

I. Intro

mapping \rightarrow encoding \rightarrow decoding

Google uses sentence piece (encoder)

- Sub-word units

GPT uses tik token

batch dimension (mini-batches)

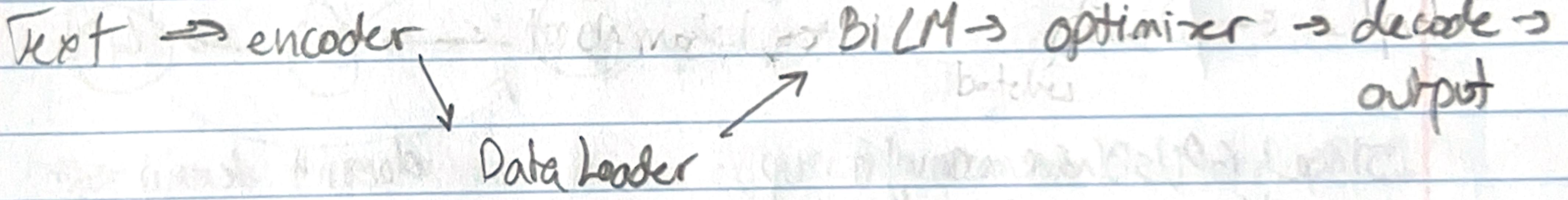
logits = scores

loss = F. cross-entropy (logits, targets)

Bigram Language Model \star - look at the previous character

logits loss

II. Diagram



III. Cool Math

Data should not be influenced by future, just past.

inefficient: take the rolling averages of tensor using torch.mean()

efficient: torch.tril() \rightarrow $\begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \end{matrix}$

$\begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{matrix} \Rightarrow \begin{matrix} 1 & 1 & 0 \\ 1 & 1 & 1 \end{matrix}$

normalized(torch.tril(a)) @ b

↑
weighted Sums

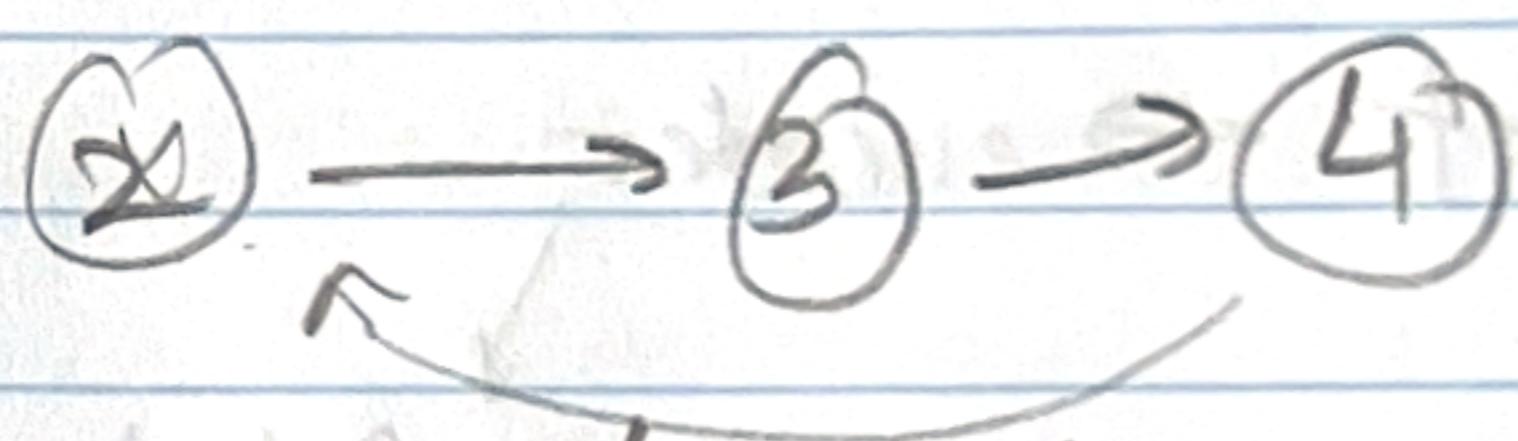
more efficient: use softmax

softmax → turns attention scores to probabilities
softmax is a normalization technique

