

A Survey of LLMs

Zhao et. al

Anthony J. Dsouza

Date - 23.09.2024

Table of Contents

1. Scaling & Need for scaling
2. The GPT family
3. Data preparation
4. Model architecture

Table of Contents - Scaling and Need for scaling

1. Why do we need scaling?
2. The scaling laws
3. Emergent abilities & scaling laws

What are LLMs?

What are LLMs?

**What factors to consider
when calling an LM
“Large”?**



What are LLMs?

**What factors to consider
when calling an LM
“Large”?**

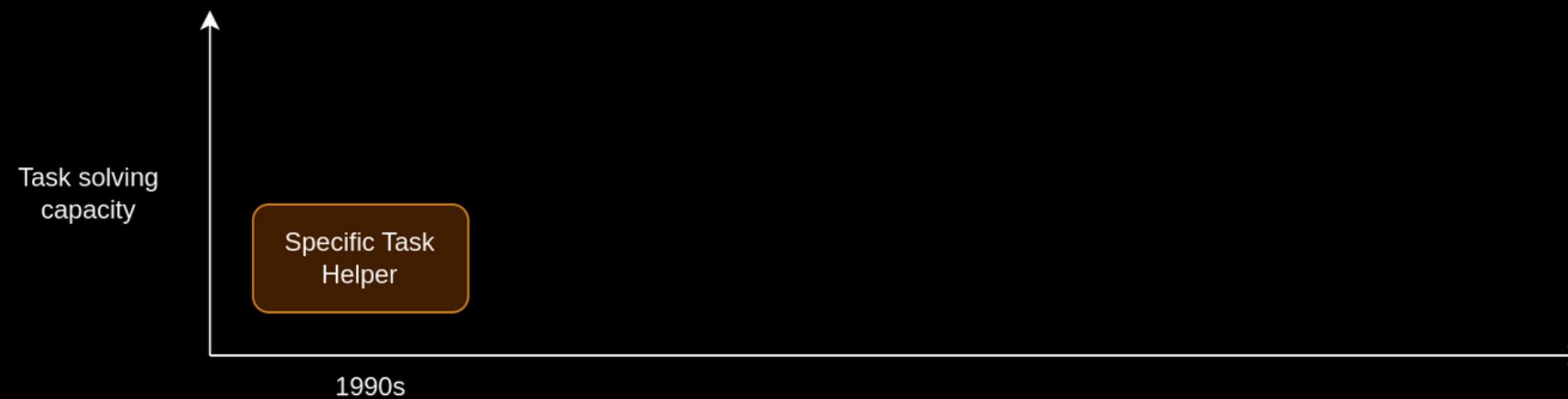
- Size

What are LLMs?

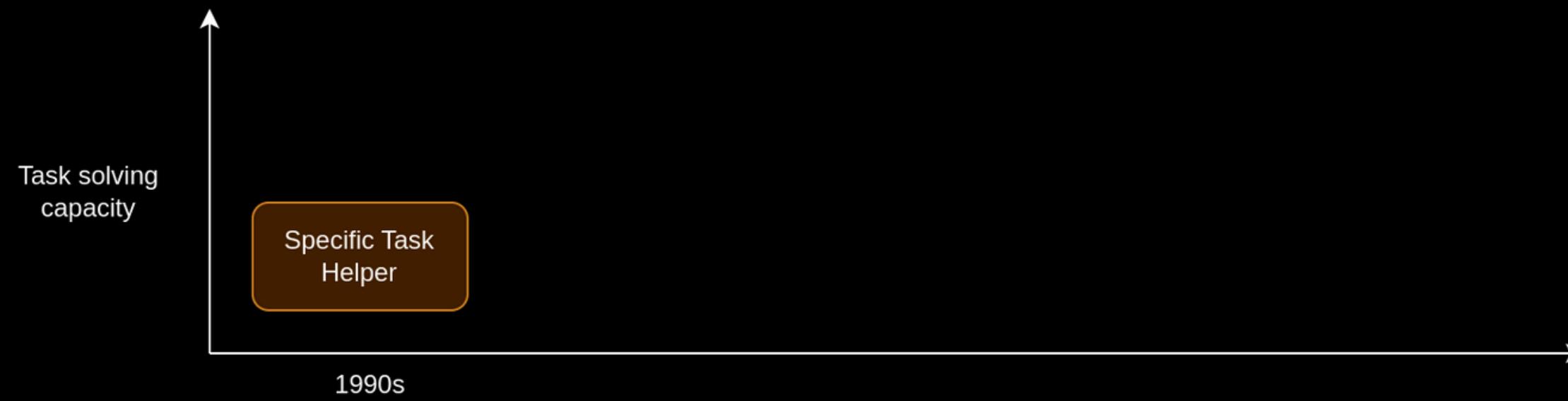
- **Size**
- **Abilities**

**What factors to consider
when calling an LM
“Large”?**

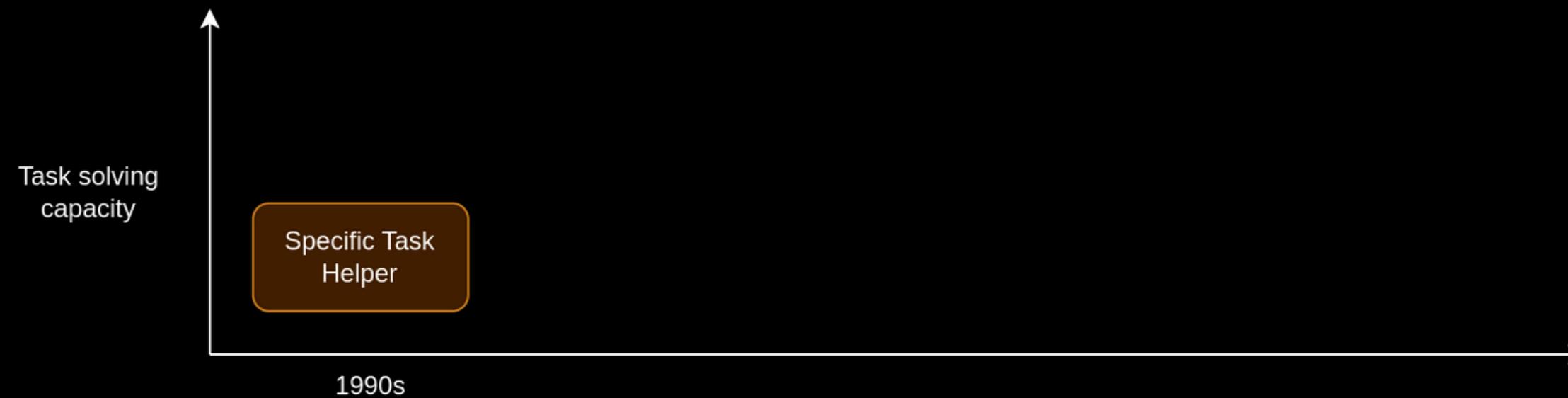
ngram models
based on probability estimation
based on _____ assumption



ngram models
based on probability estimation
based on markov assumption



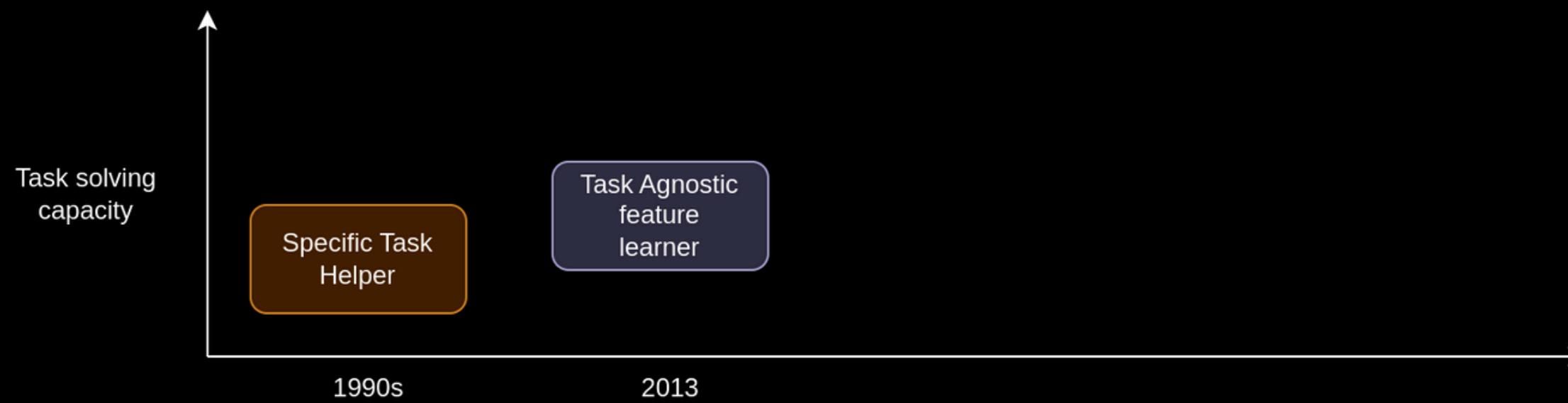
ngram models
based on probability estimation
based on markov assumption
suffer from curse of dimensionality - but why?



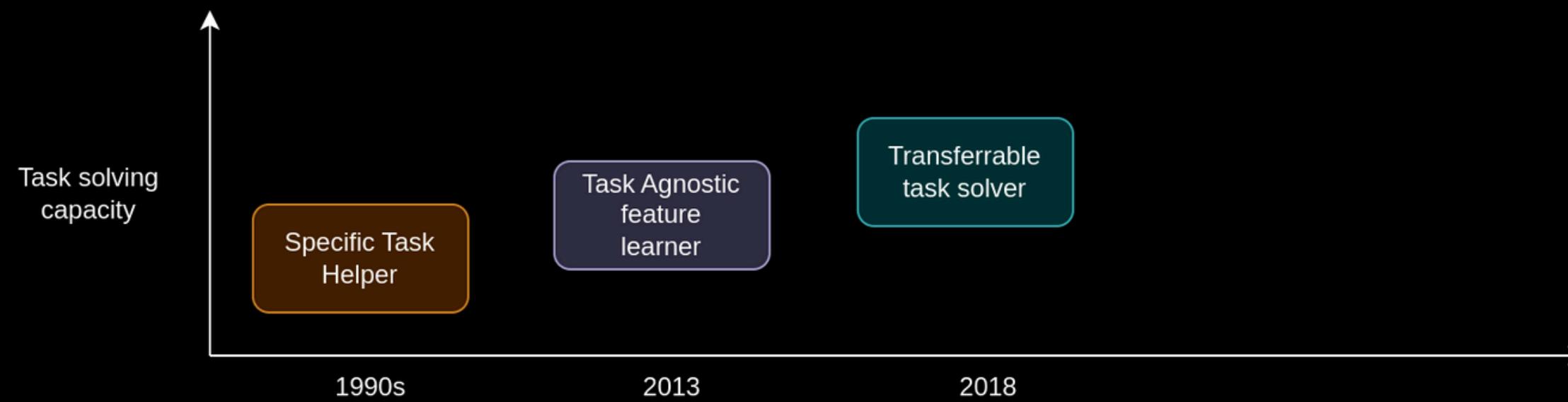
ngram models
based on probability estimation
based on markov assumption
suffer from curse of dimensionality - difficult to model all transition probabilities



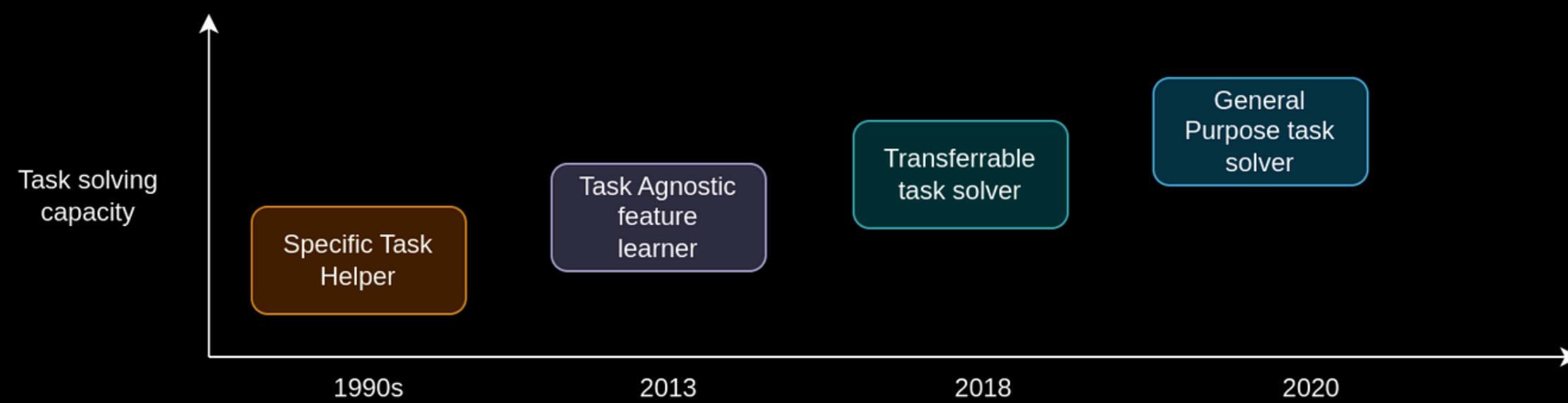
Static word representations
Neural context modelling
word2vec / GloVe
Distributed representation of words
Predicting words based upon surrounding words



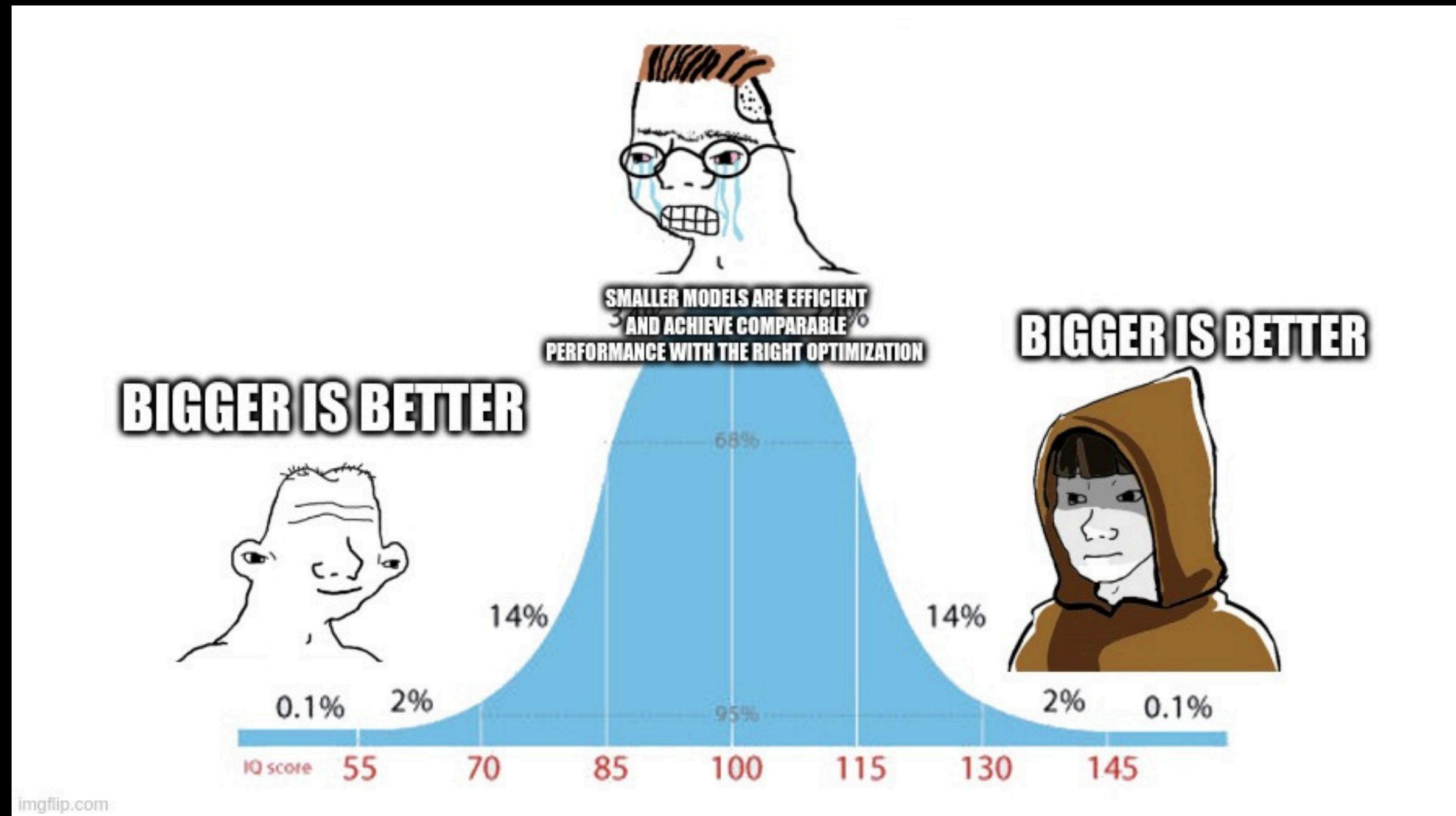
Context aware representations
pretrain on large corpora
finetune on task specific corpora



