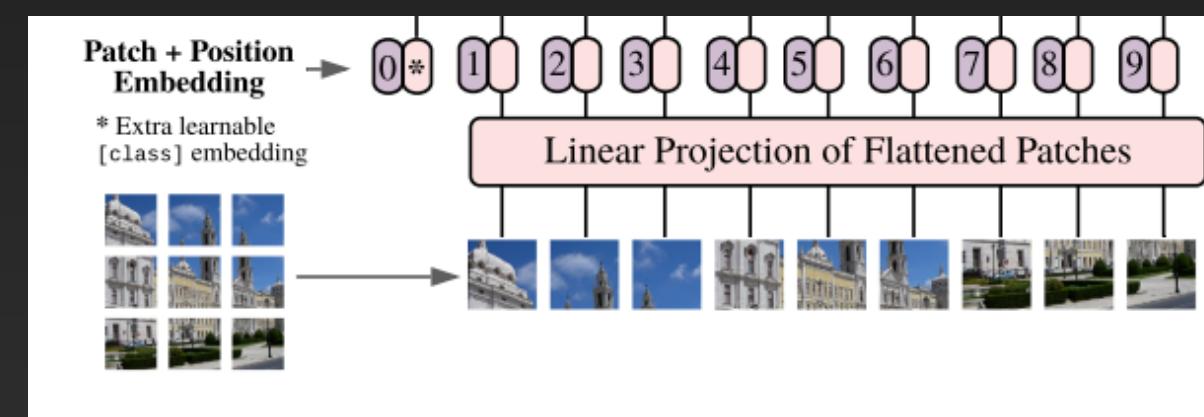
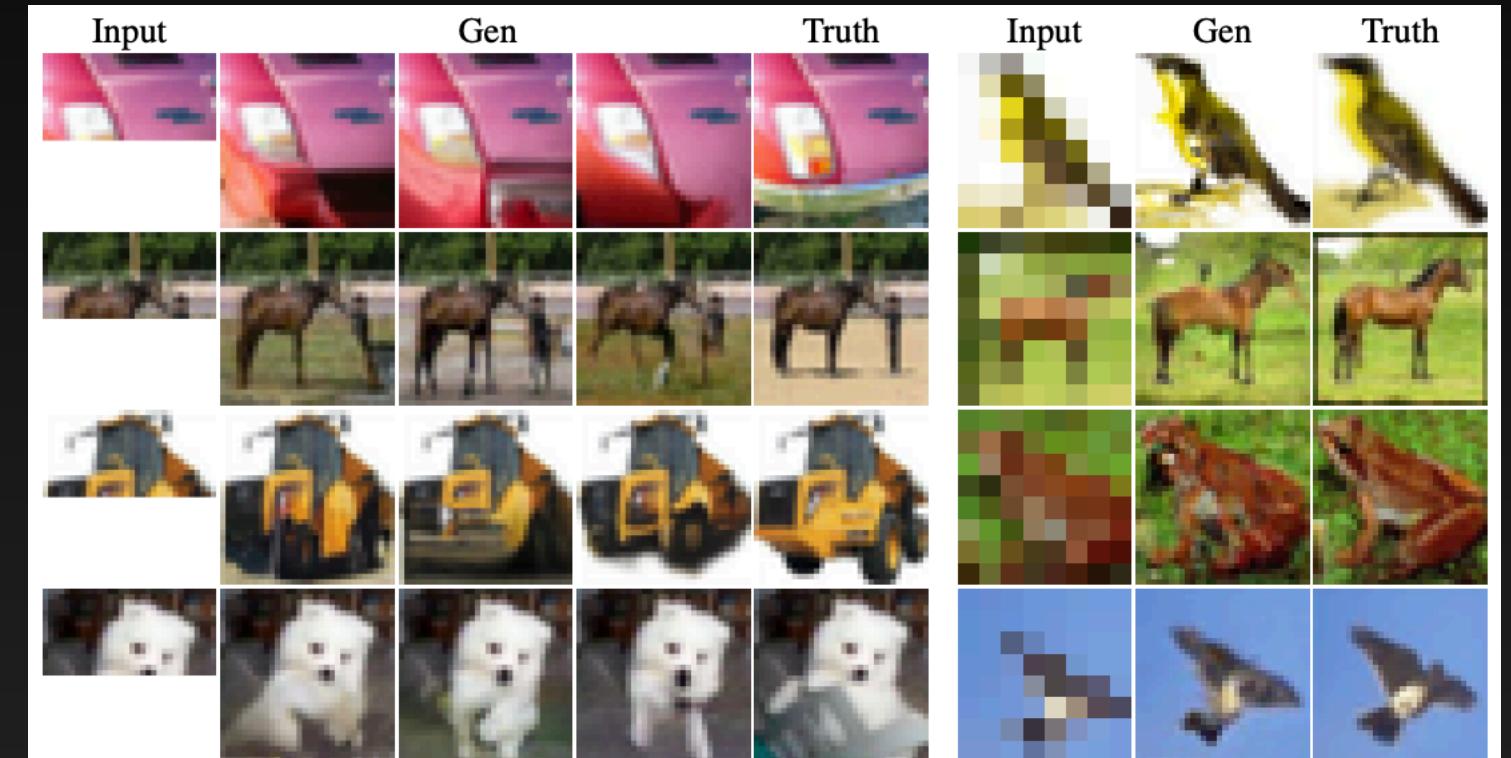
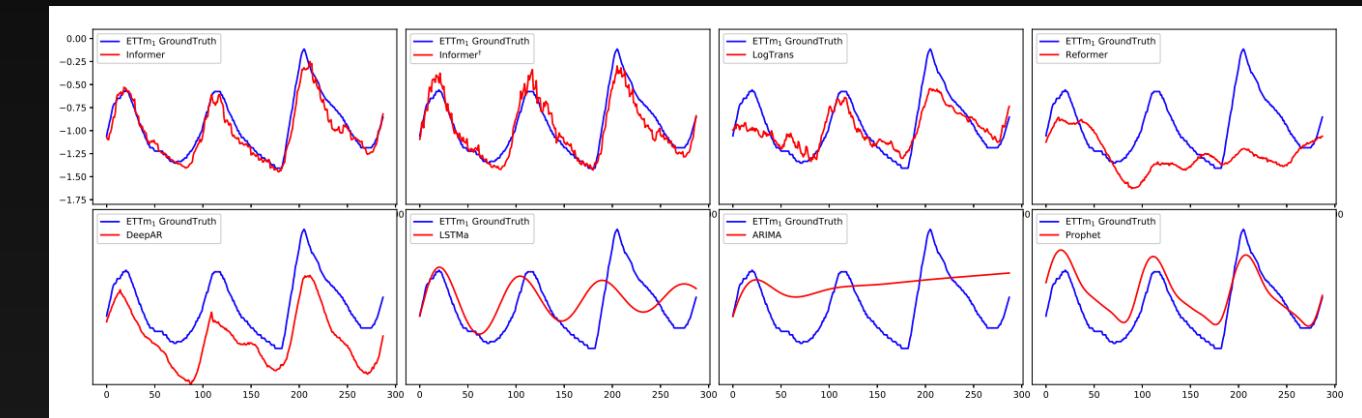
Image from

<https://quizlet.com/539807260/chapter-4-sensing-and-perceiving-our-world-flash-cards/>

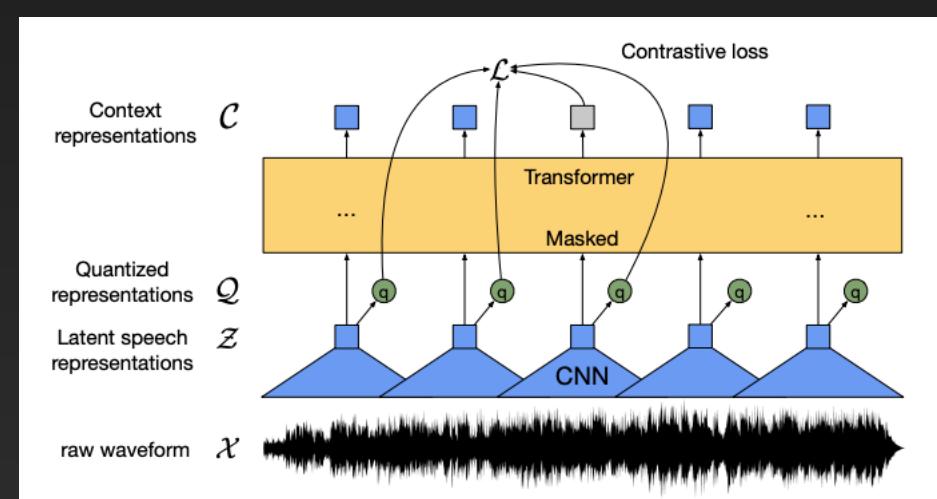
Progress towards Multimodality in the Transformer World

- One Architecture for many things:

- T5¹: Language
- Image Transformer² / ViT³: Image (&Video)
- Wave2Vec 2.0⁴: Audio
- Informer⁵: Long Range Timeseries (Power Grid)
- ECG/Stress Transformer⁶



Modality Specific
Customisations or
Training



1. Raffel, C., Shazeer, N.M., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P.J. (2020). "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". ArXiv, abs/1910.10683.
2. Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N.M., Ku, A., & Tran, D. (2018). "Image Transformer". ArXiv, abs/1802.05751.
3. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houlsby, N. (2021). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". ArXiv, abs/2010.11929.
4. Baevski, A., Zhou, H., Mohamed, A., & Auli, M. (2020). "wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations". ArXiv, abs/2006.11477.
5. Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., & Zhang, W. (2021). "Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting". AAAI.
6. Behnaein, B., Bhatti, A., Rodenburg, D., Hungler, P.C., & Etemad, A. (2021). "A Transformer Architecture for Stress Detection from ECG". 2021 International Symposium on Wearable Computers.

Progress towards Multimodality in the Transformer World

- DALL-E¹, CLIP², UC2³, MUM
(based on T5⁴+MoE⁵), Wu Dao
2.0⁶



1. Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, Ilya Sutskever. (2021) "Zero-Shot Text-to-Image Generation", Proceedings of the 38th International Conference on Machine Learning, PMLR 139:8821-8831.
2. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., & Sutskever, I. (2021). "Learning Transferable Visual Models From Natural Language Supervision". ICML.
3. Mingyang Zhou Luwei Zhou Shuhang Wang Yu Cheng Linjie Li Zhou Yu Jingjing Liu. (2021) "UC2: Universal Cross-lingual Cross-modal Vision-and-Language Pre-training" IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2021)
4. Raffel, C., Shazeer, N.M., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P.J. (2020). "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". ArXiv, abs/1910.10683.
5. Fedus, W., Zoph, B., & Shazeer, N.M. (2021). "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity". ArXiv, abs/2101.03961.
6. He, J., Qiu, J., Zeng, A., Yang, Z., Zhai, J., & Tang, J. (2021). "FastMoE: A Fast Mixture-of-Expert Training System". ArXiv, abs/2103.13262.

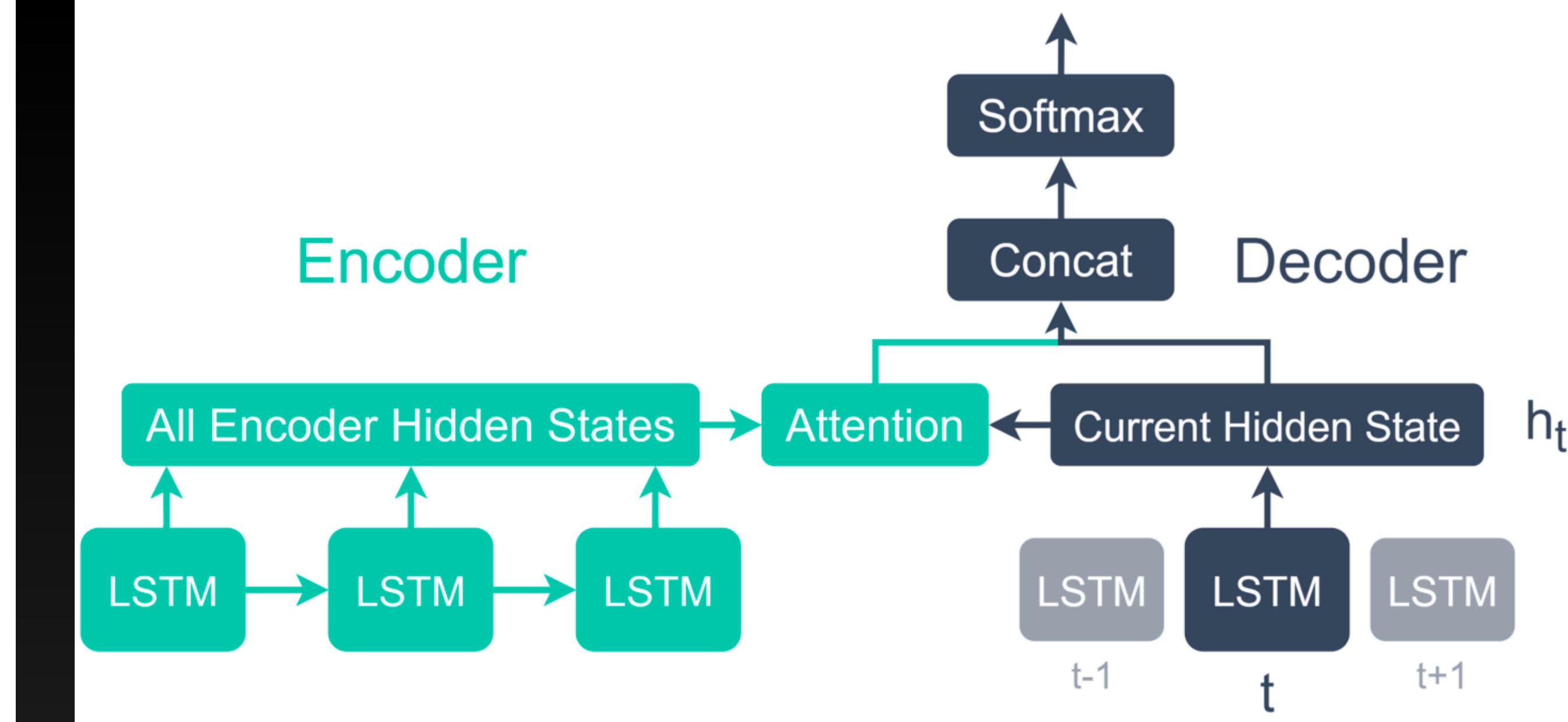
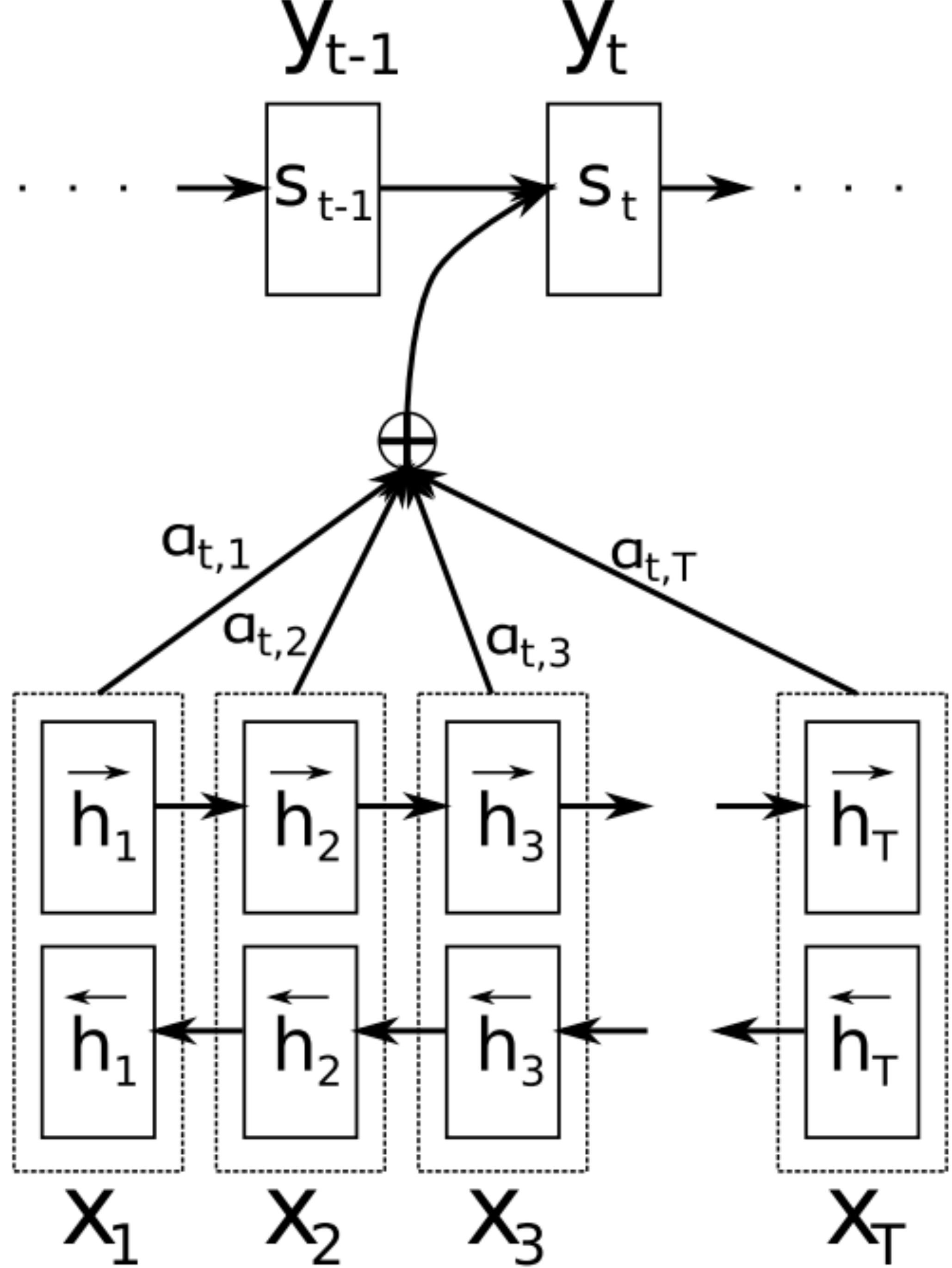
Progress towards Multimodality

