

[논문리뷰]

"LoRA: Low-Rank Adaptation of Large Language Models"

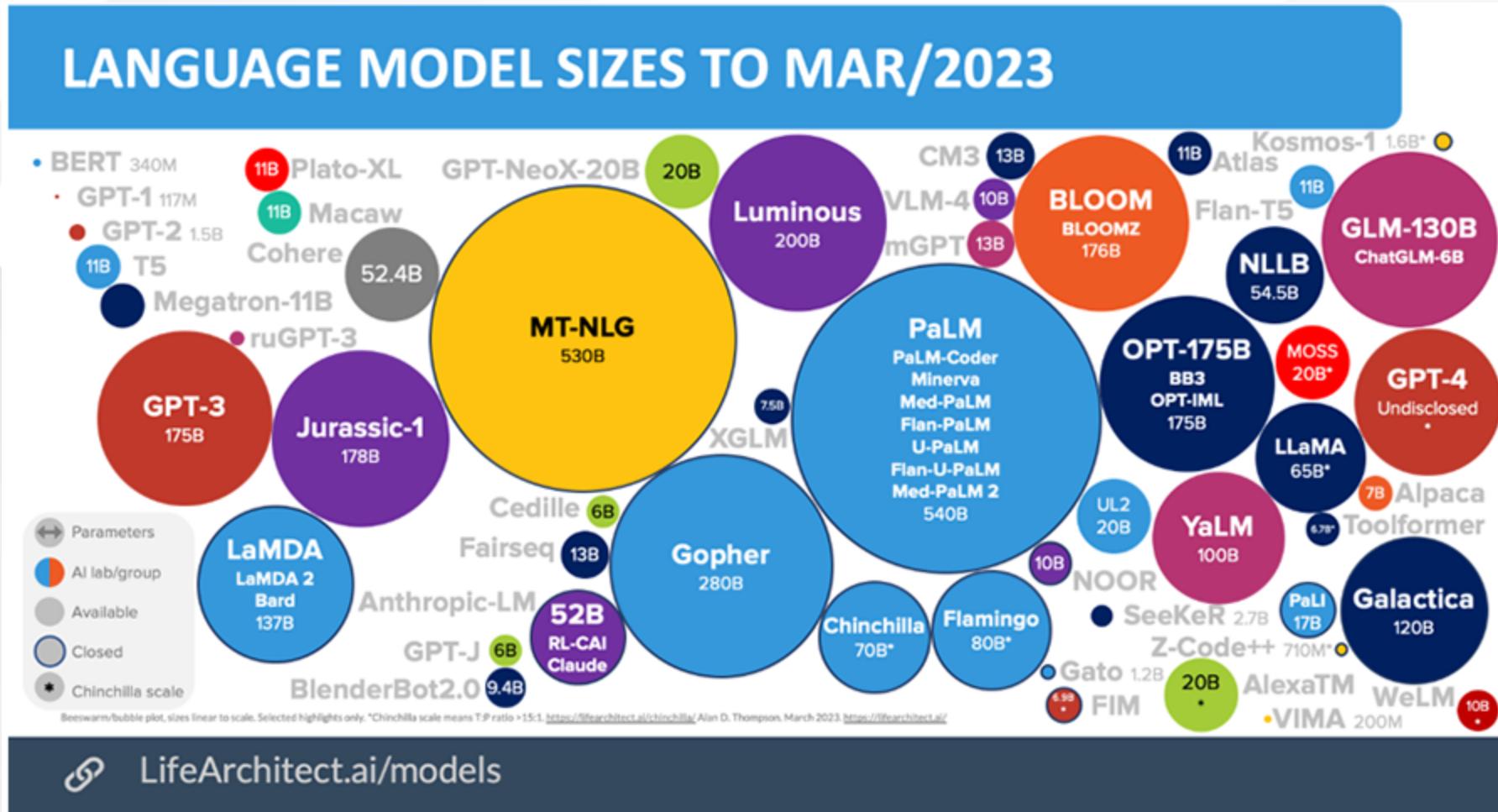
LoRA

-Microsoft Corporation-

2023.10.12

발표자 신중현

situation



motivation

pretrained model - 다양한 downstream task adaptation

규모를 키운 pretrained 모델의 성능이 경험적으로, 성능적으로 확인된 가운데, 이 모델들을 활용해서 downstream task에 adaptation해서 여러문제를 풀겠다.