Scale PLM
Use diverse data

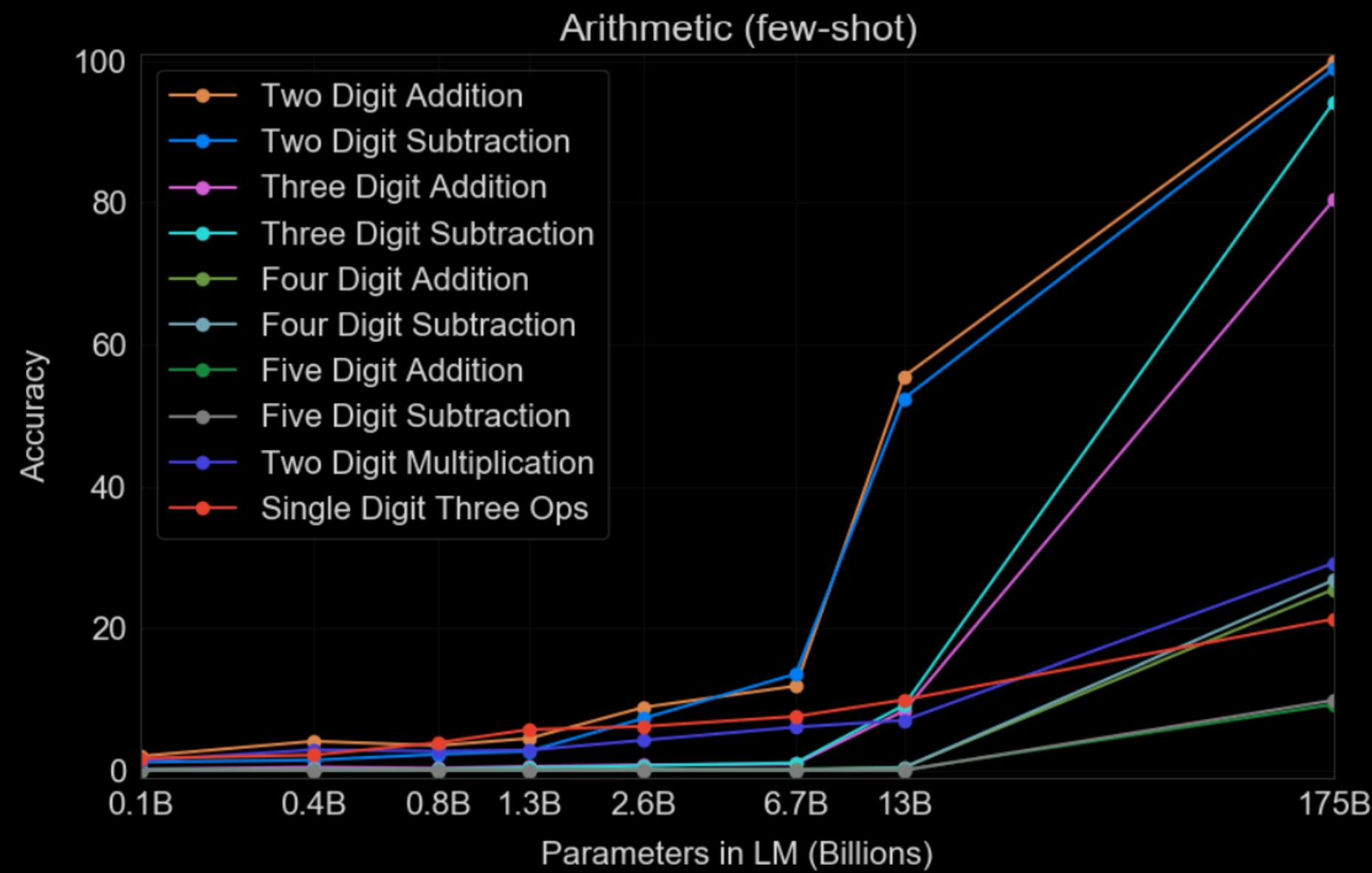


Scaling Laws



Scaling Laws

Why do we need scaling laws?



KM Scaling Law

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha N}, \alpha N \sim 0.076, N_c \sim 8.8 \times 10^{23}$$

KM Scaling Law

$$L(D) = \left(\frac{D_c}{D}\right)^{\alpha D}, \alpha D \sim 0.095, D_c \sim 5.4 \times 10^{13}$$

$$L(C) = \left(\frac{C_c}{C}\right)^{\alpha C}, \alpha C \sim 0.050, C_c \sim 3.1 \times 10^8$$

$$L(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

Chinchilla Scaling Law

$$N_{opt}(C) = G \left(\frac{C}{6} \right)^a, D_{opt}(C) = G^{-1} \left(\frac{C}{6} \right)^b$$

So, which one would you use?

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

1. Predictable Scaling

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

1. Predictable Scaling

2. Task-level predictability

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

1. Predictable Scaling

2. Task-level predictability

3. Emergent Abilities

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

1. Predictable Scaling

2. Task-level predictability

3. Emergent Abilities

a. In-Context Learning

Scaling Laws

**So, back to the question, why
do we need scaling laws?**

1. Predictable Scaling

2. Task-level predictability

3. Emergent Abilities

a. In-Context Learning

b. Instruction Following

Scaling Laws

So, back to the question, why
do we need scaling laws?

1. Predictable Scaling

2. Task-level predictability

3. Emergent Abilities

a. In-Context Learning - @Nicholas

b. Instruction Following

c. Step-by-step reasoning

In-Context Learning

Please see the following text and annotation. You need to identify the first name, last name and age from the text. the format of data is provided to you

{

- 1: I am Eric Dawson, 31M, living in Miami,
- 2: I'm Fiona Gonzales, living in SF. I am 24 years old.

}

Your response should be like the following:

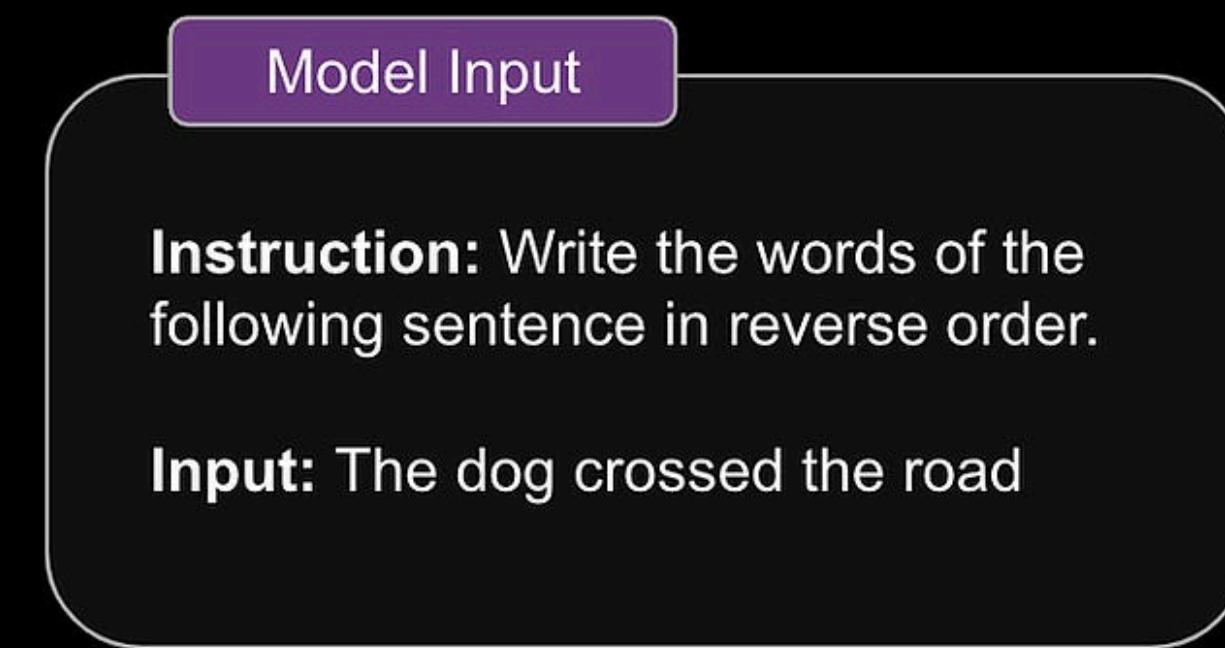
{

- 1: I am Eric [FN] Dawson [LN], 31 [age] M, living in Miami,
- 2: I'm Fiona [FN] Gonzales [LN], living in SF. I am 24 [age] years old.

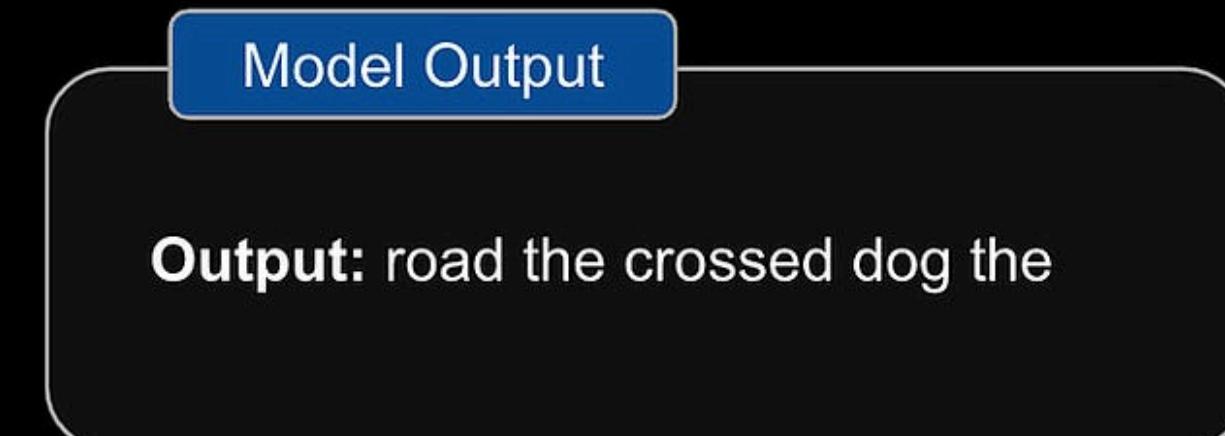
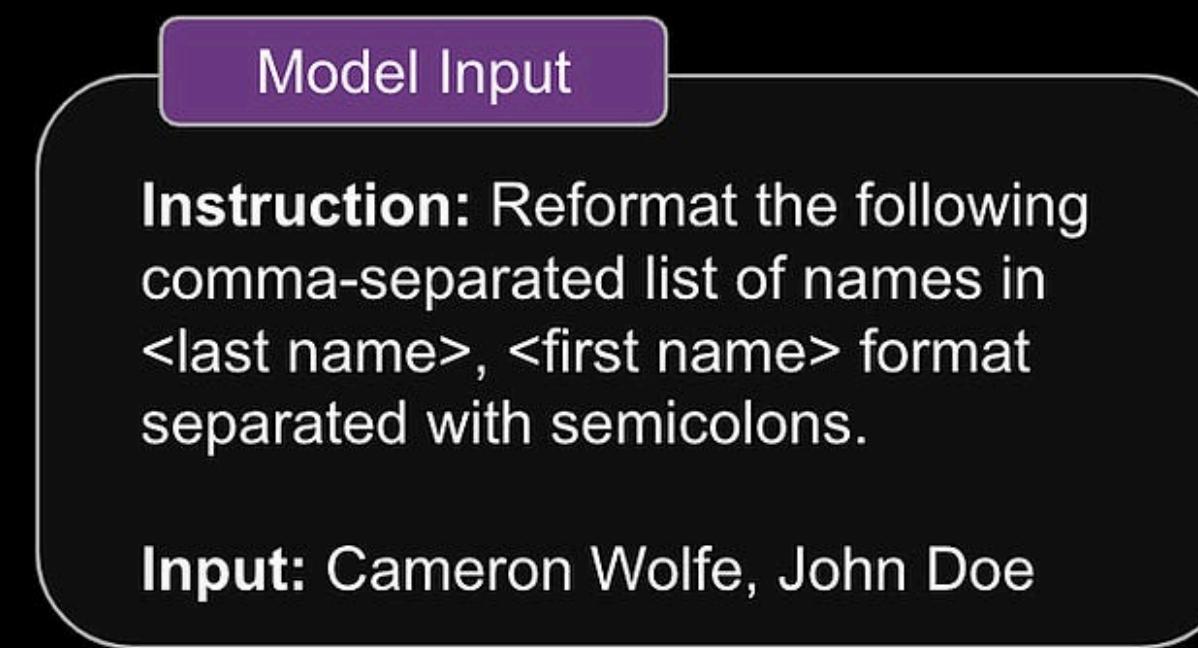
}

