

Gradient-Free Structured Pruning with Unlabeled Data



Azade Nova



Hanjun Dai



Dale Schuurmans

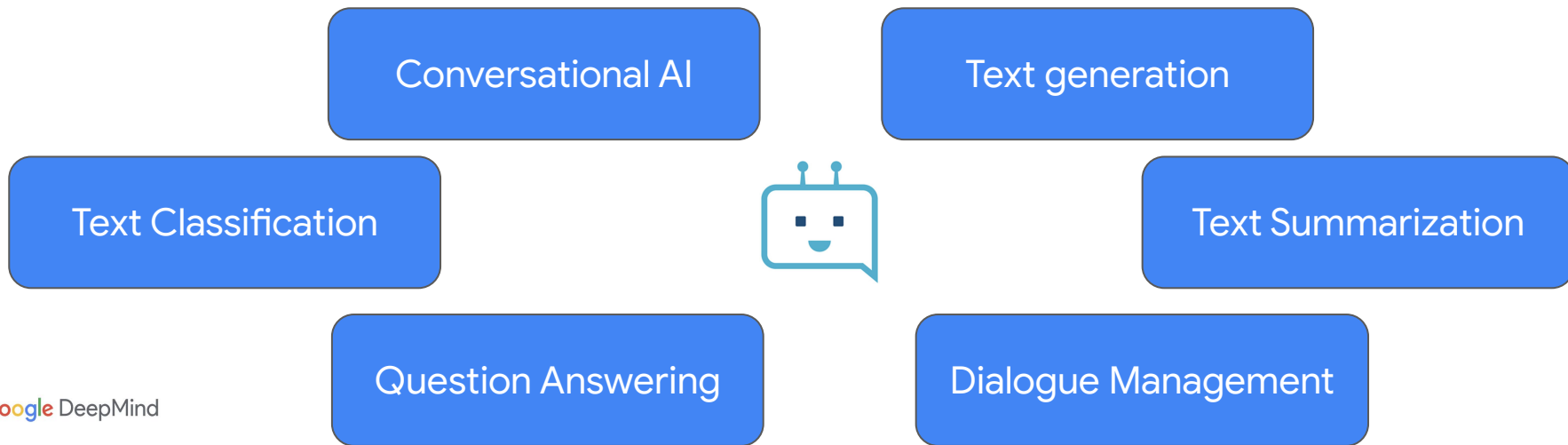
Agenda

- 01 Motivation
- 02 Existing Methods
- 03 Kernelized Convex Masking
- 04 Experimental Results
- 05 Q&A

Motivation

LLMs have achieved great success in solving difficult tasks across many domains.

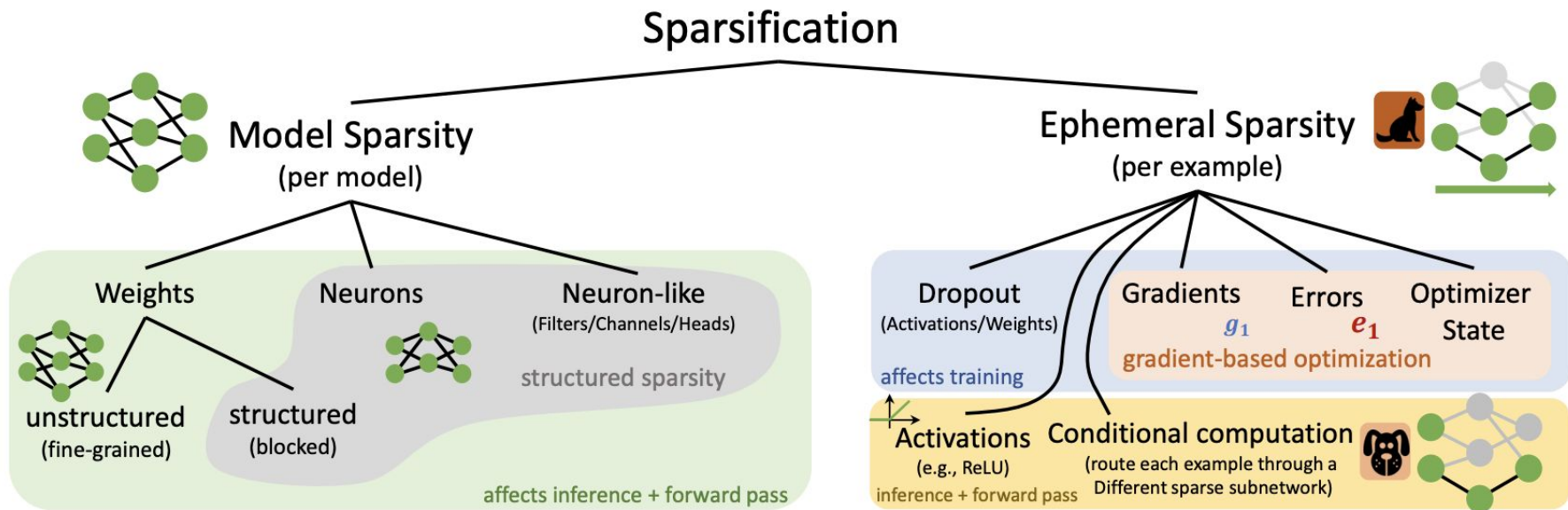
- Cost of high parameter counts
- Significant computational overhead
- Inference latency



