

Lecture 08:

Efficient DNN Training, Parameter Efficient Finetuning, Speculative Decoding

Notes

- Lab 2 is due this Friday, Lab 1 grade will post this weekend.
- Grade of Lab 1 will be posted this weekend.
- Think about the project, discuss with me during office hours or after class.
- Midterm
 - Oct 29, in class.
 - Will cover materials up to this lecture (Oct 22)
- Final presentation
 - Virtual
 - Dec 16 and Dec 17
 - Final report due at Dec 19



Notes

- Lab 3 will be posted this weekend.
 - Speculative decoding
 - Yunhai will hold extra office hours weekly about Lab 3 on Wednesday 2-3pm
 - https://nyu.zoom.us/j/98317554792
- I will traveling this Friday, so office hour will hold online. Additional office hours can be arranged upon request.
- Quiz today on Efficient LLM.



Recap

- Large Model Data Distribution
- Large Model Quantization
- Large Model Pruning
- Low-rank Decomposition for LLM

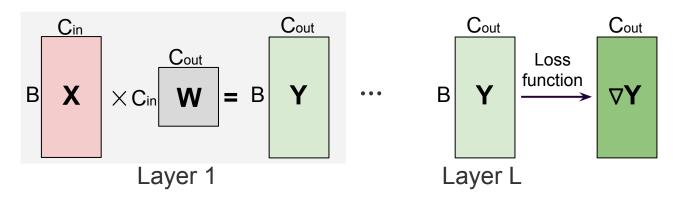


Topics

- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning
- Speculative Decoding



Forward Pass for Linear Layer



B: batch size

Cin: input channels

Y: output maps

Cout: output channels

X: input maps

W: weight filters

 The fully-connected layer during the forward propagation can be converted into matrix multiplications.



Backward Pass for Linear Layer

Weight Gradient Computation

Data Gradient Computation

$$\begin{array}{c|c}
C_{\text{out}} & C_{\text{in}} \\
\hline
C_{\text{in}} & B
\end{array}$$

X: input maps **∇X**: input gradient **W**: weight filters

Y: output maps

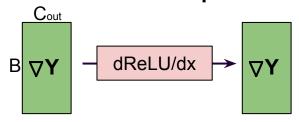
∇W: weight gradient **∀Y**: output gradient

DNN backward propagation involves two matrix multiplications



Backward Pass for Linear Layer

Data Gradient Computations



Weight Gradient Updates

$$C_{in} \boxed{\mathbf{W}} - \eta \times \nabla \mathbf{W} = \boxed{\mathbf{W'}}$$

X: input maps

∇X: input gradient

W: weight filters

∇W: weight gradient

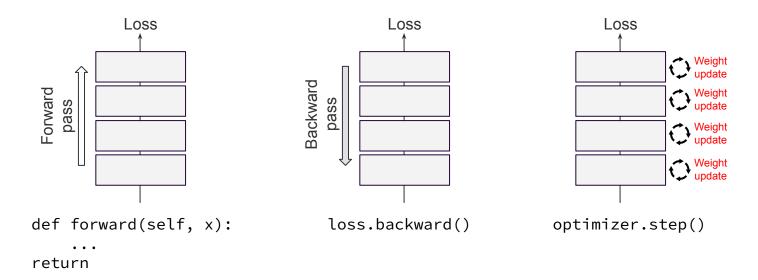
Y: output maps

∀Y: output gradient

DNN backward propagation involves two matrix multiplications

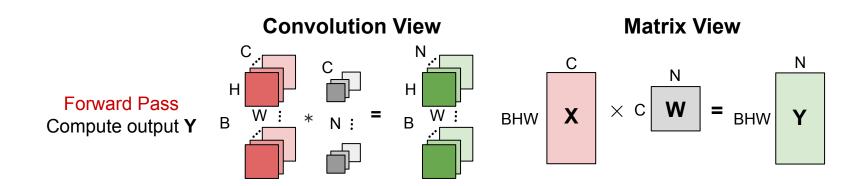


Training Process





Forward Pass for Convolutional Layer

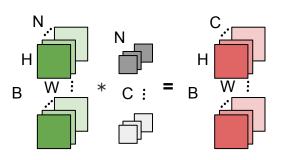


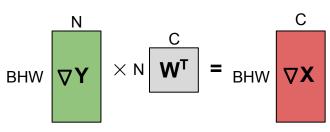
Assume a weight kernel size of 1x1.



Backward Pass for Convolutional Layer

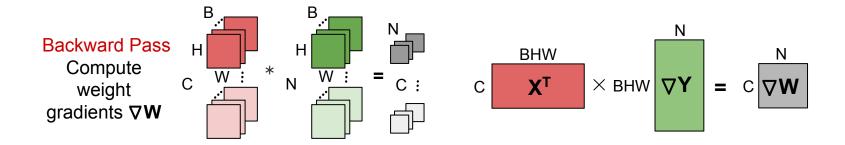
Backward Pass
Compute Activation
gradients ∇X







Backward Pass for Convolutional Layer



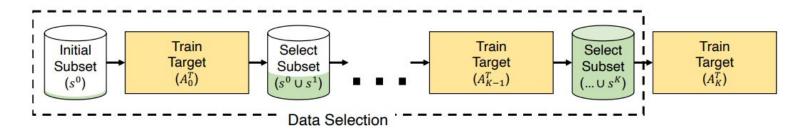


Efficient Computing during Training

- To reduce the training cost of DNN, we can proceed from the following dimensions:
 - Training data sampling
 - Parameter sampling
 - Pruning during training
 - Quantization during training



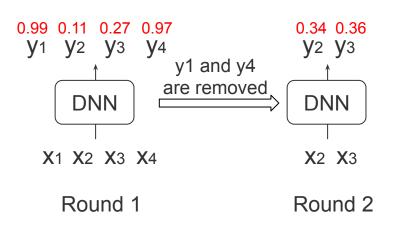
Training Data Sampling for Efficiency



- Assume a total b samples are targeted to be picked. We consider a batch setting with K rounds where we select b/K points in every round.
- Training the target model with b/K samples, then evaluate the rest of the sample over the model. Find the batch with the least confidence score. Append it to the training dataset.