-ex) summarization, reading comprehension,

finetuning

이러한 down stream task들은 training data of context-target pairs 로 표현

$$Z = (x_i, y_i)_{i=1, \dots, N}$$

fully finetuning

많은 language model의 objective function을 보면,
Maximum likelihood estimation MLE로 푼다.

$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi} (y_t \mid x, y_{<t}))$$

이런 문제를 Backward propagation을 통해서 parameters를 업데이트를 하면서 풁니다.
근데 이럴때, optimizer가 사용되는데, 어떻게 gradient 값들이 지정될까요?

SGD

SGD는 각 배치(데이터 샘플의 작은 그룹)마다 모델 가중치를 업데이트합니다. 가중치(W)를 업데이트하기 위한 수식은 다음과 같습니다:

$$W_{t+1} = W_t - \eta \nabla L(W_t)$$

여기서, W_t 는 현재 가중치, W_{t+1} 는 업데이트된 가중치,

- η 는 학습률(learning rate).
- $\nabla L(W_t)$ 는 비용 함수 L 에 대한 그래디언트(gradients)입니다.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla f(x_{t-1})$$

$$g_t = \beta_2 g_{t-1} + (1 - \beta_2) (\nabla f(x_{t-1}))^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \hat{g}_t = \frac{g_t}{1 - \beta_2^t}$$

$$x_t = x_{t-1} - \frac{\eta}{\sqrt{\hat{g}_t + \epsilon}} \cdot \hat{m}_t$$

- β_1 : Momentum의 지수이동평균 ≈ 0.9
- β_2 : RMSProp의 지수이동평균 ≈ 0.999
- \hat{m}, \hat{g} : 학습 초기 시 m_t, g_t 가 0이 되는 것을 방지하기 위한 보정 값
- ϵ : 분모가 0이 되는 것을 방지하기 위한 작은 값 $\approx 10^{-8}$
- η : 학습률 ≈ 0.001

optimizer	VRAM 사용량
SGD	가중치의 크기
Adam	가중치의 크기 + m + g

Inspiration

over-parameterized model이 실제로 낮은 low intrinsic dimension에 있다. 가정.

[Measuring the Intrinsic Dimension of Objective Landscapes](#)

[Intrinsic Dimensionality Explains the Effectiveness of Language Model Fine-Tuning](#)

[Aghajanyan](#)

LoRA: Low-Rank Adaptation

LoRA를 사용하면?

Pre-trained weight를 고정된 상태(freeze)로 유지하면서 Adaptation 중 dense layer의 변화에 대한 rank decomposition matrices를 최적화

이를 통해 신경망의 일부 dense layer를 간접적으로 훈련시키는 것이 가능

- LoRA는 trainable parameter의 수가 적고 학습 처리량이 높으며 inference latency가 이전 연구 대비 적음
- 그럼에도 불구하고 ROBERTa, DeBERTa, GPT-2, GPT-3에서 fine-tuning보다 같거나 더 나은 성능을 보여줌

Terminologies and Conventions

d_{model} : input/output dimension size 768, 1024

W_q, W_k, W_v, Wo : query/key/value/output projection matrices in self-attention module

W, W_o : pre-trained weight matrix

ΔW : accumulated gradient update during adaptation

r : rank of a LORA module

Transformers 논문의 setting을 따름

e.g., use Adam optimizer

= 4x dmodel



Model Weights

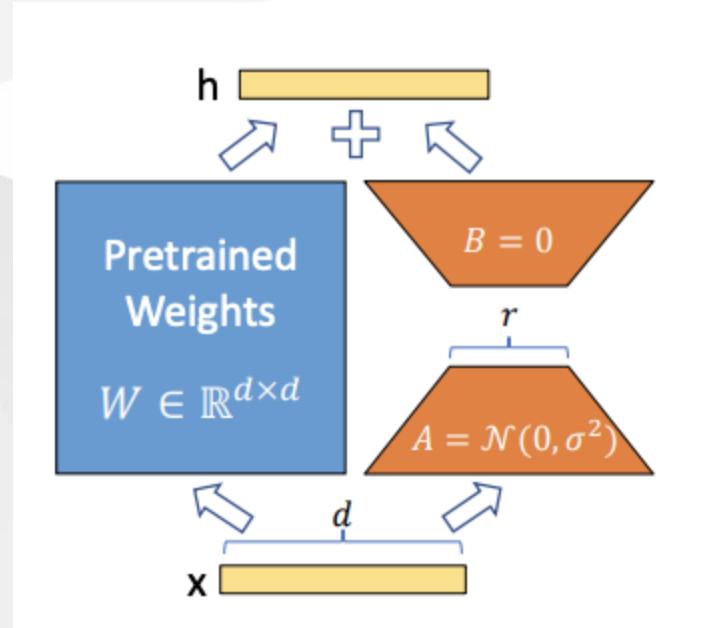
$w_{11} \ w_{12} \ w_{13}$
 $w_{21} \ w_{22} \ w_{23}$
 $w_{31} \ w_{32} \ w_{33}$

