



SCHOOL OF
**COMPUTING &
DATA SCIENCE**
The University of Hong Kong



An In-depth Investigation of Sparse Rate Reduction in Transformer-like Models

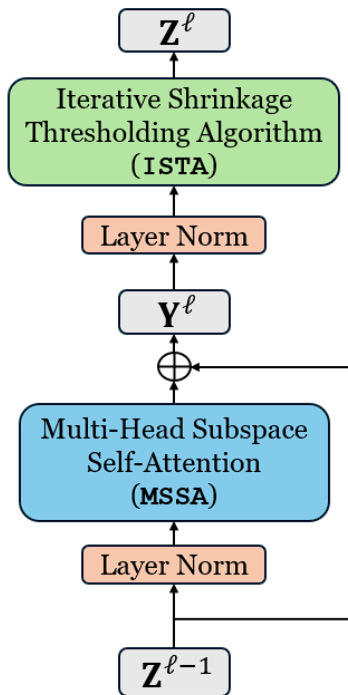
Yunzhe Hu, Difan Zou, and Dong Xu

The University of Hong Kong

Outline

- **Background:** Designing Transformer-like models via Sparse Rate Reduction (SRR)
- **Key Investigations:**
 - Analysis of behaviors of self-attention operator
 - Correlation between learning objective SRR and generalization
- **Main Results:**
 - SRR measure can be a strong predictor of generalization, better than sharpness
 - SRR measure can be incorporated as regularization for improved performance
- **Takeaways**

Background



White-box Transformer **CRATE**
[Yu et al. NeurIPS 2023]

Sparse Rate Reduction (SRR):

$$\max_{\mathbf{Z} \in \mathbb{R}^{d \times N}} R(\mathbf{Z}) - R^c(\mathbf{Z}; \mathbf{U}) - \lambda \|\mathbf{Z}\|_0$$

$$R(\mathbf{Z}) \doteq \frac{1}{2} \log \det(\mathbf{I} + \frac{d}{N\epsilon^2} \mathbf{Z}^T \mathbf{Z})$$

$$R^c(\mathbf{Z}; \mathbf{U}) \doteq \sum_{k=1}^K R(\mathbf{U}_k^T \mathbf{Z})$$

Alternating Minimization:

$$\mathbf{Y}^\ell = \underbrace{\mathbf{Z}^{\ell-1} - \alpha \nabla R^c(\mathbf{Z}^{\ell-1}; \mathbf{U}^\ell)}_{\text{minimize } R^c(\mathbf{Z}; \mathbf{U})} \approx \mathbf{Z}^{\ell-1} + \alpha \gamma^2 \text{MSSA}(\mathbf{Z}^{\ell-1}; \mathbf{U}^\ell)$$

minimize $R^c(\mathbf{Z}; \mathbf{U})$

$$\mathbf{Z}^\ell = \underbrace{\text{ReLU}(\mathbf{Y}^\ell + \beta (\mathbf{D}^\ell)^T (\mathbf{Y}^\ell - \mathbf{D}^\ell \mathbf{Y}^\ell) - \beta \lambda \mathbf{1})}_{\text{minimize } \lambda \|\mathbf{Z}\|_0 - R(\mathbf{Z})}$$

minimize $\lambda \|\mathbf{Z}\|_0 - R(\mathbf{Z})$

$$\text{MSSA}(\mathbf{Z}; \mathbf{U}) = \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z} \text{softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z}))$$

$$= \underbrace{[\mathbf{U}_1, \dots, \mathbf{U}_K]}_{\mathbf{U}} \begin{bmatrix} \mathbf{U}_1^T \mathbf{Z} \text{softmax}((\mathbf{U}_1^T \mathbf{Z})^T (\mathbf{U}_1^T \mathbf{Z})) \\ \vdots \\ \mathbf{U}_K^T \mathbf{Z} \text{softmax}((\mathbf{U}_K^T \mathbf{Z})^T (\mathbf{U}_K^T \mathbf{Z})) \end{bmatrix}$$

