

Model Compression

Mohammadreza Banaei

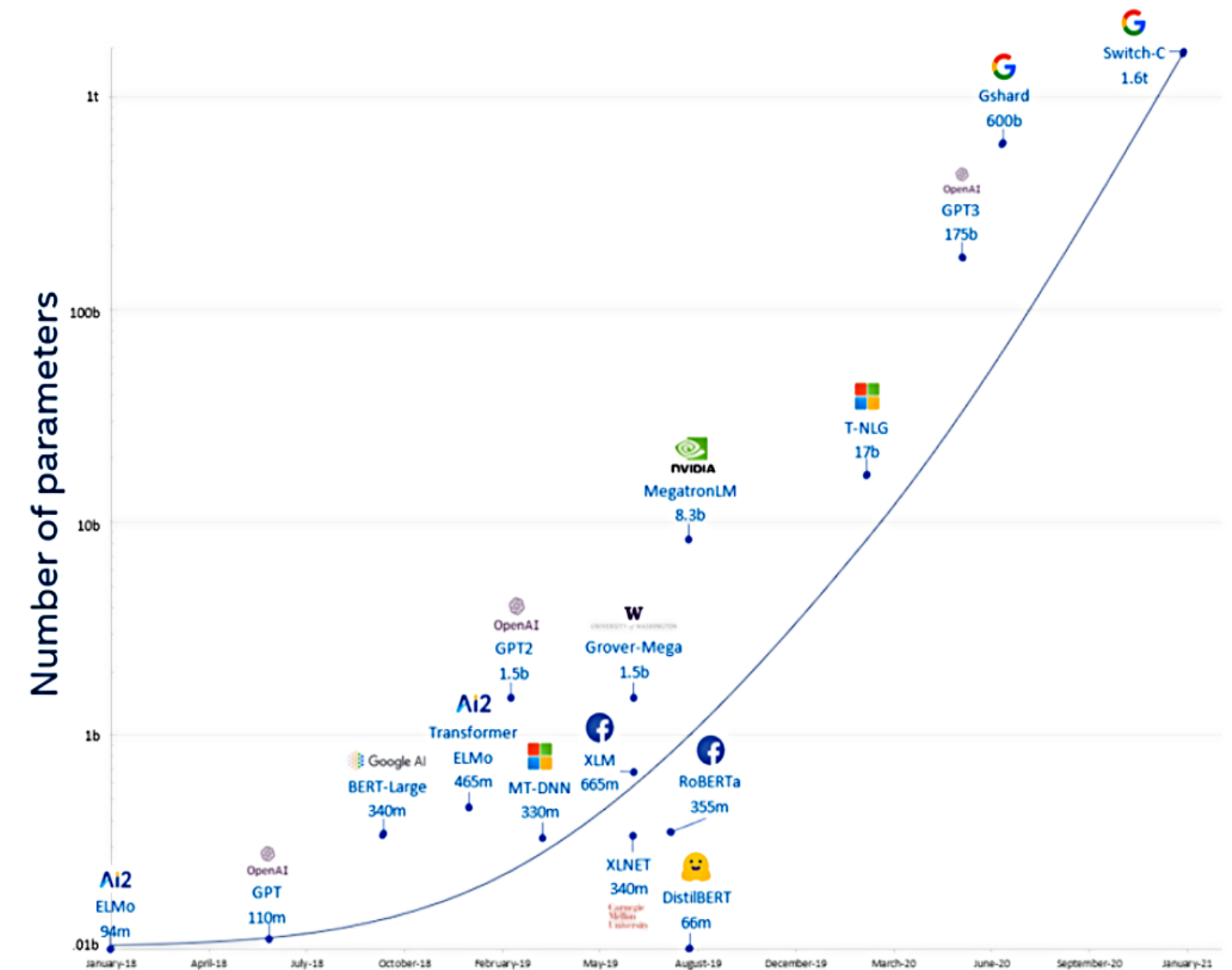
Outline

- Motivation
- Compression methods
 - Pruning
 - Quantization
 - Weight factorization
 - Knowledge distillation
 - Weight Sharing
- Sub-quadratic Transformers

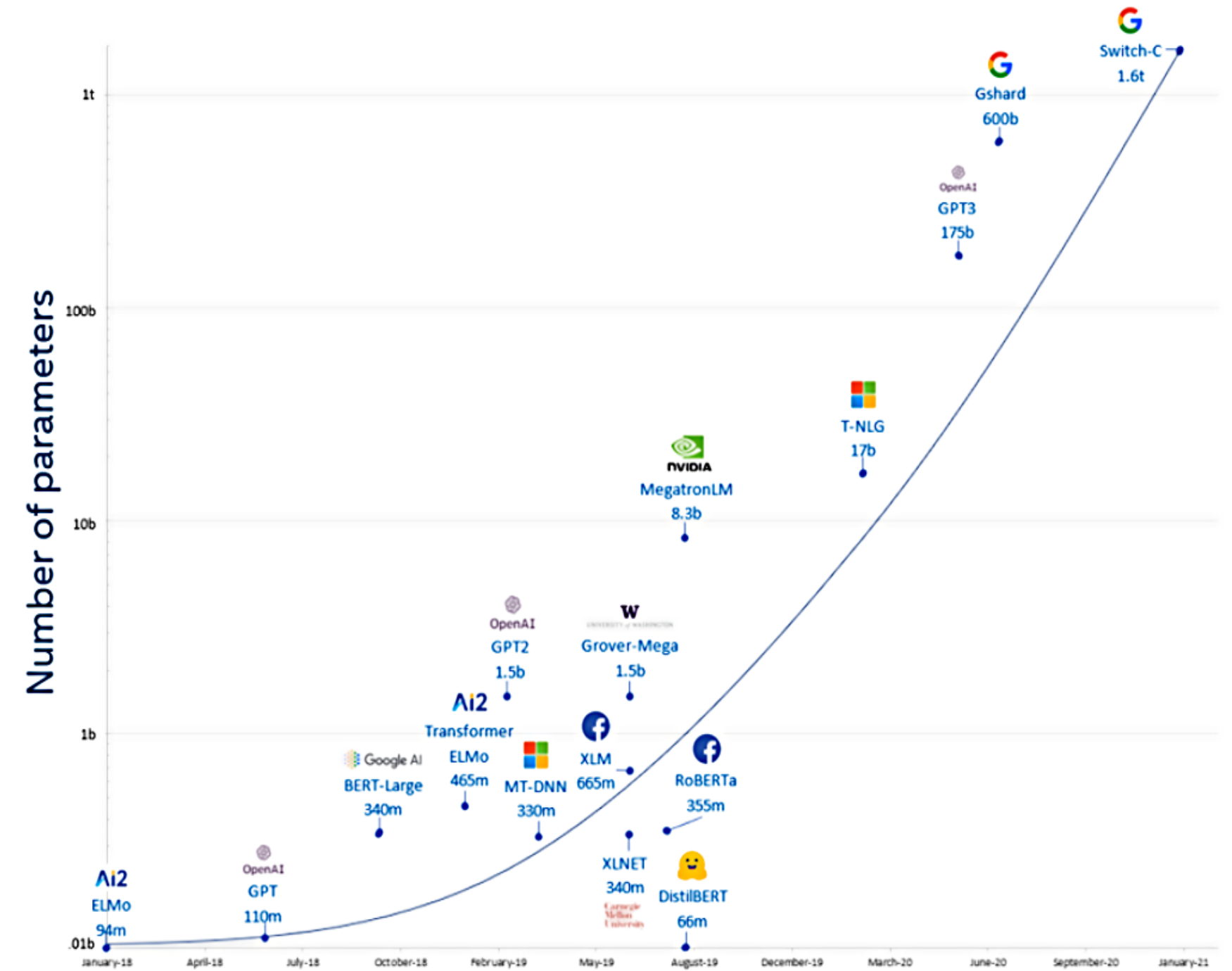
Why do we need compression?

Growth of model parameters

- Exponential growth in model parameters

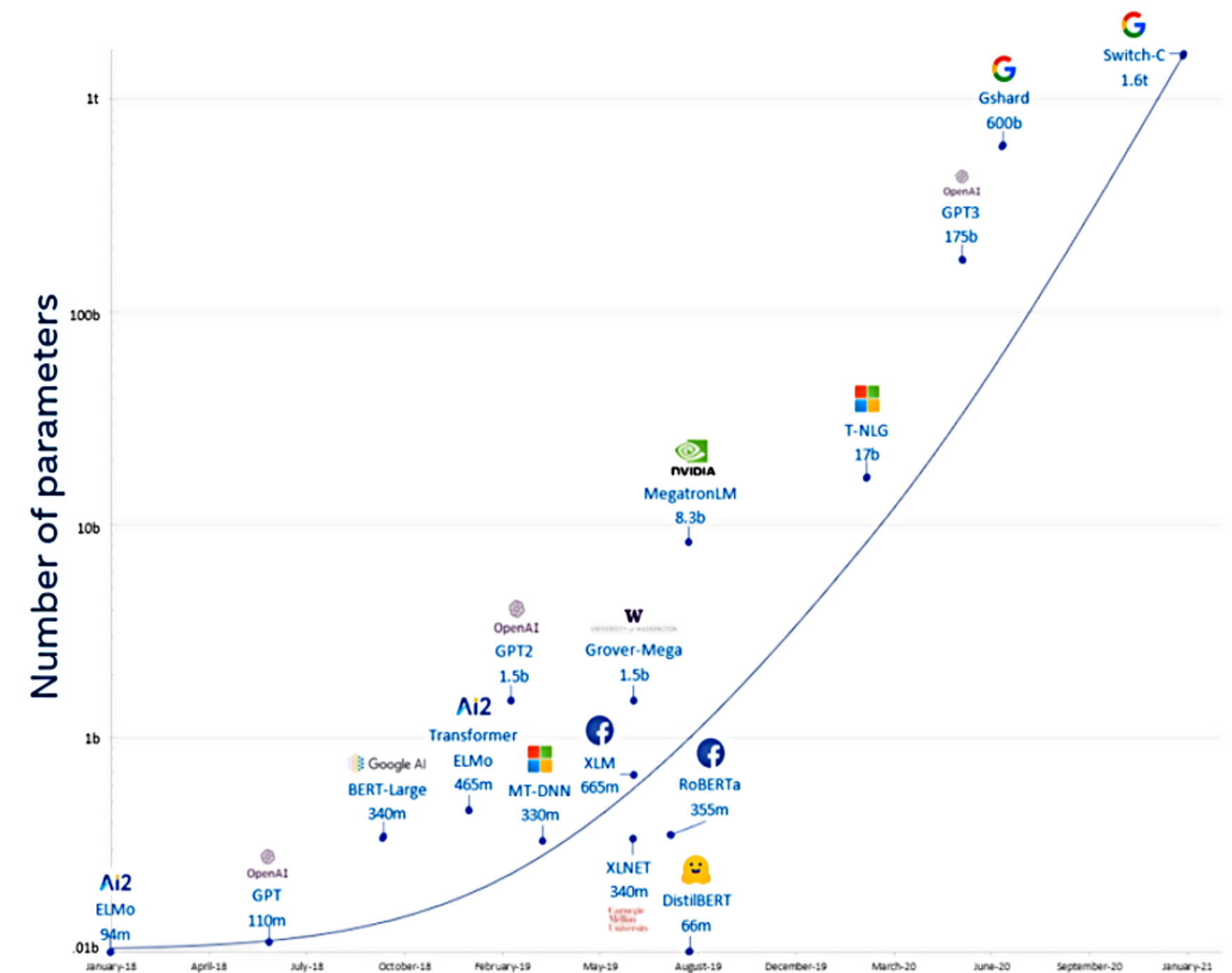
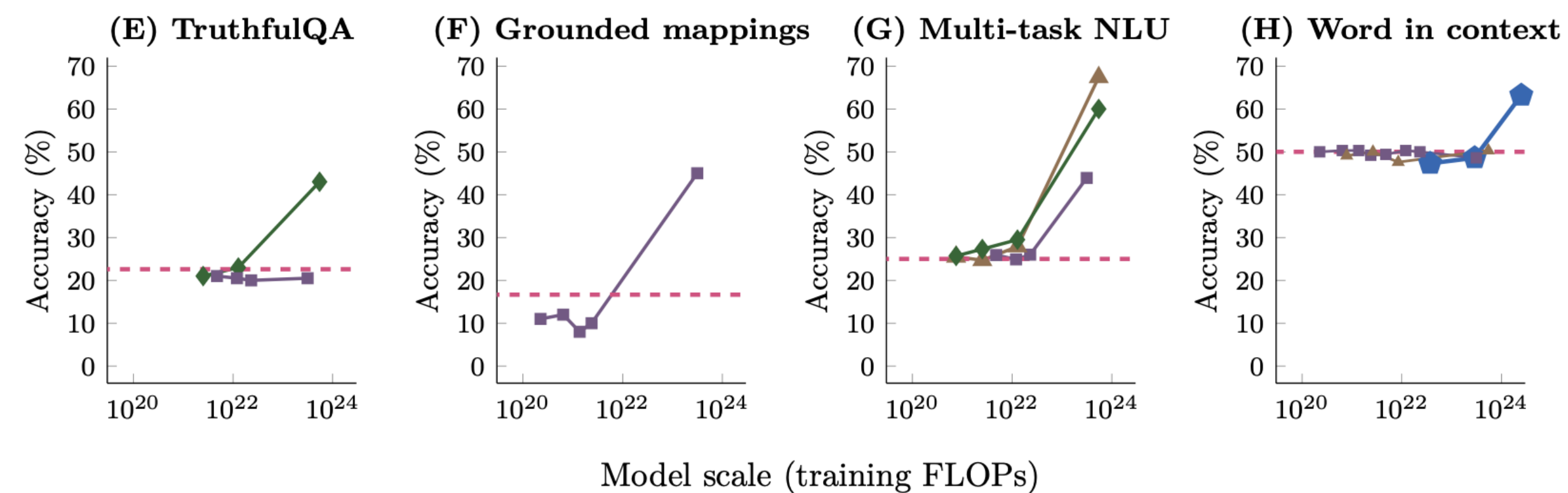


- Exponential growth in model parameters
 - Scaling laws



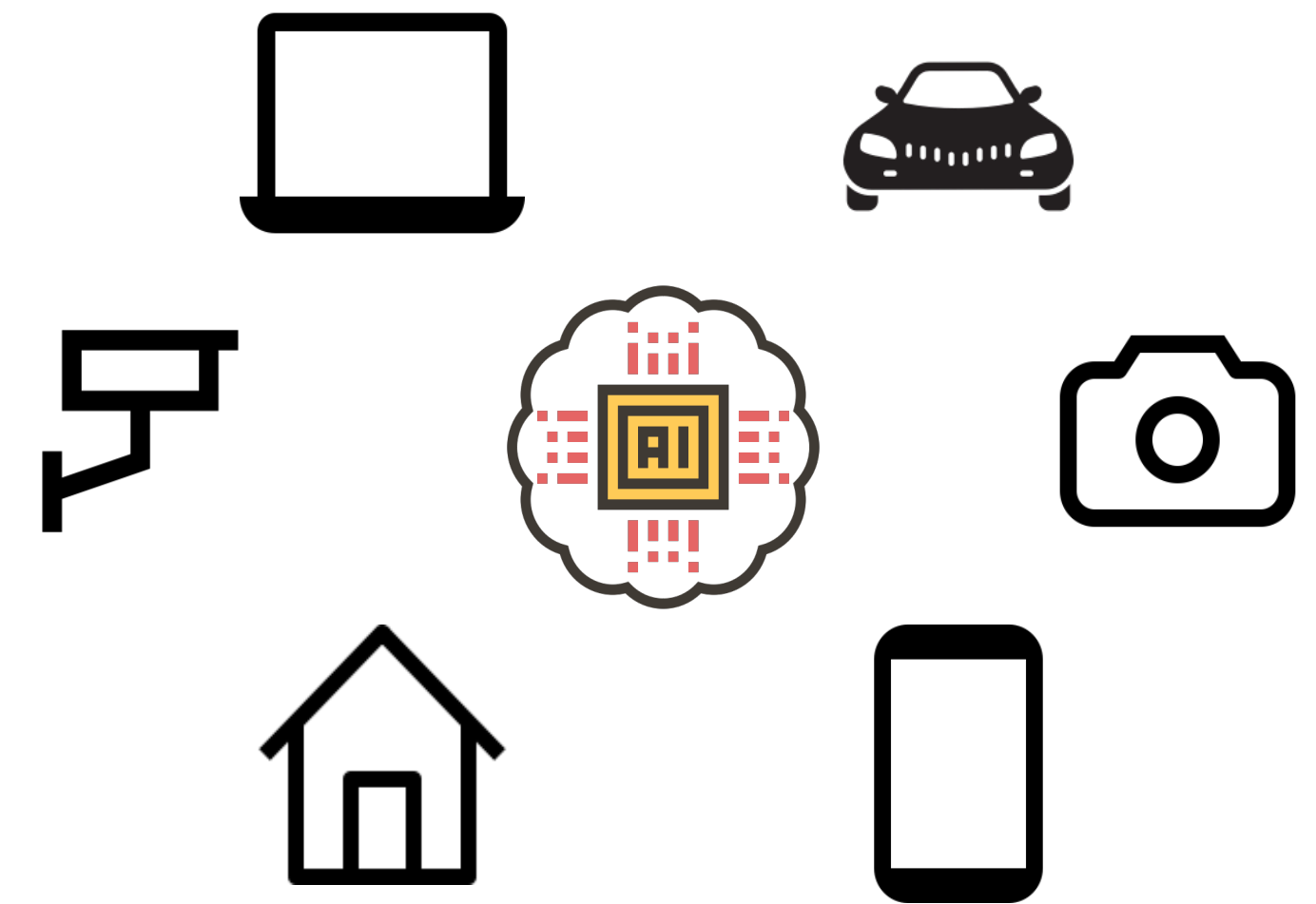
Growth of model parameters

- Exponential growth in model parameters
 - Scaling laws
 - Emergent abilities of LLMs



