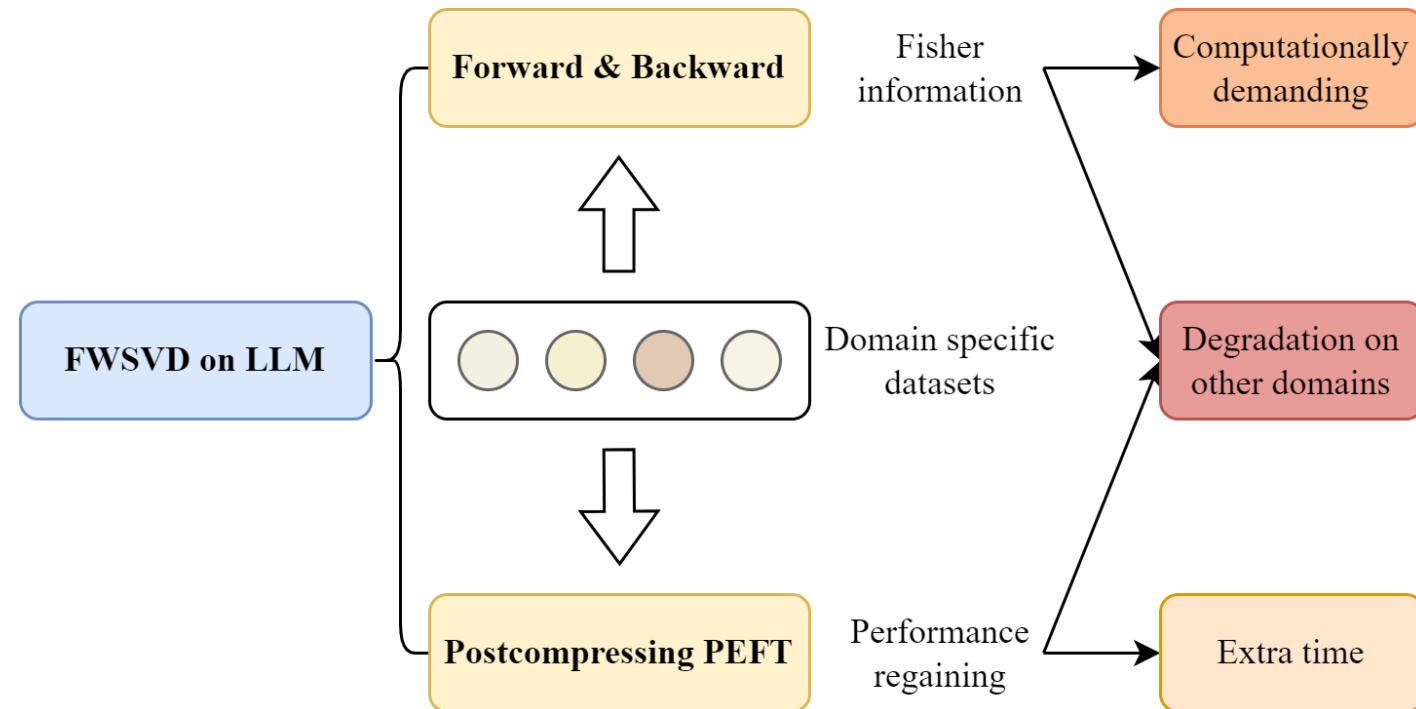


Motivation

Compression based on SVD

- FWSVD
 - Preliminary study: competitive with structure pruning on language modeling
 - Compatible with quantization
- To be improved:
 - Computationally demanding
 - Need peft for calibration

