



# 10-423/10-623 Generative AI

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

## Pretraining vs. finetuning + Modern Transformers (RoPE, GQA, Longformer) + CNNs

Matt Gormley & Henry Chai

Lecture 4

Sep. 9, 2024

# Reminders

- **Homework 0: PyTorch + Weights & Biases**
  - Out: Wed, Aug 28
  - Due: Mon, Sep 9 at 11:59pm
  - unique policy for this assignment: we will grant (essentially) any and all extension requests
- **Quiz 1: Wed, Sep 11**
- **Homework 1: Generative Models of Text**
  - Out: Mon, Sep 9
  - Due: Mon, Sep 23 at 11:59pm

# Recap So Far

## Deep Learning

- AutoDiff
  - is a tool for **computing gradients** of a differentiable function,  $b = f(a)$
  - the key building block is a **module** with a `forward()` and `backward()`
  - sometimes define `f` as **code** in `forward()` by chaining existing modules together
- Computation Graphs
  - are another way to define `f` (more conducive to slides)
  - so far, we saw two (deep) computation graphs
    - 1) RNN-LM
    - 2) Transformer-LM
    - (Transformer-LM was kind of complicated)

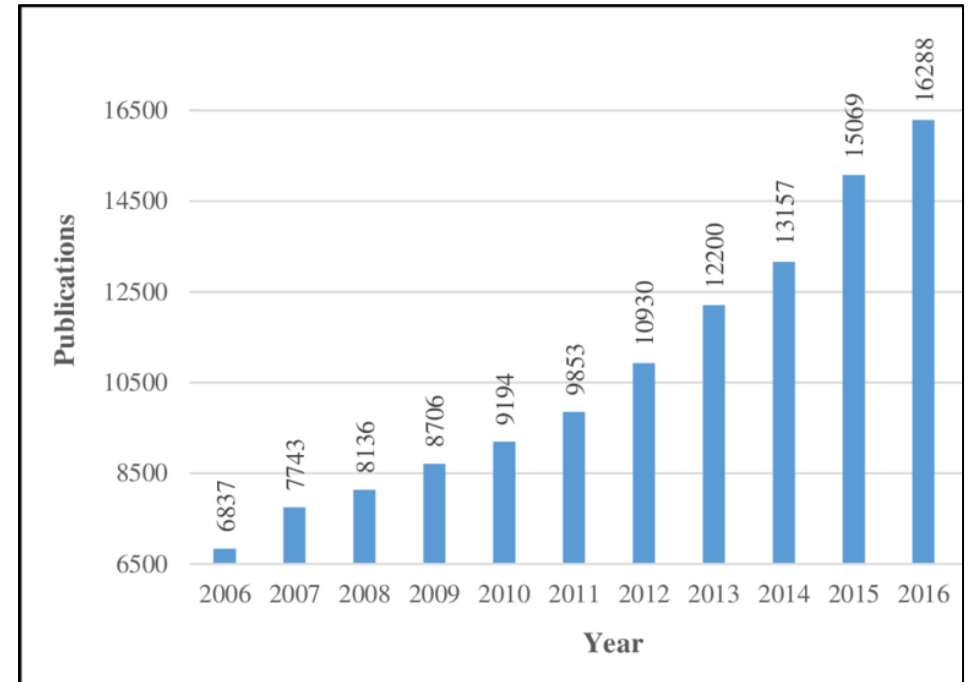
## Language Modeling

- key idea: condition on previous words to **sample the next word**
- to define the **probability** of the next word...
  - ... n-gram LM uses collection of massive 50k-sided **dice**
  - ... RNN-LM or Transformer-LM use a **neural network**
- Learning an LM
  - n-gram LMs are easy to learn: just **count** co-occurrences!
  - a RNN-LM / Transformer-LM is trained by **optimizing an objective function with SGD; compute gradients with AutoDiff**

# **PRE-TRAINING VS. FINE-TUNING**

# The Start of Deep Learning

- The architectures of modern deep learning have a long history:
  - 1960s: Rosenblatt's 3-layer multi-layer perceptron, ReLU )
  - 1970-80s: RNNs and CNNs
  - 1990s: linearized self-attention
- The spark for deep learning came in 2006 thanks to **pre-training** (e.g., Hinton & Salakhutdinov, 2006)



# Deep Network Training

- **Idea #1:**

1. Supervised fine-tuning only

- **Idea #2:**

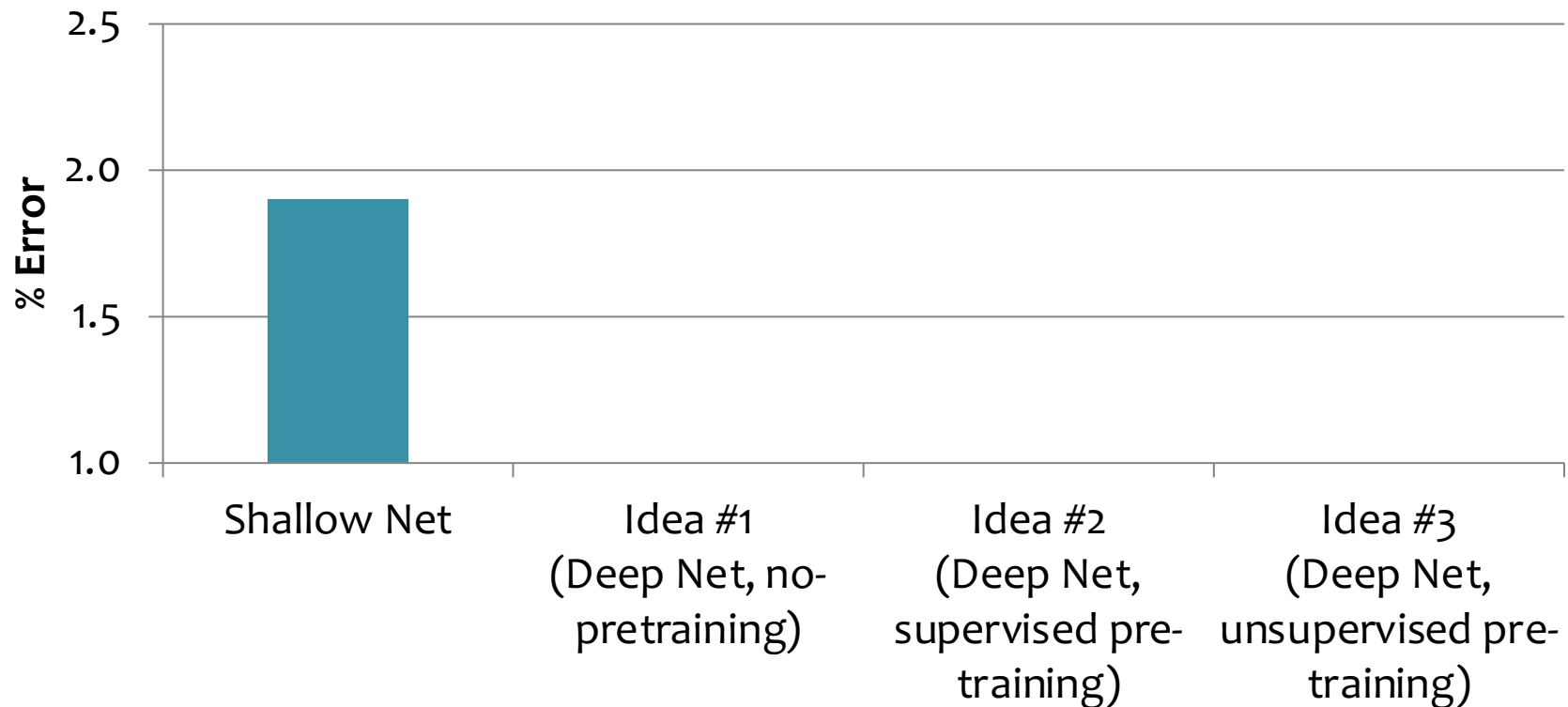
1. Supervised layer-wise pre-training
2. Supervised fine-tuning

