# Integrating Independent Layer-Wise Rank Selection with Low-Rank SVD Training for Model Compression: A Theory-Driven Approach

Yifan Guo[1] and Alyssa Yu[2]

[1] Towson University, Towson, Maryland, USA

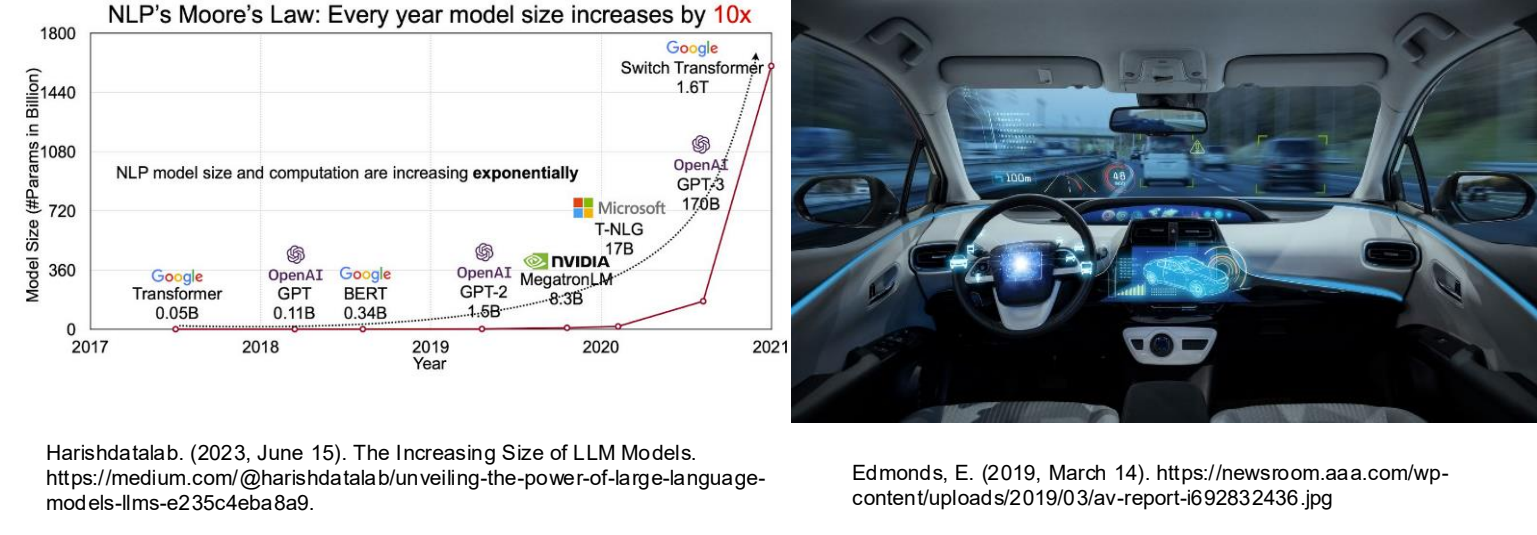[2] Poolesville High School, Poolesville, Maryland, USA

## Introduction

### Need for Model Compression

- In recent years, neural networks have grown drastically in complexity with billions or trillions of parameters.
- Reducing model size is crucial for scenarios demanding fast re-training and inference on resource-constrained edge devices.

### Uses of Model Compression

- Large language models require substantial computing resources to accommodate for their large model sizes
- Model compression allows for fast real-time decision making in self-driving cars.



Harshdatalabs (2023, June 15). The Increasing Size of LLM Models. https://medium.com/@harshdatalab/unveiling-the-power-of-large-language-models-llms-e230c4eba9a9.

Edmonds, E. (2019, March 14). https://newsroom.aaa.com/wp-content/uploads/2019/03/av-report-692832436.jpg

### Low-Rank Factorization and Truncation

Given $W$ an $m \times n$ matrix, its Singular Value Decomposition (SVD) is:

$$W = U\Sigma V^\top,$$

- $U$ is an orthogonal $m \times m$ matrix whose columns are the left singular vectors.
- $\Sigma$ is an $m \times n$ diagonal matrix with singular values $\sigma_1 \geq \cdots \geq \sigma_r > 0$ on the diagonal (denoted as singular values and $r$ is the rank of $W$ and $r \leq \min\{m,n\}$.) and the rest entries are all zero.
- $V^\top$ is the transpose of an orthogonal $n \times n$ matrix whose rows are the right singular vectors.