- Other options like CNN per modality and then some top level network or mechanism
- Majority Voting, Boosting and Bagging and other Ensemble Methods exist
- These are all late fusions of Modalities, something the brain most probably does not do

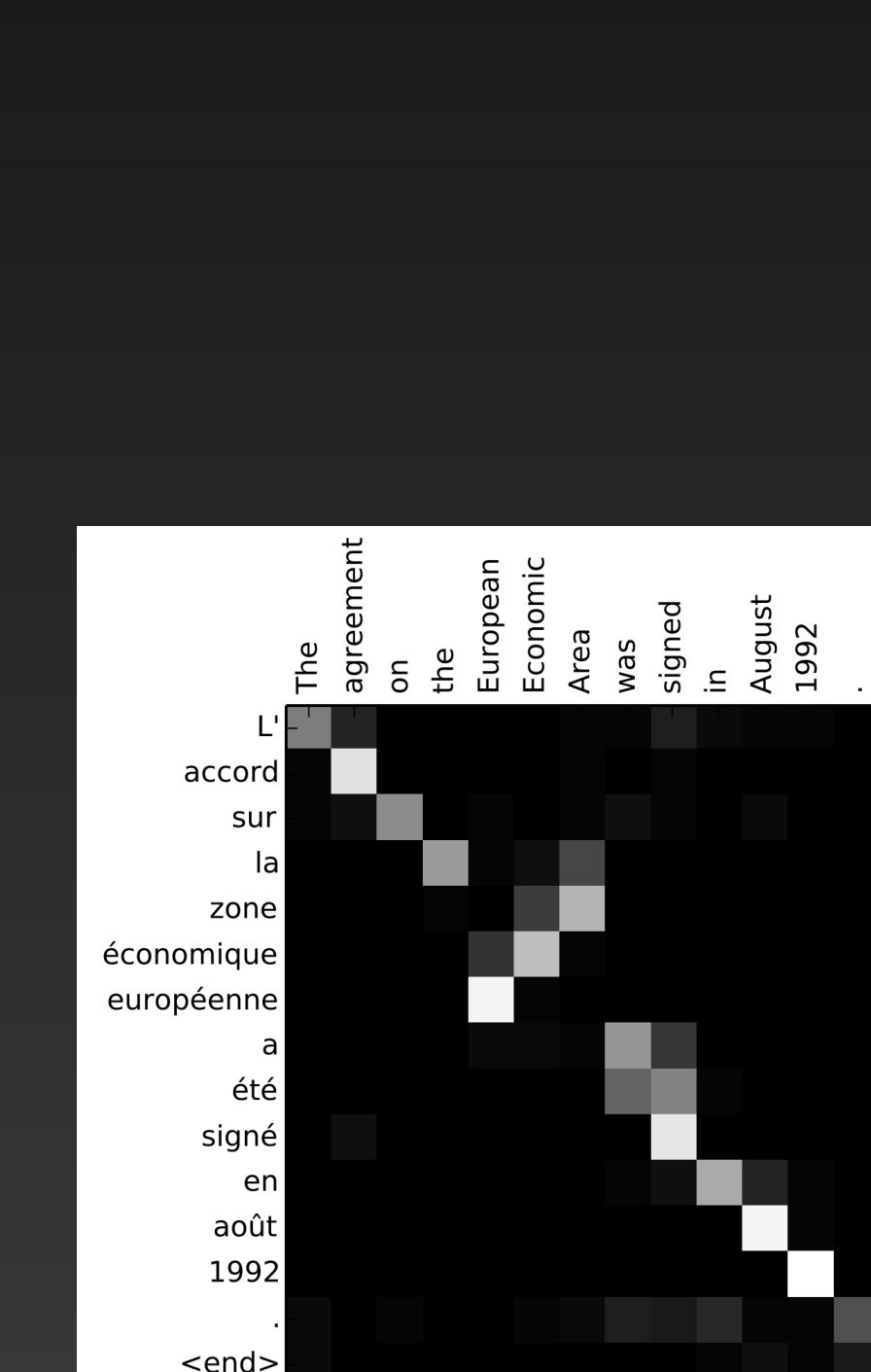
Problem with Transformers for Multimodality

- They either are still tuned towards a modality (like ViT, Wave2Vec...)
- The multimodal examples above have **extreme Hardware** requirements. You basically need an unlimited hardware budget.
- Why: Standard Multi-Head Attention has $O(n^2 \times d)$ time and memory complexity **per layer** with n being the sequence length
- Means that a 512x512 picture (not a video) and a hidden dimension of 512 and 10 layers needs at least **3.52e+14 operations** to classify
- Modern GPUs: 15 Tera Flops per second = **1.5e+13/s**
- You wait about **10s** for **one picture** to be processed on a kick ass GPU

Attention



from: <https://towardsdatascience.com/how-transformers-work-6cb4629506df>



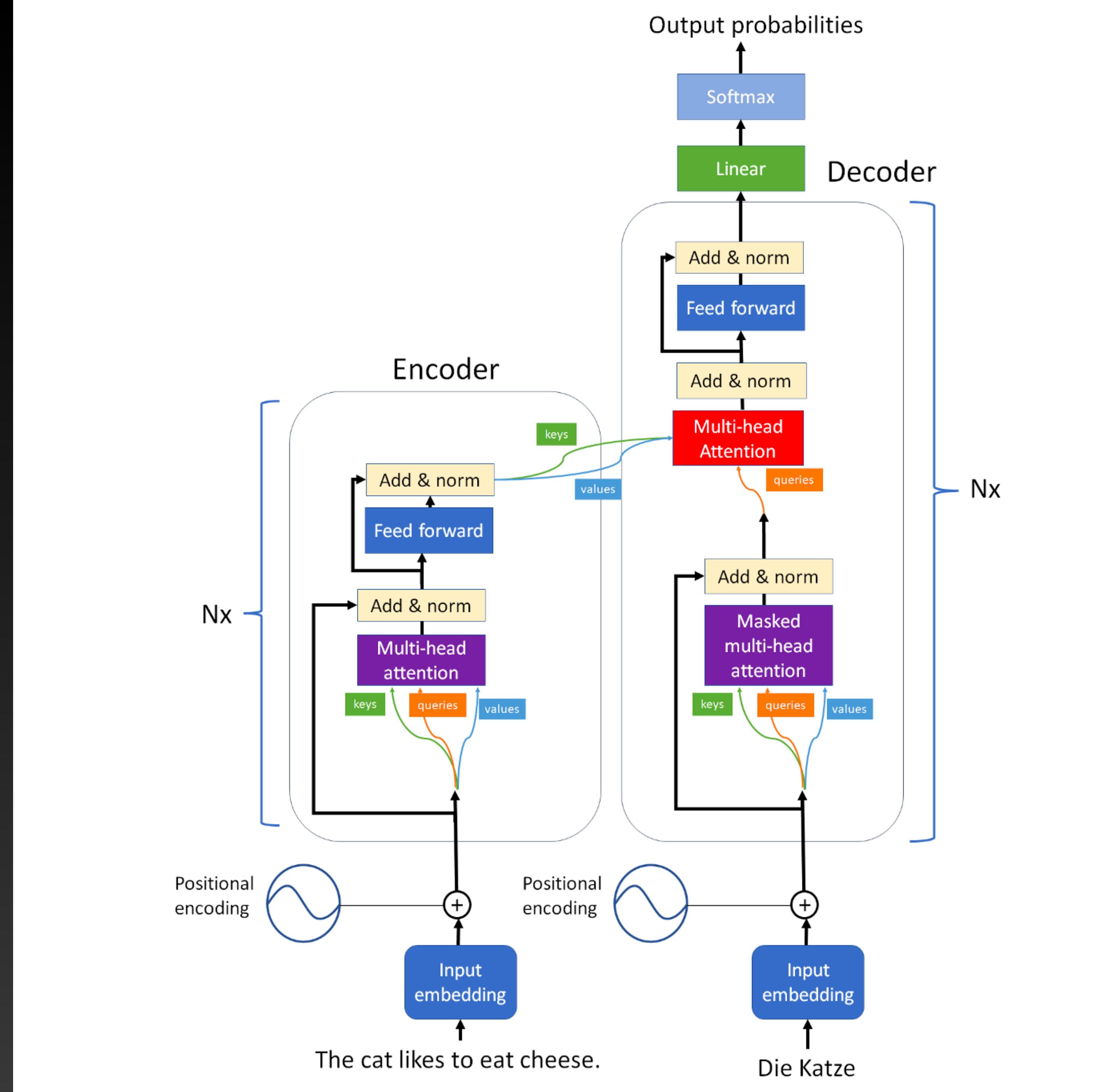
$$e_{ij} = a(s_{i-1}, h_j) \quad (1)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (2)$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (3)$$

What is the main Problem?

Sequential Training: sequentially encode source (n^2), then compute attention for each output (m)
NO PARALLEL TRAINING!



$$\mathbf{X} \times \mathbf{WQ} = \mathbf{Q}$$

The diagram illustrates matrix multiplication. On the left, a green matrix \mathbf{X} is shown as a 2x4 grid of squares. In the center, a multiplication sign (\times) is positioned between the green matrix and a purple matrix \mathbf{WQ} . The purple matrix \mathbf{WQ} is depicted as a 4x3 grid of squares. To the right of the multiplication sign is an equals sign (=). To the right of the equals sign is a purple matrix \mathbf{Q} , which is also a 2x3 grid of squares.

$$\begin{matrix} X \\ \times \\ \begin{matrix} W \\ K \end{matrix} \end{matrix} = \begin{matrix} K \end{matrix}$$

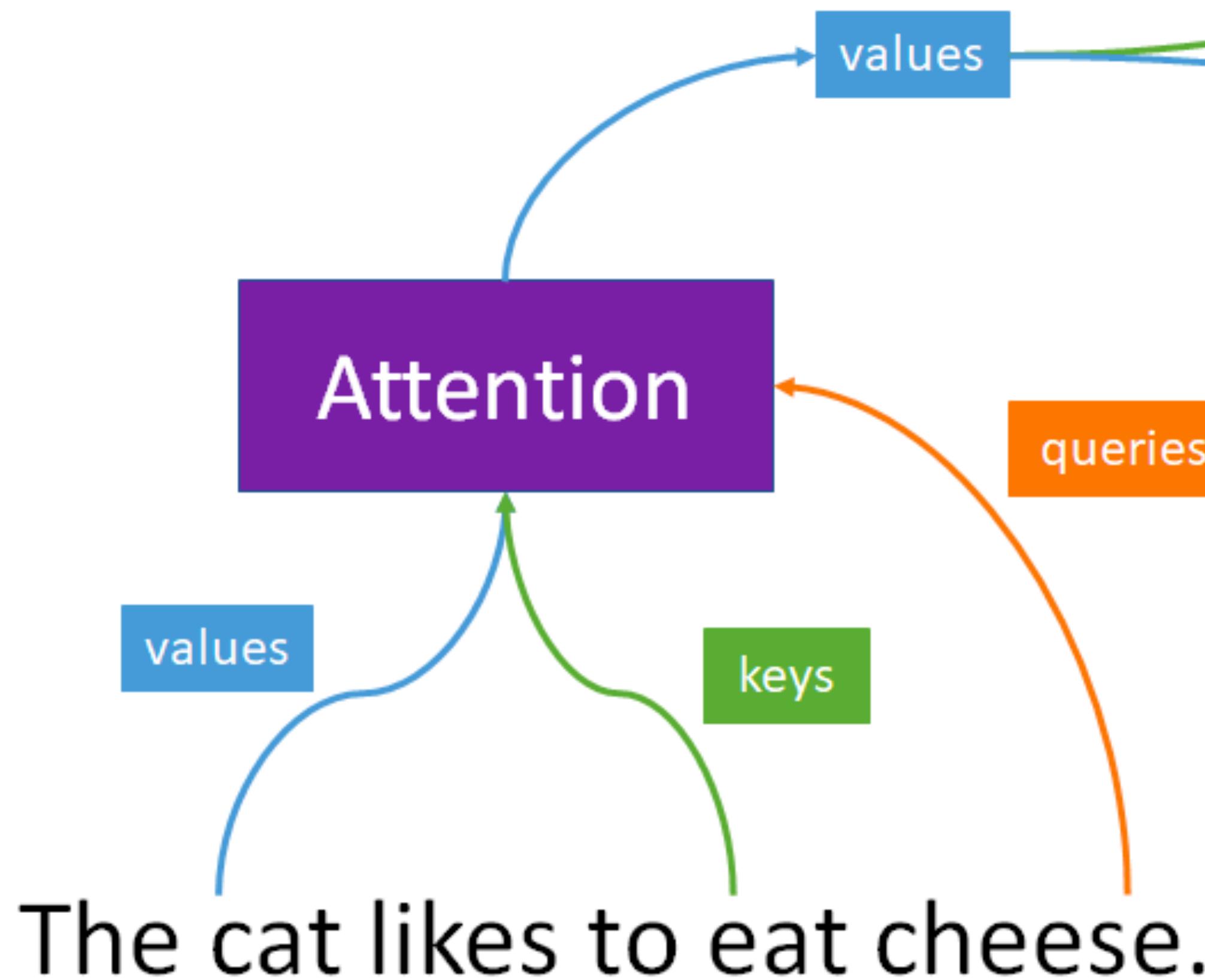
The diagram illustrates matrix multiplication. On the left, a green matrix labeled X is shown as a 2x4 grid. In the center, a multiplication sign (\times) is positioned between X and another matrix labeled W and K , which is represented as a 4x3 grid with orange borders. To the right of the multiplication sign is an equals sign (=), followed by a result matrix labeled K , which is a 2x3 grid filled with orange squares.

$$\begin{matrix} \mathbf{X} \\ \times \\ \begin{matrix} \mathbf{W} \\ \mathbf{V} \end{matrix} \end{matrix} = \mathbf{V}$$

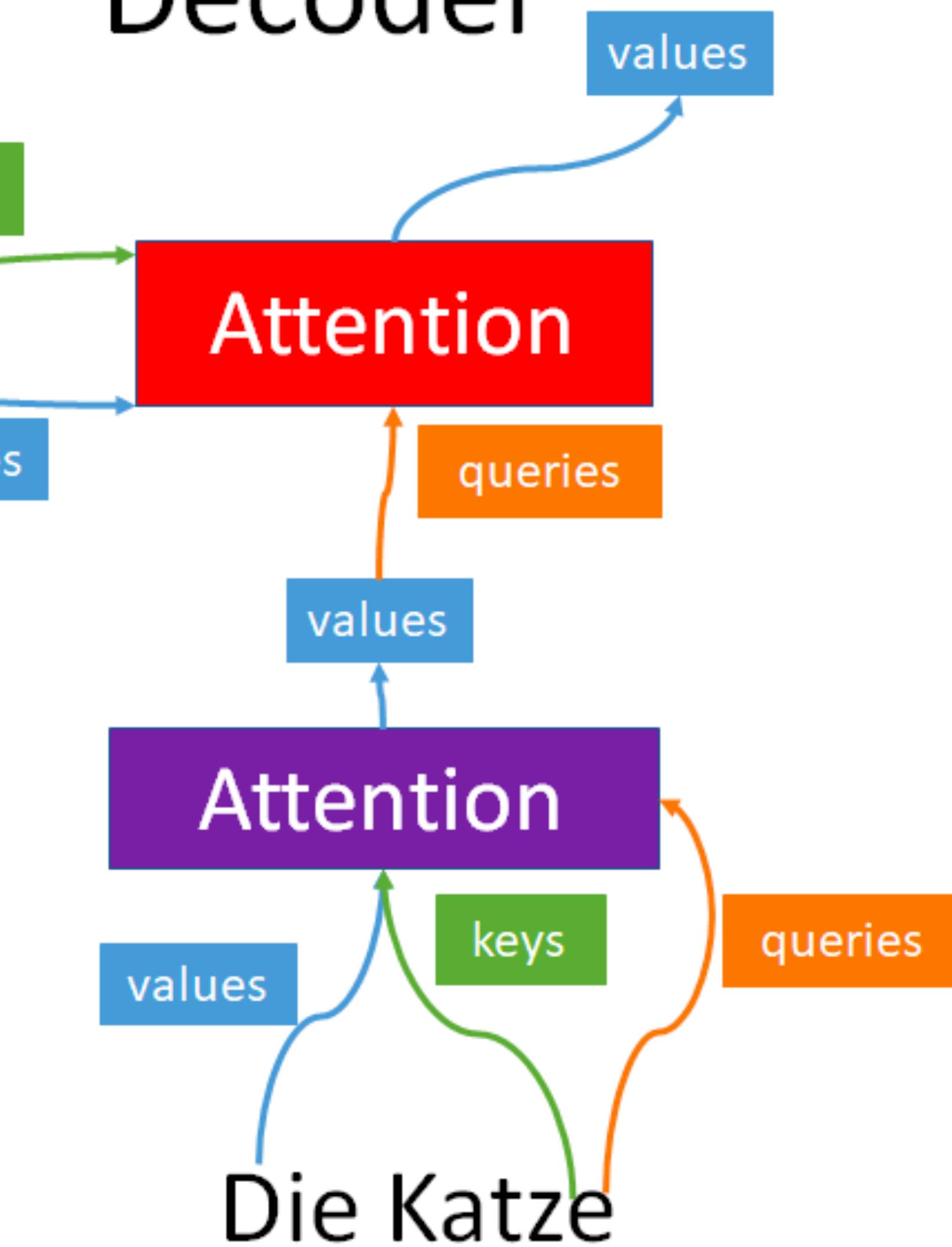

$$z = \text{softmax} \left(\frac{\begin{array}{c} \text{Q} \\ \times \\ \text{K}^T \\ \hline \text{V} \end{array}}{\sqrt{d_k}} \right)$$

The diagram illustrates the computation of the attention matrix z . It shows the multiplication of two matrices, Q and K^T , followed by the division by $\sqrt{d_k}$ and the application of the softmax function. The resulting matrix z is shown as a pink 3x3 grid.