Overview: Compressing and Optimizing Models

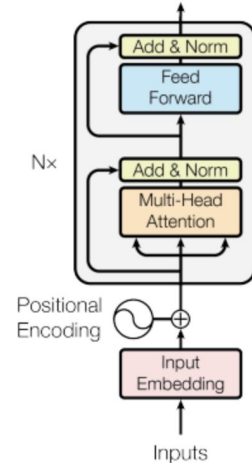
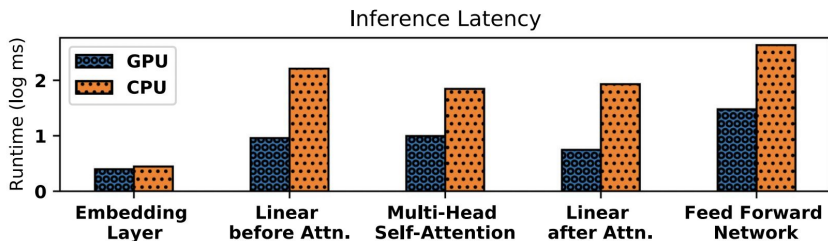
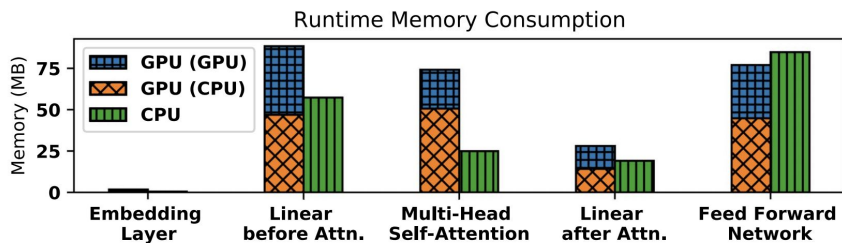
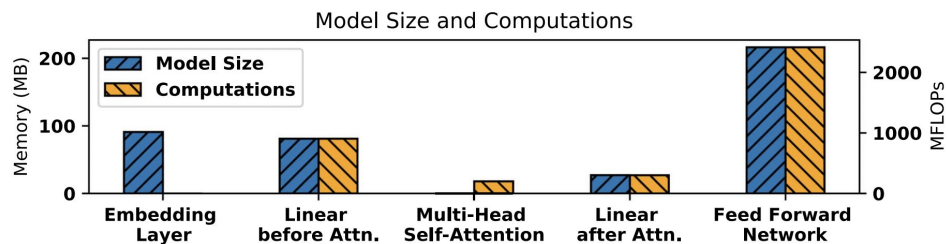
- Downsizing models
 - creates smaller dense networks, e.g through distillation or neural architecture search
- Operator factorization
 - Represent/approximate matrices by factored decomposition $W_{[n \times m]} = A_{[n \times r]} B_{[r \times m]}$
- Value quantization
 - Low precision value representation- weights, activation, etc
- Value Compression
 - Compress values, e.g., entropy-based (e.g., Huffman) or correlation based (e.g., gzip).
- Parameter sharing
 - Reuse parameters across neurons, e.g., Shapeshifter networks or CNNs
- Sparsification / pruning
 - Reduce the representational complexity using only a subset of the dimensions at a time