**Truncation:** The full-rank matrix $W$ has rank $r = \min\{m, n\}$. With a chosen $k \leq r$, we obtain the corresponding low-rank matrix $W_k = U_k \Sigma_k V_k^T$, where $U_k, \Sigma_k, V_k$ are the top-$k$ components truncated from $U$, $\Sigma$, and $V$.
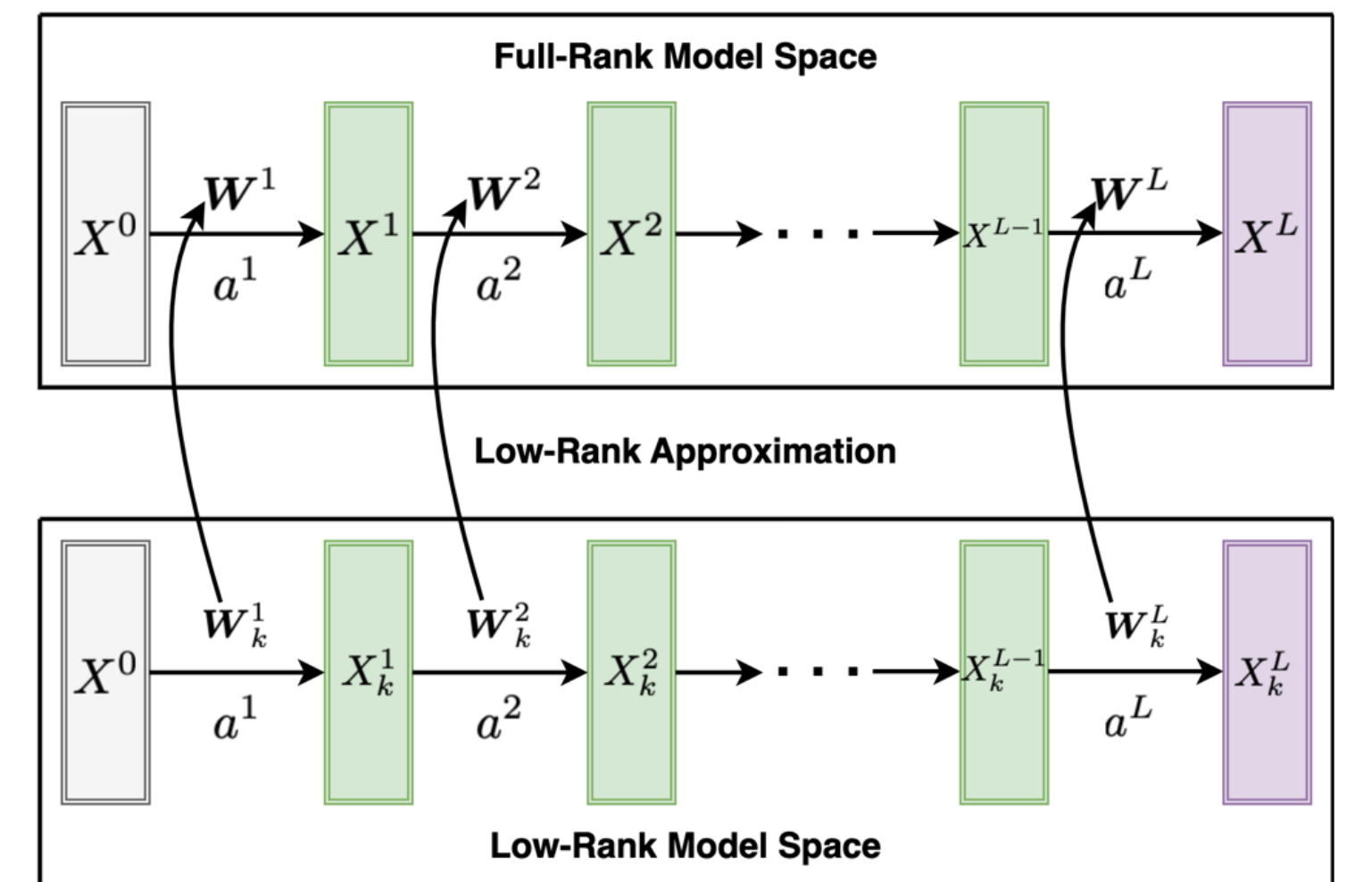


## Framework

### Uniqueness of Our Approach

- Low-rank training incorporates various penalty terms in the loss function to reduce the rank of weight matrices while preserving high accuracy.
- Three aspects considered in optimal rank selection:
  - **Independent** vs. Dependent Layer-Wise Rank Selection
  - **Theory-Driven** vs. Heuristic Rank Search Strategy
  - **During-Training** vs. Post-Training Rank Selection.
- Our approach includes independent layer-wise rank selection, is theory-driven, and conducts rank selection during training; no previous works include *all* three optimal aspects.