Encoder

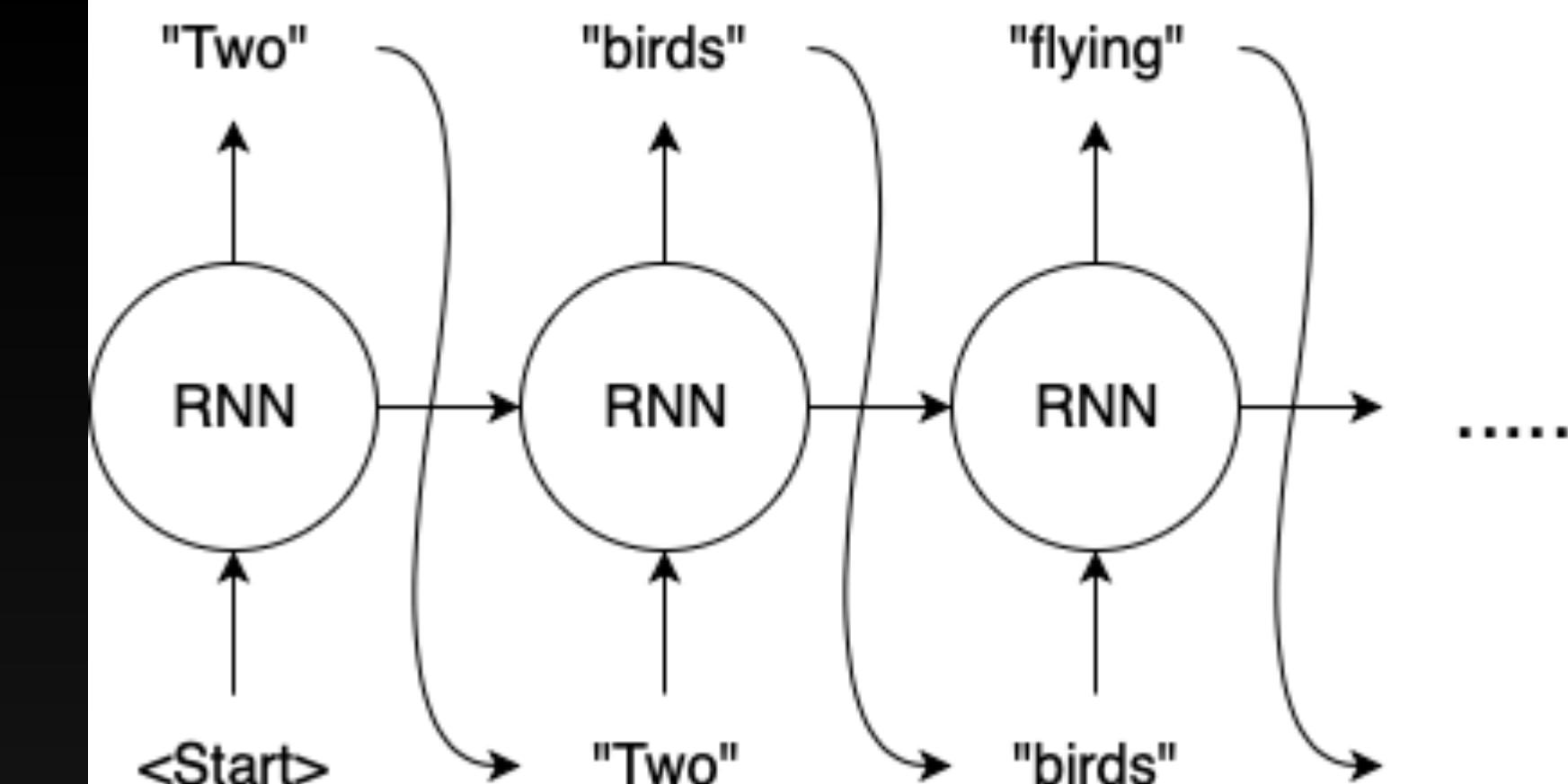


Decoder

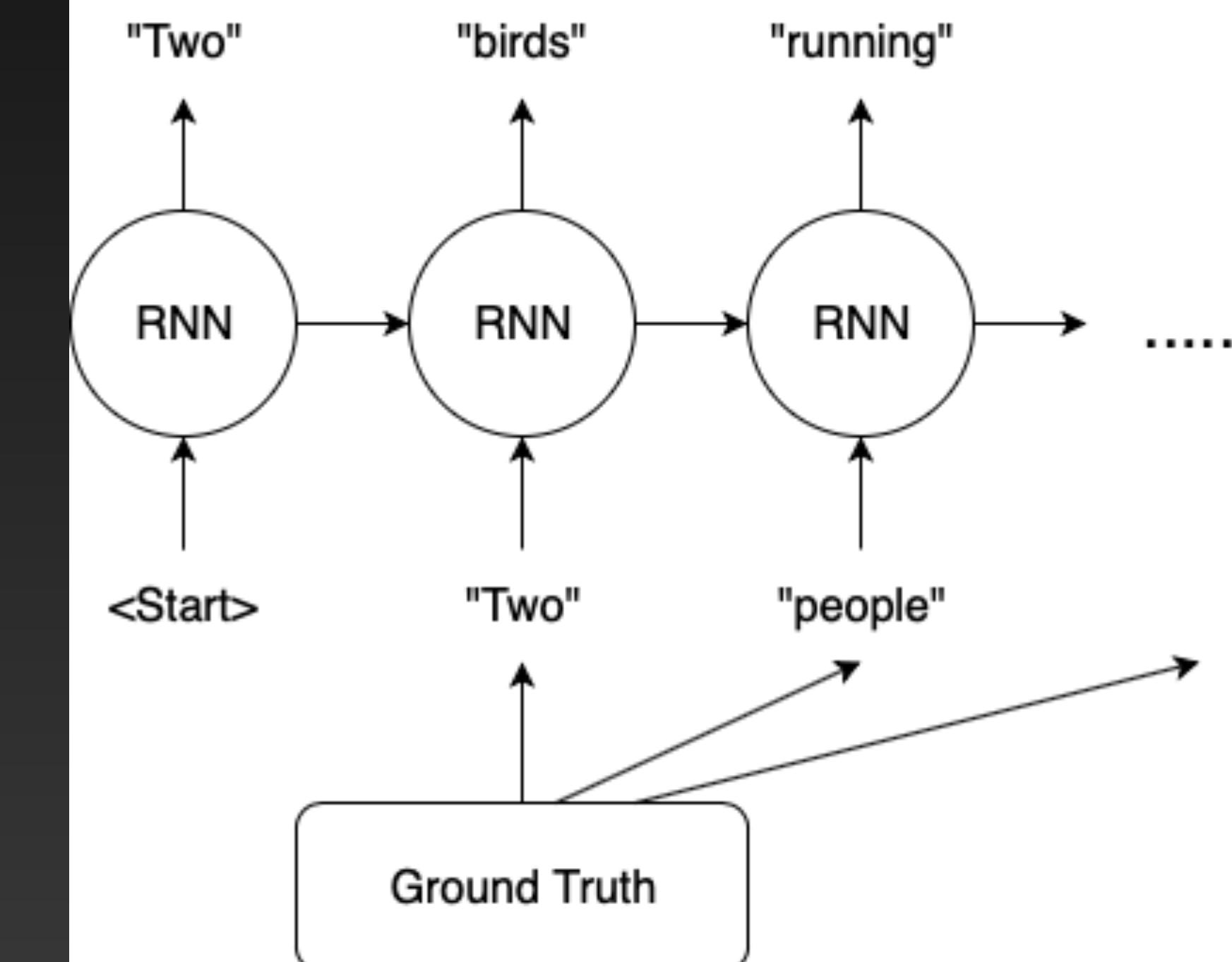


Teacher Forcing

- Originally used to improve training of RNNs to avoid cumulative error
- For Transformer: Also allow to create all training queries at once