Overview: Sparsity



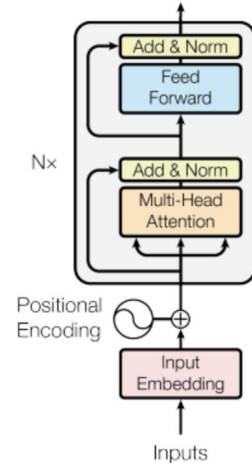
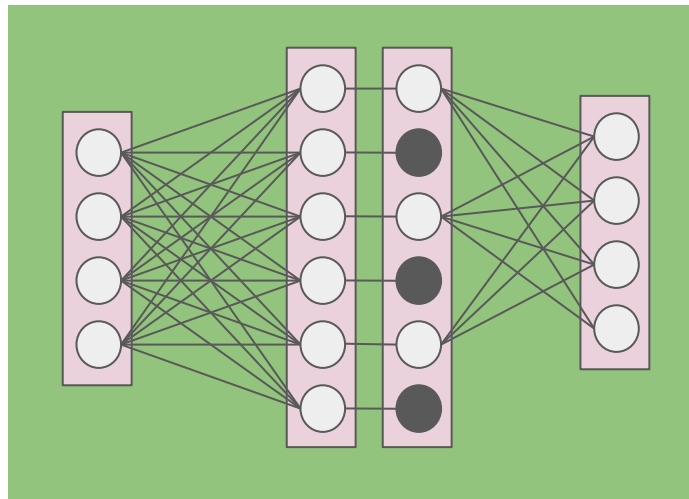
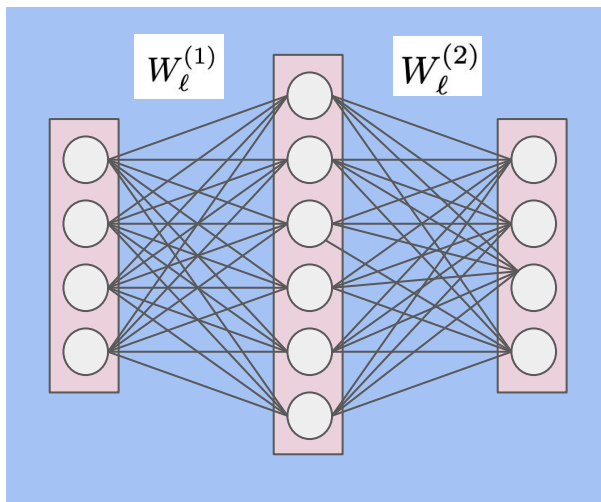
Existing approaches are complex, require a lot of labeled data, and require significant engineering effort to implement.

Breakdown Analysis of BERT_{BASE}



FFN sub-units are the parts consuming the most memory in terms of model size and executing the highest number of FLOPs

Structured Pruning by Masking



$$\widehat{FFN}_\ell(x) = \sum_{i=1}^N (\sigma(xW_\ell^{(1)}[:, i] + b_\ell^{(1)})W_\ell^{(2)}[i, :] \circ m_i) + b_\ell^{(2)}$$

$$\operatorname{argmin}_{\mathcal{M}} \mathcal{L}(\mathcal{M}) \quad s.t. \quad \operatorname{Cost}(\mathcal{M}) \leq \mathcal{C}$$

supervised setting with respect to minimizing the accuracy loss of the original model

3072 filters in BERT_{BASE}

$$W_\ell^{(1)} \in \mathbb{R}^{768 \times 3072}$$

$$W_\ell^{(2)} \in \mathbb{R}^{3072 \times 768}$$

Baselines from Structured Pruning Methods

Method	Gradient-free (! ∇)	Retrain/Finetune-free	Supervision-free	Pruning time $\leq 7min$
FLOP (Wang et al., 2019)	✗	✗	✗	✗
SLIP (Lin et al., 2020)	✗	✗	✗	✗
Sajjad et al. (Sajjad et al., 2023)	✗	✗	✗	✗
DynaBERT (Hou et al., 2020)	✗	✗	✗	✗
EBERT (Liu et al., 2021b)	✗	✗	✗	✗
Mask-Tuning (Kwon et al., 2022)	✗	✓	✗	✓
Weight-Magnitude (Li et al., 2016)	✓	✓	N/A	✓
Weight-Magnitude-Scale	✓	✓	✓	✓
KCM (ours)	✓	✓	✓	✓

$$\widehat{FFN}_\ell(x) = \sum_{i=1}^N (\sigma(xW_\ell^{(1)}[:, i] + b_\ell^{(1)}))W_\ell^{(2)}[i, :] \circ m_i + b_\ell^{(2)}$$