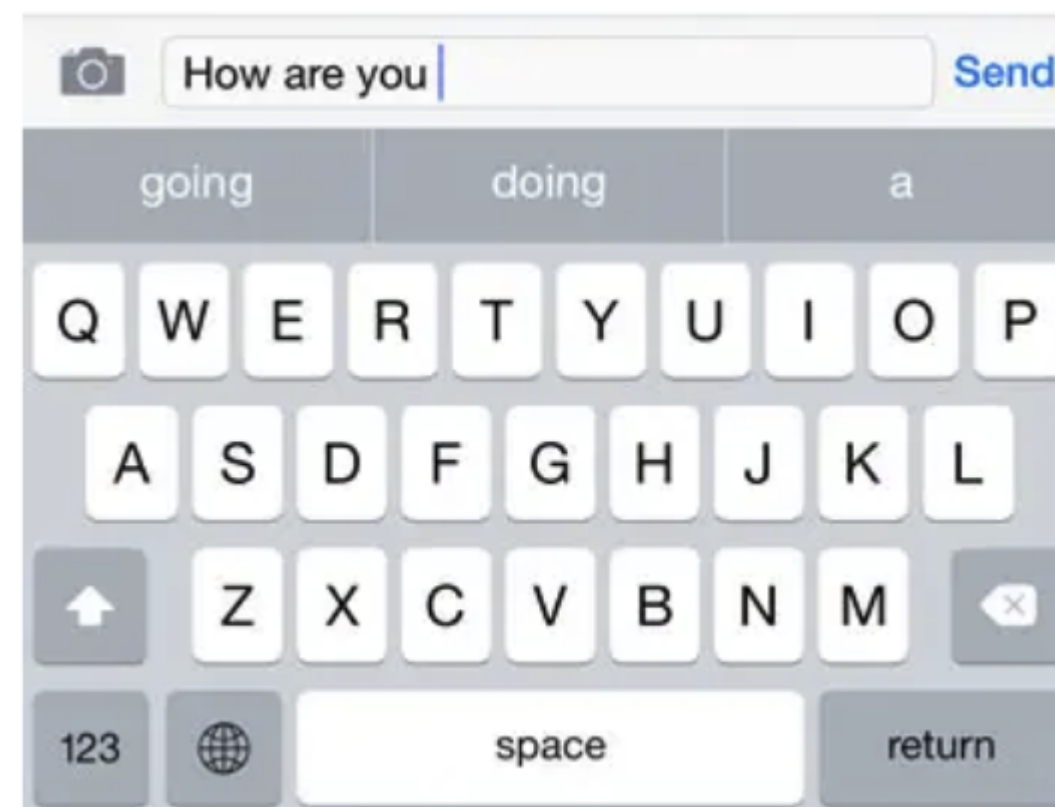
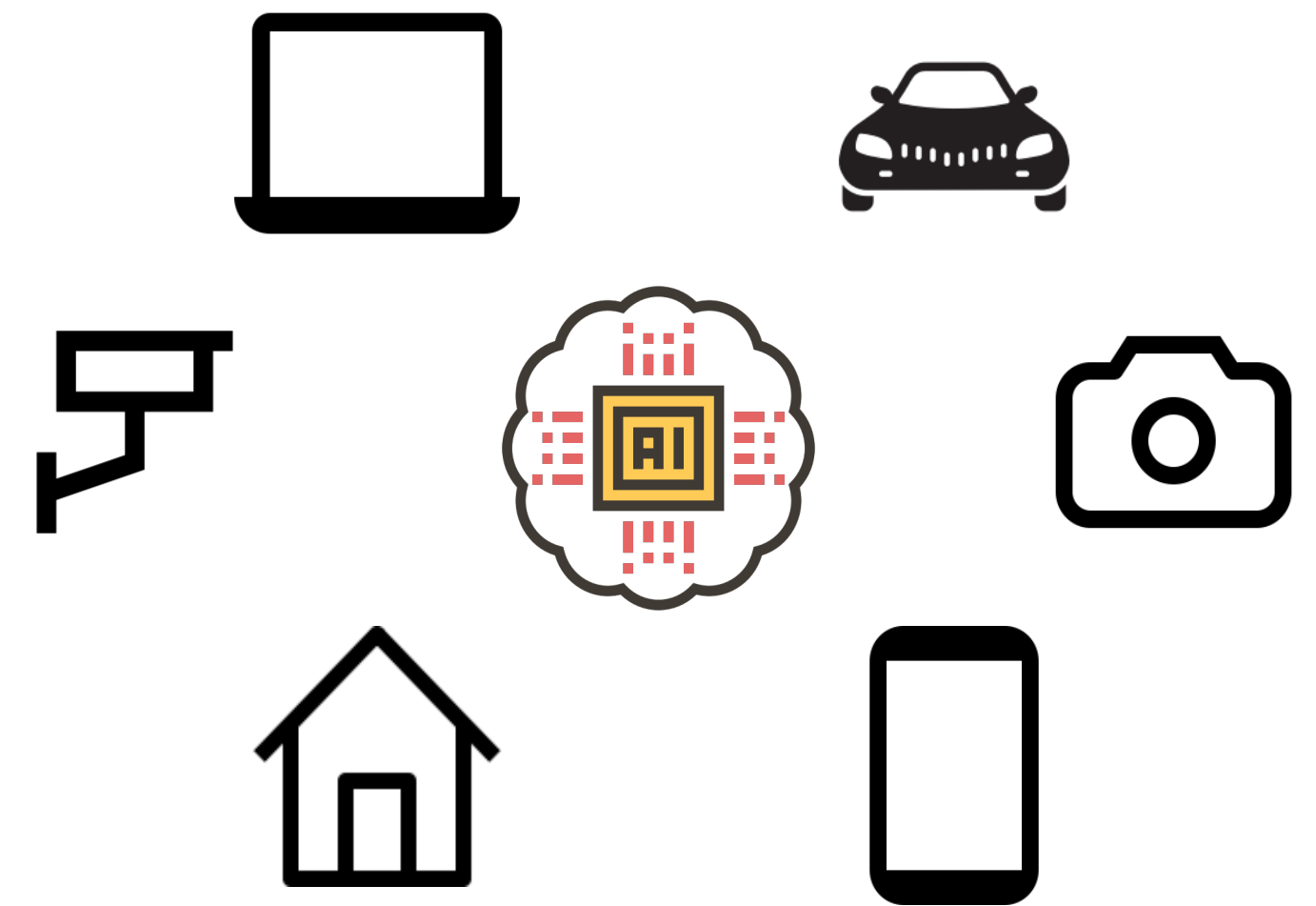
LLM Deployment in Production

- Cloud processing not always possible
 - Latency issue
 - Data privacy



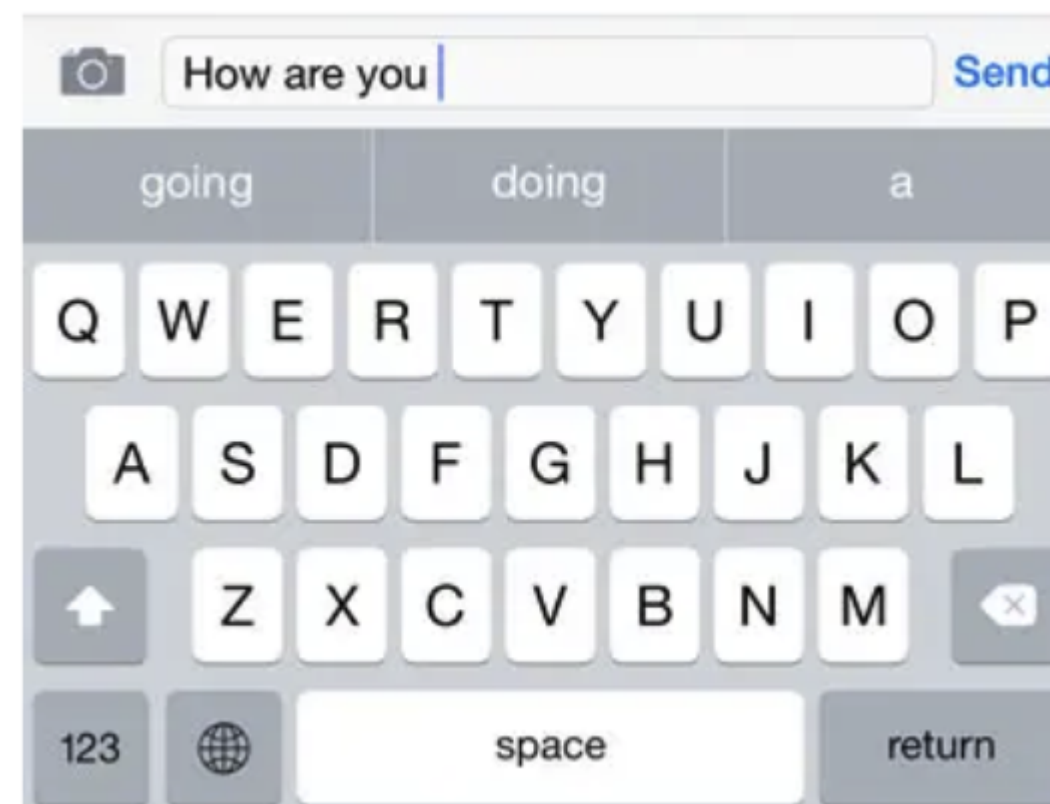
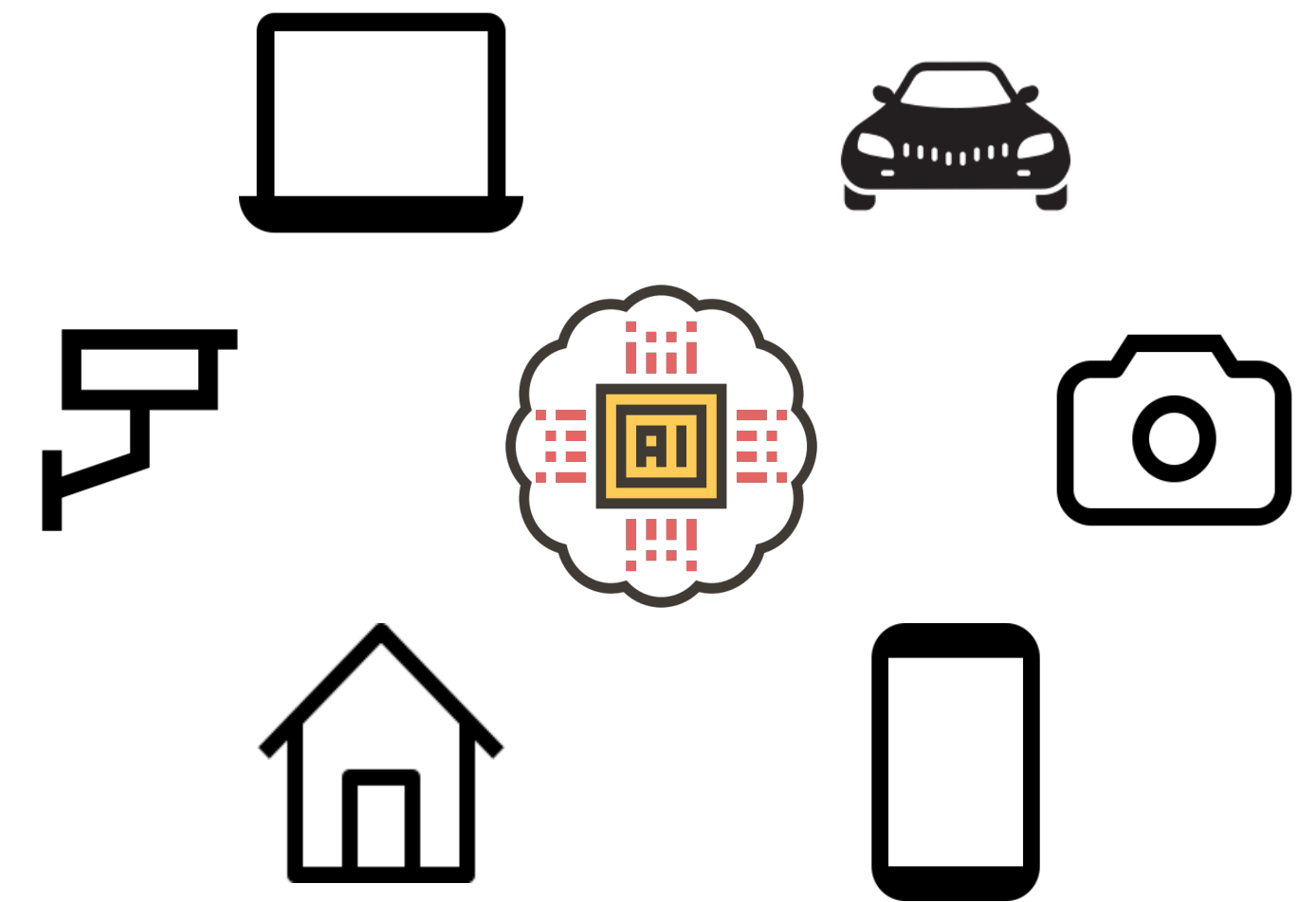
LLM Deployment in Production

- Cloud processing not always possible
 - Latency issue
 - Data privacy
- Inference time for edge devices



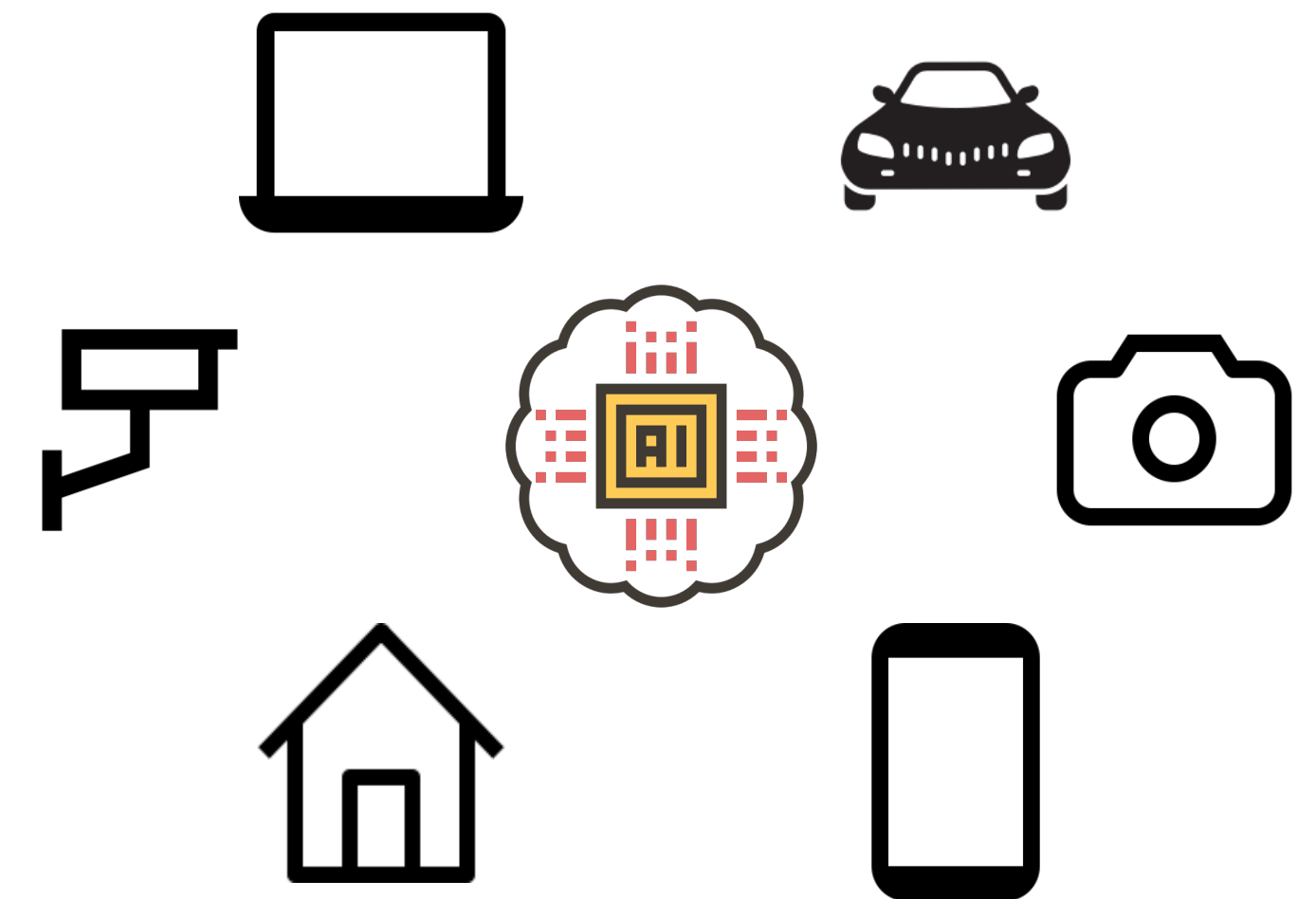
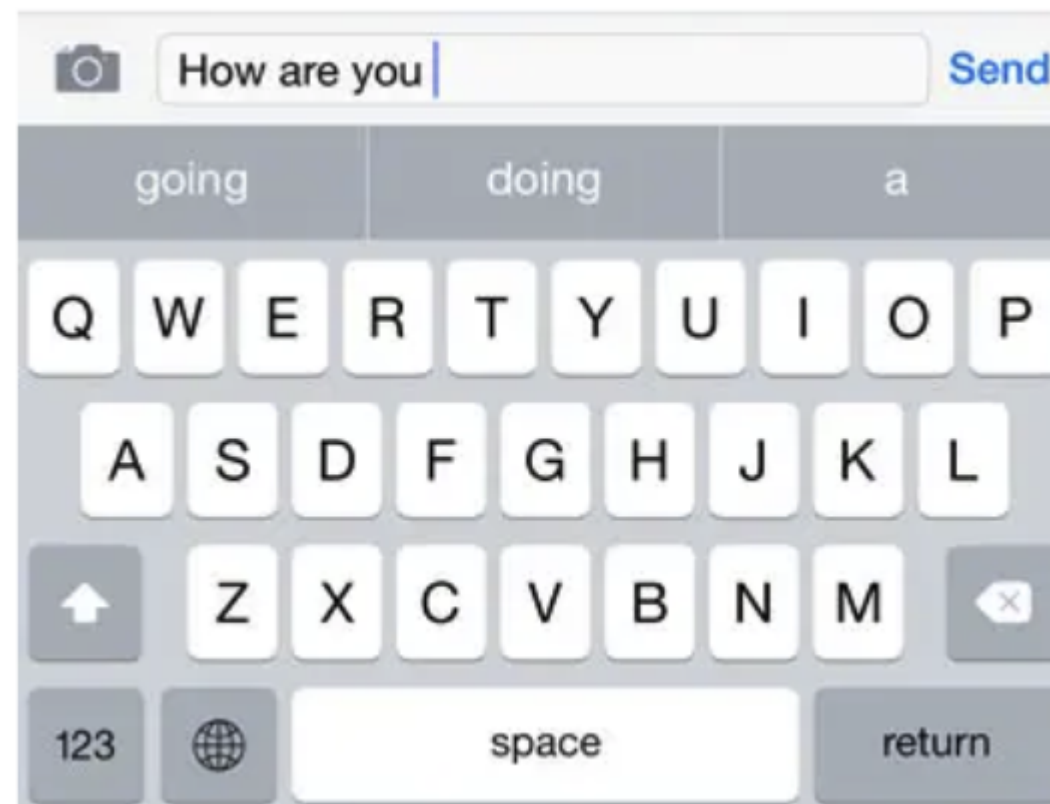
LLM Deployment in Production

- Cloud processing not always possible
 - Latency issue
 - Data privacy
- Inference time for edge devices
- Memory issue
 - ~350 GB just for **storing** a LLM weights!



