

Natural Language Processing



Efficiency

Kevin Lin – UC Berkeley

April 24, 2023

(Many slides credit to GEMLP 2020 High Performance NLP tutorial)

Efficiency

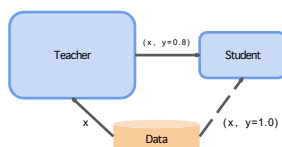


Today

- Knowledge Distillation
- Quantization
- Pruning
- Efficient Attention
- Efficient Architectures

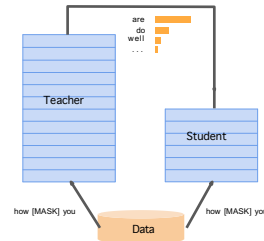
Knowledge Distillation

Hinton et al., 2015
Distilling the Knowledge in a Neural Network



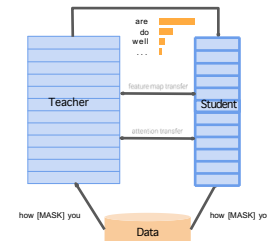
Knowledge Distillation for Pre-training

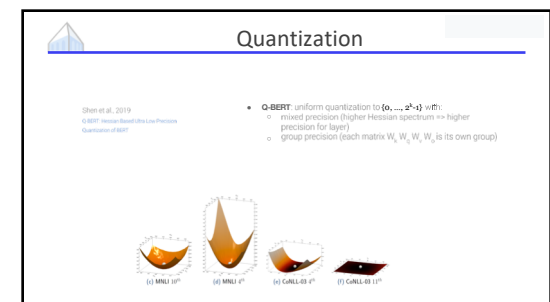
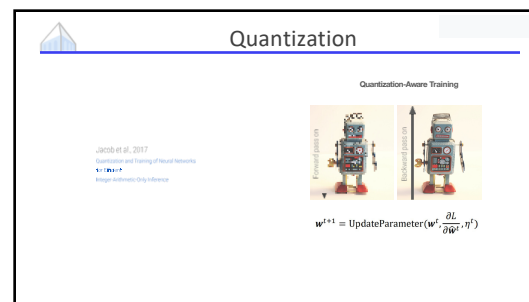
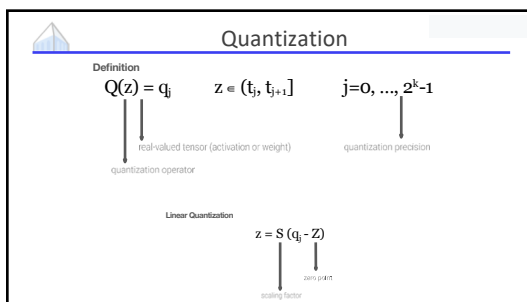
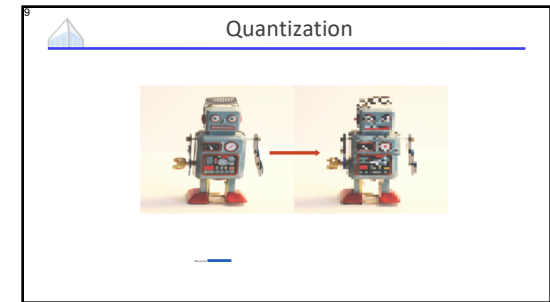
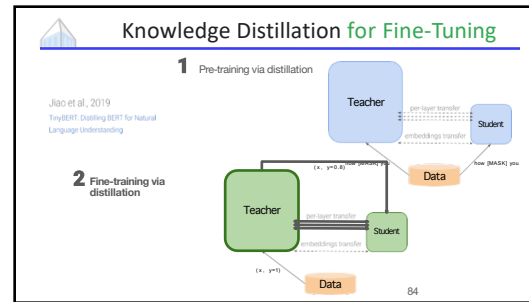
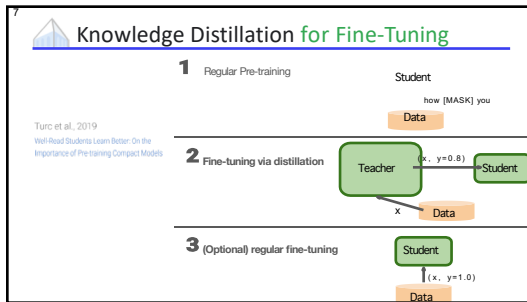
Santh et al., 2019
DisBERT: a distilled version of BERT: smaller, faster, cheaper and lighter



Knowledge Distillation for Pretraining

Sun et al., 2019
MobileBERT: a Compact Task-Agnostic BERT for Resource Limited Devices





Quantization with Distillation

Zhang et al., 2020
Towards BERT: Distillation-aware Quantization of LLMs

Figure 2: Depiction of the proposed distillation-aware quantization of BERT model.