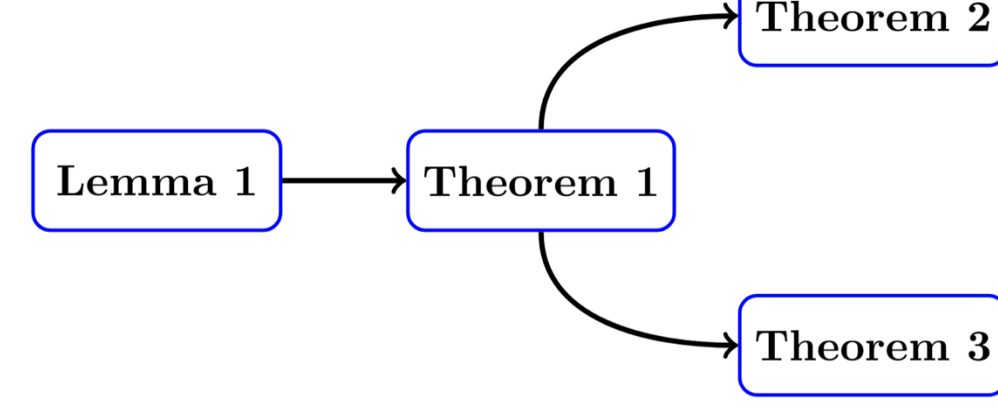| Approach | Layer-Wise Rank Selection (Independent vs. Dependent) | Search Strategy (Theory-Driven vs. Heuristic) | Rank Selection Timing (During- vs. Post-Training) |
|---|---|---|---|
| [Cheng *et al.*, 2020] | Dependent | Heuristic | Post-Training |
| [Kim *et al.*, 2021] | Dependent | Theory-Driven | Post-Training |
| [Idelbayev *et al.*, 2020] | Independent | Heuristic | Post-Training |
| [Sobolev *et al.*, 2022] | Dependent | Heuristic | Post-Training |
| [Xiao *et al.*, 2023] | Dependent | Heuristic | Post-Training |
| [Wang *et al.*, 2023] | Independent | Heuristic | During-Training |
| [Cao *et al.*, 2024] | Dependent | Heuristic | Post-Training |
| Ours | Independent | Theory-Driven | During-Training |

### Mathematical Formulation of Neural Networks



- Neural network of $L$ layers denoted by $f_\mathcal{W}(\cdot)$ and parameterized by $\mathcal{W} = \{W^1, W^2, \ldots, W^L\}$.
- $X^l = a^l(W^l X^{l-1})$ for $l = 1, 2, \ldots, L$, where $W^l$, $a^l$, and $X^l$ represent the weight matrix, non-linear activation function, and output of the $l$-th layer.

## Our Theoretical Findings

**Research Problem:** Design a theory-driven approach to search for truncations $(W_k^1, \ldots, W_k^L)$ of $(W^1, \ldots, W^L)$ with layer wise ranks $(k^1, k^2, \ldots, k^L)$ to achieve optimal compression and accuracy concurrently.



**Lemma 1** (Eckart–Young–Mirsky Theorem [Golub *et al.*, 1987]). *Let $W \in \mathbb{R}^{m \times n}$ be a matrix with rank $r$ and let $\|\cdot\|_2$ denote the spectral norm. Following the same definitions in Proposition 1, we define $W_k$ to be the best rank-$k$ approximation of $W$ in the spectral norm, i.e., $W_k = U_k \Sigma_k V_k^T = \sum_{i=1}^k \sigma_i U_{:,i} V_{i,:}^T$, where $U_k$, $\Sigma_k$, and $V_k$ are the top-$k$ components truncated from $U$, $\Sigma$, and $V$, respectively. Then, we have $\|W - W_k\|_2 = \sigma_{k+1}$.*

### How Low-Rank Truncation Impact the Output

**Theorem 1** (The output difference bound under rank-$k$ approximation for $L$-layer neural networks). *Let $a^l$ be the activation function for the $l$-th layer, where $a^l$ is $\rho_l$-Lipschitz and satisfies $a^l(0) = 0$ for all $l \in [1, L]$. Let $W_l$ be the full-rank matrix at layer $l$, and let $W_l^k$ be its rank $k^l$ approximation from keeping only the top $k^l$ singular values in the SVD decomposition of $W_l^l$. Let $\sigma_i^l$ denote the $i$th singular value of $W_l$, $X^0$ be the initial input vector, and $X^l$ and $X_k^l$ be the output vectors at the $l$-th layer after applying $W^l$ and $W_k^l$, respectively. Then, the output difference $\|X^L - X_k^L\|_2$ from a rank-$k$ approximation over an $L$-layer feed-forward network is upper-bounded by $\left(\prod_{l=1}^L \rho_l \sigma_1^l\right)\left(\sum_{l=1}^L \frac{\sigma_{k^l+1}^l}{\sigma_1^l}\right)\|X^0\|_2$.*

### How Low-Rank Truncation Impact the Loss Error

- Let $X_i^L = f_\mathcal{W}(X_i^0) \in \mathbb{R}^C$ and $X_{i,k}^L = f_{\mathcal{W}_k}(X_{i,k}^0) \in \mathbb{R}^C$ be the output logits after feeding an input $X_i^0$ sampled from the training dataset $\mathcal{D}^{tr} = \{X_i^0, y_i\}_{i=1}^R$, from the full-rank parameter space $\mathcal{W} = \{W^1, W^2, \ldots, W^L\}$ and low-rank parameter space $\mathcal{W}_k = \{W_k^1, W_k^2, \ldots, W_k^L\}$, respectively.
- Let $B$ be a constant such that $\|X_i^0\|_2 \leq B$ for all $i \in [1, R]$.