LLM Deployment in Production

- Cloud processing not always possible
 - Latency issue
 - Data privacy
- Inference time for edge devices
- Memory issue
 - ~350 GB just for **storing** a LLM weights!
- Finetuning LLMs
 - Time-consuming
 - Expensive



LLM Deployment in Production

- Cloud processing not always possible

- Latency issue
- Data privacy

- Inference on edge devices

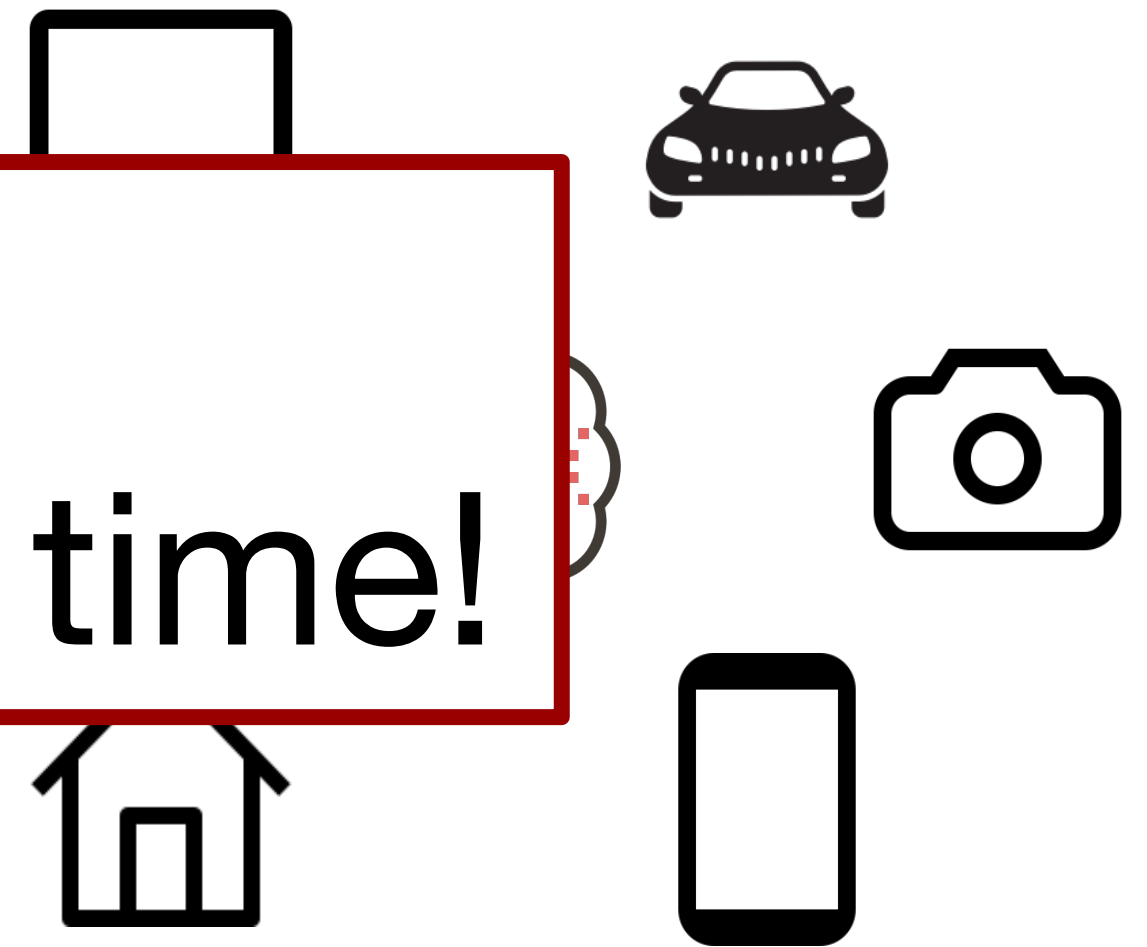
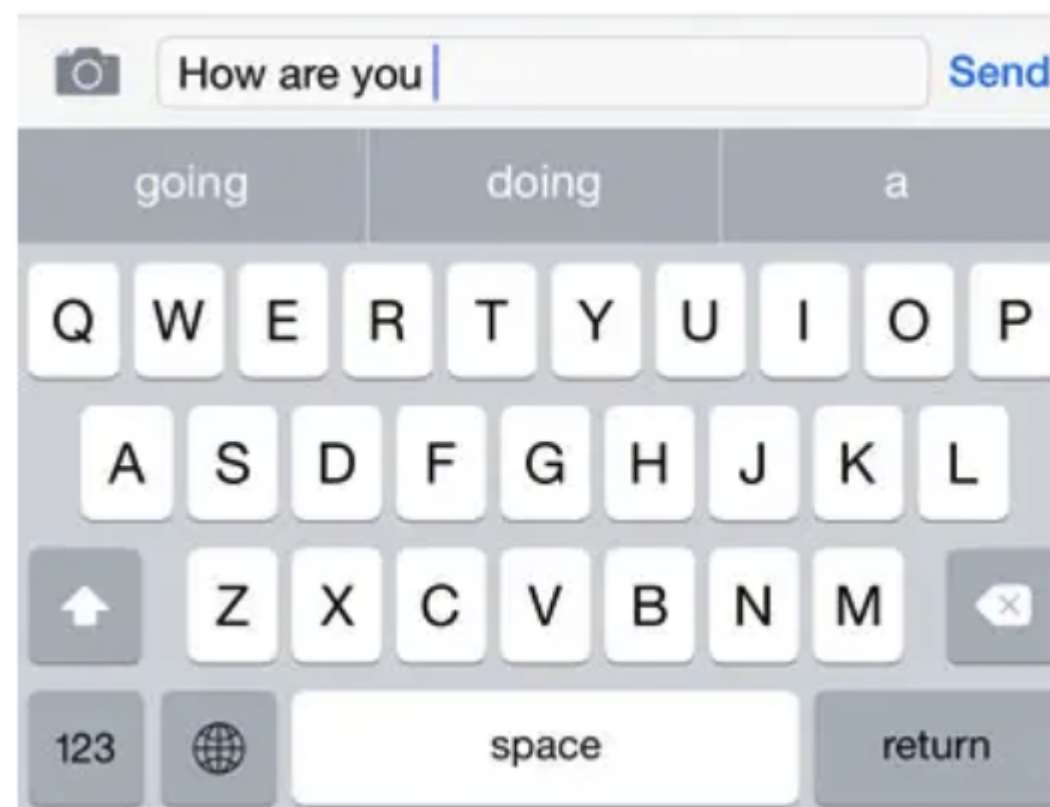
- Memory

- ~350 GB

- Finetuning

- Time-consuming
- Expensive

Compression could reduce
#parameters and inference time!



Train Large, then Compress!

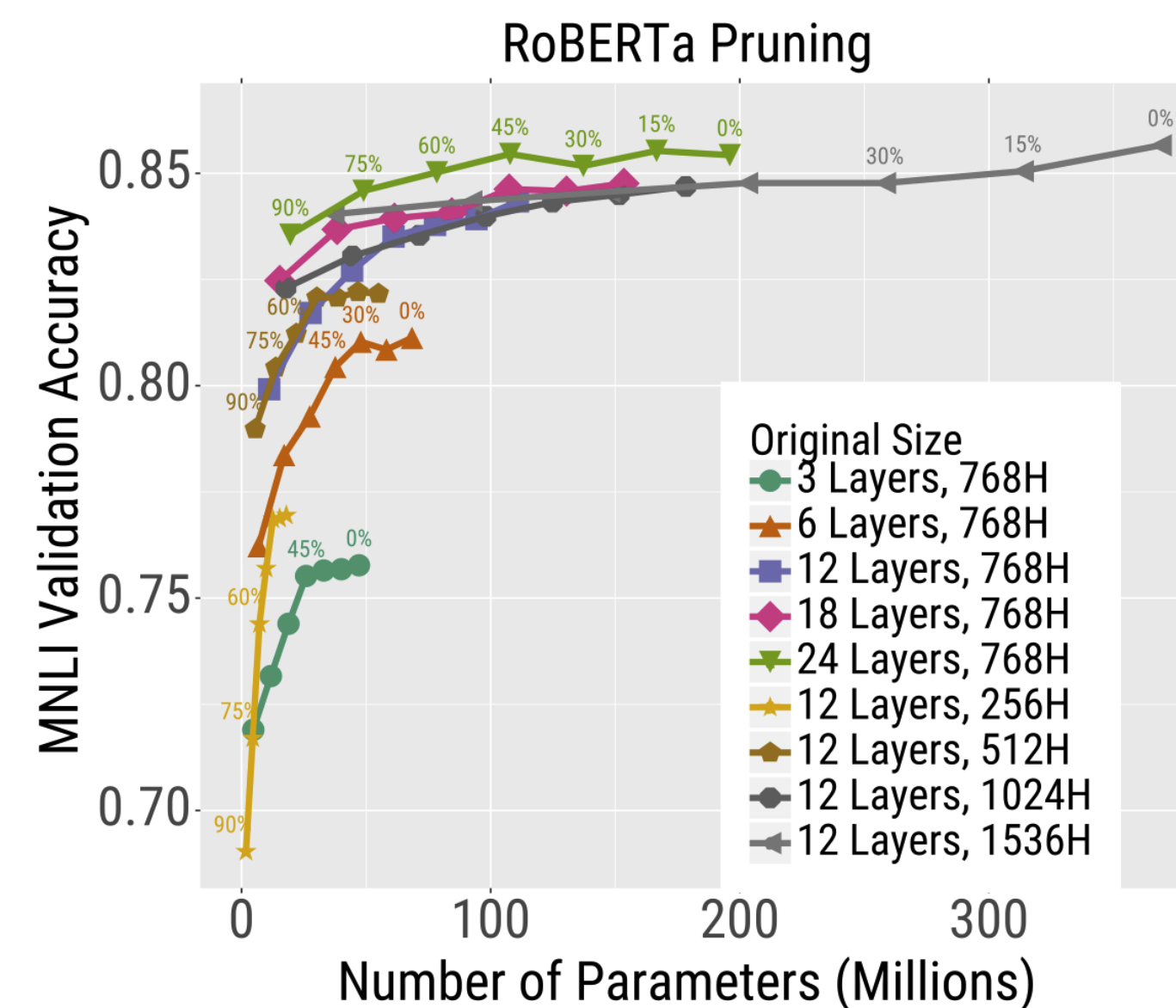
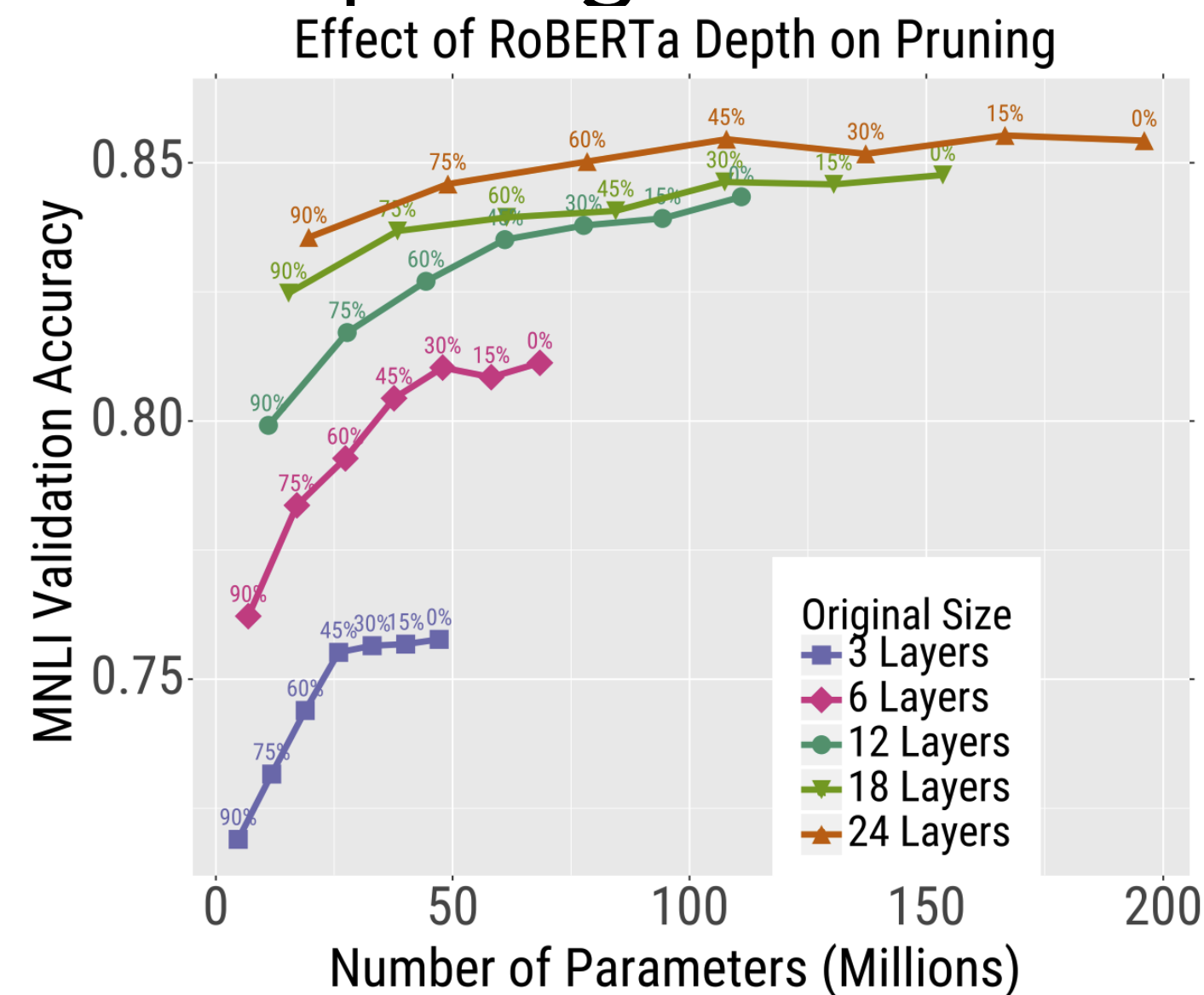
- Large models are more robust to compression techniques than small models

Train Large, then Compress!

- Large models are more robust to compression techniques than small models
- For given test-time constraints (e.g., inference time, #parameter)
 - heavily compressed, large models > lightly compressed, small models

Train Large, then Compress!

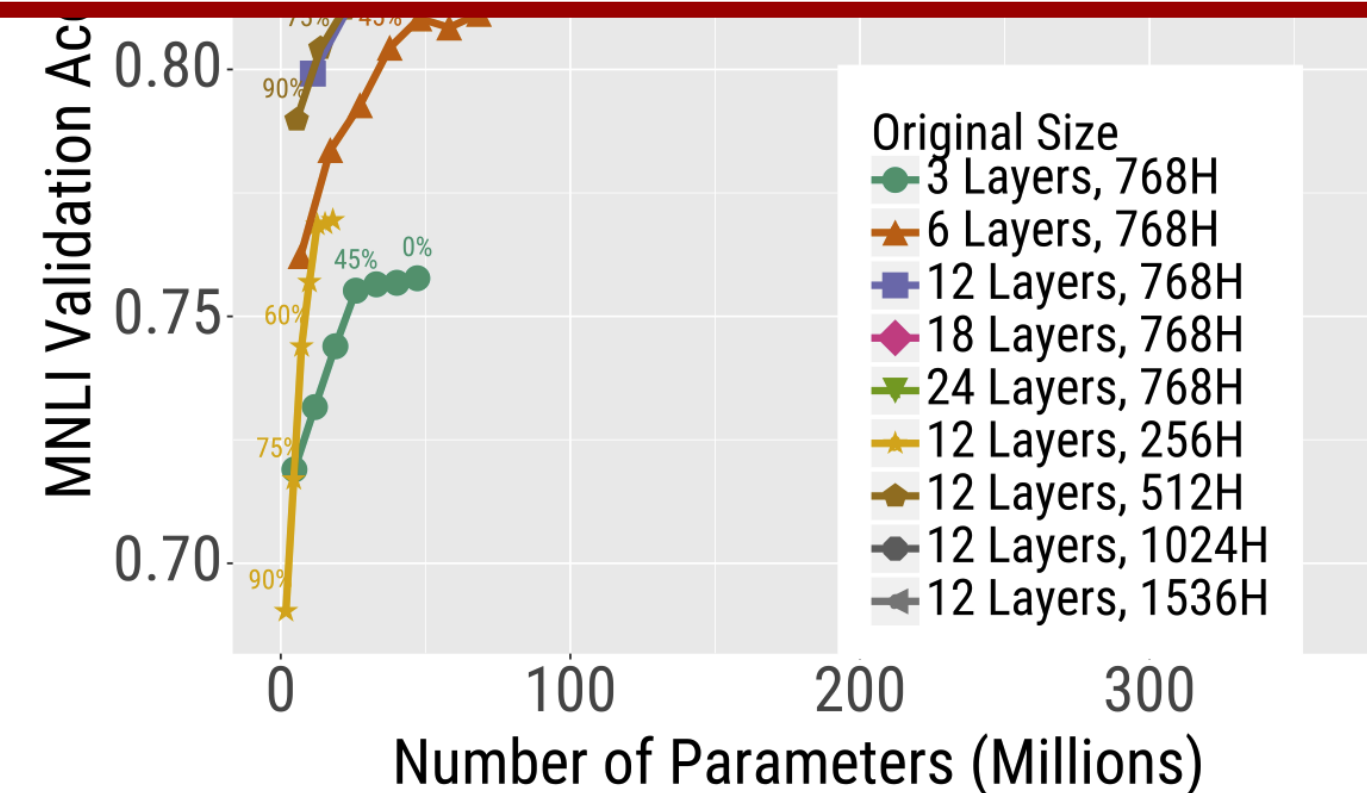
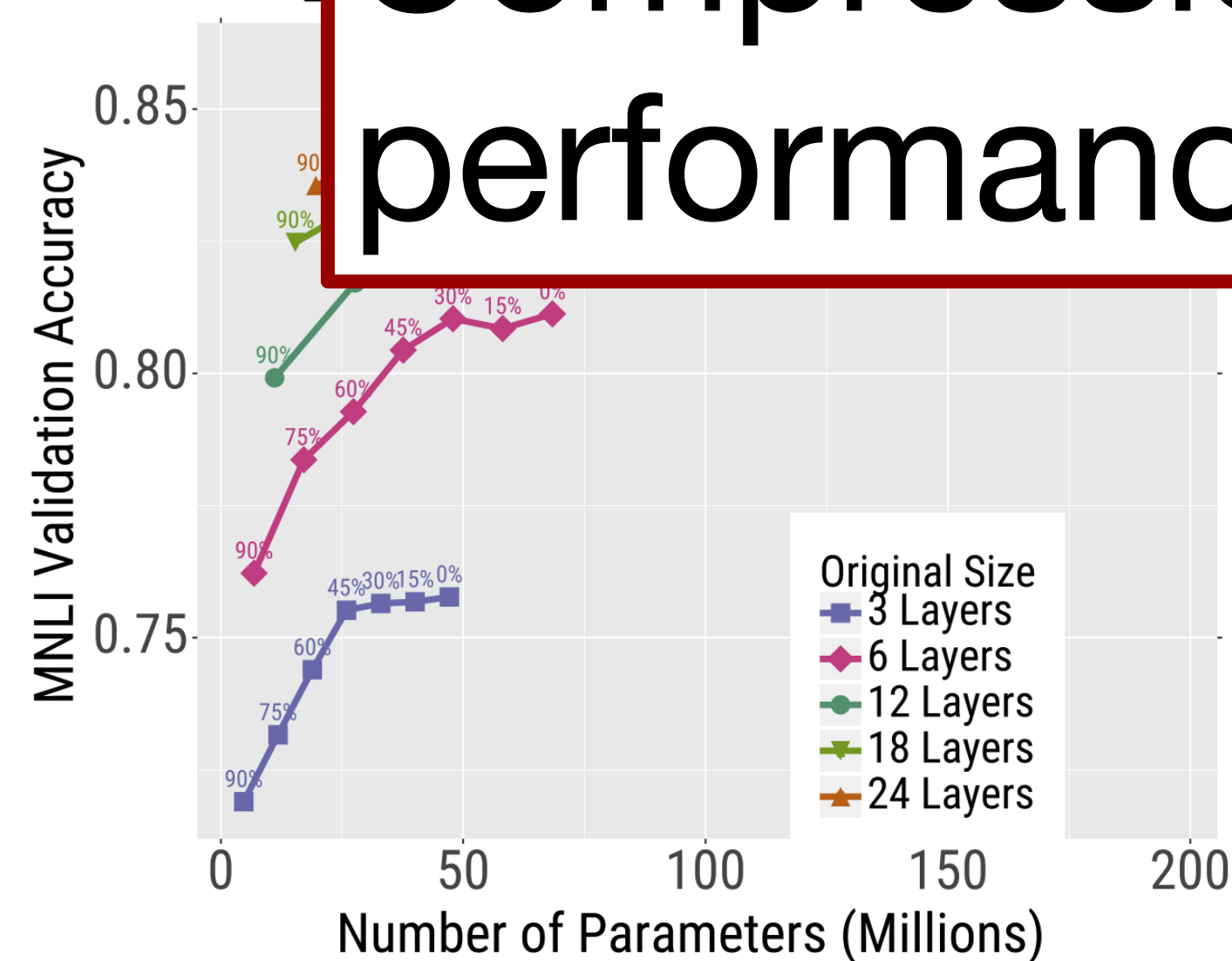
- Large models are more robust to compression techniques than small models
- For given test-time constraints (e.g., inference time, #parameter)
 - heavily compressed, large models > lightly compressed, small models
- Comparing downstream task performance for discussed scenarios



Train Large, then Compress!