- **Idea #3:**

1. Unsupervised layer-wise pre-training
2. Supervised fine-tuning

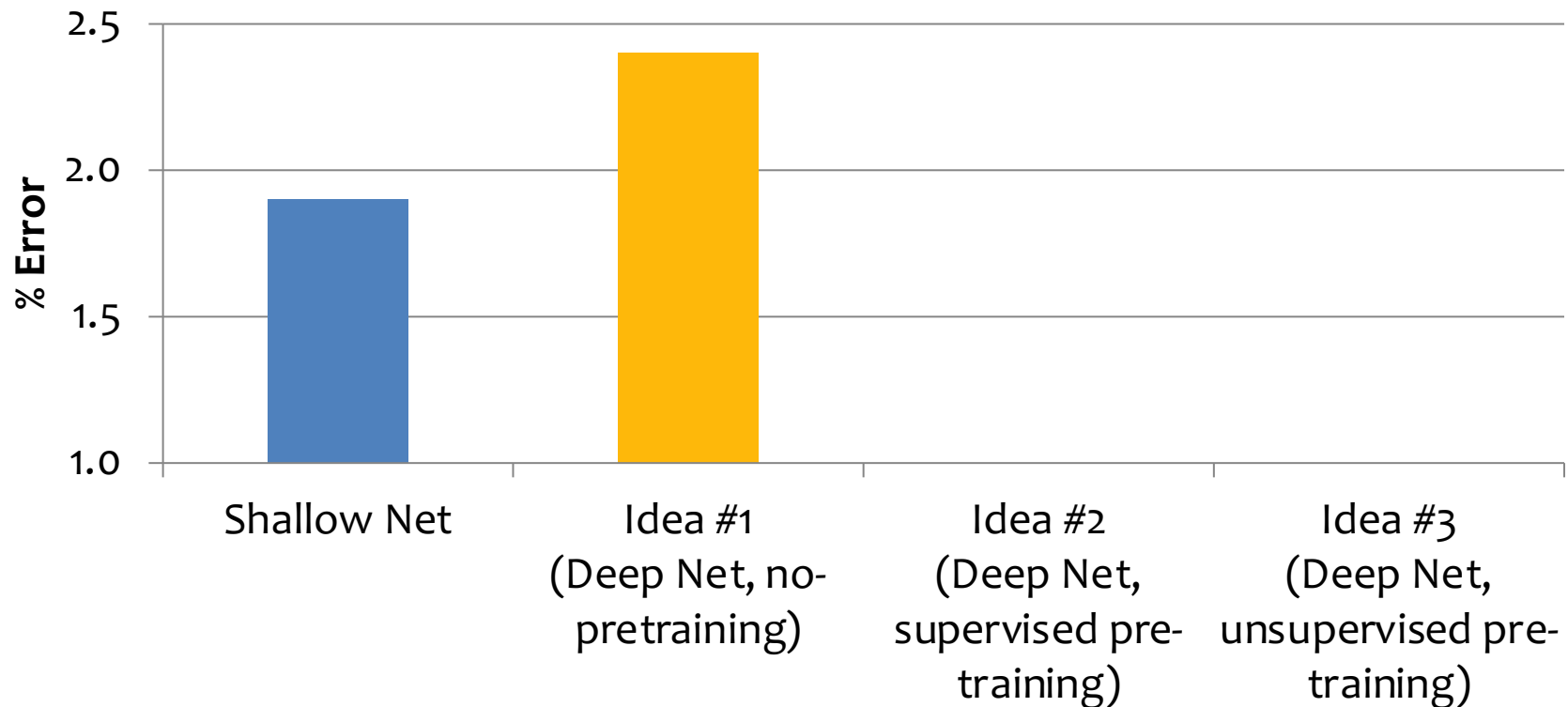
# Comparison on MNIST

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# Comparison on MNIST

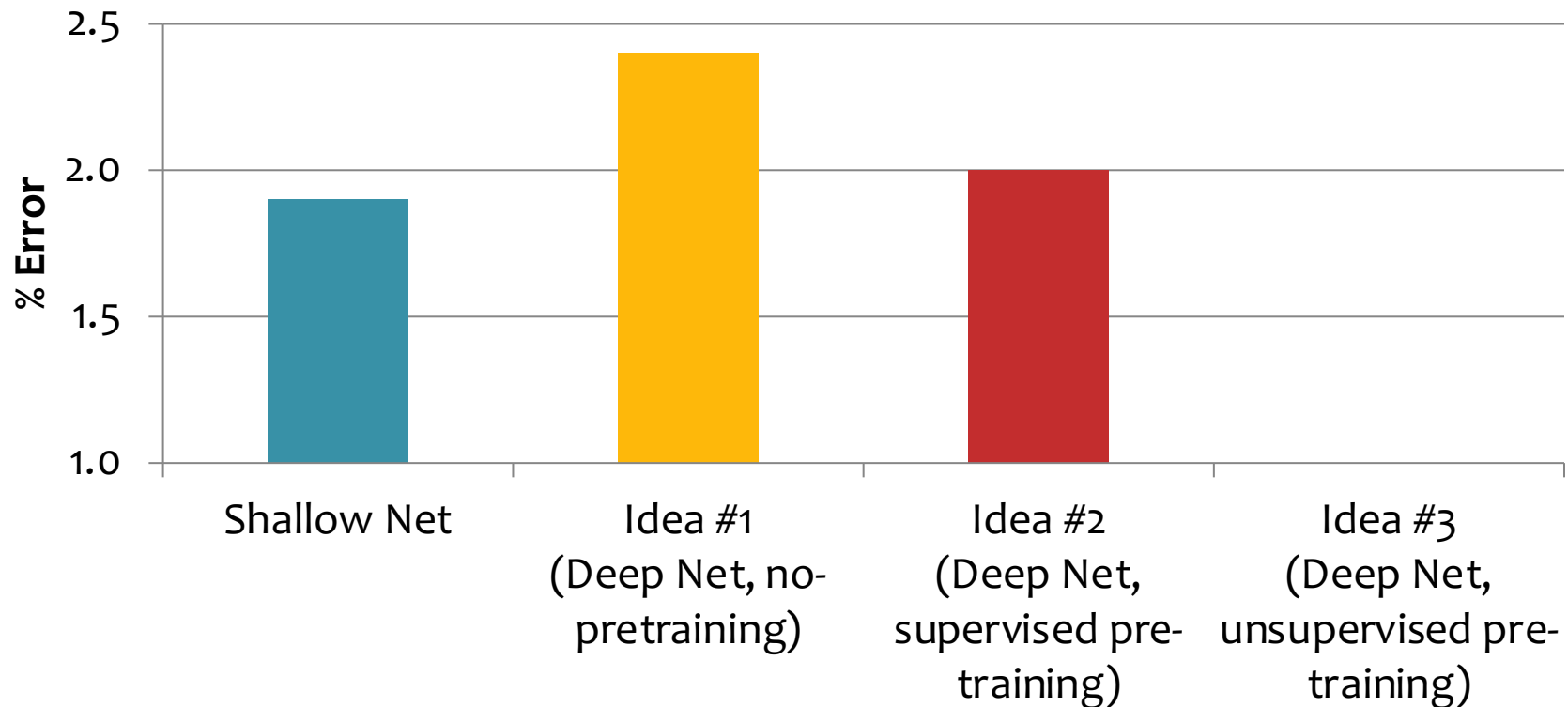
- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)





# Comparison on MNIST

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# Idea #3: Unsupervised Pre-training

- **Idea #3: (Two Steps)**

- Use our original idea, but **pick a better starting point**
- **Train each level** of the model in a **greedy way**

1. Unsupervised Pre-training

- Use **unlabeled** data
- Work bottom-up
  - Train hidden layer 1. Then fix its parameters.
  - Train hidden layer 2. Then fix its parameters.
  - ...
  - Train hidden layer n. Then fix its parameters.

2. Supervised Fine-tuning

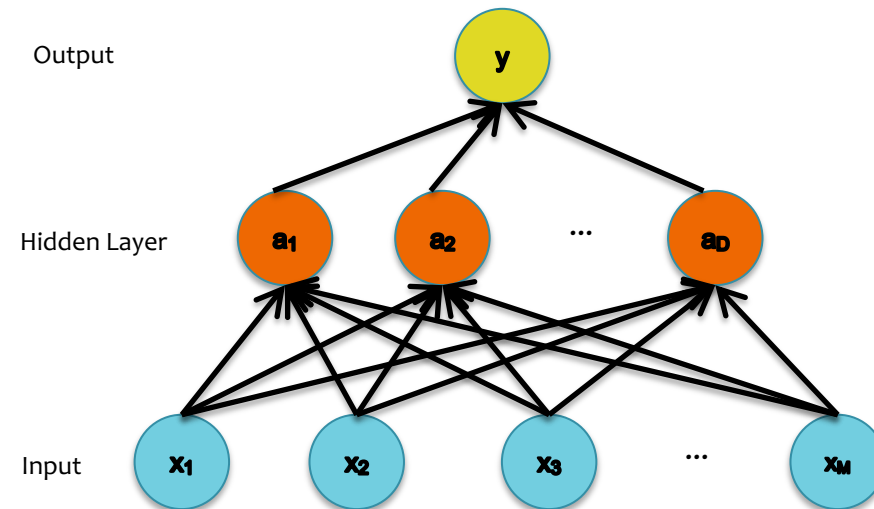
- Use **labeled** data to train following “Idea #1”
- Refine the features by backpropagation so that they become tuned to the end-task

# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training of the first layer:

- What should it predict?
- What else do we observe?
- **The input!**



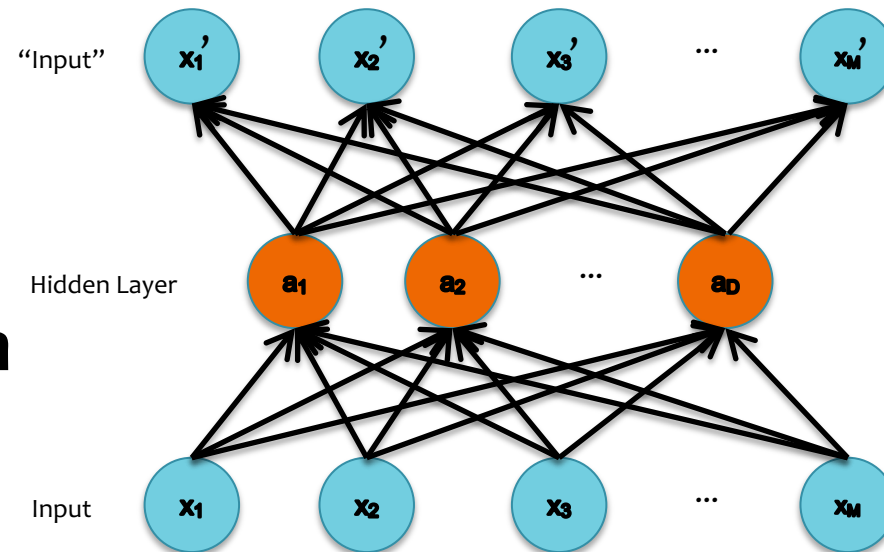
# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training of the first layer:

- What should it predict?
- What else do we observe?
- **The input!**

**This topology defines an Auto-encoder.**



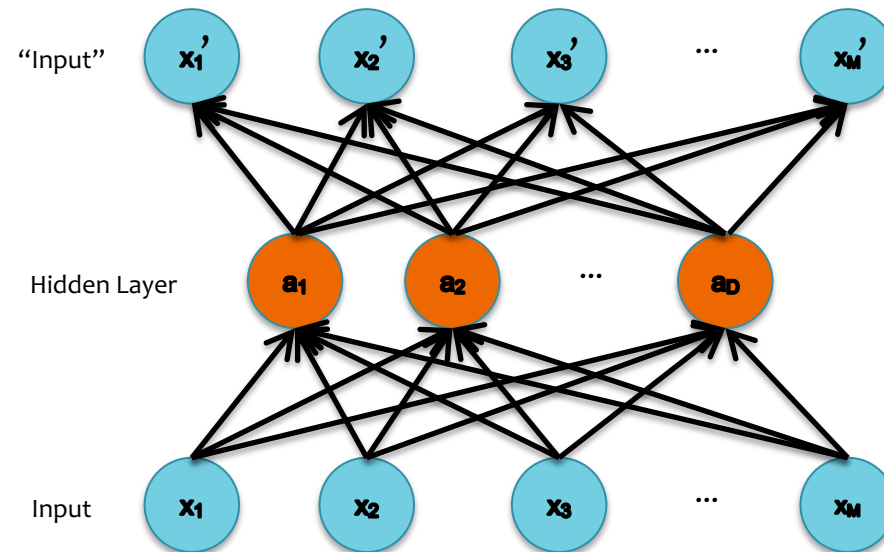
# Auto-Encoders

Key idea: Encourage  $z$  to give small reconstruction error:

- $x'$  is the *reconstruction* of  $x$
- Loss =  $\|x - \text{DECODER}(\text{ENCODER}(x))\|^2$
- Train with the same backpropagation algorithm for 2-layer Neural Networks with  $x_m$  as both input and output.

$$\text{DECODER: } x' = h(W'z)$$

$$\text{ENCODER: } z = h(Wx)$$

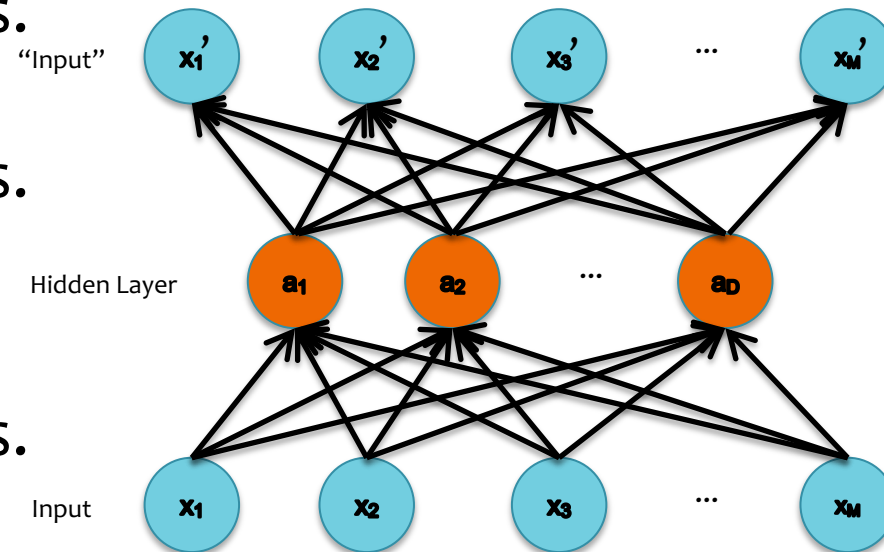


# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training

- Work bottom-up
  - Train hidden layer 1. Then fix its parameters.
  - Train hidden layer 2. Then fix its parameters.
  - ...
  - Train hidden layer n. Then fix its parameters.

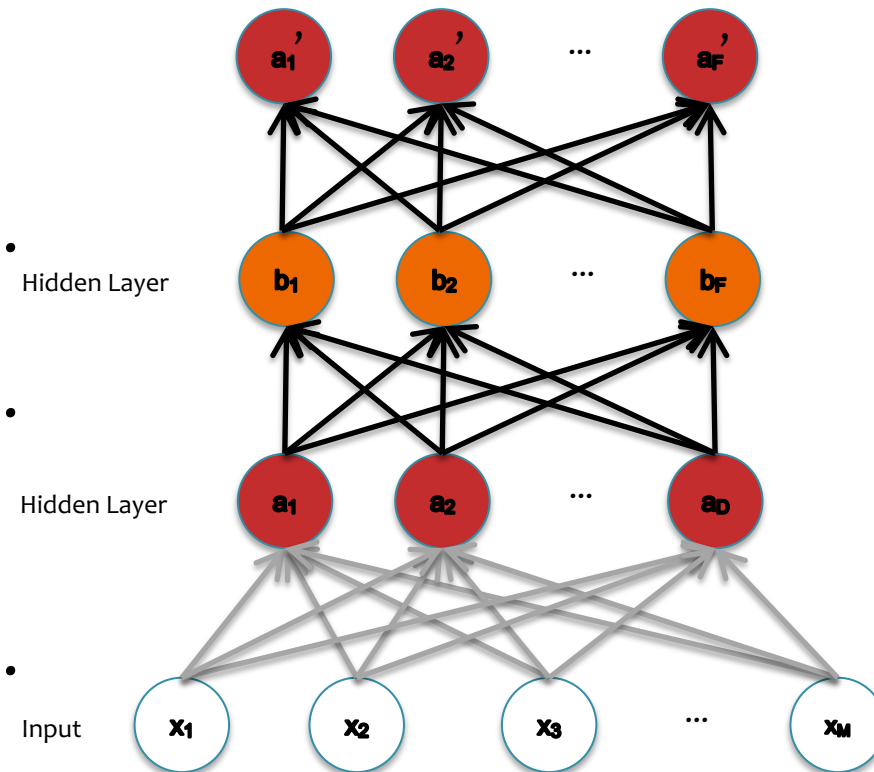


# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training

- Work bottom-up
  - Train hidden layer 1. Then fix its parameters.
  - Train hidden layer 2. Then fix its parameters.
  - ...
  - Train hidden layer n. Then fix its parameters.

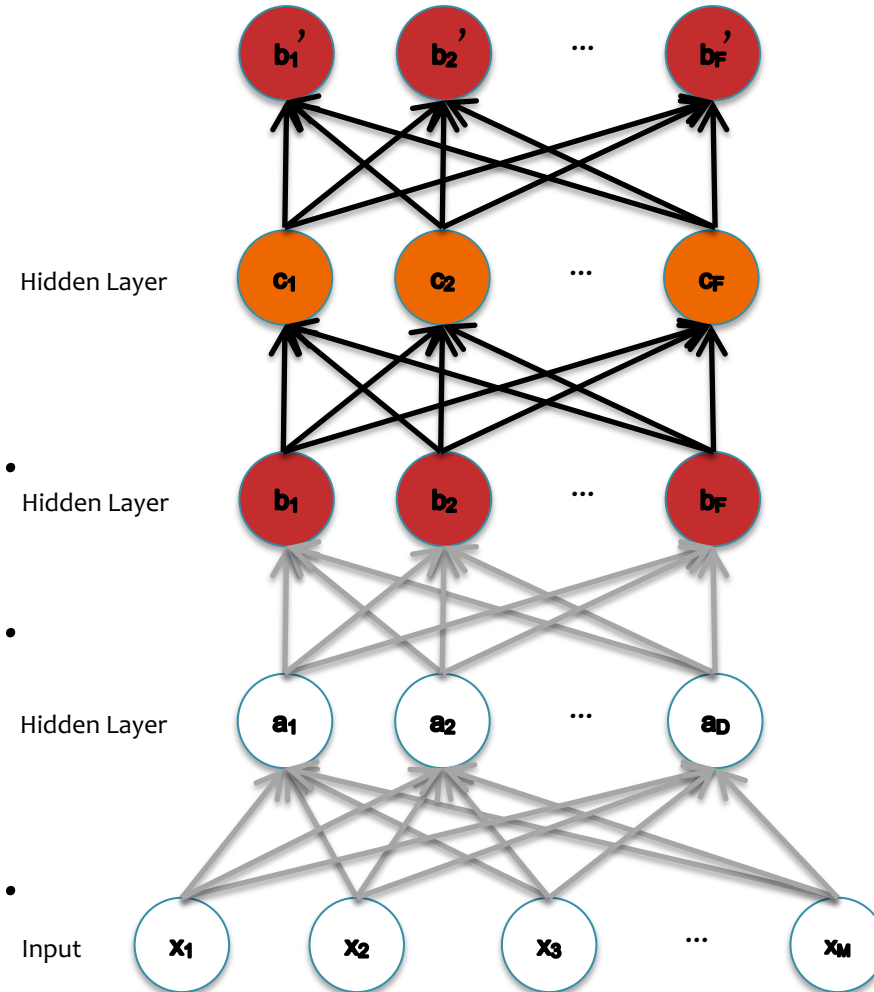


# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training

- Work bottom-up
  - Train hidden layer 1. Then fix its parameters.
  - Train hidden layer 2. Then fix its parameters.
  - ...
  - Train hidden layer n. Then fix its parameters.





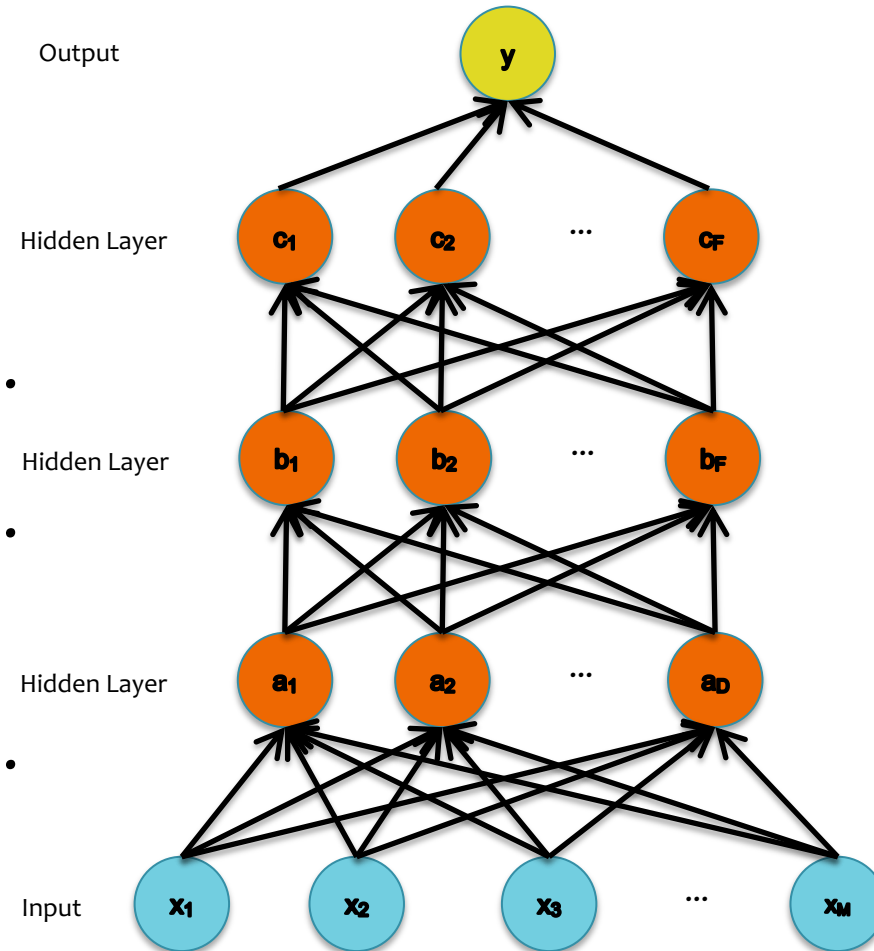
# The solution:

## *Unsupervised pre-training*

### Unsupervised pre-training

- Work bottom-up
  - Train hidden layer 1. Then fix its parameters.
  - Train hidden layer 2. Then fix its parameters.
  - ...
  - Train hidden layer  $n$ . Then fix its parameters.