Training Data Sampling for Efficiency



- Assume a total b samples are targeted to be picked. We consider a batch setting with K rounds where we select b/K points in every round.
- Training the target model with b/K samples, then evaluate the rest of the sample over the model.
 Find the batch with the least confidence score.
 Append it to the training dataset.

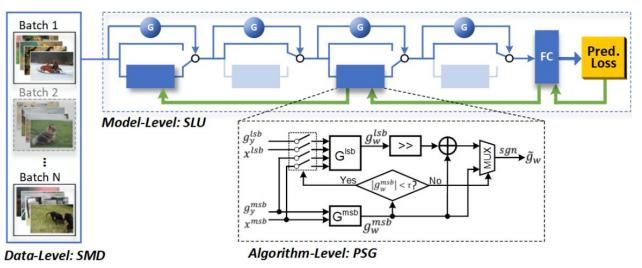


Efficient Computing during Training

- To reduce the training cost of DNN, we can proceed from the following dimensions:
 - Training data sampling
 - Parameter sampling
 - Pruning during training
 - Quantization during training



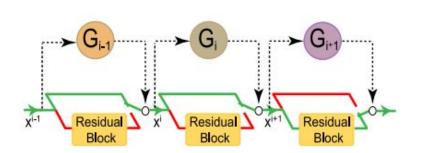
E2-Train



- A stochastic mini-batch dropping strategy is proposed.
- Stochastic minibatch dropping simply skips every mini-batch with a default probability of 0.5.
- For some easy dataset, this will generate negligible impact on performance.



Dynamically Layer Skipping



- $G_i(x_i) \in \{0,1\}$ is the gating function for layer i.
- It determines whether to skip to current residual block or not.
- During the training, G and residual blocks are trained together.
- Loss = acc_loss + computation_loss
- We will skip different layers adaptively based on inputs.



Wang, Xin, et al. "Skipnet: Learning dynamic routing in convolutional networks." *Proceedings of the European conference on computer vision (ECCV)*. 2018.

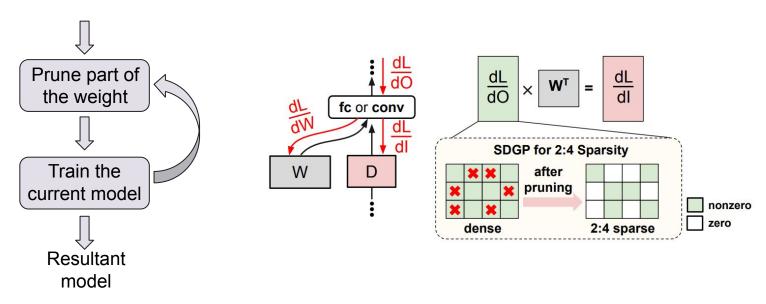
Wang, Yue, et al. "E2-train: Training state-of-the-art cnns with over 80% energy savings." *Advances in Neural Information Processing Systems* 32 (2019).

Efficient Computing during Training

- To reduce the training cost of DNN, we can proceed from the following dimensions:
 - Training data sampling
 - Parameter sampling
 - Pruning during training
 - Quantization during training



Pruning during Training



We can remove the unnecessary weight during the DNN training process.



How to Find the Winning Tickets?

Iterative Magnitude Pruning (IMP):

- Initialized DNN with random weights wo.
- While the sparsity level has not reached:
 - Train the DNN with k epochs until convergence
 - prune p% of the nonzero weights.
 - Reinitialize the remaining weights using the values in wo, finetune the remaining weights for k epochs (Rewind).
- Return the weights.
- Later work has shown that rewind to wi (i is small) works better for larger networks.

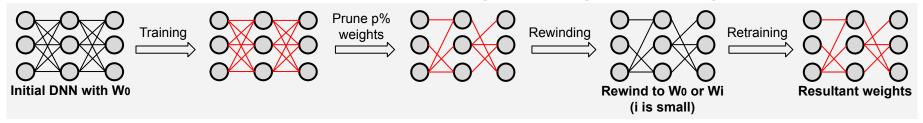


Weight Rewinding

Conventional iterative pruning



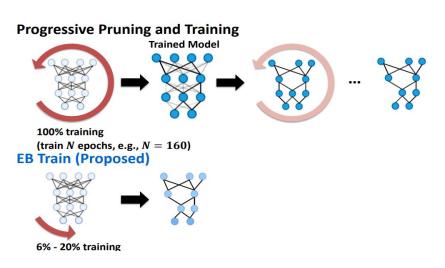
Conventional iterative pruning with weight rewinding



• The pruned architecture itself, rather than a set of inherited "important" weights, is more crucial to the accuracy in the final model, which suggests that in some cases pruning can be useful as an architecture search paradigm.