- Large models are more robust to compression techniques than small models
- For given test-time constraints (e.g., inference time, #parameter)
 - heavily compressed, large models > lightly compressed, small models

• Compression improves the model's performance given a test-time budget!



How is the compression done?

Compression Methods

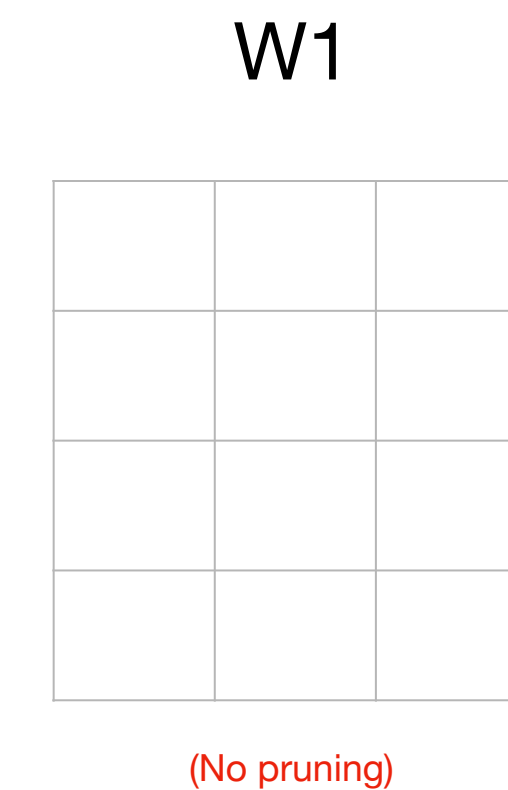
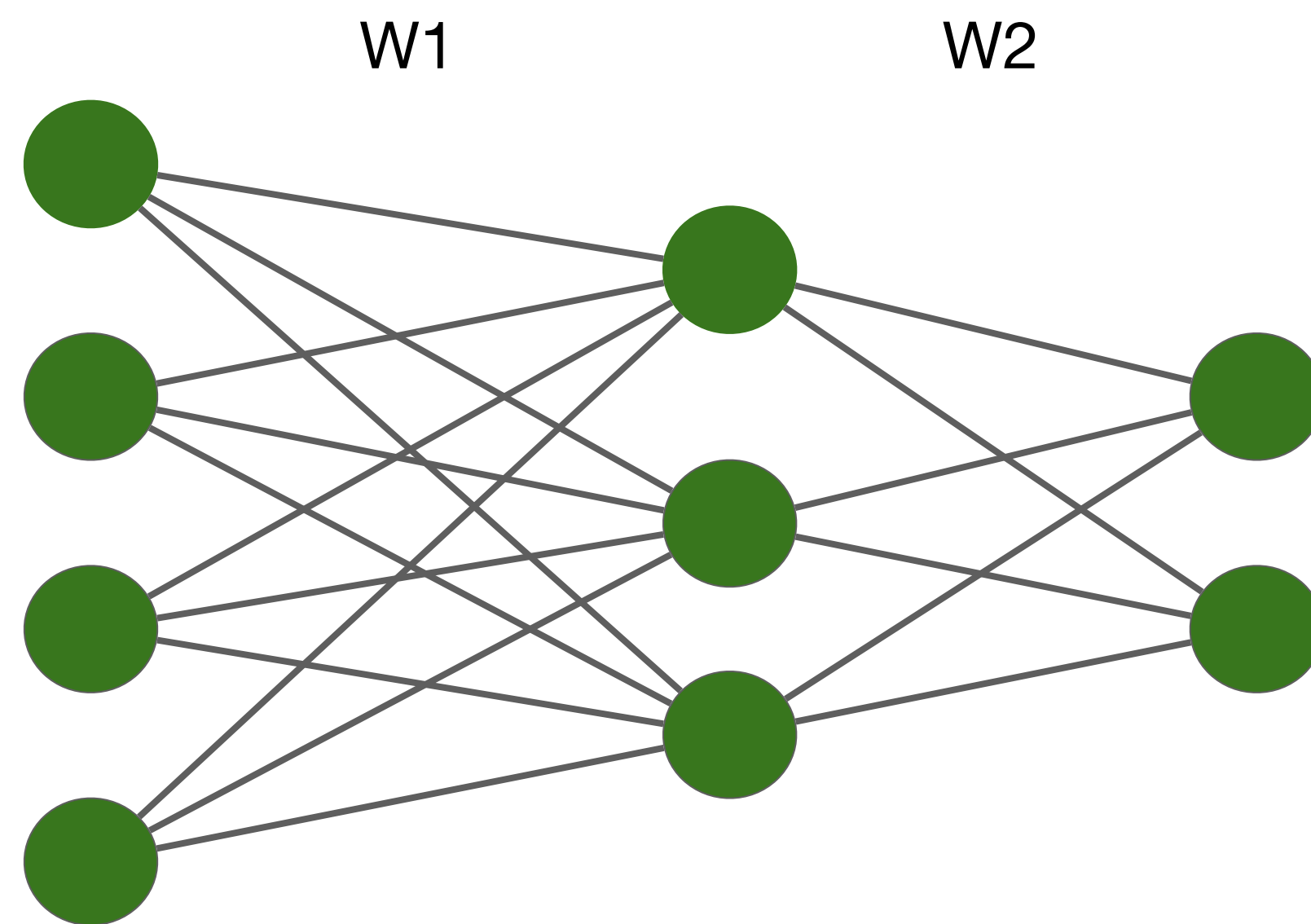
- Pruning
- Quantization
- Weight factorization
- Knowledge Distillation
- Weight sharing

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning		
Quantization		
Weight Factorization		
Weight Sharing		
Knowledge distillation		
Sub-quadratic Transformer		

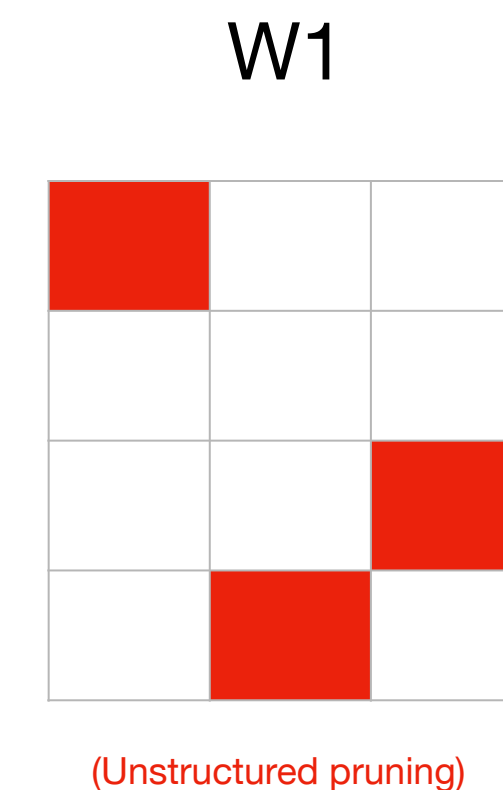
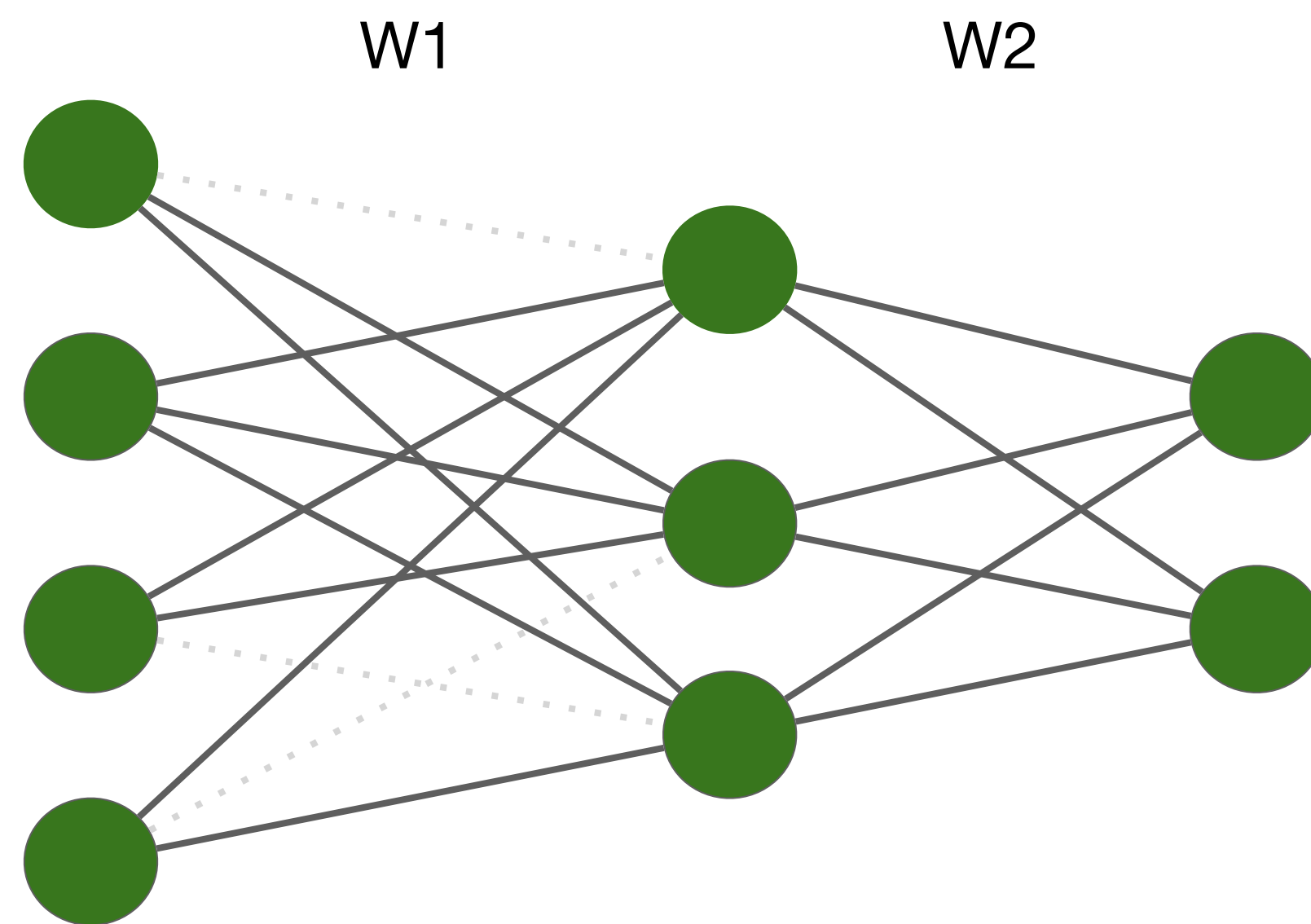
Pruning

- Sparse connectivity inspired by biological neural networks
- Unstructured pruning Vs. structured pruning



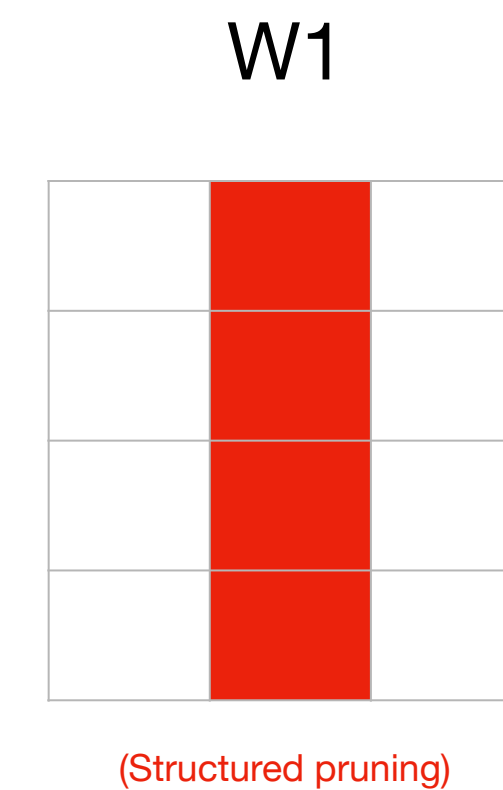
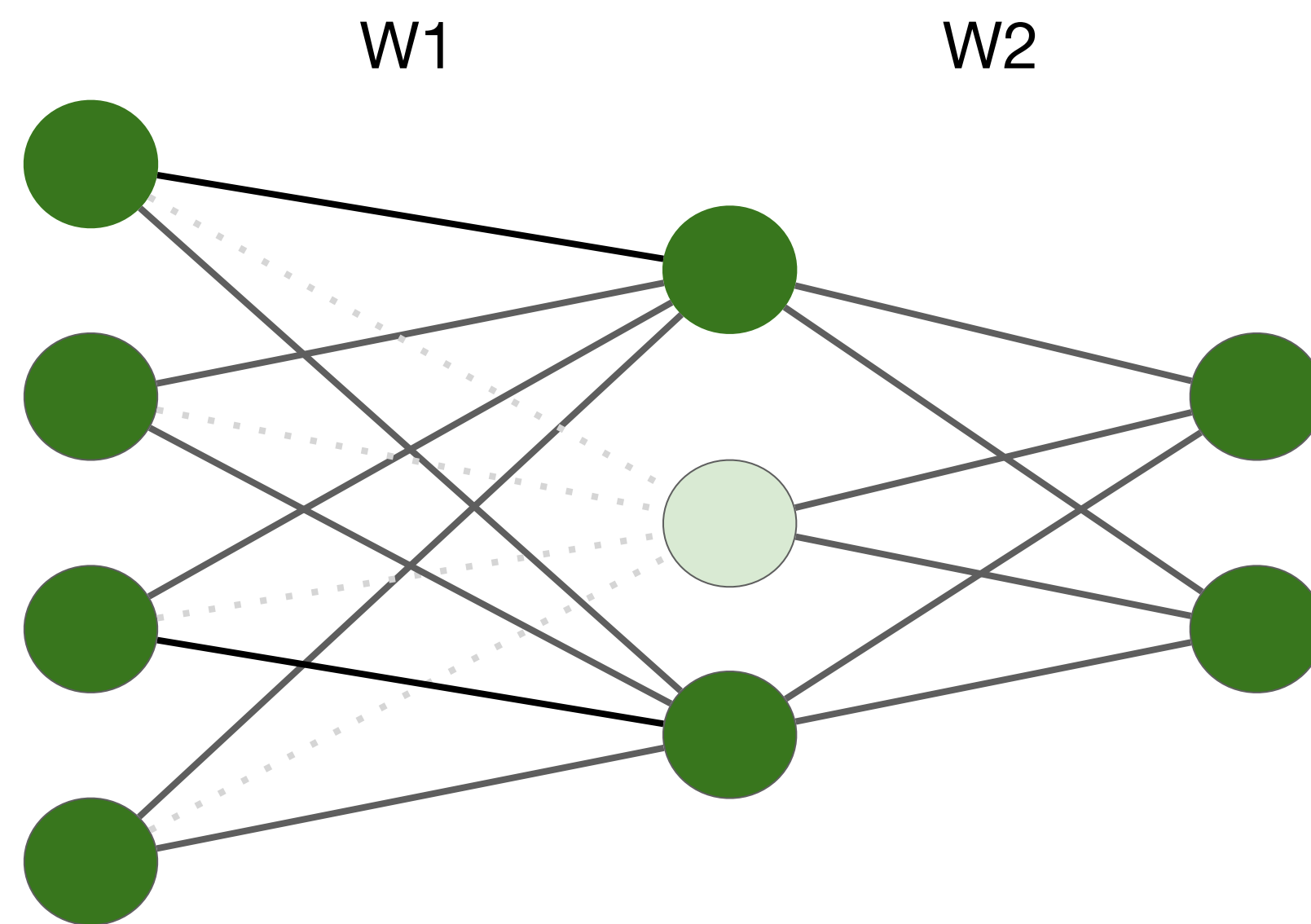
Pruning

- Sparse connectivity inspired by biological neural networks
- **Unstructured pruning** Vs. structured pruning



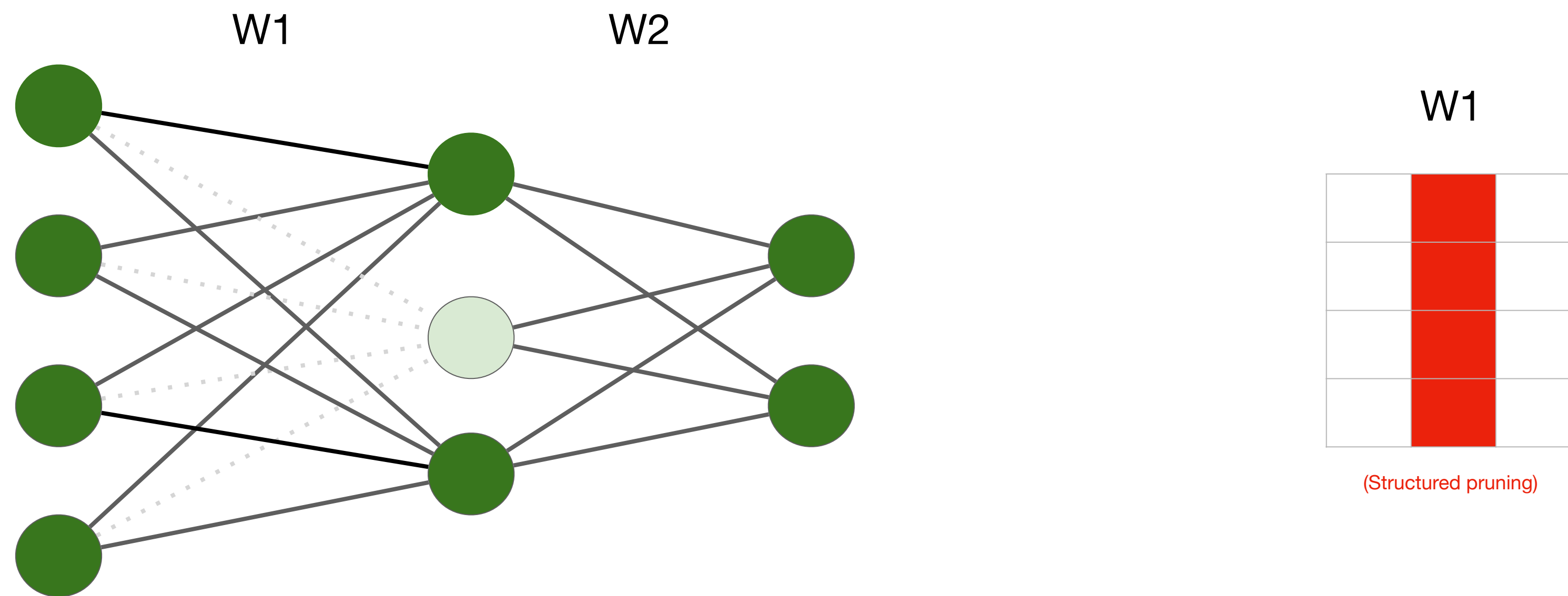
Pruning

- Sparse connectivity inspired by biological neural networks
- Unstructured pruning Vs. structured pruning



Pruning

- Sparse connectivity inspired by biological neural networks
- Unstructured pruning Vs. structured pruning

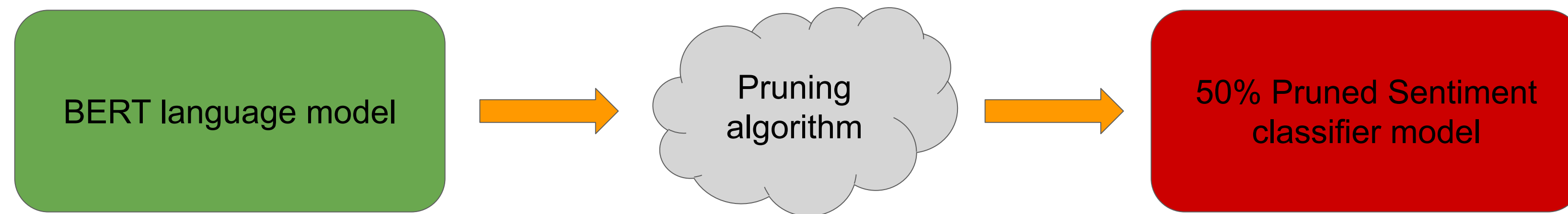


What is the potential benefit of structured pruning?