**Supervised fine-tuning**  
Backprop and update all parameters



# Deep Network Training

- **Idea #1:**

1. Supervised fine-tuning only

- **Idea #2:**

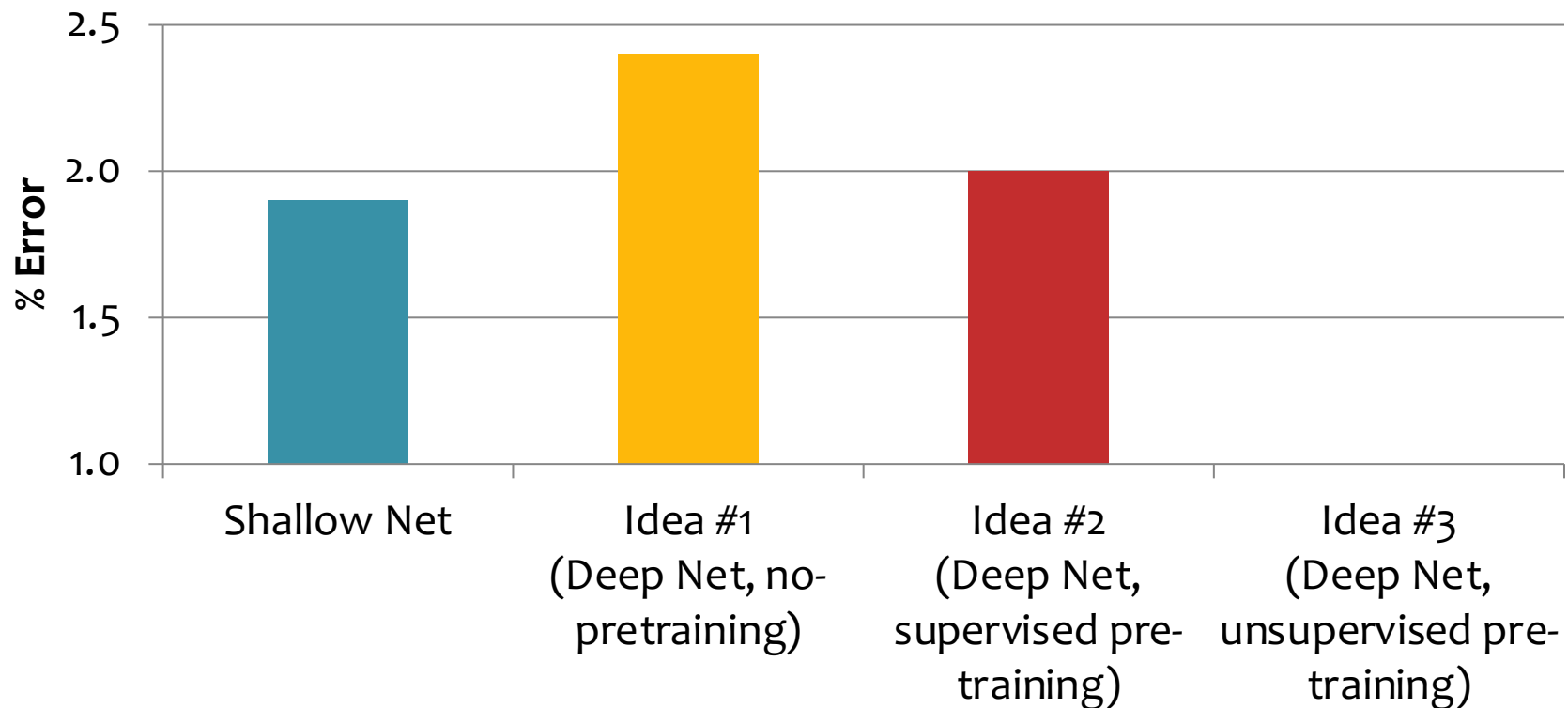
1. Supervised layer-wise pre-training
2. Supervised fine-tuning

- **Idea #3:**

1. Unsupervised layer-wise pre-training
2. Supervised fine-tuning

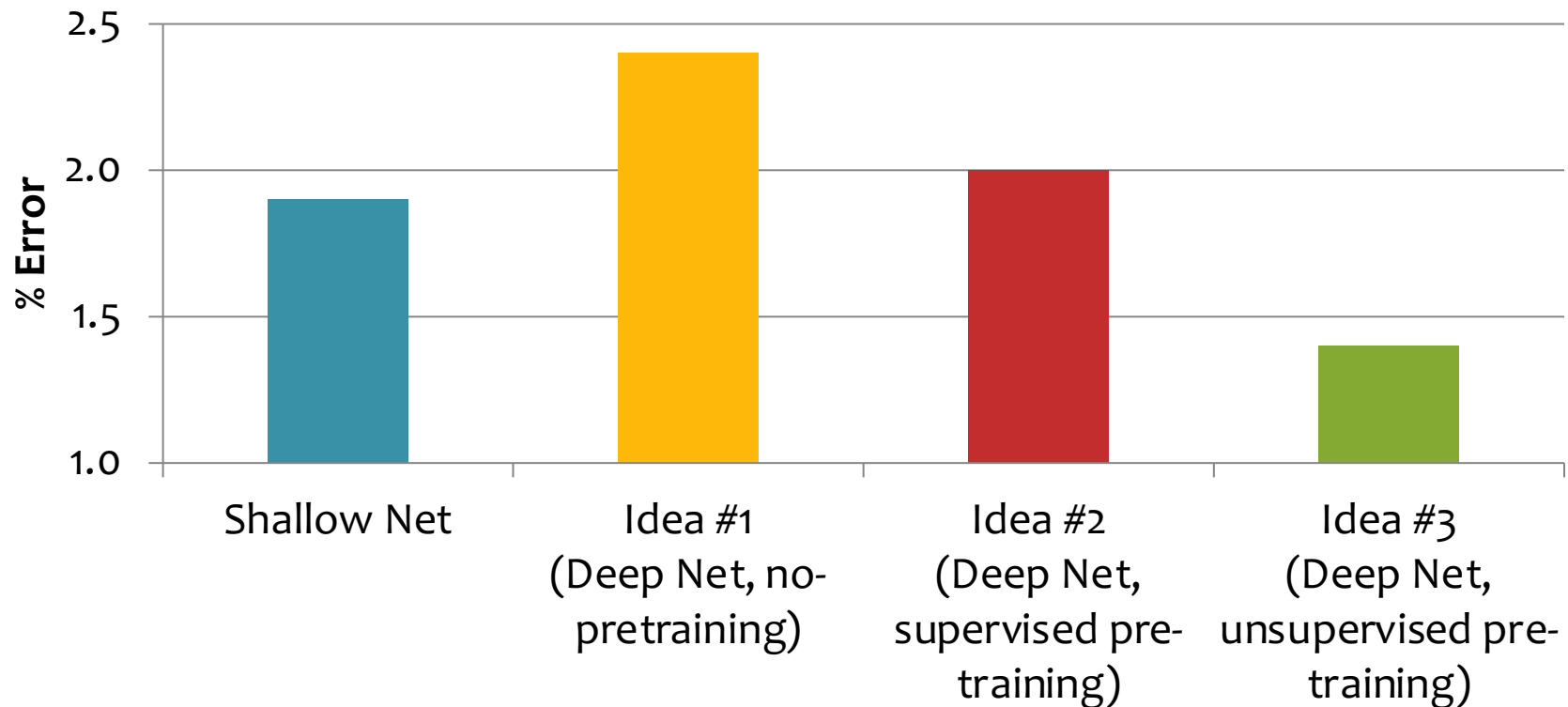
## Comparison on MNIST

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)

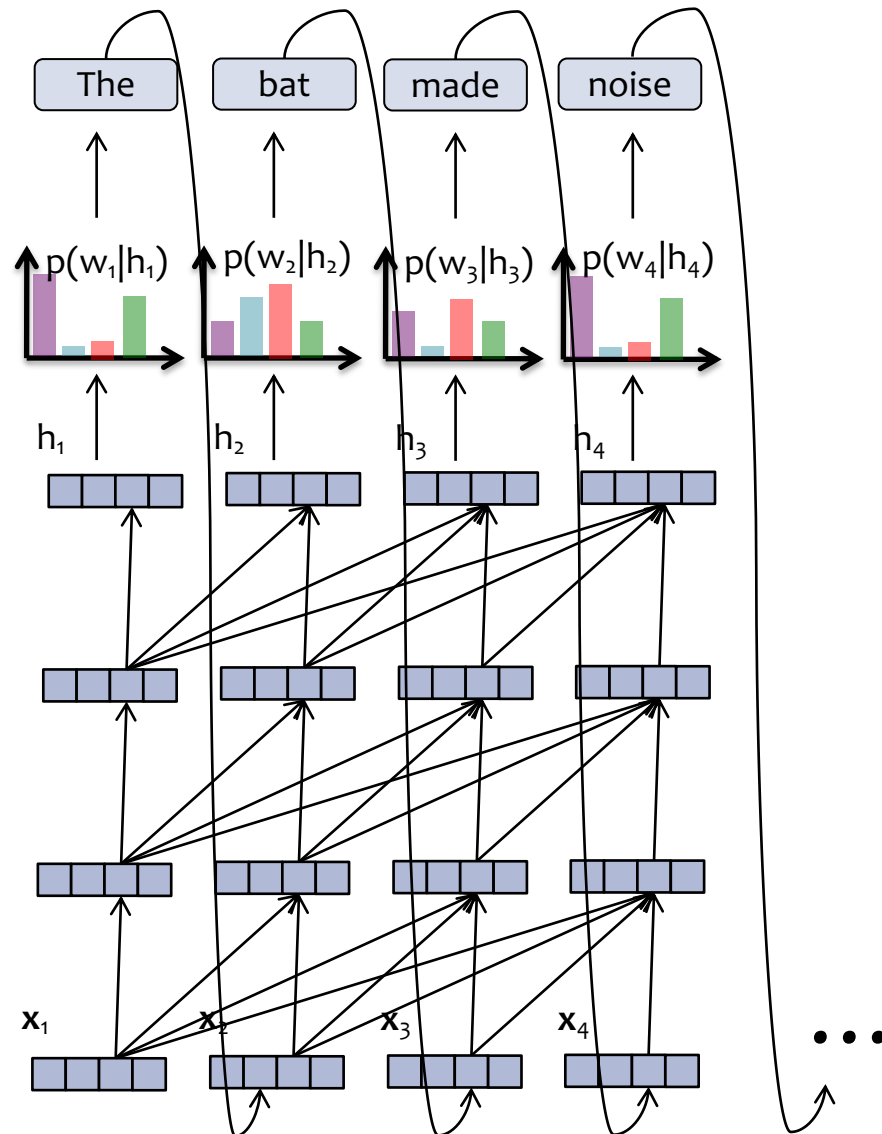


## Comparison on MNIST

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# Transformer Language Model