Early-bird Ticket



- LTH shows that there exist winning tickets (small but critical subnetworks) for dense, randomly initialized networks, that can be trained alone to achieve a comparable accuracy to the latter in a similar number of iterations.
- The winning tickets can be drawn very early in training and with aggressively low-cost training algorithms.
- Early-bird tickets can be founded via low-cost training schemes (e.g., early stopping and low-precision training) at large learning rates



Early-bird Ticket

Algorithm 1: The Algorithm for Searching EB Tickets

To search for the lottery ticket, we can early stop the DNN training.



Efficient Computing during Training

- To reduce the training cost of DNN, we can proceed from the following dimensions:
 - Training data sampling
 - Parameter sampling
 - Pruning during training
 - Quantization during training



DoReFaNet

```
Compute g_{a_L} = \frac{\partial C}{\partial a_L} knowing a_L and a^*.
10: for k = L to 1 do
        Back-propagate g_{a_k} through activation function h
      g_{a_k}^b \leftarrow f_{\gamma}^G(g_{a_k})
13: g_{a_k} \leftarrow \text{backward\_input}(g_{a_k}^b, W_k^b)
        g_{W_{k}^{b}} \leftarrow \text{backward\_weight}(g_{a_{k}}^{b}, a_{k-1}^{b})
15:
        Back-propagate gradients through pooling layer if there is one
16: end for
     {2. Accumulating the parameters gradients:}
17: for k = 1 to L do
       g_{W_k} = g_{W_k^b} \frac{\partial W_k^b}{\partial W_k}
      W_k^{t+1} \leftarrow Update(W_k, g_{W_k}, \eta)
20: end for
```

- Linear quantize the weights and activations
- Apply stochastic quantization for the gradients.



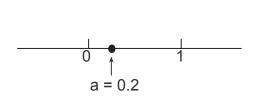
DoReFaNet

- Usually gradients requires far more bitwidth than weight and activation.
- Usually gradient requires stochastic quantization.

W	A	G	Training Complexity	Inference Complexity	Storage Relative Size	AlexNet Accuracy
1	1	6	7	1	1	0.395
1	1	8	9	1	1	0.395
1	1	32	21	1	1	0.279 (BNN)
1	1	32	-	1	1	0.442 (XNOR-Net)
1	1	32	LI.	1	1	0.401
1	1	32	£	1	1	0.436 (initialized)
1	2	6	8	2	1	0.461
1	2	8	10	2	1	0.463
1	2	32		2	1	0.477
1	2	32	5	2	1	0.498 (initialized)
1	3	6	9	3	1	0.471
1	3	32	-	3	1	0.484
1	4	6	8	4	1	0.482
1	4	32	2	4	1	0.503
1	4	32	P	4	1	0.530 (initialized)
8	8	8	FI	522	8	0.530
32	32	32	_	12	32	0.559



Deterministic and Stochastic Quantization



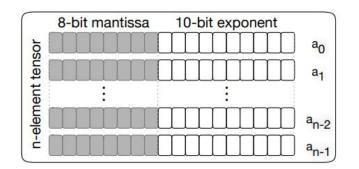
- To quantize a, conventional linear quantization will make q(a) = 0. However, this will cause a bias.
- With stochastic quantization:

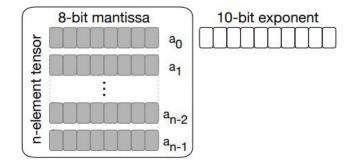
$$q(a) = egin{cases} 1 & ext{for } p = 0.2 \ 0 & ext{for } p = 0.8 \end{cases}$$

 Stochastic quantization is extremely useful when applying quantization to accelerate DNN training.



Training DNNs with Hybrid BFP





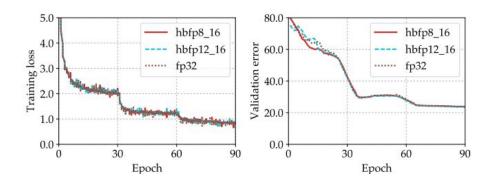
(a) FP repr. with an exponent per tensor element.

- (b) BFP repr. with an exponent per tensor.
- Block floating point format achieves a better hardware efficiency and comparable representation capability than FP.



Training DNNs with Hybrid BFP

- Use BFP in all dot-product-based operations present in DNNs (i.e., convolutions, matrix multiplications, and outer products), and floating-point representations for all other operations (i.e., activations, regularizations, etc).
- To minimize data loss in long-lasting training state, the weights are stored with wider mantissas.

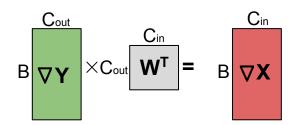


- ResNet-50 trained on ImageNet for 90 epochs.
- 8 bit mantissa, 16 bits weight seems to achieve comparable performance as FP32. A mantissa bitwidth of 12 achieves an even better performance.
- A tile size of 24.



Two Copies of Weights

Input data gradient Computation



Weight Gradient Updates

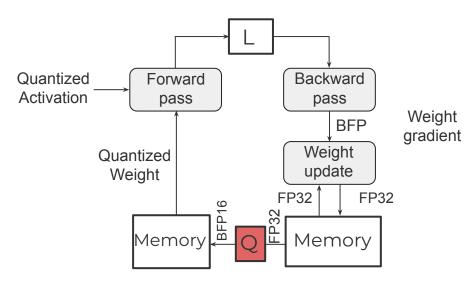
$$C_{in} \boxed{\mathbf{w}} - \eta \times \nabla \mathbf{w} = \boxed{\mathbf{w}'}$$

Weight Gradient Computation

- Gradient and forward propagation are performed using BFP.
- Weights are updated using FP.
- Two copies of weights are used.



