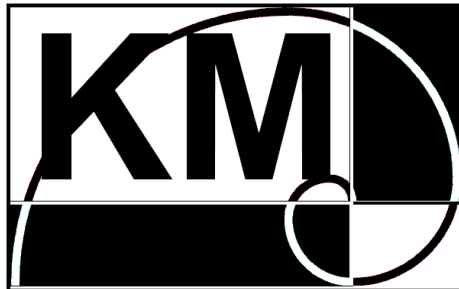


# Strojové učení II



## Sequences



Institute of Information Theory  
and Automation of the AS CR

# Sequences

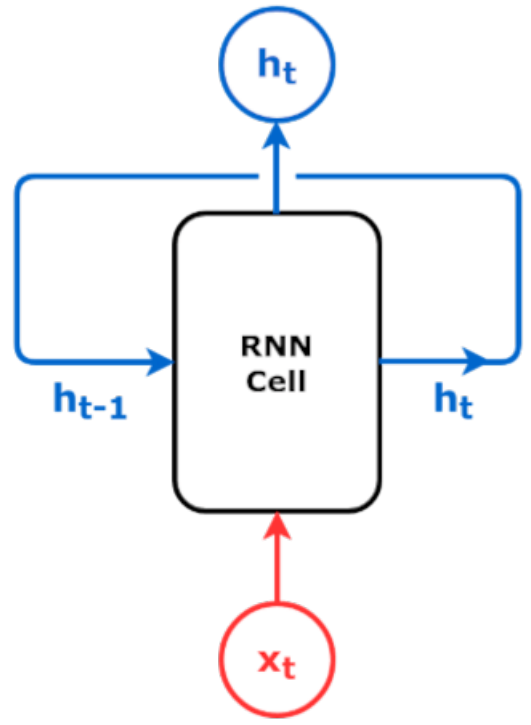


- Input:
  - sequence = an **ordered** set of data points
- Output:
  - single label: RNN
  - sequence (of different length): encoder-decoder model
- Application: time series, Natural Language Processing (NLP)

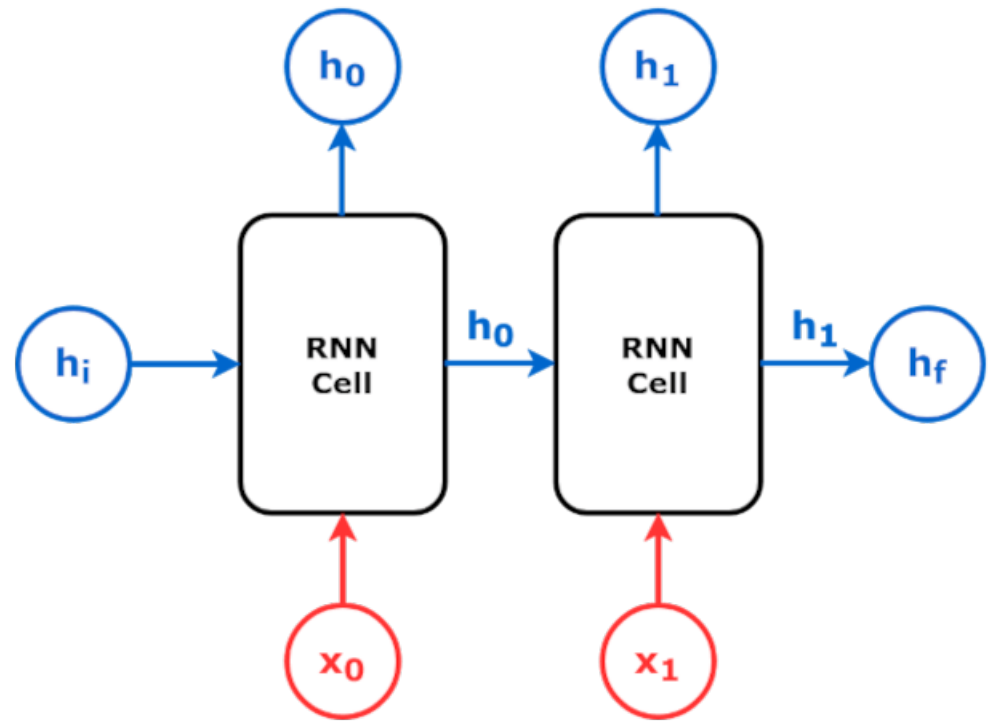


# Recurrent Neural Network (RNN)

- Basic RNN

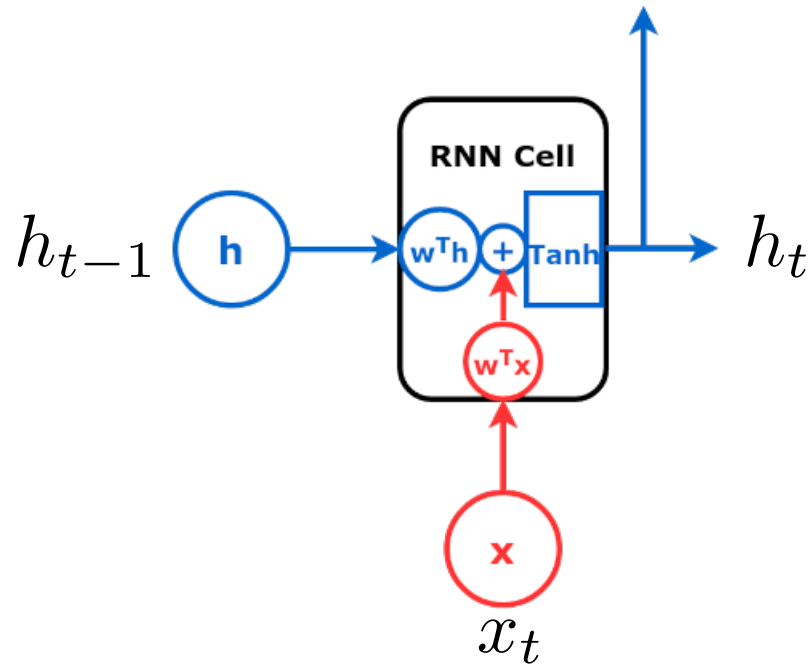


- Unfolded RNN





# RNN cell



$$t_h = W_h h_{t-1} + b_h$$

$$h_t = \tanh(t_h + t_x)$$

$$t_x = W_i x_t + b_i$$

**Problem:** Long-term dependencies

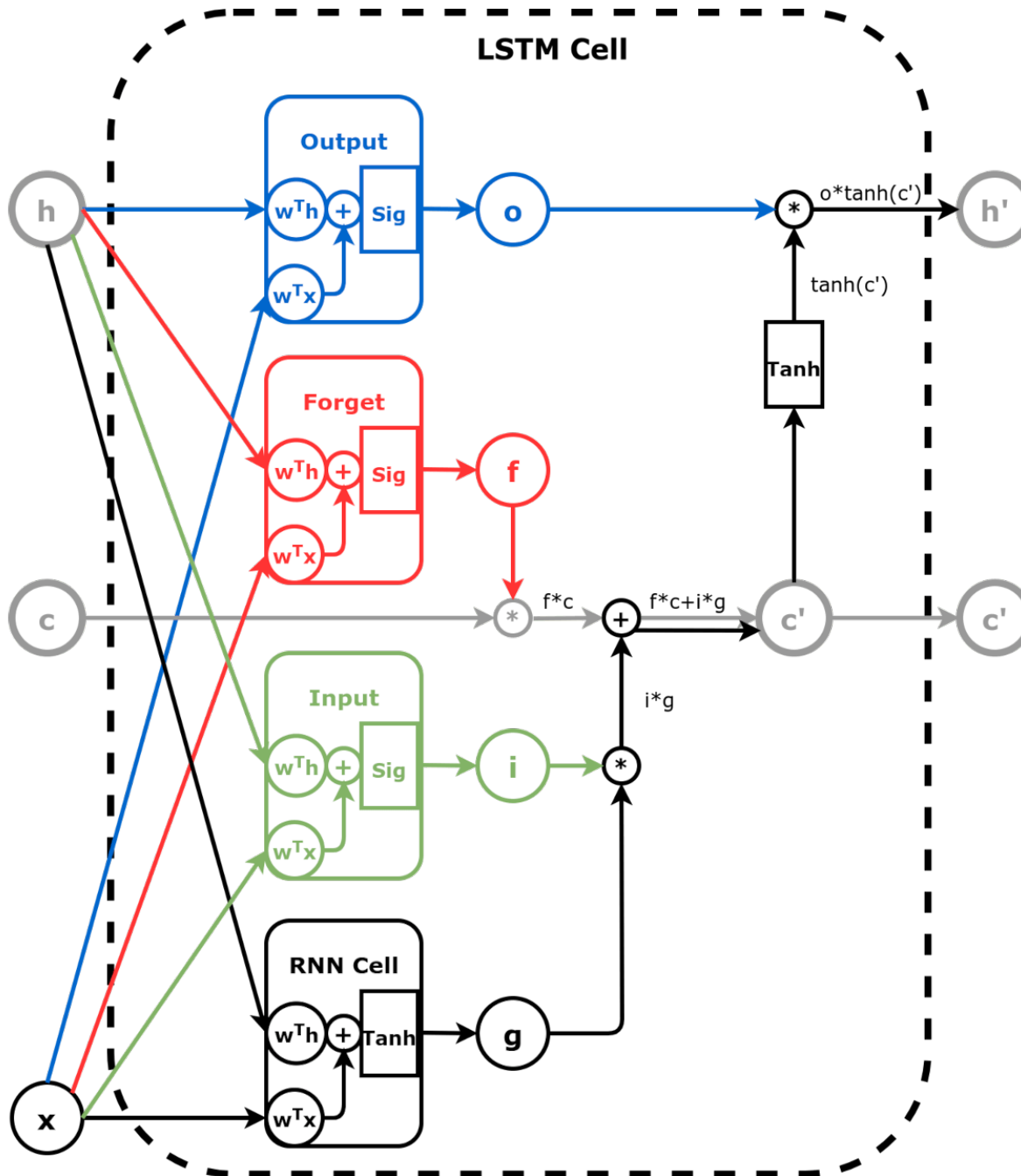


# Multiple Time Scales

- Long-term dependencies are weak in the basic RNN cell.
- Solution:
  - Skip connections through time
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Unit (GRU)



# LSTM

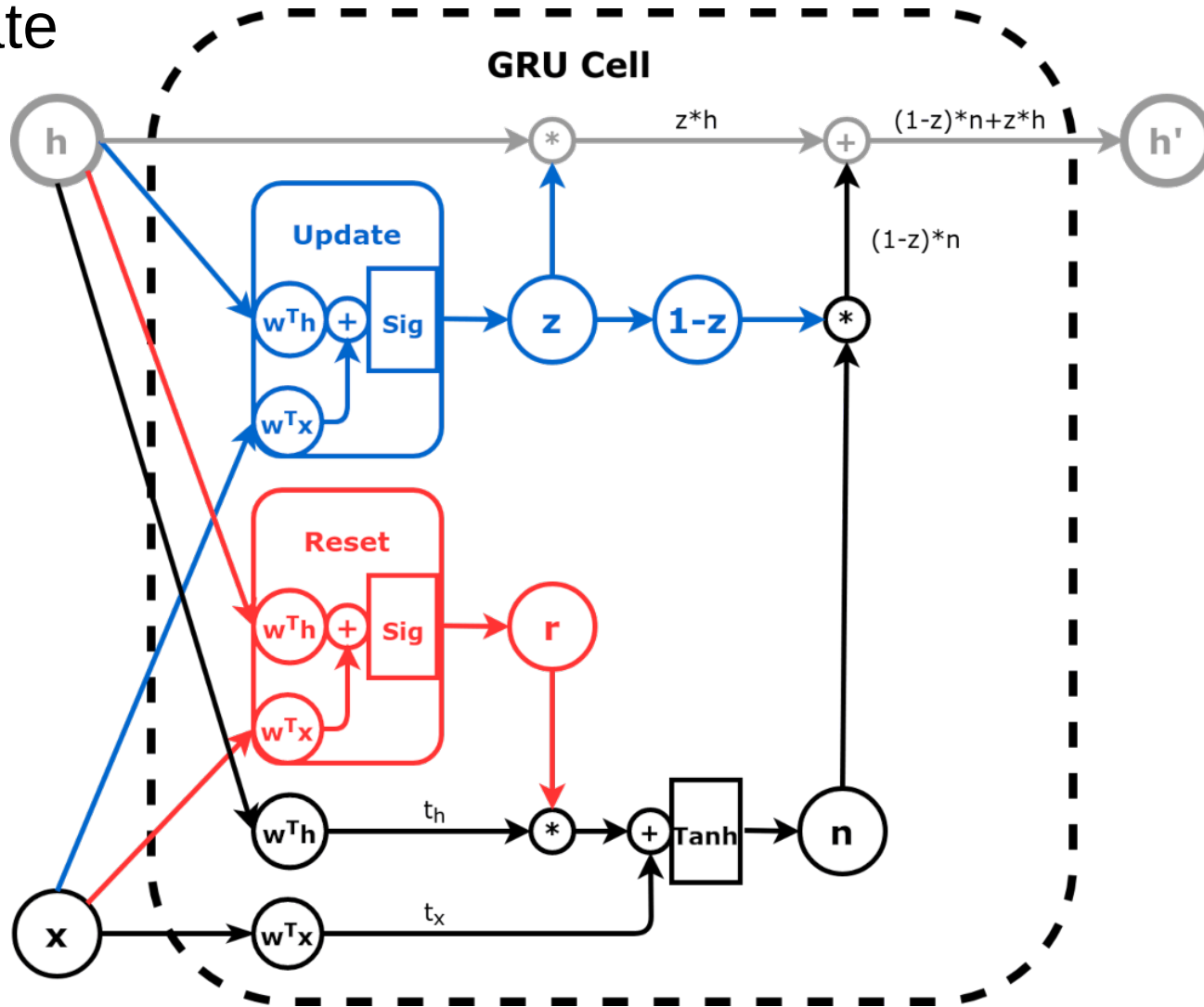


Discuss cases:  
 $i=1, f=0;$   
 $i=0, f=1;$



# GRU

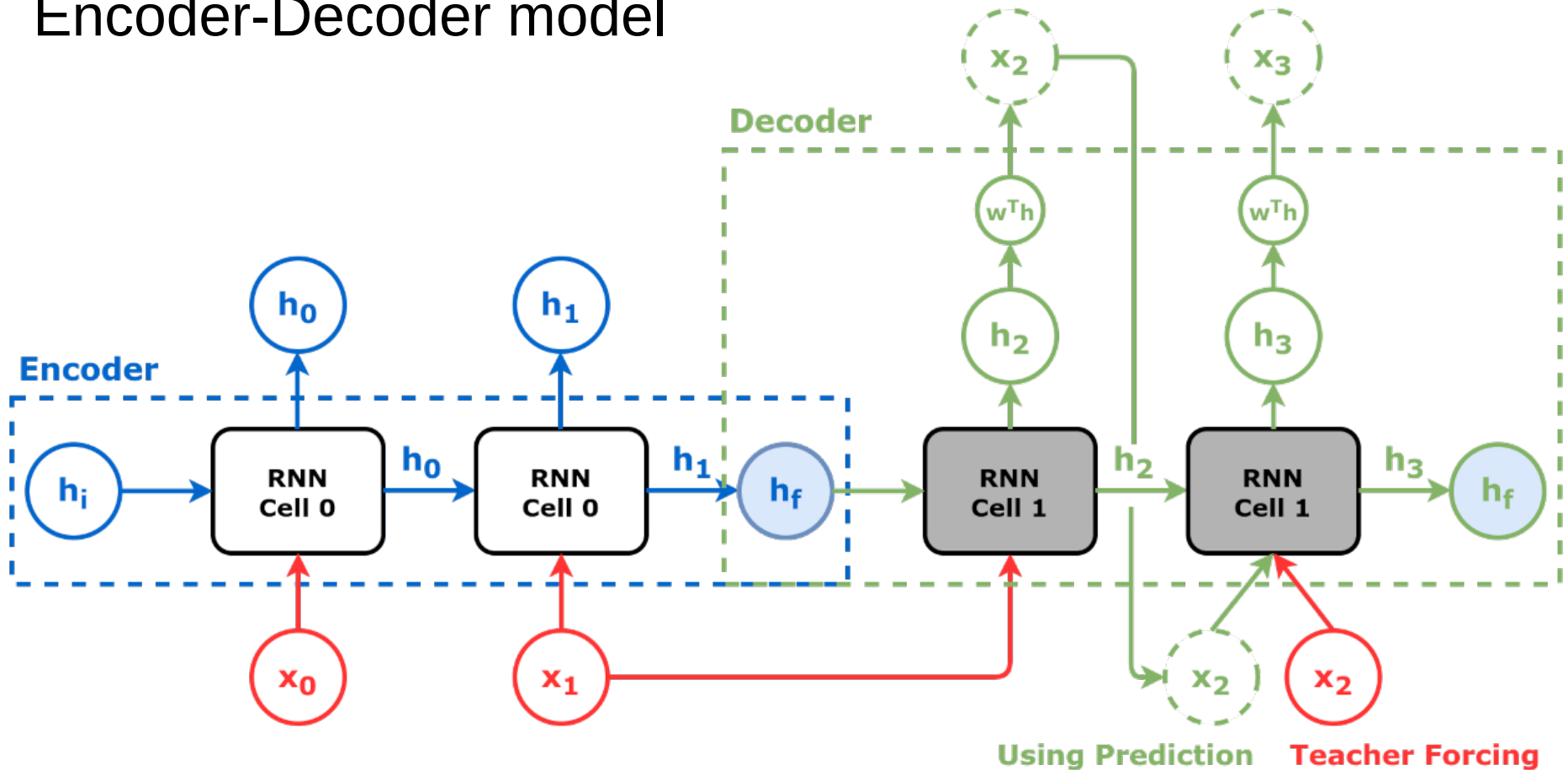
- No input gate
- No cell state





# Sequence-to-Sequence

- Encoder-Decoder model



- The European economic zone → la zona económica europea
  - Word order
  - the → el / la: To which word do we pay **attention**?