**Generative pre-training** for a deep language model:

- each training example is an (unlabeled) sentence
- the objective function is the likelihood of the observed sentence

Practically, we can **batch** together many such training examples to make training more efficient

# Training Data for LLMs

## GPT-3 Training Data:

Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

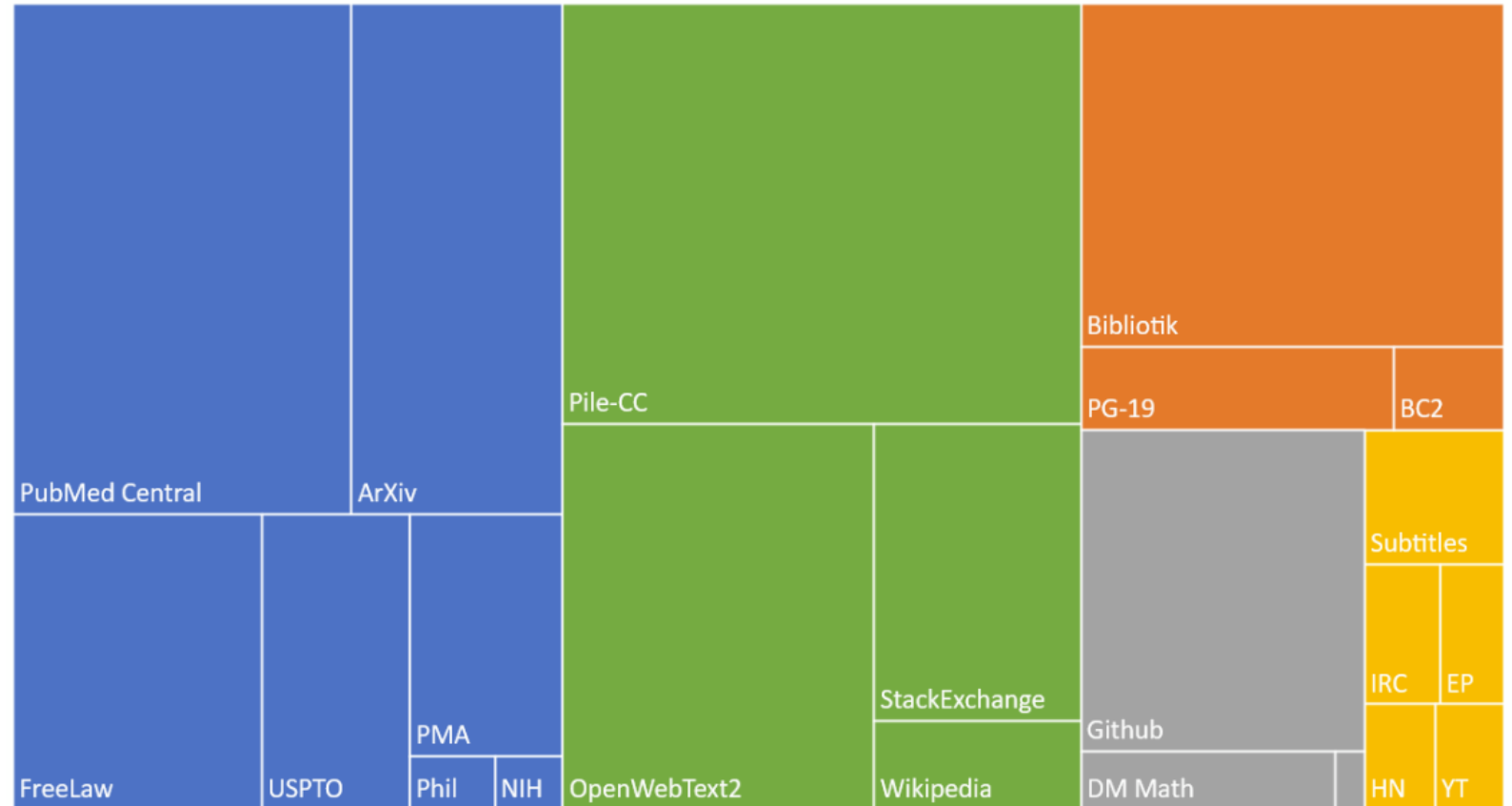
# Training Data for LLMs

## The Pile:

- An open source dataset for training language models
- Comprised of 22 smaller datasets
- Favors high quality text
- 825 Gb  $\approx$  1.2 trillion tokens

Composition of the Pile by Category

■ Academic ■ Internet ■ Prose ■ Dialogue ■ Misc



# **MODERN TRANSFORMER MODELS**



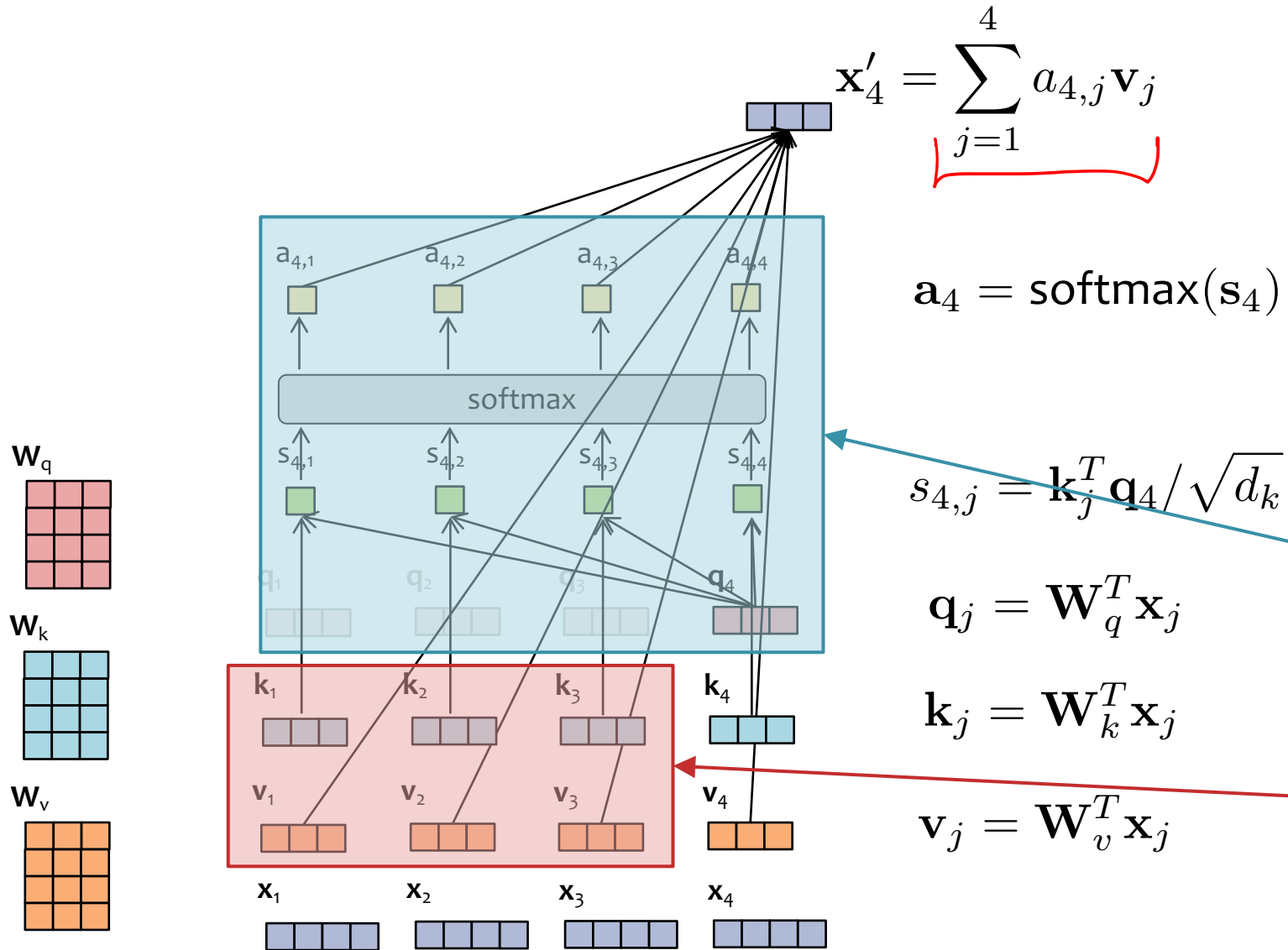
# Modern Transformer Models

- PaLM (Oct 2022)
  - 540B parameters
  - closed source
  - Model:
    - SwiGLU instead of ReLU, GELU, or Swish
    - **multi-query attention (MQA)** instead of multi-headed attention
    - **rotary position embeddings**
    - **shared input-output embeddings** instead of separate parameter matrices
  - Training: **Adafactor** on 780 billion tokens
- Llama-1 (Feb 2023)
  - collection of models of varying parameter sizes: 7B, 13B, 32B, 65B
  - semi-open source
  - Llama-13B outperforms GPT-3 on average
  - Model compared to GPT-3:
    - **RMSNorm** on inputs instead of LayerNorm on outputs
    - **SwiGLU** activation function instead of ReLU
    - **rotary position embeddings (RoPE)** instead of absolute
  - Training: **AdamW** on 1.0 – 1.4 trillion tokens
- Falcon (June - Nov 2023)
  - models of size 7B, 40B, 180B
  - first fully open source model, Apache 2.0
  - Model compared to Llama-1:
    - (GQA) instead of multi-headed attention (MHA) or **grouped query attention multi-query attention (MQA)**
    - **rotary position embeddings** (worked better than Alibi)
    - **GeLU** instead of SwiGLU
  - Training: AdamW on up to 3.5 trillion tokens for 180B model, using **z-loss** for stability and **weight decay**
- Llama-2 (Aug 2023)
  - collection of models of varying parameter sizes: 7B, 13B, 70B.
  - introduced Llama 2-Chat, fine-tuned as a dialogue agent
  - Model compared to Llama-1:
    - **grouped query attention (GQA)** instead of multi-headed attention (MHA)
    - context length of 4096 instead of 2048
  - Training: **AdamW** on 2.0 trillion tokens
- Mistral 7B (Oct 2023)
  - outperforms Llama-2 13B on average
  - introduced Mistral 7B – Instruct, fine-tuned as a dialogue agent
  - truly open source: Apache 2.0 license
  - Model compared to Llama-2
    - **sliding window attention** (with  $W=4096$ ) and grouped-query attention (GQA) instead of just GQA
    - context length of 8192 instead of 4096 (can generate sequences up to length 32K)
    - **rolling buffer cache** (grow the KV cache and the overwrite position  $i$  into position  $i \bmod W$ )
  - variant Mixtral offers a **mixture of experts** (roughly 8 Mistral models)

In this section we'll look at four techniques:

1. key-value cache (KV cache)
2. rotary position embeddings (RoPE)
3. grouped query attention (GQA)
4. sliding window attention

# Key-Value Cache



- At each timestep, we reuse all previous keys and values (i.e. we need to cache them)
- But we can get rid of the queries, similarity scores, and attention weights (i.e. we can let them fall out of the cache)

Discarded after this timestep

Computed for previous time-steps and reused for this timestep

# ROTARY POSITION EMBEDDINGS (ROPE)

# Rotary Position Embeddings (RoPE)

**Q:** Why does this slide have so many typos?

**A:** I'm really not sure. I very meticulously type up the latex for my slides myself and think carefully about all the things I put in them.