How to choose pruned weights?

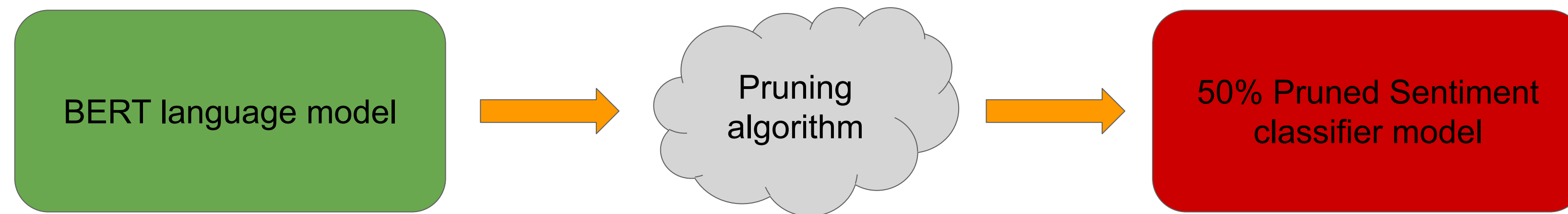
Pruning: case study

- Goal: a BERT-based sentiment classifier model
 - constraints: 50% of weights should be pruned



Pruning: case study

- Goal: a BERT-based sentiment classifier model
 - constraints: 50% of weights should be pruned

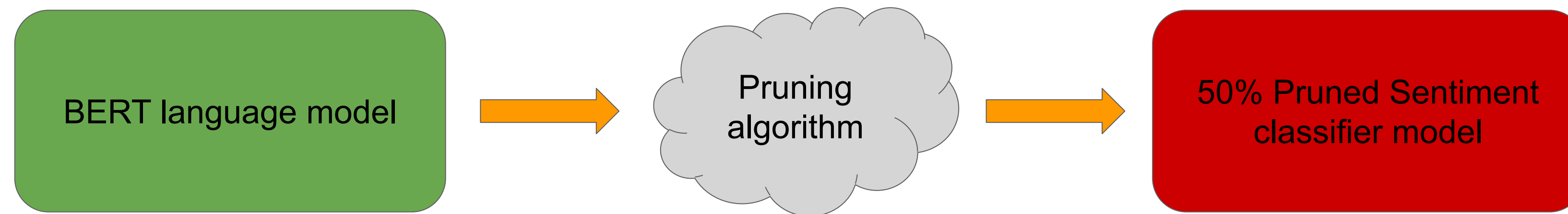


- which weights should be **pruned**?

$$\underbrace{\begin{array}{|c|c|c|c|} \hline 1 & 2 & 0.01 & 3 \\ \hline -0.01 & -1 & -2 & -0.01 \\ \hline -10 & -20 & 0.01 & -0.01 \\ \hline \end{array}}_{\text{Linear Layer}} \begin{array}{|c|} \hline X1 \\ \hline X2 \\ \hline X3 \\ \hline X4 \\ \hline \end{array}$$

Pruning: case study

- Goal: a BERT-based sentiment classifier model
 - constraints: 50% of weights should be pruned



- which weights should be **pruned**?

1	2	0.01	3	X1
-0.01	-1	-2	-0.01	X2
-10	-20	0.01	-0.01	X3
				X4

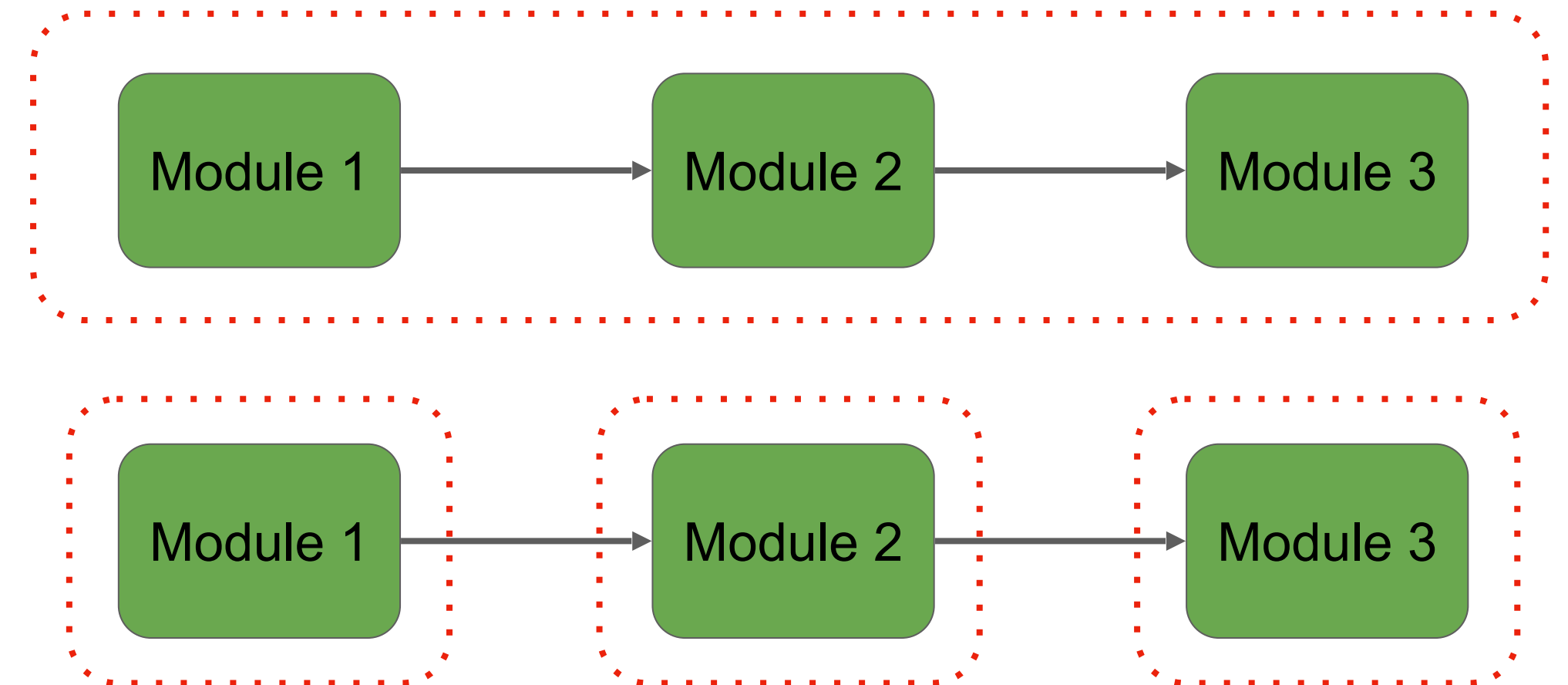
Linear Layer

Weight Pruning Methods

- Magnitude pruning
 - Pruning weights with small magnitude
 - Pruning x% at **global** Vs. **Module** level

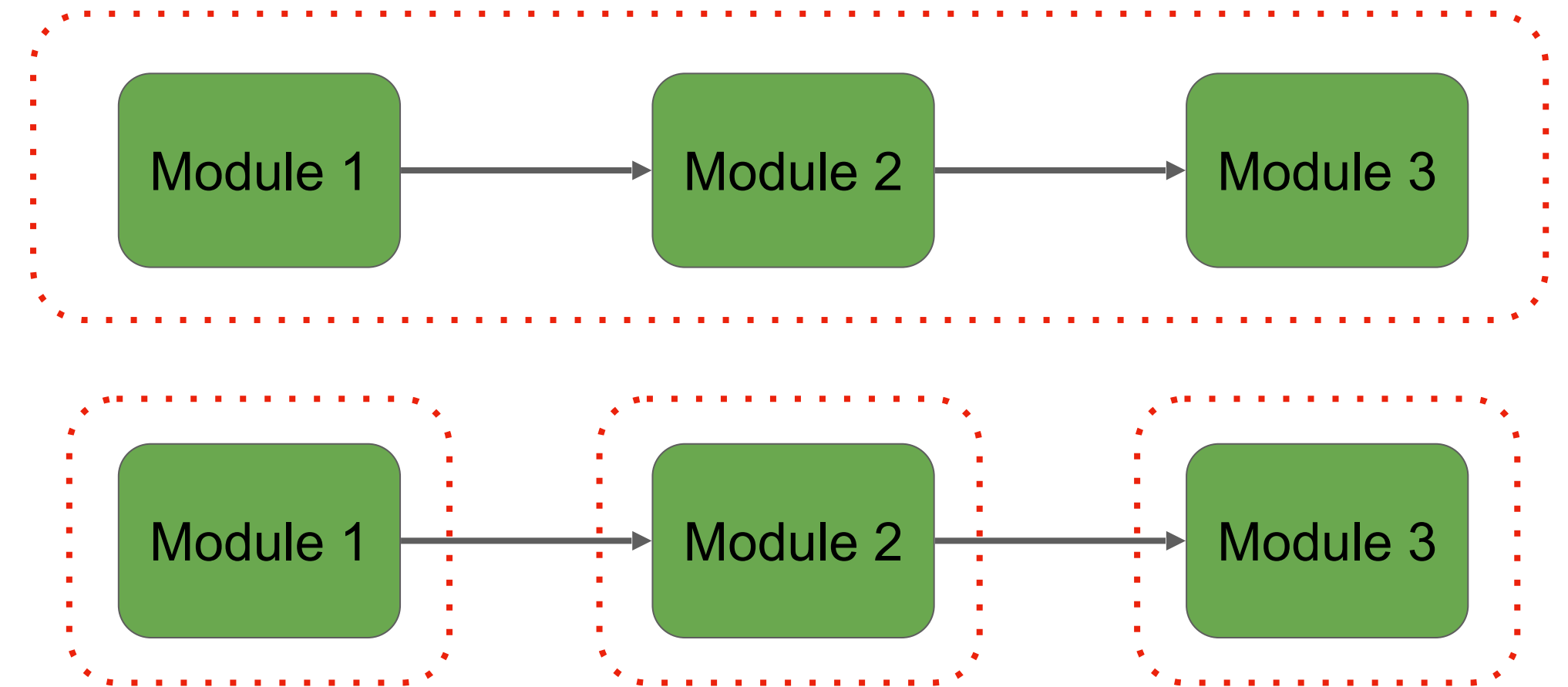
Weight Pruning Methods

- Magnitude pruning
 - Pruning weights with small magnitude
 - Pruning x% at **global** Vs. **Module** level



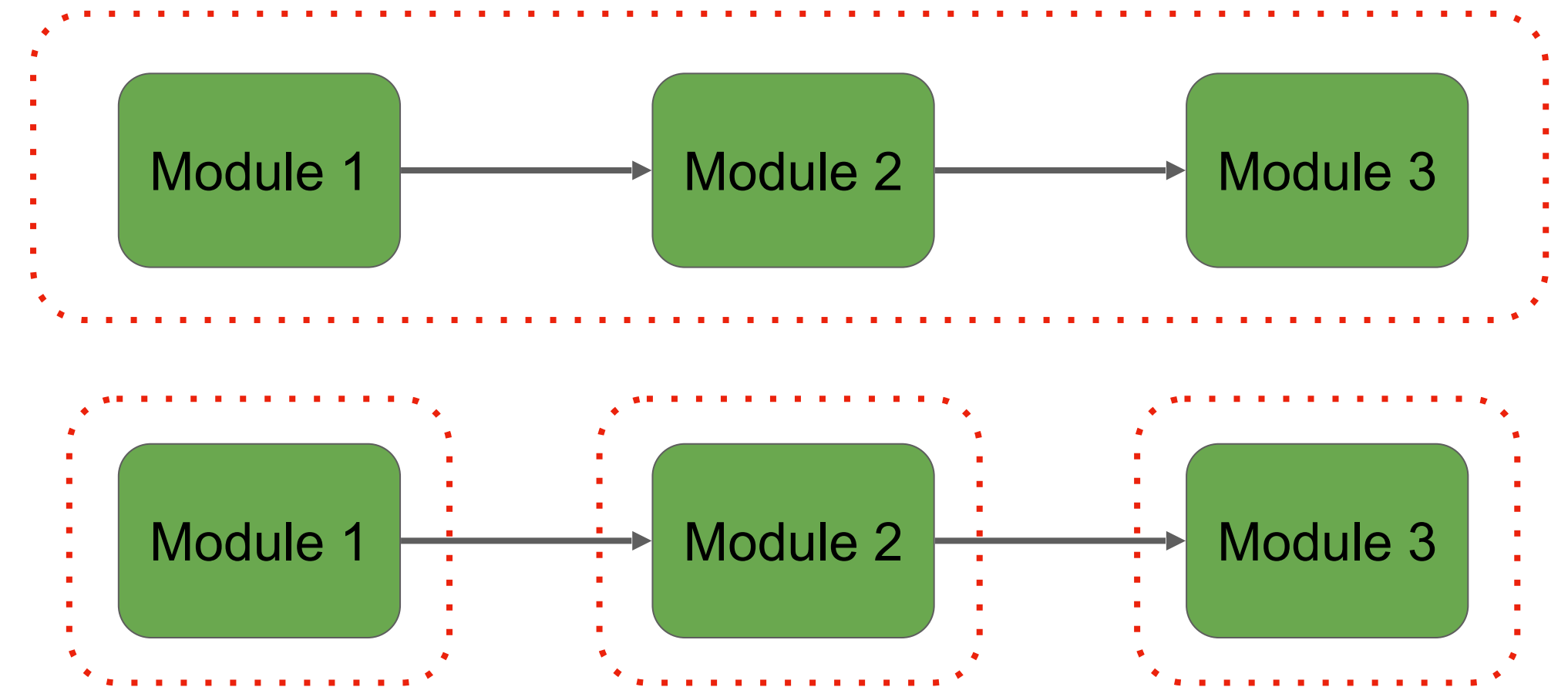
Weight Pruning Methods

- Magnitude pruning
 - Pruning weights with small magnitude
 - Pruning x% at **global** Vs. **Module** level
- Iterative magnitude pruning
 - pruning gradually during training



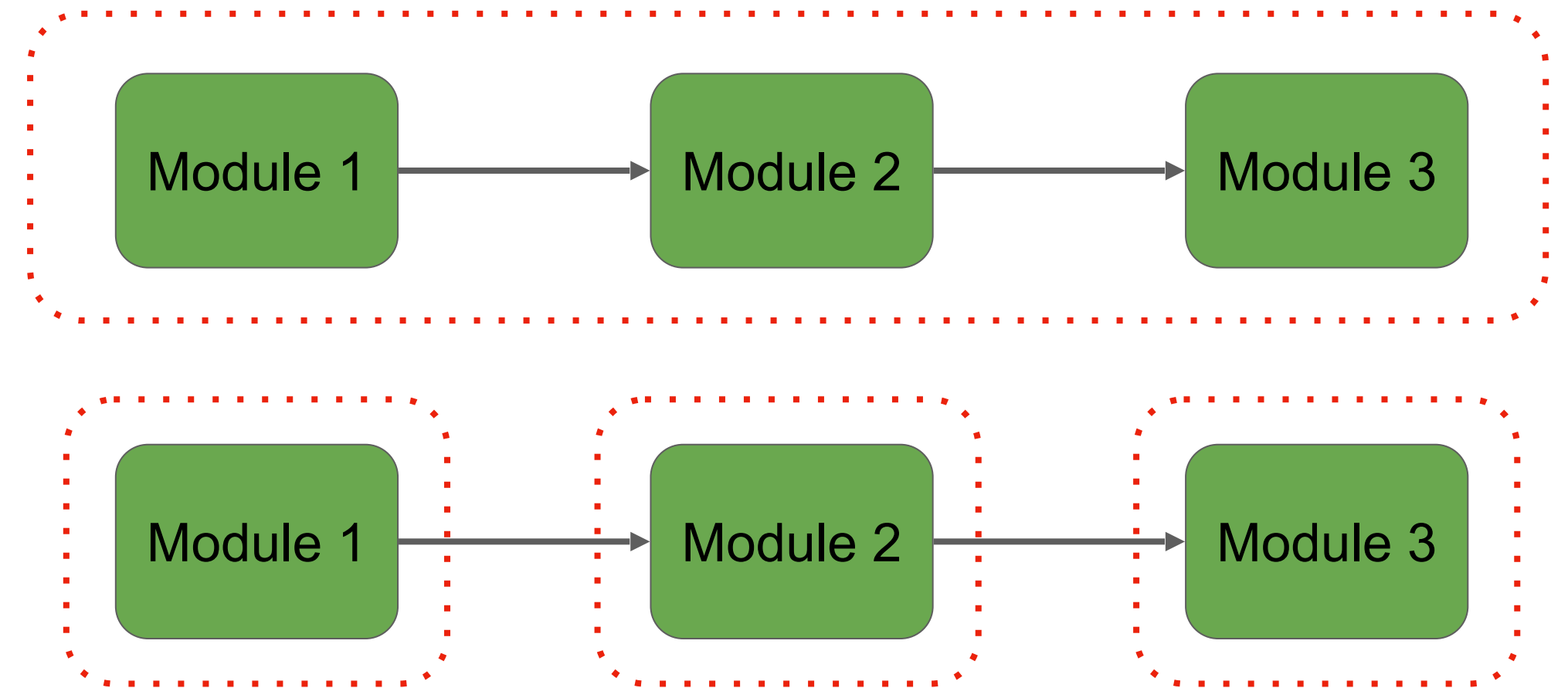
Weight Pruning Methods

- Magnitude pruning
 - Pruning weights with small magnitude
 - Pruning x% at **global** Vs. **Module** level
- Iterative magnitude pruning
 - pruning gradually during training
- Movement pruning