Background

CRATE further introduces more parameters \mathbf{W} at the expense of interpretability.

$$\text{MSSA}(\mathbf{Z}; \mathbf{U}, \mathbf{W}) = \mathbf{W} \begin{bmatrix} \mathbf{U}_k^T \mathbf{Z} \text{ softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \\ \vdots \\ \mathbf{U}_k^T \mathbf{Z} \text{ softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \end{bmatrix}$$

To differentiate, we refer to the theoretically-driven framework as **CRATE-C(onceptual)**

$$\text{MSSA}(\mathbf{Z}; \mathbf{U}) = [\mathbf{U}_1, \dots, \mathbf{U}_K] \begin{bmatrix} \mathbf{U}_k^T \mathbf{Z} \text{ softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \\ \vdots \\ \mathbf{U}_k^T \mathbf{Z} \text{ softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \end{bmatrix}$$

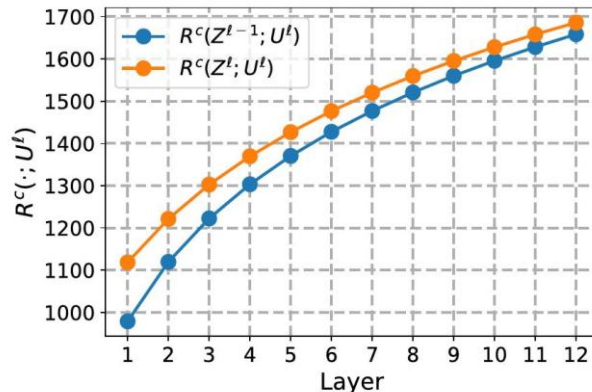
Pitfalls of Building CRATE-C

Given that MSSA operator is designed to minimize $R^c(Z; U)$, it is supposed to **decrease** monotonically as layer goes deeper.

$$\mathbf{Z}^\ell = \mathbf{Z}^{\ell-1} + \sum_{k=1}^K \mathbf{U}_k^\ell \mathbf{U}_k^{\ell T} \mathbf{Z}^{\ell-1} \text{softmax}((\mathbf{U}_k^{\ell T} \mathbf{Z}^{\ell-1})^T (\mathbf{U}_k^{\ell T} \mathbf{Z}^{\ell-1}))$$

Does the operation really achieve its design goal?

- ❖ **Empirically, No.** We can show in an isolated toy experiment that the update actually yields a counterproductive effect.

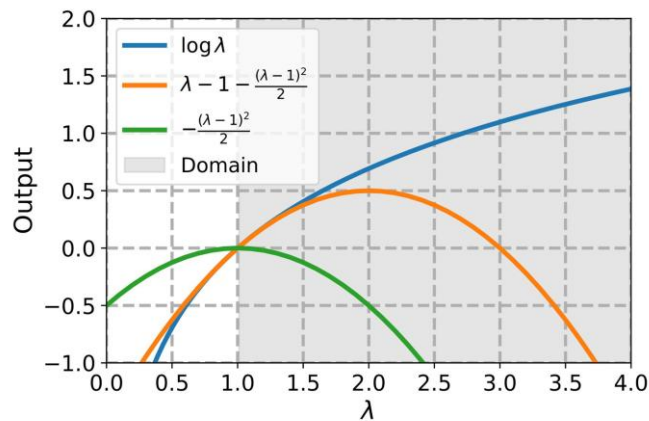


- ❖ **Theoretically, No.** We can also reveal this derivation artifacts from the eigenvalue perspective (see slides later).

Revisit and Interpret the Derivation

We first rewrite R^c with eigenvalues λ_i^k ($i = 1, \dots, N$) of $\mathbf{I} + \gamma(\mathbf{U}_k^T \mathbf{Z})^T \mathbf{U}_k^T \mathbf{Z}$. Note that $\lambda_i^k \geq 1$. Then we show R^c can be lower bounded by its Taylor expansions.