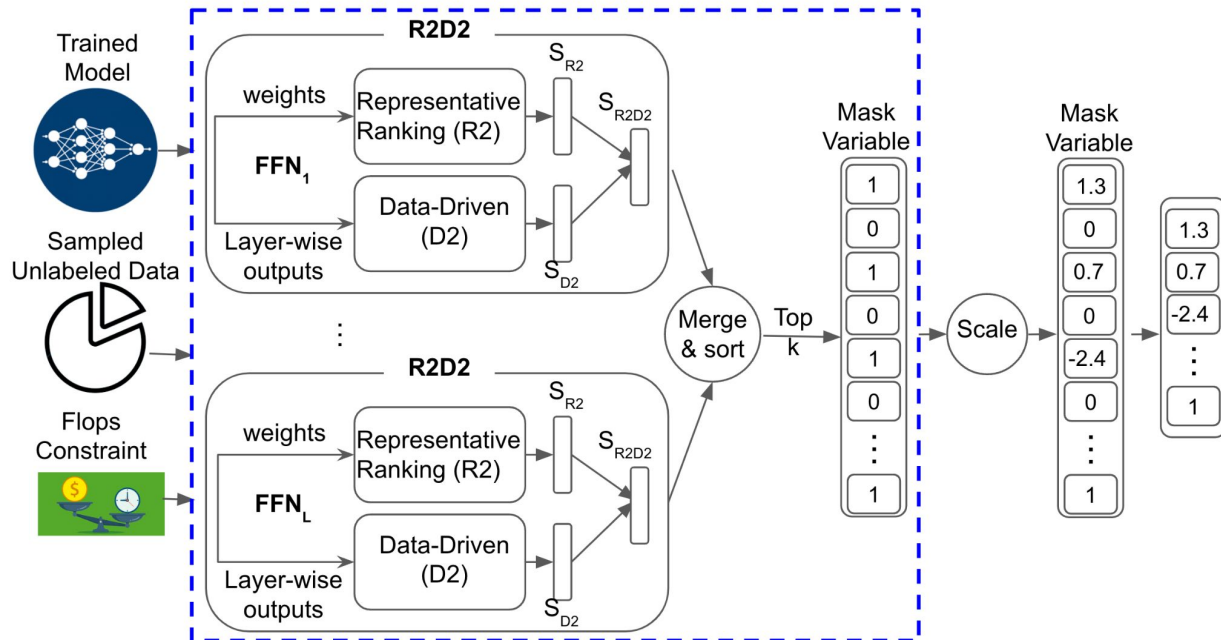
$$\operatorname{argmin}_{\mathcal{M}} \mathcal{L}_{FMT}(\mathcal{M}) \quad s.t. \quad Cost(\mathcal{M}) \leq C$$

unsupervised setting

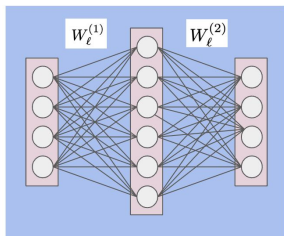
$$\mathcal{L}_{(\ell)FMT}(m) = \|FFN_{(\ell)}(x) - \widehat{FFN}_{(\ell)}(x)\|_2$$

Kernelized Convex Masking (KCM)

1. R2D2
 - a. Representative Ranking (R2)
 - b. Data-Driven (D2)
2. Merge and Sort
3. Scale

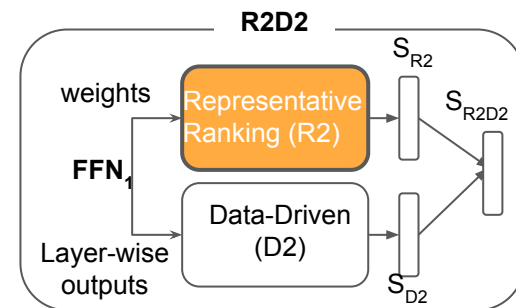


Representative Ranking (R2)



$$FFN_\ell(H_\ell^{(1)}) = H_\ell^{(1)} W_\ell^{(2)} + b_\ell^{(2)}$$

N points in d-dim



- Structured pruning goal : Select a subset of data points as representative.
- For linear functions, this problem can be reduced to finding a convex hull.
 - Complexity: $\mathcal{O}(N^{d/2})$
 - Number of convex hull points radically increases with d.

Kernelized Convex Hull Approximation^[1]

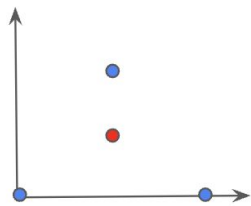
Representative Ranking (R2)

Find a positive coefficient matrix

$$C \in \mathbb{R}^{N \times N}$$

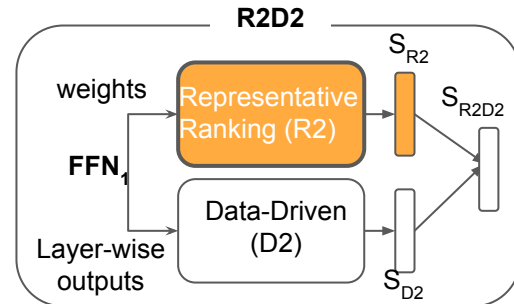
that minimizes

$$\|W_\ell^{(2)} - W_\ell^{(2)}C\|_2$$



$$\begin{bmatrix} 0 & 2 & 1 & 1 \\ 0 & 0 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 2 & 1 & 1 \\ 0 & 0 & 2 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 & 0.25 \\ 0 & 1 & 0 & 0.25 \\ 0 & 0 & 1 & 0.5 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

