

Memory-Statistics Tradeoff in Continual Learning with Structural Regularization

Haoran Li¹, Jingfeng Wu², Vladimir Braverman³
¹Shenzhen University, ²UC Berkeley, ³Johns Hopkins University

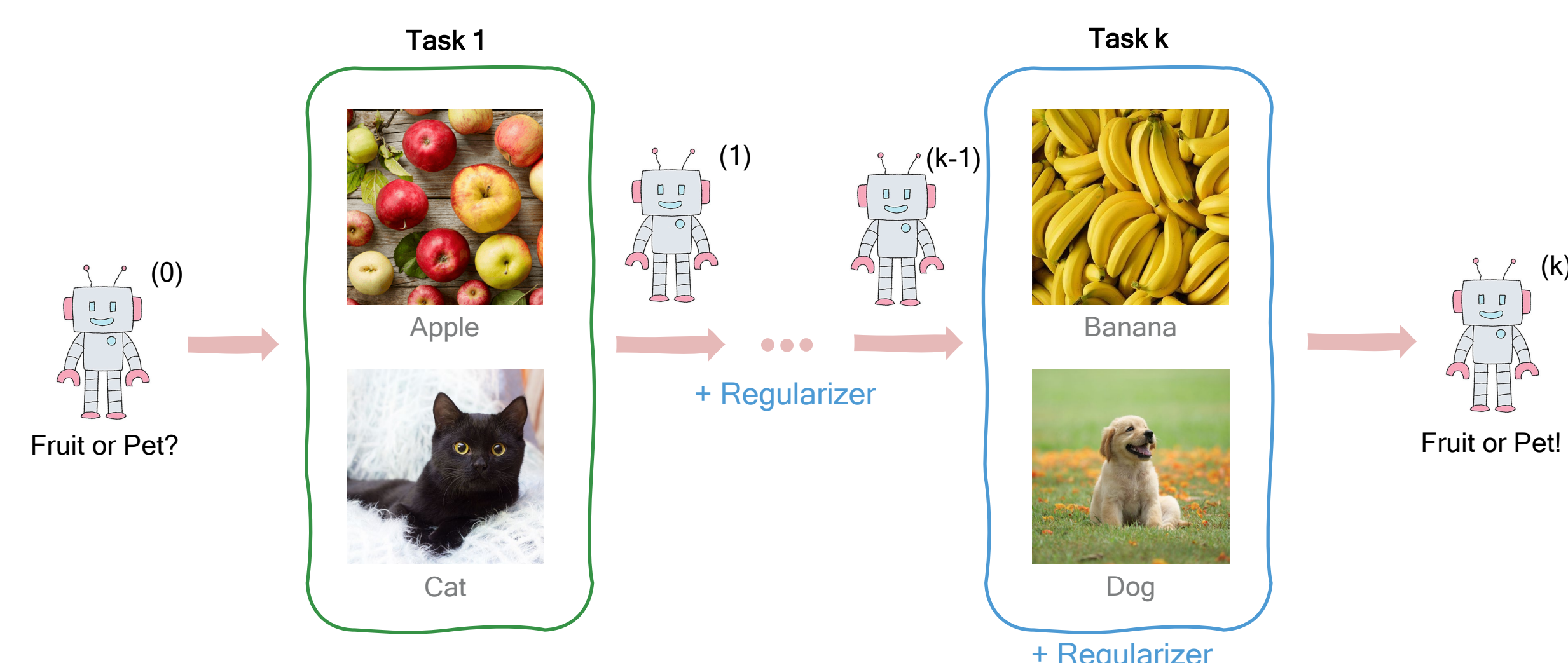


Regularized-Based CL

Sequentially learn tasks $(\mathbf{x}_t^{(1)}, y_t^{(1)})_{t=1}^n \sim \mathcal{D}^{(1)}, (\mathbf{x}_t^{(2)}, y_t^{(2)})_{t=1}^n \sim \mathcal{D}^{(2)}$:

Minimize $L_1(\mathbf{w}) = R_1(\mathbf{w})$;

Minimize $L_2(\mathbf{w}) = R_2(\mathbf{w}) + L_{\text{regu}}(\mathbf{w}, \mathbf{w}^{(1)}; \Sigma)$



Problem Formulation

Two-task Linear Regression in Random Design

For tasks $k = 1, 2$,

• $\mathbf{x}^{(k)}$ sampled from distributions $\mathcal{D}^{(k)}$ (one-hot/Gaussian)

• Shared optimal $\mathbf{y}^{(k)} = \mathbf{x}^{(k)\top} \mathbf{w}^* + \mathcal{N}(0, \sigma^2)$