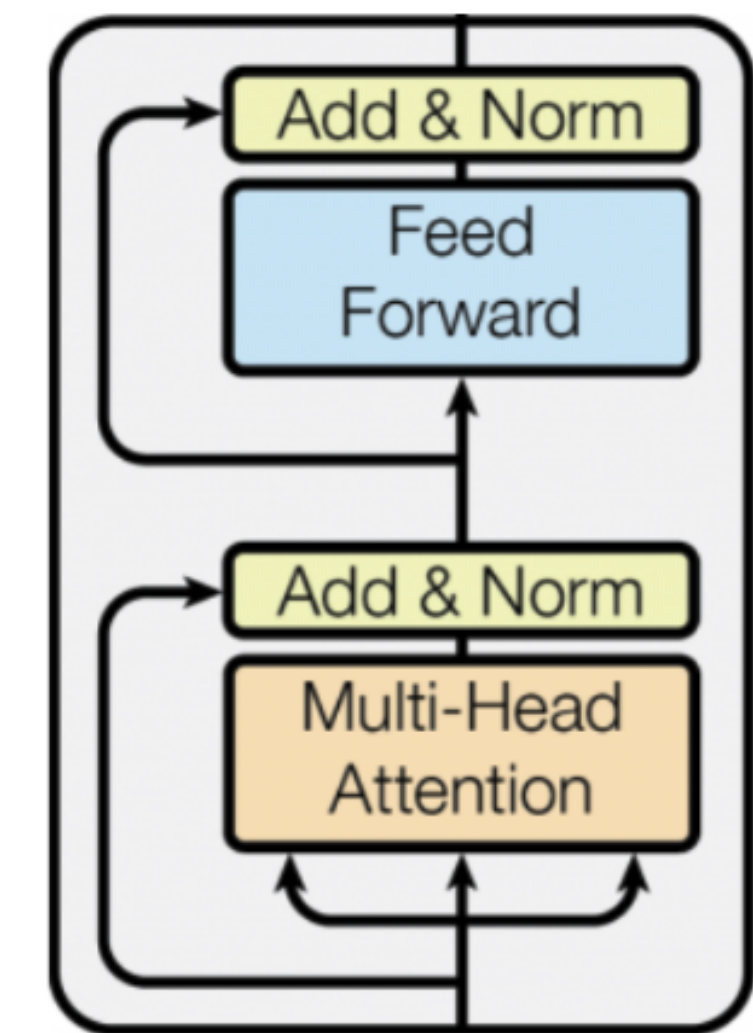
Weight Pruning Methods

- Magnitude pruning
 - Pruning weights with small magnitude
 - Pruning x% at **global** Vs. **Module** level
- Iterative magnitude pruning
 - pruning gradually during training
- Movement pruning
- (Differentiable) masking as a pruning method
 - Example: attention head masking



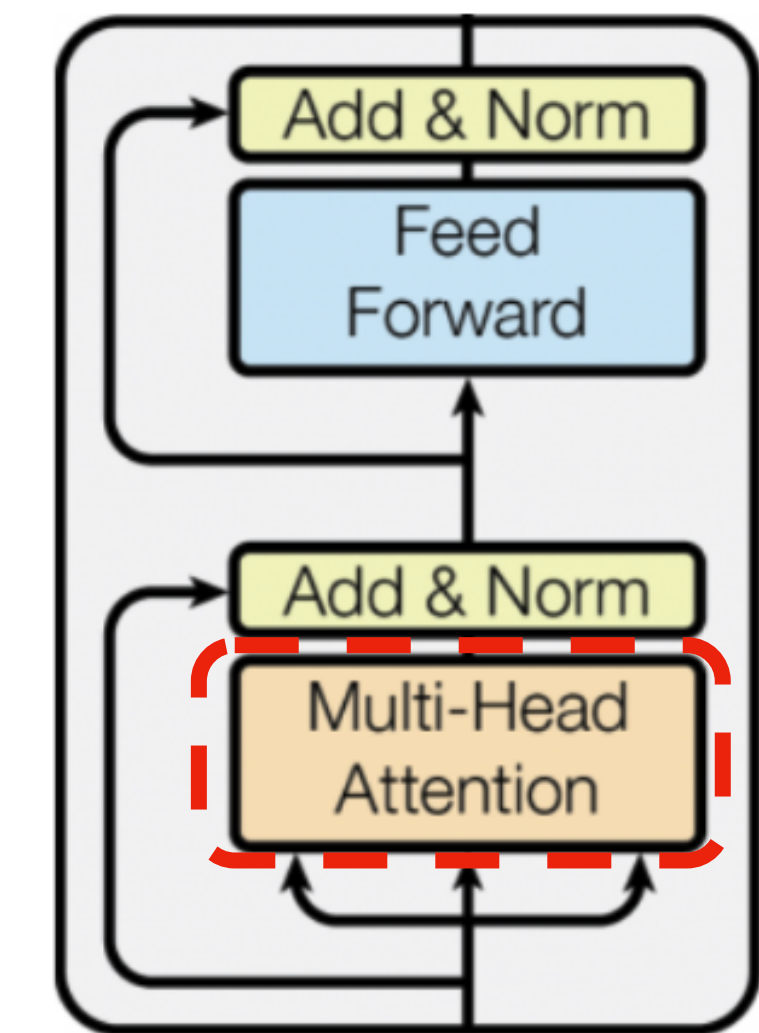
Structured Pruning

- Structured pruning for Transformer language models
 - Pruning neurons



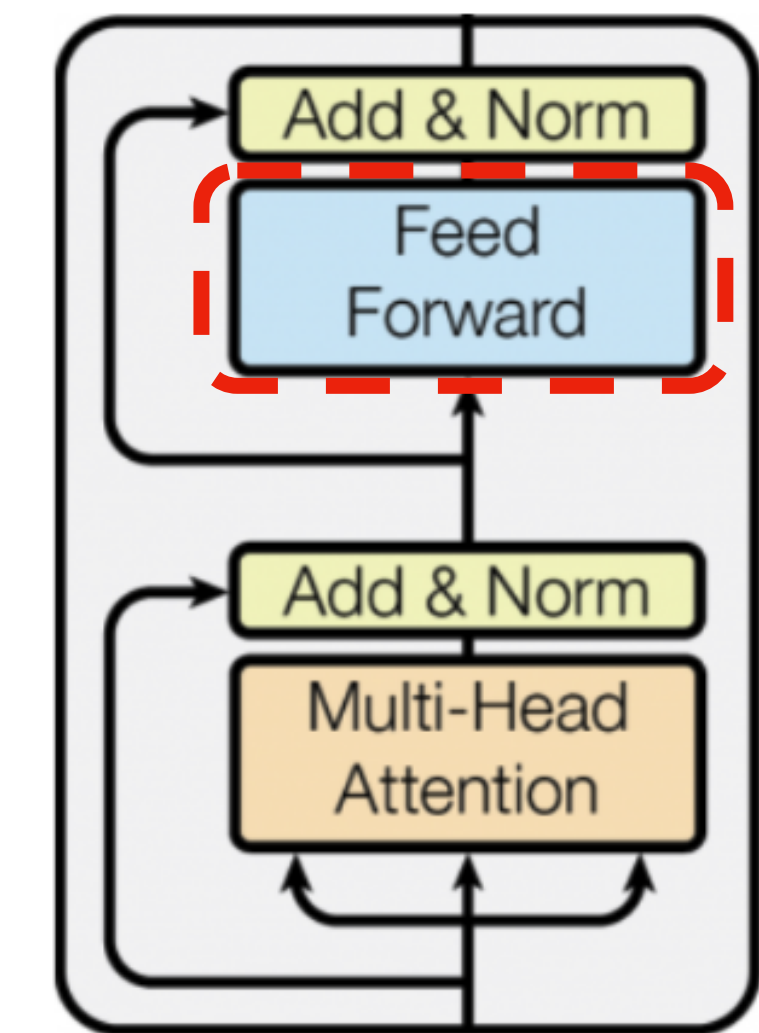
Structured Pruning

- Structured pruning for Transformer language models
 - Pruning neurons
 - Pruning attention heads



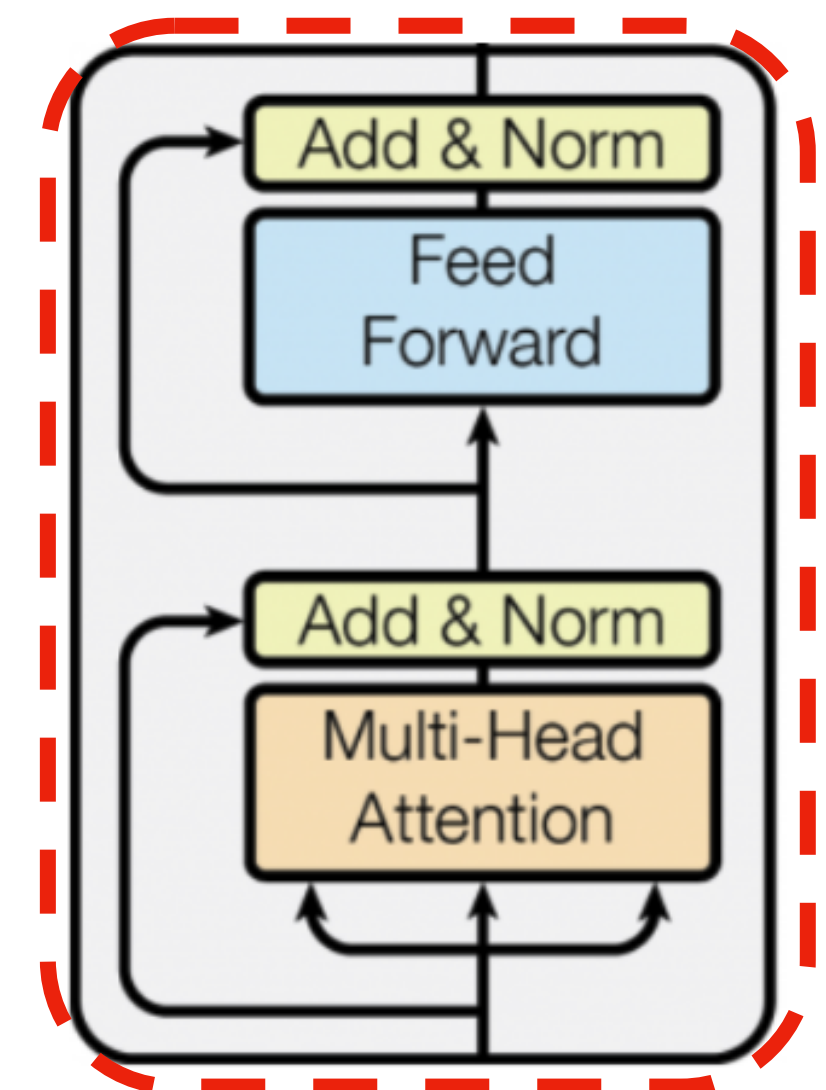
Structured Pruning

- Structured pruning for Transformer language models
 - Pruning neurons
 - Pruning attention heads
 - Pruning sub-layers
 - Example: pruning feed-forward sub-layer



Structured Pruning

- Structured pruning for Transformer language models
 - Pruning neurons
 - Pruning attention heads
 - Pruning sub-layers
 - Example: pruning feed-forward sub-layer
 - Pruning layers
 - Example: pruning the last K layers



Pruning Attention Heads

- How can we prune attention heads?

$$\text{MultiHead}(Q, K, V) = \text{Concat}_i(\text{head}_i)W^O$$

Pruning Attention Heads

- How can we prune attention heads?

$$\text{MultiHead}(Q, K, V) = \text{Concat}_i(\text{head}_i) W^O$$



$$\text{MultiHead}(Q, K, V) = \text{Concat}_i(\boxed{g_i} \cdot \text{head}_i) W^O$$

Pruning Attention Heads

- How can we prune attention heads?

$$\text{MultiHead}(Q, K, V) = \text{Concat}_i(\text{head}_i) W^O$$



$$\text{MultiHead}(Q, K, V) = \text{Concat}_i(\boxed{g_i} \cdot \text{head}_i) W^O$$

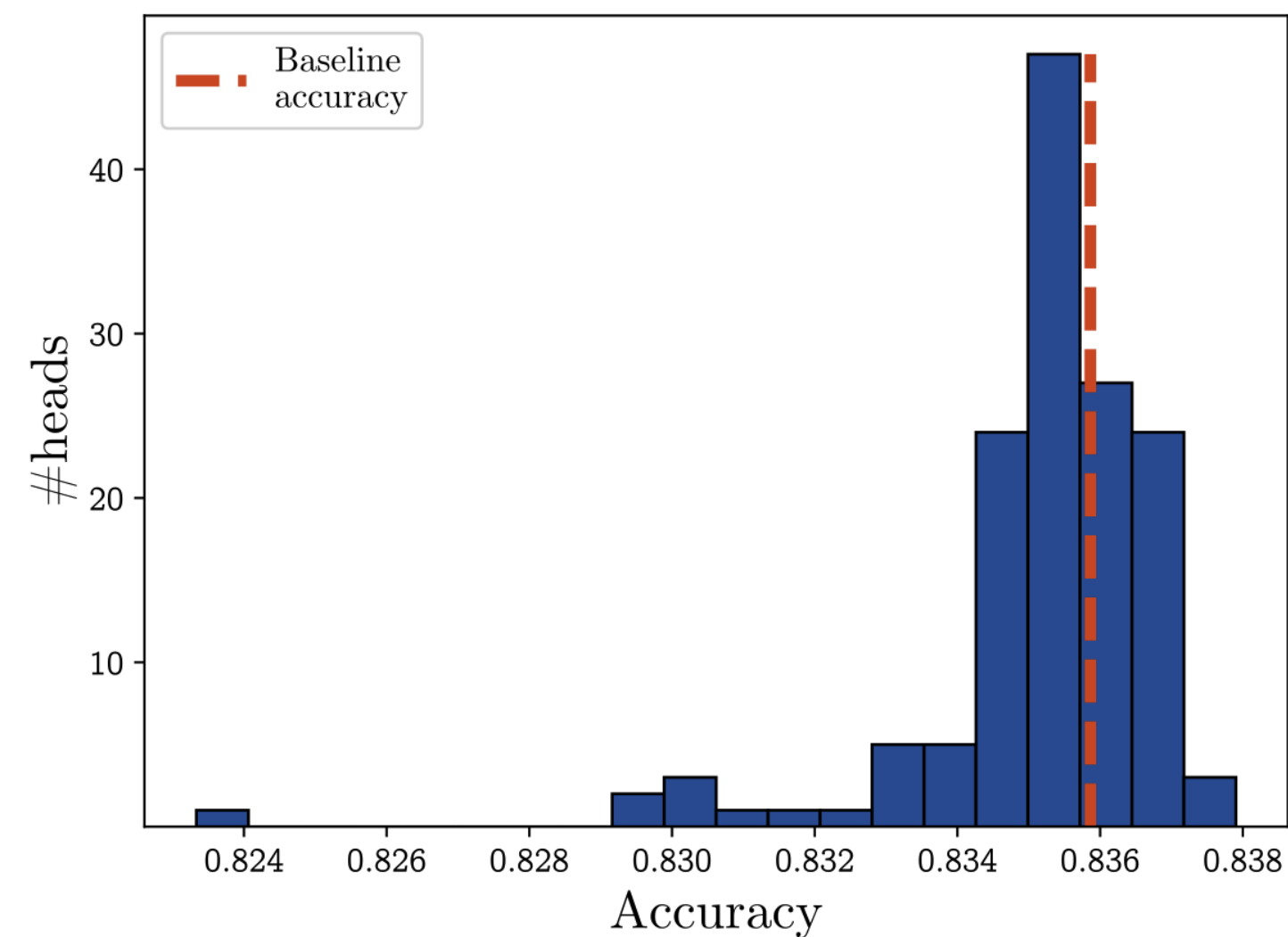
- L0 regularization over attention heads' mask parameter

- Example: Translation task

$$L = L_{xent} + \lambda L_C \quad \lambda = 0.01$$

Pruning Attention Heads

- Large fraction of Transformer attention heads can be removed at test time!



(BERT finetuned on MNLI dataset)

Layer		Layer	
1	-0.01%	7	0.05%
2	0.10%	8	-0.72%
3	-0.14%	9	-0.96%
4	-0.53%	10	0.07%
5	-0.29%	11	-0.19%
6	-0.52%	12	-0.12%

(Delta accuracy by layer when only one head is kept for MNLI BERT model)

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization		
Weight Factorization		
Weight Sharing		
Knowledge distillation		
Sub-quadratic Transformer		

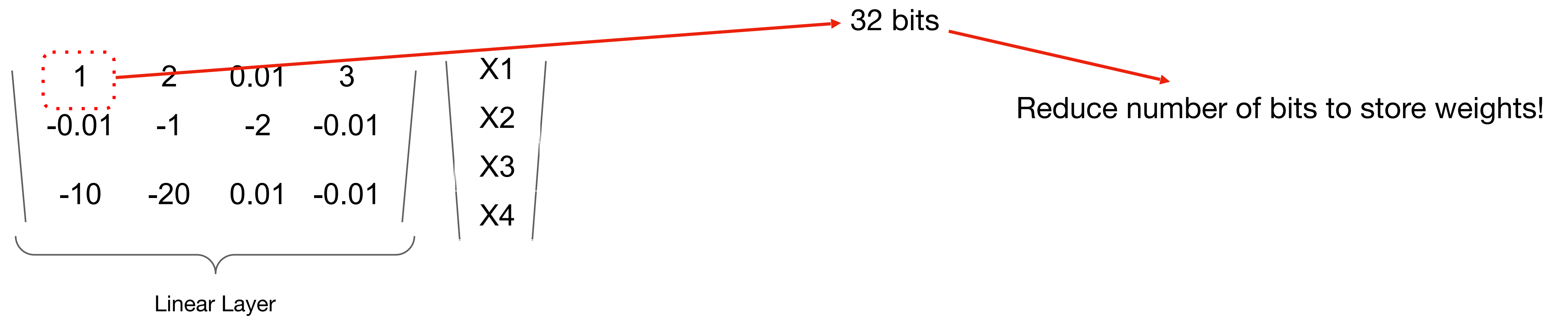
Quantization

- How else can we compress a given neural module?

$$\underbrace{\begin{pmatrix} 1 & 2 & 0.01 & 3 \\ -0.01 & -1 & -2 & -0.01 \\ -10 & -20 & 0.01 & -0.01 \end{pmatrix}}_{\text{Linear Layer}} \begin{pmatrix} X1 \\ X2 \\ X3 \\ X4 \end{pmatrix}$$

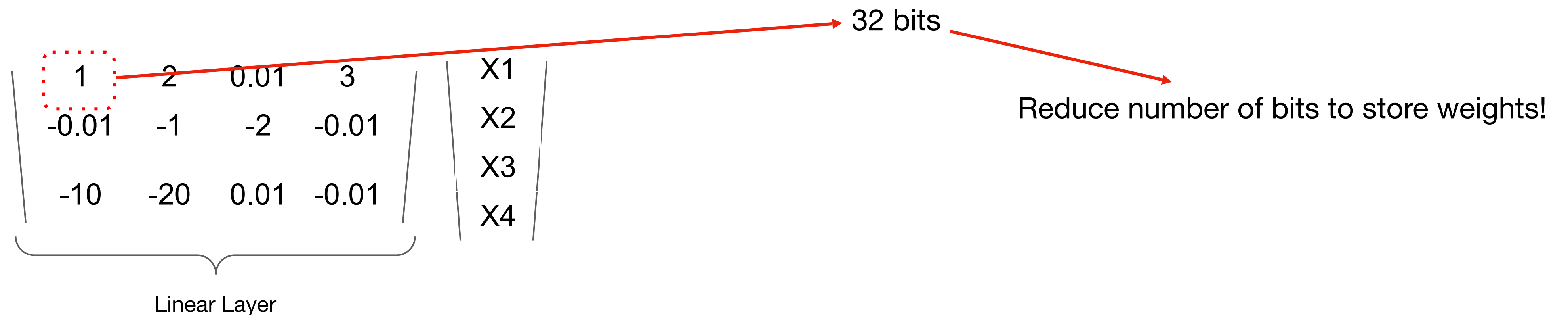
Quantization

- How else can we compress a given neural module?