C



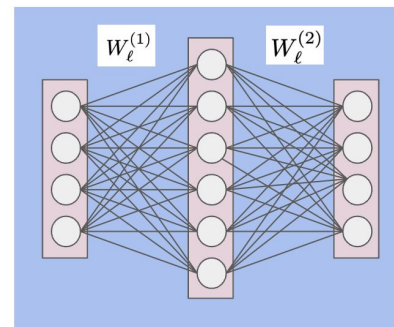
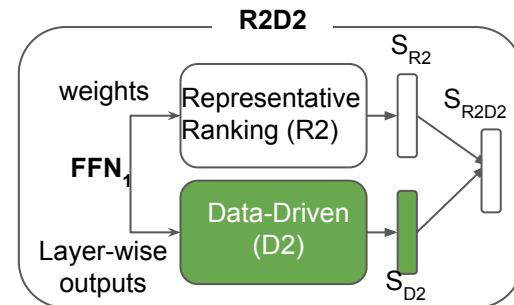
Algorithm 2 Representative Ranking (R2)

- 1: **Input:** Trained model: Model, Gaussian Kernel K with width σ , convergence rate α
- 2: **Output:** S_{R2} that represents importance of Filters in all layers.
- 3: **for** layer ℓ in layers of the Model **do**
- 4: $W_\ell^{(2)} \in \mathbb{R}^{N \times d}$ of FFN_ℓ
- 5: Initialize coefficient matrix: $C_0 \in \mathbb{R}^{N \times N} = \frac{1}{N}$
- 6: **repeat**
- 7:
$$C_{i+1} = C_i \circ \sqrt{\frac{K(W_\ell^{(2)}, W_\ell^{(2)})}{K(W_\ell^{(2)}, W_\ell^{(2)})C_i}}$$
- 8:
$$\delta = \frac{(C_{i+1} - C_i).sum()}{C_i.sum}$$
- 9: $C_i = C_{i+1}$
- 10: **until** convergence i.e. $\delta \leq \alpha$
- 11: $S_{R2}[\ell] = \text{diagonal}(C_i)$
- 12: **end for**
- 13: **return** S_{R2}

Data-Driven (D2)

Algorithm 1 Kernelized Convex Masking (KCM)

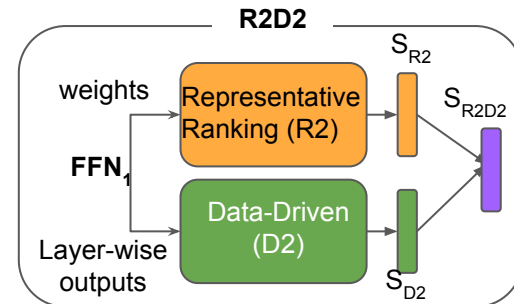
- 1: **Input:** Trained model: Model, FLOPs constraint \mathcal{C} , Gaussian Kernel K , convergence rate α
 - 2: **Output:** Mask \mathcal{M}
 - 3: Initialize mask \mathcal{M} as $\mathbf{0}$
//Call Representative Ranking (R2) Algorithm 2
 - 4: $S_{R2} = R2(\text{Model}, K, \alpha)$
// Data-Driven (D2) Ranking
 - 5: **for** batch in sample-data **do**
 - 6: for each layer ℓ in Model collect $H_\ell^{(1)}$
 - 7: $S_{D2}[\ell] = \text{average over } H_\ell^{(1)}$ for each filter
 - 8: **end for**
 - 9: $S_{R2D2}[\ell] = S_{R2}[\ell] * \text{normalized}(S_{D2}[\ell])$
 - 10: $k = \text{Number of neurons to satisfy FLOPs constraint } \mathcal{C}$
 - 11: Candidates = top- k filters of the sorted S_{R2D2}
 - 12: $\mathcal{M}[\text{Candidates}] = 1.0$
 - 13: **return** \mathcal{M}
-