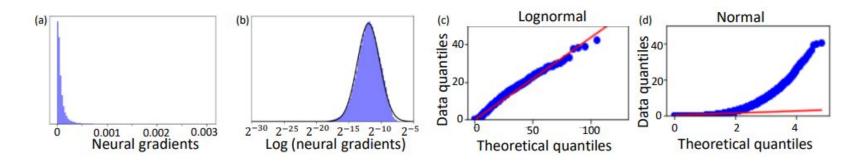
Two Copies of Weights



- Two pieces of copies are needed to be kept in the memory.
- The weight updates are usually performed with higher precision (e.g., FP16).



Neural Gradients are Near-Lognormal: Improved Quantized and Sparse Training



- The distribution of neural gradients is approximately lognormal.
- We can use lognormal regression to determine the optimal quantization setting (e.g., bitwidth, quantization interval).

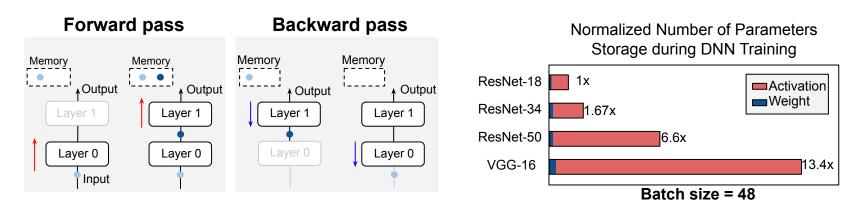


Topics

- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning
- Speculative Decoding



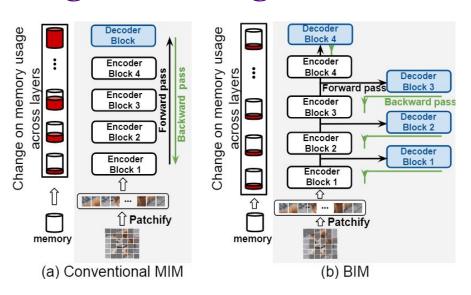
Memory Consumption During Training



- The memory footprint grows proportional with the layer depth. The activation in the early DNN layers need to be stored for a long time.
- Activations consume most of the memory space, approximately 13 times larger than the weights on average.



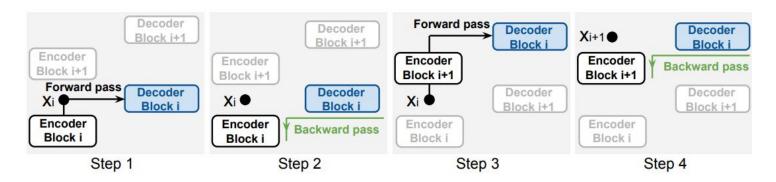
BIM: Block-Wise Local Learning with Masked Image Modeling



- Local exist is introduced during the training process.
- The intermediate results can be discarded once the training process for the current layer is complete.



BIM: Block-Wise Local Learning with Masked Image Modeling



 Once the parameter updates in encoder block i and decoder block i are finished, all intermediate features stored in the buffer, except for xi, can be cleared from memory, preserving them for future use.



Topics

- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning
- Speculative Decoding

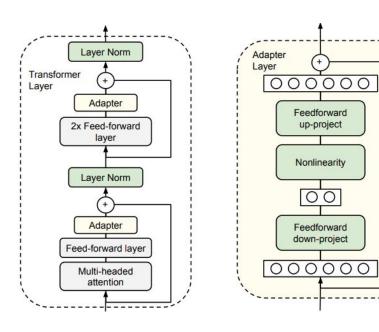


Parameter-efficient Finetuning (PEFT)

- Large models (LMs), often consisting of billions of parameters, require vast amounts of computational resources for execution.
- The expansive scale and computational demands pose considerable challenges when customizing them for particular downstream tasks.
- To better adapt the LMs over the downstream tasks, we can finetune a small portion of the LM parameters. This will make LMs achieve great performance over the downstream tasks while minimizing the training cost.
- Some of the popular PEFT Algorithms:
 - LoRA
 - Adapter
 - BitFit



Parameter-Efficient Transfer Learning for NLP



- We add the adapter module twice to each Transformer layer.
- The adapter consists of a bottleneck which contains few parameters relative to the attention and feedforward layers in the original model. The adapter also contains a skip-connection.
- The learnable parameters contributes to around 0.5 – 8% of the parameters of the original model.