**Theorem 2** (The loss error bound under rank-$k$ truncation in classification problems). *Following the definitions in Theorem 1, we consider a $C$-class classification problem. We let $z_i = \text{softmax}\left(X_i^L\right)$ and $z_{i,k} = \text{softmax}\left(X_{i,k}^L\right)$, where softmax is the softmax function. Consider the cross-entropy function $g(z, y) = -y^T \log(z)$ and let the loss functions be $L\left(\mathcal{W}; X_i^0\right) = g(z_i, y_i)$ and $L\left(\mathcal{W}_k; X_i^0\right) = g(z_{i,k}, y_i)$. We set $\sigma_{k^l+1}^l$ and $\delta$ such that $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta$, $\forall l \in [L]$. Then, for $\forall \epsilon > 0$, $\exists \delta = \frac{\epsilon}{\sqrt{2}BL\left(\prod_{l=1}^L \rho_l \sigma_1^l\right)}$ such that $|L(\mathcal{W}; \mathcal{D}^{tr}) - L(\mathcal{W}_k; \mathcal{D}^{tr})| < \epsilon$.*

**Theorem 3** (The loss error bound under rank-$k$ truncation in regression problems). *Following the definitions in Theorem 1, we consider a regression problem with loss function $g(z, y) = \|z - y\|_2$. Let $L\left(\mathcal{W}; X_i^0\right) = g\left(X_i^L, y_i\right)$ and $L\left(\mathcal{W}_k; X_i^0\right) = g\left(X_{i,k}^L, y_i\right)$. We set $\sigma_{k^l+1}^l$ and $\delta$ such that $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta$, $\forall l \in [L]$. Then, for $\forall \epsilon > 0$, $\exists \delta = \frac{\epsilon}{BL\left(\prod_{l=1}^L \rho_l \sigma_1^l\right)}$ such that $|L(\mathcal{W}; \mathcal{D}^{tr}) - L(\mathcal{W}_k; \mathcal{D}^{tr})| < \epsilon$.*

## Our Algorithms

- Calculating the Lipschitz constant of layer-wise parameters is computationally expensive in many CNN models, preventing the extension of our theoretical findings from simple feed-forward networks to complex CNNs.
- The rank selection is integrated in low-rank SVD training.
- The low-rank SVD training loss function is

$$L(\mathcal{U}, \Sigma, \mathcal{V}; D^{tr}) = \underbrace{L_T(\mathcal{U}, \Sigma, \mathcal{V})}_{\text{Training Loss}} + \lambda_O \underbrace{L_O(\mathcal{U}, \mathcal{V})}_{\text{Orthogonality Loss}} + \lambda_R \underbrace{L_R(\Sigma)}_{\text{Regularization Loss}}$$

- For the low-rank regularization loss, we consider the Nuclear norm $\|W^l\|_1$ and the Hoyer norm $\|W^l\|_1 / \|W^l\|_F$.

---

**Algorithm 1:** Proposed Rank Selection Algorithm

**Input:** full-rank parameters $\mathcal{W}$, loss error tolerance $\epsilon$, stop searching precision $\Delta\delta$
**Output:** low-rank parameters $\mathcal{W}_k$

1 **Function** RankSelection($\mathcal{W}$, $\epsilon$, $\Delta\delta$):
2    $floss \leftarrow \overline{L_T}(\mathcal{W}; \mathcal{D}^{tr})$;
3    $l \leftarrow 0$; $u \leftarrow 1$; $\delta \leftarrow (l+u)/2$;
4    **while** $|l - u| \geq \Delta\delta$ *or* $|floss - loss| \geq \epsilon$ **do**
5      $\mathcal{W}_k \leftarrow$ SVDLowRankApprox($\mathcal{W}, \delta$);
6      $loss \leftarrow \overline{L_T}(\mathcal{W}_k; \mathcal{D}^{tr})$;
7      **if** $|floss - loss| < \epsilon$ **then**
8        $l \leftarrow \delta$; $\delta \leftarrow (l+u)/2$;
9      **else**
10        $u \leftarrow \delta$; $\delta \leftarrow (l+u)/2$;
11    **return** $\mathcal{W}_k$ ;