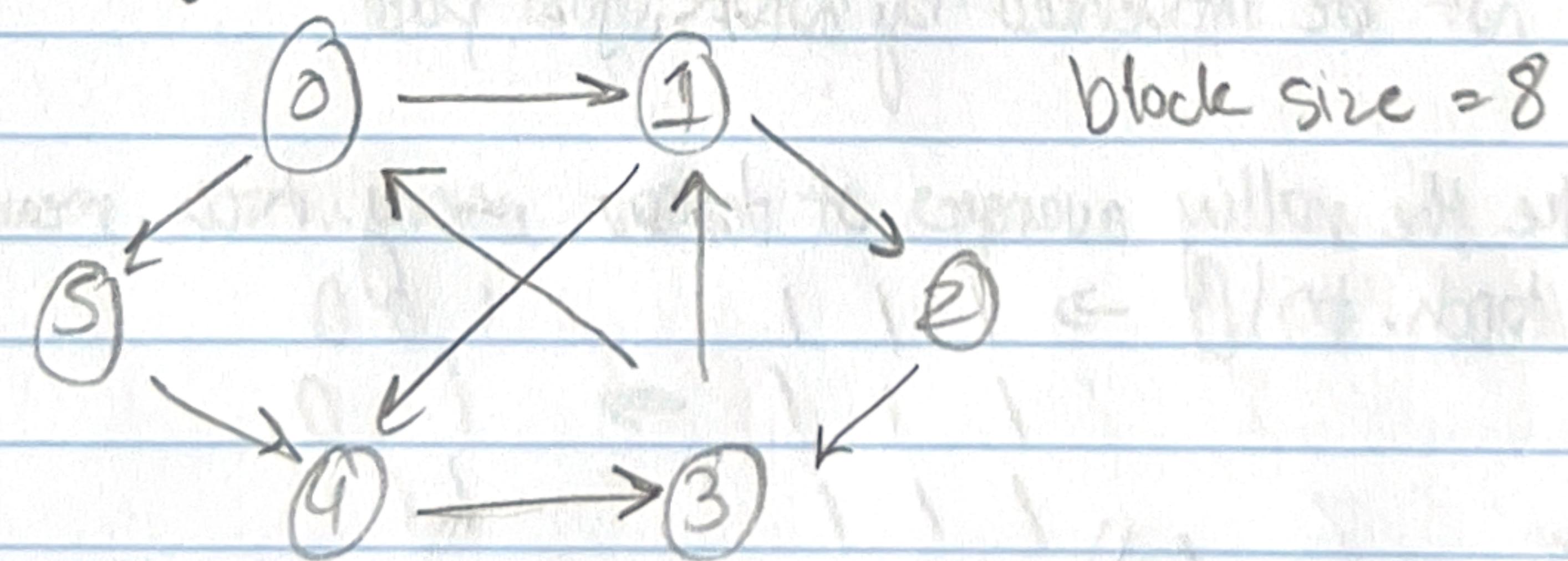
token → query (What am I looking for)
key (What do I contain) • (dot product) = learn a lot
attention scores

Transformers vs RNN

- full context
- all tokens processed at the same time
- works best w/ large datasets
- more memory
- sequential
- fuzzy and noisy
- vanishing gradients
- can only look at values thru hidden state propagation