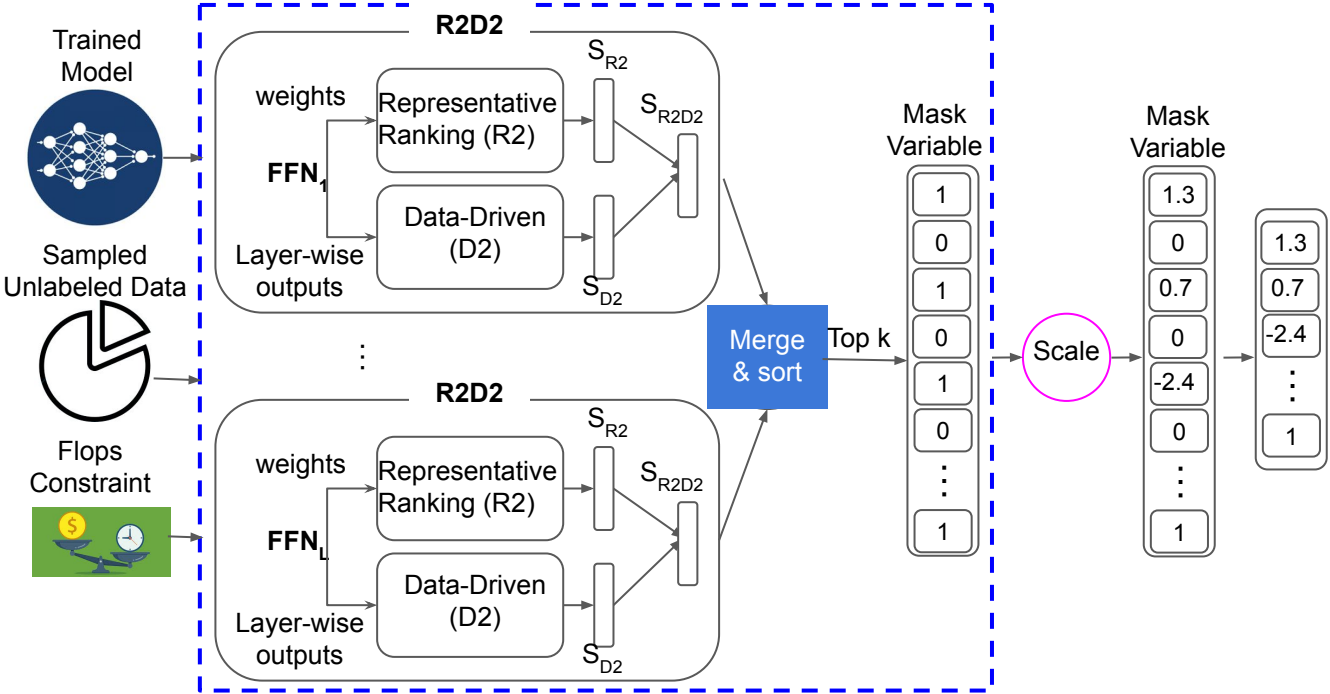
R2D2

Algorithm 1 Kernelized Convex Masking (KCM)

- 1: **Input:** Trained model: Model, FLOPs constraint \mathcal{C} , Gaussian Kernel K , convergence rate α
 - 2: **Output:** Mask \mathcal{M}
 - 3: Initialize mask \mathcal{M} as $\mathbf{0}$
//Call Representative Ranking (R2) Algorithm 2
 - 4: $S_{R2} = R2(\text{Model}, K, \alpha)$
// Data-Driven (D2) Ranking
 - 5: **for** batch in sample-data **do**
 - 6: for each layer ℓ in Model collect $H_\ell^{(1)}$
 - 7: $S_{D2}[\ell] = \text{average over } H_\ell^{(1)}$ for each filter
 - 8: **end for**
 - 9: $S_{R2D2}[\ell] = S_{R2}[\ell] * \text{normalized}(S_{D2}[\ell])$
 - 10: $k = \text{Number of neurons to satisfy FLOPs constraint } \mathcal{C}$
 - 11: Candidates = top- k filters of the sorted S_{R2D2}
 - 12: $\mathcal{M}[\text{Candidates}] = 1.0$
 - 13: **return** \mathcal{M}
-