Pruning

Definition
Pruning removes "unimportant" weights from a network.

$$a = (W \odot M) x$$

middle weights
pruning mask

Main Questions (GILMAN, 2018)

- Which weights should be eliminated?
- How should the remaining weights be adjusted?
- How can such network pruning be done in an efficient way?

Pruning

LeCun et al., 1990
OBS: Optimal Brain Damage

Hessle and Stock, 1993
OBS: Second-order derivatives for network pruning. Optimal Brain Surgeon

Main Idea:

- Start with a "reasonably large" network
- Train it to convergence
- Prune in multiple iterations, based on second-order derivatives:
 - OBS: prune and train
 - OBS: prune and update weights based on second-order statistics

Lottery Ticket Hypothesis

The Lottery Ticket Hypothesis: A randomly-initialized, dense neural network contains a subnetwork that is initialized such that—when trained in isolation—it can match the test accuracy of the original network after training for at most the same number of iterations.

Frankle and Carbin, 2018
Finding Sparse Trained Networks Using Sparse Trained Neural Networks

Searching for Tickets: $W^{(0)} \rightarrow W^{(1)} \rightarrow W^{(2)}$

Magnitude Pruning: $W^{(1)} \rightarrow W^{(2)}$

Frankle & Carbin, 2018
Top-0.001 Pruning

Movement Pruning

Barth et al., 2020
Movement Pruning: Adaptive Sparsity by $\text{Tr}(\text{Tr}(\text{Tr}))$

- First-order strategy: "Instead of selecting weights that are far from zero, we retain connections that are moving away from zero during the training process"
- The pruning mask M is learnt together with the model parameters:
 - hard version: $M = \text{Top-OS}$ where score S is learnt and γ is a hyperparameter
 - soft version: $M = (S > \gamma)$ where score S is learnt and threshold γ is a hyperparameter

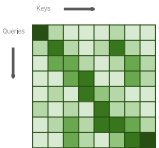
(a) Magnitude pruning (b) Movement pruning

Movement Pruning

	Unstructured Pruning	Structured Pruning
Storage	\sqrt{N}	\sqrt{N}
Inference	\sqrt{N}	\sqrt{N}
Flexibility	\sqrt{N}	\sqrt{N}

Efficient Attention

Goal:
Approximate the computation of attention
via more efficient operations



Keys →
↑
Queries

Efficient Attention

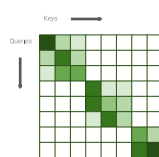
- Data-Independent
- Data-Dependent
- Kernels
- Recurrence
- I/O Aware-Attention

Data-Independent Patterns

Blockwise Patterns

Divide sequence into local blocks and restrict attention within them

Examples:
Blockwise Transformer (Gou et al., 2020)
Local Attention (Parmar et al., 2018)



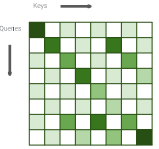
Keys →
↑
Queries

Data-Independent Patterns

Strided Patterns

Skip some query-key pairs.
Quadratic in sequence length / stride

Examples:
Sparse Transformer (Child et al., 2019)
Longformer (Beltagy et al., 2020)



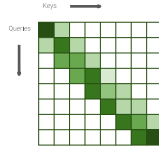
Keys →
↑
Queries

Data-Independent Patterns

Diagonal Patterns

Compute attention over the diagonal.
Linear in sequence length and window size.