---

**Algorithm 2:** Standard SVD Low-Rank Approximation
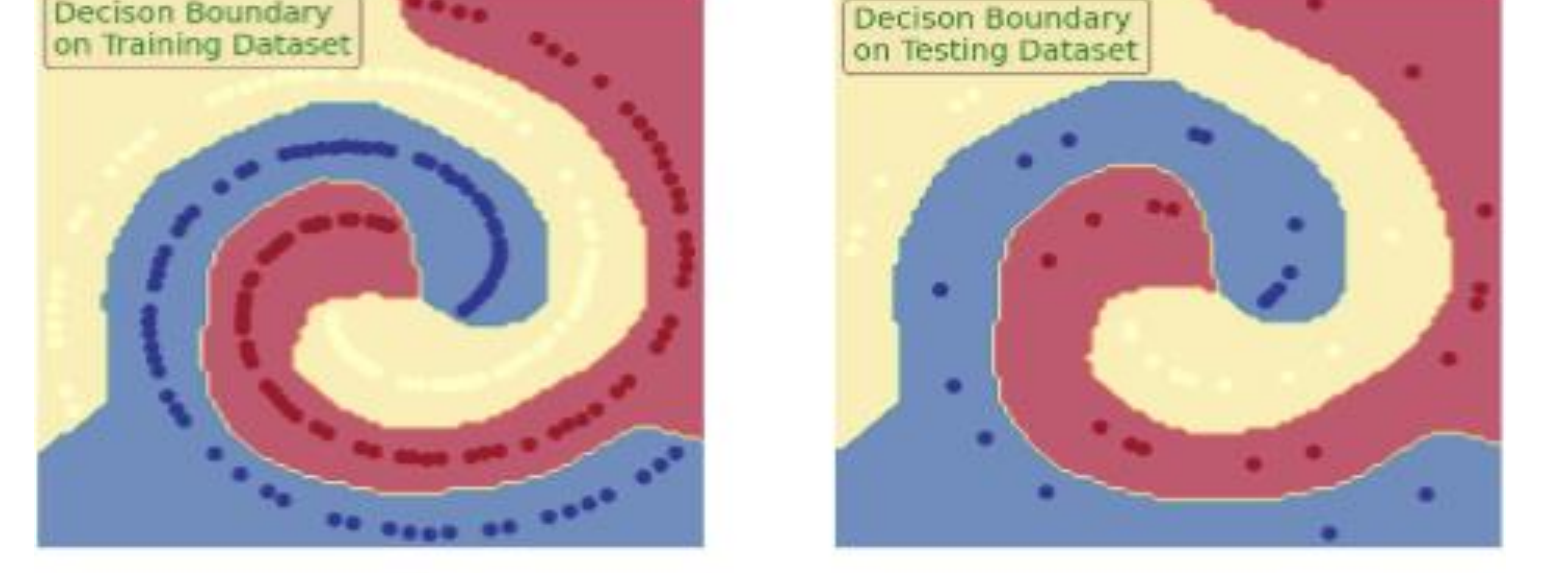
1 **Function** SVDLowRankApprox($\mathcal{W}, \delta$):
2    $\mathcal{W}_k \leftarrow \emptyset$;
3    **for** *each* $W^l \in \mathcal{W}$ **do**
4      $U^l, \Sigma^l, V^l \leftarrow$ SVD($W^l$);
5      $k^l \leftarrow \text{argmax}_k\{k|\sigma_k^l/\sigma_1^l \geq \delta\}$;
6      $W_k^l \leftarrow U_k^l \Sigma_k^l V_k^l$; $\mathcal{W}_k \leftarrow \mathcal{W}_k \cup W_k^l$;
     // $U_k^l, \Sigma_k^l, V_k^l$ are top-$k^l$ vectors truncated from $U^l, \Sigma^l, V^l$
7    **return** $\mathcal{W}_k$