RoPE attention:

$$f_q(\mathbf{x}_t, m) \triangleq \mathbf{R}_{\Theta, m} \mathbf{W}_q^T \mathbf{x}_t$$

$$f_k(\mathbf{x}_j, m) \triangleq \mathbf{R}_{\Theta, m} \mathbf{W}_k^T \mathbf{x}_j$$

$$s_{t,j} = f_k(\mathbf{x}_j, m)^T f_q(\mathbf{x}_t, m) / \sqrt{|\mathbf{k}|},$$

$$\forall j, t \text{ where } m = t - j$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

where  $\mathbf{W}_k, \mathbf{W}_q \in \mathbb{R}^{d_{model} \times d_k}$ , and the rotary matrix  $\mathbf{R}_{\Theta, m} \in \mathbb{R}^{d_k \times d_k}$  is given by:

$$R_{\Theta, m} = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \dots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos m\theta_{d_k/2} & -\sin m\theta_{d_k/2} \\ 0 & 0 & 0 & 0 & \dots & \sin m\theta_{d_k/2} & \cos m\theta_{d_k/2} \end{pmatrix}$$

The  $\theta_i$  parameters are fixed ahead of time and defined as below

$$\Theta = \{\theta_i = 10000^{-2^{i-1}/d}, i \in [1, 2, \dots, d/2]\}$$

wrong

wrong

wrong

wrong

wrong

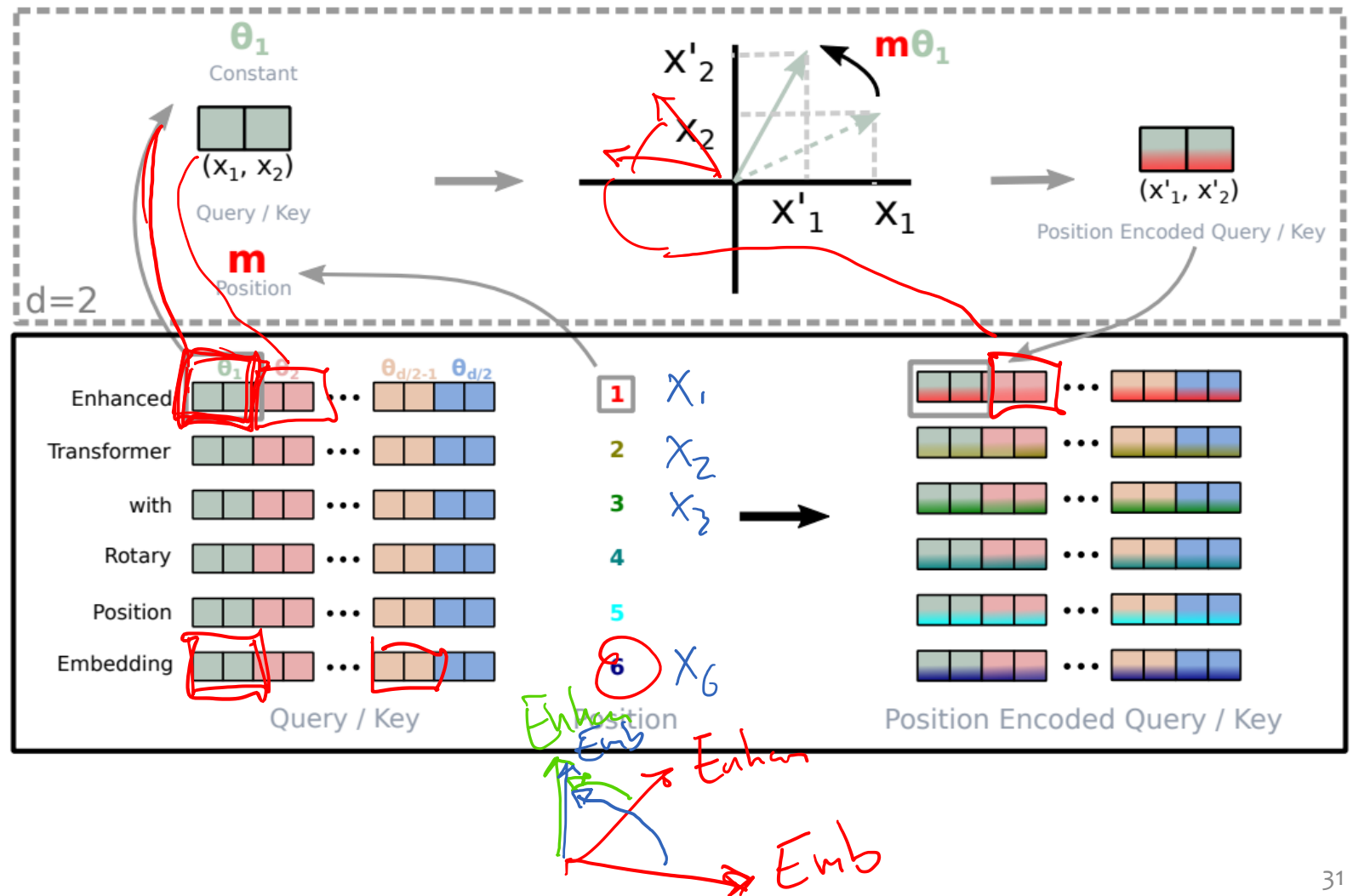
# Rotary Position Embeddings (RoPE)

**Q:** Why does this slide have so many typos?

**A:** I'm really not sure. I very meticulously type up the latex for my slides myself and think carefully about all the things I put in them.

# Rotary Position Embeddings (RoPE)

- Rotary position embeddings are a kind of **relative** position embeddings
- Key idea:
  - break each  $d$ -dimensional input vector into  $d/2$  vectors of length 2
  - rotate each of the  $d/2$  vectors by an amount scaled by  $m$
  - $m$  is the absolute position of the query or the key



# Rotary Position Embeddings (RoPE)

Standard attention:

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j, \forall j$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$

$$s_{t,j} = \mathbf{k}_j^T \mathbf{q}_t / \sqrt{|\mathbf{k}|}, \forall j, t$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

RoPE attention:

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j, \forall j$$

$$\tilde{\mathbf{q}}_j = \mathbf{R}_{\Theta,j} \mathbf{q}_j$$

$$s_{t,j} = \tilde{\mathbf{k}}_j^T \tilde{\mathbf{q}}_t / \sqrt{d_k}, \forall j, t$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$

$$\tilde{\mathbf{k}}_j = \mathbf{R}_{\Theta,j} \mathbf{k}_j$$

where  $\mathbf{W}_k, \mathbf{W}_q \in \mathbb{R}^{d_{model} \times d_k}$ . Herein we use  $d = d_k$  for brevity.

For some fixed absolute position  $m$ , the rotary matrix  $\mathbf{R}_{\Theta,m} \in \mathbb{R}^{d_k \times d_k}$  is given by:

$$R_{\Theta,m} = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \dots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos m\theta_{d_k/2} & -\sin m\theta_{d_k/2} \\ 0 & 0 & 0 & 0 & \dots & \sin m\theta_{d_k/2} & \cos m\theta_{d_k/2} \end{pmatrix}$$

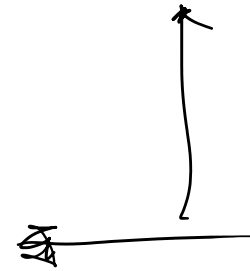
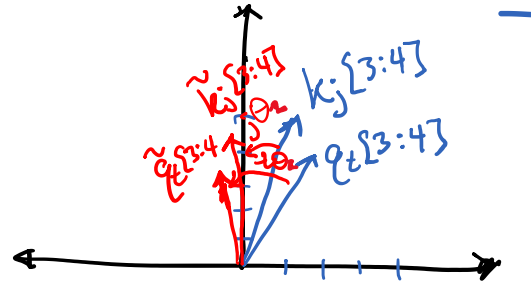
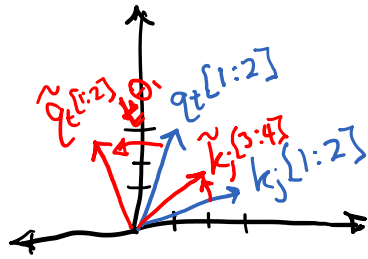
The  $\theta_i$  parameters are fixed ahead of time and defined as below.

$$\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]\}$$

# Rotary Position Embeddings (RoPE)

$$q_t = W_q^T x_t = [1, 3, 2, 4, 4, 7] \quad t=4 \quad d=6$$

$$k_j = W_k^T x_j = [3, 1, 2, 5, -2, 1] \quad j=2 \quad d/2=3$$



$$\tilde{q}_t = f(q_t, t) = R_{\theta, t} q_t$$

$$\tilde{k}_j = f(k_j, j) = R_{\theta, j} k_j$$



# Rotary Position Embeddings (RoPE)

Standard attention:

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j, \forall j$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$

$$s_{t,j} = \mathbf{k}_j^T \mathbf{q}_t / \sqrt{|\mathbf{k}|}, \forall j, t$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

RoPE attention:

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j, \forall j$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$

$$\tilde{\mathbf{q}}_j = \mathbf{R}_{\Theta,j} \mathbf{q}_j$$

$$\tilde{\mathbf{k}}_j = \mathbf{R}_{\Theta,j} \mathbf{k}_j$$

$$s_{t,j} = \tilde{\mathbf{k}}_j^T \tilde{\mathbf{q}}_t / \sqrt{d_k}, \forall j, t$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

Because of the block sparse pattern in  $\mathbf{R}_{\theta,m}$ , we can efficiently compute the matrix-vector product of  $\mathbf{R}_{\theta,m}$  with some arbitrary vector  $\mathbf{y}$  in a more efficient manner:

$$s_{t,j} = f(k_j, q_{t,j}, t)$$

$$\mathbf{R}_{\Theta,m} \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ \vdots \\ y_{d-1} \\ y_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -y_2 \\ y_1 \\ -y_4 \\ y_3 \\ \vdots \\ -y_d \\ y_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$

# Matrix Version of RoPE

## RoPE attention:

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j, \forall j$$

$$\tilde{\mathbf{q}}_j = \mathbf{R}_{\Theta, j} \mathbf{q}_j$$

$$s_{t, j} = \tilde{\mathbf{k}}_j^T \tilde{\mathbf{q}}_t / \sqrt{d_k}, \forall j, t$$

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t), \forall t$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$

$$\tilde{\mathbf{k}}_j = \mathbf{R}_{\Theta, j} \mathbf{k}_j$$

## Matrix Version:

$$\mathbf{Q} = \mathbf{XW}_q$$

$$\tilde{\mathbf{Q}} = g(\mathbf{Q}; \Theta)$$

$$\mathbf{S} = \tilde{\mathbf{Q}}\tilde{\mathbf{K}}^T / \sqrt{d_k}$$

$$\mathbf{A} = \text{softmax}(\mathbf{S})$$

$$\mathbf{K} = \mathbf{XW}_k$$

$$\tilde{\mathbf{K}} = g(\mathbf{K}; \Theta)$$

**Goal:** to construct a new matrix  $\tilde{\mathbf{Y}} = g(\mathbf{Y}; \Theta)$  such that  $\tilde{\mathbf{Y}}_{m, \cdot} = \mathbf{R}_{\Theta, m} \mathbf{y}_m$

$$\mathbf{C} = \left[ \begin{array}{ccc|ccc} 1\theta_1 & \cdots & 1\theta_{\frac{d}{2}} & 1\theta_1 & \cdots & 1\theta_{\frac{d}{2}} \\ \vdots & & \vdots & \vdots & & \vdots \\ N\theta_1 & \cdots & N\theta_{\frac{d}{2}} & N\theta_1 & \cdots & N\theta_{\frac{d}{2}} \end{array} \right]$$

$$\begin{aligned} \tilde{\mathbf{Y}} &= g(\mathbf{Y}; \Theta) \\ &= \left[ \mathbf{Y}_{\cdot, 1:d/2} \mid \mathbf{Y}_{\cdot, d/2+1:d} \right] \otimes \cos(\mathbf{C}) \\ &\quad + \left[ -\mathbf{Y}_{\cdot, d/2+1:d} \mid \mathbf{Y}_{\cdot, 1:d/2} \right] \otimes \sin(\mathbf{C}) \end{aligned}$$

~~Handwritten scribbles and notes in red ink, including a large box and some illegible text.~~

Handwritten notes in red ink:  $\mathbf{Y} = \left[ \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right]$  with a bracket indicating dimensions  $1 \dots d/2, d/2+1 \dots d$ .

# Matrix Version of RoPE

**Q:** Is this slide correct?

**A:** I'm really not sure.

But I did write it myself!

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j, \forall j$$
$$\tilde{\mathbf{k}}_j = \mathbf{R}_{\Theta, j} \mathbf{k}_j$$

**Matrix Version:**

$$\mathbf{Q} = \mathbf{XW}_q$$

$$\tilde{\mathbf{Q}} = g(\mathbf{Q}; \Theta)$$

$$\mathbf{S} = \tilde{\mathbf{Q}}\tilde{\mathbf{K}}^T / \sqrt{d_k}$$

$$\mathbf{A} = \text{softmax}(\mathbf{S})$$

$$\mathbf{K} = \mathbf{XW}_k$$

$$\tilde{\mathbf{K}} = g(\mathbf{K}; \Theta)$$

**Goal:** to construct a new matrix  $\tilde{\mathbf{Y}} = g(\mathbf{Y}; \Theta)$  such that  $\tilde{\mathbf{Y}}_{m, \cdot} = \mathbf{R}_{\Theta, m} \mathbf{y}_m$

$$\mathbf{C} = \left[ \begin{array}{ccc|ccc} 1\theta_1 & \cdots & 1\theta_{\frac{d}{2}} & 1\theta_1 & \cdots & 1\theta_{\frac{d}{2}} \\ \vdots & & \vdots & \vdots & & \vdots \\ N\theta_1 & \cdots & N\theta_{\frac{d}{2}} & N\theta_1 & \cdots & N\theta_{\frac{d}{2}} \end{array} \right]$$

$$\tilde{\mathbf{Y}} = g(\mathbf{Y}; \Theta)$$

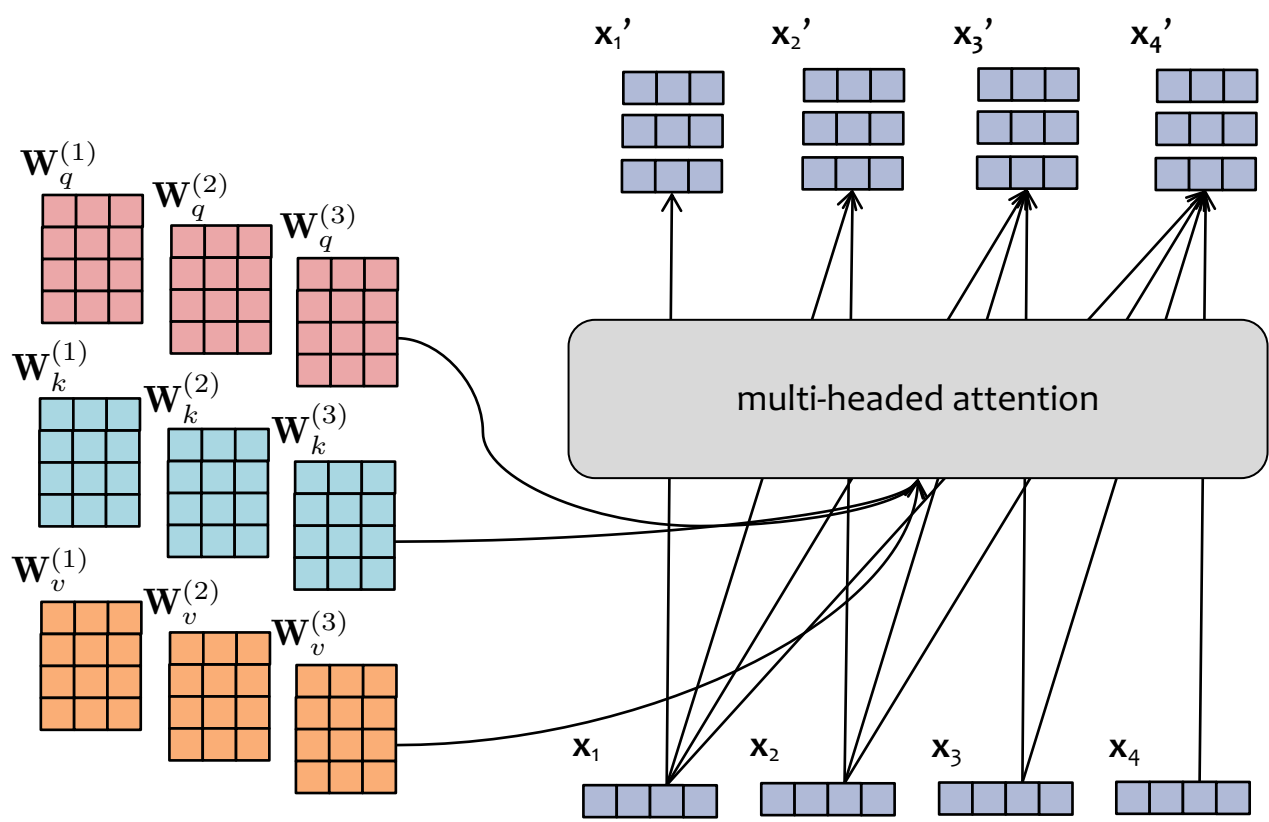
$$= \left[ \mathbf{Y}_{\cdot, 1:d/2} \mid \mathbf{Y}_{\cdot, d/2+1:d} \right] \otimes \cos(\mathbf{C})$$

$$+ \left[ -\mathbf{Y}_{\cdot, d/2+1:d} \mid \mathbf{Y}_{\cdot, 1:d/2} \right] \otimes \sin(\mathbf{C})$$

# **GROUPED QUERY ATTENTION (GQA)**

# Matrix Version of Multi-Headed (Causal) Attention

$$\mathbf{X}' = \text{concat}(\mathbf{X}'^{(1)}, \dots, \mathbf{X}'^{(h)})$$



$$\mathbf{X}'^{(i)} = \text{softmax} \left( \frac{\mathbf{Q}^{(i)} (\mathbf{K}^{(i)})^T}{\sqrt{d_k}} + \mathbf{M} \right) \mathbf{V}^{(i)}$$

$$\mathbf{Q}^{(i)} = \mathbf{X} \mathbf{W}_q^{(i)}$$

$$\mathbf{K}^{(i)} = \mathbf{X} \mathbf{W}_k^{(i)}$$

$$\mathbf{V}^{(i)} = \mathbf{X} \mathbf{W}_v^{(i)}$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_4]^T$$