since GPUs have caught up, doesn't know how they use this



Batch dimensions don't talk to each other

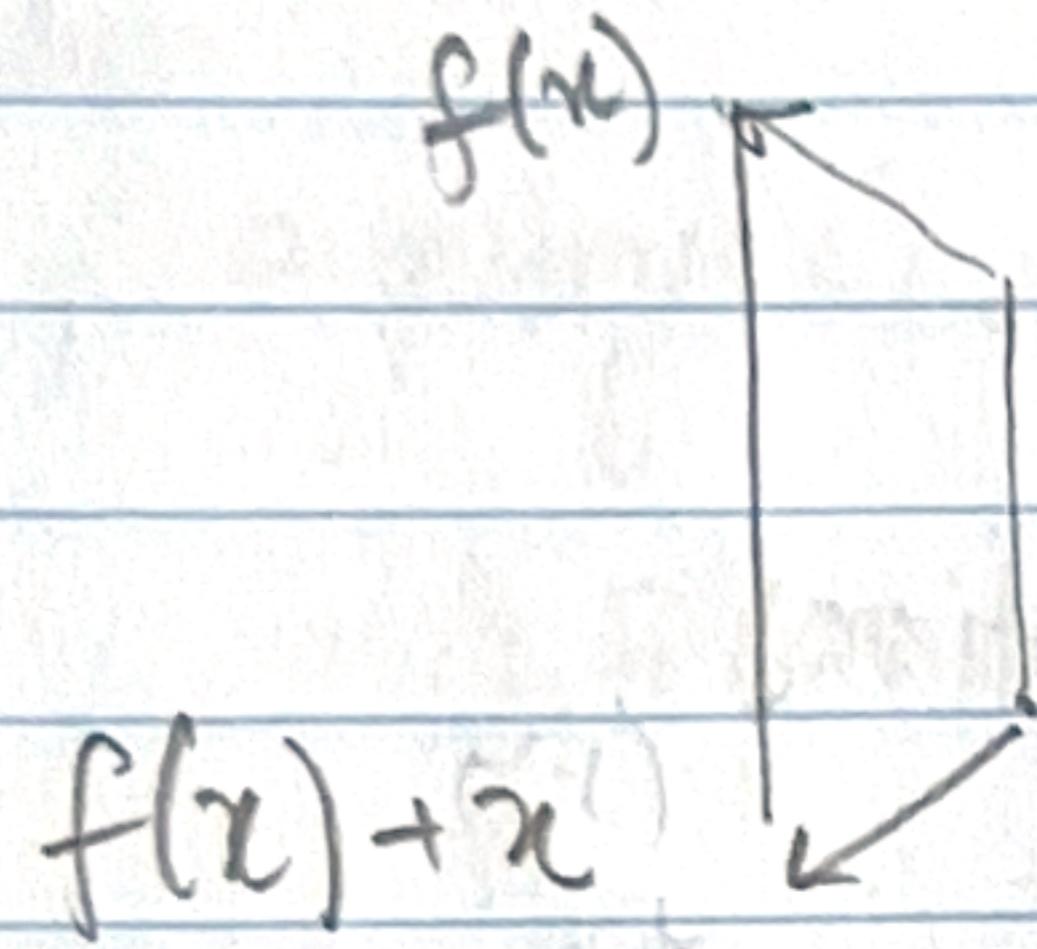
Cross attention - if you get info from another set of keys

queries come from x but keys and vals come from other source

Feed Forward: Input \rightarrow Hidden layers \rightarrow Output

For transformers: applied to each token, and helps model interpret what attention gave it

Residual connections



Layer Norm: normalizes values of vector mean of 0, std of 1

Attention \rightarrow Feed Forward \rightarrow Layer Norm

Dropout = reduces generality, drops neural net links

Gelu

pretraining: decoder

PGPT training

train w/ supervised \rightarrow collect comparisons and train reward \rightarrow \rightarrow optimize usg PPO

Proximal Policy Optimization - Slow and steady reinforcement

EX 1. C

Goal m

transformer

- Self attention
- feed forward
- Layer Norm + Residuals

Old

I. Input Embeddings + Positional Encoding
Each token is turned into a vector, add positional encodings

for head

or

final =

II. Self Attention

$$\text{attention} = \text{softmax} \left(\frac{Q \times K^t}{\sqrt{d_k}} \right) \times V$$

New

q =

III. MultiHead Attention

Repeat self attention in parallel multiple times

i. Q =

K =

V =

Concatenate all heads \rightarrow linear projection

IV. Feed Forward, multilayer perception

Before attention, MLP, token's vector goes through 2-layer MLP
Transforms info it got

ii. R

V. Layer Norm + Residuals

Before attention/MLP, Layer Norm normalizes vector

iii. S

Adds input back to output (Residual Connection)

Prevents gradient problems

VI. Linear \rightarrow Softmax to predict next token

EX1. Combine Head and Multithead attention into one class that process all the heads in parallel

Goal multihead operation becomes batch matrix operation

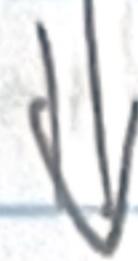
Old

for head in heads:

out = head(x)

outputs.append(out)

final = torch.cat(outputs, dim = -1)



New

$q = \text{query}(x).view(B, T, n_head, head_size).transpose(1, 2)$

$\# \rightarrow (B, n_head, T, head_size)$

i. $Q = x W^Q$

$K = x W^K \rightarrow$ All shape $R^{B \times T \times C}$

$V = x W^V$

ii. Reshape / Transform group by head

$Q \rightarrow R^{B \times n_head \times T \times d_k}$

iii. $\text{scores} = \frac{QK^T}{\sqrt{d_k}} \in R^{B \times n_head \times T \times T}$

output = softmax(scores) $\cdot V$ ← value

iv.

if, $W^Q, W^K, W^V \in \mathbb{R}^{C \times C}$

$$\text{Attention}(x) = \text{Concat}_{h=1}^n \left(\text{softmax} \left(\frac{xW_h^Q \times (xW_h^K)^T}{\sqrt{d_k}} \cdot xW_h^V \right) \right)$$