---

**Algorithm 3:** Proposed Rank Selection Enabled Low-Rank SVD Training Algorithm
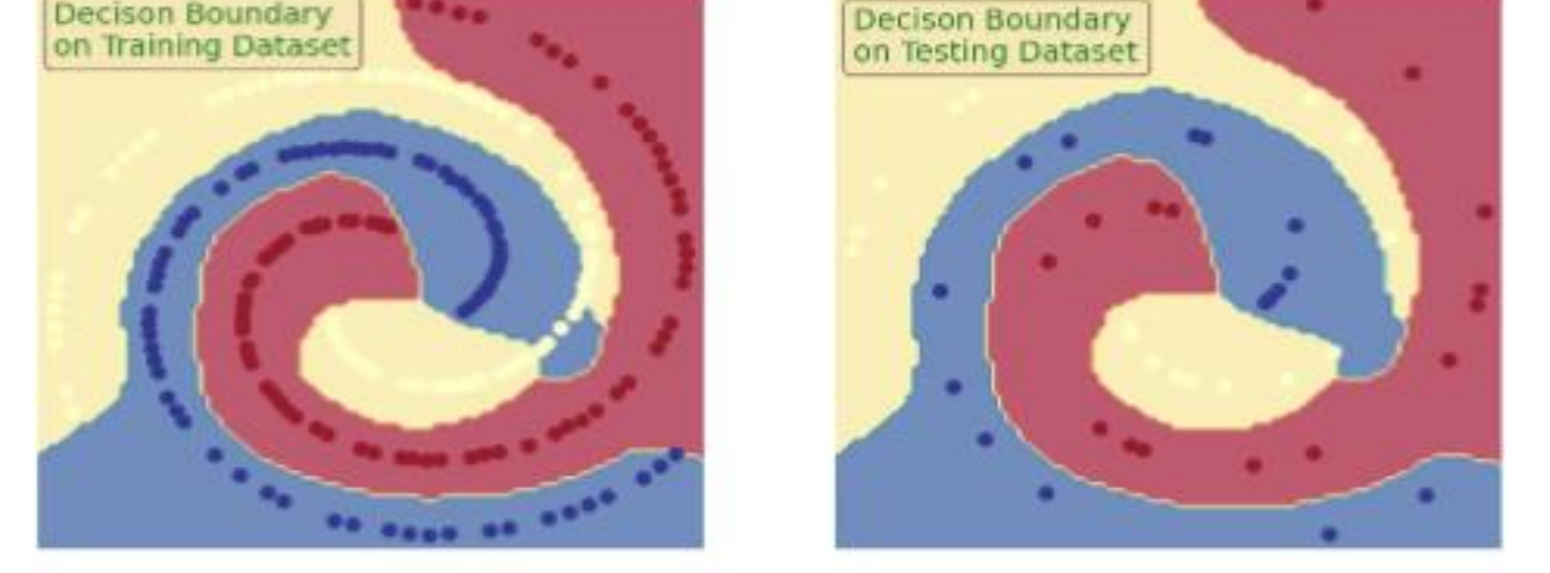
**Input:** full-rank parameters $\mathcal{U}, \Sigma, \mathcal{V}$, loss error tolerance $\epsilon$, stop searching precision $\Delta\delta$, training epoch $E$
**Output:** low-rank parameters $\mathcal{U}_k, \Sigma_k, \mathcal{V}_k$

1 Initialize parameters $\mathcal{U}, \Sigma, \mathcal{V}$ ;
2 $e \leftarrow 1$ ;
3 **while** $e \leq E$ **do**
4    Update $\mathcal{U}, \Sigma, \mathcal{V}$ based on loss function $L(\mathcal{U}, \Sigma, \mathcal{V})$ with an appropriate optimizer and extract the learning loss $L_T(\mathcal{U}, \Sigma, \mathcal{V})$; // $L_O(\mathcal{U}, \mathcal{V})$ and $L_R(\Sigma)$ are not used in the next-step rank selection
5    $\mathcal{U}_k, \Sigma_k, \mathcal{V}_k \leftarrow$ RankSelection($\mathcal{U}, \Sigma, \mathcal{V}, \epsilon, \Delta\delta, L_T(\mathcal{U}, \Sigma, \mathcal{V})$);
6    $\mathcal{U}, \Sigma, \mathcal{V} \leftarrow \mathcal{U}_k, \Sigma_k, \mathcal{V}_k$; // Use truncated models for the next round of training
7    $e \leftarrow e + 1$;
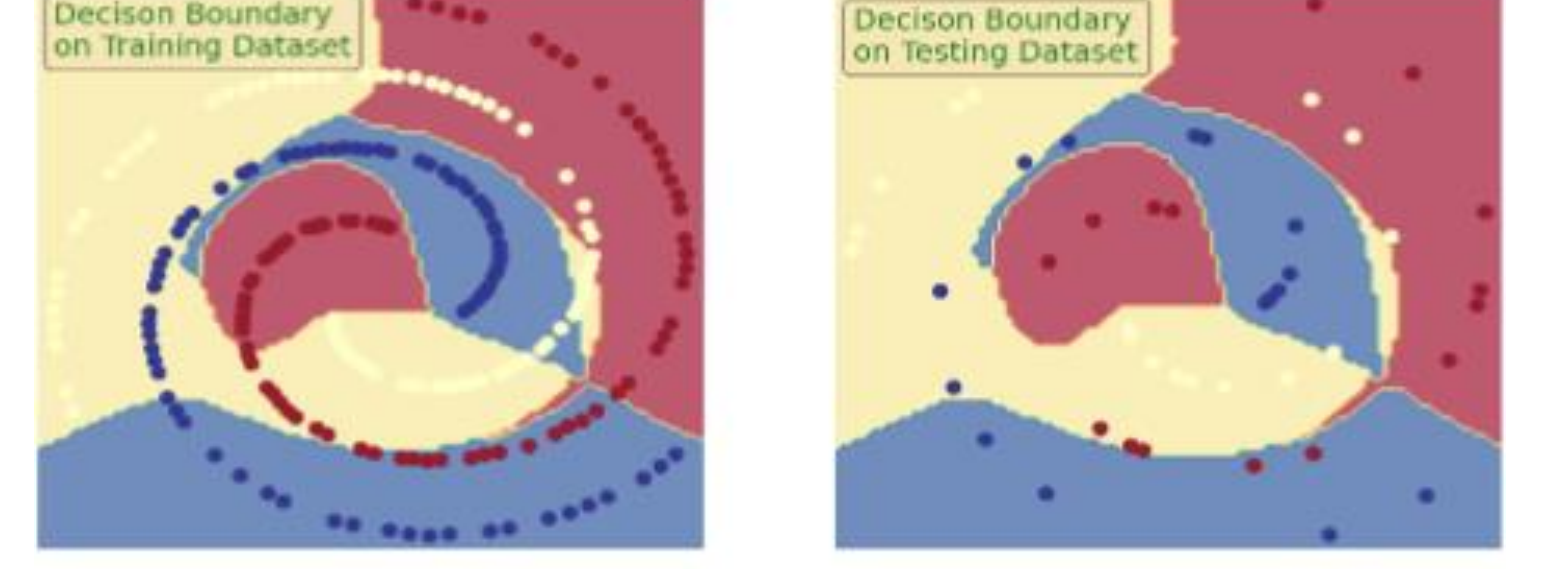8 **return** $\mathcal{U}_k, \Sigma_k, \mathcal{V}_k$ ;