# Grouped Query Attention (GQA)

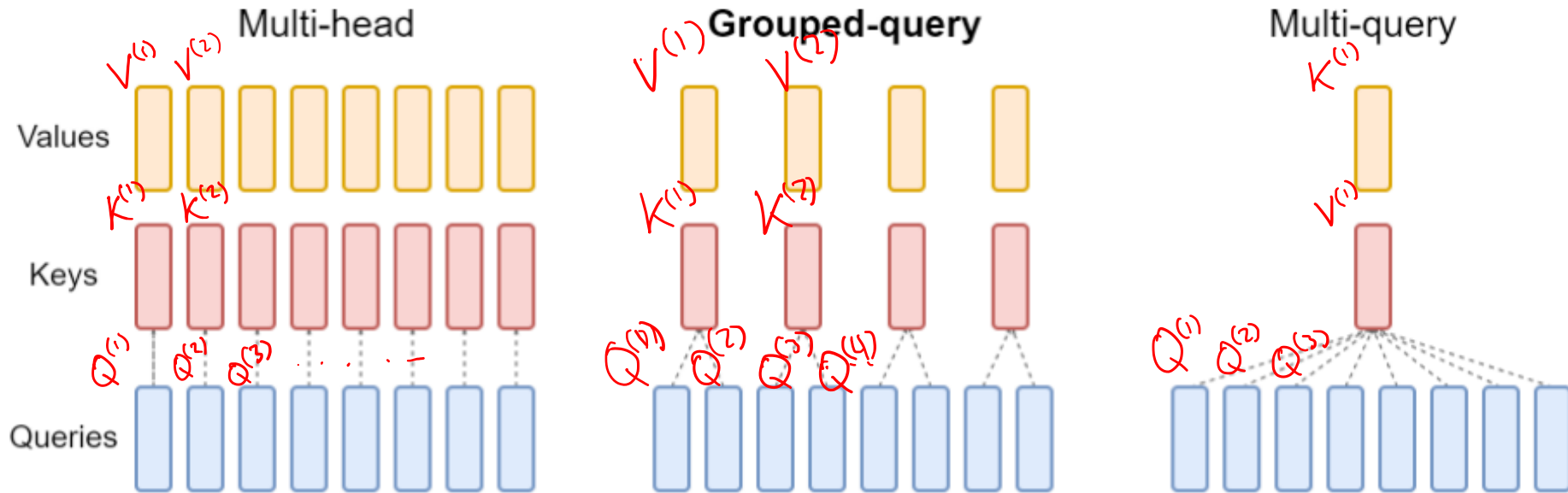
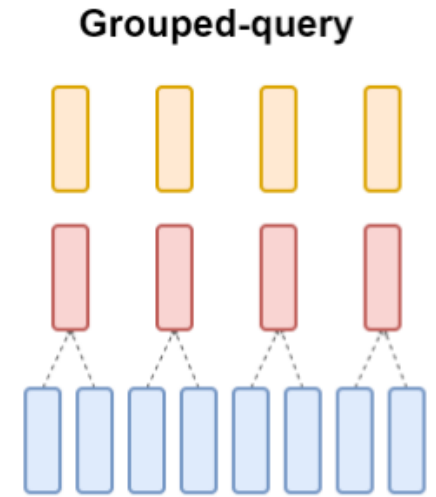


Figure 2: Overview of grouped-query method. Multi-head attention has  $H$  query, key, and value heads. Multi-query attention shares single key and value heads across all query heads. Grouped-query attention instead shares single key and value heads for each *group* of query heads, interpolating between multi-head and multi-query attention.

# Grouped Query Attention (GQA)

- **Key idea:** reuse the same key-value heads for multiple different query heads
- **Parameters:** The parameter matrices are all the same size, but we now have fewer key/value parameter matrices (heads) than query parameter matrices (heads)

- $h_q$  = the number of query heads
- $h_{kv}$  = the number of key/value heads
- Assume  $h_q$  is divisible by  $h_{kv}$
- $g = h_q/h_{kv}$  is the size of each group (i.e. the number of query vectors per key/value vector).



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T]^T$$

$$\mathbf{V}^{(i)} = \mathbf{XW}_v^{(i)}, \forall i \in \{1, \dots, \cancel{d_{kv}}\} h_{kv}$$

$$\mathbf{K}^{(i)} = \mathbf{XW}_k^{(i)}, \forall i \in \{1, \dots, \cancel{d_{kv}}\} h_{kv}$$

$$\mathbf{Q}^{(i,j)} = \mathbf{XW}_q^{(i,j)}, \forall i \in \{1, \dots, \cancel{d_{kv}}\} h_{kv}, \forall j \in \{1, \dots, g\}$$

# **SLIDING WINDOW ATTENTION**

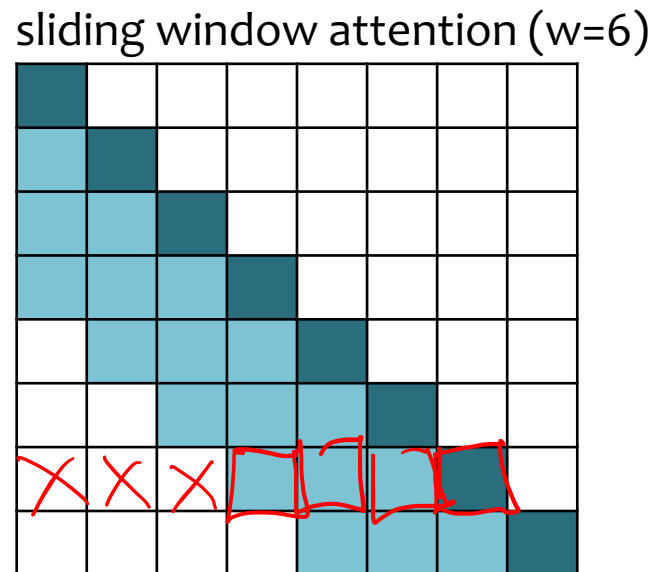
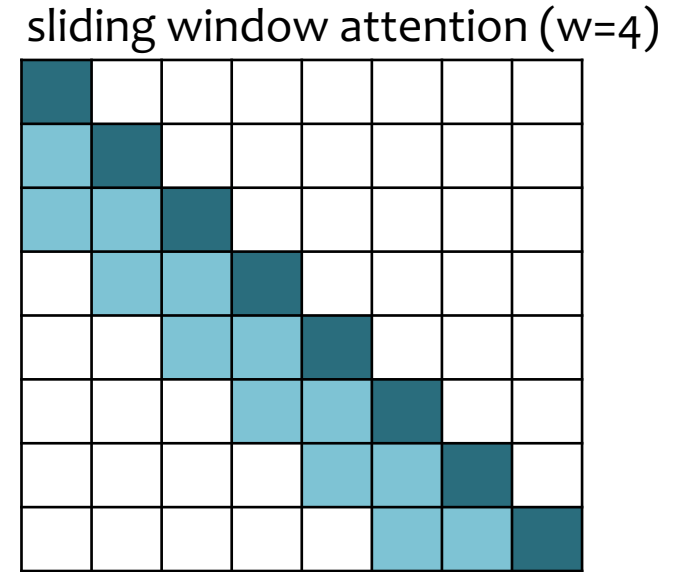
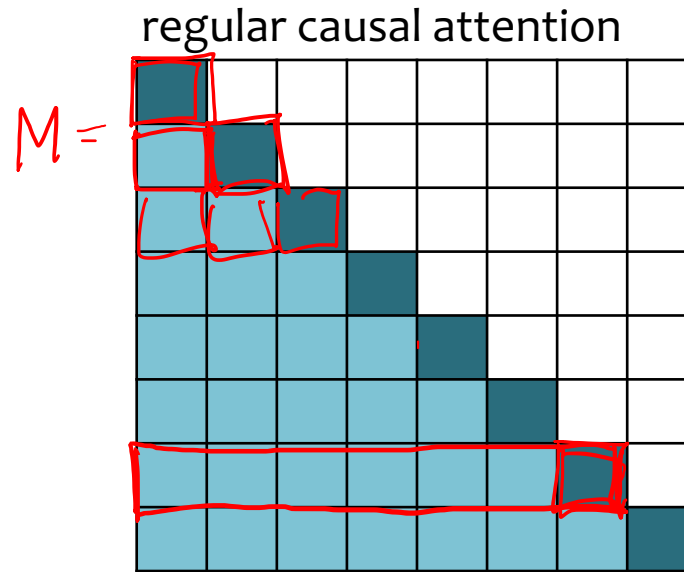


# Sliding Window Attention

## Sliding Window Attention

- also called “local attention” and introduced for the Longformer model (2020)
- **The problem:** regular attention is computationally expensive and requires a lot of memory
- **The solution:** apply a causal mask that only looks at the include a window of  $(\frac{1}{2}w+1)$  tokens, with the rightmost window element being the current token (i.e. on the diagonal)

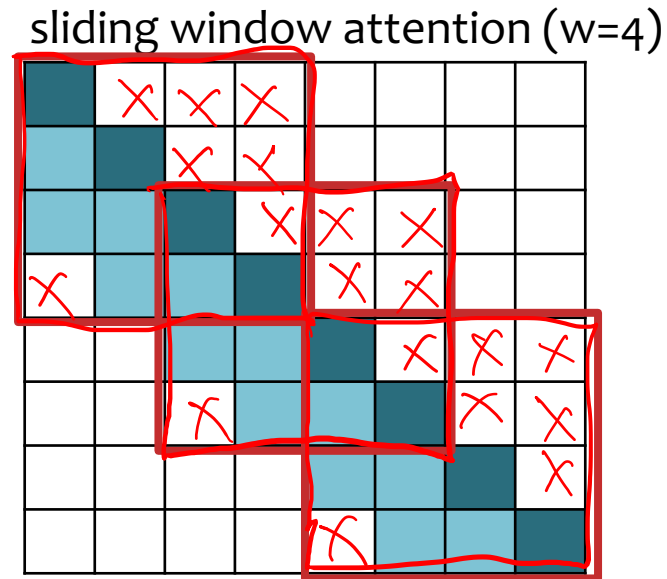
$$\mathbf{X}' = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} + \mathbf{M} \right) \mathbf{V}$$



# Sliding Window Attention

## Sliding Window Attention

- also called “local attention” and introduced for the Longformer model (2020)
- **The problem:** regular attention is computationally expensive and requires a lot of memory
- **The solution:** apply a causal mask that only looks at the include a window of  $(\frac{1}{2}w+1)$  tokens, with the rightmost window element being the current token (i.e. on the diagonal)



## 3 ways you could implement

1. *naïve implementation:* just do the matrix multiplication, but this is still slow
2. *for-loop implementation:* asymptotically faster / less memory, but unusable in practice b/c for-loops in PyTorch are too slow
3. *sliding chunks implementation:* break into Q and K into chunks of size  $w \times w$ , with overlap of  $\frac{1}{2}w$ ; then compute full attention within each chunk and mask out chunk (very fast/low memory in practice)

$$\mathbf{X}' = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} + \mathbf{M} \right) \mathbf{V}$$