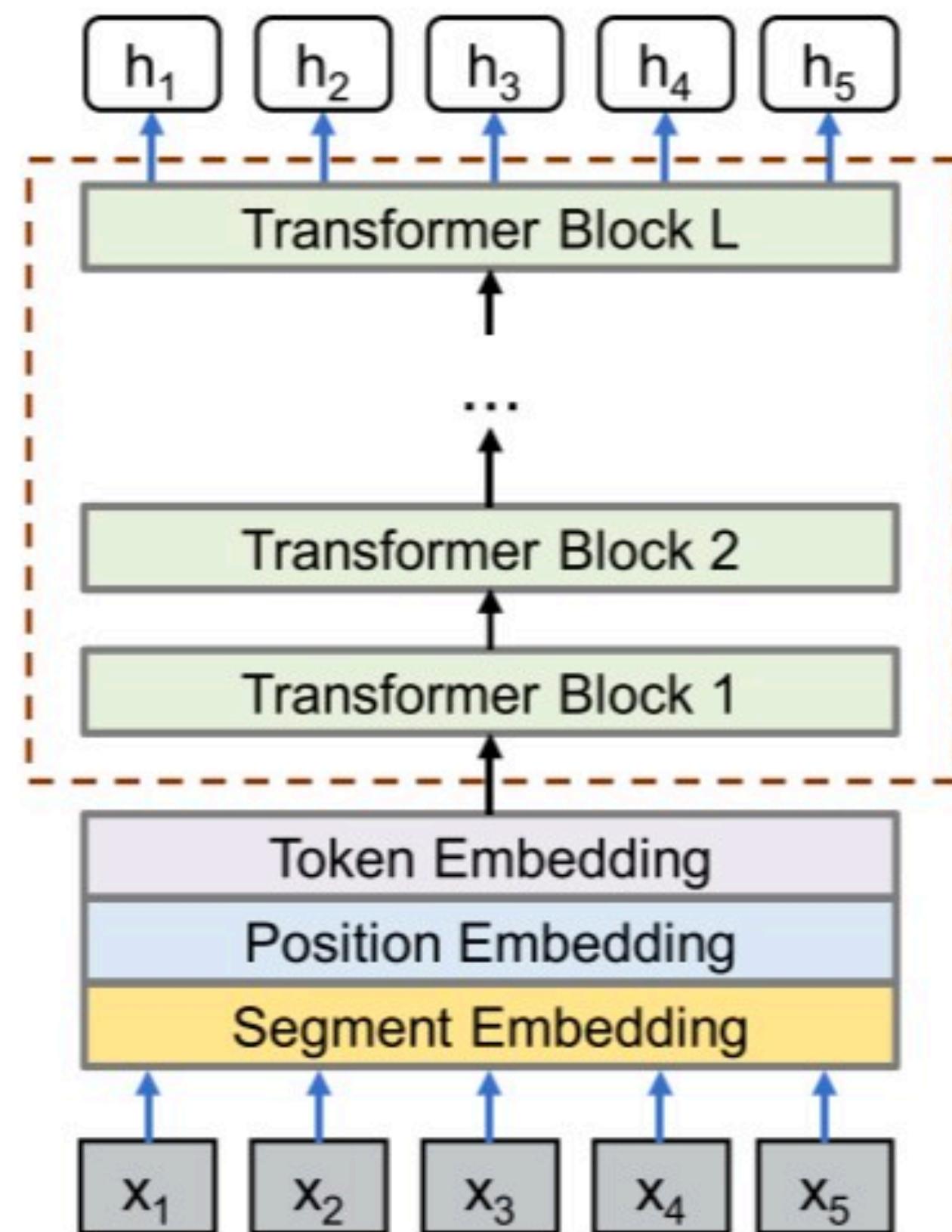


Without Teacher Forcing

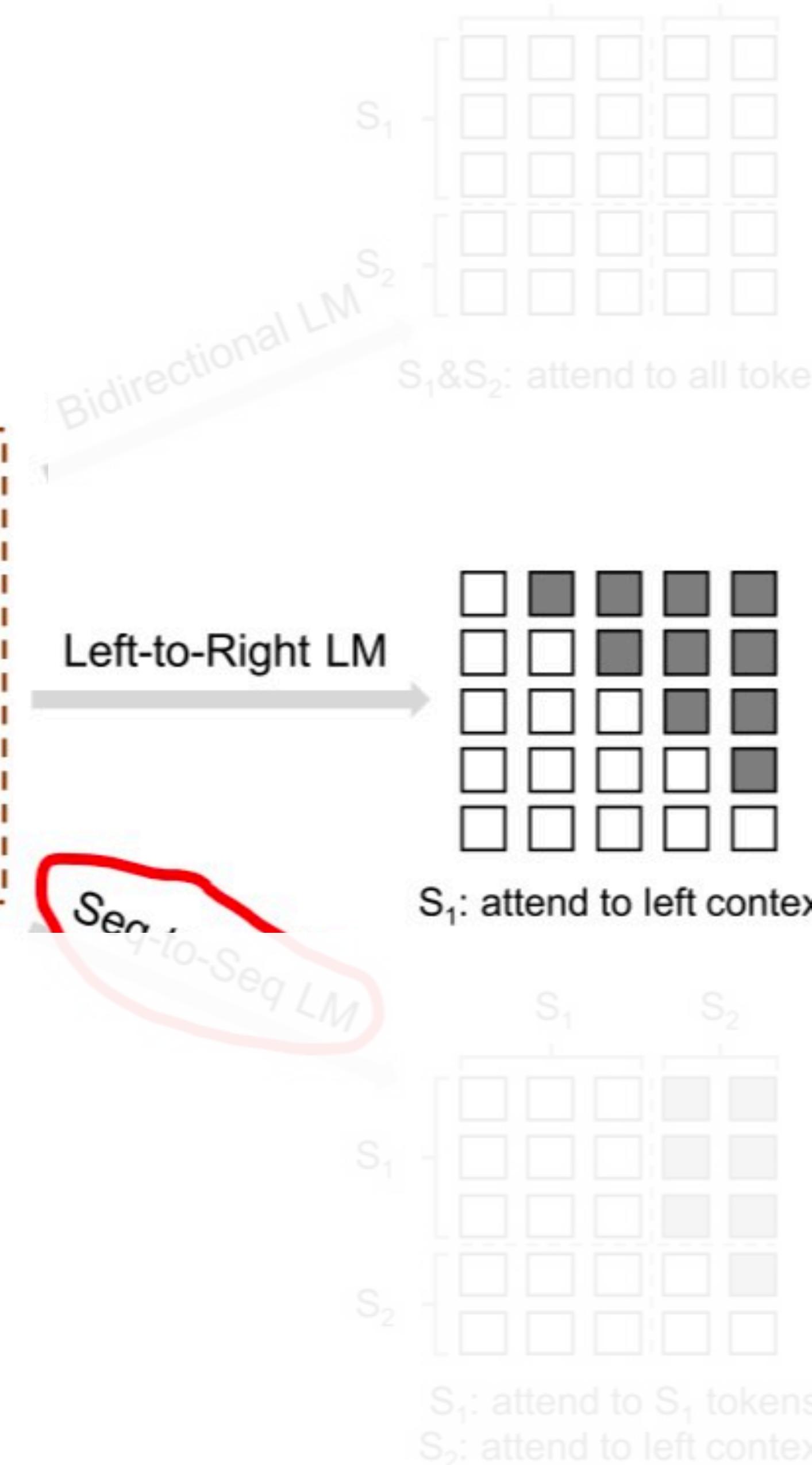


With Teacher Forcing

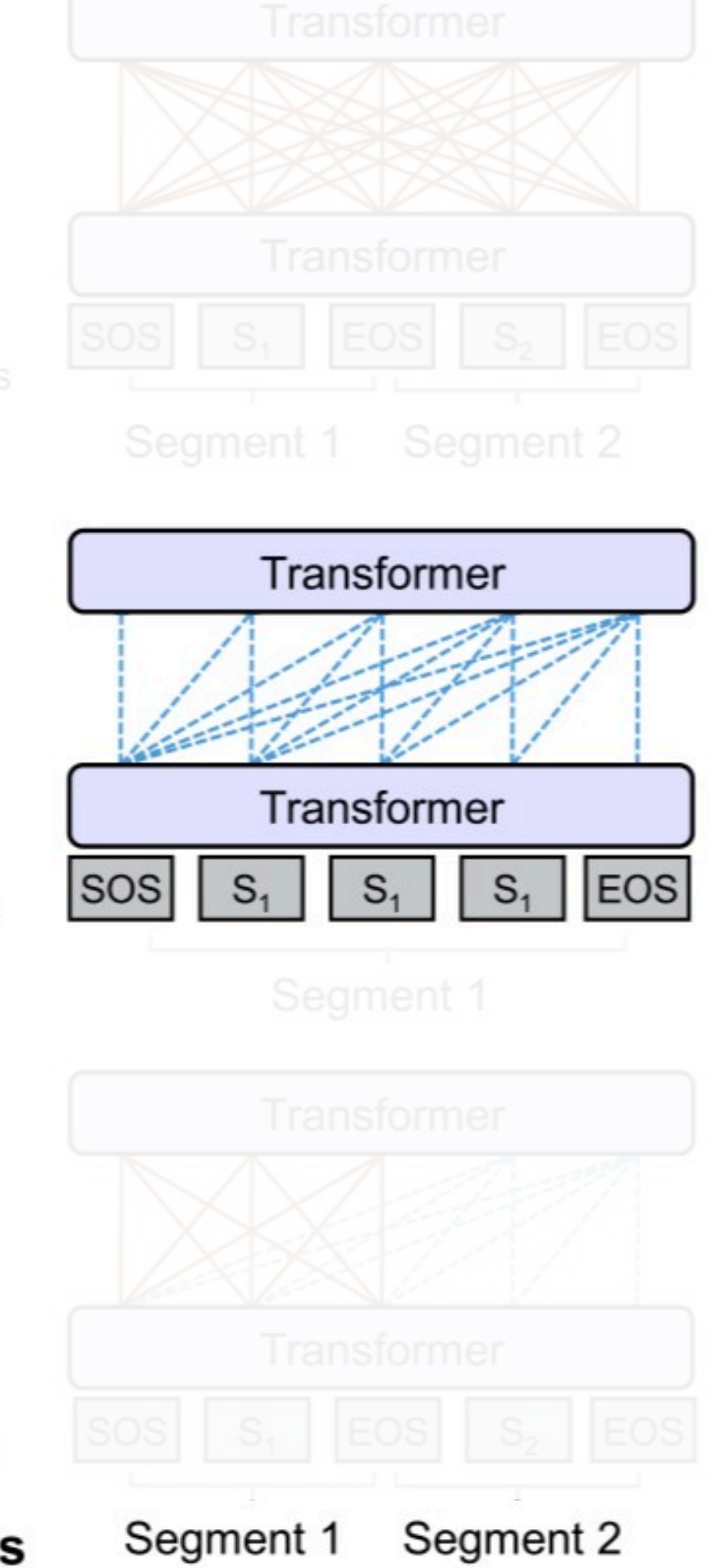
Allow to attend
 Prevent from attending



Unified LM with Shared Parameters

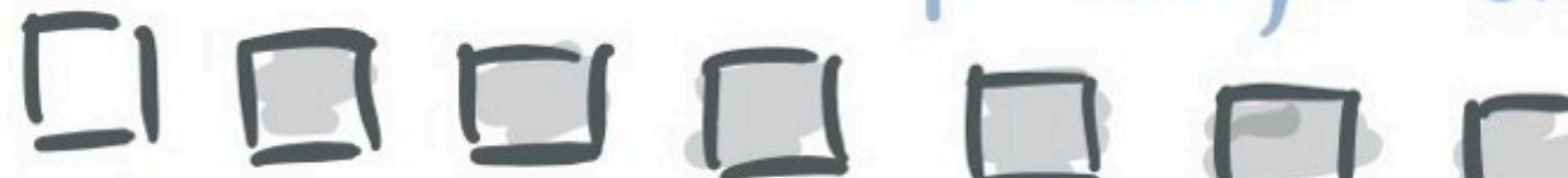


Self-attention Masks



I like to eat pizza for breakfast

I



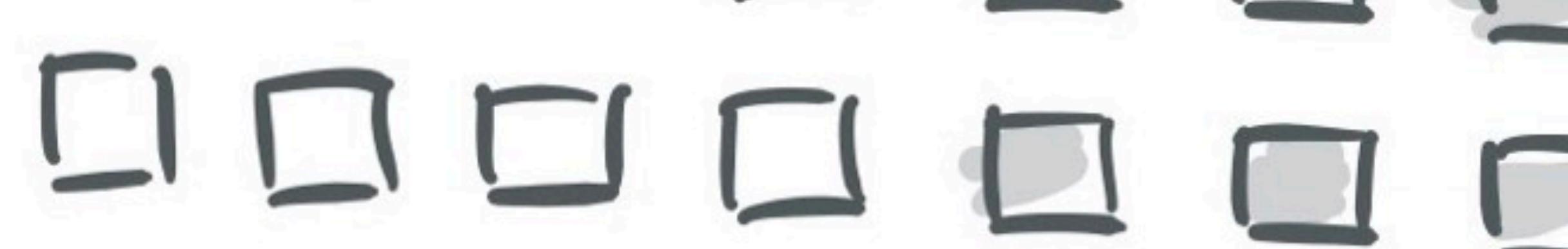
like



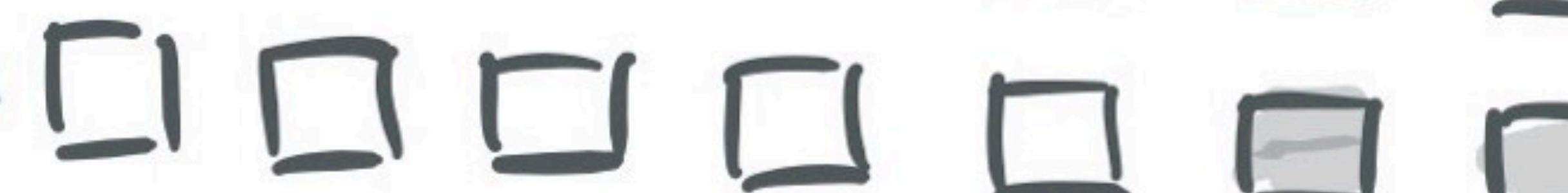
to



eat



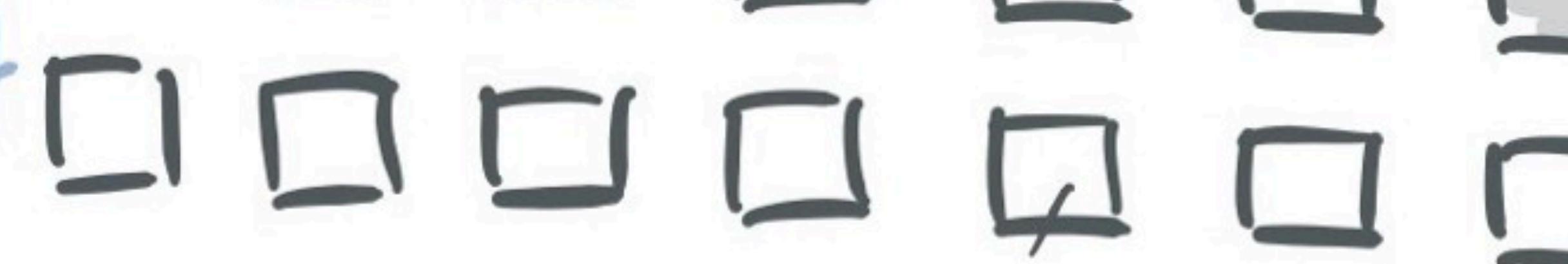
pizza



for



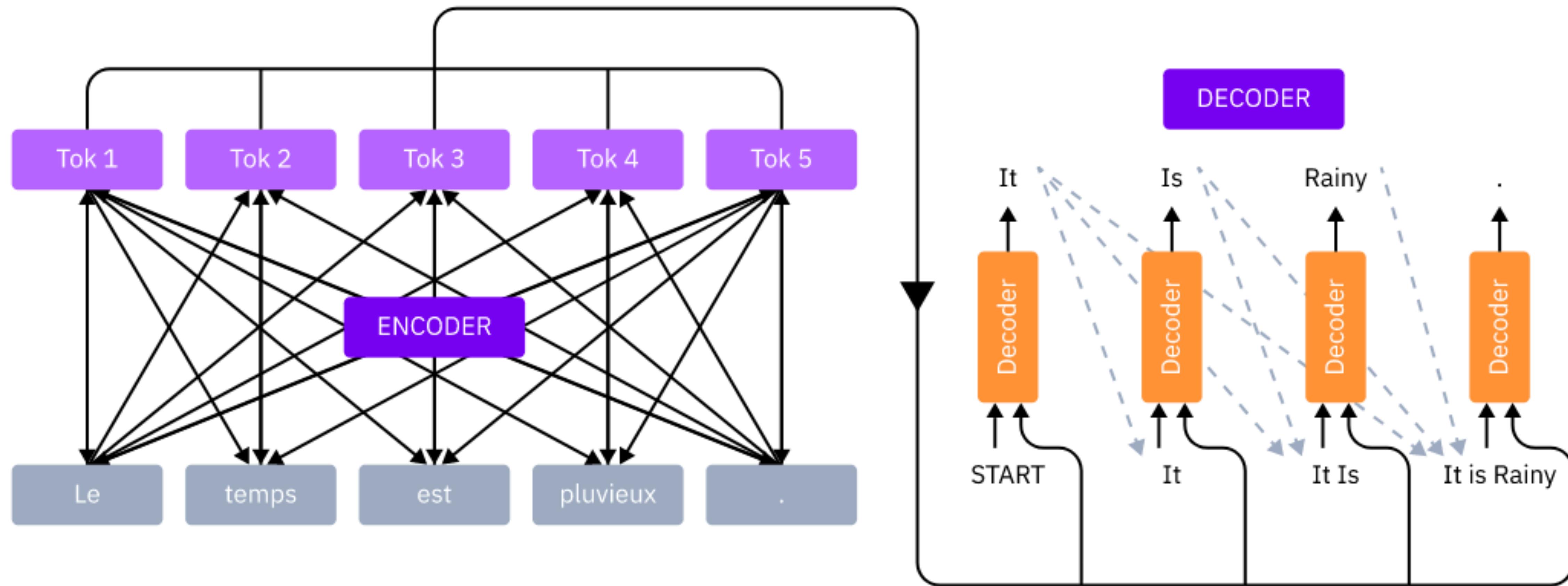
breakfast

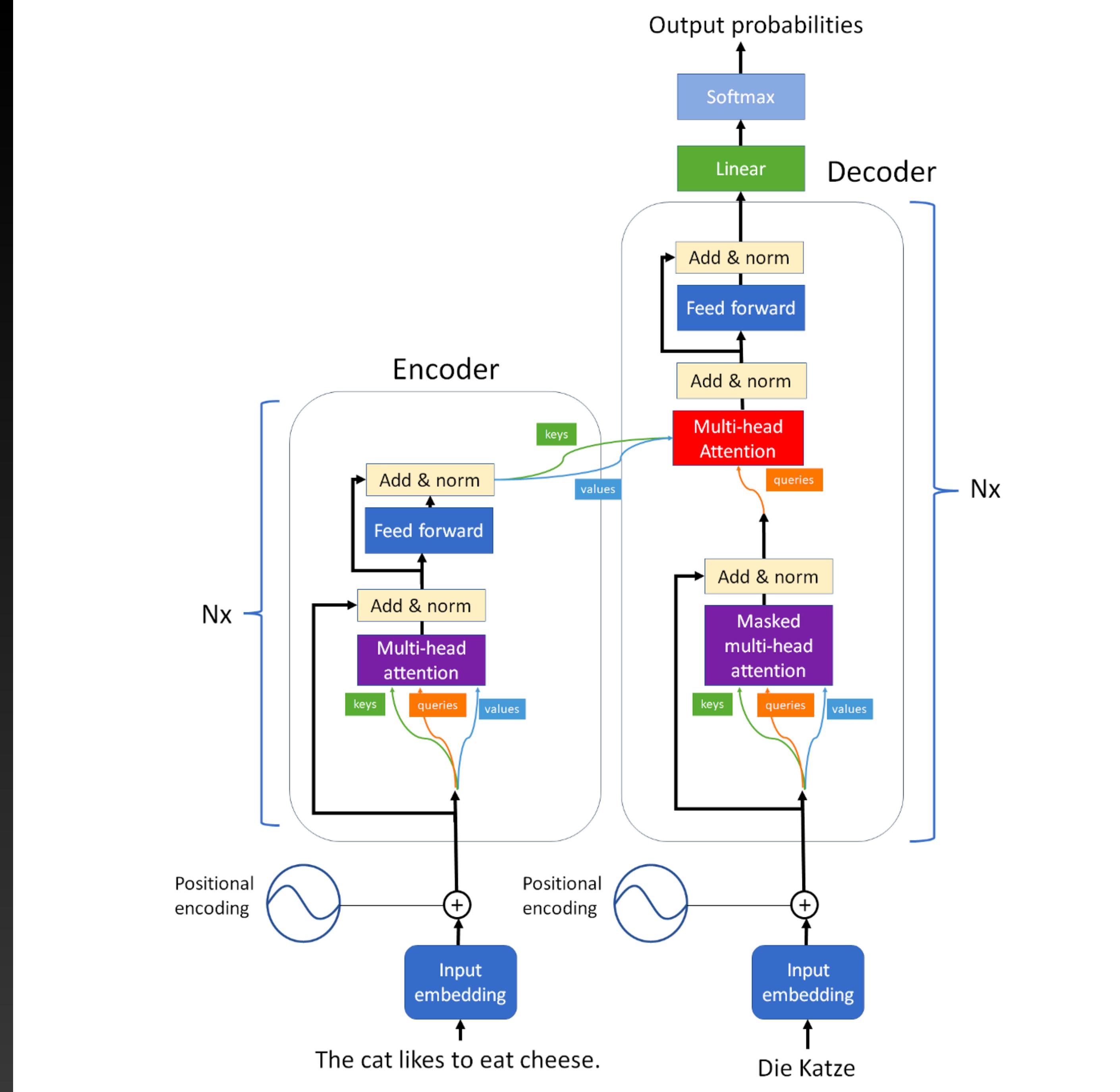


- ∞

Language Model
(left-to-right)

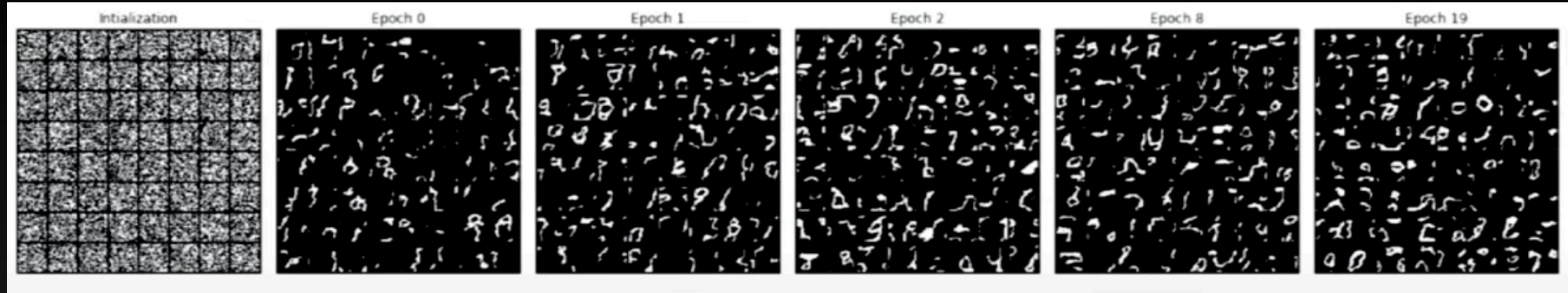
Mash





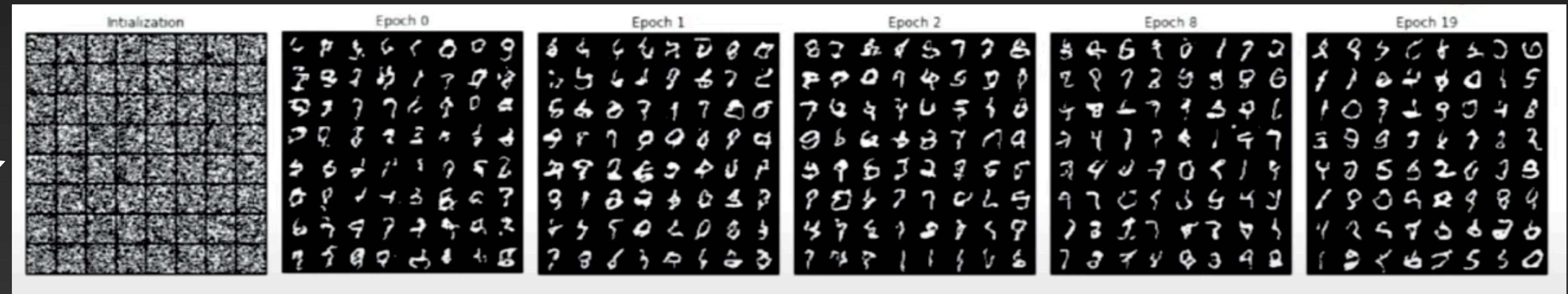
Is there still a problem?