Instruction following

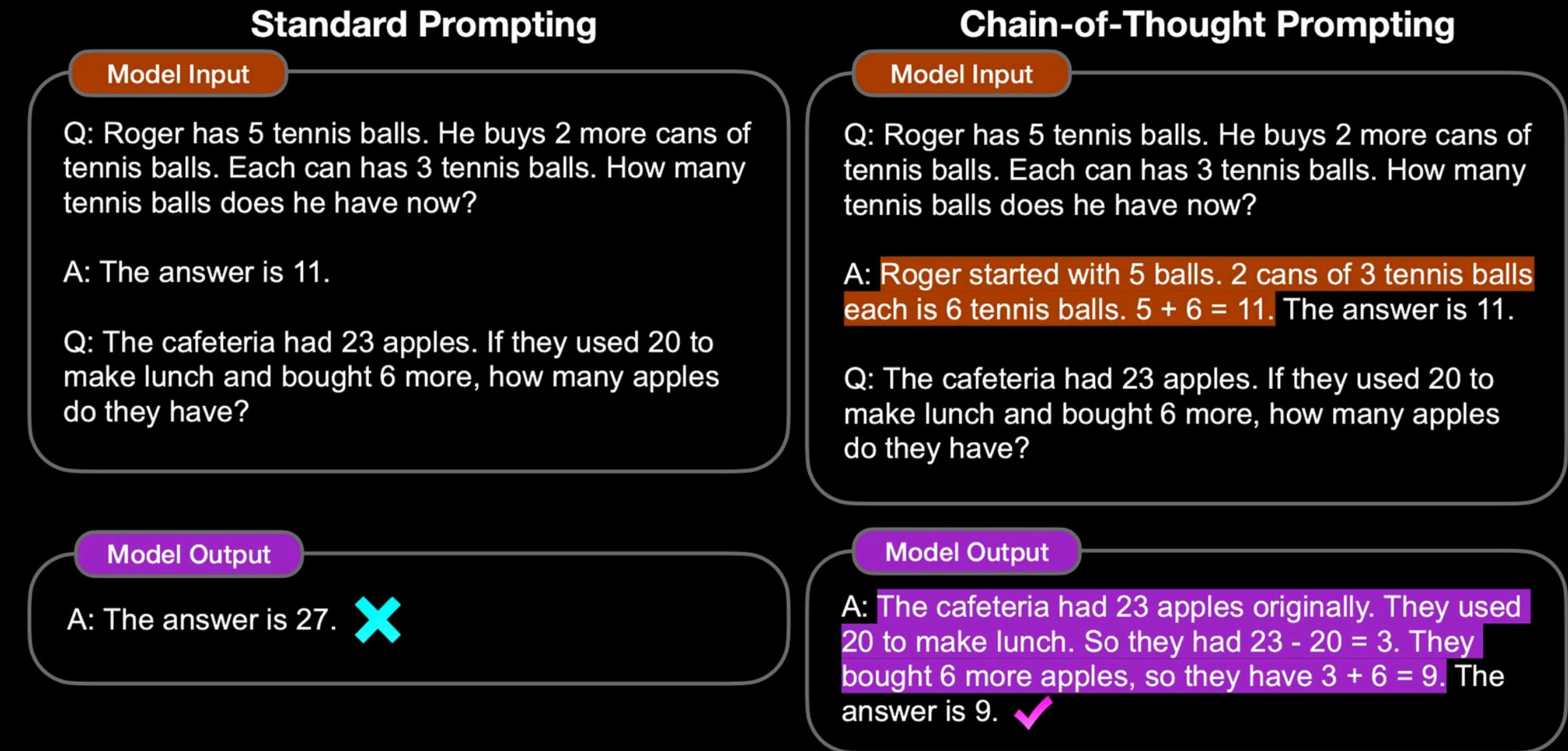
without additional instruction



with additional instruction



Chain of Thought



Section Summary

We saw

1. History of language models
2. Scaling Laws (OpenAI and Chinchilla)
 - a. OpenAI - Best possible model
 - b. Chinchilla - Best possible model with provided budget
3. Why we need scaling laws
4. Emergent abilities

Questions

Table of Contents - Evolution of GPT

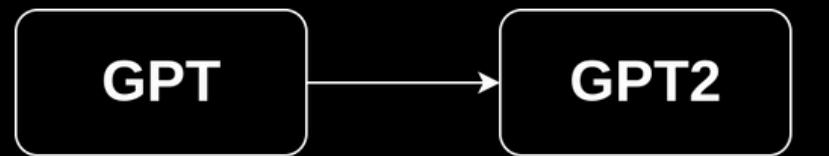
1. GPT on a timeline
2. Data preparation

Generative Pre-training

Technical Evolution

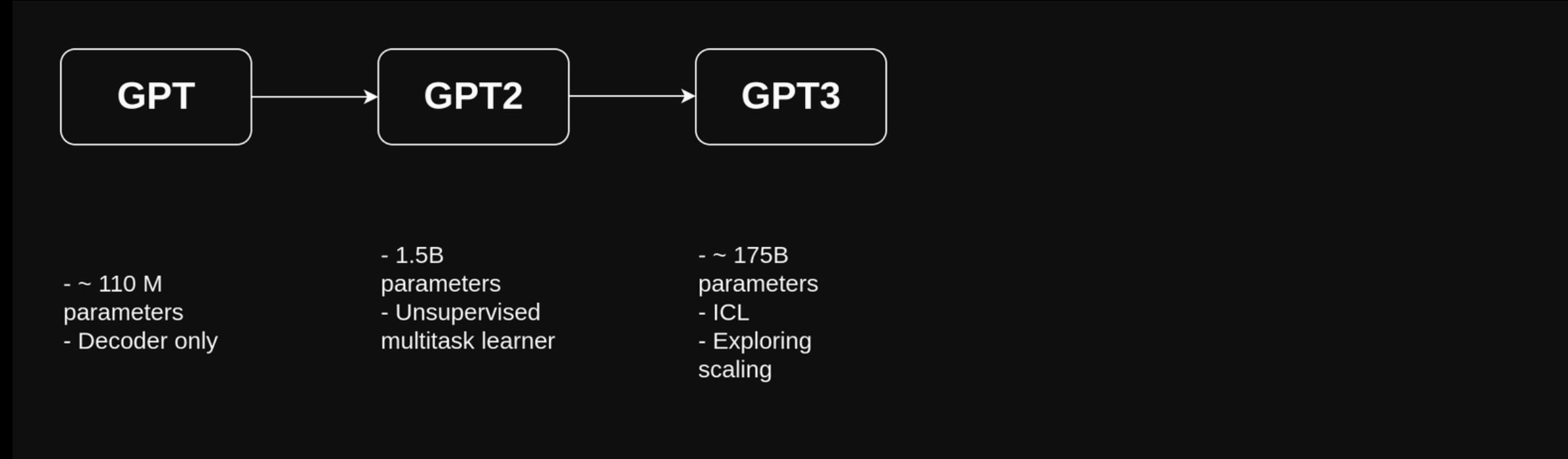
GPT

- ~ 110 M parameters
- Decoder only



- ~ 110 M parameters
- Decoder only

- 1.5B parameters
- Unsupervised multitask learner





- ~ 110 M parameters
- Decoder only

- 1.5B parameters
- Unsupervised multitask learner

- ~ 175B parameters
- ICL
- Exploring scaling

- Code + instruction finetuning + RLHF + Chat



- ~ 110 M parameters
- Decoder only

- 1.5B parameters
- Unsupervised multitask learner

- ~ 175B parameters
- ICL
- Exploring scaling

- Code + instruction finetuning + RLHF + Chat

- Multimodal

Resources

Pretraining Corpora

Resources

Pretraining Corpora

Webpages		
Books		
Conversation		
code		

Resources

Pretraining Corpora

Webpages	General text	CommonCrawl
Books		
Conversation		
code		

Resources

Pretraining Corpora

Webpages	General text	CommonCrawl
Books	Long context understanding	BookCorpus
Conversation		
code		

Resources

Pretraining Corpora

Webpages	General text	CommonCrawl
Books	Long context understanding	BookCorpus
Conversation	reply response	RedditLinks
code		

Resources

Pretraining Corpora

Webpages	General text	CommonCrawl
Books	Long context understanding	BookCorpus
Conversation	reply response	RedditLinks
code	coding & reasoning abilities	GitHub

Resources

Fine-tuning Datasets

Resources

Fine-tuning Datasets

Instruction Finetuning

1. NLP Task Datasets

Resources

Fine-tuning Datasets

Instruction Finetuning

- 1. NLP Task Datasets**
- 2. Daily Chat**

Resources

Fine-tuning Datasets

Instruction Finetuning

1. NLP Task Datasets

2. Daily Chat

3. Synthetic Datasets - @Ansh

Resources

Libraries

Transformers	
DeepSpeed	
Triton	
vLLM	

Resources

Libraries

Transformers	model architecture & weights
DeepSpeed	
Triton	
vLLM	

Resources

Libraries

Transformers	model architecture & weights
DeepSpeed	efficient distributed training
Triton	
vLLM	

Resources

Libraries

Transformers	model architecture & weights
DeepSpeed	efficient distributed training
Triton	writing autotuned cuda kernels
vLLM	

Resources

Libraries

Transformers	model architecture & weights
DeepSpeed	efficient distributed training
Triton	writing autotuned cuda kernels
vLLM	inference