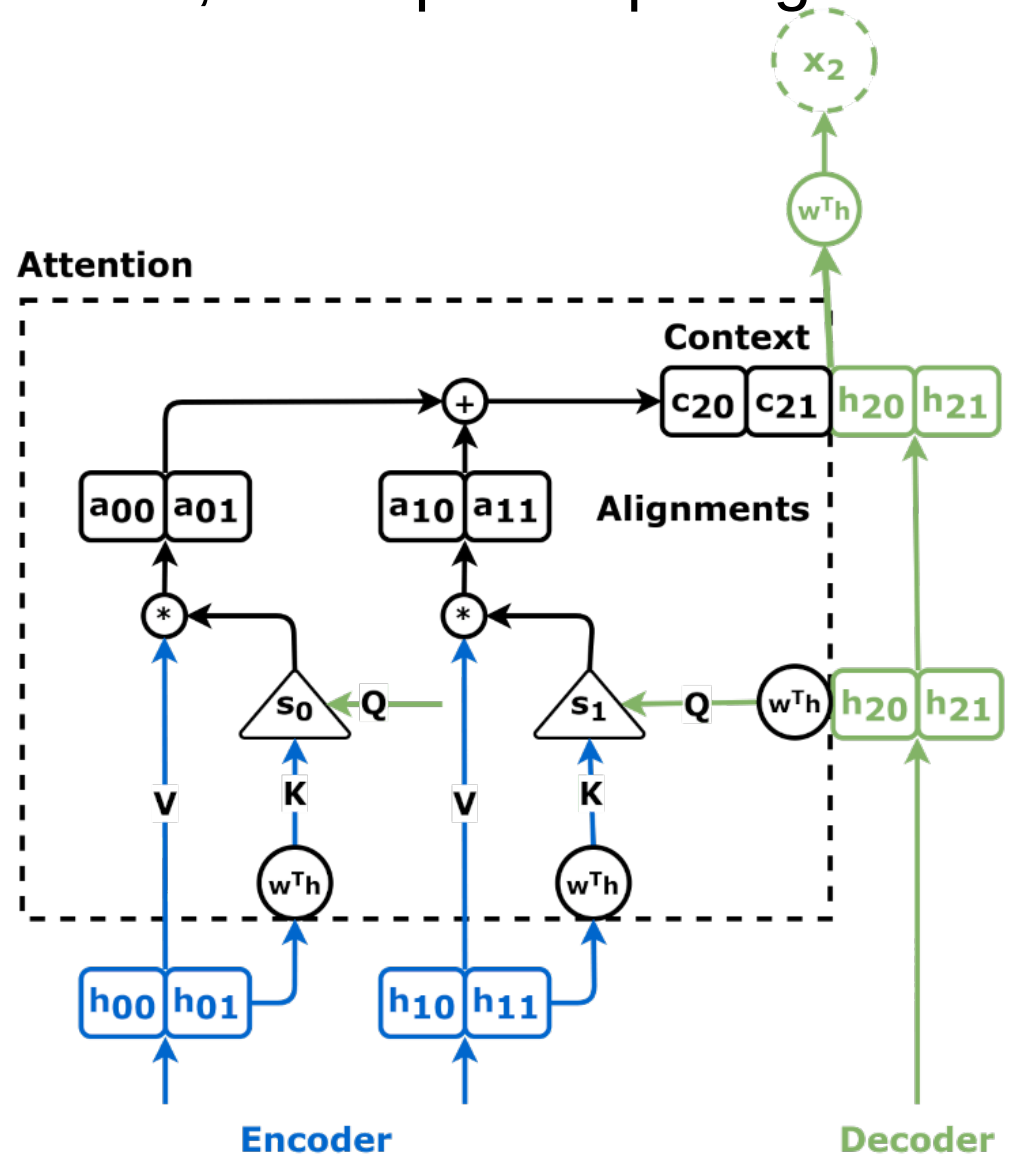
# Attention

- V... values [D x L], D... hidden dim., L... input seq. length
- K... keys [D x L]
- q... query

$$c = V (K^T q)_{\text{softmax}}$$

$$c = \sum_i \text{softmax} \left( \frac{\langle q, k_i \rangle}{\sqrt{D}} \right) v_i$$

What is D and L in the figure?





# Self-Attention

- Forget RNN!

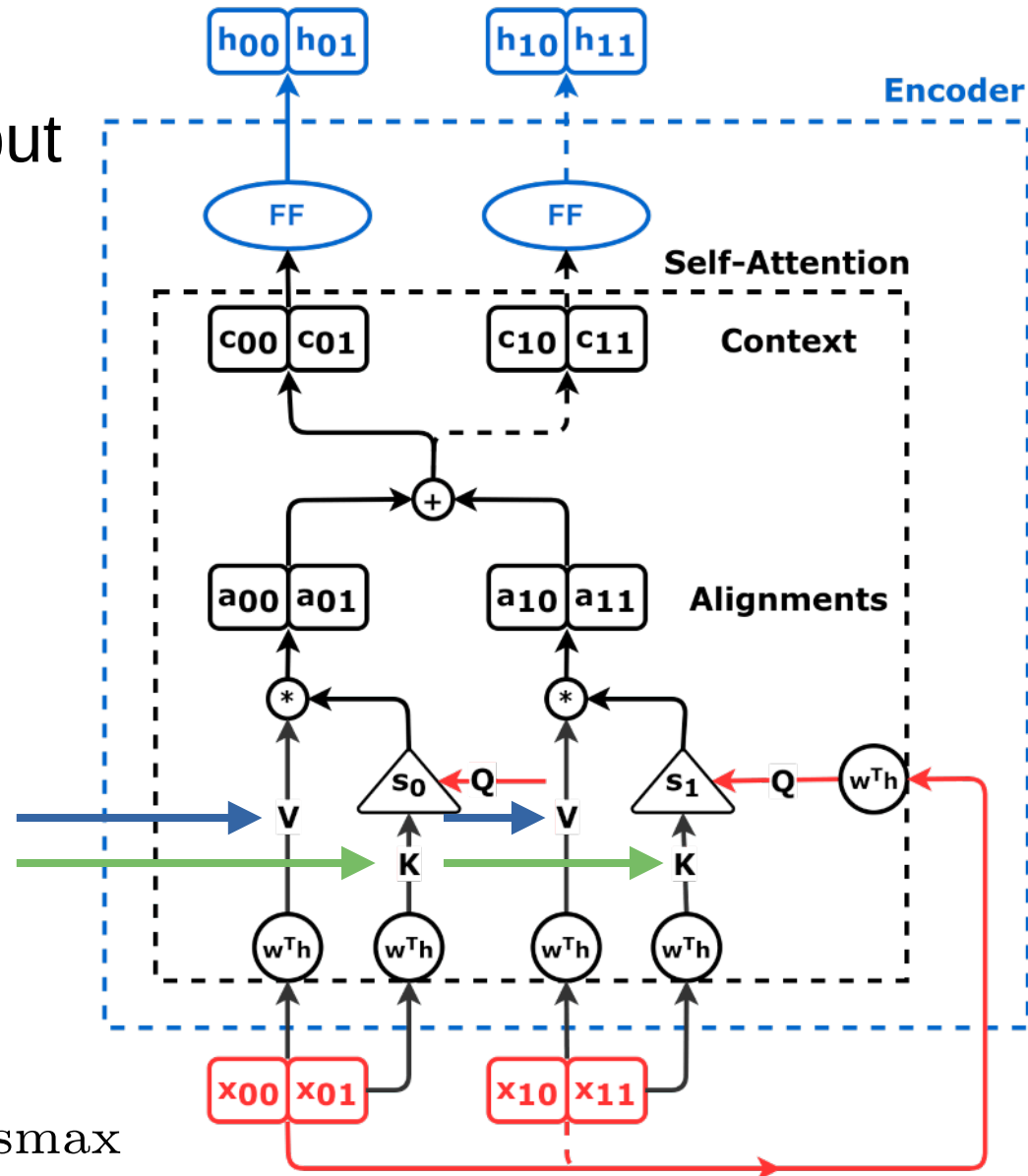
Compute V, K, Q from the input sequence  $X = [x_1, \dots, x_L]$

- Attention

$$c = \underline{V} \left( \underline{K}^T \underline{q} \right)_{\text{ssmax}}$$

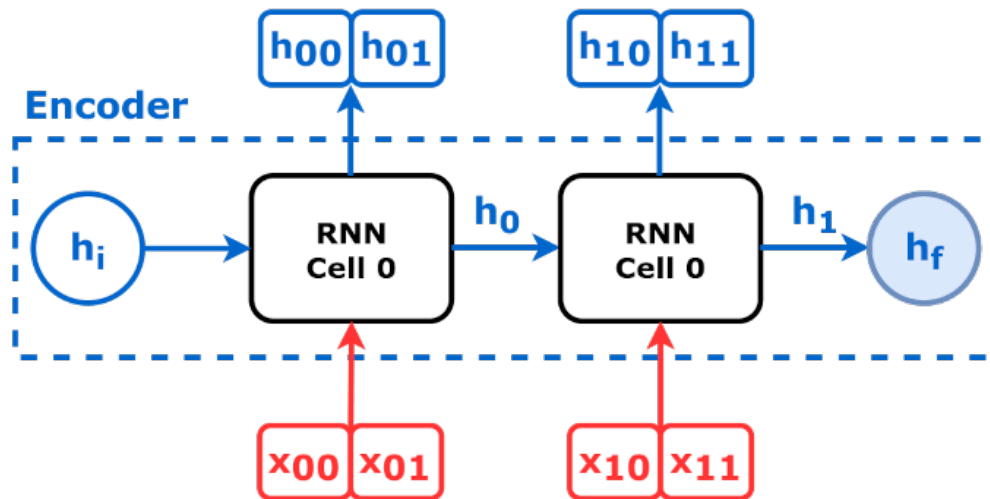
- Self-Attention

$$C = \left[ \underline{W}_v X \right] \left( \left[ \underline{W}_k X \right]^T \left[ \underline{W}_q X \right] \right)_{\text{ssmax}}$$

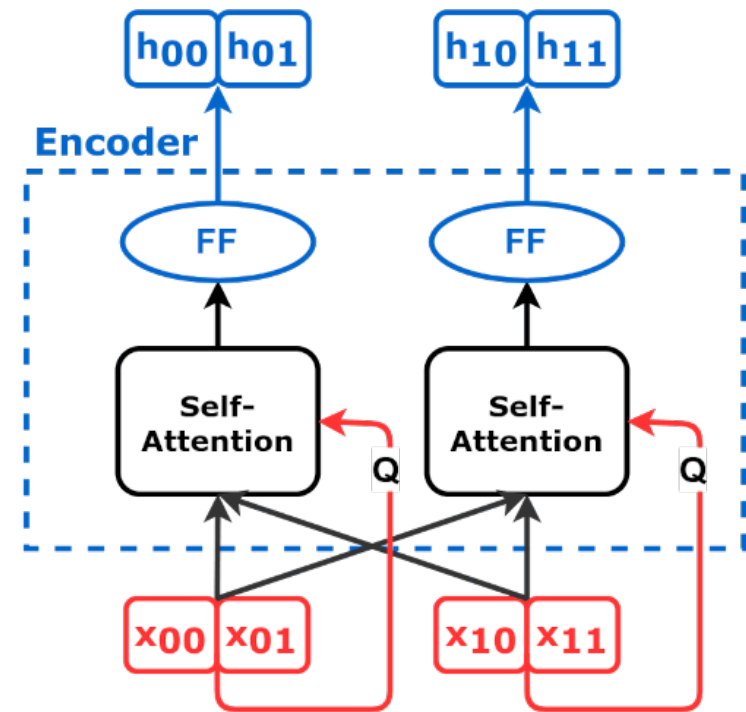




# RNN vs Self-Attention



RNN



Self-Attention



# Convolution vs Self-Attention

$$C = \underbrace{[W_v X]} \left( \underbrace{[W_k X]^T} \underbrace{[W_q X]} \right)_{\text{ssmax}}$$

$$c_i = \sum_j \langle \underbrace{q_i}, \underbrace{k_j} \rangle \underbrace{v_j}$$

$$W_v = W_k = W_q = I$$

$$C = [X] \left( [X]^T [X] \right)_{\text{ssmax}}$$

Convolution:

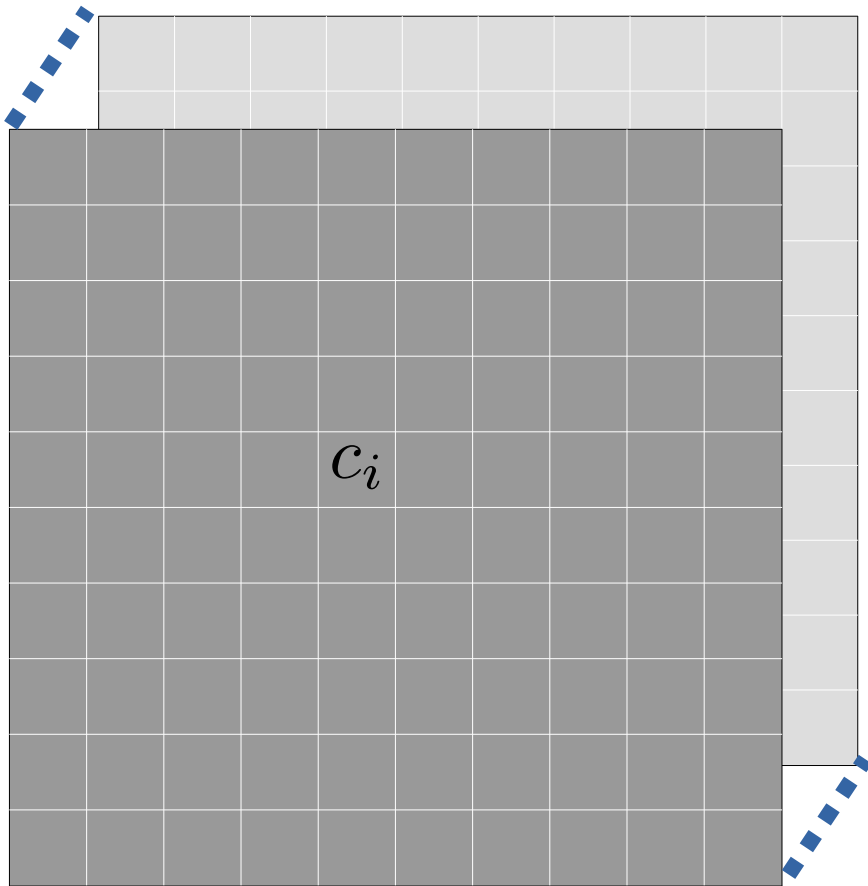
$$c_i = \sum_j h_{i-j} x_j$$

$$c_i = \sum_j \langle x_i, x_j \rangle x_j$$



# Convolution vs Self-Attention

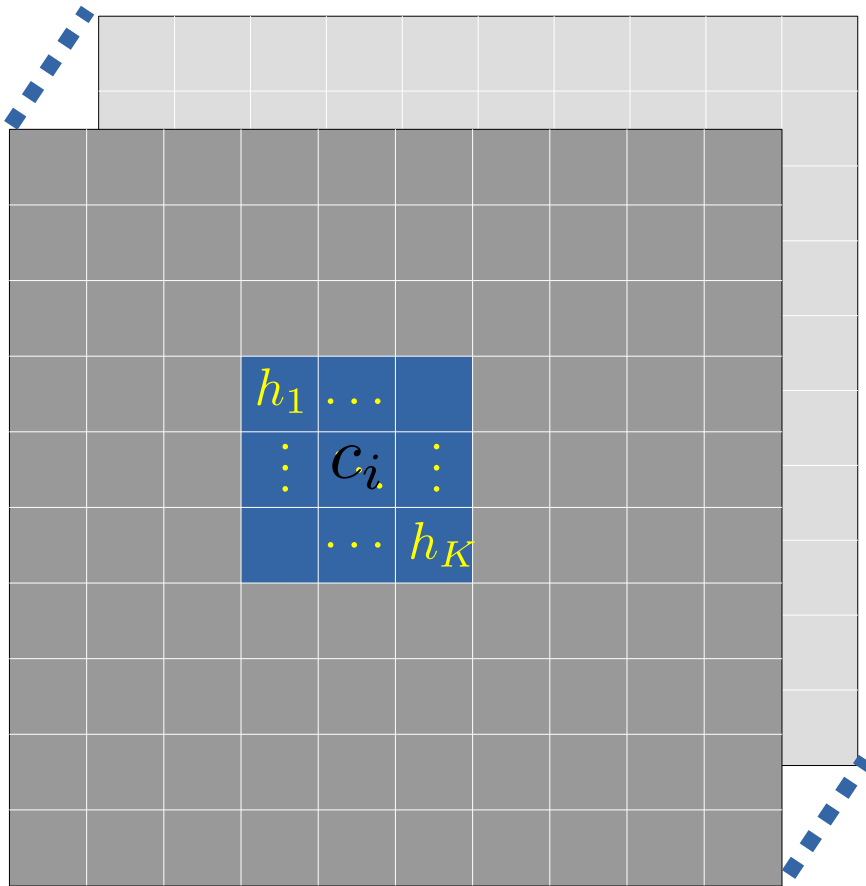
$$c_i = \sum_j h_{i-j} \underline{x_j}$$





# Convolution vs Self-Attention

$$c_i = \sum_j h_{i-j} \underline{x_j}$$



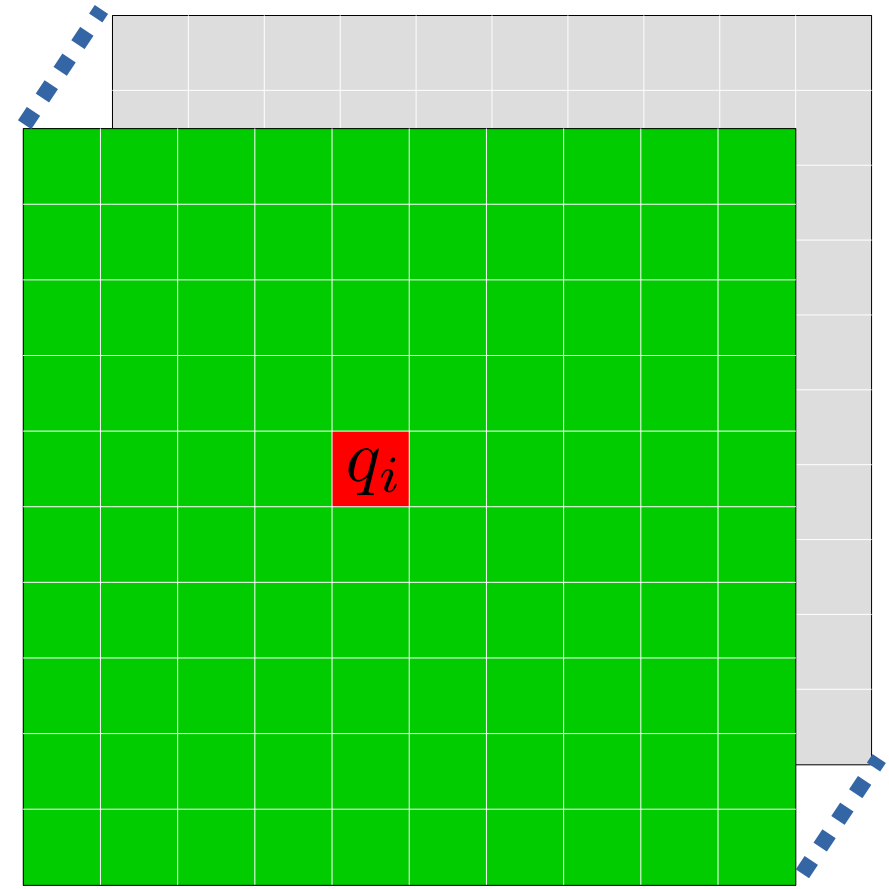
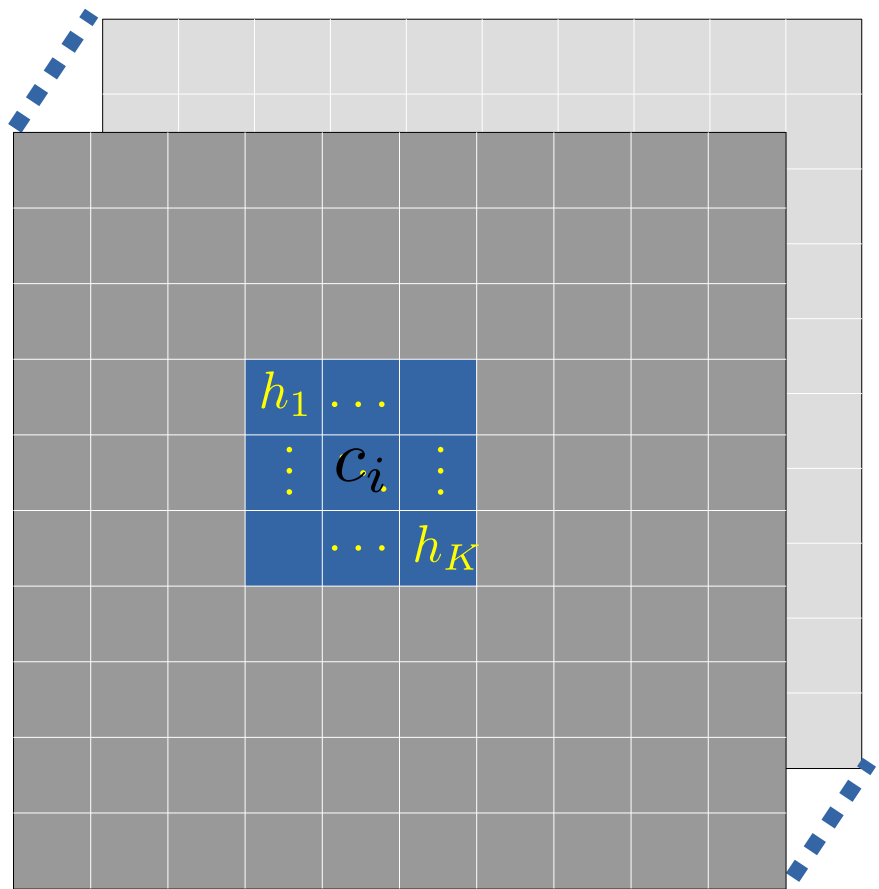
Local and fixed  
during inference



# Convolution vs Self-Attention

$$c_i = \sum_j h_{i-j} x_j$$

$$c_i = \sum_j \alpha_{ij} \langle q_i, k_j \rangle v_j$$



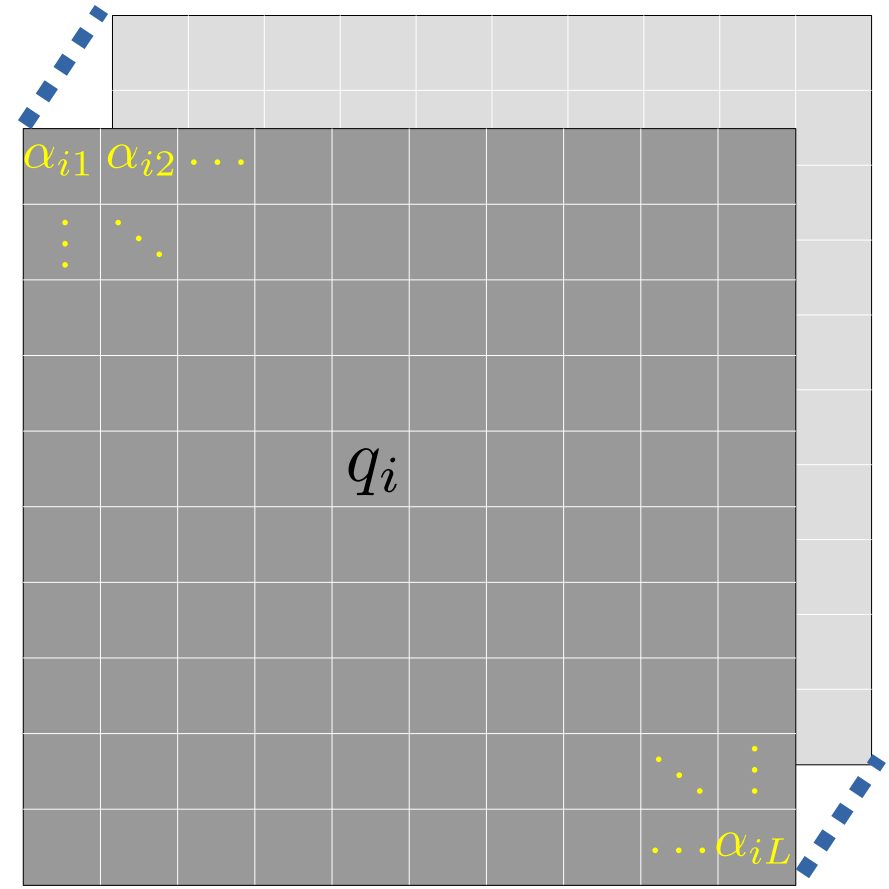
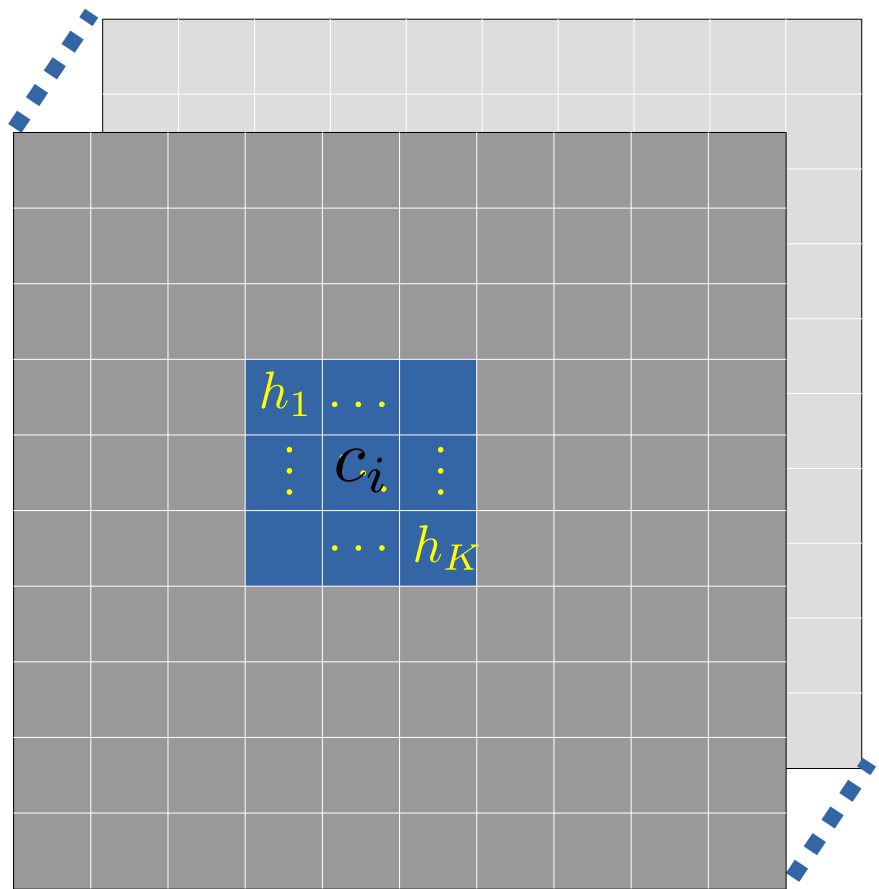
Local and fixed during inference