## Proof-of-Concept



(a) Full-rank model space ($\epsilon = 0$)

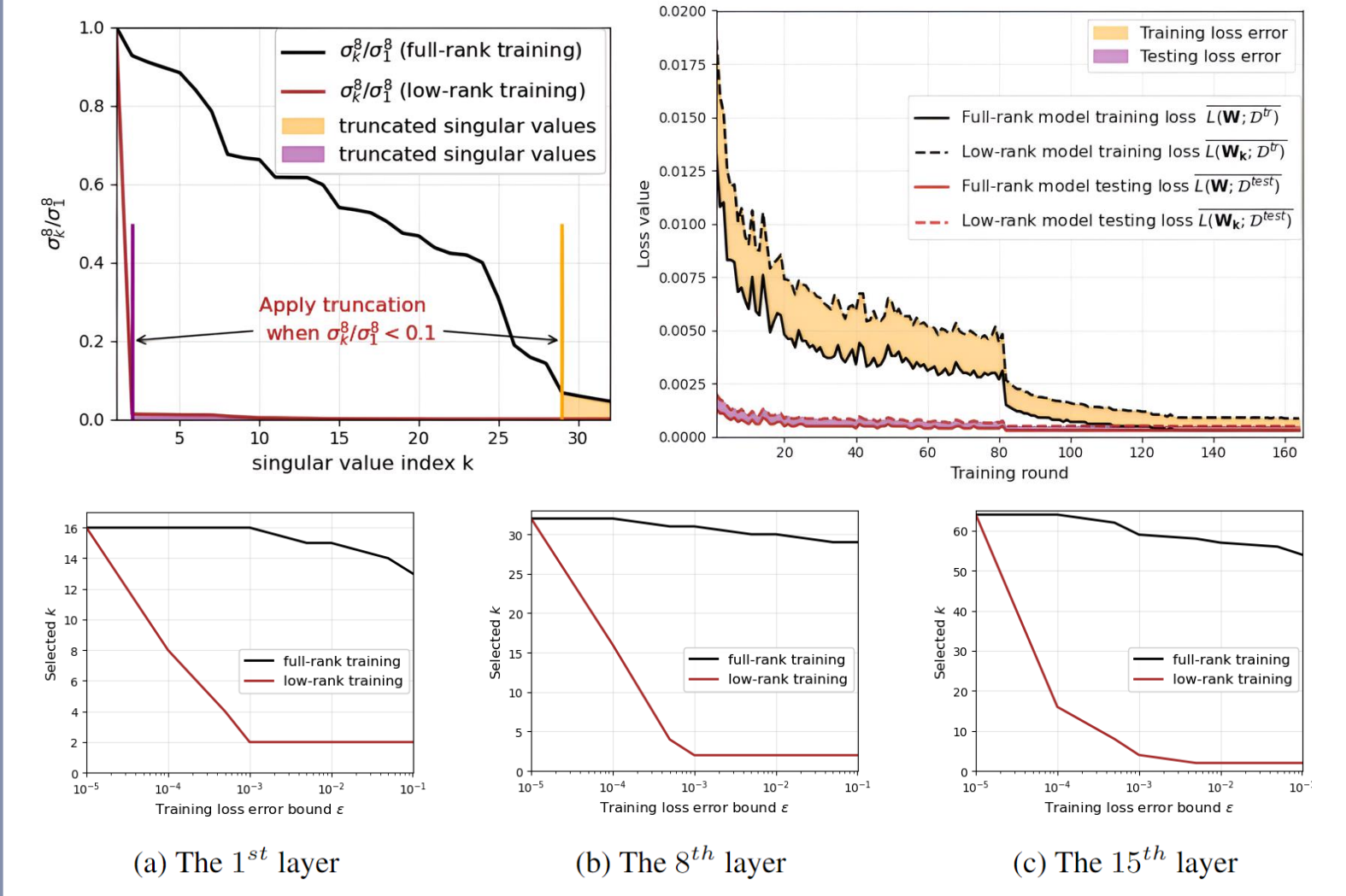(b) Low-rank model space ($\epsilon = 0.28$, $\delta = 0.025$)