BitFit

$$\begin{aligned} \mathbf{Q}^{m,\ell}(\mathbf{x}) &= \mathbf{W}_q^{m,\ell} \mathbf{x} + \mathbf{b}_q^{m,\ell} \\ \mathbf{K}^{m,\ell}(\mathbf{x}) &= \mathbf{W}_k^{m,\ell} \mathbf{x} + \mathbf{b}_k^{m,\ell} \\ \mathbf{V}^{m,\ell}(\mathbf{x}) &= \mathbf{W}_v^{m,\ell} \mathbf{x} + \mathbf{b}_v^{m,\ell} \\ \mathbf{h}_2^{\ell} &= \mathsf{Dropout} \big(\mathbf{W}_{m_1}^{\ell} \cdot \mathbf{h}_1^{\ell} + \mathbf{b}_{m_1}^{\ell} \big) \\ \mathbf{h}_3^{\ell} &= \mathbf{g}_{LN_1}^{\ell} \odot \frac{(\mathbf{h}_2^{\ell} + \mathbf{x}) - \mu}{\sigma} + \mathbf{b}_{LN_1}^{\ell} \\ \mathbf{h}_4^{\ell} &= \mathsf{GELU} \big(\mathbf{W}_{m_2}^{\ell} \cdot \mathbf{h}_3^{\ell} + \mathbf{b}_{m_2}^{\ell} \big) \\ \mathbf{h}_5^{\ell} &= \mathsf{Dropout} \big(\mathbf{W}_{m_3}^{\ell} \cdot \mathbf{h}_4^{\ell} + \mathbf{b}_{m_3}^{\ell} \big) \\ \mathsf{out}^{\ell} &= \mathbf{g}_{LN_2}^{\ell} \odot \frac{(\mathbf{h}_5^{\ell} + \mathbf{h}_3^{\ell}) - \mu}{\sigma} + \mathbf{b}_{LN_2}^{\ell} \end{aligned}$$

- BitFiT is a sparse-finetuning method where only the bias-terms of the model are being modified.
- Applying BitFit on pre-trained BERT models is competitive with (and sometimes better than) fine-tuning the entire model.
- Bias parameters make up 0.09% of the total number of parameters in BER.



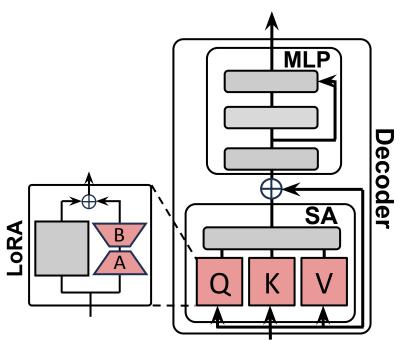
Finetune Bias is Cheap

$$B = \mathbf{X} \times C_{in} \mathbf{W} = B + \mathbf{\beta}$$

$$rac{dL}{deta} = \sum_{y \in Y} rac{dL}{dy}$$

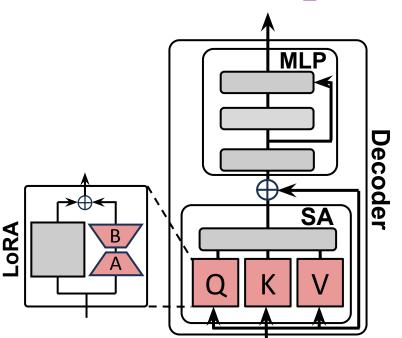
 Updating the bias does not require buffering any intermediate results during the forward pass of DNN training.





- LoRA (Low-Rank Adaptation) is a PEFT method for large pre-trained models. Instead of updating all model weights during fine-tuning, LoRA inserts small trainable low-rank matrices into specific layers (usually the attention projections).
- This dramatically reduces memory and compute requirements while maintaining near full fine-tuning performance.

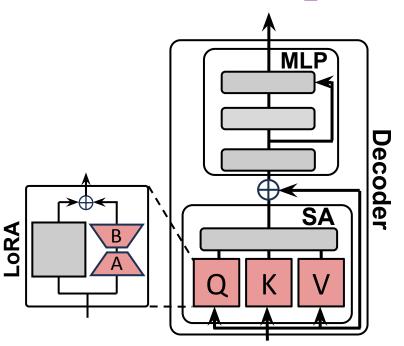




$$h = W_0 x + \Delta W x = W_0 x + BAx$$

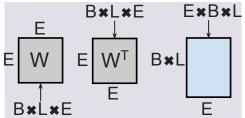
- Only the weights within the red blocks are updated.
- Assume the weight matrix has a dimension of E×E, A and B have a size of E×r and r×E, where r << k (e.g., r=4).
- BA can be merged with the original weight W₀, leading to no additional computational and storage cost.

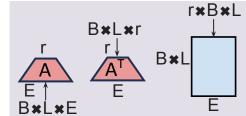




$$h = W_0 x + \Delta W x = W_0 x + BAx$$

- Compared with finetuning the entire WQ, WK and W√, this will lead to great compute savings:
 - o 3BLE²
 - o 6BLrE







Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB _{base} (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
RoB _{base} (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
$RoB_{base} (Adpt^{D})^*$	0.3M	$87.1_{\pm .0}$	$94.2_{\pm.1}$	$88.5_{\pm 1.1}$	$60.8_{\pm .4}$	$93.1_{\pm .1}$	$90.2_{\pm .0}$	$71.5_{\pm 2.7}$	$89.7_{\pm .3}$	84.4
$RoB_{base} (Adpt^{D})^*$	0.9M	$87.3_{\pm .1}$	$94.7_{\pm .3}$	$88.4_{\pm .1}$	$62.6_{\pm.9}$	$93.0_{\pm .2}$	$90.6_{\pm .0}$	$75.9_{\pm 2.2}$	$90.3_{\pm .1}$	85.4
RoB _{base} (LoRA)	0.3M	$87.5_{\pm .3}$	$95.1_{\pm .2}$	$89.7_{\pm .7}$	$63.4_{\pm 1.2}$	$93.3{\scriptstyle \pm .3}$	$90.8_{\pm.1}$	$86.6_{\pm .7}$	$\textbf{91.5}_{\pm .2}$	87.2
RoB _{large} (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB _{large} (LoRA)	0.8M	$\textbf{90.6}_{\pm.2}$	$96.2_{\pm.5}$	$90.9_{\pm 1.2}$	$\textbf{68.2}_{\pm 1.9}$	$\textbf{94.9}_{\pm.3}$	$91.6_{\pm .1}$	87.4 $_{\pm 2.5}$	$\textbf{92.6}_{\pm .2}$	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2±.3	96.1±.3	90.2±.7	68.3 ±1.0	94.8±.2	91.9 ±.1	83.8 _{±2.9}	92.1±.7	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5±.3	96.6±.2	$89.7_{\pm 1.2}$	$67.8_{\pm 2.5}$	94.8±.3	$91.7_{\pm .2}$	$80.1_{\pm 2.9}$	$91.9_{\pm .4}$	87.9
RoB _{large} (Adpt ^H)†	6.0M	$89.9_{\pm .5}$	96.2±.3	$88.7_{\pm 2.9}$	$66.5_{\pm 4.4}$	$94.7_{\pm .2}$	$92.1_{\pm .1}$	$83.4_{\pm 1.1}$	$91.0_{\pm 1.7}$	87.8
RoB _{large} (Adpt ^H)†	0.8M	$90.3_{\pm .3}$	$96.3_{\pm .5}$	$87.7_{\pm 1.7}$	$66.3_{\pm 2.0}$	$94.7_{\pm .2}$	$91.5_{\pm .1}$	$72.9_{\pm 2.9}$	$91.5_{\pm .5}$	86.4
RoB _{large} (LoRA)†	0.8M	90.6 _{±.2}	$96.2_{\pm .5}$	$90.2_{\pm 1.0}$	$68.2_{\pm1.9}$	$94.8_{\pm .3}$	$91.6_{\pm .2}$	$85.2_{\pm 1.1}$	$\textbf{92.3}_{\pm.5}$	88.6
DeB _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB _{XXL} (LoRA)	4.7M	$91.9_{\pm .2}$	$96.9_{\pm .2}$	92.6 _{±.6}	72.4 $_{\pm 1.1}$	$96.0_{\pm.1}$	92.9 $_{\pm .1}$	94.9 _{±.4}	$\textbf{93.0}_{\pm.2}$	91.3

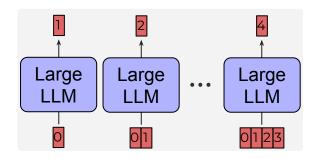


LoRA achieves better results than Adapter and BitFit.

Topics

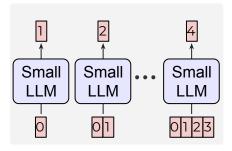
- Efficient training of DNNs
 - Efficient computing
 - Efficient storage
- Parameter efficient finetuning
- Speculative Decoding





Accurate but slow

$$T_{tot} = NT_{p,1}$$

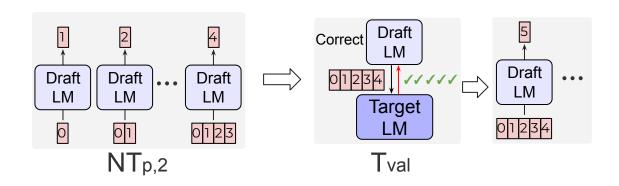


Fast but inaccurate

$$T_{tot} = NT_{p,2}$$

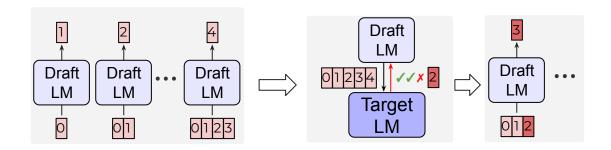
Speculative decoding enables lossless token generation with low latency.





$$T_{tot} = NT_{p,2} + T_{val} < NT_{p,1}$$

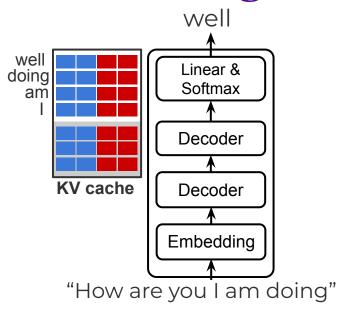


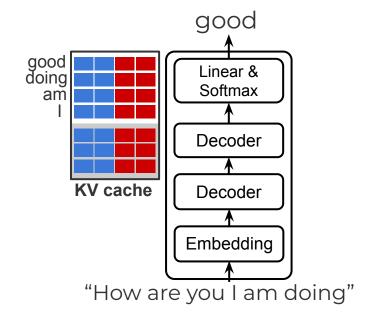


- If the token is incorrect, the target model provides the correct token to the draft model to help it generate subsequent tokens more accurately.
- If the amount of tokens that pass the verification is too low, then it is possible that speculative decoding is slower than autoregressive baseline.



LLM Decoding





• We can simply select the token with the highest score. But better results are achieved if the model considers other words as well. So a better strategy is to sample a word from the entire list using the score as the probability of selecting that word.



- To increase the diversity of the LLM output, a better strategy is to sample a word from the entire list using the score as the probability of selecting that word.
- Let p(x), q(x) denote the probability density function specified by the target and draft LLM
- To sample x ~ p(x), we instead sample x ~ q(x), keeping it if q(x) ≤ p(x), and in case q(x) > p(x) we reject the sample with probability 1- p(x)/q(x) and sample x again from an adjusted distribution p'(x) = norm(max(0, p(x) q(x))) instead.