# Convolution vs Self-Attention

$$c_i = \sum_j h_{i-j} x_j$$

$$c_i = \sum_j \langle q_i, k_j \rangle v_j$$



Local and fixed during inference

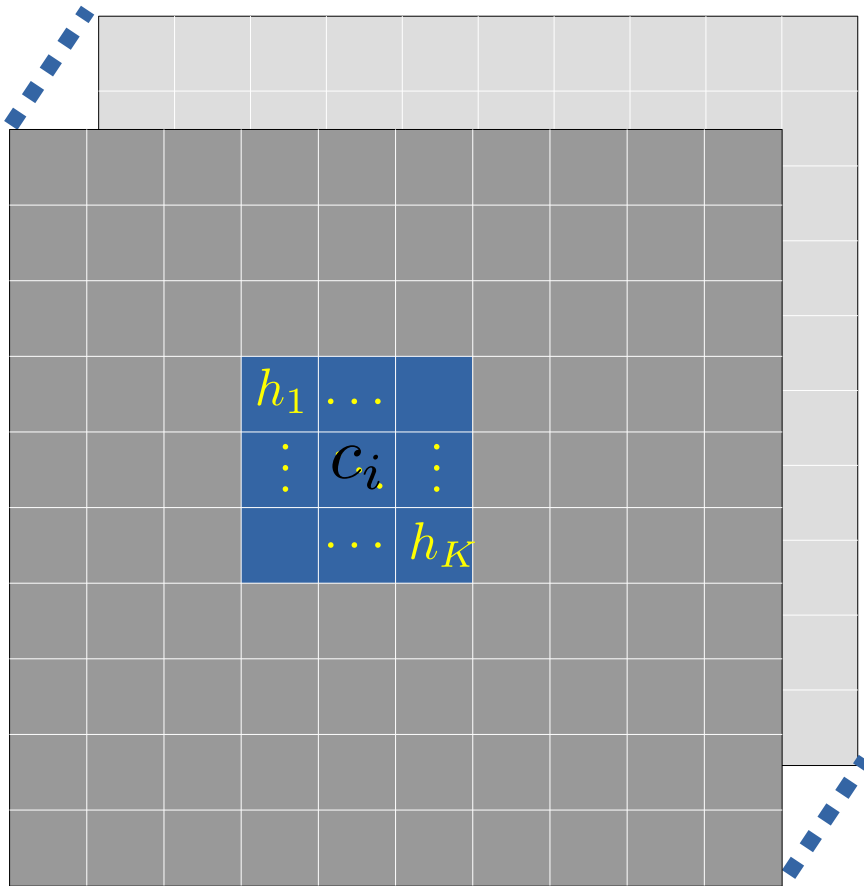




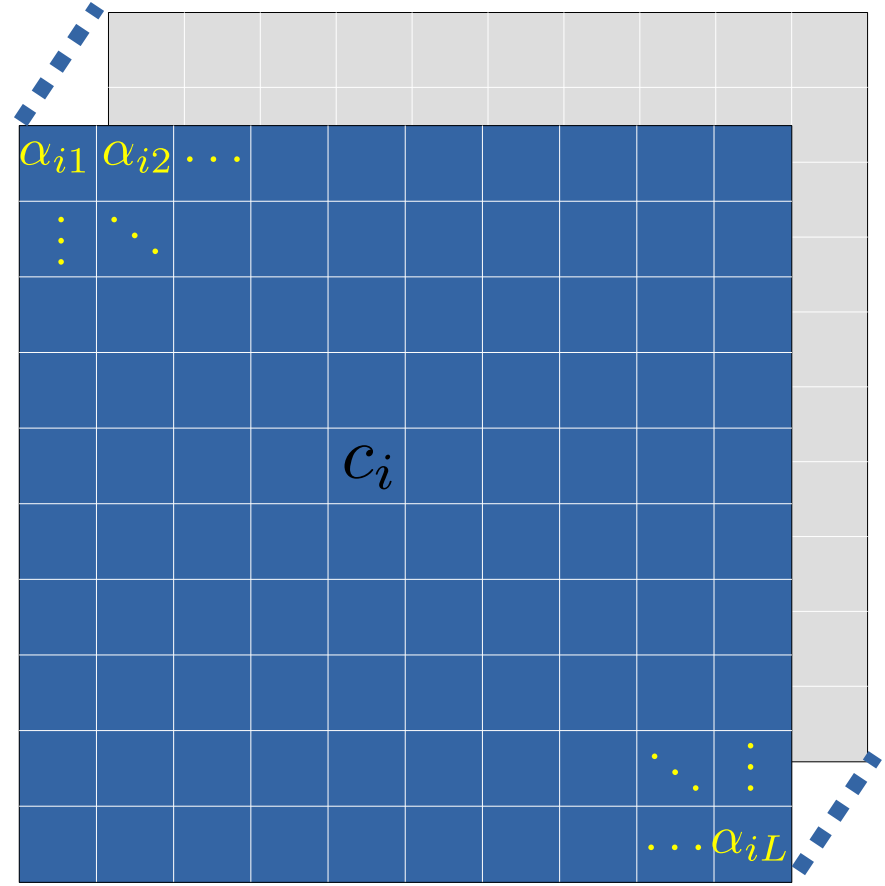
# Convolution vs Self-Attention

$$c_i = \sum_j h_{i-j} x_j$$

$$c_i = \sum_j \langle q_i, k_j \rangle v_j$$





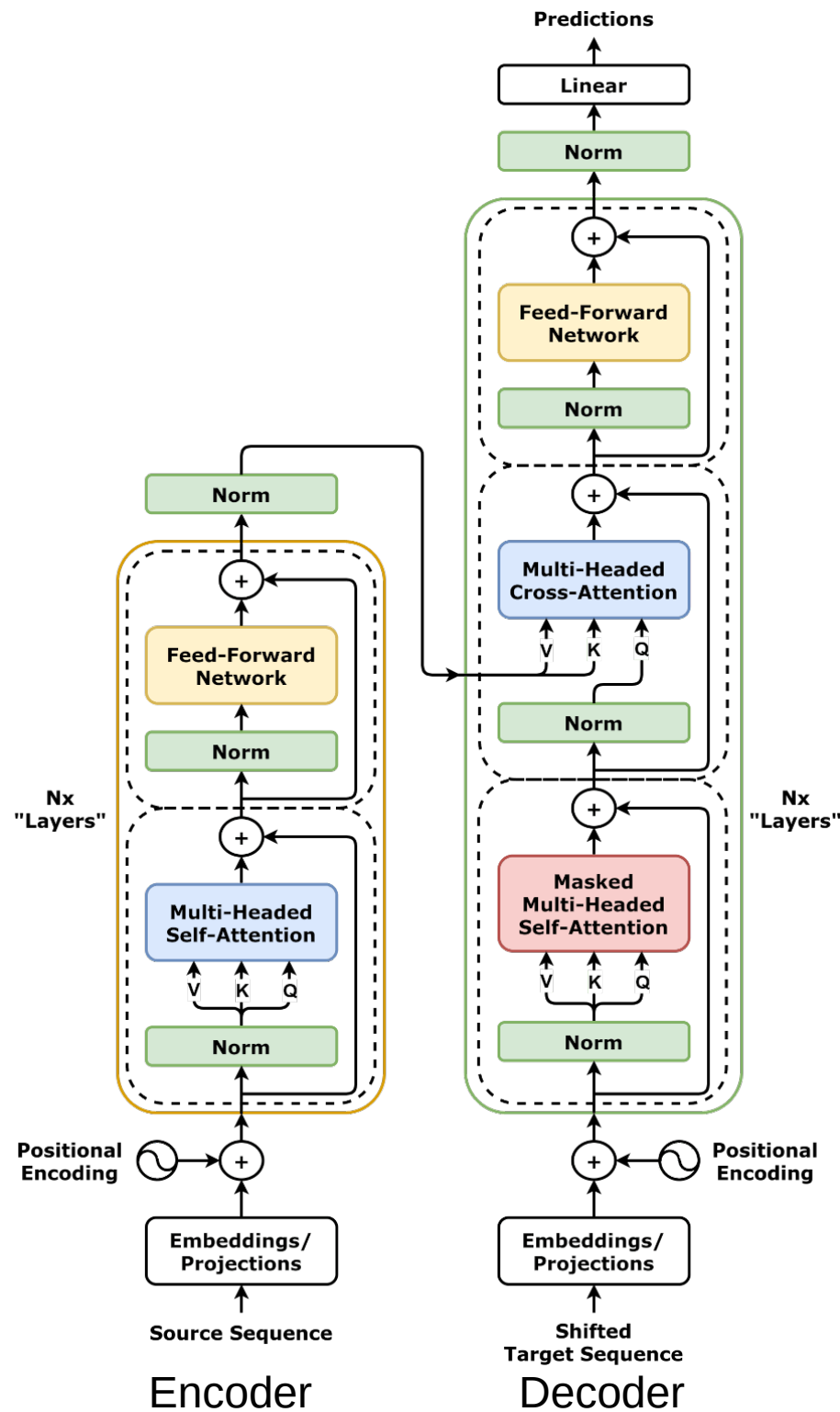
Local and fixed during inference



Global and variable during inference

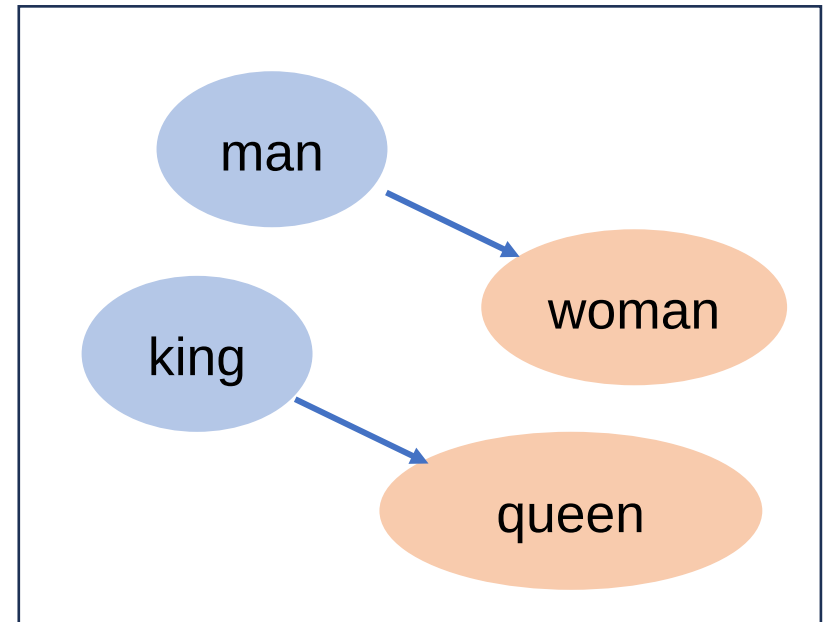
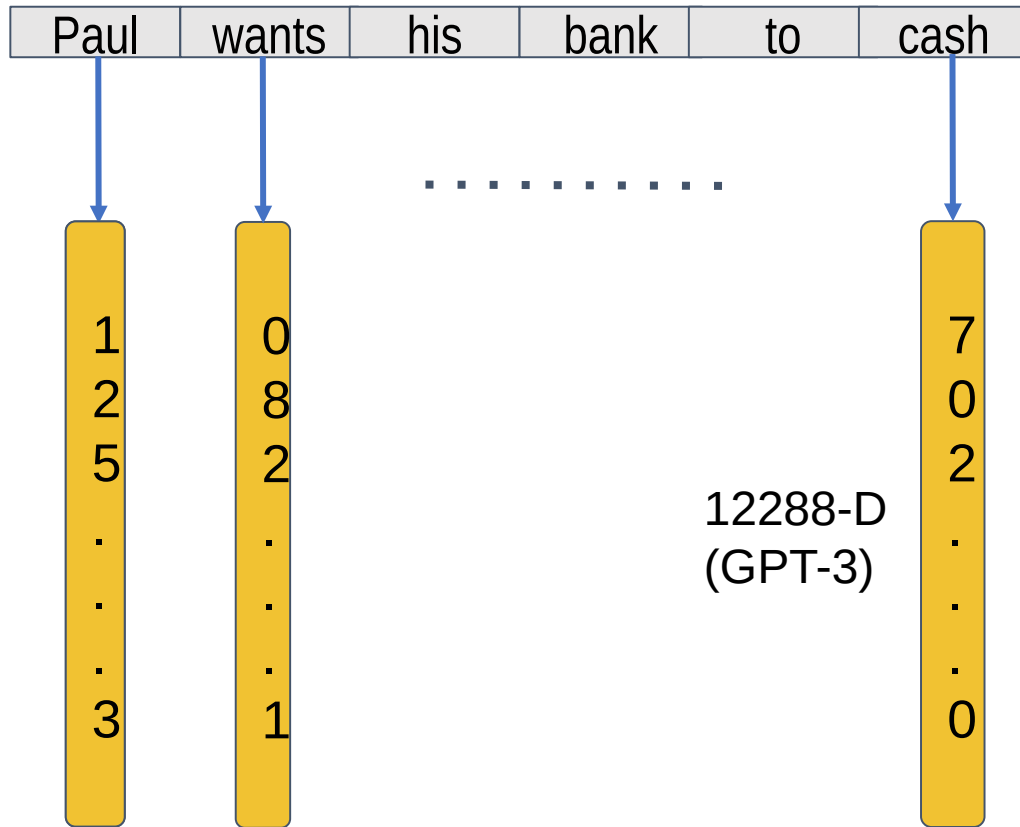
# The Transformer

- Multi-Headed Self- and Cross-Attention
- Masked Multi-Headed Self-Attention
- Layer Normalization
- Linear + ReLU
- Positional Encoding 
- Residual Connection 
- Dropouts





# Embedding



Byte pair encoding → tokens

50257 tokens in GPT-3

king – man + woman = queen

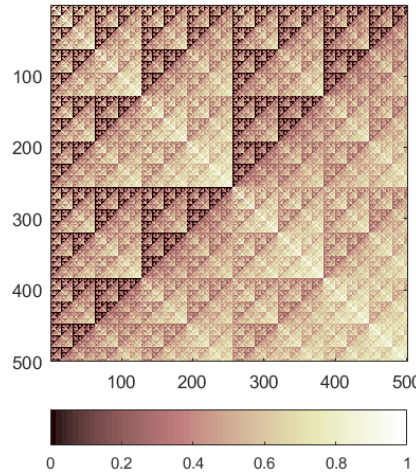
Mikolov et al., *Word2Vec*, 2013



# Positional Encoding

- **Unique** encoding for each position
- Encoding distance between two positions **consistent**
- **Generalize** to any sequence length

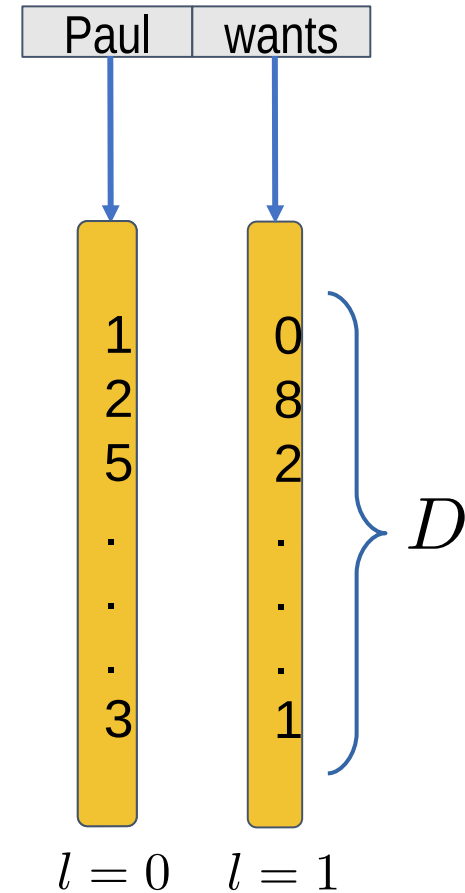
- Binary code?



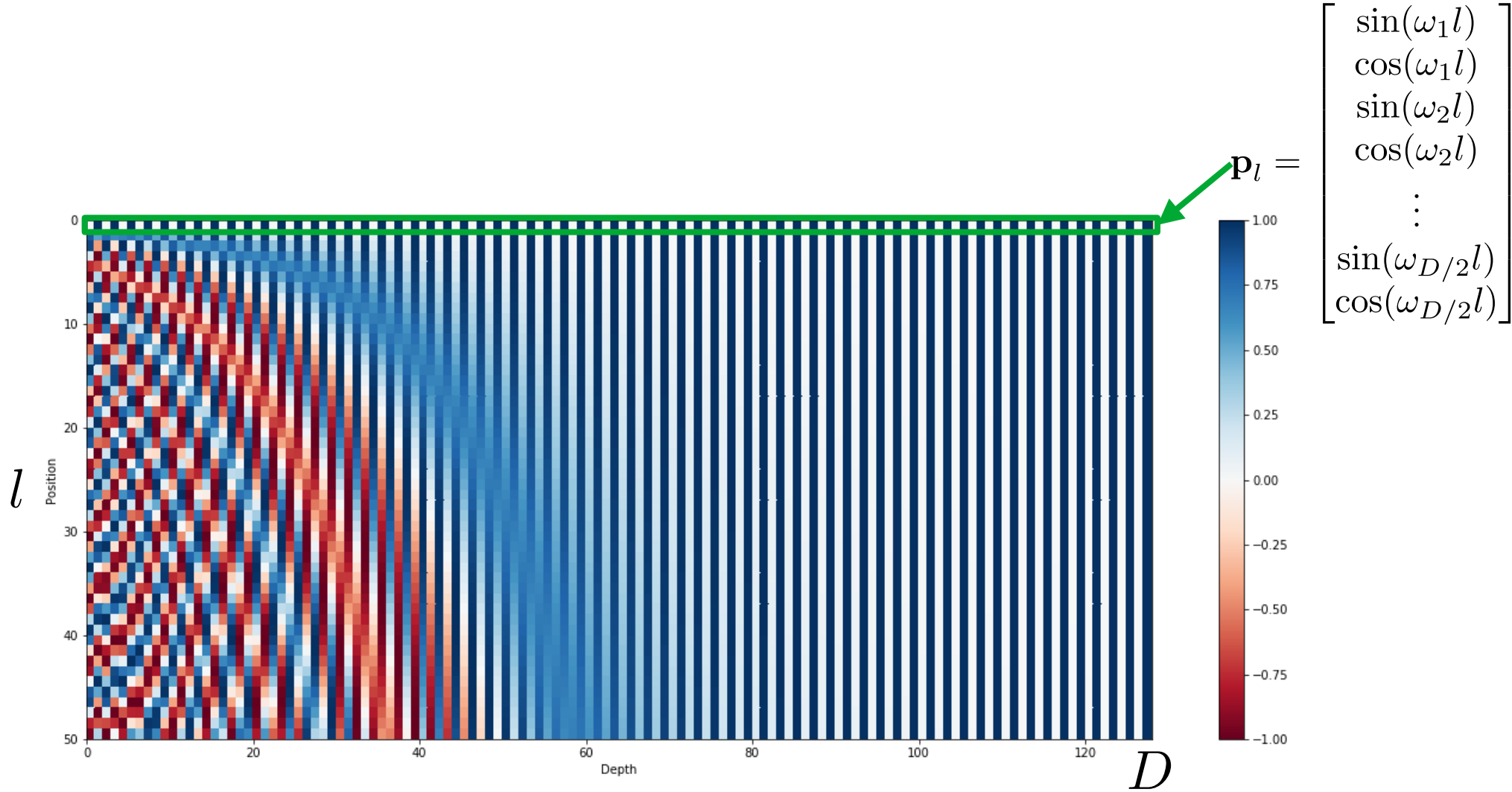
$$\omega_k = \frac{1}{10000^{2k/D}}$$

$$\mathbf{p}_l =$$

$$\begin{bmatrix} \sin(\omega_1 l) \\ \cos(\omega_1 l) \\ \sin(\omega_2 l) \\ \cos(\omega_2 l) \\ \vdots \\ \sin(\omega_{D/2} l) \\ \cos(\omega_{D/2} l) \end{bmatrix}$$