Pre-Training

Data Preparation



Section Summary

We saw

1. Evolution of GPT
2. Pretraining corpora
3. Finetuning Corpora
4. Data Preparation

Questions

Table of Contents - Scaling and Need for scaling

1. Various architectures
2. Positional encodings
3. Scaling beyond 3000 tokens

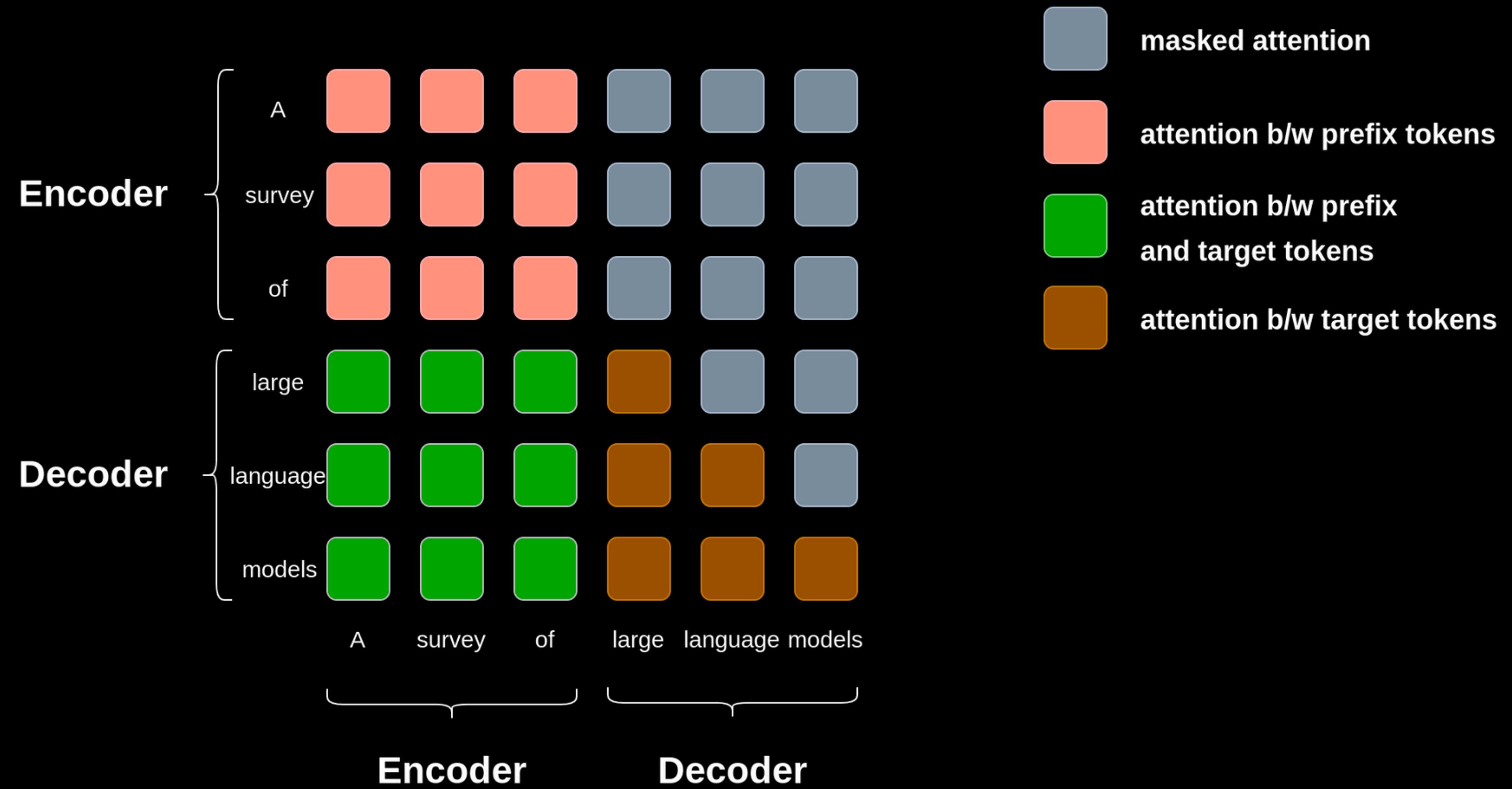
Pre-Training

Choosing the Architecture

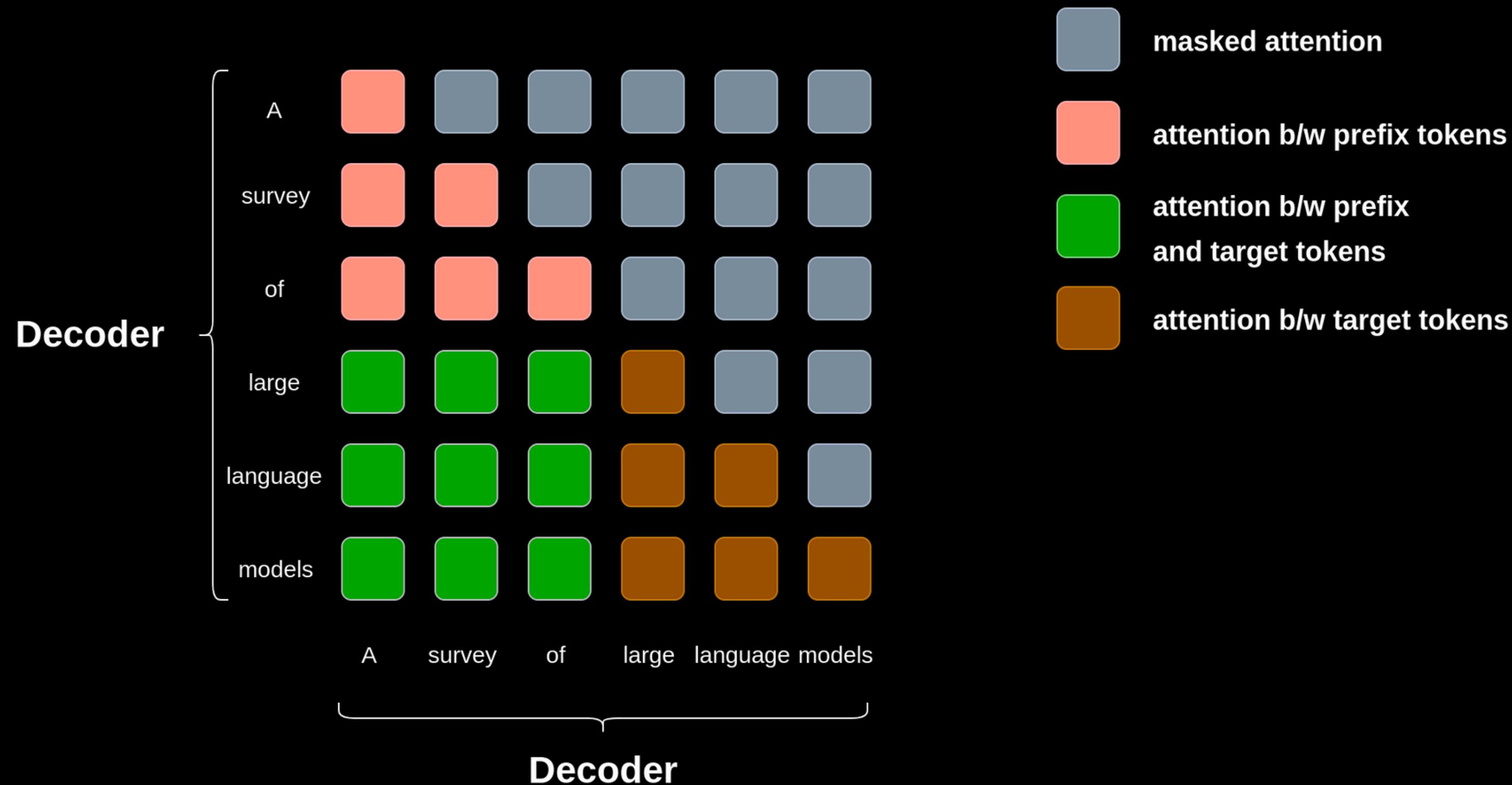
Pre-Training

Choosing the Architecture

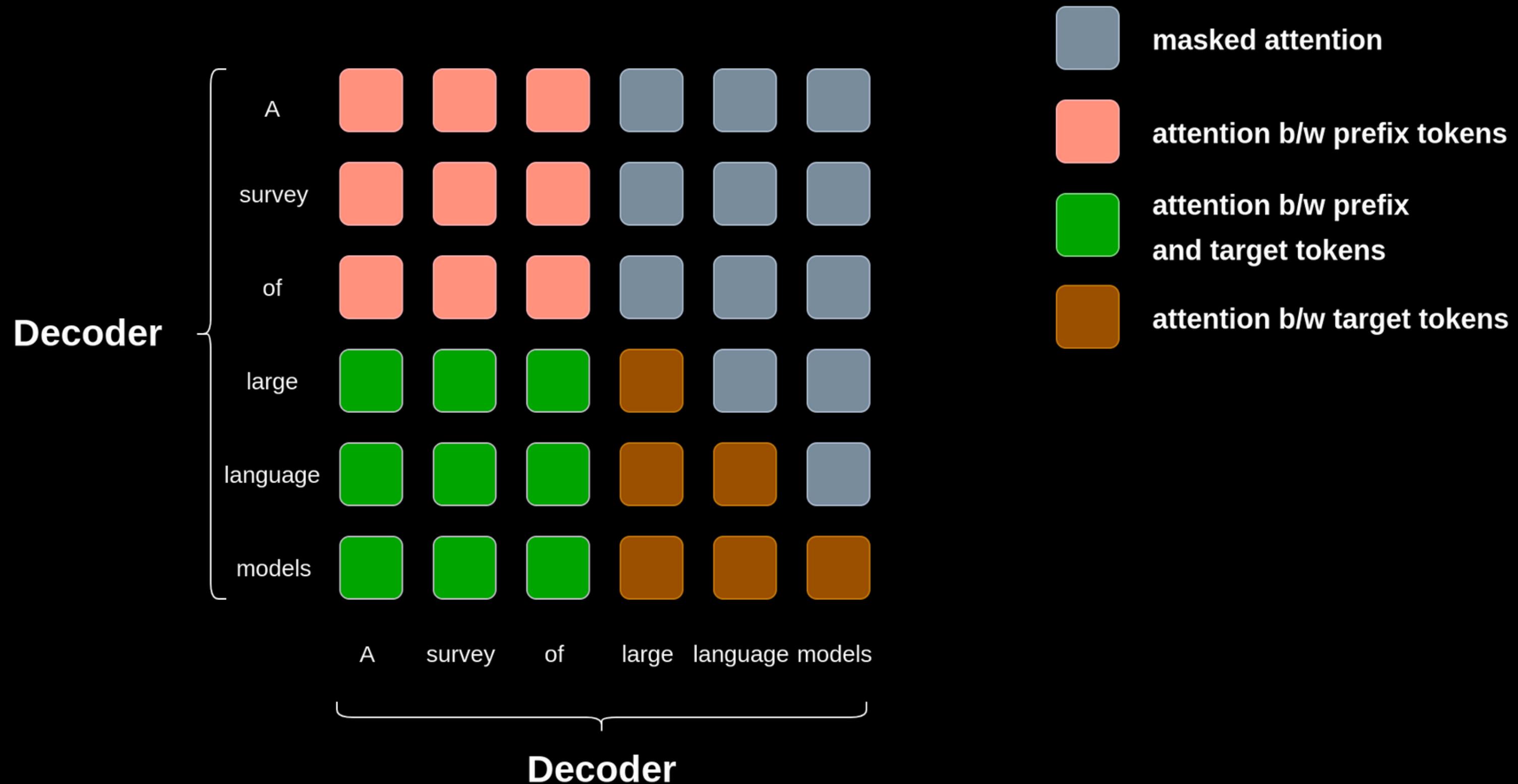
Encoder-Decoder Model



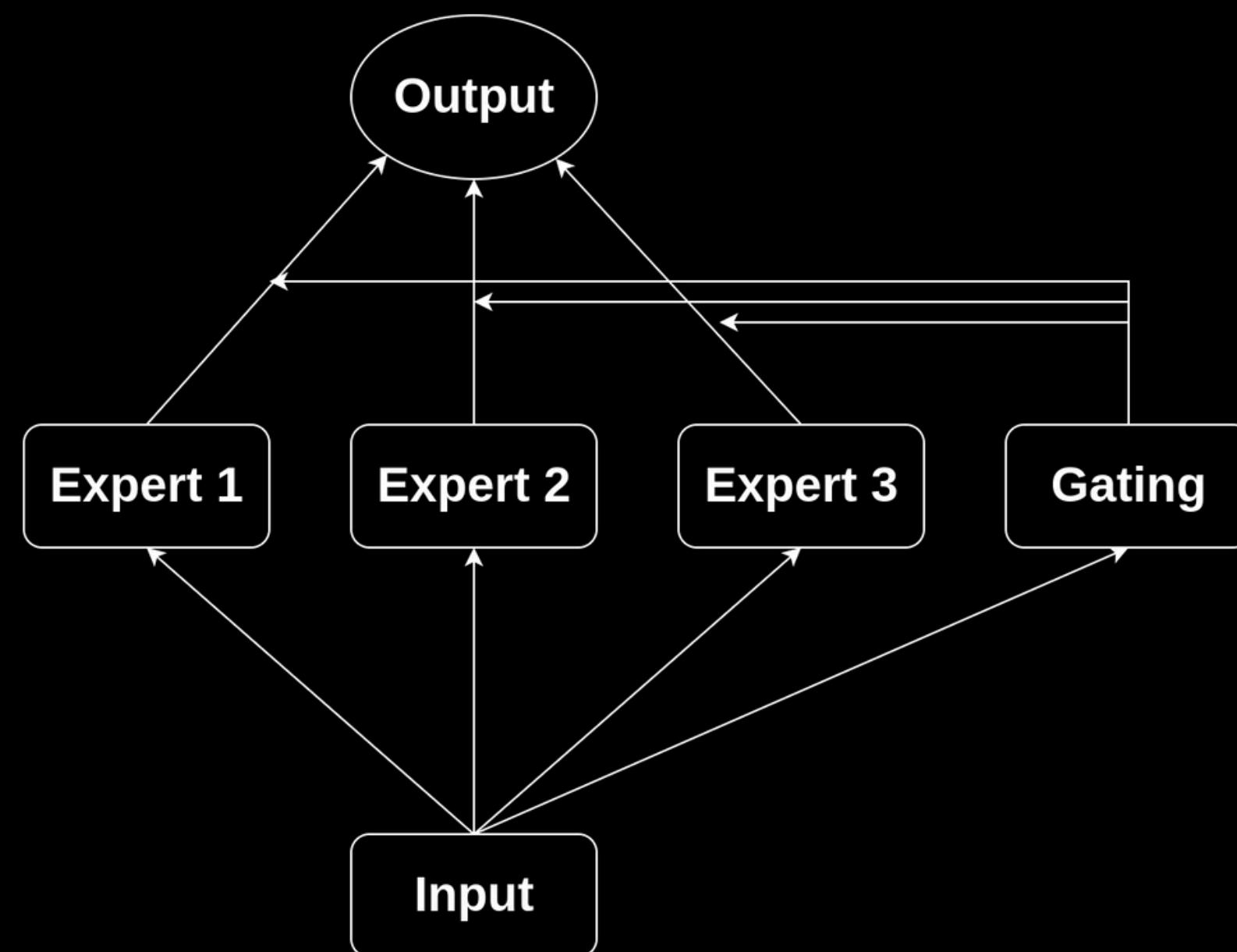
Causal Decoder Model



Non-causal Decoder Model



Mixture of Experts



Emergent Architectures

RWKV
State Space Models

H3
XLSTM
S4

And many others...!