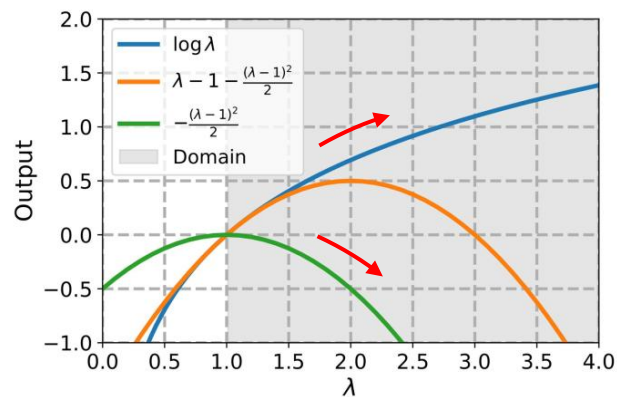
$$\begin{aligned} R^c(\mathbf{Z}; \mathbf{U}) &= \sum_{k=1}^K \sum_{i=1}^N \frac{1}{2} \log \lambda_i^k \geq \sum_{k=1}^K \sum_{i=1}^N \frac{1}{2} \left(\lambda_i^k - 1 - \frac{(\lambda_i^k - 1)^2}{2} \right) \\ &= \sum_{k=1}^K \left(\underbrace{\frac{\gamma}{2} \|\mathbf{U}_k^T \mathbf{Z}\|_F^2}_{\text{First-order term}} - \underbrace{\frac{\gamma^2}{4} \|(\mathbf{U}_k^T \mathbf{Z})^T \mathbf{U}_k^T \mathbf{Z}\|_F^2}_{\text{Second-order term}} \right) \end{aligned}$$



Revisit and Interpret the Derivation

MSSA operator with skip connection is constructed by performing an **approximation** of gradient descent on R^c .

$$\begin{aligned}
 \mathbf{Z} - \alpha \nabla_{\mathbf{Z}} R^c(\mathbf{Z}; \mathbf{U}) &= \mathbf{Z} - \alpha \gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z} \left(\mathbf{I} + \gamma (\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z}) \right)^{-1} \\
 &\approx \mathbf{Z} - \alpha \left(\underbrace{\gamma \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z}}_{\nabla \text{ of first-order term}} - \underbrace{\gamma^2 \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z} (\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})}_{\nabla \text{ of second-order term}} \right) \\
 &\approx \mathbf{Z} + \alpha \gamma^2 \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z} \text{softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})).
 \end{aligned}$$



Why this construction produces the opposite effect, i.e., increasing R^c ? Only utilizing the second-order term of its gradient !

Variants of CRATE

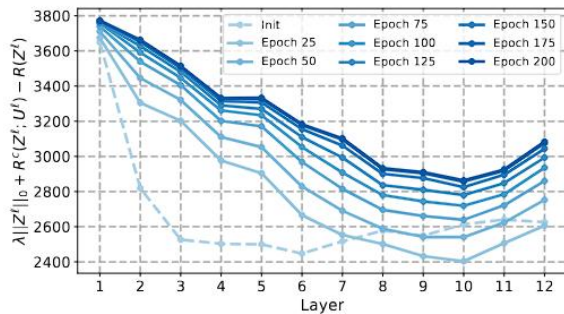
To implement the design purpose more faithfully, the sign before MSSA operator can be naturally reversed, performing ascent method. We name this framework **CRATE-N(egative)**.

$$\mathbf{Z} - \alpha\gamma^2 \sum_{k=1}^K \mathbf{U}_k \mathbf{U}_k^T \mathbf{Z} \text{softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z}))$$