$+ \eta.$

Gradients for each Weight

$$\begin{array}{lll} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \frac{\partial L}{\partial w_{13}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} & \frac{\partial L}{\partial w_{23}} \\ \frac{\partial L}{\partial w_{31}} & \frac{\partial L}{\partial w_{32}} & \frac{\partial L}{\partial w_{33}} \end{array}$$

LoRA fintuning method



$$\begin{aligned} W(\text{finetuning}) &\rightarrow W(\text{freeze}) + \nabla W \\ W &\in R^{d \times k} \\ \nabla W &= BA \\ A &\in R^{d \times r}, B \in R^{r \times k} \end{aligned}$$

Fully fintuning

$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi} (y_t \mid x, y_{<t}))$$

LoRA

$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Theta_0 + \Delta\Phi(\Theta)} (y_t \mid x, y_{<t}))$$

- $\Delta\Phi = \Delta\Phi(\Theta)$: the task-specific parameter increment
- 훨씬 작은 set of parameters Θ s.t $\Theta \ll |\Phi_0|$

fintunig으로 점점 커지는 LLM을 감당하기 힘들기에, 그걸 효율적으로 파인튜닝 시키는 방법.

what is Rank?

- numbers how many linearly independent columns in matrix. = rank of columns = rank of rows
- rank : row reduced echelon form의 1의 개수

rref

elementary matrix를 곱해서 만들수있음.

Goal: Put matrix in REF or RREF

$$\left(\begin{array}{cccc|ccccc} * & * & \dots & \dots & & & & & \\ 0 & * & \dots & \dots & \dots & & & & \\ 0 & 0 & * & \dots & \dots & \dots & & & \\ 0 & 0 & 0 & * & \dots & \dots & & & \\ 0 & 0 & 0 & 0 & \dots & 0 & & & \end{array} \right)$$

REF

$$\left(\begin{array}{cccc|ccccc} 1 & 0 & 0 & 0 & * & \dots & \dots & & \\ 0 & 1 & 0 & 0 & * & \dots & \dots & & \\ 0 & 0 & 1 & 0 & * & \dots & \dots & & \\ 0 & 0 & 0 & 1 & * & \dots & \dots & & \\ 0 & 0 & 0 & 0 & 0 & \dots & \dots & 0 & \end{array} \right)$$

RREF

$$A = \begin{pmatrix} 2 & -2 & 9 \\ 2 & 1 & 10 & 7 \\ -4 & 4 & -8 & 4 \\ 4 & -1 & 14 & 6 \end{pmatrix}$$

RREF Z.

$$\bar{E}_1 \cdot \bar{E}_2 \bar{E}_1 A = \begin{pmatrix} 1 & 0 & 4 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$A \sim_r B$$

$$A \sim_r C$$

$$B = C$$

rref의 유일성과 존재성 증명됨.

LoRA

$$\underline{W_0} + \underline{\Delta W} = W_0 + BA, \quad B \in \mathbb{R}^{d \times r}$$

freeze pretraind train LoRA params

if. $A \in \mathbb{R}^{r \times k}$ $r \leq \min(d, k)$

$$A \approx_r r \begin{pmatrix} 1 & & & & & \\ & 1 & & & & \\ & & \ddots & & & \\ & & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{pmatrix}^T$$