Pre-Training

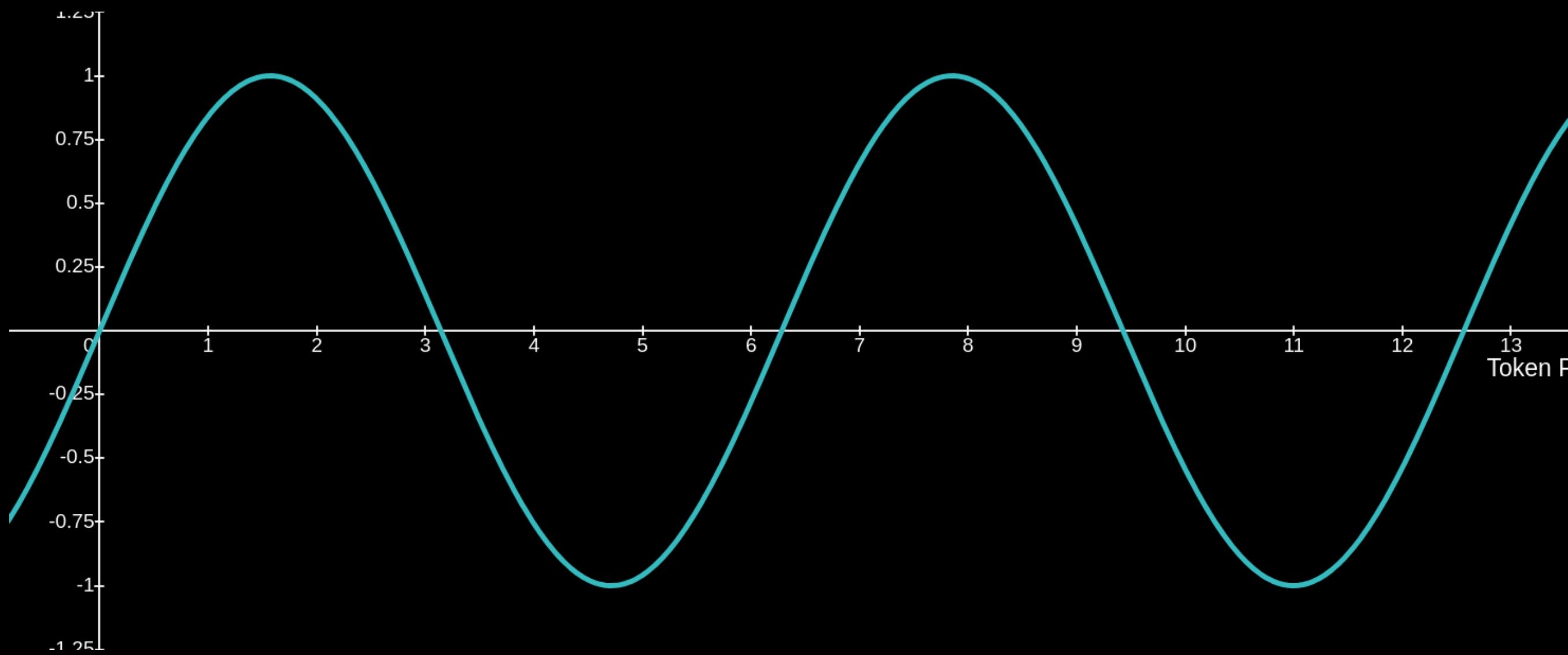
Position Embeddings

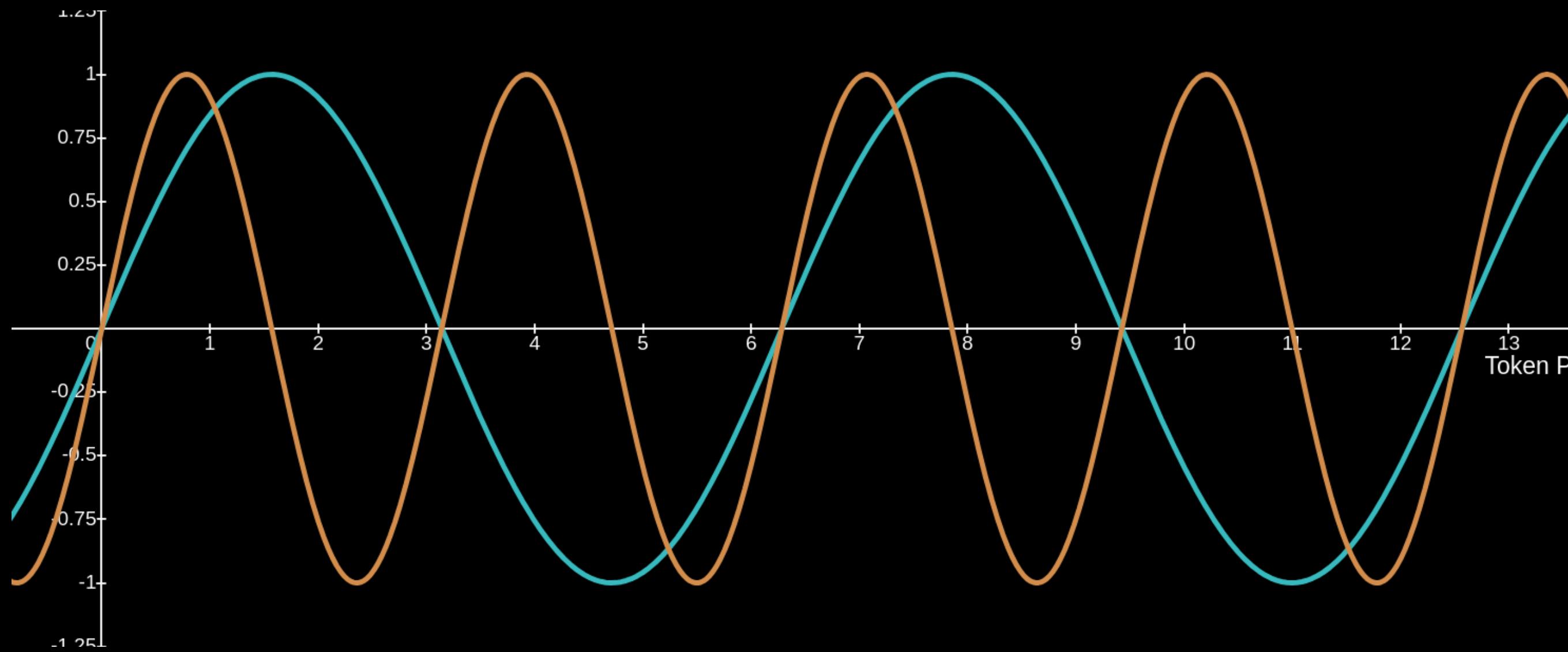
Absolute Position Embeddings

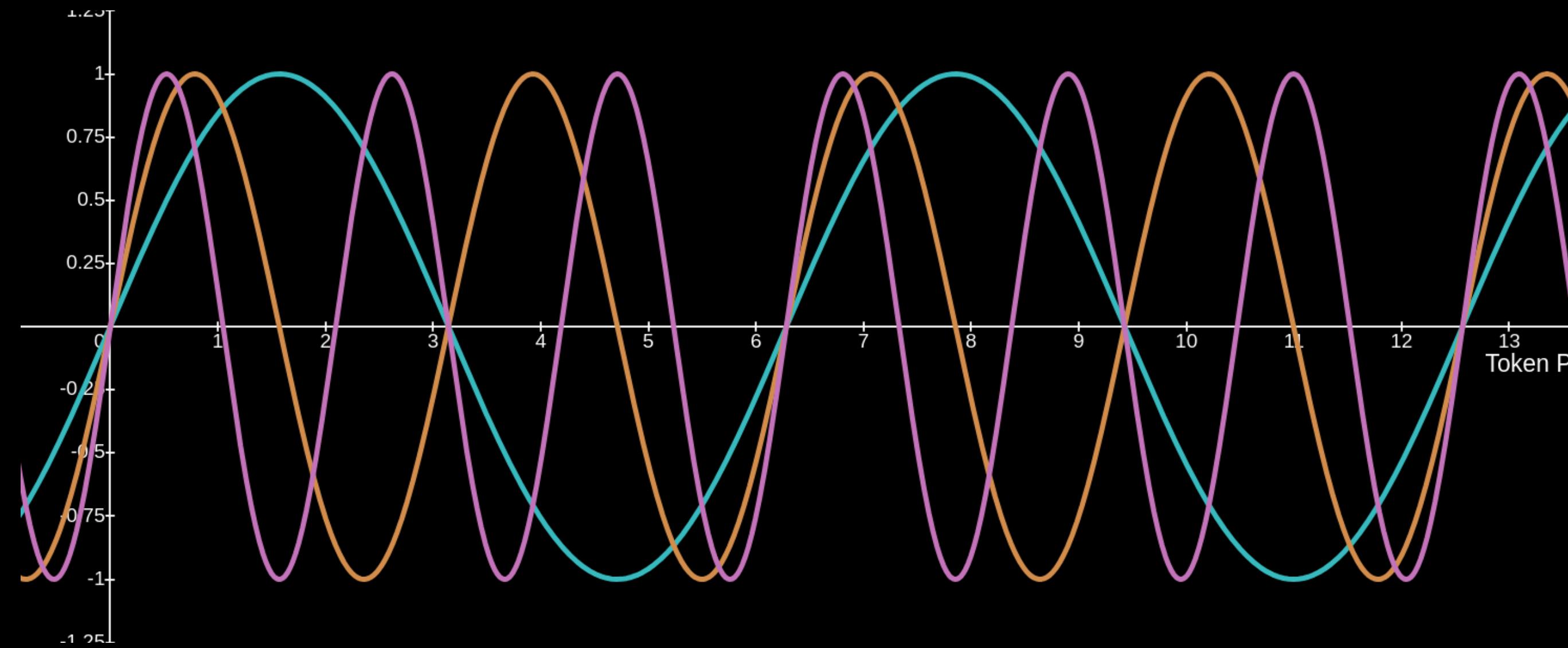
Sinusoidal

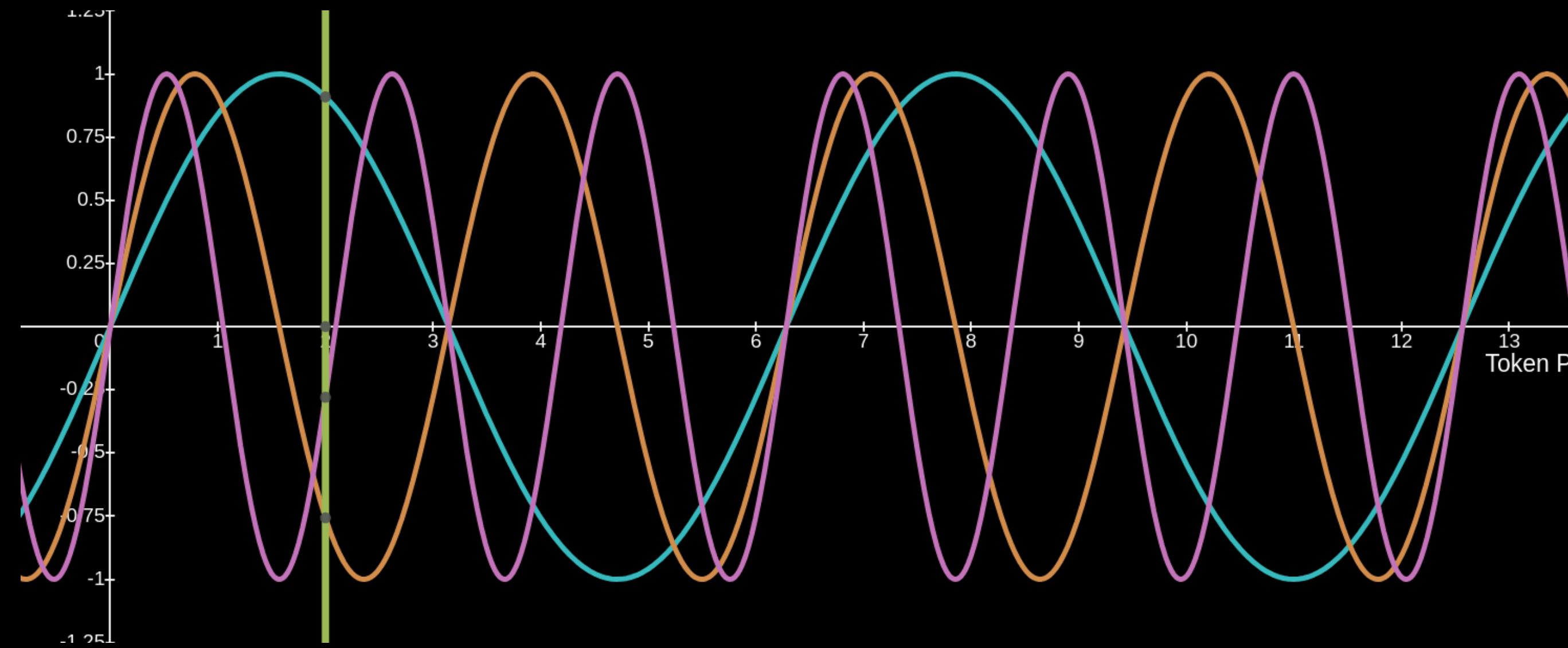
Learned

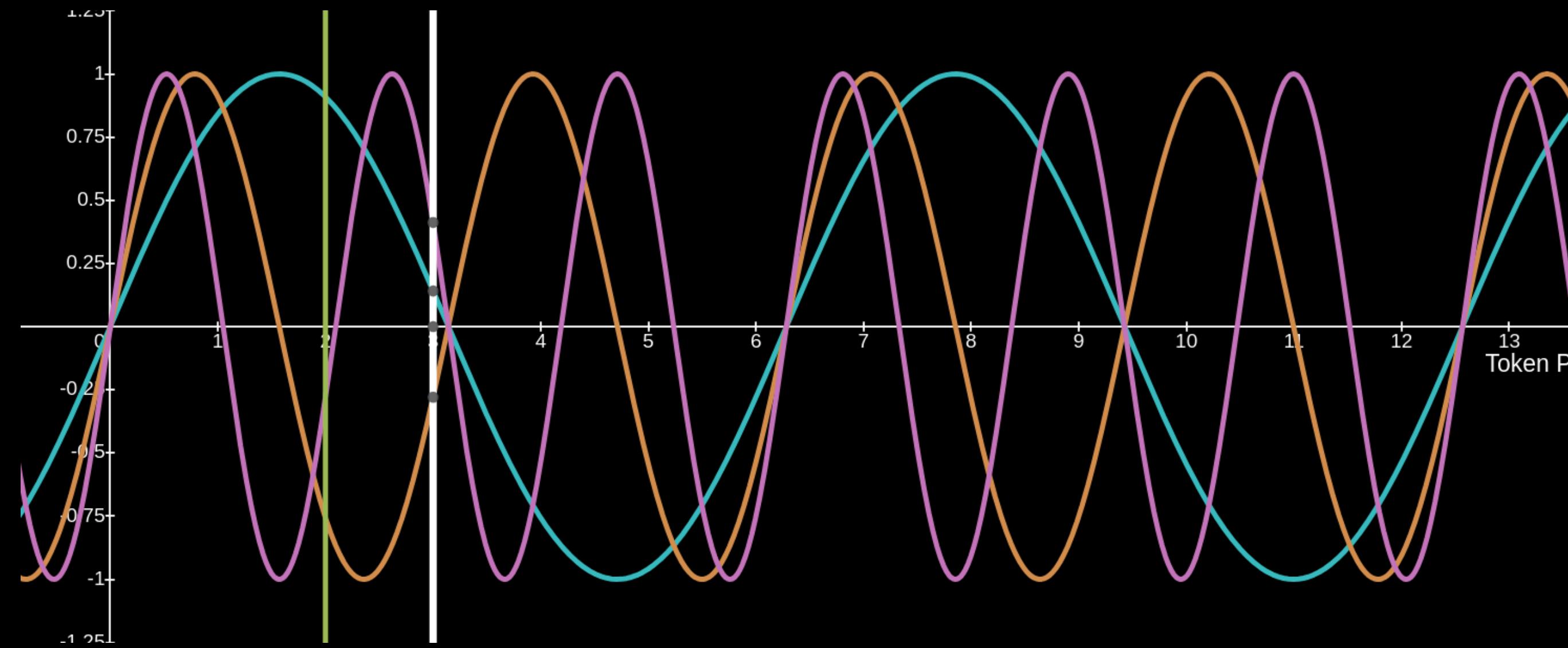
An Intuition for Sinusoidal Embeddings

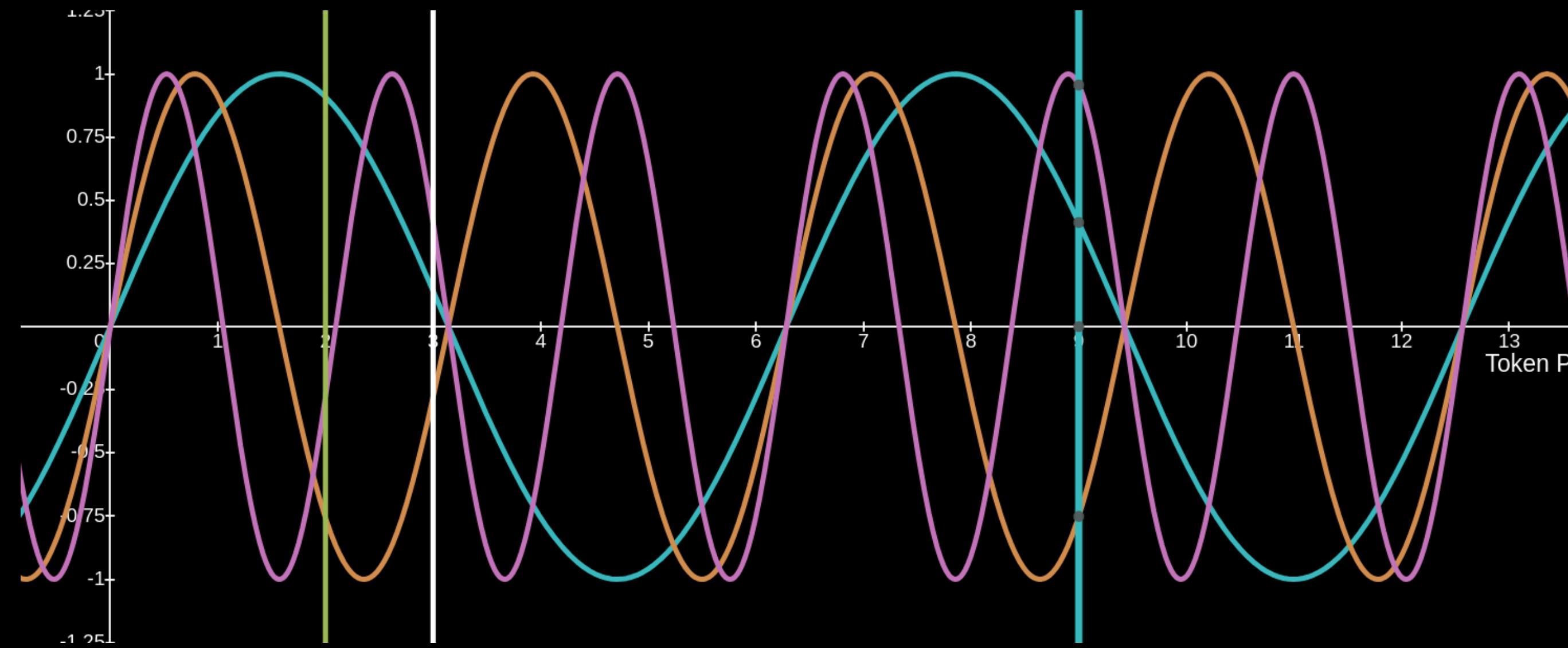


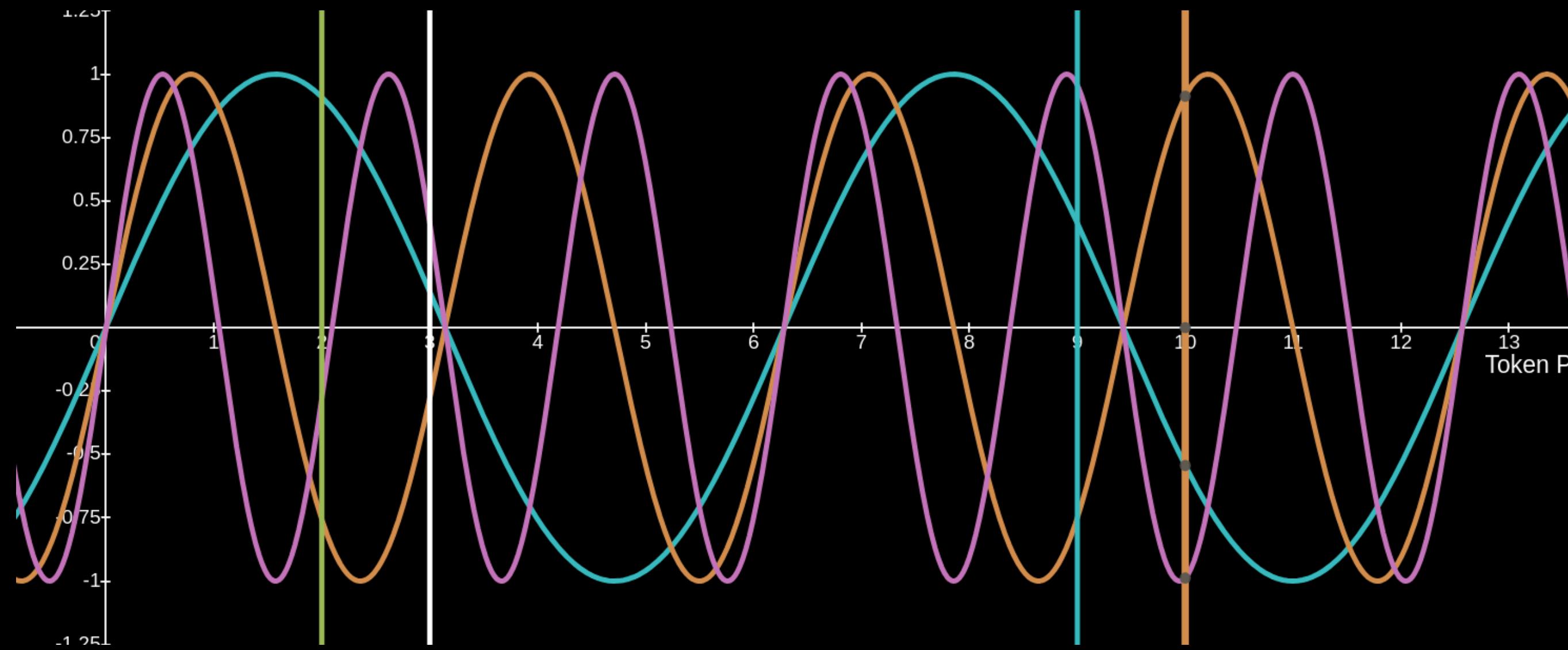












$$\begin{bmatrix} \overrightarrow{P_1} \\ \overrightarrow{P_2} \\ \dots \\ \dots \\ \overrightarrow{P_{maxlen}} \end{bmatrix} + \begin{bmatrix} \overrightarrow{W_1} \\ \overrightarrow{W_2} \\ \dots \\ \dots \\ \overrightarrow{W_{maxlen}} \end{bmatrix} = \begin{bmatrix} \overrightarrow{E_1} \\ \overrightarrow{E_2} \\ \dots \\ \dots \\ \overrightarrow{E_{maxlen}} \end{bmatrix}$$