Attention Mechanism actually allow one token to see all others especially in encoder: They loose all notion of position! (It is not a RNN)

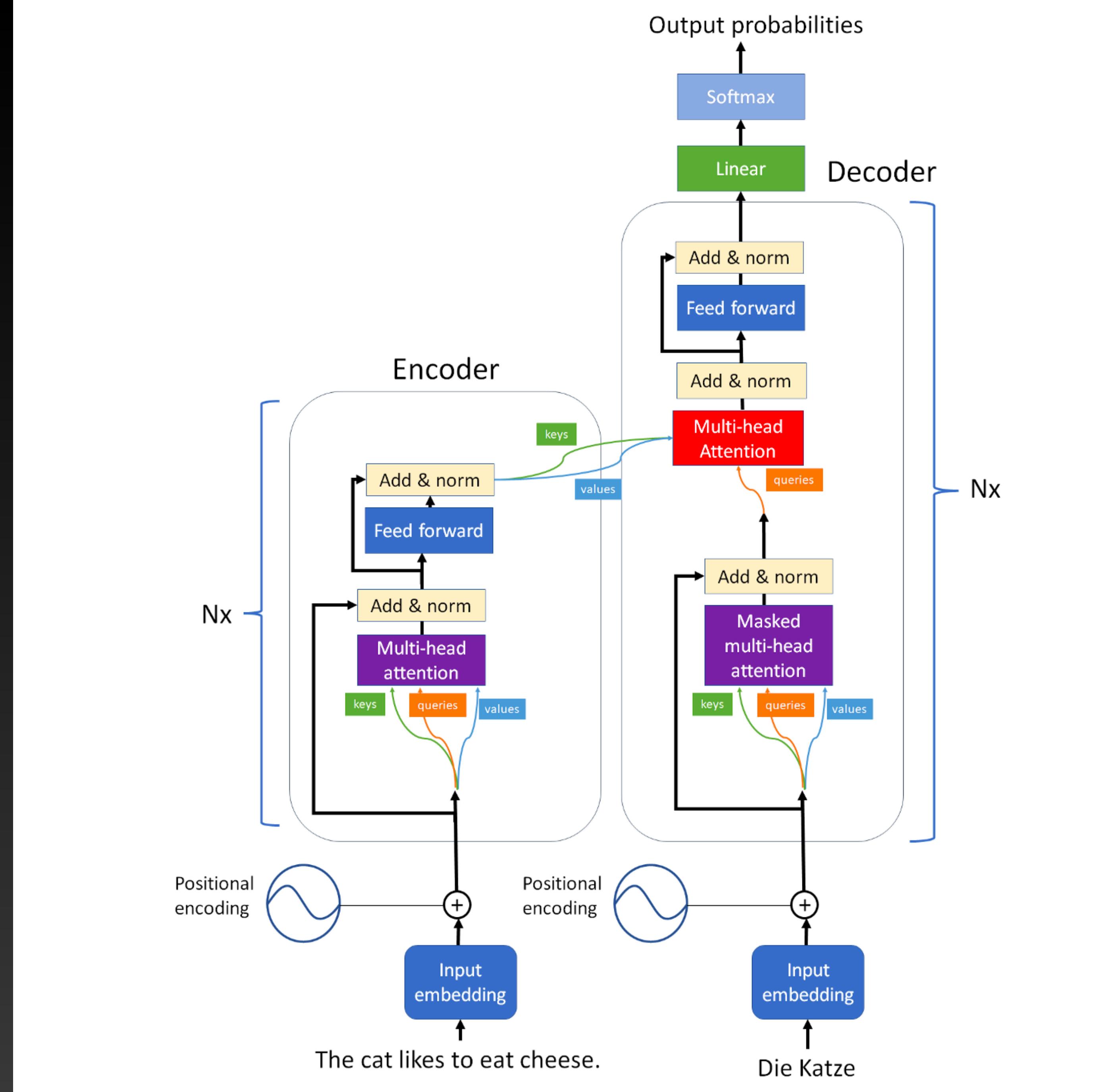


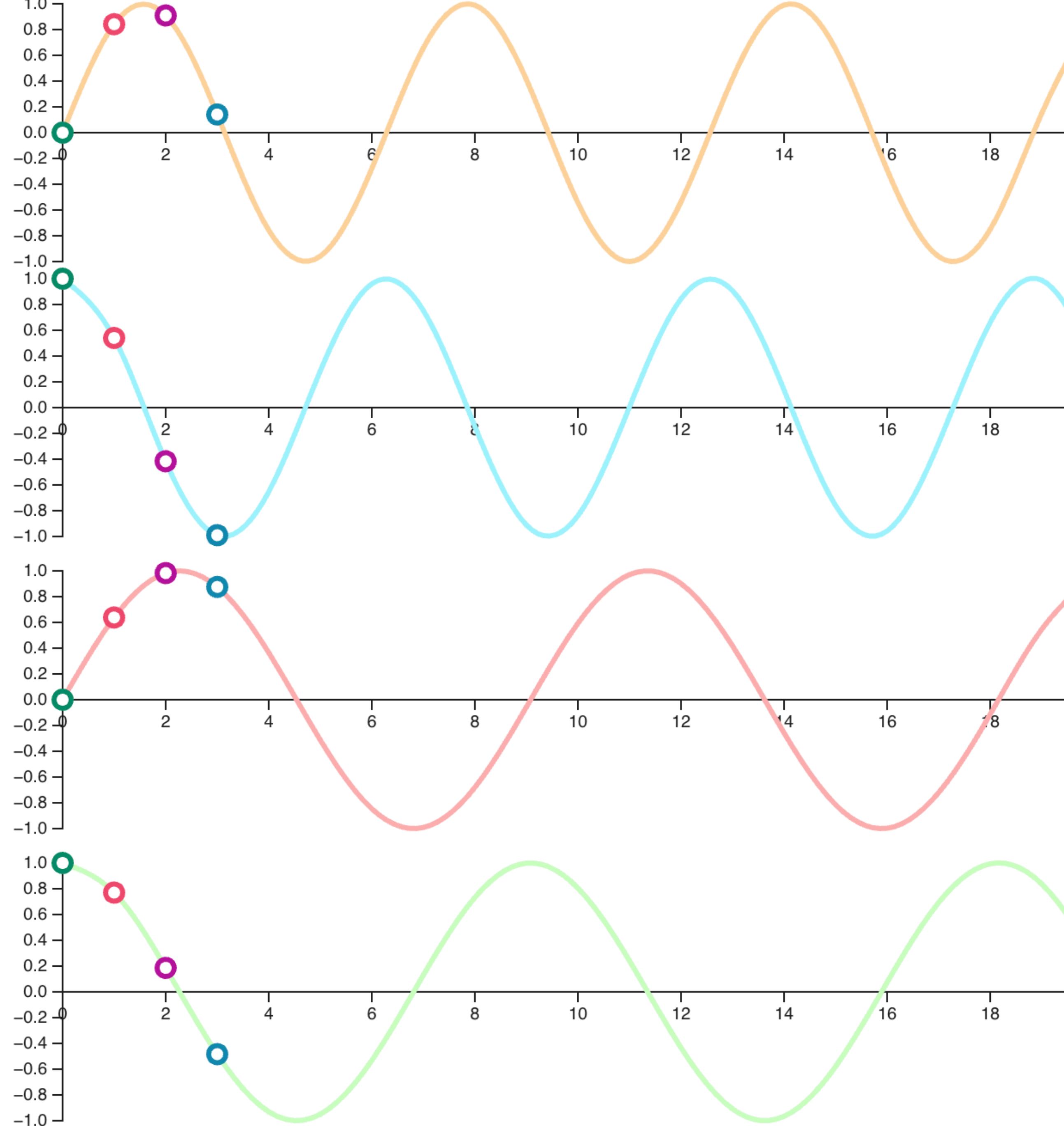
WaveNet

with position
encoding



Here simply append (row, column) num to each pixel value (so add two dimensions)





	p0	p1	p2	p3	i=0
0.000	0.841	0.909	0.141		i=0
1.000	0.540	-0.416	-0.990		i=1
0.000	0.638	0.983	0.875		i=2
1.000	0.770	0.186	-0.484		i=3

Positional Encoding

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

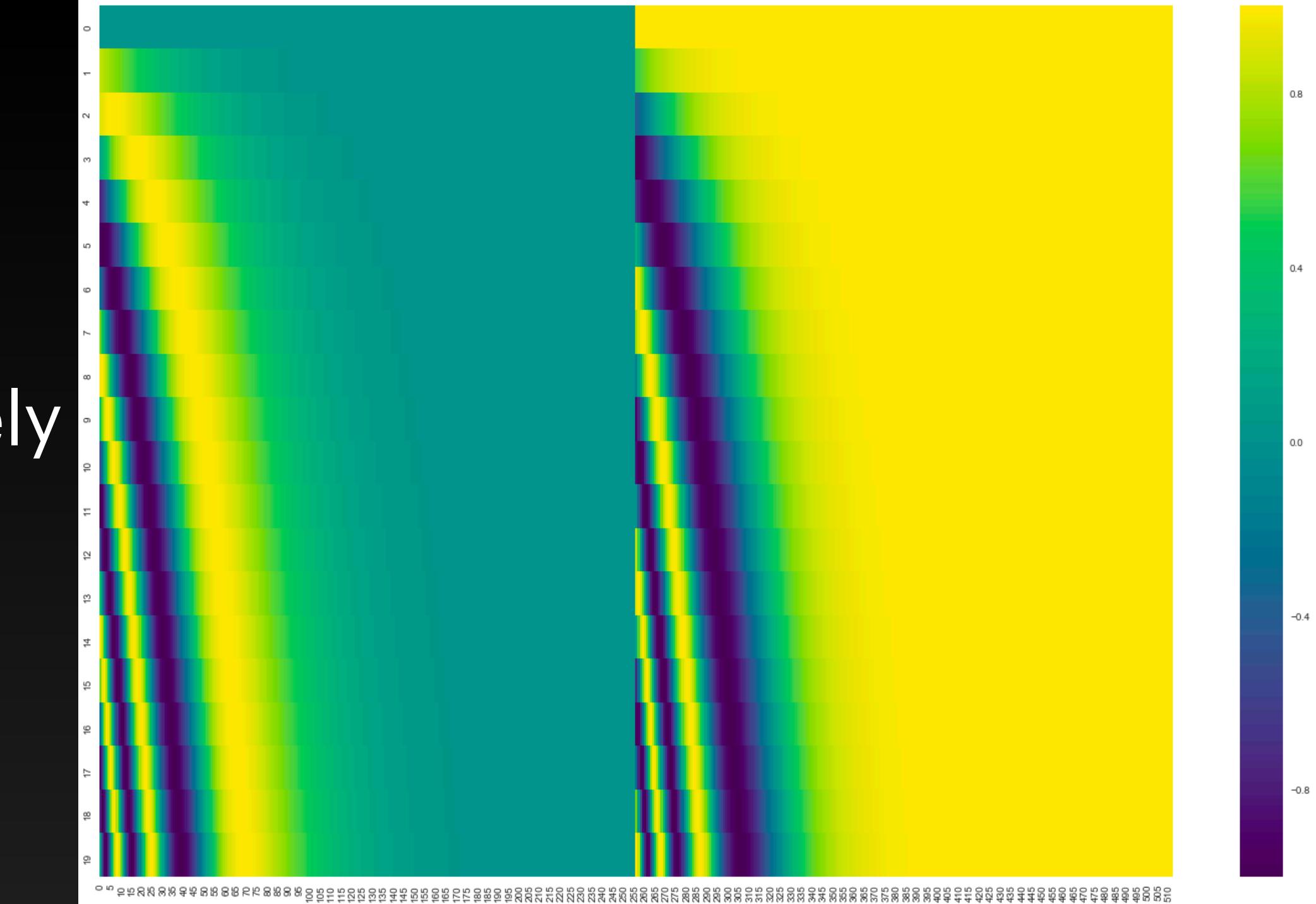
Settings: $d = 50$

The value of each positional encoding depends on the *position* (*pos*) and *dimension* (*d*). We calculate result for every *index* (*i*) to get the whole vector.