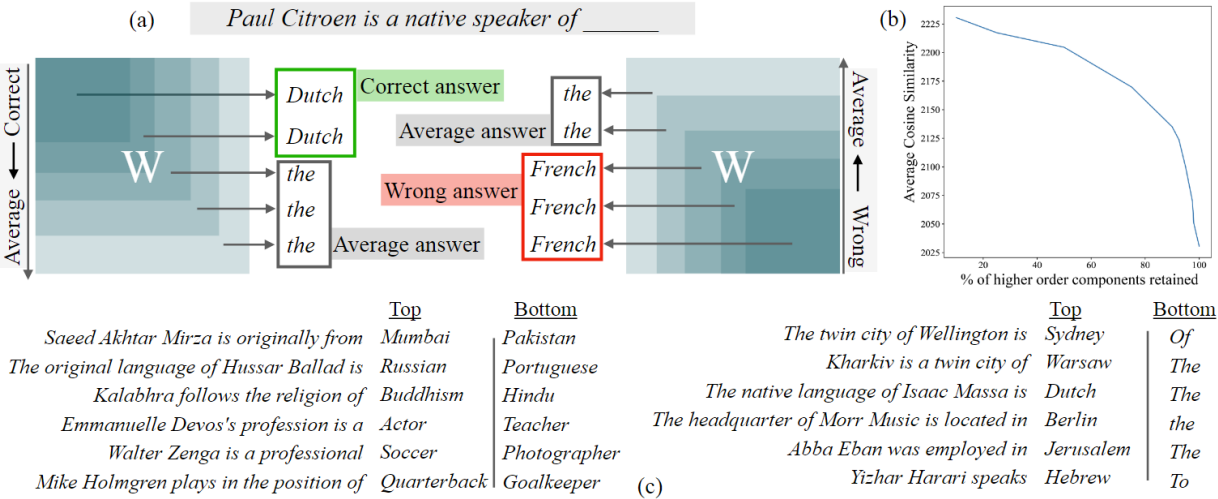


LoRA Weighted SVD

- Degradation on other domains
 - Retuning-free base model with competitive LM ability
- Memory footprint
 - Utilize LoRA to avoid computing full params' gradient, saving 25% GPU mem
 - QLoRA is also compatible
- Further improvement
 - Layer-wise compressing strategy
 - Kernel for efficient SVDLinear OP

Theoretical Research

- Current trend: low-rank + quant/sparse
 - LQ-LoRA: low-bit weights + LoRA
 - LoSparse: SP weight matrix + LoRA
 - Layer-Selective Rank Reduction: components with lower singular values may introduce “noise”
- Topic:
 - Low order/coherent parts的作用



values. Therefore, the coherent parts of neurons can be well approximated by the low-rank matrix computed by singular value thresholding.

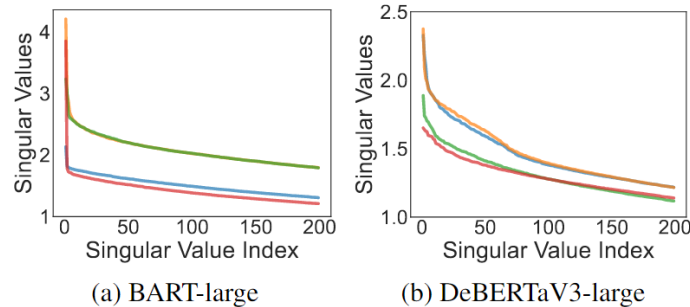
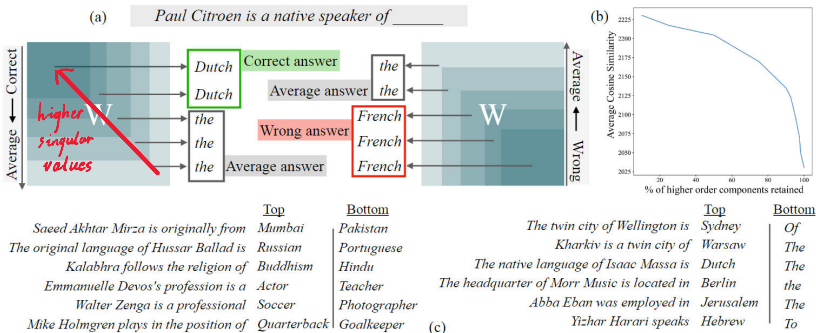


Figure 3. Singular values in language models. (a) Singular values of weight matrices of the 10th decoder layer in BART-large; (b) Singular values of weight matrices of the 14th encoder layer in DeBERTaV3-large.

日期: /



LASER indicates that ①头部奇异值已经可以做出语言建模
②拖尾部分实际会给出一些相似词 (noise)

values. Therefore, the coherent parts of neurons can be well approximated by the low-rank matrix computed by singular value thresholding.

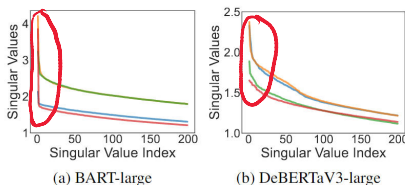


Figure 3. Singular values in language models. (a) Singular values of weight matrices of the 10th decoder layer in BART-large; (b) Singular values of weight matrices of the 14th encoder layer in DeBERTaV3-large.

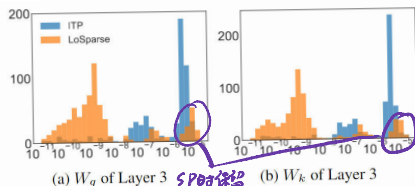


Figure 4. Neuron importance scores of selected linear projections when compressing DeBERTaV3-base on SST-2 with ITP (blue) and LoSparse (orange). It shows LoSparse successfully separates incoherent parts of neurons and make it easy to prune the non-expressive components.

LoSparse: 头部奇异值 (coherent parts) 使用 SVD 保留, 拖尾部分用结构化剪枝 (SP)

Topic: coherent part 在 LM 中有什么作用

Hypo: 形成 global pattern, 为 output feature 定一个 ground, 由拖尾部分进行微调

Prove: compare activate feature map ① coherent parts
分布是否有明显区别 ② coherent parts + tail

Further explore: SVD 为什么不 work → 没考虑拖尾

FWSVD, SP 为什么不 work → LoRA 与 coherent part 的关系