$$\max(rk(A)) = r$$

only adapting the attention weights

W_q, W_k, W_v, W_o

freeze MLP modules

Alternatives

1. adding Adapter Layer in each transformer block
2. Prefix tuning

reference

youtube channels : 딥러닝논문읽기모임

<https://www.youtube.com/watch?v=BJqwmDpa0wM>

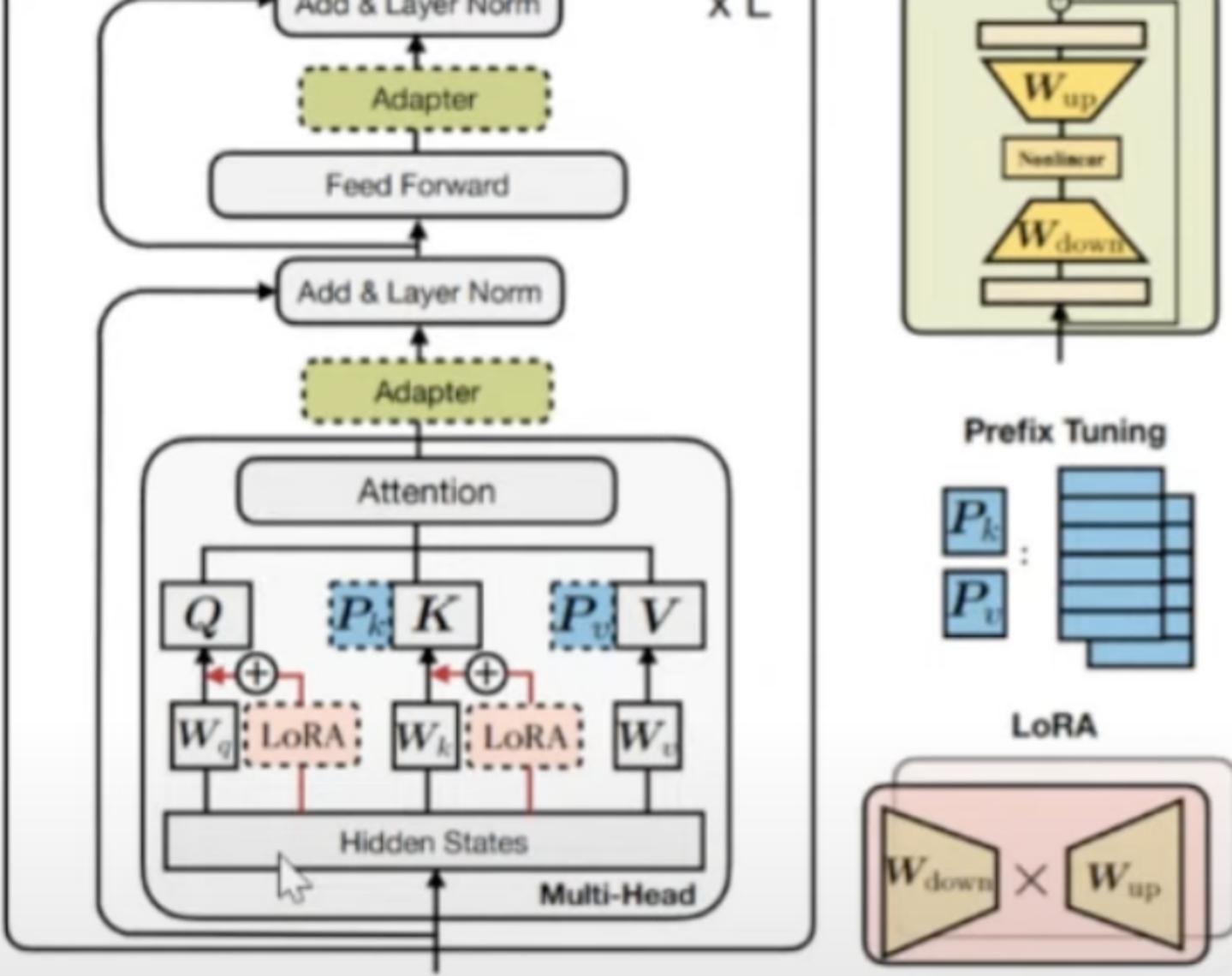
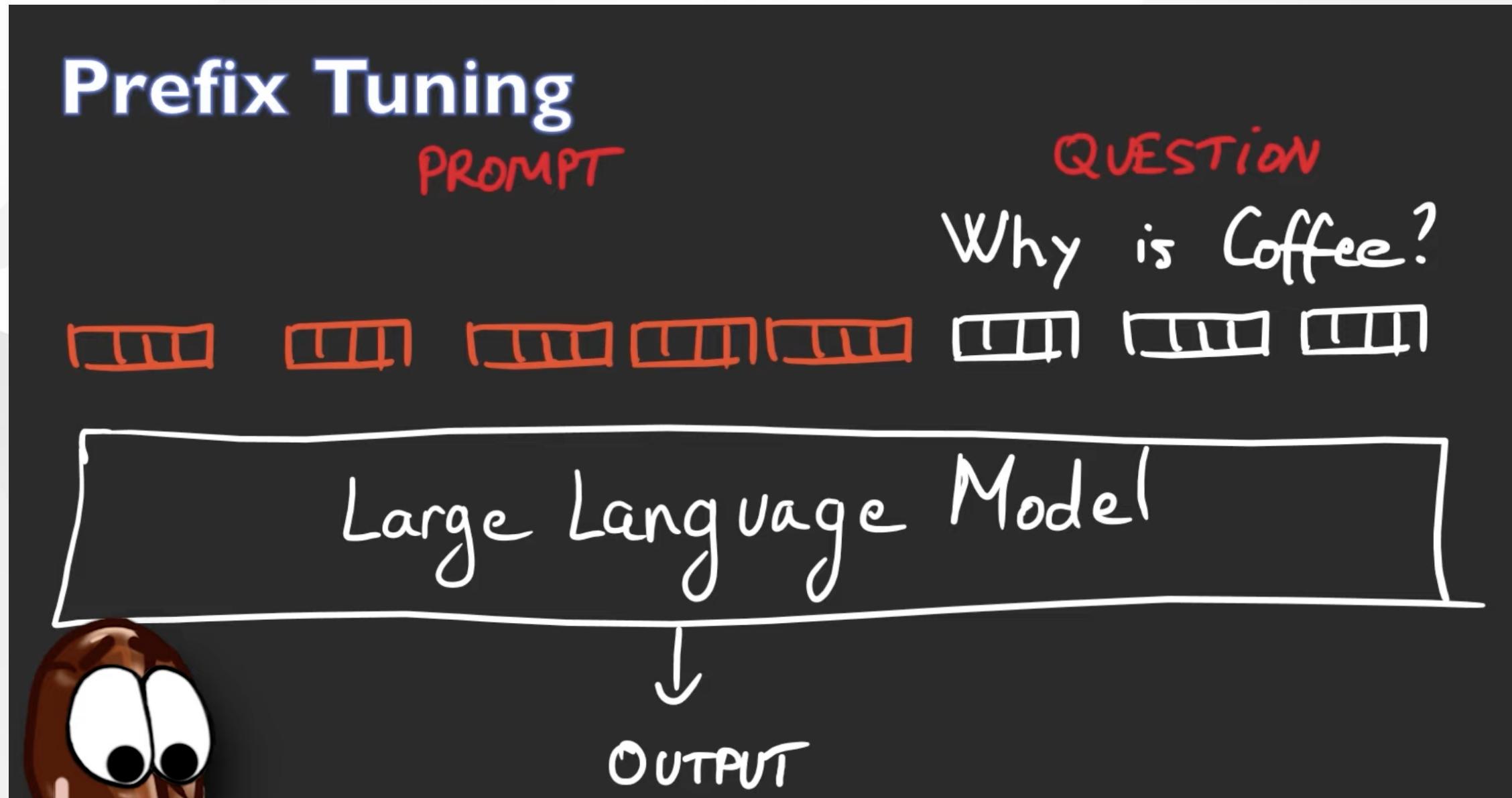