Merge and Scale

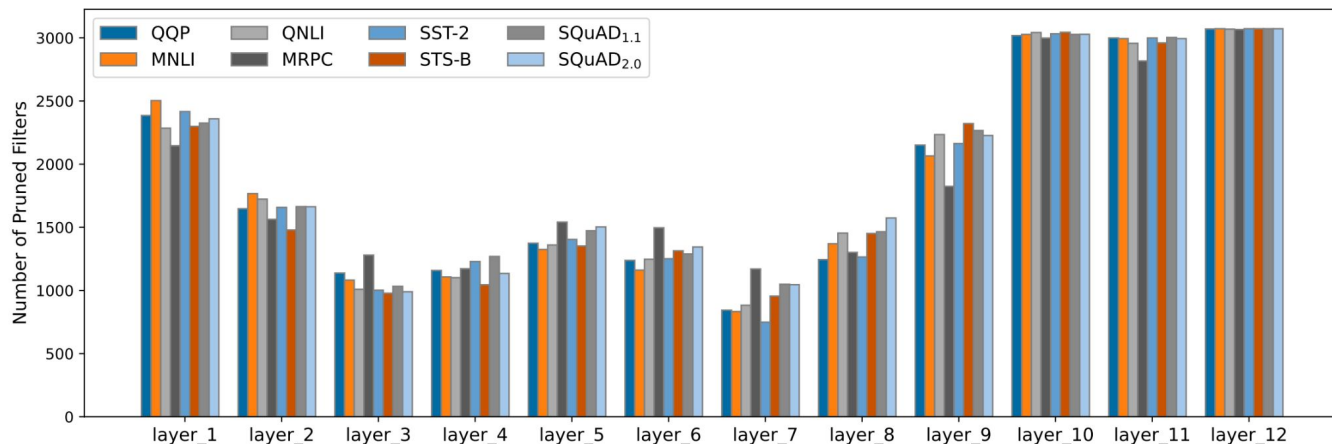


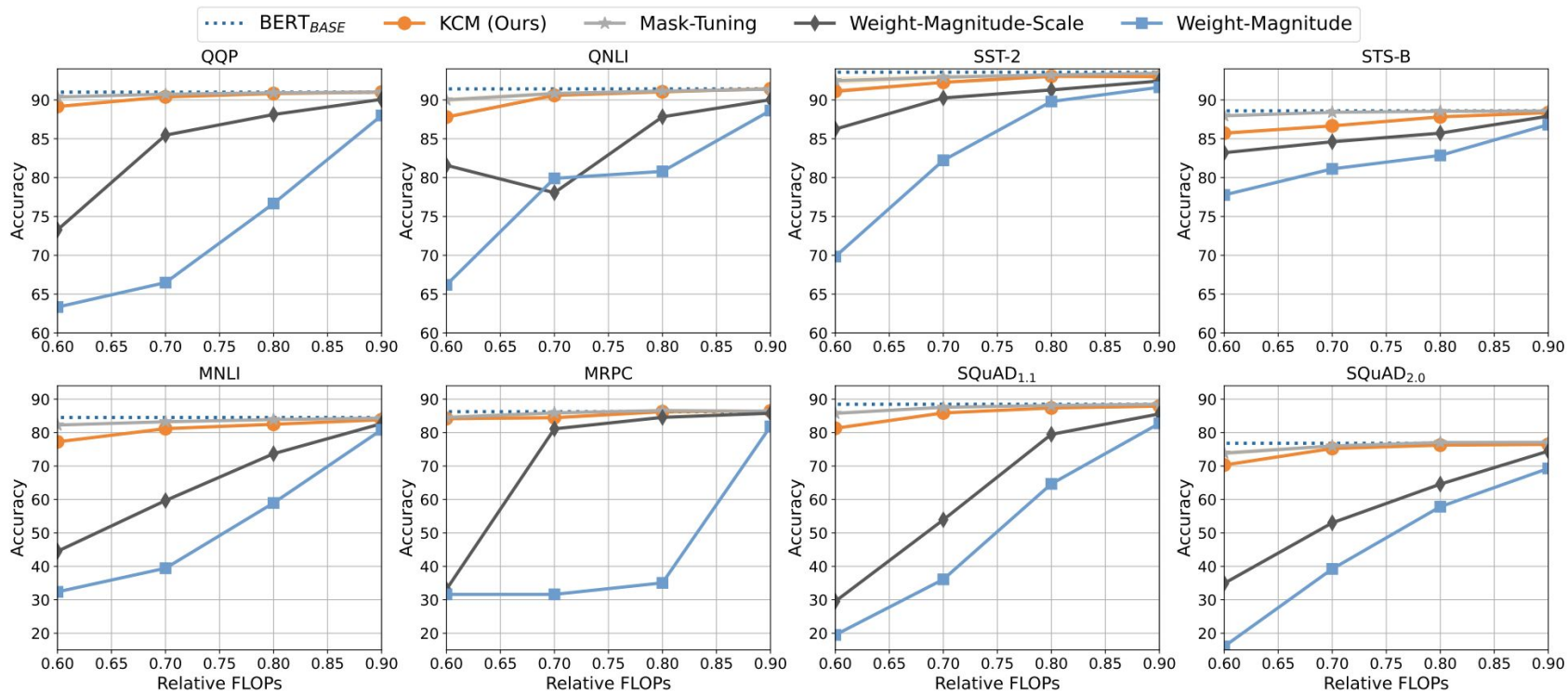
Experiments on BERT_{BASE}

!∇	Method	QQP		MNLI		MRPC		QNLI	
		60%	70%	60%	70%	60%	70%	60%	70%
	baseline	91.00		84.53		86.27		91.41	
✗	Mask-Tuning	90.38 ± 0.07	90.74 ± 0.07	82.26 ± 0.21	83.24 ± 0.16	84.51 ± 0.63	85.91 ± 0.40	90.00 ± 0.26	90.83 ± 0.16
✓	KCM (Ours)	89.15 ± 0.04	90.39 ± 0.04	77.24 ± 0.10	81.18 ± 0.10	84.19 ± 0.44	84.46 ± 0.29	87.79 ± 0.15	90.58 ± 0.08

!∇	Method	SST-2		STS-B		SQuAD _{1.1}		SQuAD _{2.0}	
		60%	70%	60%	70%	60%	70%	60%	70%
	baseline	93.57		88.59		88.48		76.82	
✗	Mask-Tuning	92.47 ± 0.41	92.92 ± 0.26	87.95 ± 0.12	88.40 ± 0.05	85.77 ± 0.41	87.57 ± 0.11	73.86 ± 0.55	76.00 ± 0.29
✓	KCM (Ours)	91.11 ± 0.23	92.26 ± 0.09	85.72 ± 0.12	86.66 ± 0.05	81.29 ± 0.06	85.89 ± 0.04	70.30 ± 0.13	75.24 ± 0.10

Dynamic Pruning Size





+ Limited labeled data

- One forward-backward pass and gather the gradient over the mask variables
- Refine top-k results

!∇	Method	QQP		MNLI		MRPC		QNLI	
		60%	70%	60%	70%	60%	70%	60%	70%
	baseline	89.99		82.11		84.80		88.56	
✗	Mask-Tuning	88.71 ± 0.22	89.66 ± 0.06	80.51 ± 0.19	81.65 ± 0.09	84.73 ± 0.71	84.83 ± 0.35	87.72 ± 0.38	88.43 ± 0.07
✓	KCM	88.16 ± 0.03	89.28 ± 0.03	78.05 ± 0.08	80.60 ± 0.05	79.66 ± 0.27	83.01 ± 0.16	85.93 ± 0.09	86.93 ± 0.13
✗	Extension(512 labeled data)	88.76 ± 0.25	89.45 ± 0.07	80.02 ± 0.25	81.37 ± 0.11	83.70 ± 1.40	84.49 ± 0.49	87.21 ± 0.54	88.21 ± 0.15