 Speculative decoding does not save computation, but greatly reduce the memory traffic by reducing the number of memory reads, further reducing the overall latency.

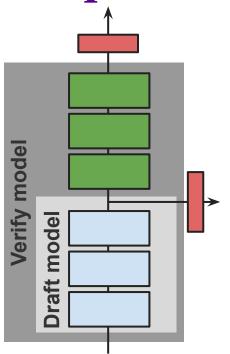
Algorithm 1 SpeculativeDecodingStep

return $prefix + [x_1, \ldots, x_n, t]$

```
Inputs: M_p, M_q, prefix.
\triangleright Sample \gamma guesses x_{1,...,\gamma} from M_q autoregressively.
for i=1 to \gamma do
   q_i(x) \leftarrow M_q(prefix + [x_1, \dots, x_{i-1}])
   x_i \sim q_i(x)
end for
\triangleright Run M_p in parallel.
p_1(x), \ldots, p_{\gamma+1}(x) \leftarrow
       M_p(prefix), \ldots, M_p(prefix + [x_1, \ldots, x_{\gamma}])
\triangleright Determine the number of accepted guesses n.
r_1 \sim U(0,1), \ldots, r_{\gamma} \sim U(0,1)
n \leftarrow \min(\{i-1 \mid 1 \le i \le \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})
\triangleright Adjust the distribution from M_p if needed.
p'(x) \leftarrow p_{n+1}(x)
if n < \gamma then
   p'(x) \leftarrow norm(max(0, p_{n+1}(x) - q_{n+1}(x)))
end if
\triangleright Return one token from M_p, and n tokens from M_q.
t \sim p'(x)
```



Self-Speculative Decoding



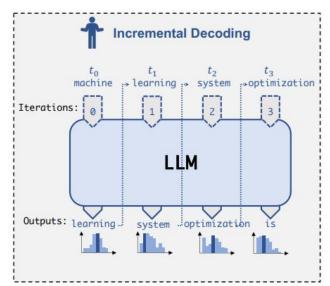
- Self-Speculative decoding the draft model is a subnetwork of the verify model. All the intermediate results from the draft model are reusable.
- No additional network needs to be trained, except a simple classification layer.

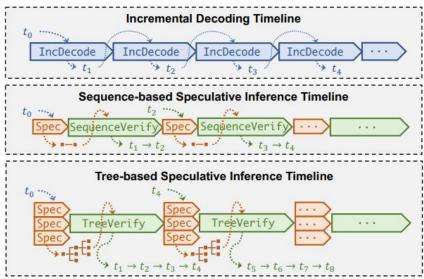


Zhang, Jun, et al. "Draft & verify: Lossless large language model acceleration via self-speculative decoding." arXiv preprint arXiv:2309.08168 (2023).

Elhoushi, Mostafa, et al. "Layer skip: Enabling early exit inference and self-speculative decoding." *arXiv preprint arXiv:2404.16710* (2024).

SpecInfer





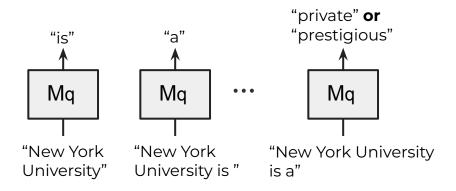
(a) Incremental decoding.

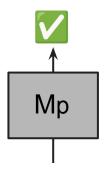
(b) Timeline Comparison.



Miao, Xupeng, et al. "SpecInfer: Accelerating Generative Large Language Model Serving with Tree-based Speculative Inference and Verification." arXiv preprint arXiv:2305.09781 (2023).

SpecInfer



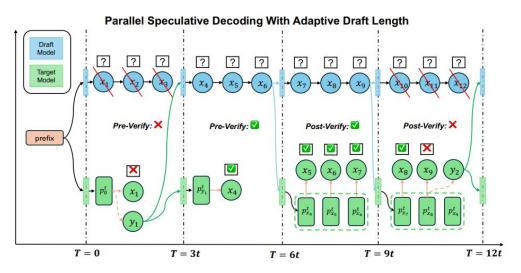


"New York University is a private research university"

or

"New York University is a prestigious research university"

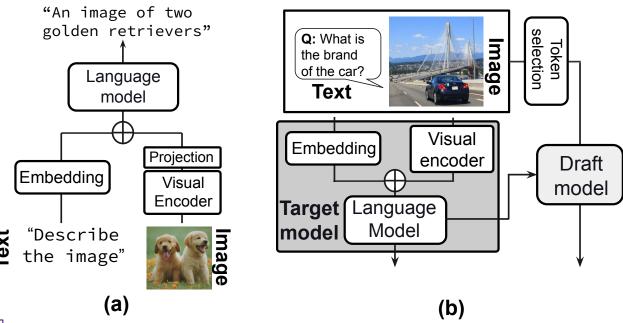
Parallel Speculative Decoding



- PEARL is a parallel inference framework based on speculative decoding which utilizes pre-verify and post-verify to achieve adaptive draft length.
- The draft model continues to decode during the verification stage.
- If the verification fails, the windows size will become 1 in the next cycle.