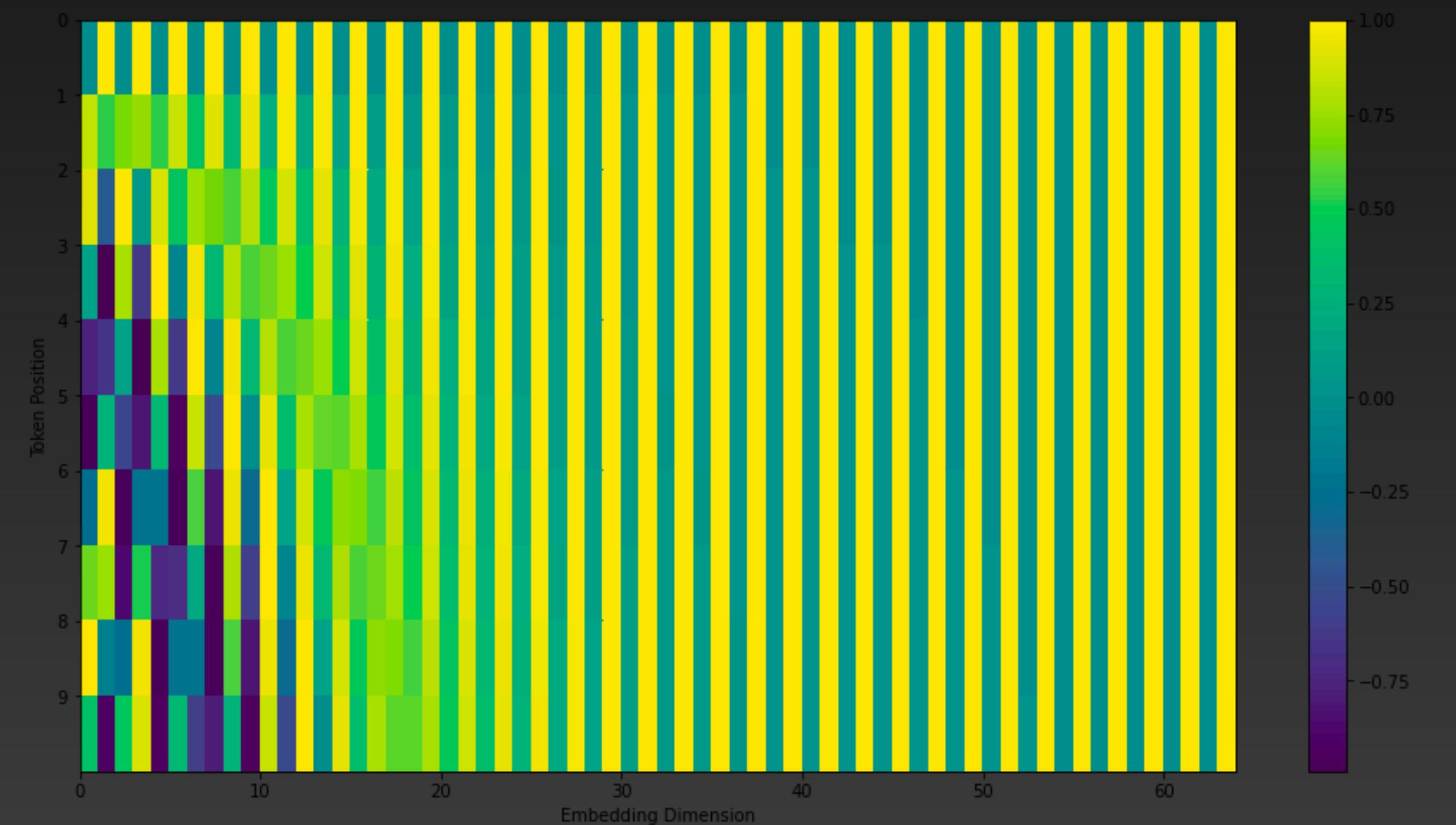
Real Example from Text

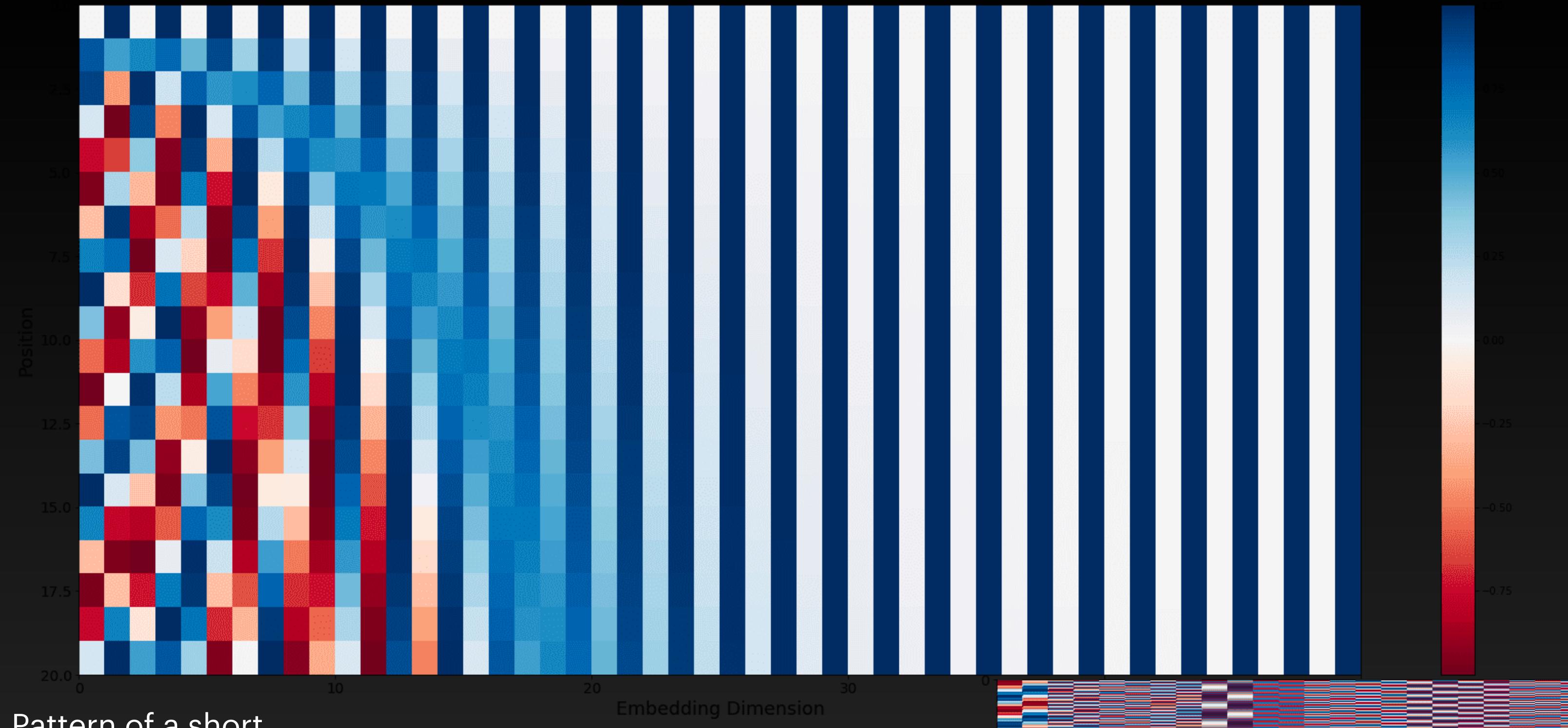
20 words text, 512 embedding
Dimension

sine and cosine separately

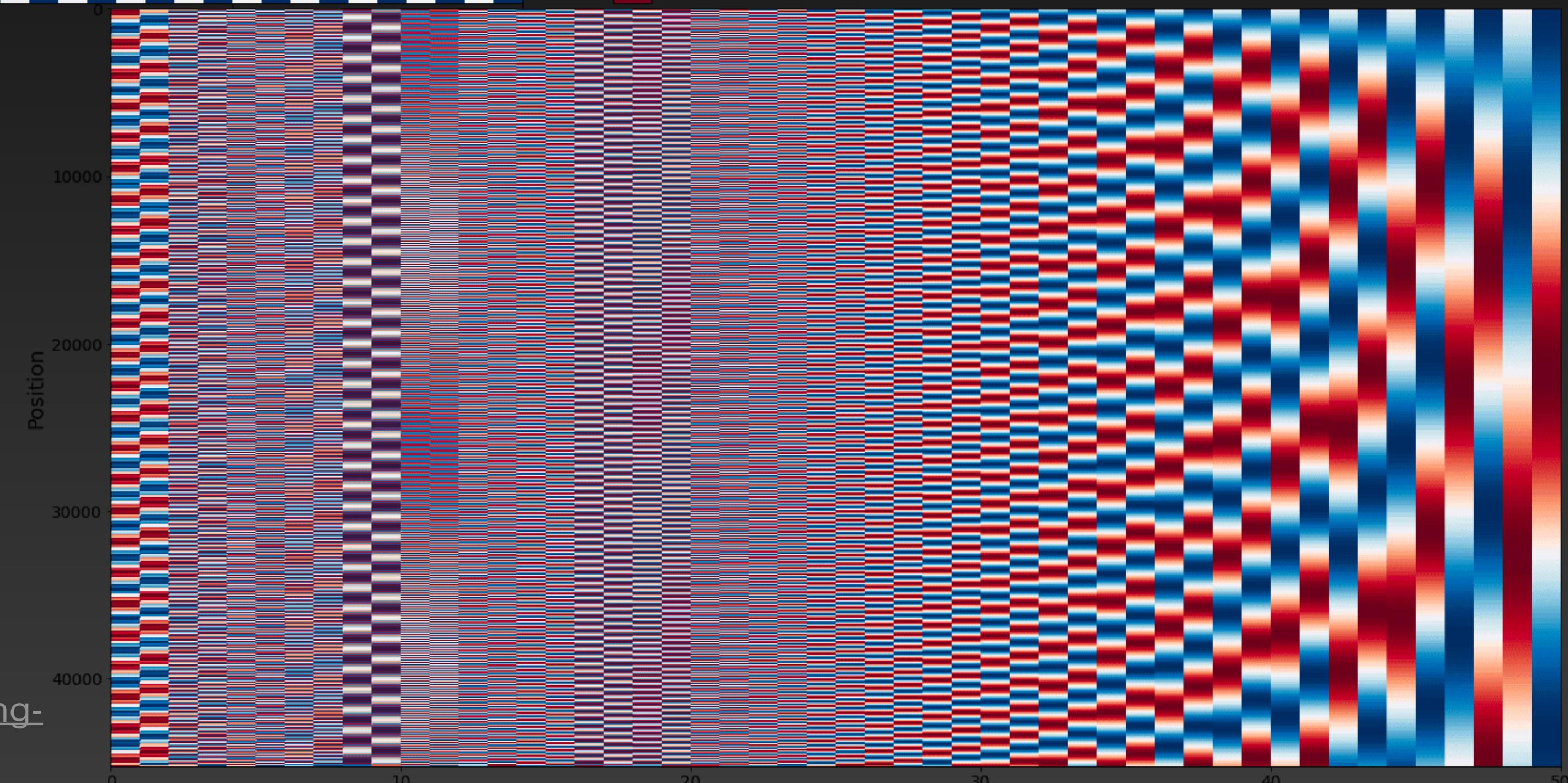


combined position embedding

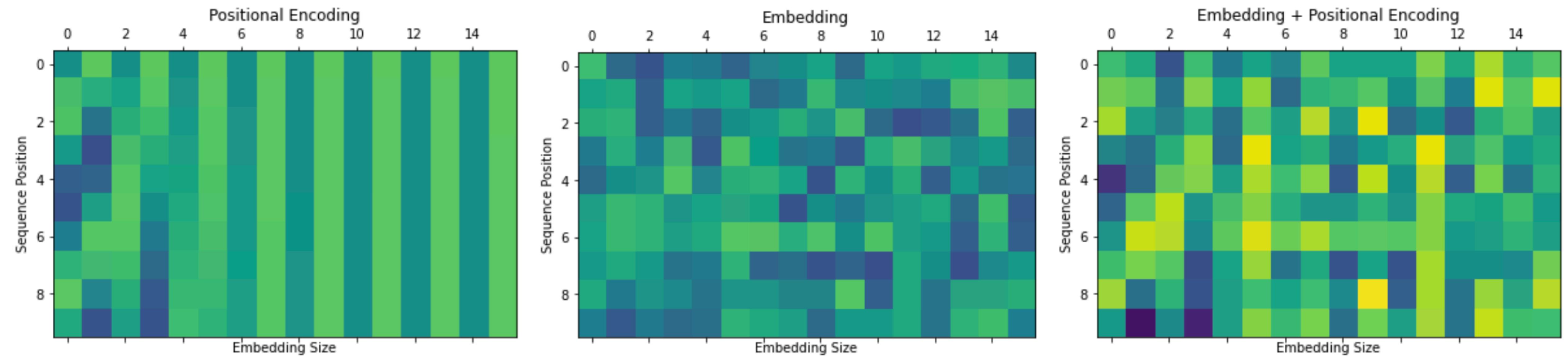




Pattern of a short
sequence (20 steps)

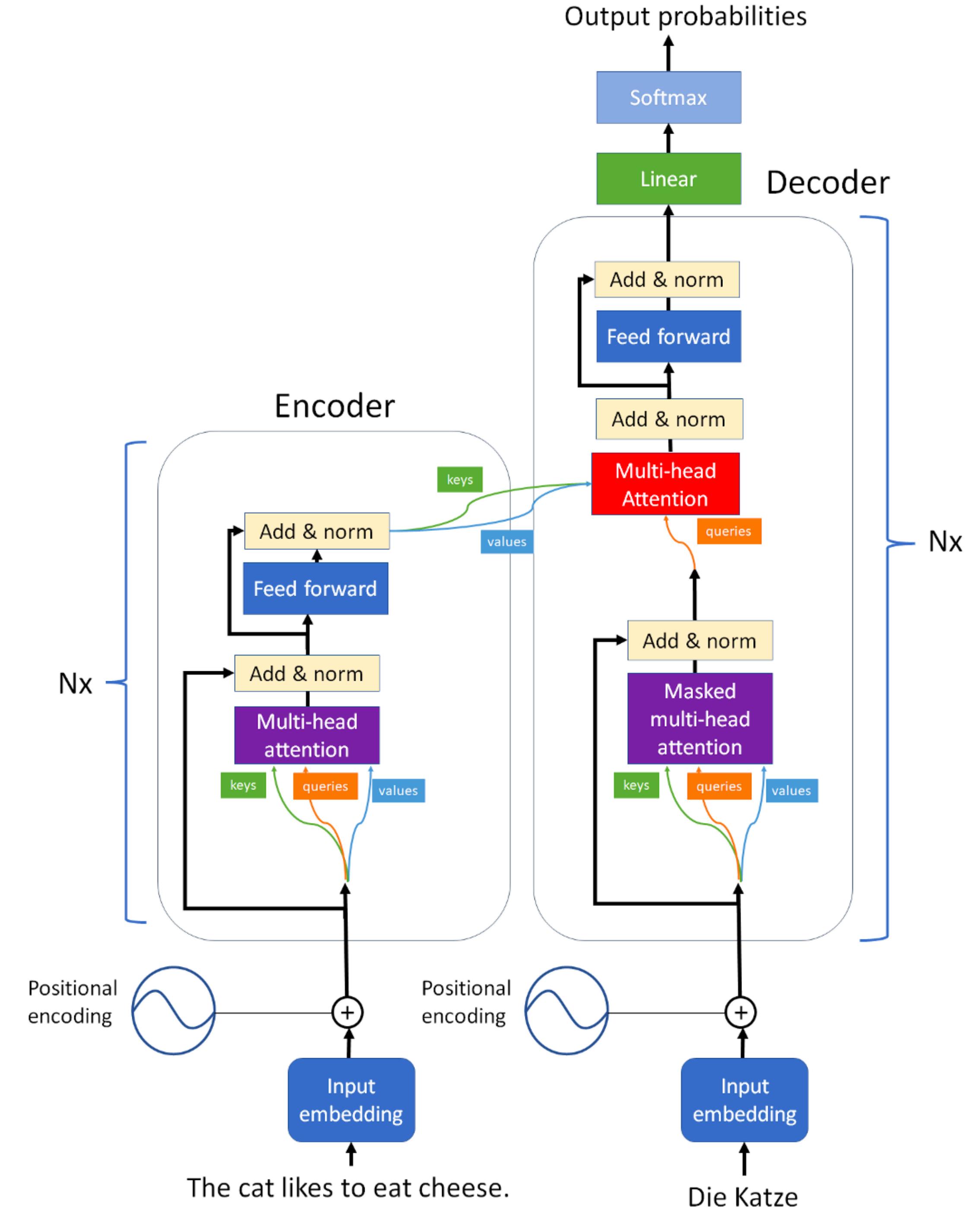


Pattern of a long
sequence (50k steps)



Full Picture Vanilla Transformer

- Data is flattened and embedded into Embeddings
- Position Encodings are created and plainly added to the embeddings
- Key, Query, Value Multi Head Self Attention is encoded
- Causal Self-Attention is used to create all decoder queries at once
- The Encoder and Decoder are combined as keys and values from the encoder and queries from the decoder for the output

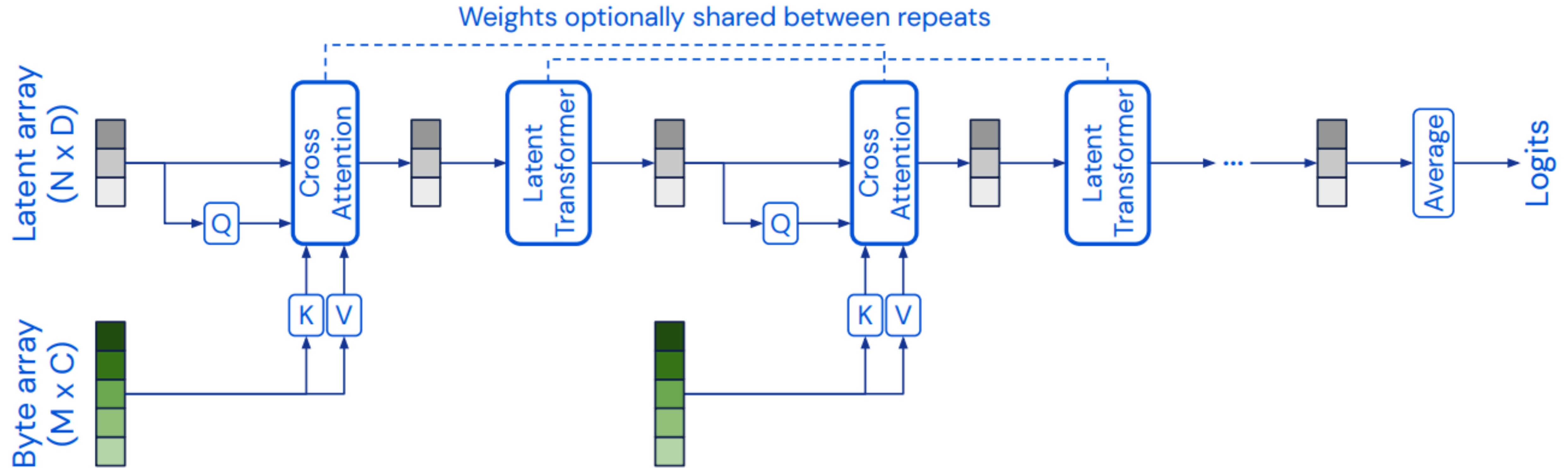


Some other aspect

We don't have time to do that in detail too

- The classic Transformer has an encoder and decoder part
 - Depending on the use case you only need one
 - GPT: Decoder only / focus on generation
 - BERT: Encoder only / focus on classification
 - T5: All of it
 - Physiological Transformers often classify: Often only Transformer
- Training Types matter:
 - Mainly Masked Self-Supervised Pre Training. Especially in the encoder in both directions
 - BERT
 - Wave2Vec 2.0