An Intuition for Learned Embeddings

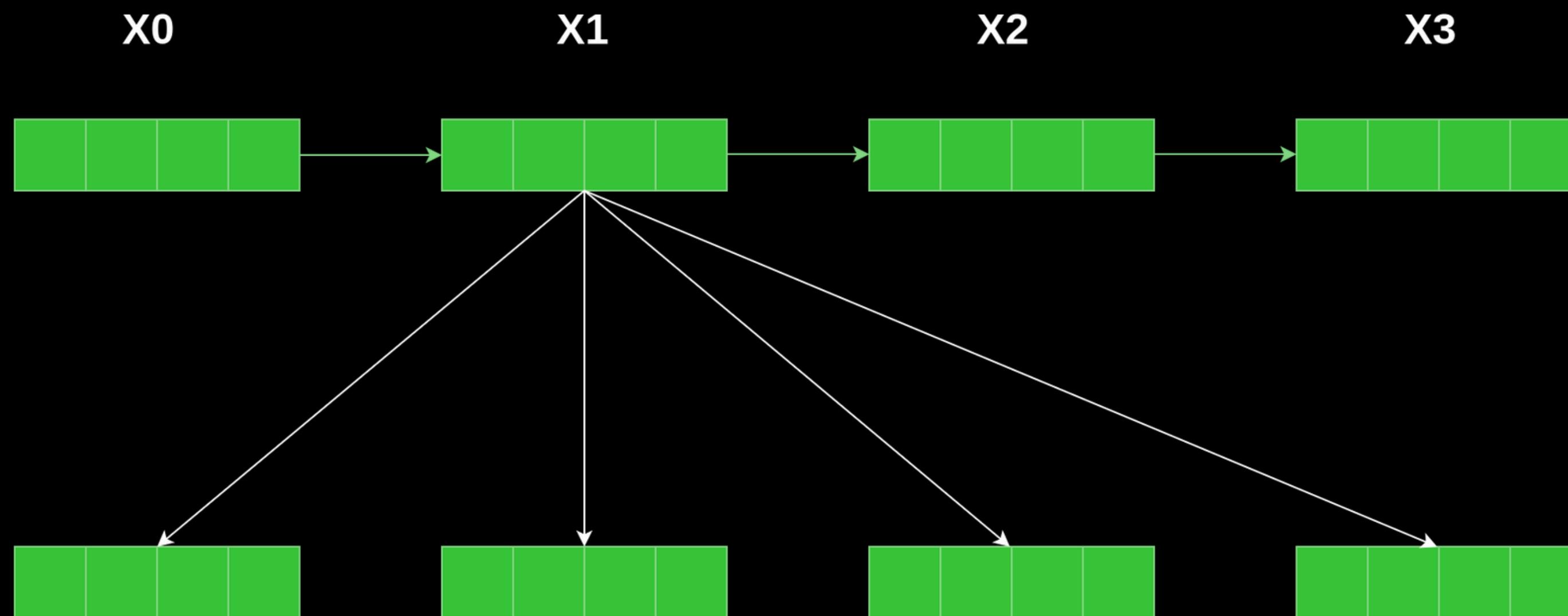
$$\begin{bmatrix} \overrightarrow{P_1} \\ \overrightarrow{P_2} \\ \dots \\ \dots \\ \overrightarrow{P_{maxlen}} \end{bmatrix} + \begin{bmatrix} \overrightarrow{W_1} \\ \overrightarrow{W_2} \\ \dots \\ \dots \\ \overrightarrow{W_{maxlen}} \end{bmatrix} = \begin{bmatrix} \overrightarrow{E_1} \\ \overrightarrow{E_2} \\ \dots \\ \dots \\ \overrightarrow{E_{maxlen}} \end{bmatrix}$$

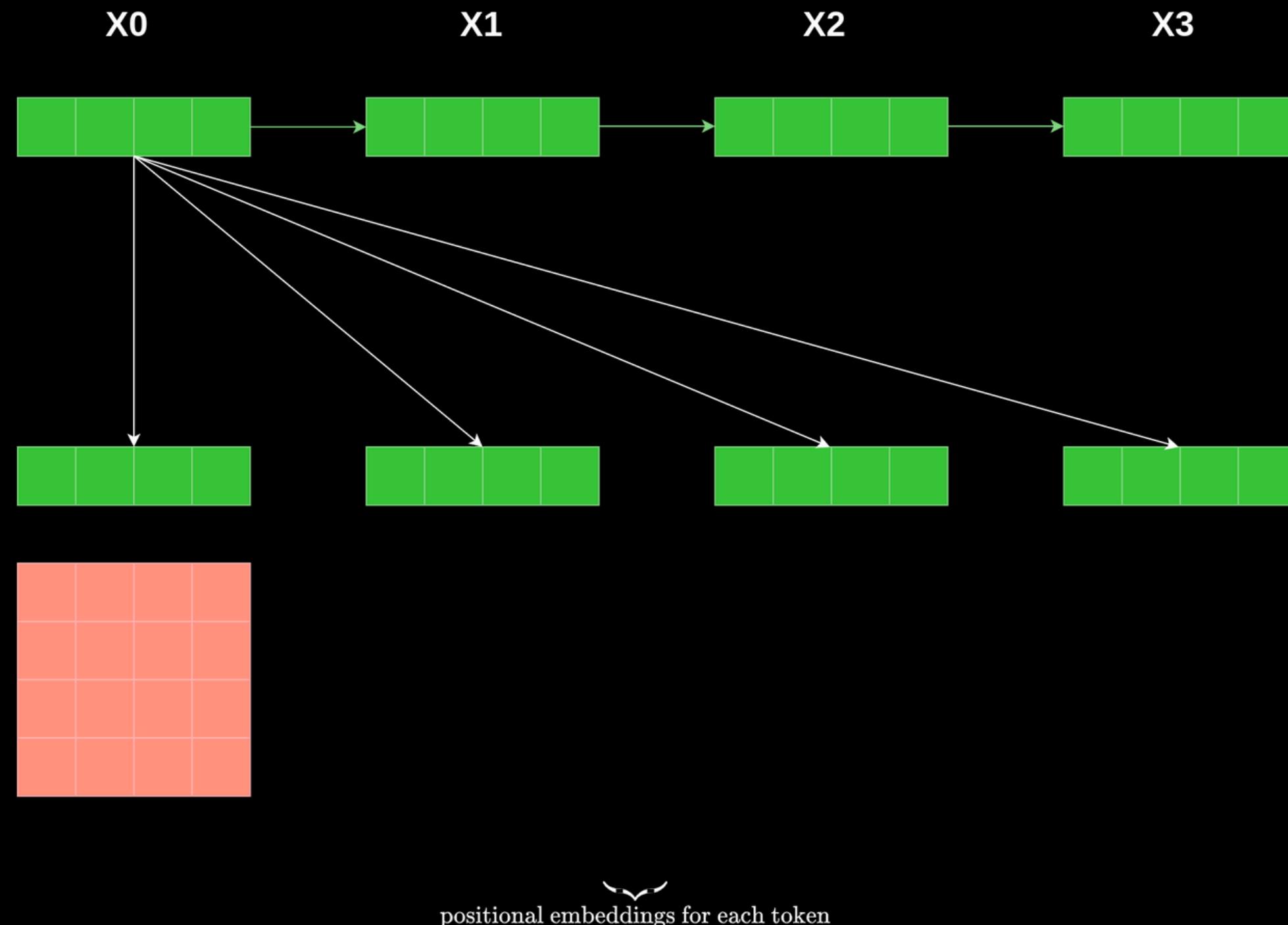
But these are learned along
with other params of the model

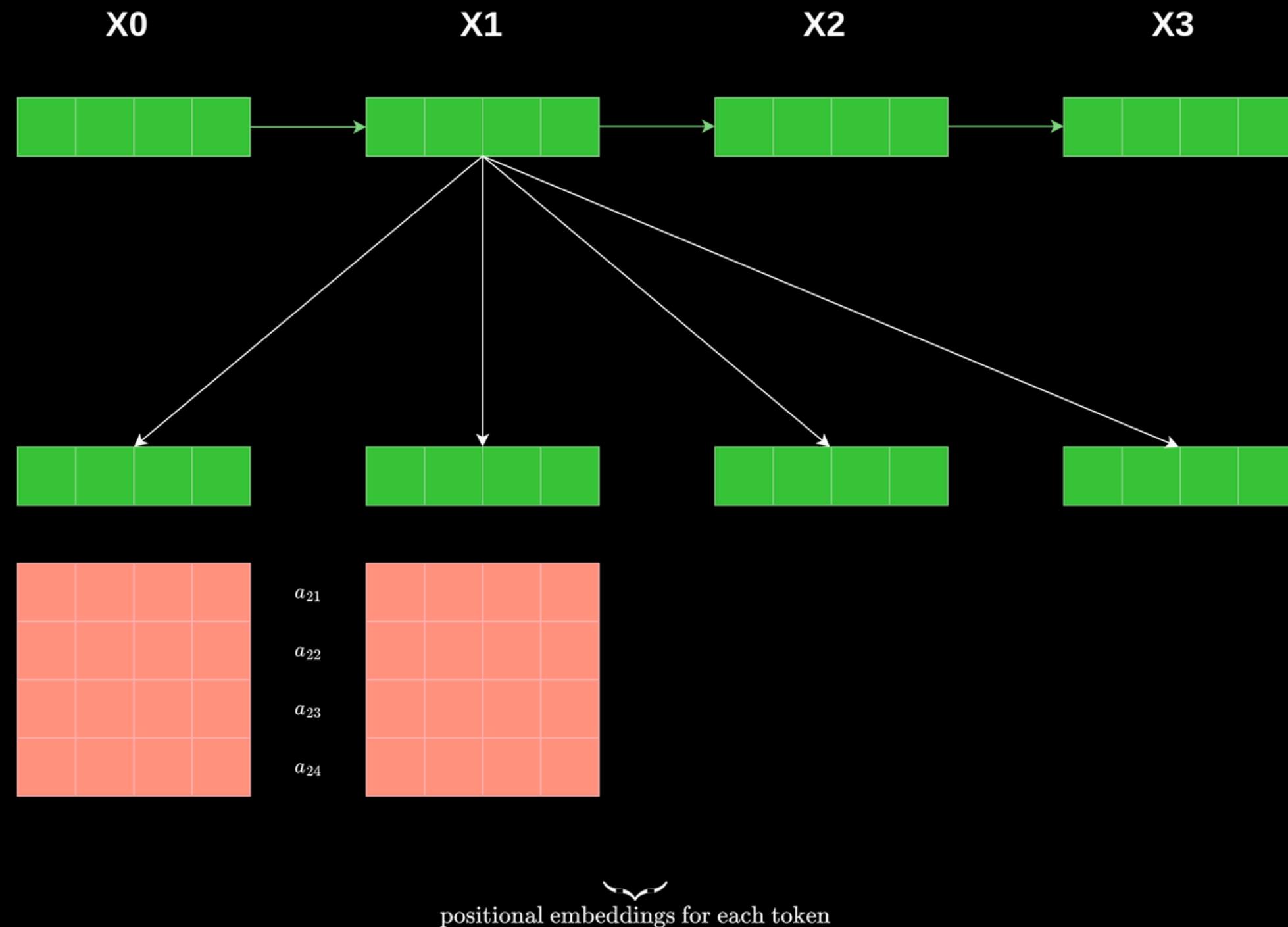
Pre-Training

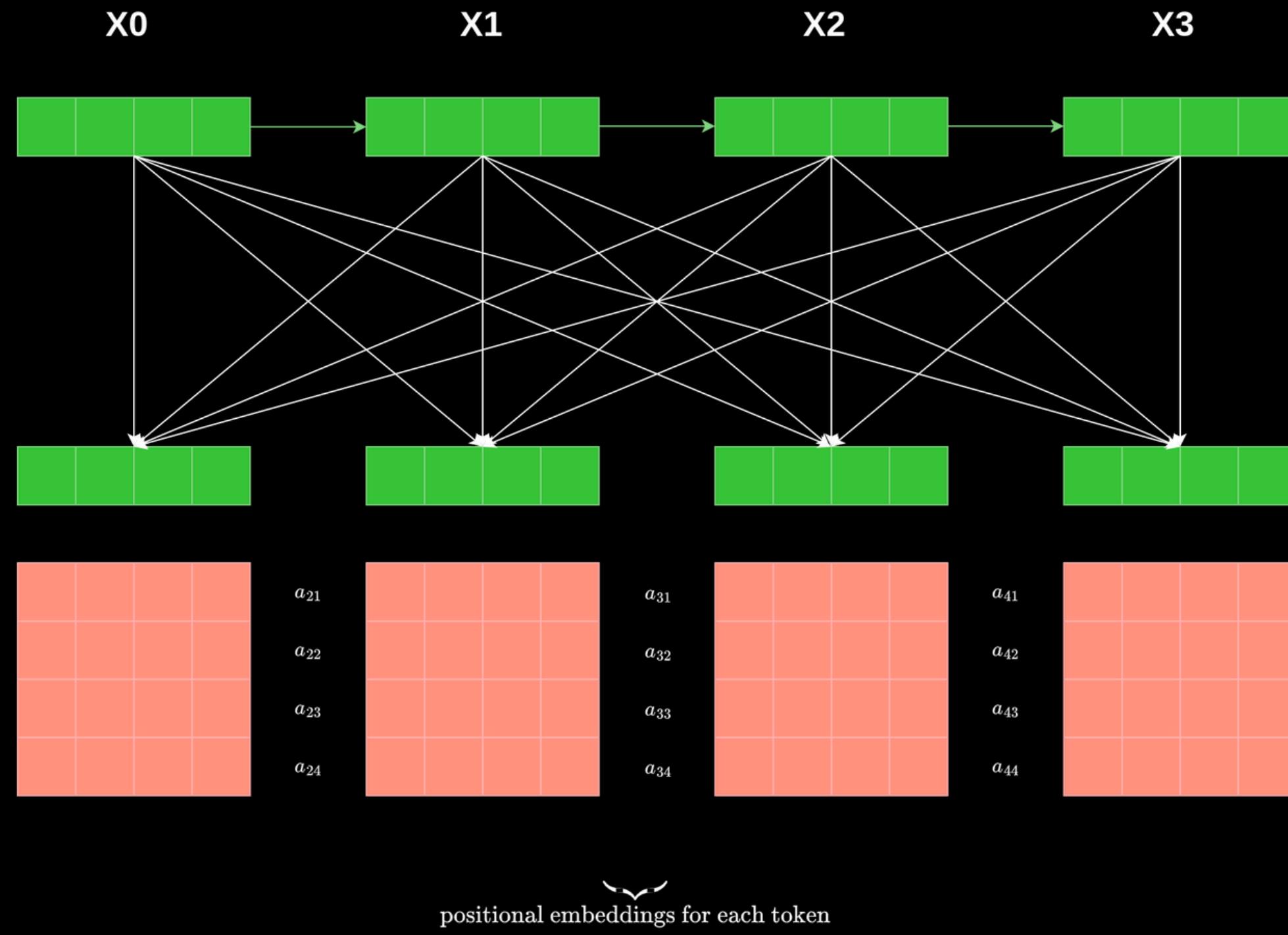
Relative Position Embeddings

Position Embeddings



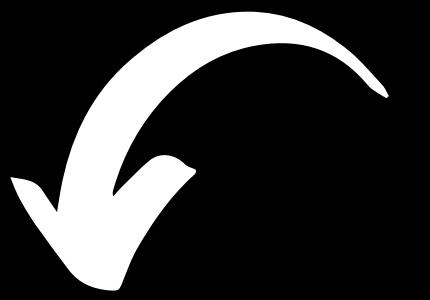
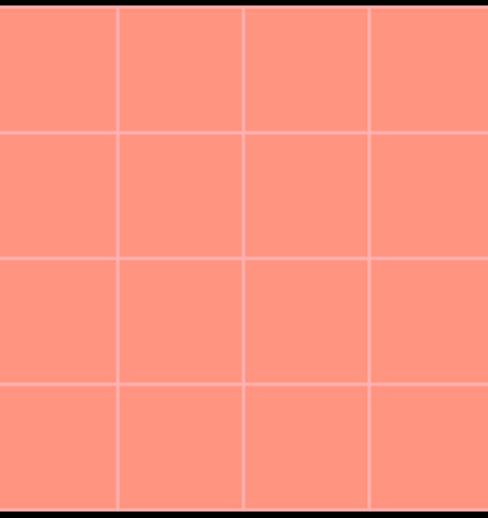