(c) Low-rank model space ($\epsilon = 0.33$, $\delta = 0.03$)

**A visualization of the decision boundaries on the training dataset (left column) and testing dataset (right column) for different rank selections.**

- We conduct a pilot study using a simple 3-layer feedforward neural network for a ternary classification problem.
- Validate the feasibility of identifying $\delta$ based on our derived $\epsilon - \delta$ and determining the optimal $(k^1, k^2, \ldots, k^L)$

## Experimental Results



Correlation between training loss error and rank selection on ResNet-20 model

| Approach | ResNet-20 Test Acc | ResNet-20 CR ↓ | ResNet-32 Test Acc | ResNet-32 CR ↓ | ResNet-56 Test Acc | ResNet-56 CR ↓ | ResNet-110 Test Acc | ResNet-110 CR ↓ |
|---|---|---|---|---|---|---|---|---|
| [Yang *et al.*, 2020] (Channel, Squared Hoyer) | 0.887 | 0.202 | 0.893 | 0.367 | 0.924 | 0.485 | 0.924 | 0.511 |
| [Yang *et al.*, 2020] (Spatial, Squared Hoyer) | 0.866 | 0.242 | 0.889 | 0.413 | 0.918 | 0.522 | 0.917 | 0.543 |
| [Wang *et al.*, 2023] | 0.822 | 0.337 | 0.834 | 0.398 | 0.844 | 0.441 | 0.849 | 0.462 |
| Ours (Channel, Nuclear) | 0.870 | **0.155** | 0.879 | **0.312** | 0.892 | **0.422** | 0.892 | **0.445** |
| Ours (Spatial, Nuclear) | 0.867 | 0.221 | 0.873 | 0.319 | 0.899 | 0.438 | 0.898 | 0.465 |
| Ours (Channel, Squared Hoyer) | **0.892** | **0.155** | **0.899** | **0.312** | **0.929** | **0.422** | **0.928** | **0.445** |
| Ours (Spatial, Squared Hoyer) | 0.873 | 0.221 | 0.881 | 0.319 | 0.912 | 0.438 | 0.912 | 0.465 |

**Overall performance comparisons on CIFAR-10 dataset**

| Approach | ResNet-18 Test Acc | ResNet-18 CR ↓ | ResNet-50 Test Acc | ResNet-50 CR ↓ |
|---|---|---|---|---|
| [Yang *et al.*, 2020] (Channel, Squared Hoyer) | 0.684 | 0.204 | 0.691 | 0.392 |
| [Yang *et al.*, 2020] (Spatial, Squared Hoyer) | 0.670 | 0.221 | 0.678 | 0.411 |
| Ours (Channel, Squared Hoyer) | **0.690** | **0.181** | **0.696** | **0.366** |
| Ours (Spatial, Squared Hoyer) | 0.672 | 0.207 | 0.681 | 0.398 |

**Overall performance comparisons on ImageNet dataset**

## Conclusions

- In-depth theoretical analysis that quantitatively measures how low-rank approximation affects training losses.
- Rank selection enabled low-rank training inspired by our theoretical findings.
- Our algorithm, paired with channel decomposition and Hoyer regularization, achieves better results than the previous state-of-the-art studies.

## Future Work

- A stricter bound of results from Theorems 1-3 can give us better control over the balance between deep learning model compression and accuracy.
- The rank-selection algorithm could be based off a generalization bound for the predicted accuracy on unseen data.
- Rank-selection-based model compression can be implemented in other models other than ResNet.

## Key References

- [Idelbayev et al., 2020] Y. Idelbayev and M. Á. Carreira-Perpiñán, Low-rank compression of neural nets: Learning the rank of each layer, in Proc. of 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 8046-8056.
- [Wang et al., 2023] Hongyi Wang, Saurabh Agarwal, Yoshiki Tanaka, Eric Xing, Dimitris Papailiopoulos, et al. Cuttlefish: Low-rank model training without all the tuning. In Proc. of Machine Learning and Systems, vol. 5, pp: 578–605, 2023.
- [Yang et al., 2020] H. Yang, M. Tang, W. Wen, F. Yan, D. Hu, A. Li, H. Li, Y. Chen. Learning low-rank deep neural networks via singular vector orthogonality regularization and singular value sparsification. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp: 2899–2908.