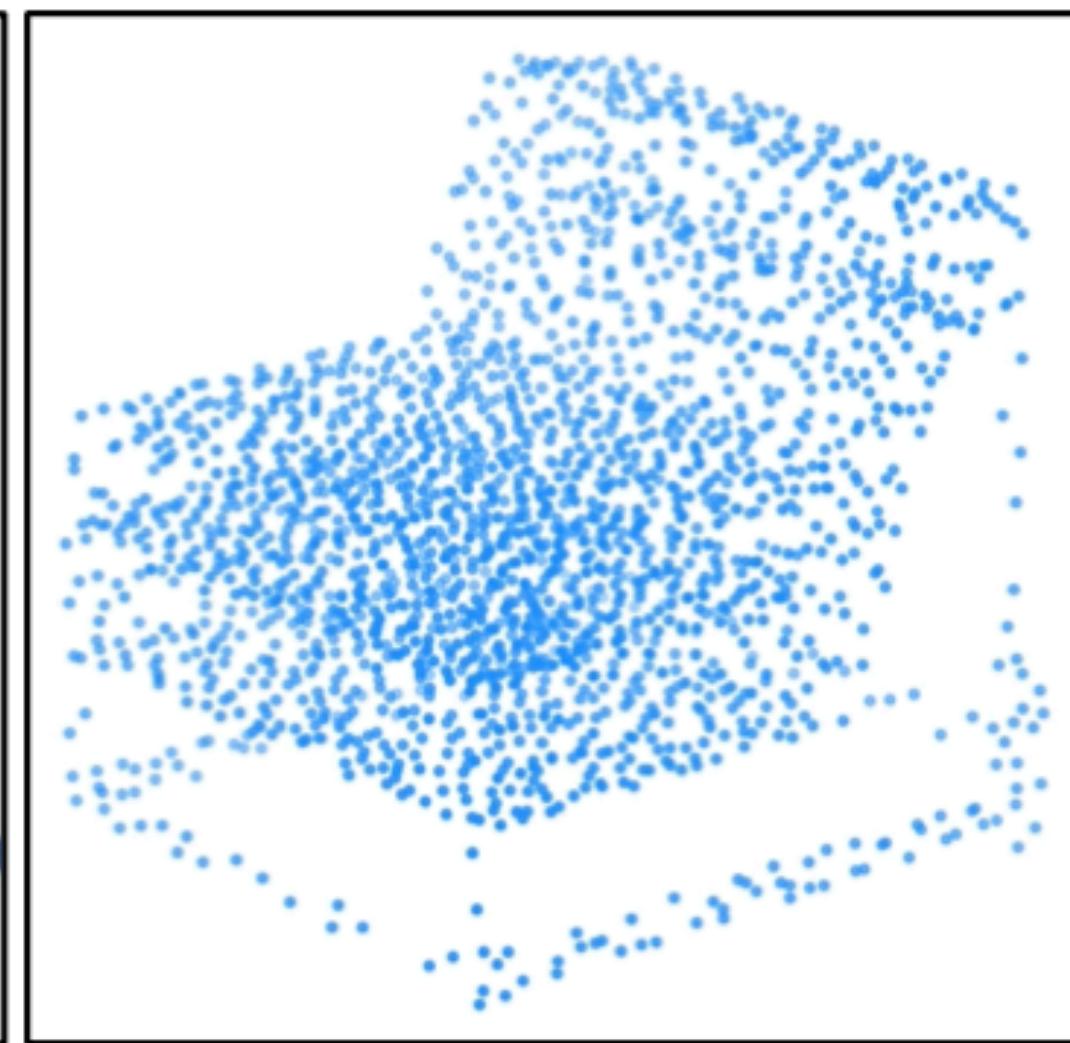
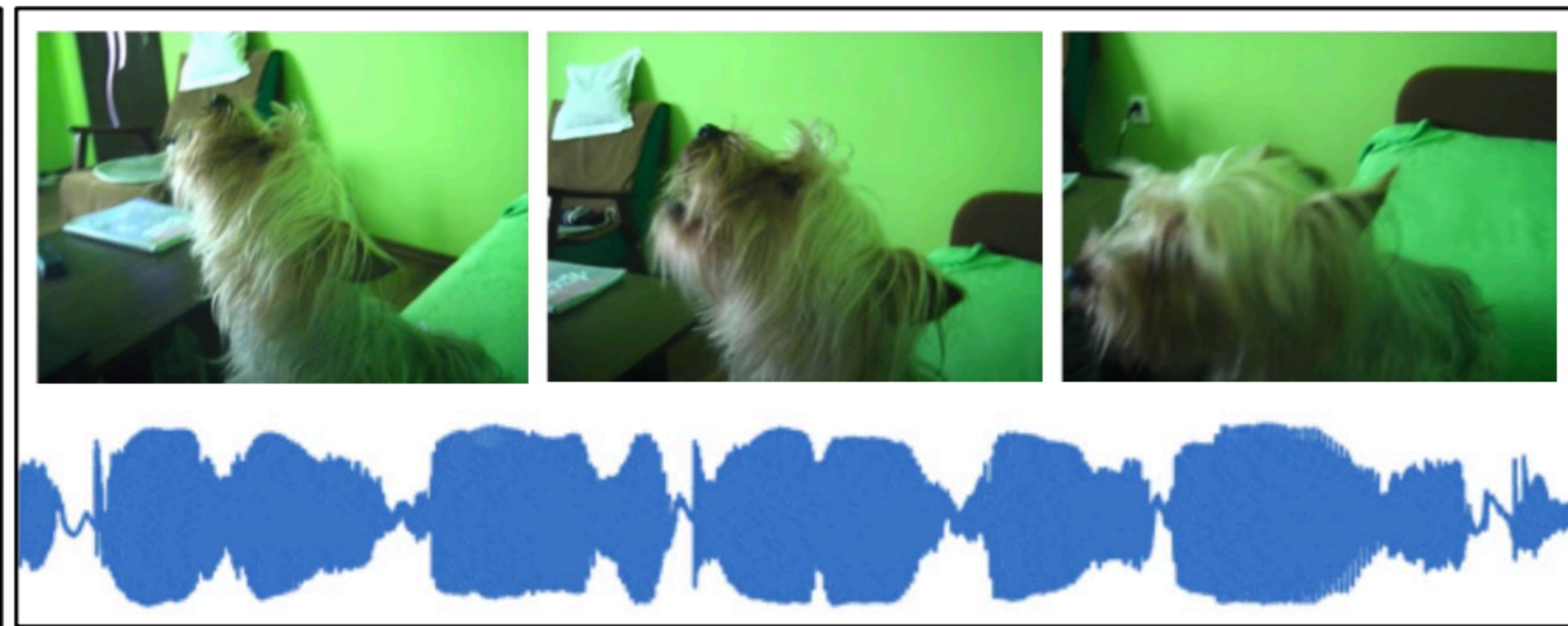
Perceiver



Pervceiver

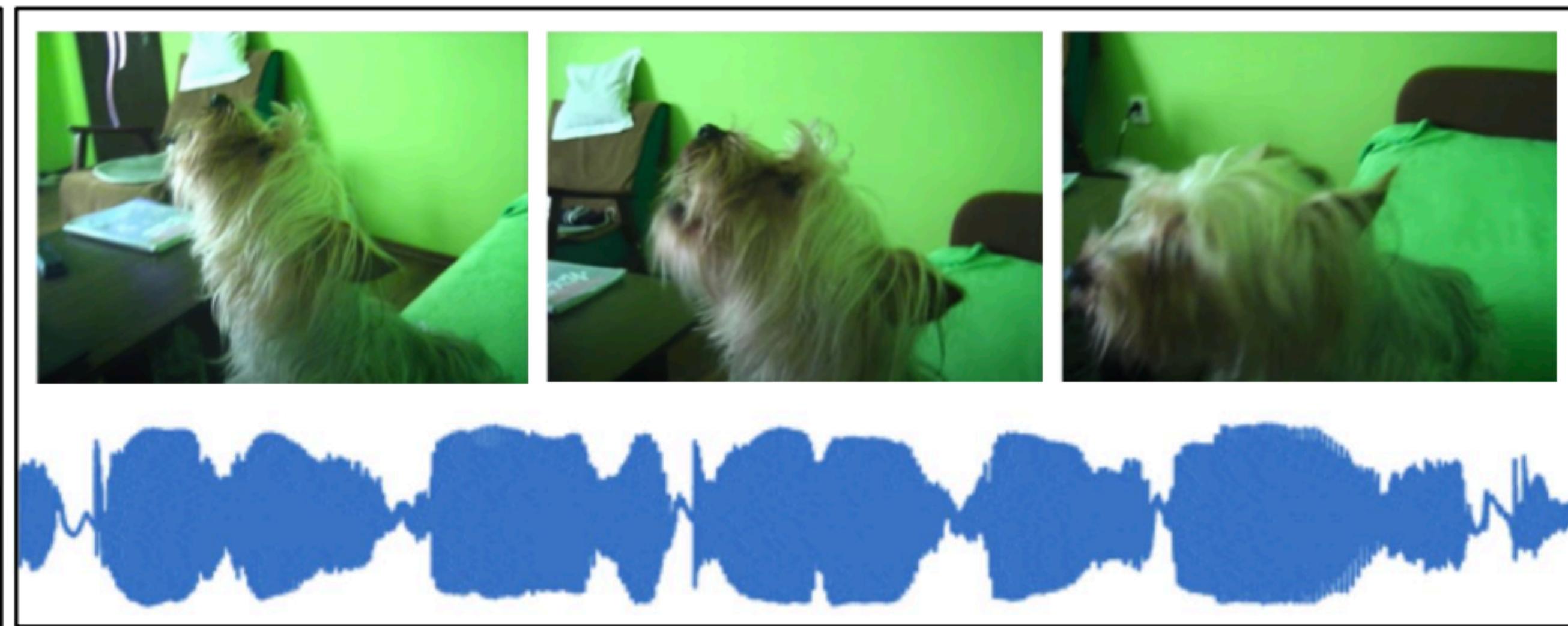
General Idea

- We have less timepoints than the original input in the latent dimension for most layers (Latent Transformer)
- We occasionally attend to the whole input again (Cross Attention)
- Since the latent is much smaller than the sequence length, we gain speedups
- We only allow classification for now / only encoder of a Transformer
- Instead of Modality specific Architectures we always keep the Architecture the same and:
 - Always just flatten the input into one vector (or in patches: matrix)
 - Use modality specific position encodings
- Some other tricks like weight sharing exist, but are again modality specific