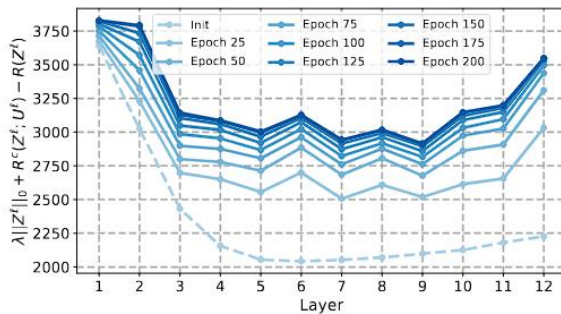
Can we further find a variant that performs competitively with CRATE without new parameters? Well, a simple transpose could do (see experiments later). We term this one **CRATE-T(ranspose)**.

$$\mathbf{Z} + \alpha\gamma^2 [\mathbf{U}_1, \dots, \mathbf{U}_K]^T \begin{bmatrix} \mathbf{U}_k^T \mathbf{Z} \text{softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \\ \vdots \\ \mathbf{U}_k^T \mathbf{Z} \text{softmax}((\mathbf{U}_k^T \mathbf{Z})^T (\mathbf{U}_k^T \mathbf{Z})) \end{bmatrix}$$

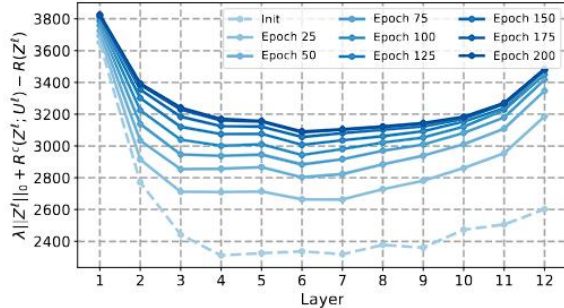
Behaviors of Sparse Rate Reduction



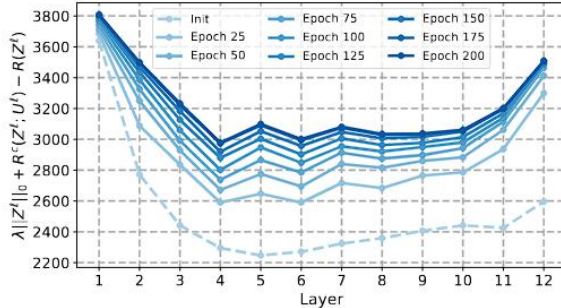
(a) Sparse rate reduction measure of CRATE-C (7)



(b) Sparse rate reduction measure of CRATE-N (8)



(c) Sparse rate reduction measure of CRATE-T (9)



(d) Sparse rate reduction measure of CRATE

Figure 2: Sparse rate reduction measure $\lambda\|\mathbf{Z}\|_0 + R^c(\mathbf{Z}; \mathbf{U}) - R(\mathbf{Z})$ of CRATE and its variants evaluated at different layers and epochs on CIFAR-10.

Whether Sparse Rate Reduction Benefits Generalization?

So far, we have partially confirmed the validity of different implementations of transformer-like models.

However, there are still some lingering questions

- **Whether** this SRR objective is beneficial or principled for these models to generalize ?
- If so, **how much** is the benefit ?

Whether Sparse Rate Reduction Benefits Generalization?

We will explore its causal relationship to the generalization and adopt SRR as an empirical predictor of generalization (measure of complexity).

$$\mu_{\text{SRR}}(\mathbf{w}; \mathbf{Z}) = \frac{1}{L} \sum_{\ell=1}^L \mu_{\text{SRR}}^{\ell}(\mathbf{w}^{\ell}; \mathbf{Z}^{\ell}) = \frac{1}{L} \sum_{\ell=1}^L \left(\lambda \|\mathbf{Z}^{\ell}\|_0 + R^c(\mathbf{Z}^{\ell}; \mathbf{U}^{\ell}) - R(\mathbf{Z}^{\ell}) \right)$$