Figure 1: Illustration of the transformer architecture and several state-of-the-art parameter-efficient tuning methods. We use blocks with dashed borderlines to

1. Prefix tuning



LoRA

원래 weight에 BA 를 통해서 weight_0에서 variants를 가하겠다.

CODE

<https://github.com/microsoft/LoRA>

A,B Definition

```
class Linear(nn.Linear, LoRALayer):
    # LoRA implemented in a dense layer

        # Actual trainable parameters

#####
#####      Definition !!!!!      #####
    if r > 0:
        self.lora_A = nn.Parameter(self.weight.new_zeros((r, in_features)))
        self.lora_B = nn.Parameter(self.weight.new_zeros((out_features, r)))
        self.scaling = self.lora_alpha / self.r
        # Freezing the pre-trained weight matrix
        self.weight.requires_grad = False
    self.reset_parameters()
    if fan_in_fan_out:
        self.weight.data = self.weight.data.transpose(0, 1)
```

A,B initialize

```
def reset_parameters(self):
    self.conv.reset_parameters()
    if hasattr(self, 'lora_A'):
        # initialize A the same way as the default for nn.Linear and B to zero
        nn.init.kaiming_uniform_(self.lora_A, a=math.sqrt(5))
        nn.init.zeros_(self.lora_B)
```

```
# LoRA implemented in a dense layer
def reset_parameters(self):
    nn.Embedding.reset_parameters(self)
    if hasattr(self, 'lora_A'):
        # initialize A the same way as the default for nn.Linear and B to zero
        nn.init.zeros_(self.lora_A)
        nn.init.normal_(self.lora_B)
```


elementwise sum

```
def train(self, mode: bool = True):
    def T(w):
        return w.transpose(0, 1) if self.fan_in_fan_out else w
    nn.Linear.train(self, mode)
    if mode:
        if self.merge_weights and self.merged:
            # Make sure that the weights are not merged
            if self.r > 0:
                self.weight.data -= T(self.lora_B @ self.lora_A) * self.scaling
        self.merged = False
```

Experiments

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 _{±.0}	94.2 _{±.1}	88.5 _{±1.1}	60.8 _{±.4}	93.1 _{±.1}	90.2 _{±.0}	71.5 _{±2.7}	89.7 _{±.3}	84.4
RoB _{base} (Adpt ^D)*	0.9M	87.3 _{±.1}	94.7 _{±.3}	88.4 _{±.1}	62.6 _{±.9}	93.0 _{±.2}	90.6 _{±.0}	75.9 _{±2.2}	90.3 _{±.1}	85.4
RoB _{base} (LoRA)	0.3M	87.5 _{±.3}	95.1 _{±.2}	89.7 _{±.7}	63.4 _{±1.2}	93.3 _{±.3}	90.8 _{±.1}	86.6 _{±.7}	91.5 _{±.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	90.6 _{±.2}	96.2 _{±.5}	90.9 _{±1.2}	68.2 _{±1.9}	94.9 _{±.3}	91.6 _{±.1}	87.4 _{±2.5}	92.6 _{±.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 _{±1.2}	67.8 _{±2.5}	94.8 _{±.3}	91.7 _{±.2}	80.1 _{±2.9}	91.9 _{±.4}	87.9
RoB _{large} (Adpt ^H)†	6.0M	89.9 _{±.5}	96.2 _{±.3}	88.7 _{±2.9}	66.5 _{±4.4}	94.7 _{±.2}	92.1 _{±.1}	83.4 _{±1.1}	91.0 _{±1.7}	87.8
RoB _{large} (Adpt ^H)†	0.8M	90.3 _{±.3}	96.3 _{±.5}	87.7 _{±1.7}	66.3 _{±2.0}	94.7 _{±.2}	91.5 _{±.1}	72.9 _{±2.9}	91.5 _{±.5}	86.4
RoB _{large} (LoRA)†	0.8M	90.6 _{±.2}	96.2 _{±.5}	90.2 _{±1.0}	68.2 _{±1.9}	94.8 _{±.3}	91.6 _{±.2}	85.2 _{±1.1}	92.3 _{±.5}	88.6
DeBERT _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeBERT _{XXL} (LoRA)	4.7M	91.9 _{±.2}	96.9 _{±.2}	92.6 _{±.6}	72.4 _{±1.1}	96.0 _{±.1}	92.9 _{±.1}	94.9 _{±.4}	93.0 _{±.2}	91.3

Table 2: RoBERTa_{base}, RoBERTa_{large}, and DeBERTa_{XXL} with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all metrics. * indicates numbers published in prior works. † indicates runs configured in a setup similar to Houlsby et al. (2019) for a fair comparison.