Quantization

- How else can we compress a given neural module?



- Number of parameters remains the same!
 - Improvement in memory footprint + inference time
- Quantization is mostly applied on a **trained** model

Binarized Network

- Essentially using 1 bit per parameter!
- Deterministic Binarization
 - c_1 and c_2 from K-means over the weights
 - c_1 and c_2 **tuned** on downstream task

$$w_b = \begin{cases} c_1 & \text{if } w \geq (c_1 + c_2)/2 \\ c_2 & \text{if } w < (c_1 + c_2)/2 \end{cases}$$

Binarized Network

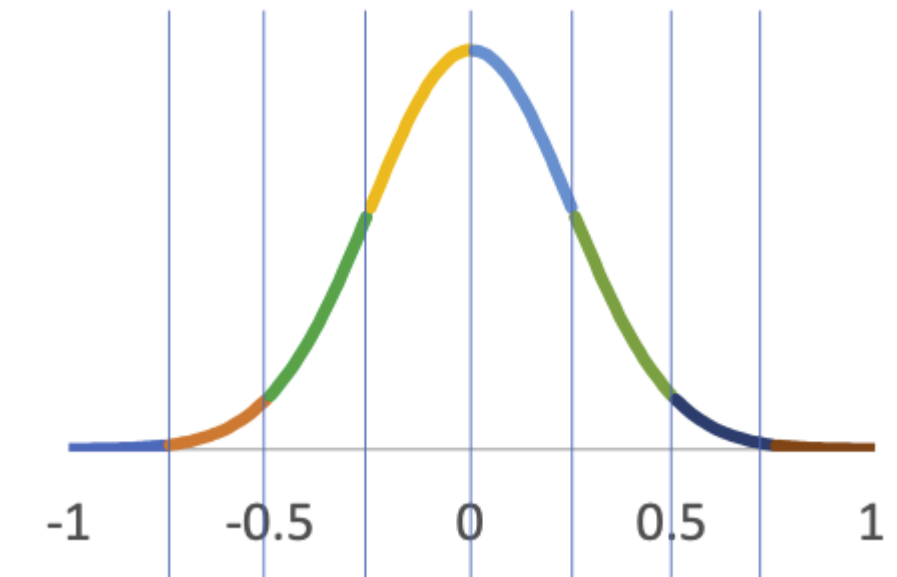
- Essentially using 1 bit per parameter!
- Deterministic Binarization
 - c_1 and c_2 from K-means over the weights
 - c_1 and c_2 **tuned** on downstream task

$$w_b = \begin{cases} c_1 & \text{if } w \geq (c_1 + c_2)/2 \\ c_2 & \text{if } w < (c_1 + c_2)/2 \end{cases}$$

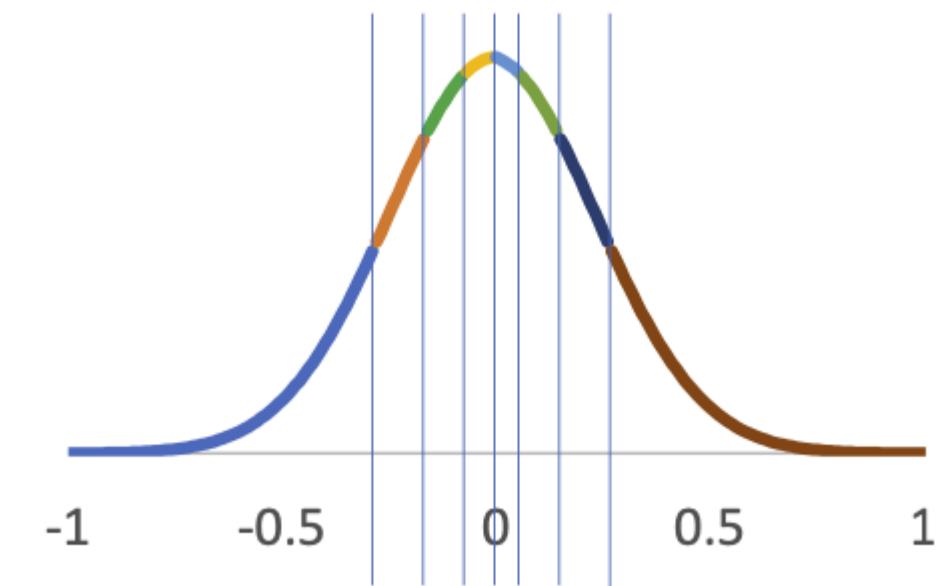
- Question: How can we improve the binarized network performance?

General Quantized Networks

- Uniform Quantization
 - Not necessarily optimal
- Balanced Quantization
 - Better fitted for non-uniform weights!
 - Example: Decide bin boundaries using clustering!



Quantization with 3 bits



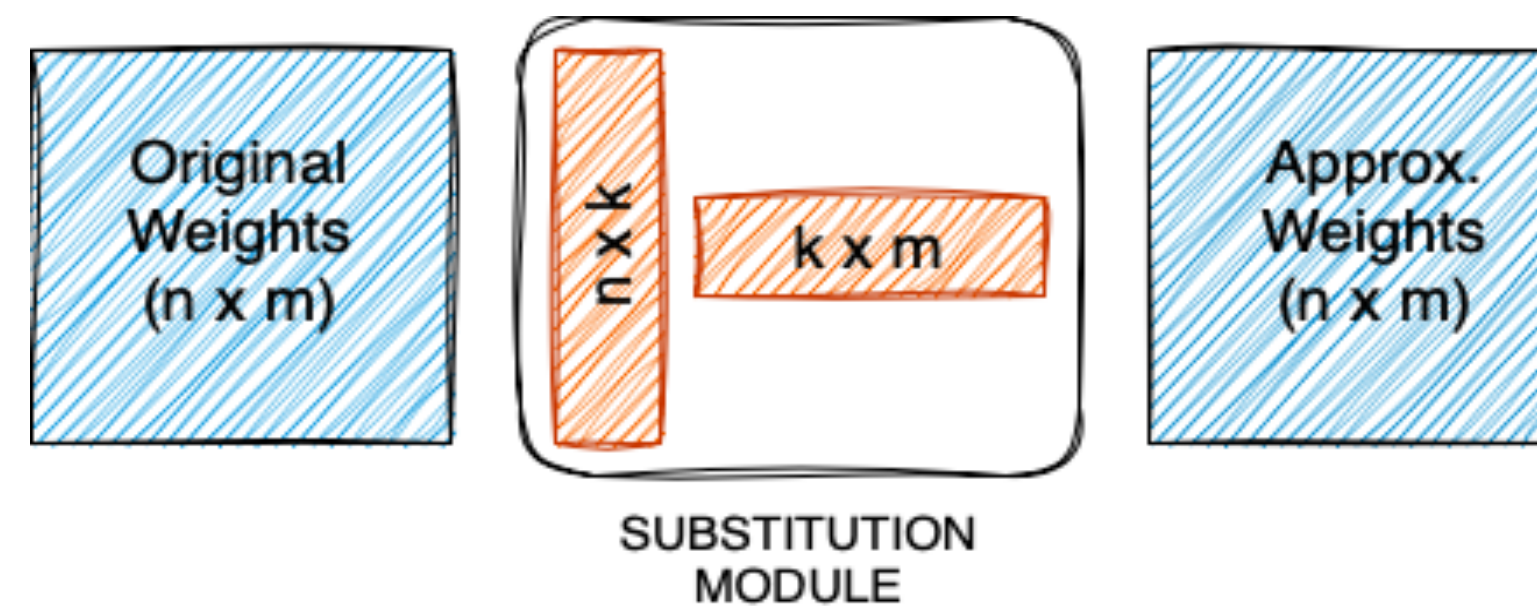
Quantization with 3 bits

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization	Yes	Yes
Weight Factorization		
Weight Sharing		
Knowledge distillation		
Sub-quadratic Transformer		

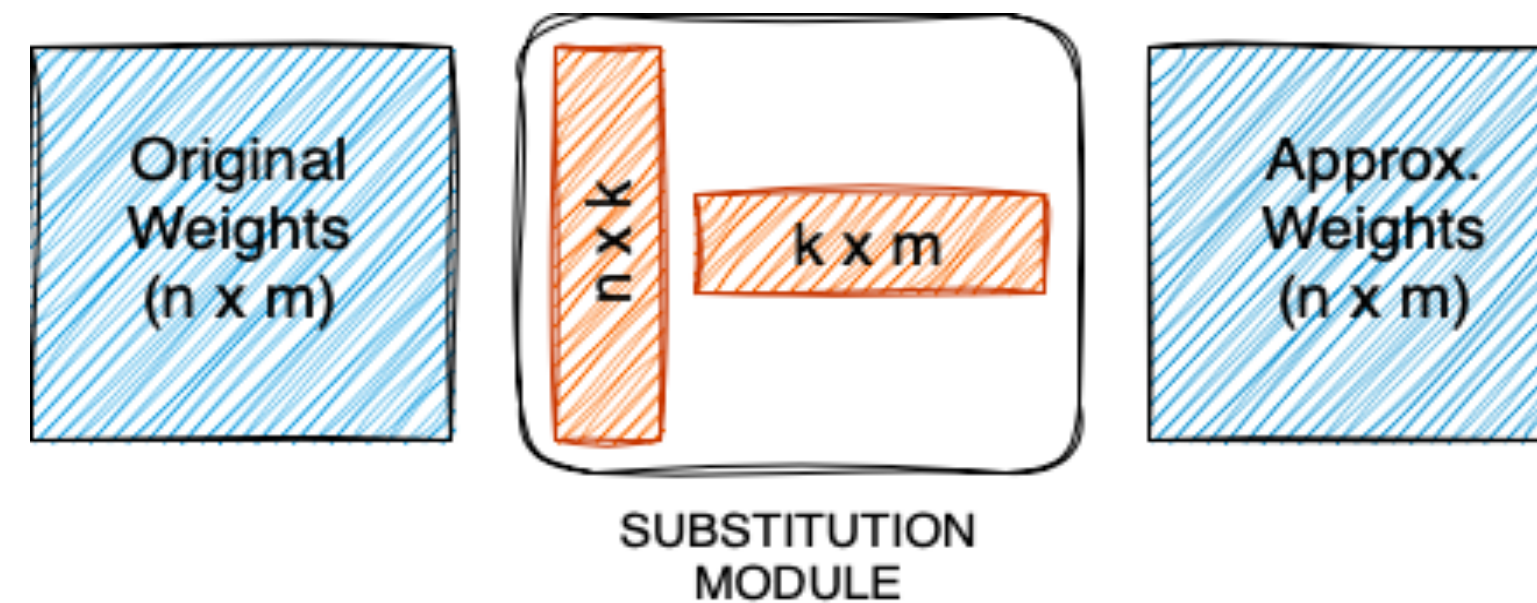
Weight Factorization

- The weight modules are replaced by their factorized matrices



Weight Factorization

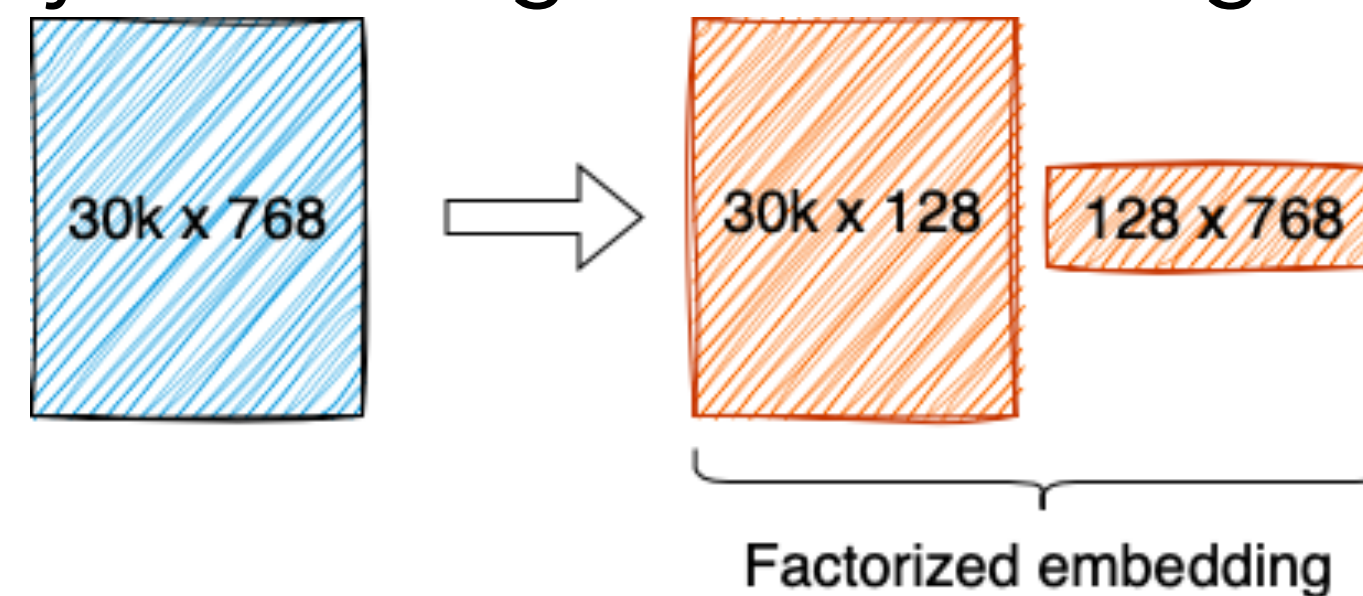
- The weight modules are replaced by their factorized matrices



- Factorization methods
 - Two low-rank matrices (similar to SVD)
 - Tensor decomposition
 - Non-linear factorization by using Auto-encoders

Case Study: ALBERT

- #parameters issue in token embedding matrix
 - ~23M out of 110M parameters in BERT-base
 - More than half of the parameters in mBERT
- Token embedding dimension are generally tied to the hidden dimension
 - What if we disentangle them by factorizing the embedding matrix?



Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization	Yes	Yes
Weight Factorization	Yes	No
Weight Sharing		
Knowledge distillation		
Sub-quadratic Transformer		

Weight Sharing

- A common example of parameter compression/efficiency
 - Finding weight blocks that can share the same weight

Weight Sharing

- A common example of parameter compression/efficiency
 - Finding weight blocks that can share the same weight
- Examples of weight sharing
 - Sharing token embedding and LM decoder head
 - Parameter sharing in the embedding matrix
 - Cross-layer parameter sharing (e.g., ALBERT)

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization	Yes	Yes
Weight Factorization	Yes	No
Weight Sharing	Yes	No
Knowledge distillation		
Sub-quadratic Transformer		

Knowledge Distillation

- Training a smaller **student** network by distilling a large **teacher** model
 - The student's goal is to **imitate** teacher's behavior!
- Can we have the best of the two worlds?
 - Good **performance** of teacher model + **faster** & **parameter**-efficient student model

Knowledge Distillation

- Training a smaller **student** network by distilling a large **teacher** model
 - The student's goal is to **imitate** teacher's behavior!
- Can we have the best of the two worlds?
 - Good **performance** of teacher model + **faster** & **parameter**-efficient student model
- Knowledge distillation Vs. Transfer learning
 - Transfer learning → deals with shared architecture/layers
 - Knowledge distillation → often the student model has a different smaller architecture

How can we distill the teacher's knowledge?

Knowledge Distillation

- Intuition behind knowledge distillation

Knowledge Distillation

- Intuition behind knowledge **distillation**
- Consider a 3-class sentiment analysis dataset
 - We pass the following 2 samples to the teacher model to get class probabilities

Sample #1

	Positive	Negative	Neutral
Groundtruth	1	0	0
Teacher prob.	0.94	0.01	0.05

Sample #2

	Positive	Negative	Neutral
Groundtruth	1	0	0
Teacher prob.	0.67	0.02	0.31

Soft Labels

Knowledge Distillation

- How to leverage **soft** labels for the student model?

- Additional cross-entropy to soft labels (**soft loss**)
- Cross-entropy loss to ground-truth labels → **hard** loss

$$\mathcal{L} = \underbrace{\alpha \cdot \mathcal{L}_{\text{CE}}}_{\text{Hard Loss}} + (1 - \alpha) \cdot \underbrace{\mathcal{L}_{\text{distill}}}_{\text{Soft Loss}}$$

Knowledge Distillation

- How to leverage **soft** labels for the student model?

- Additional cross-entropy to soft labels (**soft loss**)
- Cross-entropy loss to ground-truth labels → **hard** loss

$$\mathcal{L} = \underbrace{\alpha \cdot \mathcal{L}_{\text{CE}}}_{\text{Hard Loss}} + (1 - \alpha) \cdot \underbrace{\mathcal{L}_{\text{distill}}}_{\text{Soft Loss}}$$

- Teacher making confident prediction for easy downstream tasks

- Solution: increase softmax temperature to get suitably soft targets!

Knowledge Distillation

- How to leverage **soft** labels for the student model?

- Additional cross-entropy to soft labels (**soft loss**)
- Cross-entropy loss to ground-truth labels → **hard** loss

$$\mathcal{L} = \underbrace{\alpha \cdot \mathcal{L}_{\text{CE}}}_{\text{Hard Loss}} + (1 - \alpha) \cdot \underbrace{\mathcal{L}_{\text{distill}}}_{\text{Soft Loss}}$$

- Teacher making confident prediction for easy downstream tasks

- Solution: increase softmax temperature to get suitably soft targets!

#	Model	SST-2	QQP	MNLI-m	MNLI-mm
		Acc	F ₁ /Acc	Acc	Acc
1	BERT _{LARGE} (Devlin et al., 2018)	94.9	72.1/89.3	86.7	85.9
2	BERT _{BASE} (Devlin et al., 2018)	93.5	71.2/89.2	84.6	83.4
3	OpenAI GPT (Radford et al., 2018)	91.3	70.3/88.5	82.1	81.4
4	BERT ELMo baseline (Devlin et al., 2018)	90.4	64.8/84.7	76.4	76.1
5	GLUE ELMo baseline (Wang et al., 2018)	90.4	63.1/84.3	74.1	74.5
6	Distilled BiLSTM _{SOFT}	90.7	68.2/88.1	73.0	72.6
7	BiLSTM (our implementation)	86.7	63.7/86.2	68.7	68.3

Case Study: distilBERT

- 6-layer student model distilled from BERT-base (i.e., teacher)
 - Initialize the student from the teacher by taking one layer out of two

Case Study: distilBERT

- 6-layer student model distilled from BERT-base (i.e., teacher)
 - Initialize the student from the teacher by taking one layer out of two
- Distillation on MLM loss
 - Improving LM generalization

I absolutely [MASK] natural language processing field. 

Case Study: distilBERT

- 6-layer student model distilled from BERT-base (i.e., teacher)
 - Initialize the student from the teacher by taking one layer out of two
- Distillation on MLM loss
 - Improving LM generalization

I absolutely [MASK] natural language processing field. → BERT-base

0.241	I absolutely hate natural language processing field.
0.154	I absolutely love natural language processing field.
0.045	I absolutely need natural language processing field.
0.041	I absolutely mean natural language processing field.
0.040	I absolutely missed natural language processing field.
0.034	I absolutely hated natural language processing field.
0.032	I absolutely understand natural language processing field.
0.024	I absolutely loved natural language processing field.
0.023	I absolutely like natural language processing field.
0.020	I absolutely miss natural language processing field.