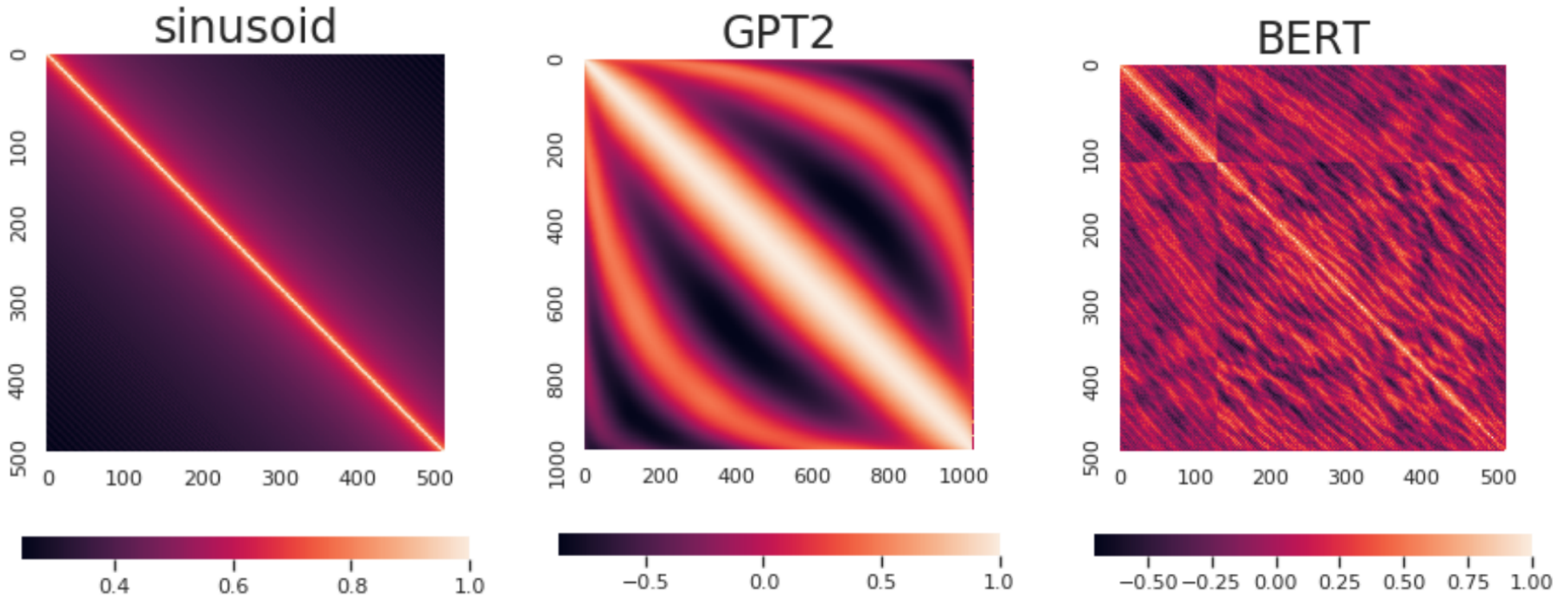
# Handcrafted Positional Encoding





# Learned Positional Encoding

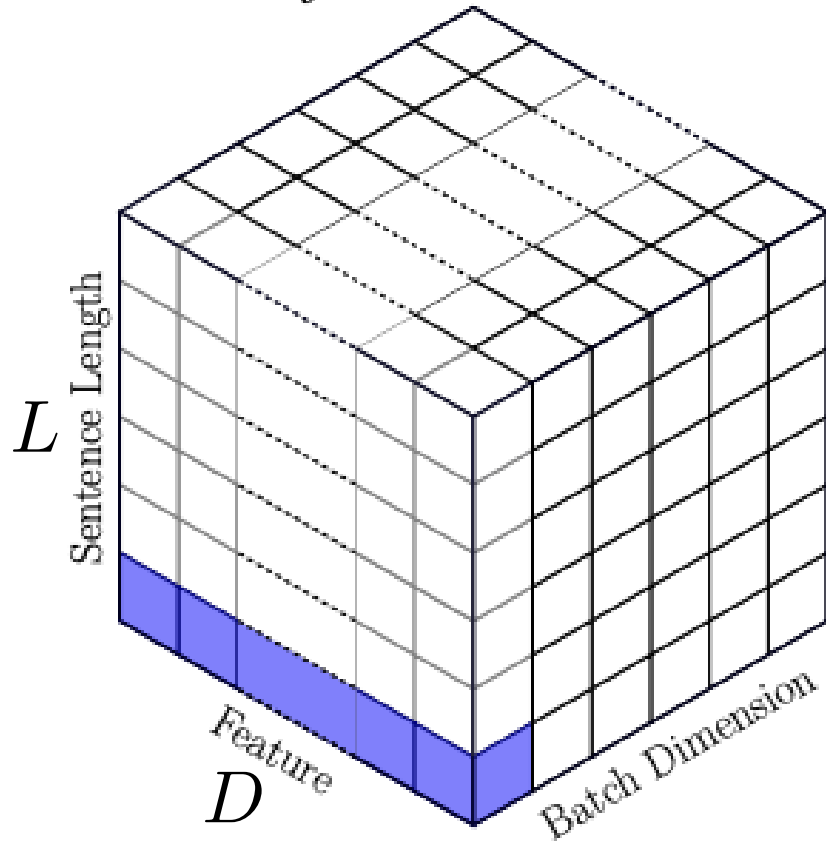
- PE Correlation between different positions



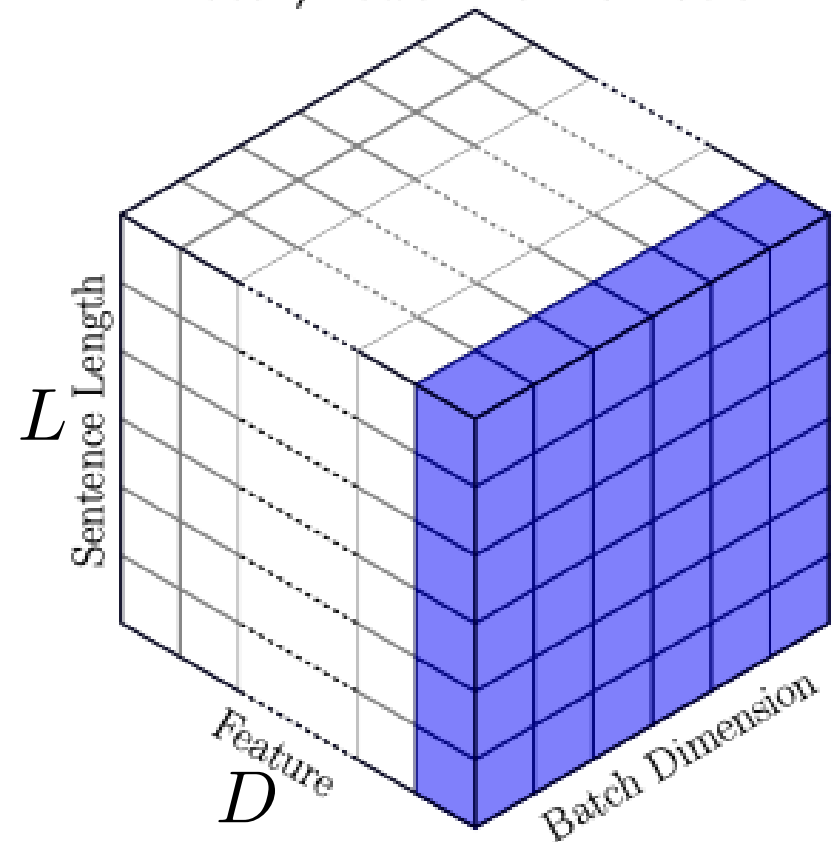


# Layer Normalization

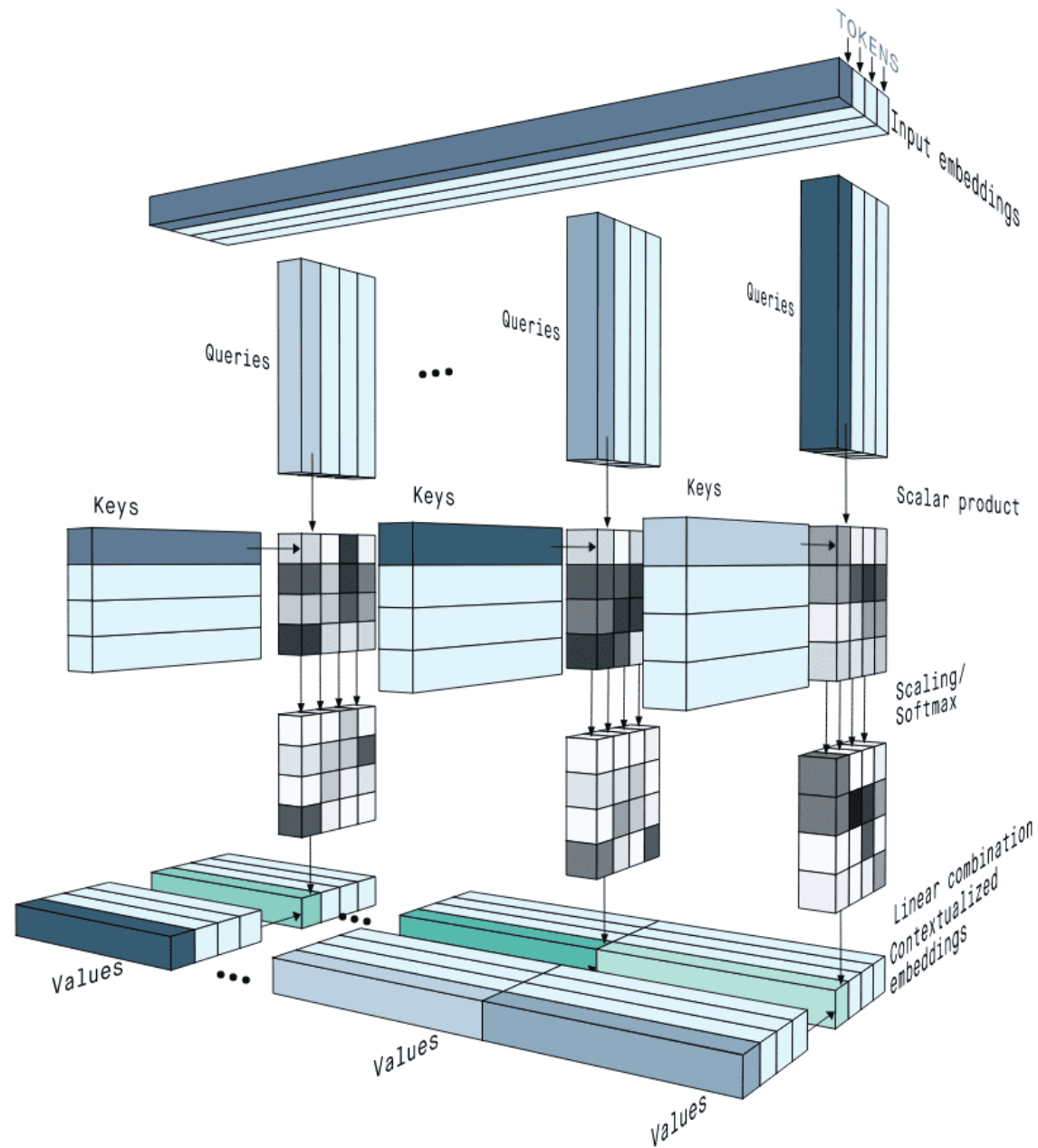
Layer Normalization



Batch/Power Normalization



# Multi-headed Self-Attention









# Large Language Models

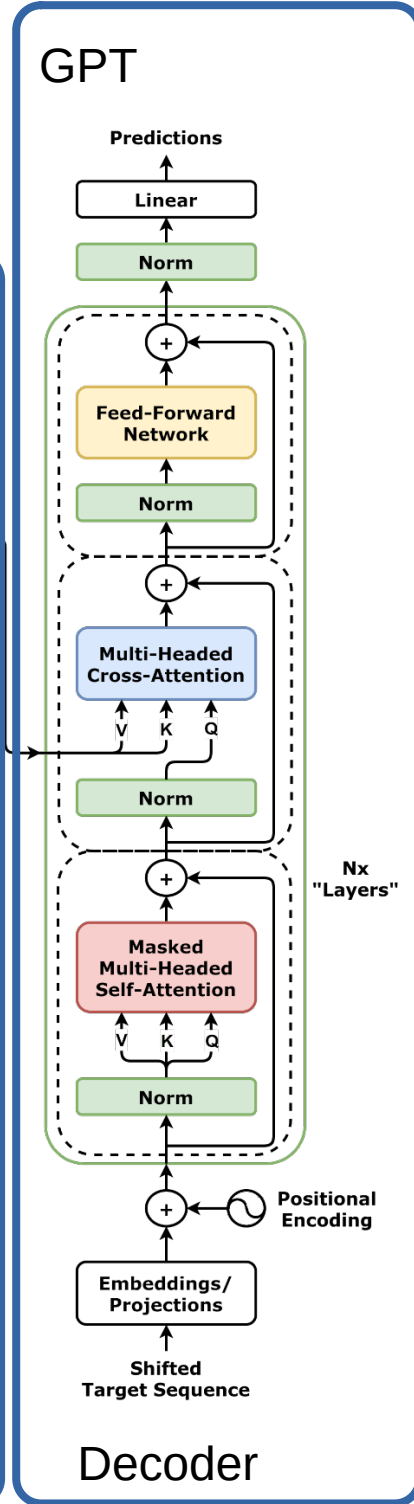
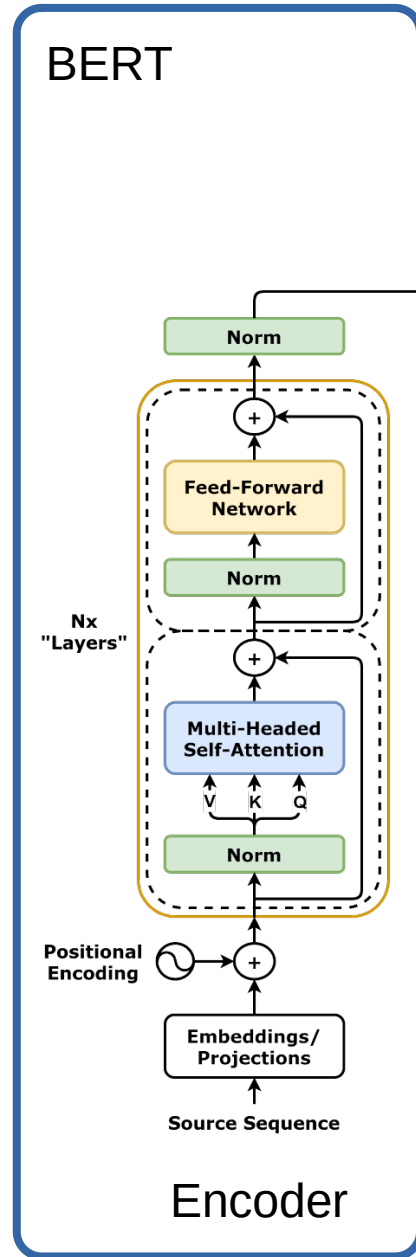
## 1) Pretraining