DREAM

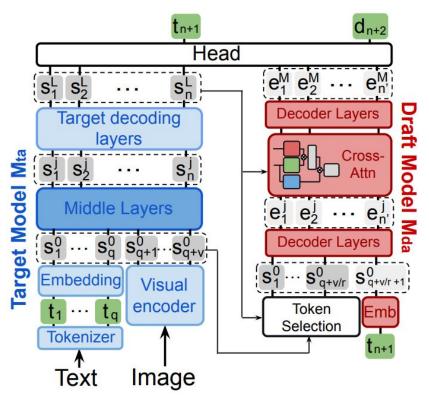




Hu, Yunhai, et al. "DREAM: Drafting with Refined Target Features and Entropy-Adaptive Cross-Attention Fusion for Multimodal Speculative Decoding." *arXiv preprint arXiv:2505.19201* (2025).

DREAM

- During operation, the target model will send their intermediate results to the draft model to better guide the generation of the draft model.
- The visual tokens will also be pruned to remove the redundant tokens to reduce the processing latency of the draft model.





DREAM

		MMT		SEED		ScienceQA		OCRBench		ChartQA		MathVista		Average	
Models	Methods	S	au	S	au	S	au	S	au	S	au	S	au	S	au
					Ten	nperatui	re = 0								
LLaVA-v1.6 Vicuna-7B	SPD [10]	1.10	1.88	0.81	1.17	1.08	1.87	0.89	1.25	0.91	1.24	1.06	1.76	0.97	1.53
	Kangaroo [28]	1.32	2.11	1.33	2.12	1.31	2.09	1.17	1.89	1.18	1.98	1.15	1.86	1.24	2.01
	Medusa [4]	1.58	2.88	1.59	3.01	1.44	2.77	1.22	2.33	1.25	2.41	1.22	2.34	1.38	2.62
	Hydra [2]	1.78	3.86	1.72	3.88	1.68	3.79	1.41	3.21	1.35	3.11	1.42	3.25	1.56	3.52
	EAGLE [25]	2.10	5.04	2.09	5.01	1.98	4.88	1.72	4.13	1.56	3.98	1.78	4.25	1.87	4.55
	EAGLE-2 [24]	2.31	5.48	2.31	5.61	2.15	5.22	1.92	4.88	1.77	4.22	1.87	4.67	2.05	5.01
	DREAM	2.52	6.40	2.48	6.20	2.33	5.82	2.05	4.88	1.89	4.44	2.11	5.32	2.23	5.51
LLaVA-v1.6 Vicuna-13B	SPD	1.07	1.78	1.06	1.79	1.09	1.88	0.86	1.12	0.89	1.25	0.87	1.22	1.00	1.58
	Kangaroo	1.43	1.77	1.51	1.87	1.22	1.55	1.21	1.54	1.27	1.61	1.53	2.01	1.36	1.72
	Medusa	1.99	2.67	1.96	2.76	1.93	2.77	1.40	2.92	1.51	2.82	1.51	2.62	1.72	2.76
	Hydra	2.12	2.87	2.08	2.99	2.21	3.12	1.49	3.07	1.65	3.03	1.66	2.87	1.87	2.99
	EAGLE	2.45	3.56	2.19	3.24	2.63	3.98	1.65	3.31	1.85	3.27	1.8	3.09	2.10	3.41
	EAGLE-2	2.89	4.05	3.18	4.33	3.09	4.97	2.20	4.12	2.41	4.15	2.39	3.76	2.69	4.23
	DREAM	3.68	5.58	3.51	5.34	3.36	5.29	2.69	4.64	2.59	4.20	2.53	4.18	3.06	4.87
Pixtral-12B	SPD	1.08	1.51	1.03	1.47	1.05	1.49	1.05	1.49	1.04	1.43	1.04	1.46	1.05	1.47
	Kangaroo	1.26	1.54	1.09	1.39	1.14	1.51	1.16	1.52	1.12	1.47	1.13	1.49	1.15	1.49
	Medusa	1.37	1.81	1.37	1.81	1.35	1.87	1.24	1.69	1.22	1.68	1.16	1.47	1.28	1.72
	Hydra	1.58	2.24	1.47	2.04	1.53	2.06	1.38	1.81	1.34	1.79	1.36	1.78	1.44	1.95
	EAGLE	2.38	3.47	1.97	2.53	2.31	3.64	1.69	2.73	1.78	2.84	1.64	2.47	1.96	2.95
	EAGLE-2	2.81	3.95	2.31	3.07	2.64	4.03	2.12	3.25	2.14	3.17	1.81	2.73	2.31	3.37
	DREAM	2.93	4.52	2.61	3.67	2.98	4.33	2.38	3.55	2.35	3.49	2.36	3.42	2.65	3.78

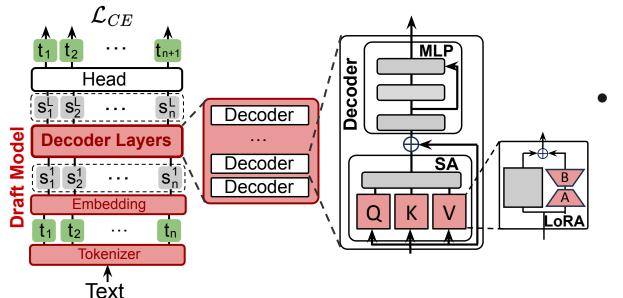


What makes an ideal draft model?

- Ideally, the draft model should have:
 - High acceptance rate
 - Low execution latency
- This is exactly the goal of DNN pruning, quantization, knowledge distillation, dynamic computing...



Speculative Decoding with Finetuning



The draft model will be trained using the dataset with cross-entropy loss to achieve a better acceptance ratio.