Case Study: distilBERT

- 6-layer student model distilled from BERT-base (i.e., teacher)
 - Initialize the student from the teacher by taking one layer out of two
- Distillation on MLM loss
 - Improving LM generalization

I absolutely [MASK] natural language processing field.

BERT-base

0.241	I absolutely hate natural language processing field.
0.154	I absolutely love natural language processing field.
0.045	I absolutely need natural language processing field.
0.041	I absolutely mean natural language processing field.
0.040	I absolutely missed natural language processing field.
0.034	I absolutely hated natural language processing field.
0.032	I absolutely understand natural language processing field.
0.024	I absolutely loved natural language processing field.
0.023	I absolutely like natural language processing field.
0.020	I absolutely miss natural language processing field.

- Proposed Loss: MLM + distilling BERT ML

Case Study: distilBERT

- 6-layer student model distilled from BERT-base (i.e., teacher)
 - Initialize the student from the teacher by taking one layer out of two
- Distillation on MLM loss
 - Improving LM generalization

I absolutely [MASK] natural language processing field.

BERT-base

0.241 I absolutely hate natural language processing field.
0.154 I absolutely love natural language processing field.
0.045 I absolutely need natural language processing field.
0.041 I absolutely mean natural language processing field.
0.040 I absolutely missed natural language processing field.
0.034 I absolutely hated natural language processing field.
0.032 I absolutely understand natural language processing field.
0.024 I absolutely loved natural language processing field.
0.023 I absolutely like natural language processing field.
0.020 I absolutely miss natural language processing field.

- Proposed Loss: MLM + distilling BERT ML
- Competitive performance to the teacher

Model	IMDb (acc.)	SQuAD (EM/F1)
BERT-base	93.46	81.2/88.5
DistilBERT	92.82	77.7/85.8
DistilBERT (D)	-	79.1/86.9

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization	Yes	Yes
Weight Factorization	Yes	No
Weight Sharing	Yes	No
Knowledge distillation	Yes...	Yes...
Sub-quadratic Transformer		

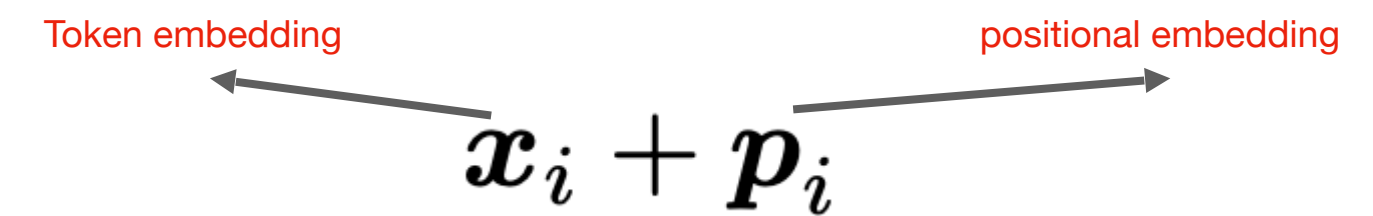
Processing Long Contexts

- Issue with **trained** position embeddings
 - Example: BERT model

The diagram illustrates the decomposition of the sum of token and positional embeddings. It features the mathematical expression $x_i + p_i$ in the center. A red arrow points from the left side of the expression to the label "Token embedding". Another red arrow points from the right side of the expression to the label "positional embedding".

Processing Long Contexts

- Issue with **trained** position embeddings
 - Example: BERT model
- Sinusoidal position embedding
 - Example: (original) Transformer paper



A diagram illustrating the combination of token and positional embeddings. It shows the equation $x_i + p_i$ in the center. An arrow points from the equation to the left, labeled "Token embedding" in red. Another arrow points from the equation to the right, labeled "positional embedding" in red.

$$\omega_k = \frac{1}{10000^{2k/d}} \quad \vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

Processing Long Contexts

- Issue with **trained** position embeddings
 - Example: BERT model
- Sinusoidal position embedding
 - Example: (original) Transformer paper
- Relative positional encoding
 - Rotary Position Embedding (RoPE)
 - Attention with Linear Biases (ALiBi)

A diagram illustrating the decomposition of the sum of token and positional embeddings. The expression $x_i + p_i$ is shown in the center. A red arrow points from the expression to the left, labeled "Token embedding". Another red arrow points from the expression to the right, labeled "positional embedding".

$$\omega_k = \frac{1}{10000^{2k/d}} \quad \vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

ALiBi

- No additive position embeddings in the input layer
- adding a linear bias to each attention score

The diagram illustrates the ALiBi mechanism, showing the addition of a linear bias matrix to the attention scores.

Left matrix (Query-Key dot products):

$q_1 \cdot k_1$				
$q_2 \cdot k_1$	$q_2 \cdot k_2$			
$q_3 \cdot k_1$	$q_3 \cdot k_2$	$q_3 \cdot k_3$		
$q_4 \cdot k_1$	$q_4 \cdot k_2$	$q_4 \cdot k_3$	$q_4 \cdot k_4$	
$q_5 \cdot k_1$	$q_5 \cdot k_2$	$q_5 \cdot k_3$	$q_5 \cdot k_4$	$q_5 \cdot k_5$

Right matrix (Linear bias):

0				
-1	0			
-2	-1	0		
-3	-2	-1	0	
-4	-3	-2	-1	0

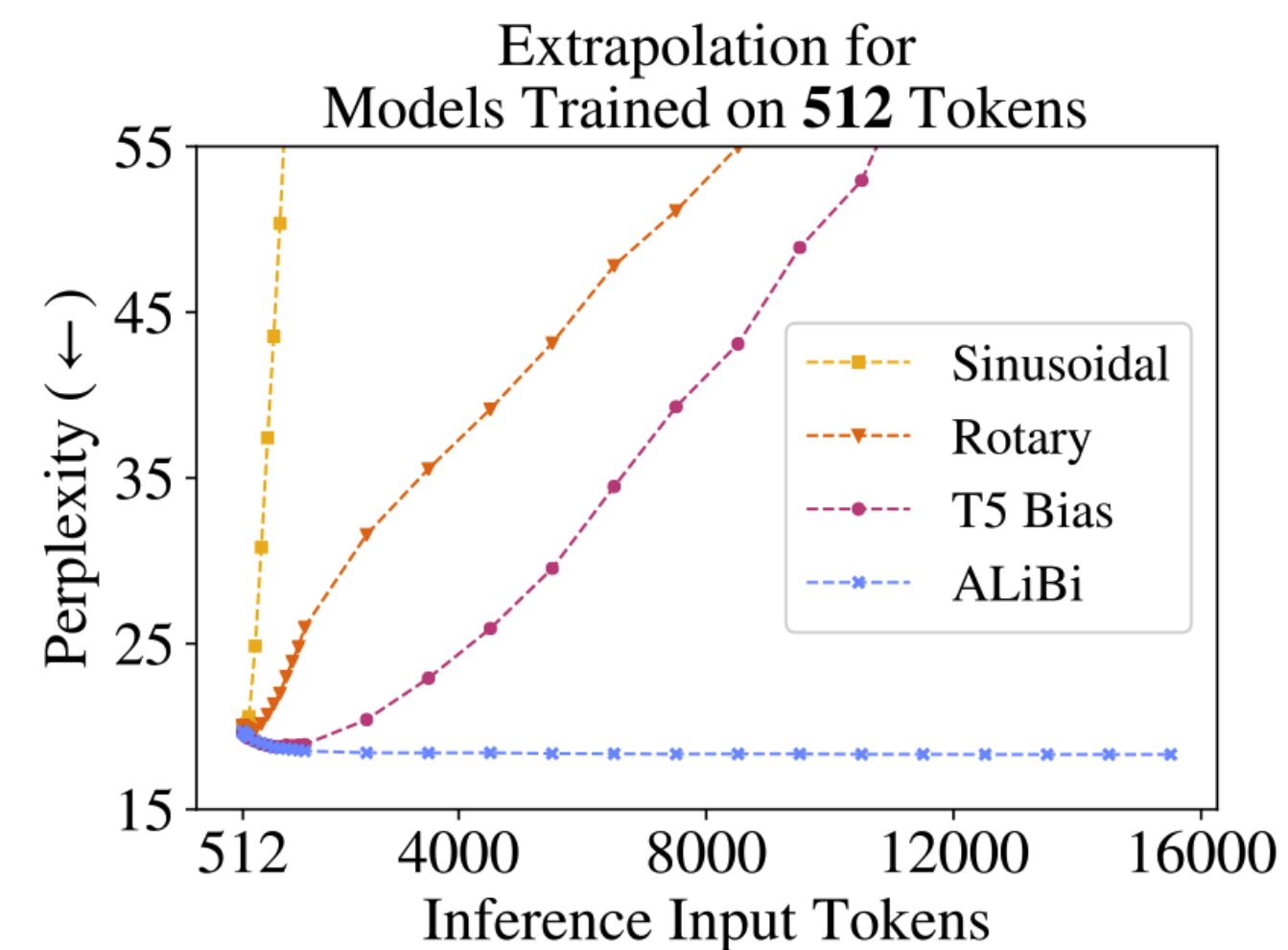
The operation is represented as:

$$\text{Left Matrix} + \text{Right Matrix} \cdot m$$

ALiBi

- No additive position embeddings in the input layer
- adding a linear bias to each attention score

$$\begin{bmatrix} q_1 \cdot k_1 & & & & \\ q_2 \cdot k_1 & q_2 \cdot k_2 & & & \\ q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 & & \\ q_4 \cdot k_1 & q_4 \cdot k_2 & q_4 \cdot k_3 & q_4 \cdot k_4 & \\ q_5 \cdot k_1 & q_5 \cdot k_2 & q_5 \cdot k_3 & q_5 \cdot k_4 & q_5 \cdot k_5 \end{bmatrix} + \begin{bmatrix} 0 & & & & \\ -1 & 0 & & & \\ -2 & -1 & 0 & & \\ -3 & -2 & -1 & 0 & \\ -4 & -3 & -2 & -1 & 0 \end{bmatrix} \cdot m$$



Inference time for long inputs?

Sub-quadratic Transformers

- Time and activation memory grows quadratically with the sequence length
 - Especially important for long sequences
 - Potentially limiting the maximum sequence length

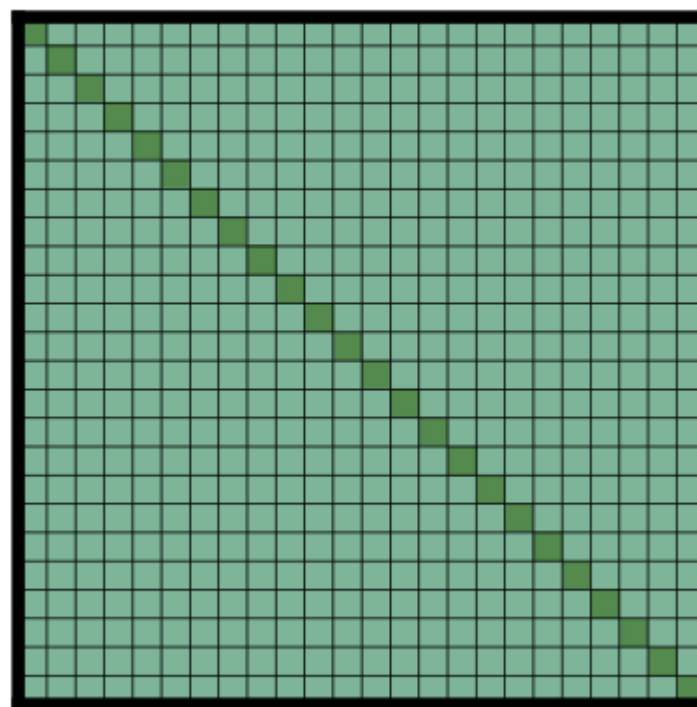
Sub-quadratic Transformers

- Time and activation memory grows quadratically with the sequence length
 - Especially important for long sequences
 - Potentially limiting the maximum sequence length
- Do tokens need to directly attend to every other token?
 - What if attention is performed more locally! → Longformer
 - Masking attention between far tokens (using M matrix)

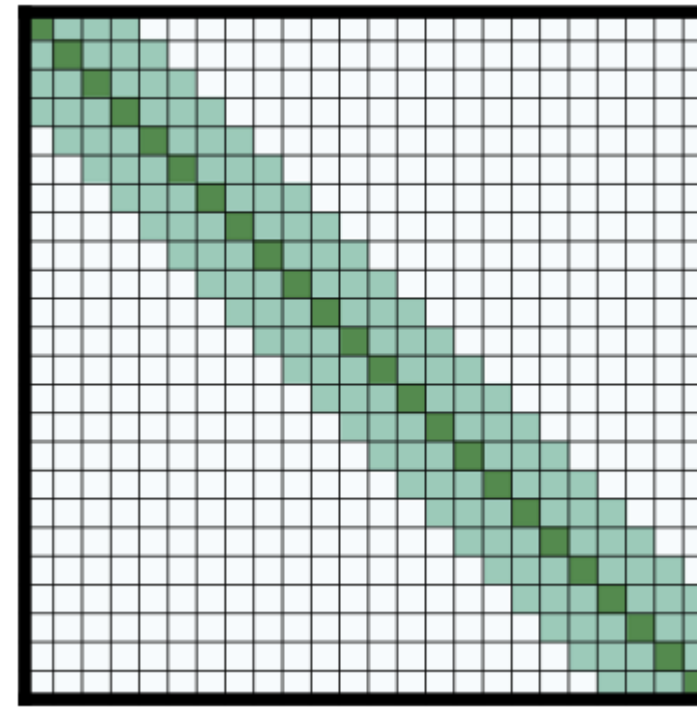
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{M}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \odot \mathbf{M}\right)\mathbf{V}$$

Longformer

- Every token should attend to its neighbor tokens



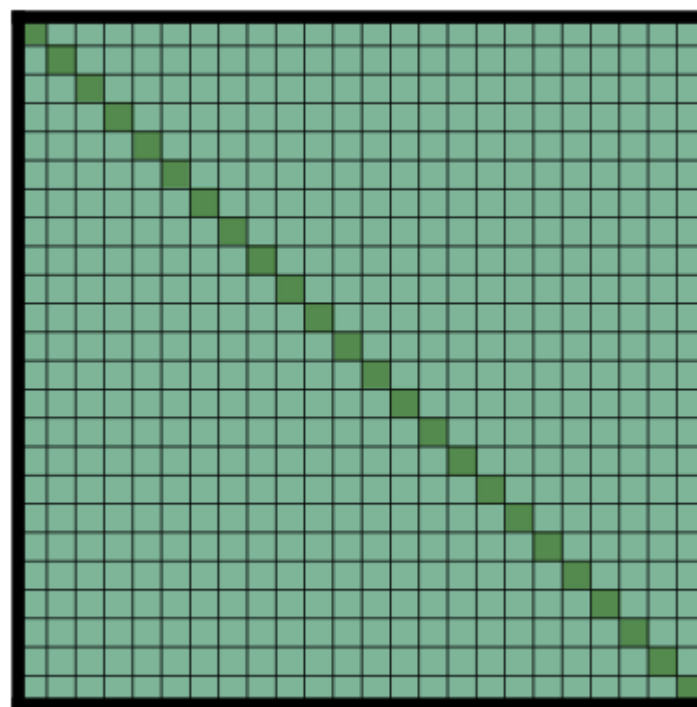
Full n^2 attention



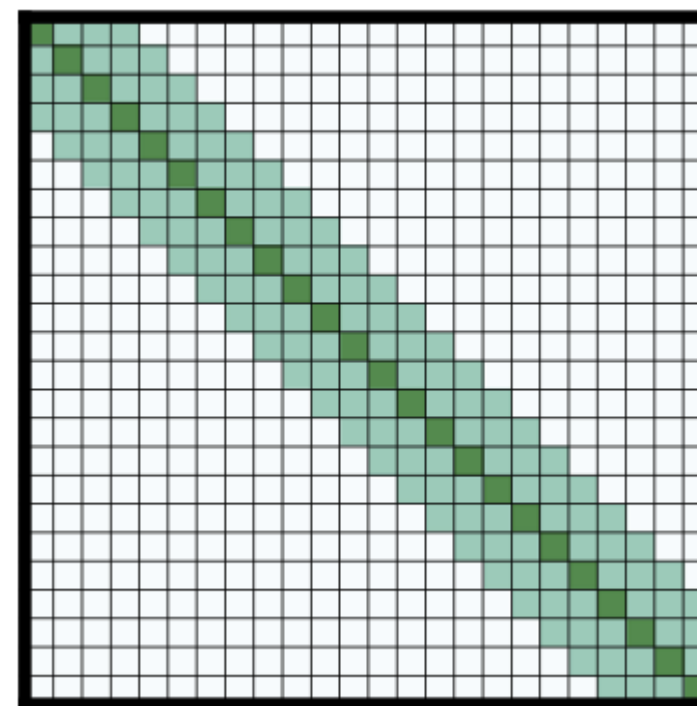
Sliding window attention

Longformer

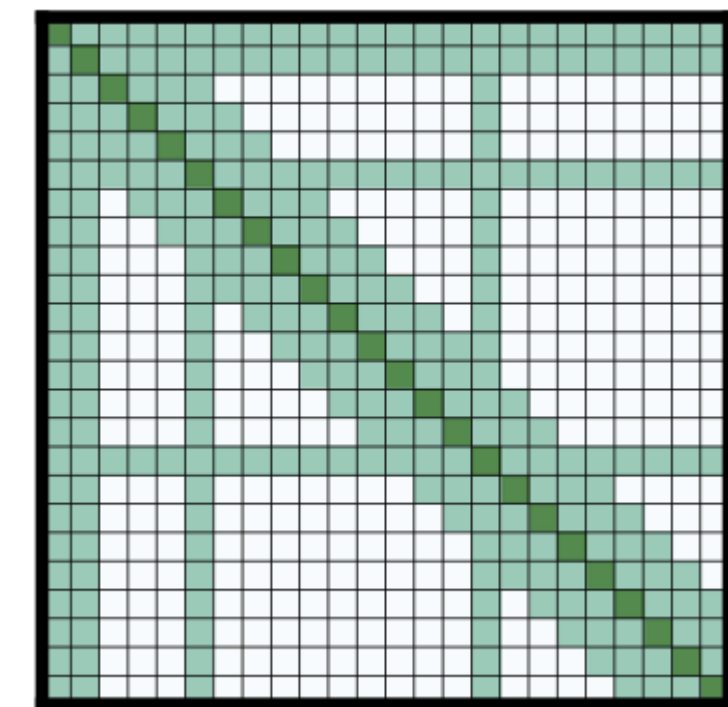
- Every token should attend to its neighbor tokens
- Need for some **global** tokens to **bridge** across the sequence
 - [CLS] for text classification
 - Question tokens for QA



Full n^2 attention



Sliding window attention



Global+sliding window

Methods Overview

Approach	Improvement on memory footprint	Improvement on inference time
Pruning	Y/N	Y/N
Quantization	Yes	Yes
Weight Factorization	Yes	No
Weight Sharing	Yes	No
Knowledge distillation	Yes...	Yes...
Sub-quadratic Transformer	No	Yes

Recap

- Compression leads to improving:
 - Number of parameters
 - Inference time
 - Size-performance trade-off
 - **heavily** compressed large models > **lightly** compressed **small** models
- Different compression techniques
 - Pruning, quantization, factorization, weight sharing, knowledge distillation
- Improvement over quadratic attention mechanism