Model & Method	# Trainable Parameters	E2E NLG Challenge				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 M (FT)*	354.92M	68.2	8.62	46.2	71.0	2.47
GPT-2 M (Adapter ^L)*	0.37M	66.3	8.41	45.0	69.8	2.40
GPT-2 M (Adapter ^L)*	11.09M	68.9	8.71	46.1	71.3	2.47
GPT-2 M (Adapter ^H)	11.09M	67.3 _{±.6}	8.50 _{±.07}	46.0 _{±.2}	70.7 _{±.2}	2.44 _{±.01}
GPT-2 M (FT ^{Top2})*	25.19M	68.1	8.59	46.0	70.8	2.41
GPT-2 M (PreLayer)*	0.35M	69.7	8.81	46.1	71.4	2.49
GPT-2 M (LoRA)	0.35M	70.4_{±.1}	8.85_{±.02}	46.8_{±.2}	71.8_{±.1}	2.53_{±.02}
GPT-2 L (FT)*	774.03M	68.5	8.78	46.0	69.9	2.45
GPT-2 L (Adapter ^L)	0.88M	69.1 _{±.1}	8.68 _{±.03}	46.3 _{±.0}	71.4 _{±.2}	2.49_{±.0}
GPT-2 L (Adapter ^L)	23.00M	68.9 _{±.3}	8.70 _{±.04}	46.1 _{±.1}	71.3 _{±.2}	2.45 _{±.02}
GPT-2 L (PreLayer)*	0.77M	70.3	8.85	46.2	71.7	2.47
GPT-2 L (LoRA)	0.77M	70.4_{±.1}	8.89_{±.02}	46.8_{±.2}	72.0_{±.2}	2.47 _{±.02}

Table 3: GPT-2 medium (M) and large (L) with different adaptation methods on the E2E NLG Challenge. For all metrics, higher is better. LoRA outperforms several baselines with comparable or fewer trainable parameters. Confidence intervals are shown for experiments we ran. * indicates numbers published in prior works.

Model&Method	# Trainable Parameters	WikiSQL	MNLI-m	SAMSum
		Acc. (%)	Acc. (%)	R1/R2/RL
GPT-3 (FT)	175,255.8M	73.8	89.5	52.0/28.0/44.5
GPT-3 (BitFit)	14.2M	71.3	91.0	51.3/27.4/43.5
GPT-3 (PreEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5
GPT-3 (PreLayer)	20.2M	70.1	89.5	50.8/27.3/43.5
GPT-3 (Adapter ^H)	7.1M	71.9	89.8	53.0/28.9/44.8
GPT-3 (Adapter ^H)	40.1M	73.2	91.5	53.2/29.0/45.1
GPT-3 (LoRA)	4.7M	73.4	91.7	53.8/29.8/45.9
GPT-3 (LoRA)	37.7M	74.0	91.6	53.4/29.2/45.1

Table 4: Performance of different adaptation methods on GPT-3 175B. We report the logical form validation accuracy on WikiSQL, validation accuracy on MultiNLI-matched, and Rouge-1/2/L on SAMSum. LoRA performs better than prior approaches, including full fine-tuning. The results on WikiSQL have a fluctuation around $\pm 0.5\%$, MNLI-m around $\pm 0.1\%$, and SAMSum around $\pm 0.2/\pm 0.2/\pm 0.1$ for the three metrics.

		# of Trainable Parameters = 18M						
Weight Type		W_q	W_k	W_v	W_o	W_q, W_k	W_q, W_v	W_q, W_k, W_v, W_o
Rank r		8	8	8	8	4	4	2
WikiSQL ($\pm 0.5\%$)		70.4	70.0	73.0	73.2	71.4	73.7	73.7
MultiNLI ($\pm 0.1\%$)		91.0	90.8	91.0	91.3	91.3	91.3	91.7

Table 5: Validation accuracy on WikiSQL and MultiNLI after applying LoRA to different types of attention weights in GPT-3, given the same number of trainable parameters. Adapting both W_q and W_v gives the best performance overall. We find the standard deviation across random seeds to be consistent for a given dataset, which we report in the first column.

	Weight Type	$r = 1$	$r = 2$	$r = 4$	$r = 8$	$r = 64$
WikiSQL($\pm 0.5\%$)	W_q	68.8	69.6	70.5	70.4	70.0
	W_q, W_v	73.4	73.3	73.7	73.8	73.5
	W_q, W_k, W_v, W_o	74.1	73.7	74.0	74.0	73.9
MultiNLI ($\pm 0.1\%$)	W_q	90.7	90.9	91.1	90.7	90.7
	W_q, W_v	91.3	91.4	91.3	91.6	91.4
	W_q, W_k, W_v, W_o	91.2	91.7	91.7	91.5	91.4

Table 6: Validation accuracy on WikiSQL and MultiNLI with different rank r . To our surprise, a rank as small as one suffices for adapting both W_q and W_v on these datasets while training W_q alone needs a larger r . We conduct a similar experiment on GPT-2 in Section H.2.

This suggests the update matrix ΔW could have a very small “intrinsic rank”

$$\phi(A_{r=8}, A_{r=64}, i, j) = \frac{\|U_{A_{r=8}}^{i\top} U_{A_{r=64}}^j\|_F^2}{\min(i, j)} \in [0, 1] \quad (4)$$

where $U_{A_{r=8}}^i$ represents the columns of $U_{A_{r=8}}$ corresponding to the top- i singular vectors.