Query Projection Key Projection

$$e_{ij} = \frac{\overbrace{x_i W^Q}^{\text{Query Projection}} \left(\overbrace{x_j W^K}^{\text{Key Projection}} \right)^\top}{\sqrt{d_z}}$$



$$e_{ij} = \frac{\text{Query Projection } \widehat{x_i W^Q} \quad (\text{Key Projection } \widehat{x_j W^K} + \text{Key position vector } \widehat{a_{ij}^K})^\top}{\sqrt{d_z}}$$

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

**What is this operation
called?**

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

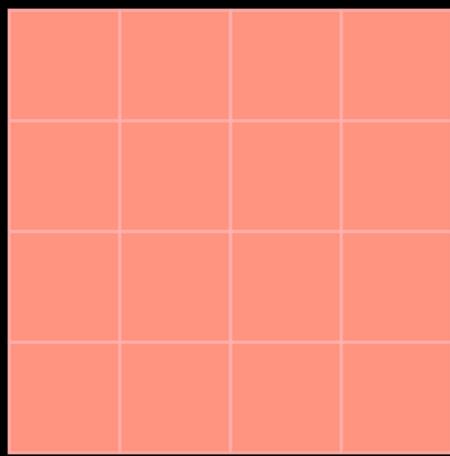


$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^n \exp(e_{ik})}$$



What is this operation
called?
SoftMax!

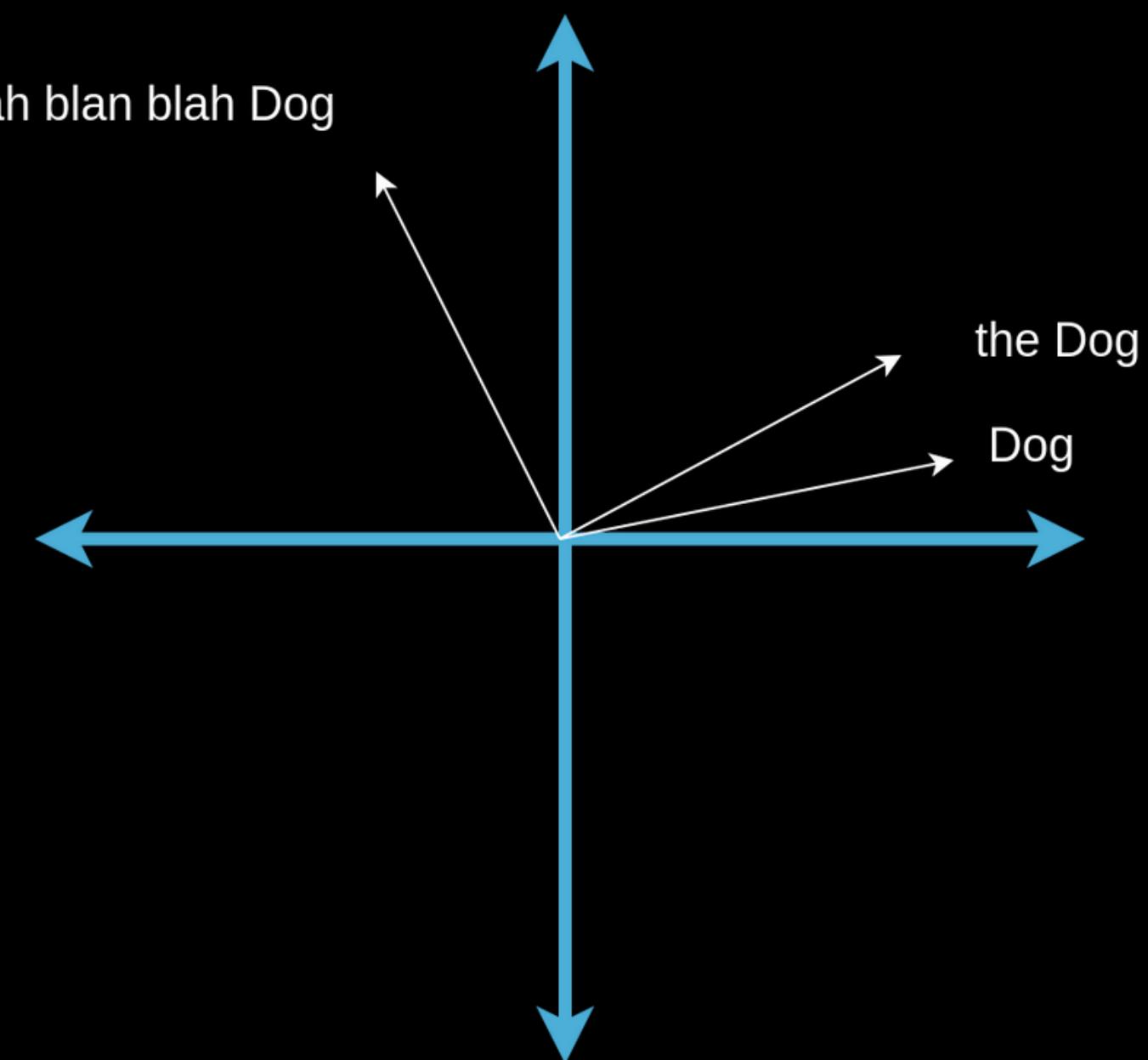
$$z_i = \sum_{j=1}^n \alpha_{ij} \left(\underbrace{x_j W^V}_{\text{Value Projection}} + \underbrace{a_{ij}^V}_{\text{Value position vector}} \right)$$

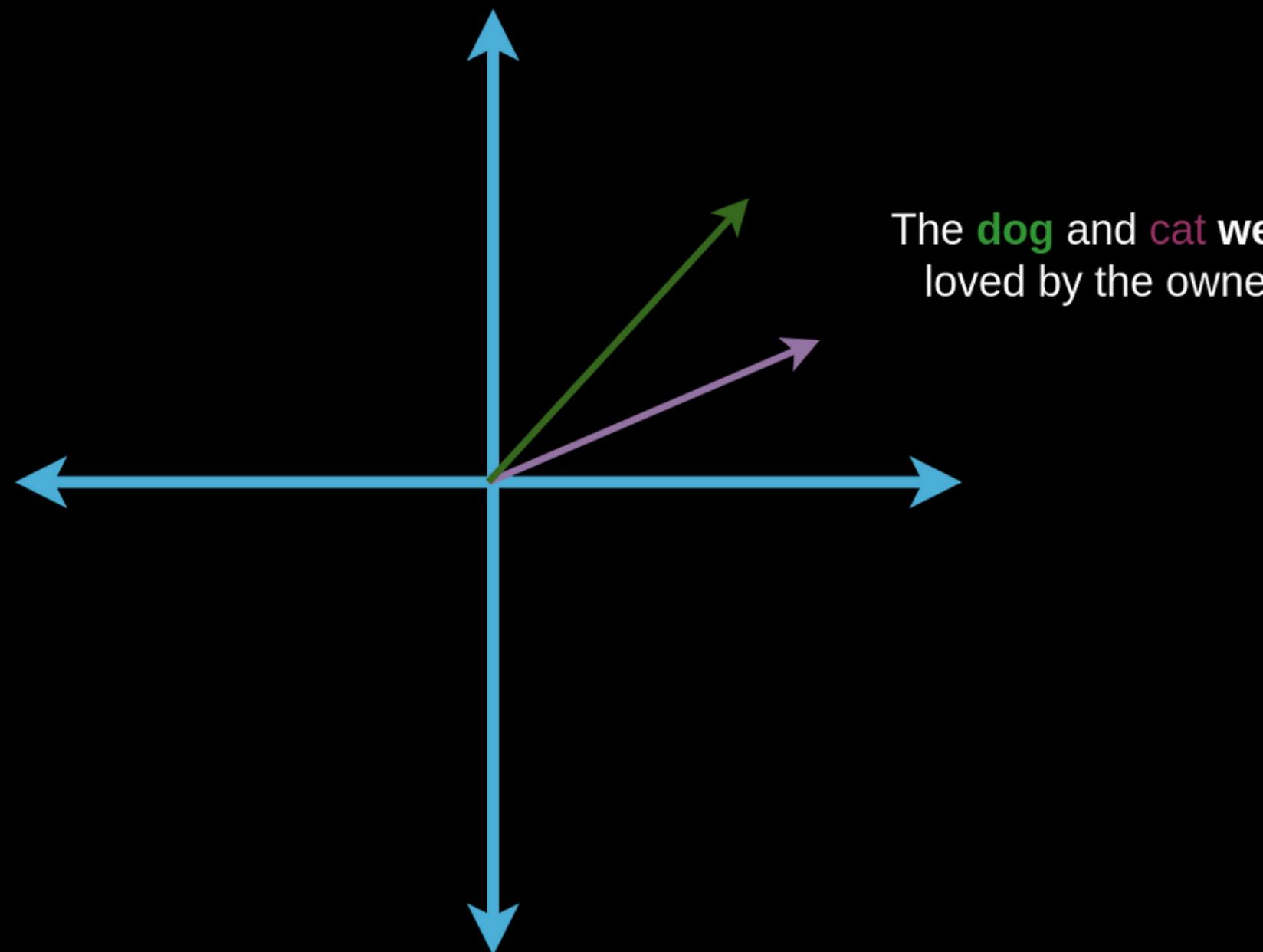


Pre-Training

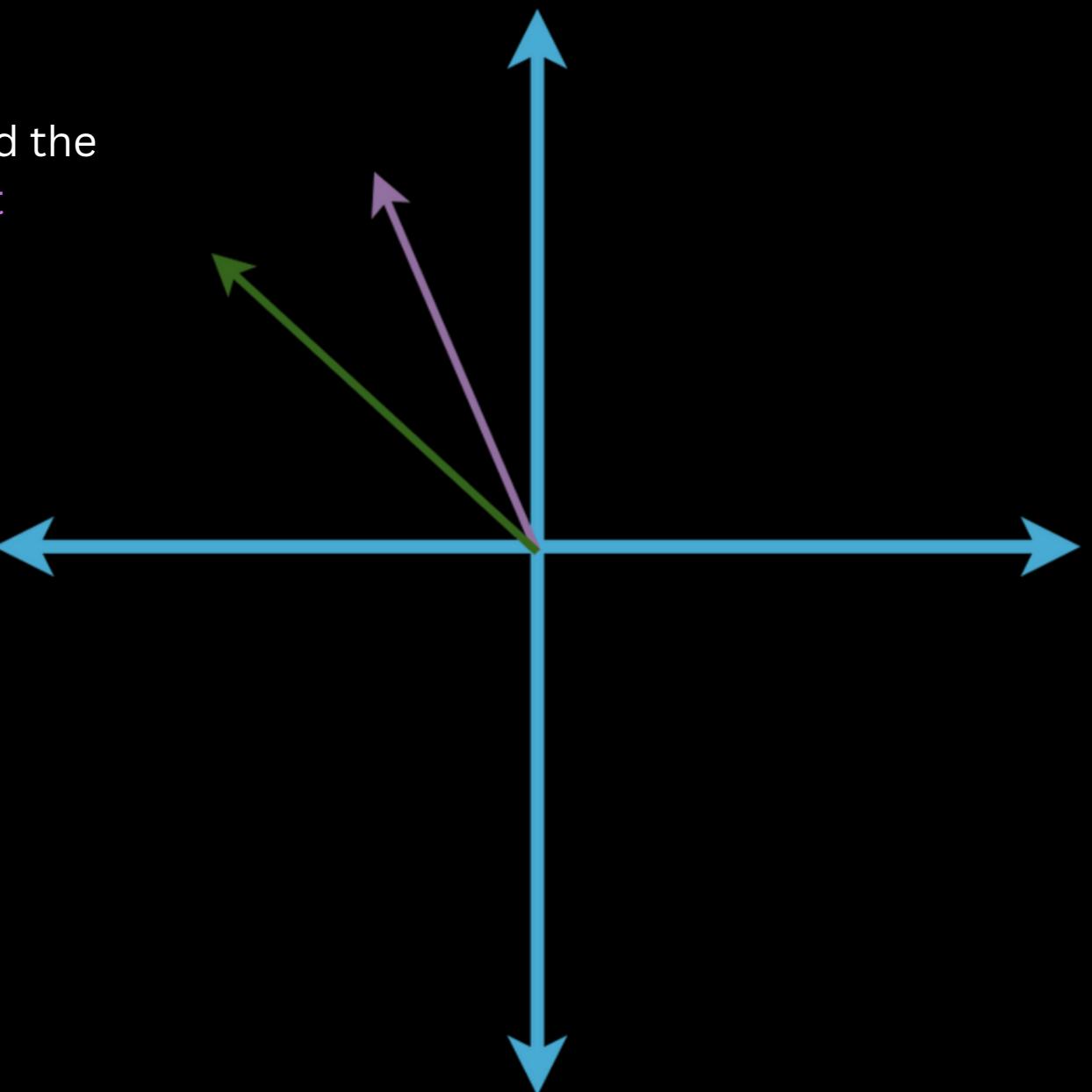
Rotary Position Embeddings (RoPE)

Position Embeddings





The owner loved the
dog and cat



Rotation Matrix

$$f_{q,k}(x_m, m) = \overbrace{\begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix}}^{\text{Vector we want to rotate}} \begin{pmatrix} W_{qk}^{(11)} & W_{qk}^{(12)} \\ W_{qk}^{(21)} & W_{qk}^{(22)} \end{pmatrix} \overbrace{\begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}}^{\text{Rotation Matrix}}$$

So, what would you use?