Examples:
Longformer (Beltagy et al., 2020)
Big Bird (Zoph et al., 2020)



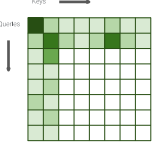
Keys →
↑
Queries

Data-Independent Patterns

Global Attention

Applied to one or a few special tokens, often prepended to the sequence.
Usually combined with other patterns

Examples:
Big Bird (Zoph et al., 2020)
Longformer (Beltagy et al., 2020)
ETC (Vaswani et al., 2020)



Keys →
↑
Queries

Data-Dependent Patterns

Buckets

Create buckets/clusters and compute attention within.

Ideally, buckets should contain the highest attention weights in the matrix

Examples:
[Reformer \(Kitaev et al., 2020\)](#)
[Routing Transformer \(Bayer et al., 2020\)](#)

Data-Dependent Patterns

Buckets: Hashing

[Locality Sensitive Hashing \(LSH\)](#)

Key idea: take a random projection matrix R , compute hash for a vector x through:

$$h(x) = \arg \max_j (x R_j - \sigma R_j)$$

Examples:
[Reformer \(Kitaev et al., 2020\)](#)

Data-Dependent Patterns

Compression

E.g. pooling, strided convolution, low-rank projections with learnable weights

Examples:
[Compressed Attention \(Liu et al., 2019\)](#)
[Lionformer \(Wang et al., 2020\)](#)
[Synthesizers \(Tay et al., 2020\)](#)

Kernel

Kernels

Recap: attention in its general form uses a similarity function $\phi(Q_i, K_j)$

However, we can simplify things with a decomposable **kernel**:

$$\phi(Q_i, K_j) = \phi(Q_i)^T \phi(K_j)$$

Kernel

Kernels

Recap: attention in its general form uses a similarity function $\phi(Q_i, K_j)$

However, we can simplify things with a decomposable **kernel**:

$$\phi(Q_i, K_j) = \phi(Q_i)^T \phi(K_j)$$

Independent of query!

Kernel

Kernels

In vectorized form:

$$O = \phi(Q) \phi(K)^T V$$

Compute this $d' \times d$ matrix first

This allows us to compute attention in **linear** time with respect to sequence length!

In [Katharopoulos et al., 2020](#)

Recurrence

Compressive Transformers
[\(Vaswani et al., 2020\)](#)

Dual memory system:

- Primary mem. contains activations from previous segment
- Secondary mem. is compresses activations from all previous segments

I/O Awareness

FlashAttention (Dao et al., 2019)

Mixture-of-Experts

FFN(x) = $\max(0, xW_1 + b_1)W_2 + b_2$

(Dao et al., 2017)
 Laptev et al., 2020

Mixture-of-Experts

Version 2 @ [Laptev et al., 2020](#)

$$\mathcal{L} = \ell_{all} + k * \ell_{aux}$$

where k is a constant loss weight (a good value is 0.1, usually between 0.01 and 1.0)

- Random dispatch (use 2nd expert proportionally to the softmax gate probability)
- Have a frequency cutoff – a token budget – for each expert. If this budget is exceeded the expert degenerates to a zero matrix. This effectively reduces the output of the MoE layer to zero and thus only the residual connection output around the MoE layer is fed to the next layer.

Mixture-of-Experts

- Works well on diverse data like multilingual machine translation
- Can be difficult to train due to balancing specialization issues
- Only faster than transformers if you can run it with a large enough batch size to saturate distributed experts
- If you scale the model across a cluster, you will need excellent interconnect performance (TPU v4 Pod, NVIDIA SuperPod)

Structure State Spaces

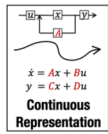
Continuous-time (CTM)	Recurrent (RNN)	Convolutional (CNN)
✓ Continuous data Irregular sampling	✓ Unbounded context Stateful inference	✓ Easy optimization Parallelizable training
✗ Complex, very inefficient Vanishing gradients	✗ Inefficient training Vanishing gradients	✗ Inefficient inference Bounded context

Existing model families have clear tradeoffs
 All struggle with long-range dependencies (LRD)

S4 (Su et al., 2022)

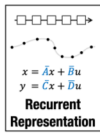


Structure State Spaces



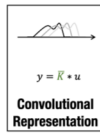
- ✓ mathematically tractable
- ✓ handles irregular data

Discretize



- ✓ unbounded context
- ✓ efficient inference

Unroll



- ✓ easy optimization
- ✓ parallelizable training