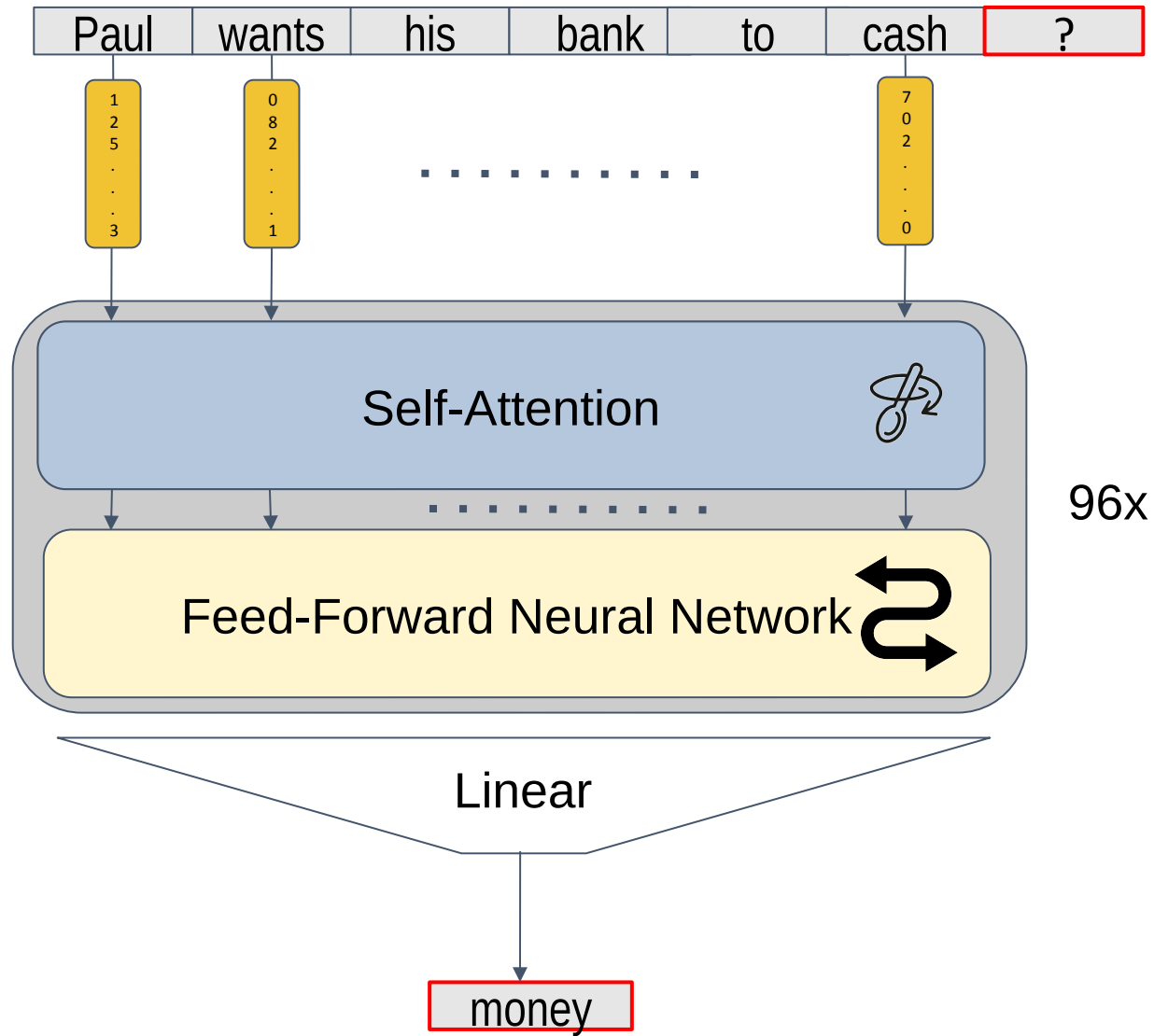
- Annotation is a curse
- Task - “Predict a next (masked) word in an incomplete sentence.”

## 2) Fine-tuning on a downstream task

- Small annotation dataset
- Reinforcement learning from Human Feedback
- Immitate human preferences - Reward Model

GPT (OpenAI) , BERT (Google),  
LLaMA (Meta AI), Titan (Amazon)







Self-Attention 

- Understanding, relations, reasoning

Paul	wants	his	bank	to	cash
------	-------	-----	------	----	------

Self-Attention 

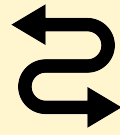
			financial		
	Paul	Paul	institution		
Paul	wants	his	bank	to	cash

Self-Attention 

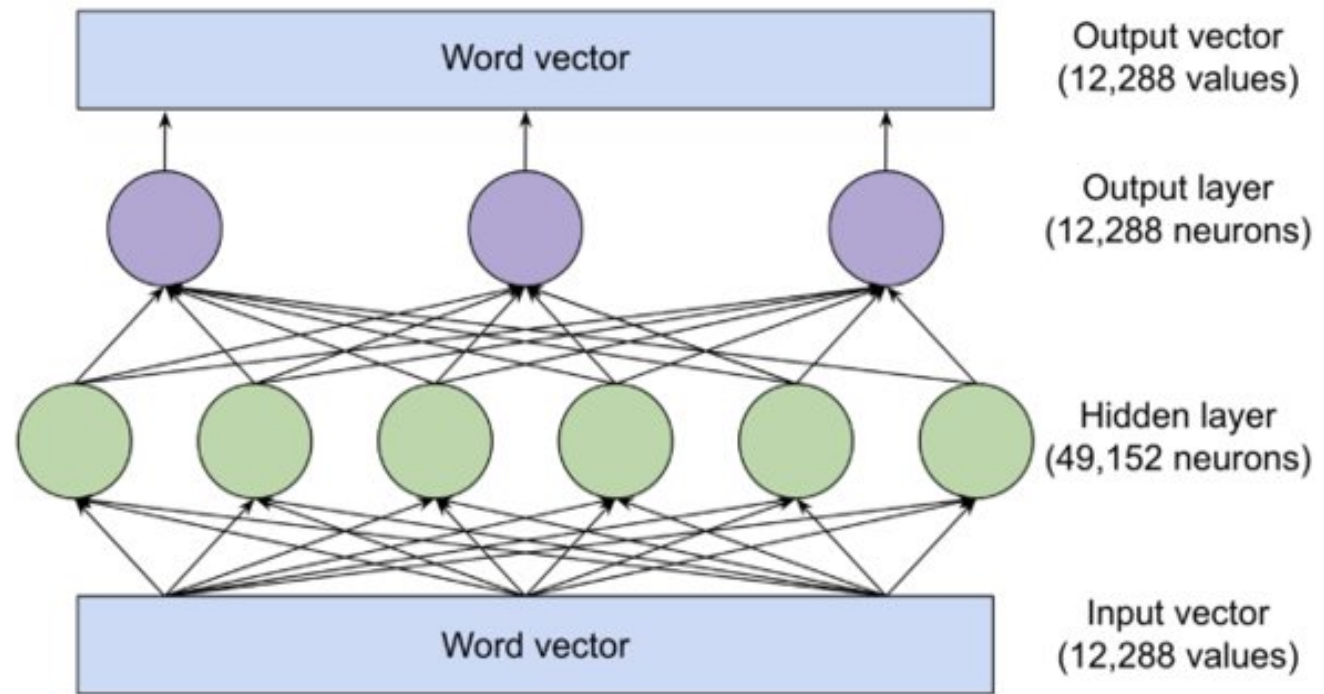
	Paul	Paul	financial		
			institution		
				bank	
				transaction	
Paul	wants	his	bank	to	cash

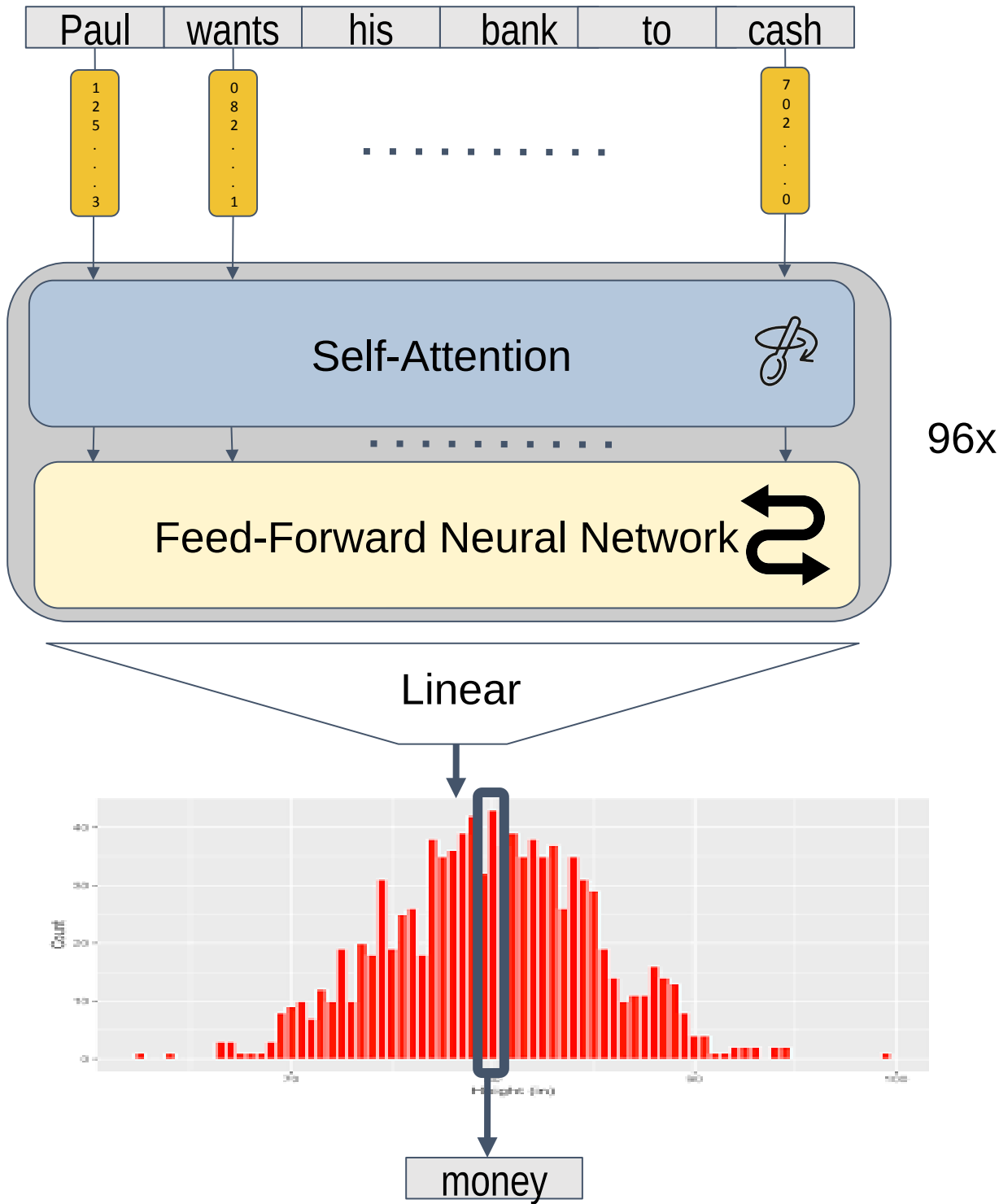
Paul	wants	his	bank	to	cash
------	-------	-----	------	----	------

# Feed-Forward Neural Network



- Long-term Memory – facts from the training data
- GPT-3:
  - 175 (x10<sup>9</sup>) billion of parameters
- GPT-4:
  - 1 700 billion
- Human brain:
  - 700 000 billion





Non-deterministic!

Mode

Model

Temperature

Maximum length 2050

Stop sequences  
Enter sequence and press Tab

Top P


Frequency penalty 0


Presence penalty 0

# Context in Prompting

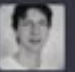



Self-Attention can infer opposite implication using the context.

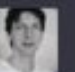
 Představ si, že Tom má matku Mary. Jak se jmenuje syn Mary?


 Syn Mary by se jmenoval Tom, protože v úvodu tvé otázky jsi řekl, že Tom má matku Mary.

Long-term memory in FF does not have this ability.

 Jak se jmenuje syn Mary Lee Pfeiffer?

 Bohužel nemám informace o synovi Mary Lee Pfeiffer, protože moje znalosti sahají pouze do září 2021, a nemám aktuální informace o soukromých osobách a jejich rodinných vztazích. Pro aktuální informace o této osobě byste mohli zkusit prověřit dostupné veřejné zdroje nebo vyhledat aktuální zpravodajské články.

 Jak se jmenuje matka Toma Cruise?

 Matka Toma Cruise se jmenuje Mary Lee Pfeiffer.



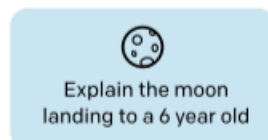


# Training with Human Feedback

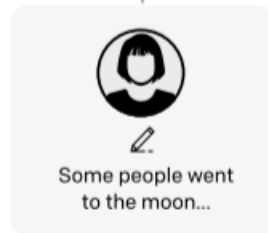
Step 1

**Collect demonstration data, and train a supervised policy.**

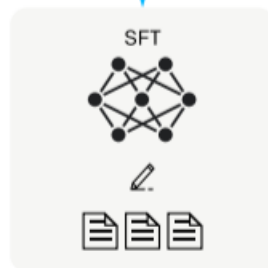
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



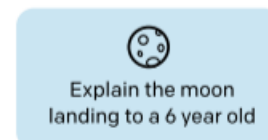
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

**Collect comparison data, and train a reward model.**

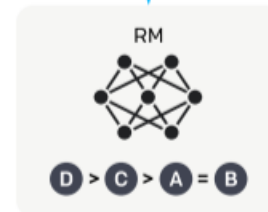
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



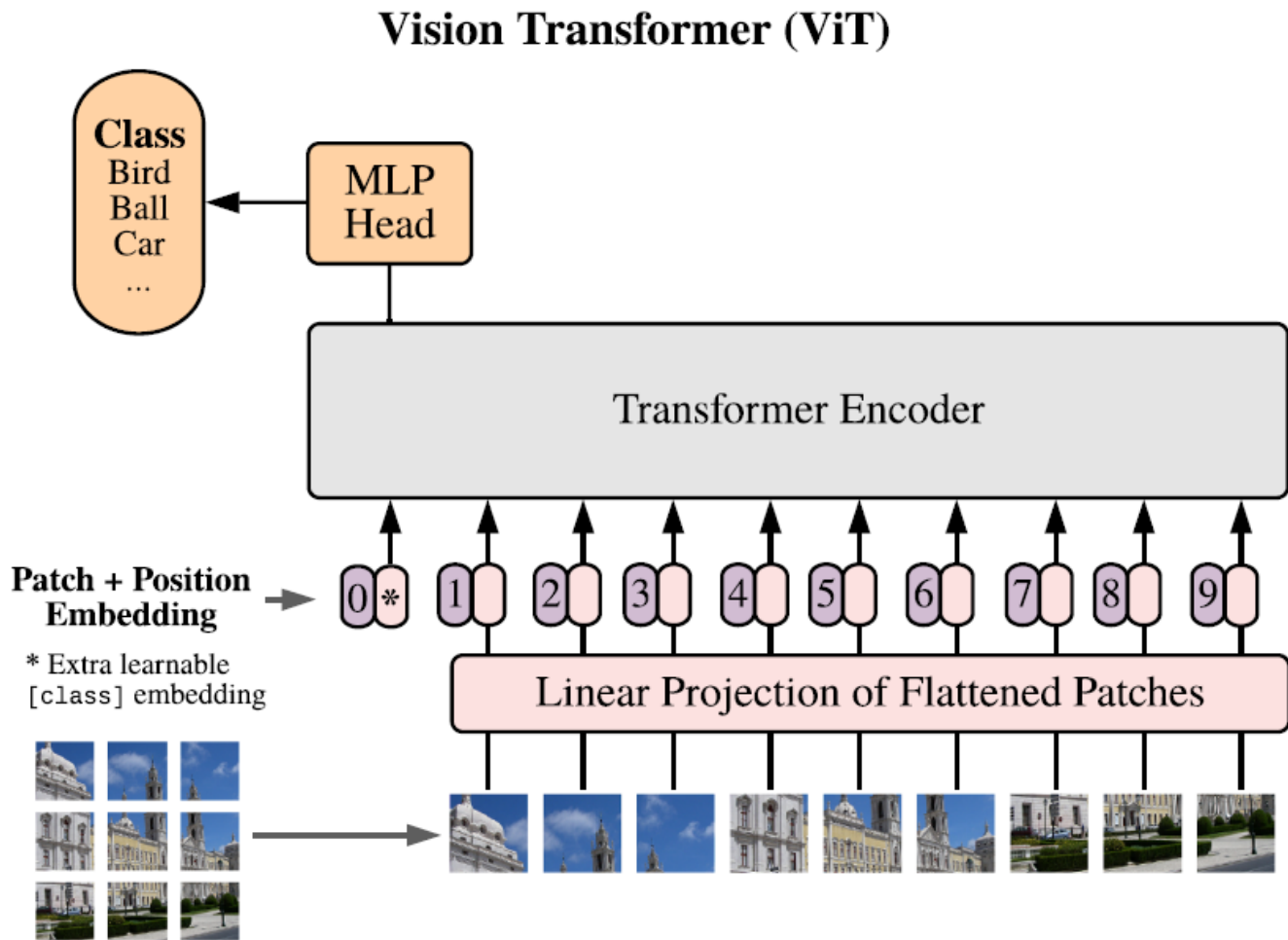


# Attention in Image Processing

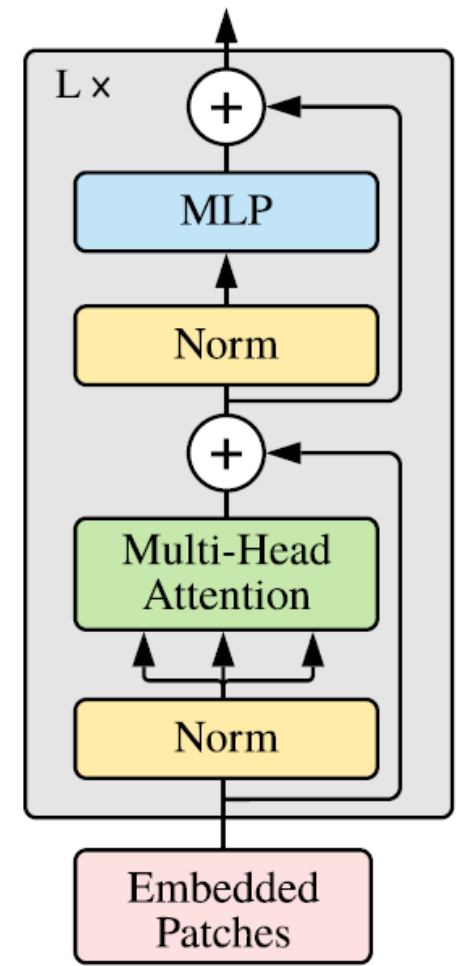
- Image Classification & Detection:
  - **CBAM** (Convolutional Block Attention Module)
  - **Dual Attention** (Spatial and Channel)
  - **ViT** (Vision Transformer)
  - **CoAtNet** (Convolution with Self-Attention)