Scaling beyond 3000 tokens

Scaling beyond 3000 tokens

Increase context length
slowly during training

Scaling beyond 3000 tokens

Increase context length
slowly during training

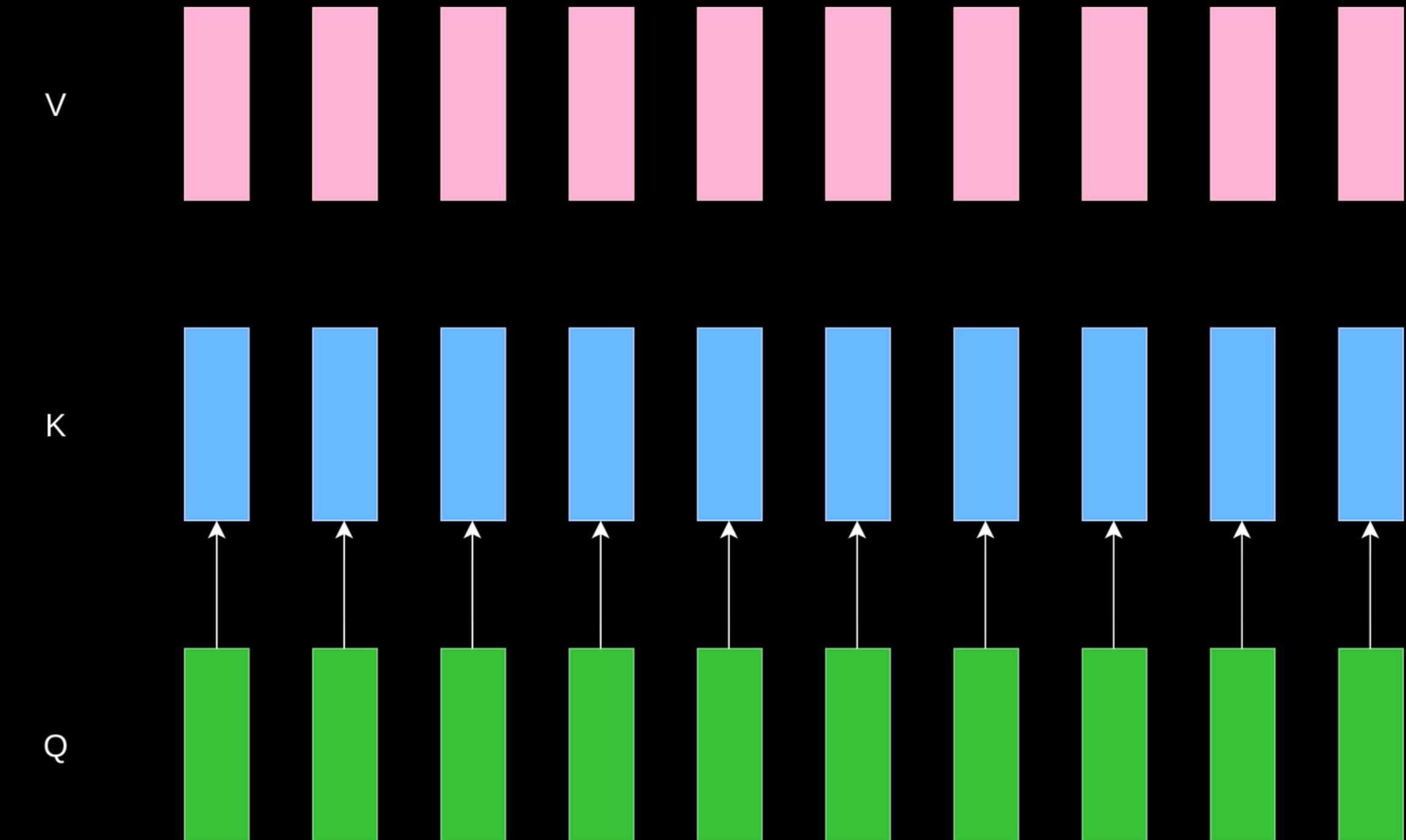
Position interpolation by
multiplying with a factor

Pre-Training

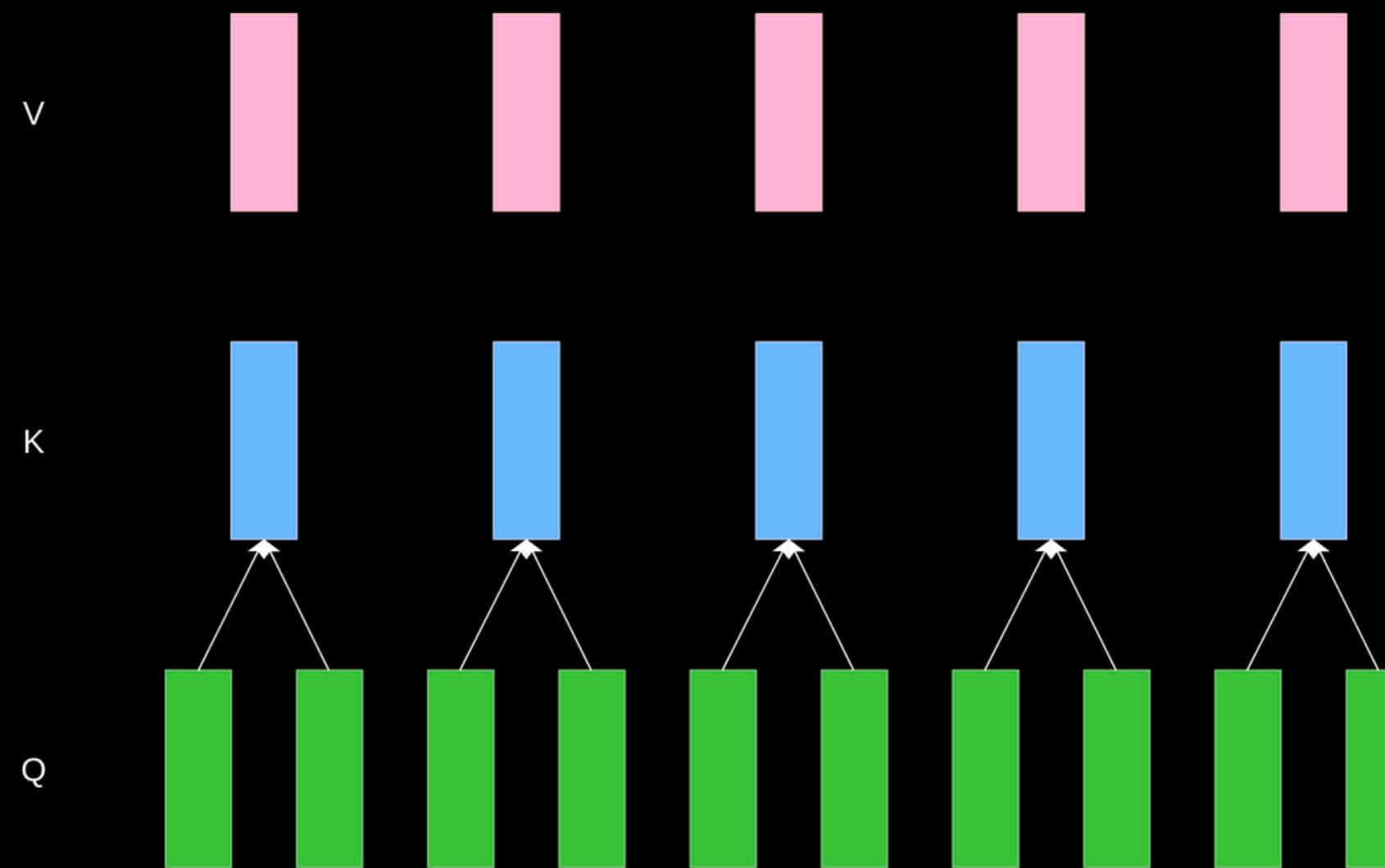
Attention Mechanism

Full Attention

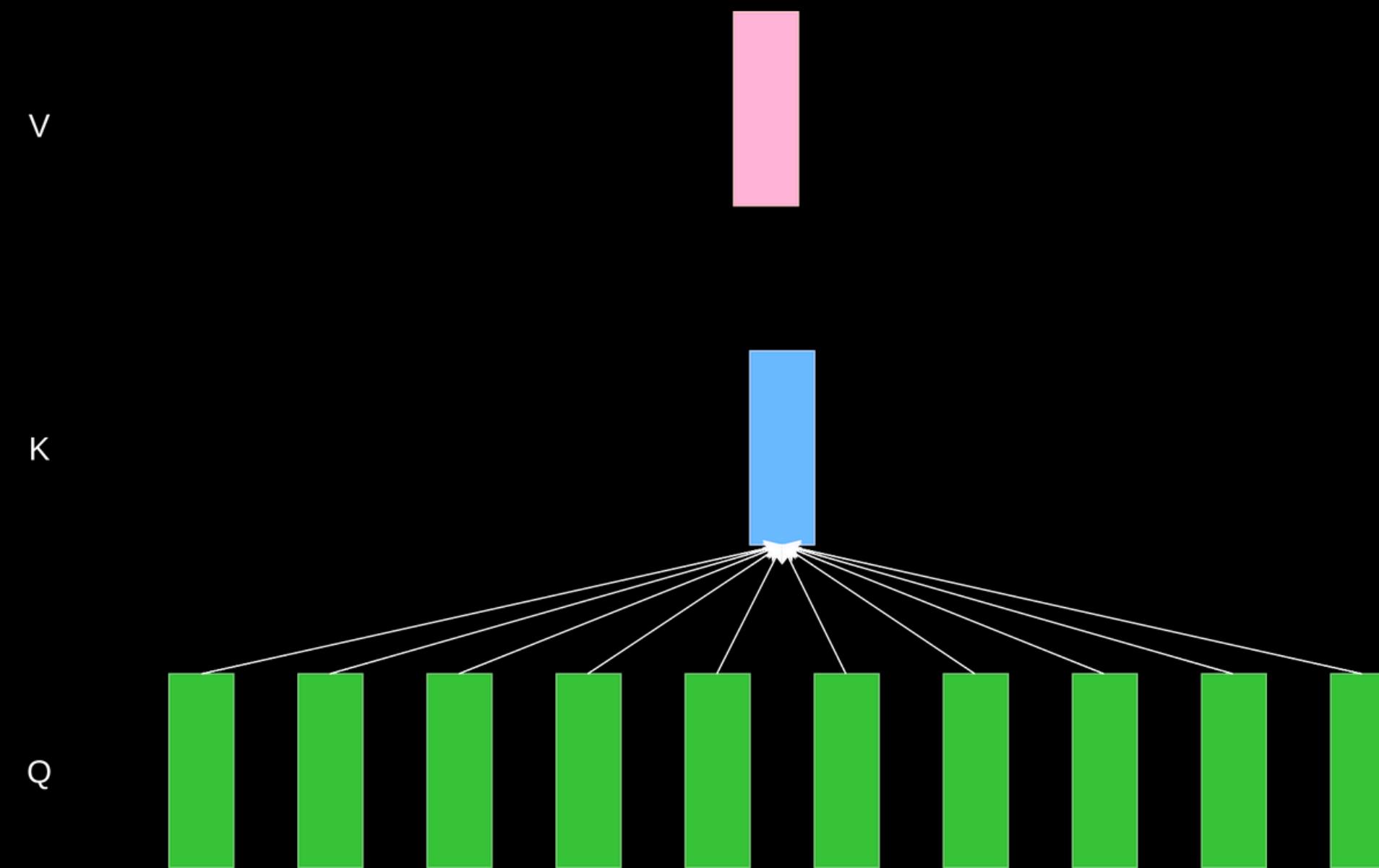
$$\text{Attention}(q, k, v) = \text{softmax} \left(\frac{\text{Query Projection } \widehat{Q} \quad \text{Key Projection } \widehat{K}^\top}{\underbrace{\sqrt{d_k}}_{\text{Scaling factor}}} \right) \text{Value Projection } \widehat{V}$$



Grouped Query Attention

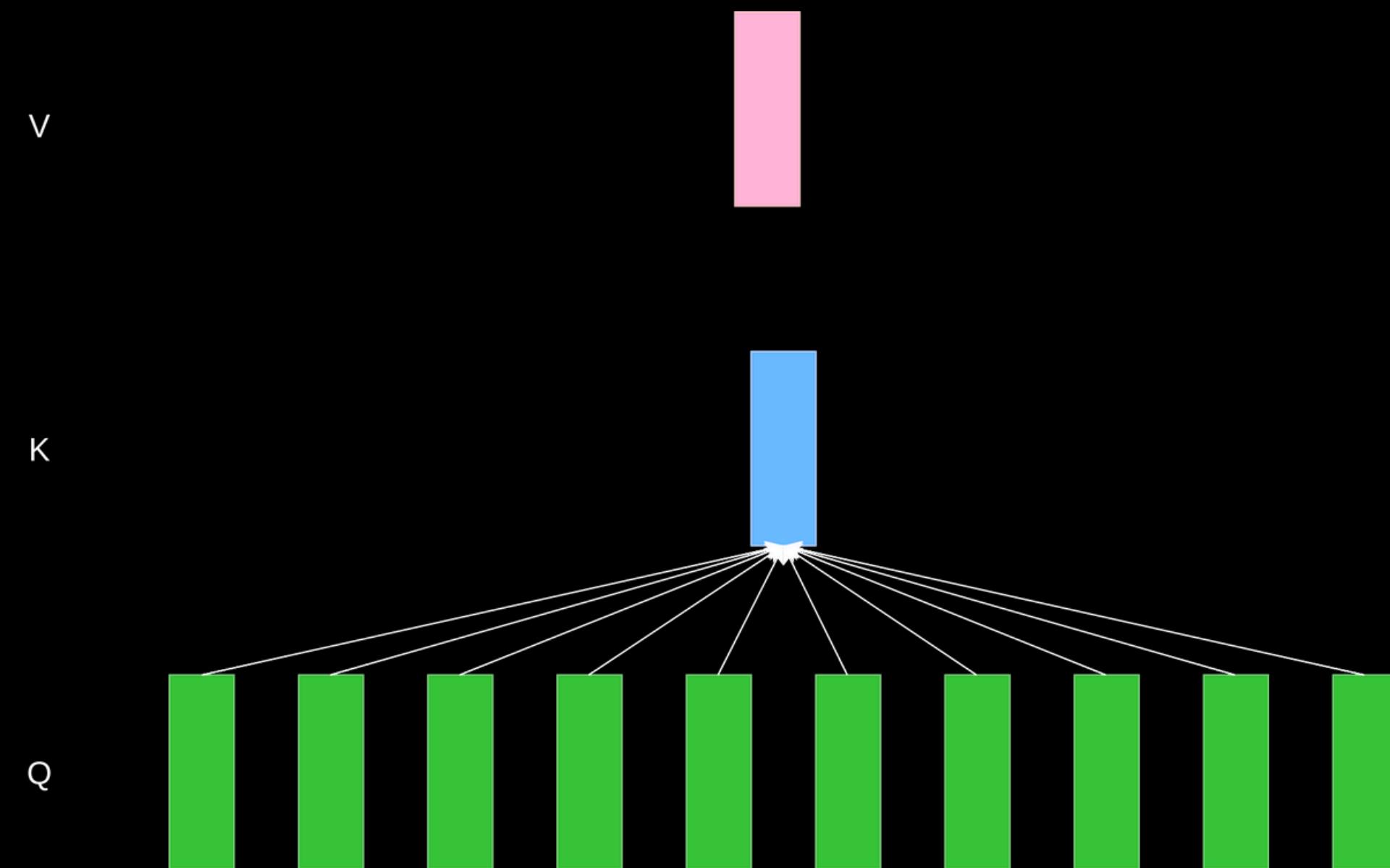


Multi Query Attention



Multi Query Attention

Do you think this reduces the quality of generation?



Section Summary

We saw

1. Various architectures for our LLM
2. Positional Embeddings and how to scale beyond 3000 tokens
3. Attention mechanisms

Questions

Conclusion

- Pretraining a LLM is a daunting task
- Scaling laws need to be derived for optimal compute usage.
- Data quality is very important!
- Choosing a good architecture is important as it determines how much compute you use!

Thank you for
listening