Typically, a good measure should have the property where lower complexity should indicate smaller generalization gap. For example, the following is true if the measure μ is described by generalization bound.

$$L_{\text{test}} - \hat{L}_{\text{train}} \leq \sqrt{\frac{\mu}{m}}$$

Correlation Analysis

Similar to [Jiang et al. ICLR 2020], we collected a set of models with varied hyperparameters trained until convergence and evaluated how well the generalization gap correlates with the measure.

$$\mathcal{T} \triangleq \cup_{\theta \in \Theta_1 \times \dots \times \Theta_n} \{(\mu(\theta), g(\theta))\}$$

Table 3: Choices of hyperparameters.

Hyperparameters	Choices
batch size	{64, 128}
initial learning rate	$\{2 \times 10^{-5}, 1 \times 10^{-4}\}$
width	{384, 768}
dropout	{0.0, 0.1}
model type	{CRATE, CRATE-C, CRATE-N, CRATE-T}

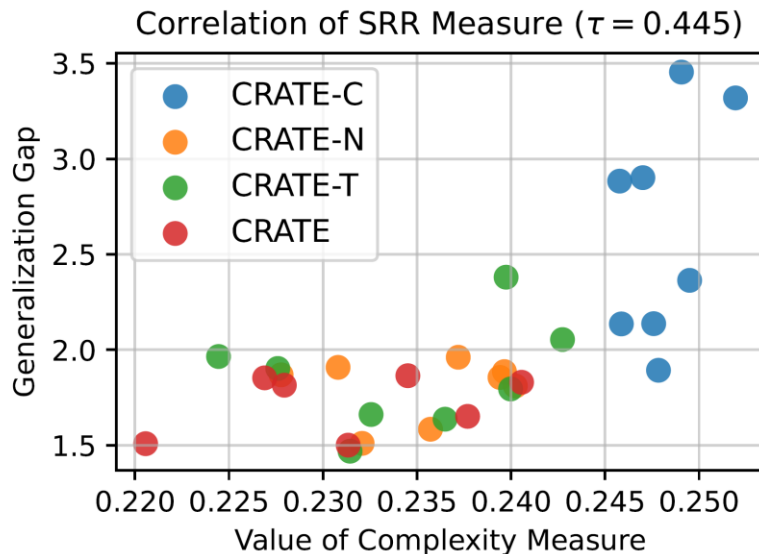
Kendall's rank-correlation coefficient: Range in $[-1, 1]$. The close to one, the stronger the positive correlation.

$$\tau(\mathcal{T}) \triangleq \frac{1}{|\mathcal{T}|(|\mathcal{T}| - 1)} \sum_{(\mu_1, g_1) \in \mathcal{T}} \sum_{(\mu_2, g_2) \in \mathcal{T} \setminus (\mu_1, g_1)} \text{sign}(\mu_1 - \mu_2) \text{sign}(g_1 - g_2)$$

Correlation Analysis Results

Table 1: Correlation of complexity measures with generalization gap (width $d = 384$).