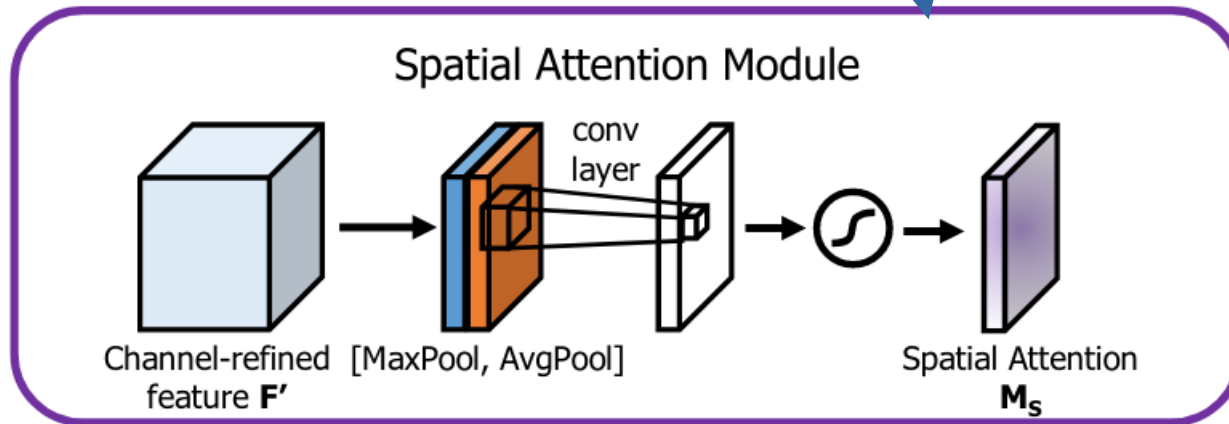
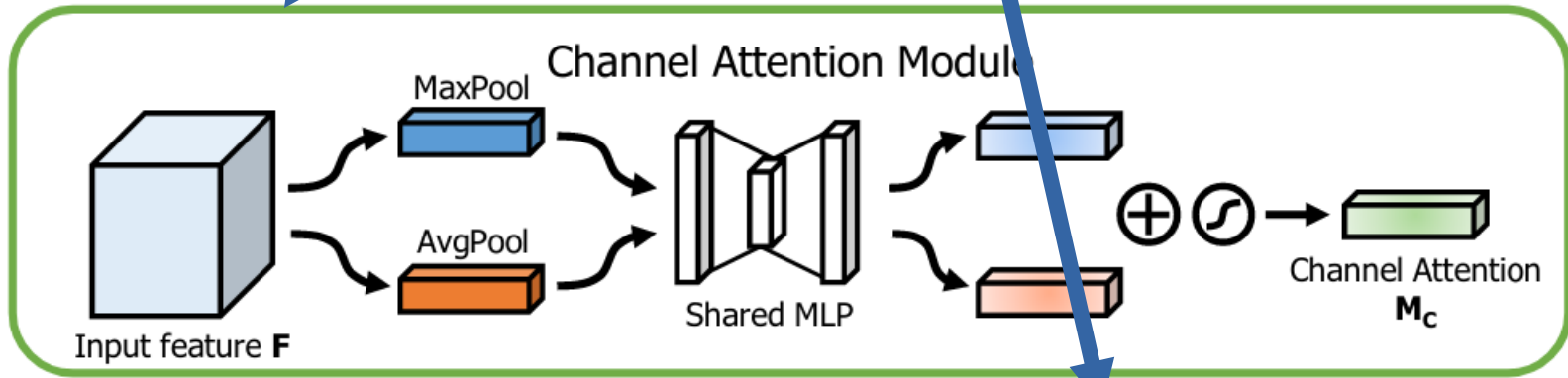
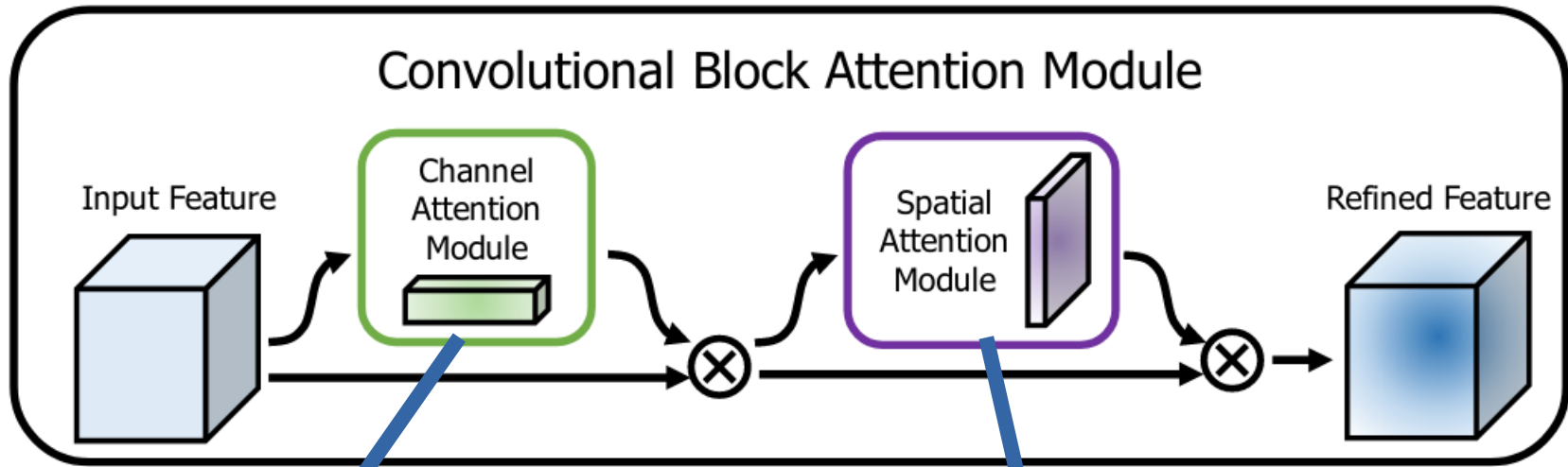
# Vision Transformer (ViT)