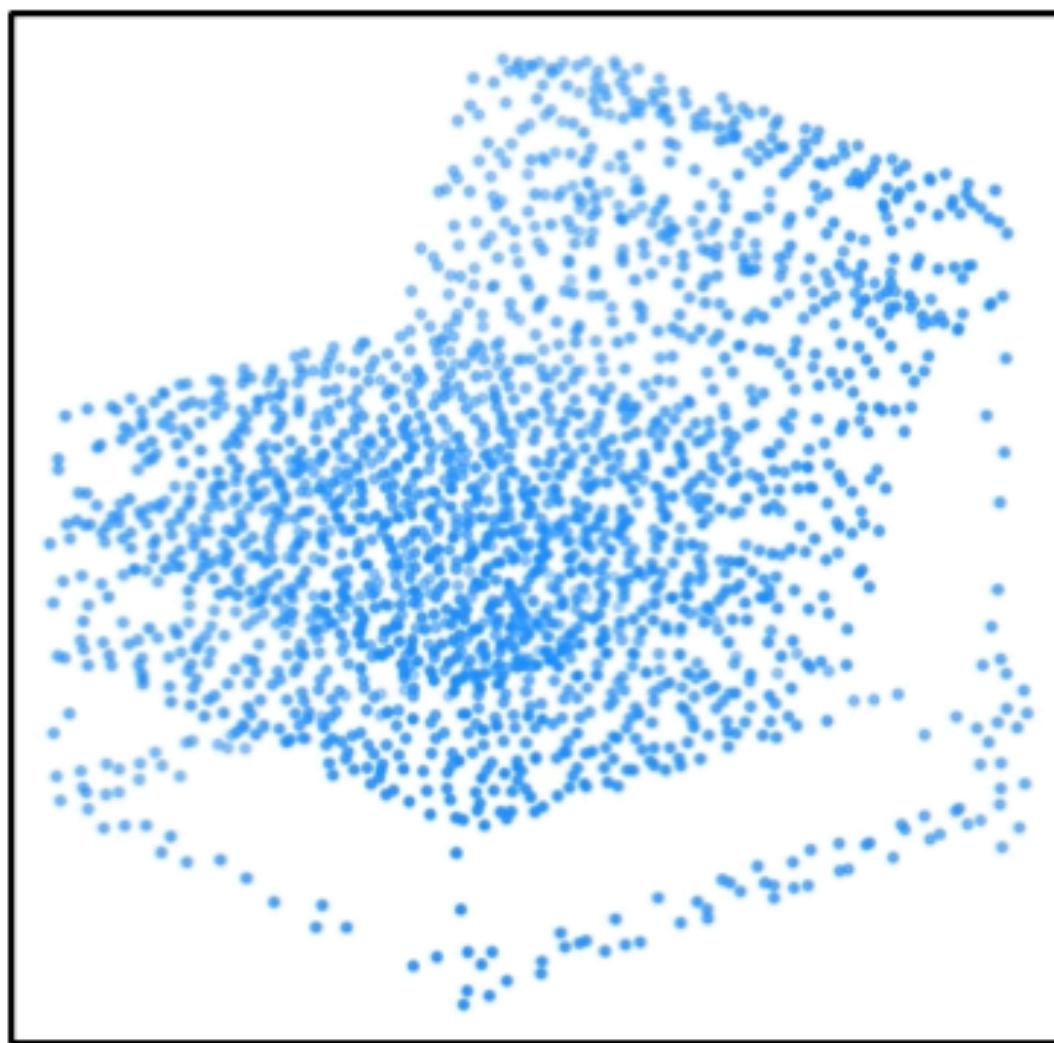




One ImageNet Model



One AudioSet Model



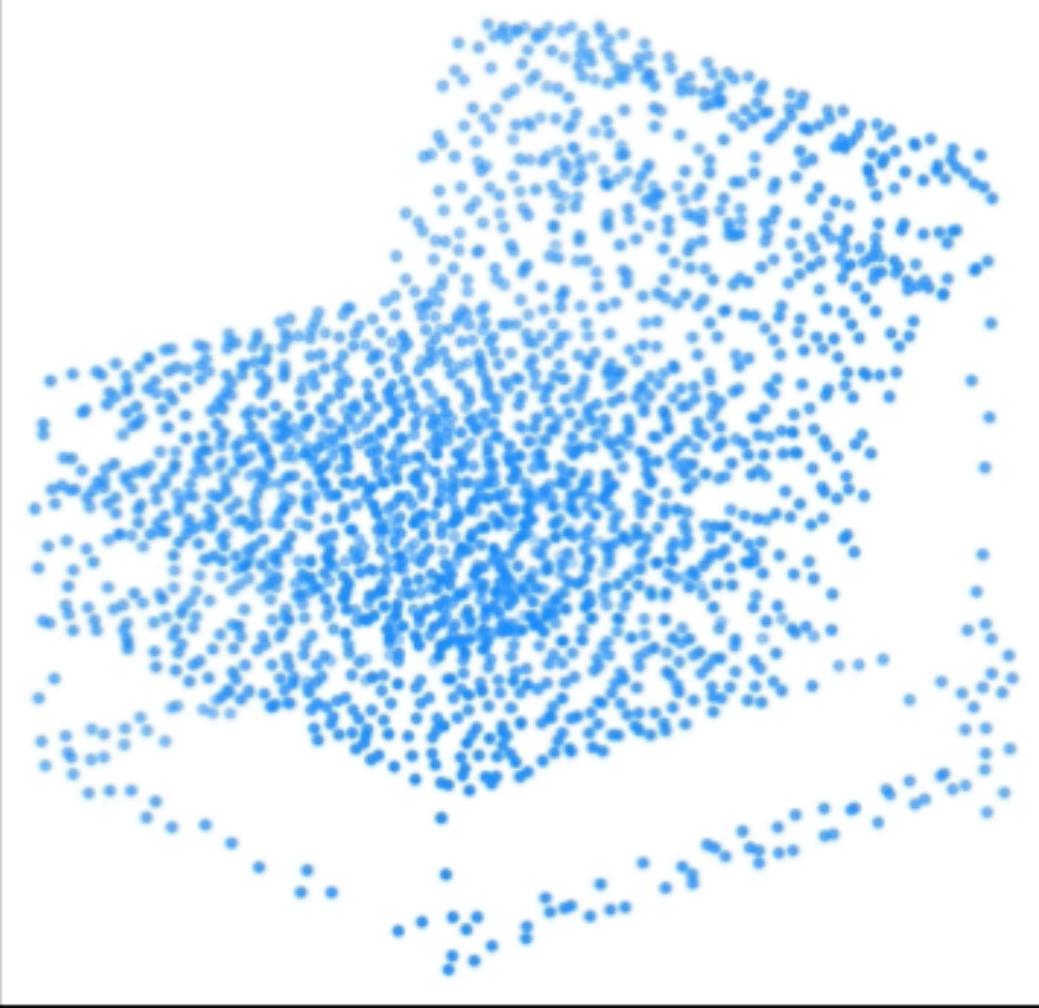
One ModelNet40 Model



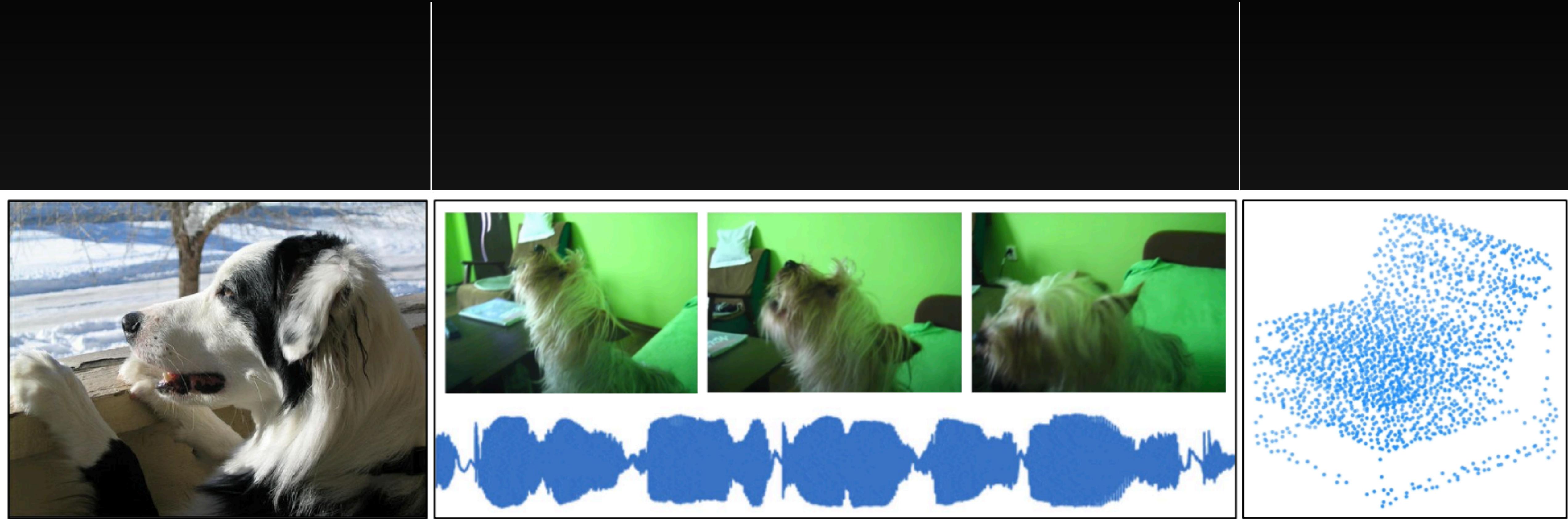
One ImageNet Model



One AudioSet Model



One ModelNet40 Model



	Raw	Perm.	Input RF
ResNet-50 (FF)	73.5	39.4	49
ViT-B-16 (FF)	76.7	61.7	256
Transformer (64x64) (FF)	57.0	57.0	4,096
Perceiver: (FF)	78.0	78.0	50,176
(Learned pos.)	70.9	70.9	50,176

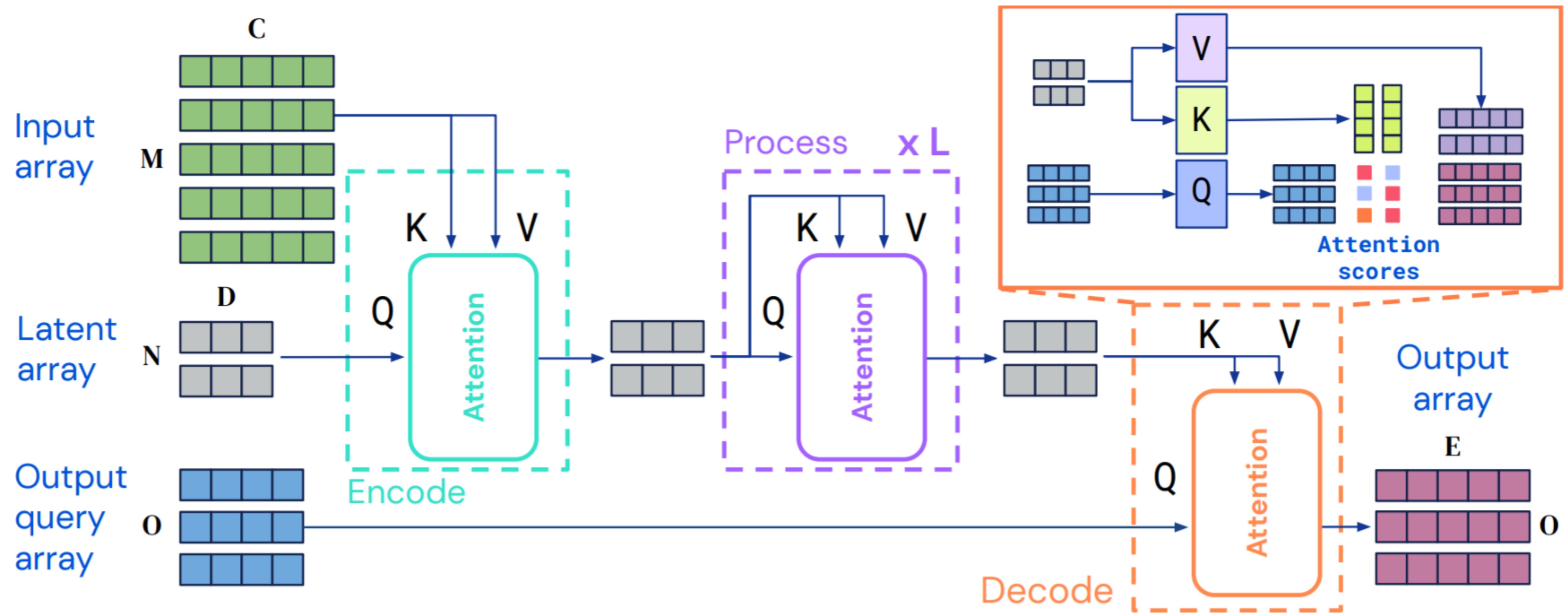
Model / Inputs	Audio	Video	A+V
Benchmark (Gemmeke et al., 2017)	31.4	-	-
Attention (Kong et al., 2018)	32.7	-	-
Multi-level Attention (Yu et al., 2018)	36.0	-	-
ResNet-50 (Ford et al., 2019)	38.0	-	-
CNN-14 (Kong et al., 2020)	43.1	-	-
CNN-14 (no balancing & no mixup) (Kong et al., 2020)	37.5	-	-
G-blend (Wang et al., 2020c)	32.4	18.8	41.8
Attention AV-fusion (Fayek & Kumar, 2020)	38.4	25.7	46.2
Perceiver (raw audio)	38.3	25.8	43.5
Perceiver (mel spectrogram)	38.4	25.8	43.2
Perceiver (mel spectrogram - tuned)	-	-	44.2

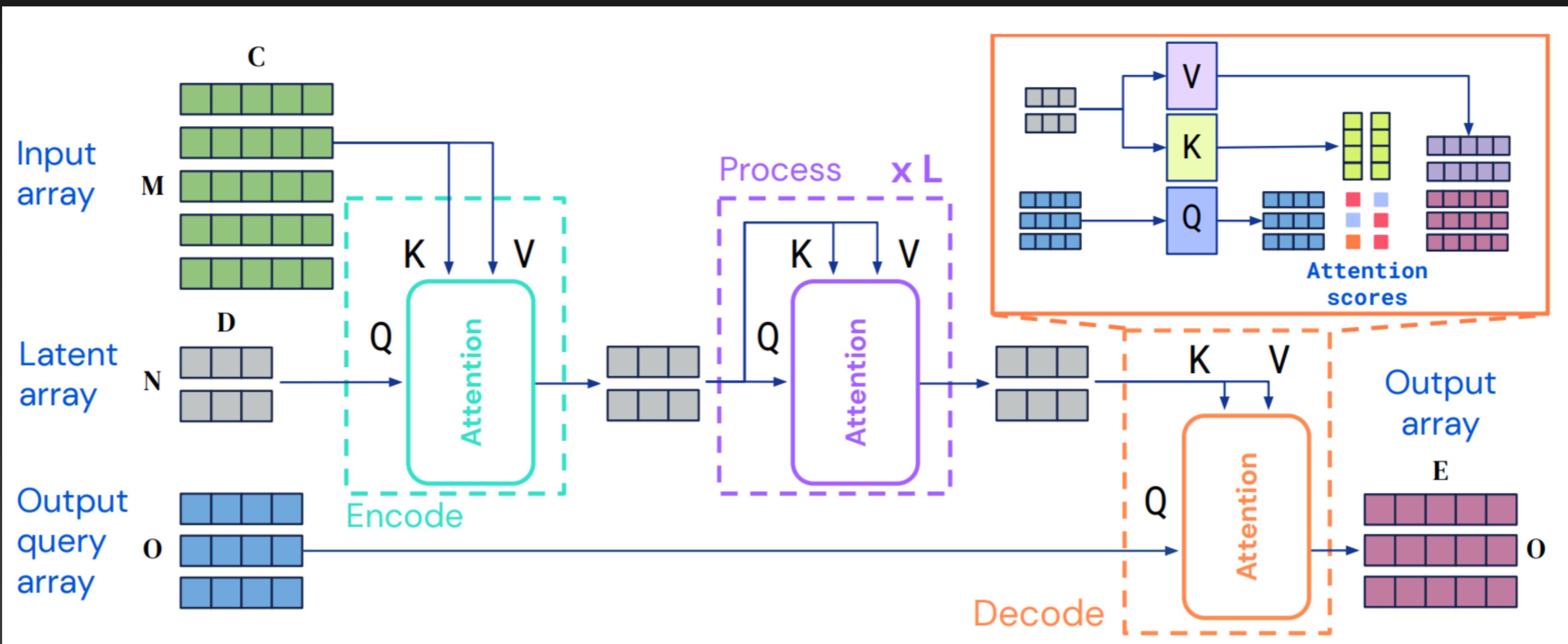
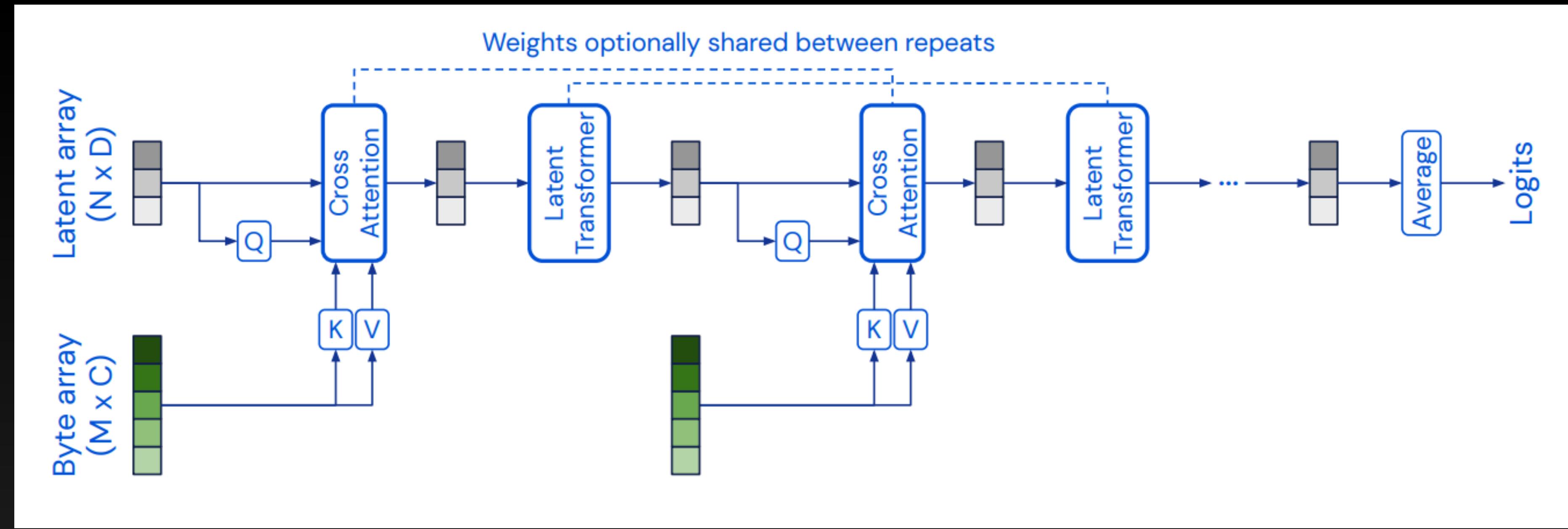
	Accuracy
PointNet++ (Qi et al., 2017)	91.9
ResNet-50 (FF)	66.3
ViT-B-2 (FF)	78.9
ViT-B-4 (FF)	73.4
ViT-B-8 (FF)	65.3
ViT-B-16 (FF)	59.6
Transformer (44x44)	82.1
Perceiver	85.7

Perceiverlo

PerceiverIO

- The original Perveiver only can do classification
- But what if you want do do all kind of tasks?
 - Decoding? Generation of Data? Segmentation?
- Key for a Transformer: Getting the queries right
- Key Idea of Perveiver: One output format to rule them all!
 - Everything is a query of position encodings or learned queries

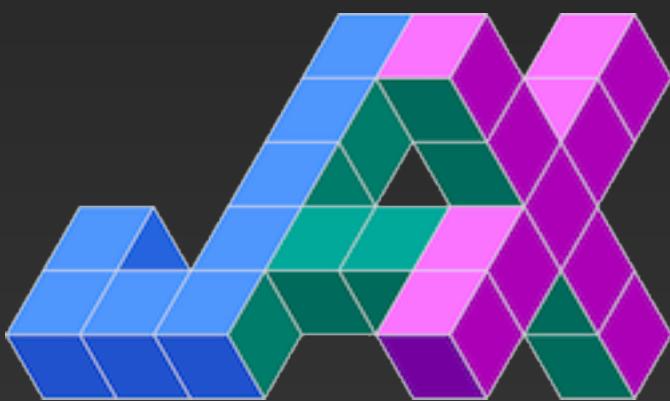
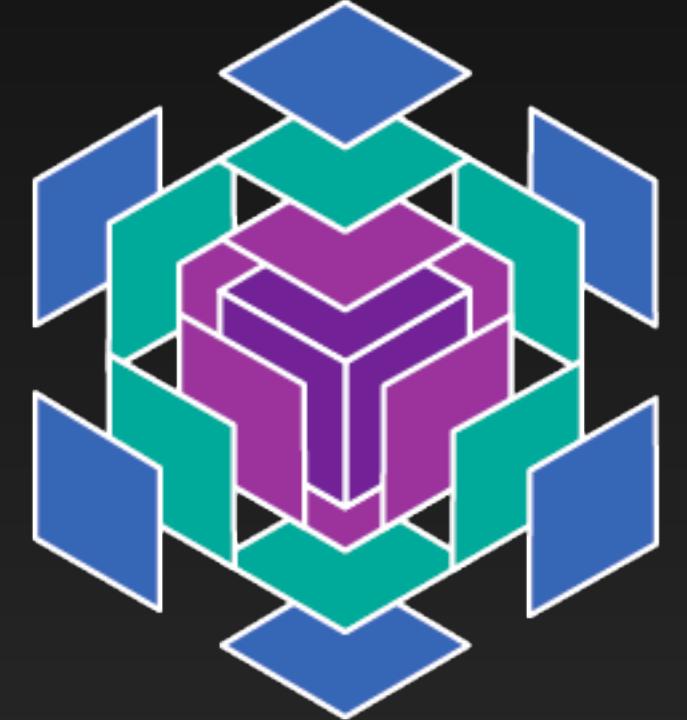




War of Frameworks

The AI Community can not decide

- Current state of Frameworks:
 - PyTorch most popular in research
 - Tensorflow is the main legacy
 - But there is JAX
 - And there is JAX (Flask)
 - TRAX
 - JAX as part of TF
 - JAX from OpenAI



Where to get Code

- PyTorch:
 - <https://github.com/krasserm/perceiver-io> (PyTorch Lightening)
 - <https://github.com/esceptico/perceiver-io/>
 - <https://github.com/lucidrains/perceiver-pytorch> (LucidRains is the PyTorch Transformer provider beside Huggingface [which is NLP focused])
 - <https://github.com/fac2003/perceiver-multi-modality-pytorch> (Multimodal Extension)
 - <https://github.com/frenkiboy/perceiver-pytorch> (Radiology Data)
- Tensorflow:
 - <https://github.com/Rishit-dagli/Perceiver>
- JAX:
 - <https://github.com/deepmind/deepmind-research/tree/master/perceiver>

Code examples