$\phi(\cdot)$ has a range of $[0, 1]$, where 1 represents a complete overlap of subspaces and 0 a complete separation. See Figure 3 for how ϕ changes as we vary i and j . We only look at the 48th layer (out of 96) due to space constraint, but the conclusion holds for other layers as well, as shown in Section H.1.

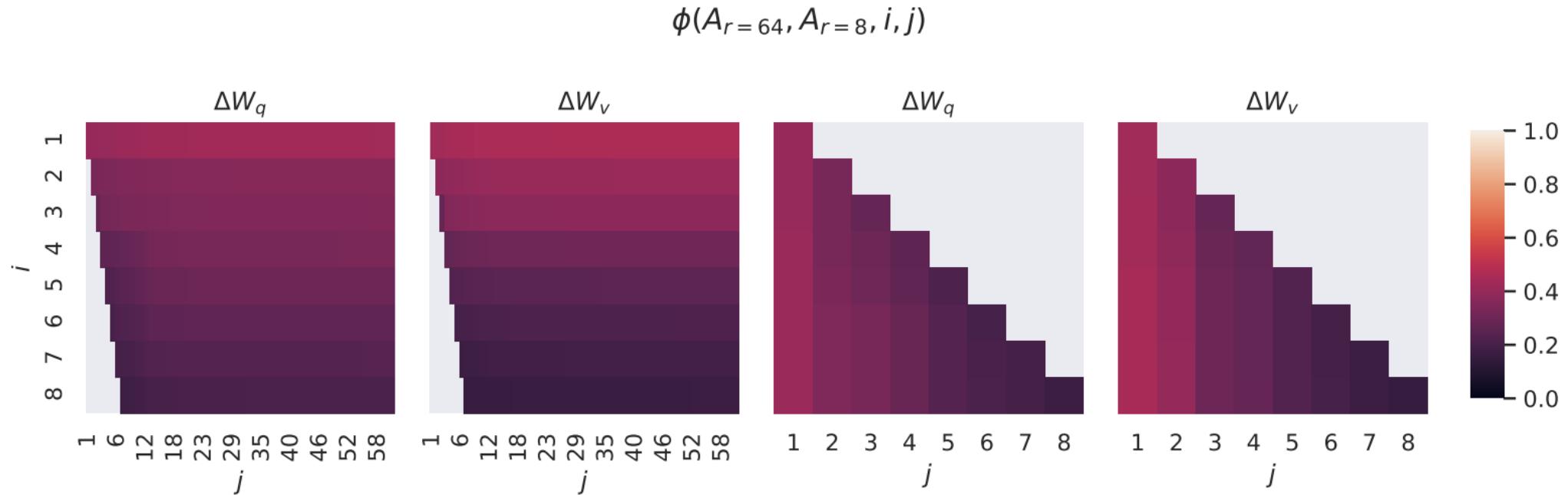


Figure 3: Subspace similarity between column vectors of $A_{r=8}$ and $A_{r=64}$ for both ΔW_q and ΔW_v . The third and the fourth figures zoom in on the lower-left triangle in the first two figures. The top directions in $r = 8$ are included in $r = 64$, and vice versa.

$$\phi(A, B, i, j) = \psi(U_A^i, U_B^j) = \frac{\sum_{i=1}^p \sigma_i^2}{p} = \frac{1}{p} \left(1 - d(U_A^i, U_B^j)^2\right)$$

This similarity satisfies that if U_A^i and U_B^j share the same column span, then $\phi(A, B, i, j) = 1$. If they are completely orthogonal, then $\phi(A, B, i, j) = 0$. Otherwise, $\phi(A, B, i, j) \in (0, 1)$.

conculusion

Conclusion

- Large-scale language model을 효율적으로 튜닝하는 LORA 제안
- Adapter 류의 기법과 다르게 inference latency가 발생하지 않음
- Prefix-tuning과 다르게 usable sequence length를 줄일 필요가 없음
- 가중치 업데이트 행렬이 low intrinsic rank를 가진다고 가정
- 논문에선 LM에 초점을 맞췄지만 이론적으로 모든 dense layer에 적용가능

END

읽어주셔서 감사합니다.

-Shin joong hyun-