✗	Extension(1k labeled data)	88.92 ± 0.20	89.53 ± 0.08	80.41 ± 0.11	81.50 ± 0.12	84.17 ± 0.45	84.68 ± 0.49	87.60 ± 0.31	88.29 ± 0.16

!∇	Method	SST-2		STS-B		SQuAD _{1.1}		SQuAD _{2.1}	
		60%	70%	60%	70%	60%	70%	60%	70%
	baseline	91.40		86.12		85.73		68.84	
✗	Mask-Tuning	90.44 ± 0.41	90.93 ± 0.24	85.73 ± 0.07	85.96 ± 0.10	83.20 ± 0.16	84.64 ± 0.09	62.36 ± 1.40	65.32 ± 0.48
✓	KCM	88.38 ± 0.25	90.61 ± 0.25	85.26 ± 0.02	85.55 ± 0.03	76.92 ± 0.11	82.65 ± 0.06	64.56 ± 0.11	68.19 ± 0.06
✗	Extension(512 labeled data)	89.32 ± 0.54	90.38 ± 0.35	85.83 ± 0.11	86.02 ± 0.07	81.41 ± 0.26	83.30 ± 0.09	66.51 ± 0.24	67.72 ± 0.18
✗	Extension(1k labeled data)	89.86 ± 0.56	90.62 ± 0.40	85.90 ± 0.11	86.04 ± 0.06	81.16 ± 0.22	83.34 ± 0.11	66.35 ± 0.32	67.74 ± 0.18

Conclusion

- We studied the problem of structured pruning with unlabeled data and no backward pass.
- We proposed a gradient-free structured pruning framework that prunes the filters with the help of our proposed R2D2 that combines two ranking techniques called Representative Ranking (R2) and Data-Driven (D2).
- We empirically evaluated our framework on GLUE and SQuAD benchmarks using BERT_{BASE} and DistilBERT. Compared to when the labeled data is available, our approach achieved up to 40% FLOPs reduction with less than 4% accuracy loss over all tasks considered.