RoPE

positional encodings tell where the word is

- i. Absolute pos enc = $e_{\text{final}}(\text{pos}) = e_{\text{token}} + e_{\text{position}}$
issues: given short, irrelevant training words not good.
don't need to know the 39th word if you are 2 words after despite

ii. Relative pos enc

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

Too many relationships, very slow, Key-Value cache, not commonly used

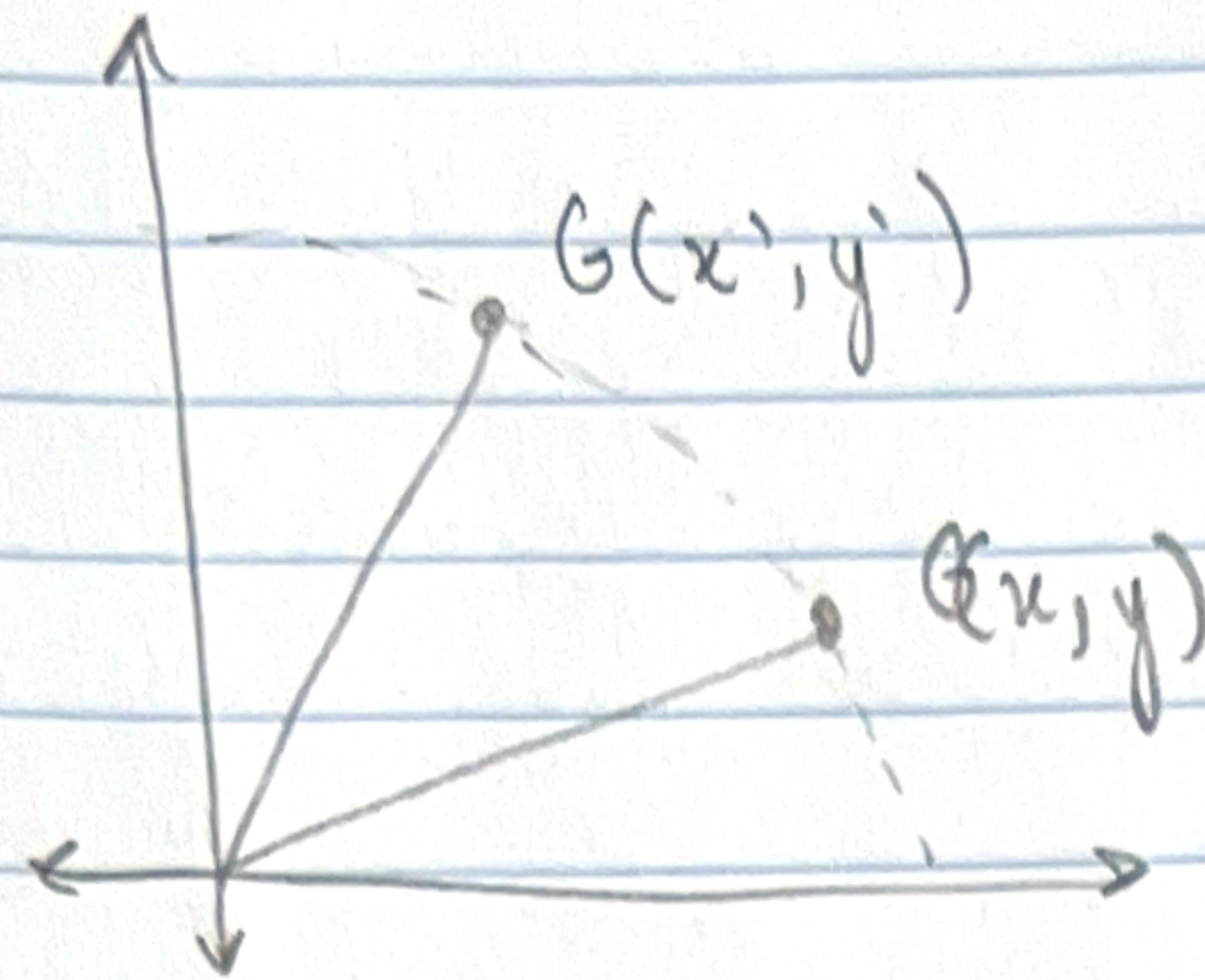
iii. RoPE!

- encodes positional info so it allows the model understand absolute and relative pos enc. does a rotational mechanism

The quick brown fox jumps
constant

more rotations →

Lower dimensions rotate faster higher dim rotate slower



$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

RoPE can handle sequences longer than it was trained on

Easy to cache next word

$$f_{\text{Eq}, k_2}(x_m, m) = R_{\theta, m}^d W_{\text{Eq}, k_2} x_m$$

\Rightarrow n-dimension corkscrew visualization