## Transformer Encoder

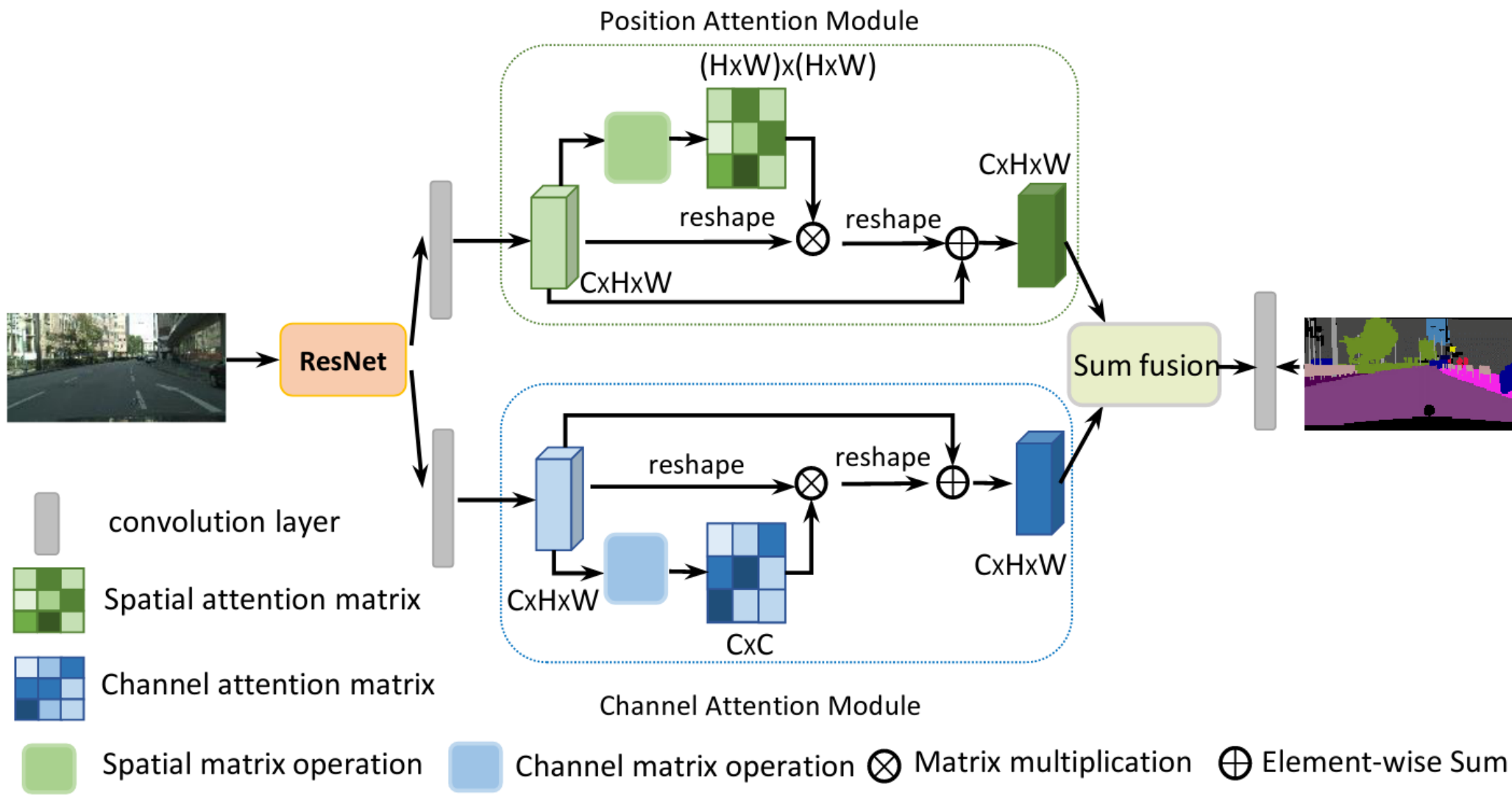


# Convolutional Block Attention Module



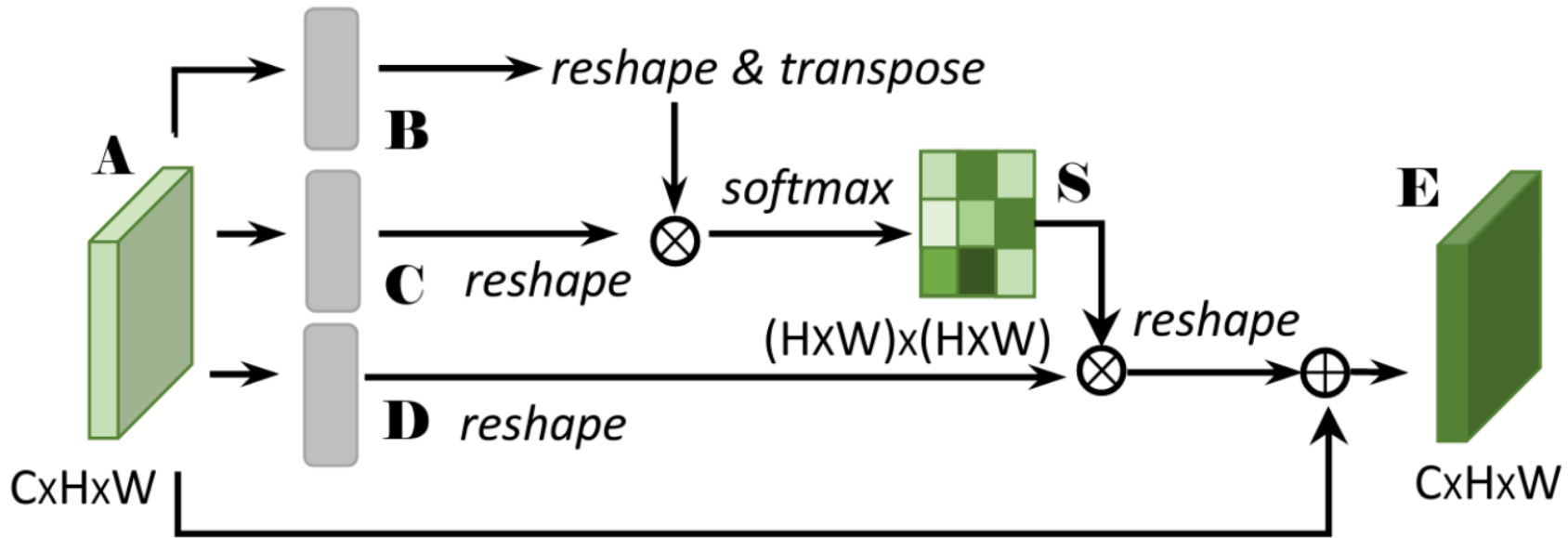


# Dual Attention

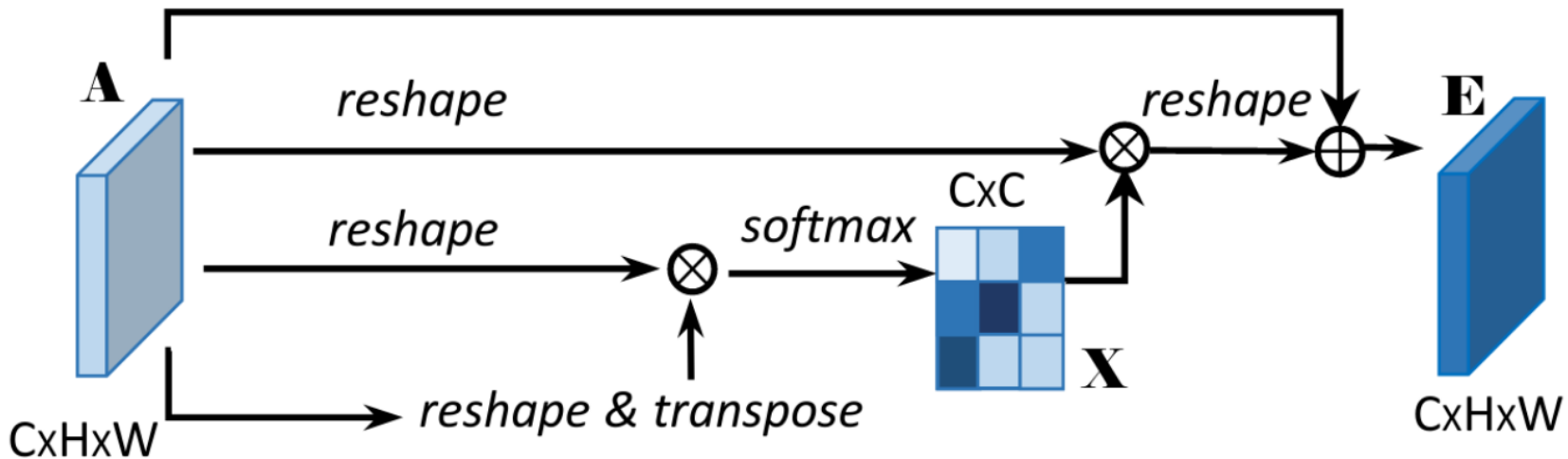




# Dual Attention



A. Position attention module



B. Channel attention module

# Self-supervised training

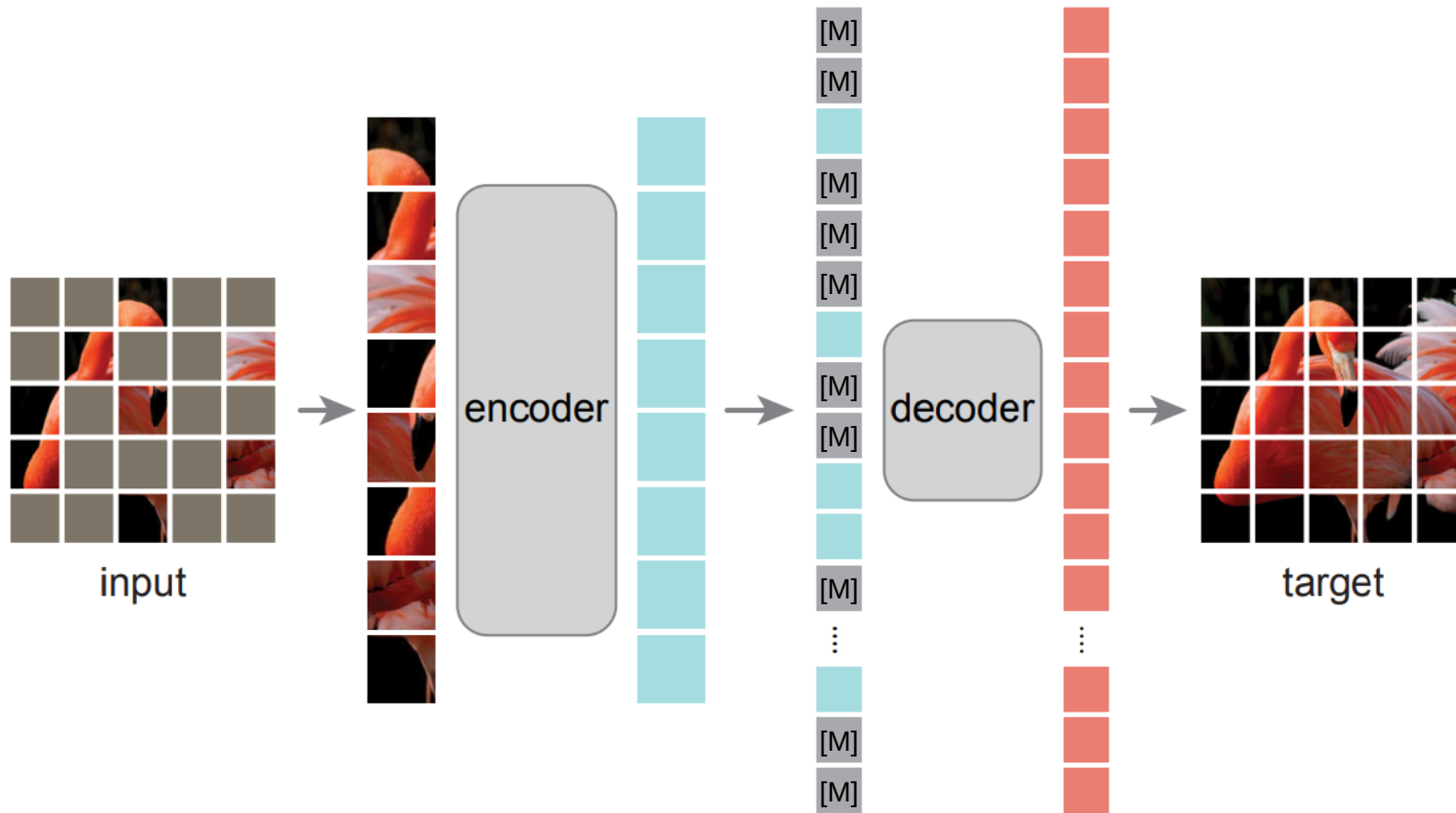


- ViT requires a lot of training data
- Pretraining + fine-tuning
- Pretraining methods:
  - Mask Autoencoder (MAE)
  - Self-distillation with no labels (DINO)



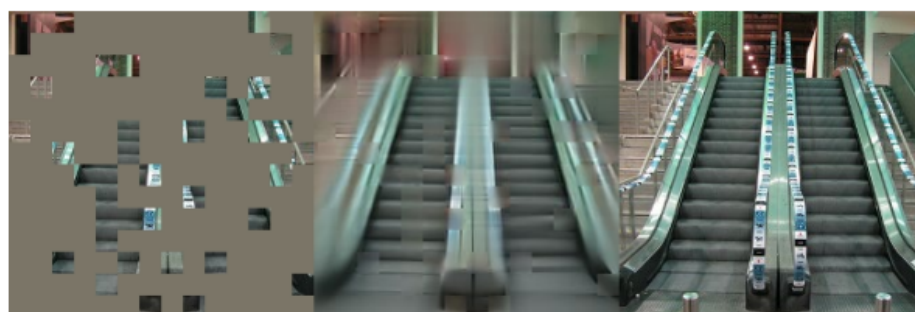
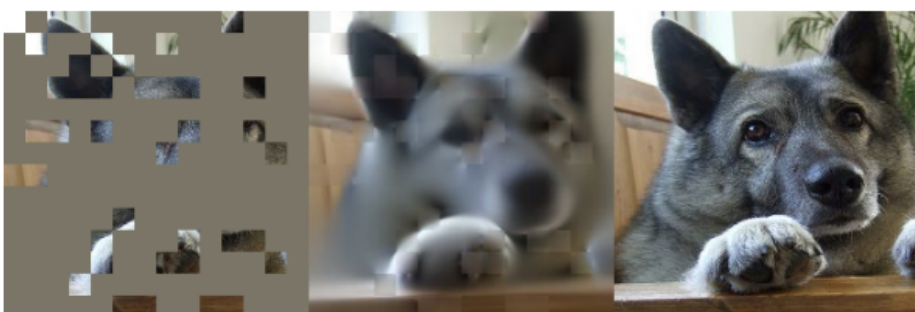
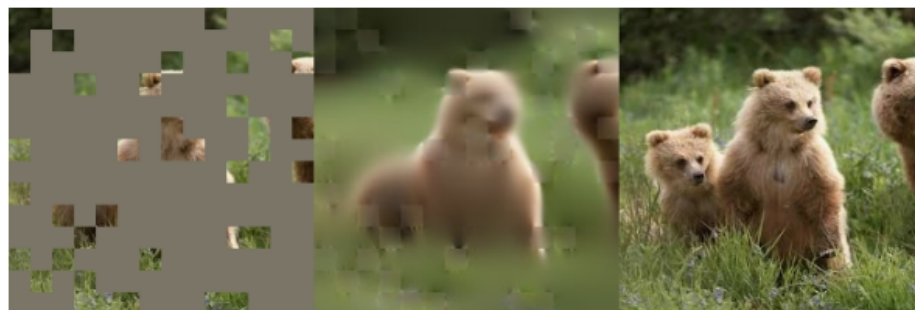
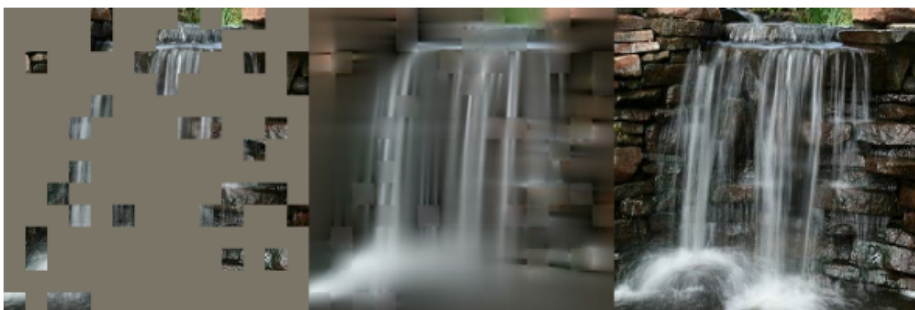
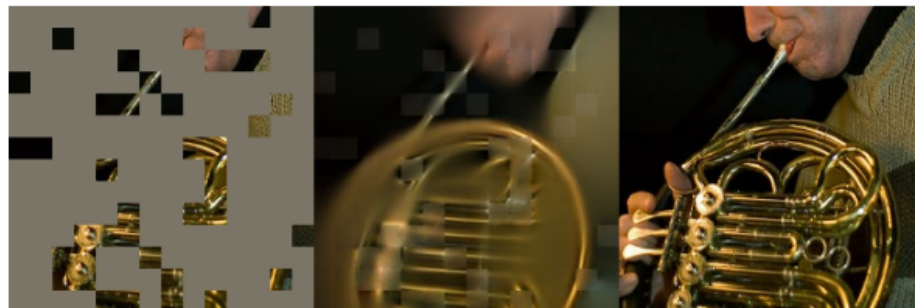
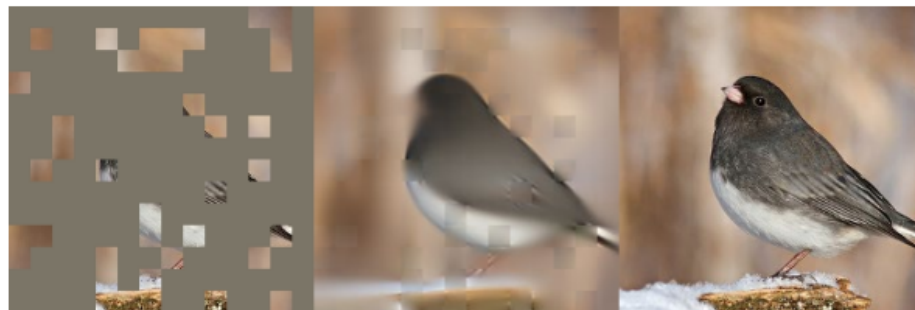
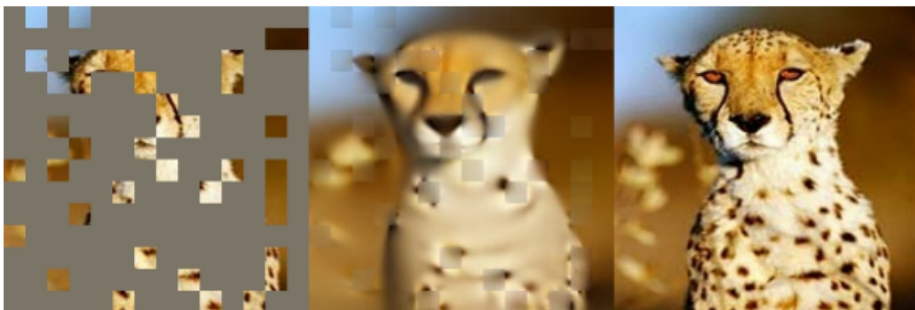
# MAE

- Similar to BART pretraining





# MAE



# DINO



- Several local views
  - small areas of the original image (0.05 - 0.4)
  - 96x96 pixels
  - For the student
- Two global views
  - large area of the original image (0.4 – 1.0)
  - 224x224 pixels
  - Only for the teacher



Student



Teacher



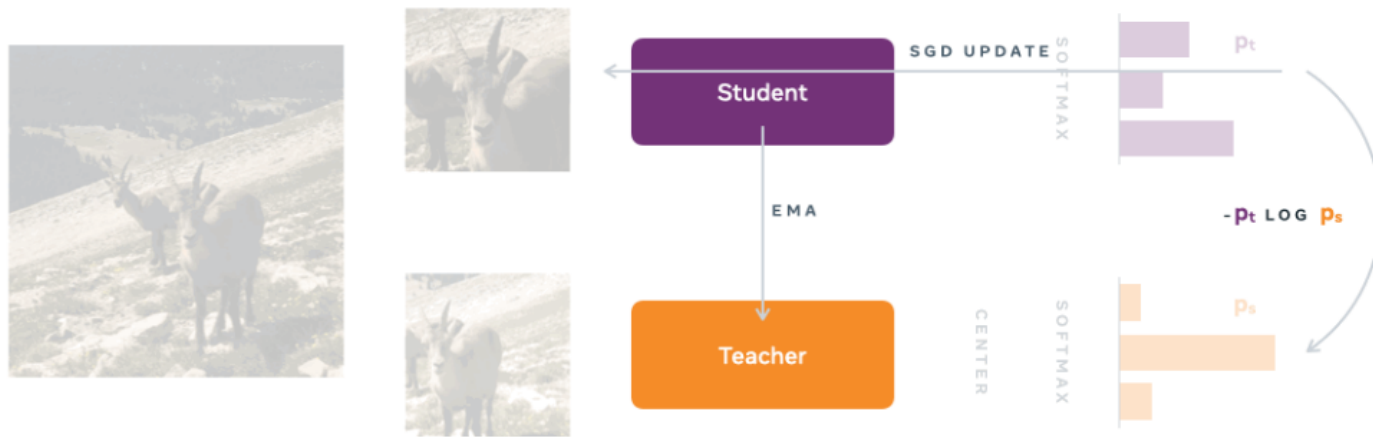


Student



Teacher





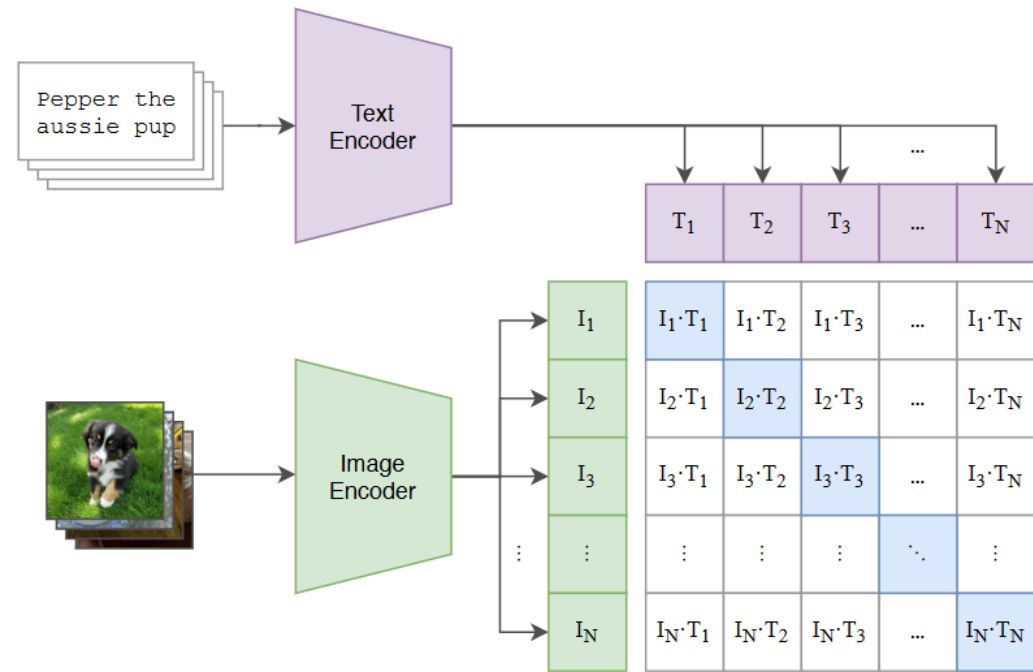
Avoid collapse: centering and sharpening



# CLIP

- Contrastive Language-Image Pretraining

(1) Contrastive pre-training



(2) Create dataset classifier from label text