Complexity measures	Batch size	Learning rate	Dropout	Model type	Overall τ	Ψ
ℓ_2 -norm	0.200	-0.333	-0.333	-0.429	-0.363	-0.224
ℓ_2 -norm-init	0.200	-0.200	-0.333	-0.286	-0.290	-0.158
# params	0.000	0.000	0.000	-0.572	-0.351	-0.143
1/margin	-0.067	0.467	0.467	0.238	0.415	0.276
sum-of-spec	0.200	-0.333	-0.467	-0.381	-0.290	-0.245
prod-of-spec	0.200	-0.333	-0.467	-0.476	-0.338	-0.269
sum-of-spec/margin	0.333	-0.333	-0.467	-0.048	-0.230	-0.129
prod-of-spec/margin	0.333	-0.333	-0.467	-0.143	-0.260	-0.152
fro/spec	-0.200	0.333	0.467	-0.476	0.019	0.031
spec-init-main	0.333	-0.333	-0.467	-0.190	-0.273	-0.164
spec-orig-main	0.200	-0.333	-0.467	-0.095	-0.252	-0.174
sum-of-fro	0.200	-0.333	-0.333	-0.381	-0.325	-0.212
prod-of-fro	0.200	-0.333	-0.333	-0.429	-0.372	-0.224
sum-of-fro/margin	0.333	-0.200	-0.467	-0.048	-0.217	-0.095
prod-of-fro/margin	0.333	-0.200	-0.467	-0.143	-0.247	-0.119
fro-distance	0.200	-0.200	-0.333	-0.286	-0.290	-0.155
spec-distance	0.200	-0.200	-0.333	-0.286	-0.290	-0.155
param-norm	0.200	-0.333	-0.333	-0.429	-0.363	-0.224
path-norm	0.333	-0.600	-0.467	-0.286	-0.191	-0.255
pac-bayes-init	0.200	0.200	-0.600	0.238	0.015	-0.009
pac-bayes-orig	-0.200	0.333	0.467	0.381	0.333	0.245
1/ σ pac-bayes-flatness	-0.267	0.333	0.333	0.455	0.333	0.213
SRR	-0.067	0.467	0.333	0.714	0.445	0.362



Better predictive power than widely investigated flatness-based measure

Sparse Rate Reduction as Regularization

Since SRR measure enjoys a strong correlation to generalization, it is reasonable to incorporate it during training and optimize it with task-specific loss simultaneously, similar to sharpness-aware minimization [Foret et al. ICLR 2021] for improved generalization.

Well, the most straightforward way is through regularization.

$$\min_{\mathbf{w}} \mathcal{L}_{\text{ce}}(\mathbf{w}) + \eta \cdot \frac{1}{L} \sum_{\ell=1}^L \mu_{\text{SRR}}^{\ell}(\mathbf{w}^{\ell}; \mathbf{Z}_{\text{StopGrad}}^{\ell})$$

Sparse Rate Reduction as Regularization Results

An efficient implementation at the last layer can already give consistent performance gain.

Table 2: Top-1 accuracy for CRATE and its variants trained with or without SRR regularization on CIFAR-10/100 from scratch (width $d = 384$).

Models	CIFAR-10		CIFAR-100	
	cross-entropy	+ SRR regularization (L=12)	cross-entropy	+ SRR regularization (L=12)
CRATE-C	76.87	77.61	43.40	44.53
CRATE-N	81.52	81.91	55.11	55.62
CRATE-T	85.49	85.52	60.59	60.69
CRATE	86.67	86.79	62.40	62.52

In fact, we also find that imposing regularization on first few layers performs best. Only as a proof-of-concept for scalable depth here.

Table 5: Top-1 accuracy for CRATE and its variants trained with efficient implementations of SRR regularization on CIFAR-10 from scratch (width $d = 384$).

Training methods	CIFAR-10			
	CRATE-C	CRATE-N	CRATE-T	CRATE
cross-entropy	76.87	81.52	85.49	86.67
+ Layer 2 reg	77.75	82.41	85.84	87.03
+ Layer 4 reg	77.95	81.57	85.46	87.03
+ Layer 6 reg	77.48	80.83	85.22	87.02
+ Layer 8 reg	77.04	81.29	85.12	86.64
+ Layer 10 reg	77.44	81.19	85.68	86.67
+ Layer 12 reg	77.61	81.91	85.52	86.79
+ Random layer reg	75.19	79.66	84.27	85.36

Takeaways

- SRR objective can be viewed as an energy function that is optimized in the forward pass of transformer-like models.
 - It *almost* monotonically decreases and hovers around the stationary point.
 - Behaviors persist across varied implementations.
- SRR objective could be a choice to design transformer-like models, but not necessarily principled.
 - We demonstrate its *positive and strong (by comparison)* correlation with generalization.
 - More faithful instantiation does not necessarily give a better model.