• Covariance matrices $\mathbf{G} := \mathbb{E} \mathbf{x}^{(1)\top} \mathbf{x}^{(1)}, \mathbf{H} := \mathbb{E} \mathbf{x}^{(2)\top} \mathbf{x}^{(2)}$

• Risk for task k $R_k(\mathbf{w}) := \mathbb{E} \|\mathbf{y}^{(k)} - \mathbf{x}^{(k)} \mathbf{w}\|_2^2$

• Joint population risk $R(\mathbf{w}) := R_1(\mathbf{w}) + R_2(\mathbf{w})$

• One-hot setting: $\mathbb{P}(\mathbf{x}^{(1)} = \mathbf{e}_i) = \mu_i, \mathbb{P}(\mathbf{x}^{(2)} = \mathbf{e}_i) = \lambda_i$

• Gaussian setting: $\mathbf{x}^{(1)} \sim \mathcal{N}(0, \mathbf{G}), \mathbf{x}^{(2)} \sim \mathcal{N}(0, \mathbf{H})$,

Generalized ℓ_2 -Regularized CL (GRCL)

$$\mathbf{w}^{(1)} = (\mathbf{X}^{(1)\top} \mathbf{X}^{(1)})^{-1} \mathbf{X}^{(1)\top} \mathbf{y}^{(1)}$$

$$\text{GRCL: } \mathbf{w}^{(2)} = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{y}^{(2)} - \mathbf{X}^{(2)} \mathbf{w}\|_2^2 + \|\mathbf{w} - \mathbf{w}^{(1)}\|_{\Sigma}^2$$

$$\cdot \ell_2\text{-RCL } (\Sigma \rightarrow \gamma \mathbf{I}): \mathbf{w}^{(2)} = \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{y}^{(2)} - \mathbf{X}^{(2)} \mathbf{w}\|_2^2 + \gamma \|\mathbf{w} - \mathbf{w}^{(1)}\|_2^2$$

$$\cdot \text{OCL } (\Sigma \rightarrow 0): \mathbf{w}^{(2)} = \mathbf{w}^{(1)} + (\mathbf{X}^{(2)\top} \mathbf{X}^{(2)})^{-1} \mathbf{X}^{(2)\top} (\mathbf{y}^{(2)} - \mathbf{X}^{(2)} \mathbf{w}^{(1)}),$$

$$\cdot \text{Joint Learning (JL): } \mathbf{w}_{\text{joint}} = (\mathbf{X}_{\text{joint}}^{\top} \mathbf{X}_{\text{joint}})^{-1} \mathbf{X}_{\text{joint}}^{\top} \mathbf{y}_{\text{joint}}$$

Main Result

One-hot Setting

$$\mathbb{E}[R(\mathbf{w}^{(2)})] - \min R(\cdot) = \text{bias} + \text{variance}$$

where

$$\begin{aligned} \text{bias} &\approx \langle (\mathbf{G} + \mathbf{H})(\mathbf{I} - \mathbf{G})^n (\Sigma^2 (\Sigma + \mathbf{H})^{-2} + (\mathbf{I} - \mathbf{H})^n), \mathbf{w}^* \mathbf{w}^{*\top} \rangle, \\ \text{variance} &\approx \sigma^2 \langle (\mathbf{G} + \mathbf{H})(\Sigma^2 (\Sigma + \mathbf{H})^{-2} + (\mathbf{I} - \mathbf{H})^n), \frac{1}{n} \mathbf{G}_{\mathbb{J}^c}^{-1} + n \mathbf{G}_{\mathbb{J}^c} \rangle \\ &\quad + \langle \mathbf{G} + \mathbf{H}, \frac{1}{n} (\mathbf{H}_{\mathbb{K}} + \Sigma_{\mathbb{K}})^{-2} \mathbf{H}_{\mathbb{K}} + n (\mathbf{I}_{\mathbb{K}^c} + n \Sigma_{\mathbb{K}^c})^{-2} \mathbf{H}_{\mathbb{K}^c} \rangle. \end{aligned}$$

[Gaussian Setting] For OCL, bias, variance satisfy

$$\begin{aligned} \text{bias} &\gtrsim \left\| \left(\frac{\text{tr} \mathbf{G}_{\mathbb{J}^c}}{n^2} \mathbf{G}_{\mathbb{J}^c}^{-2} + \mathbf{I}_{\mathbb{J}^c} \right)^{\frac{1}{2}} \cdot \left(\frac{\text{tr} \mathbf{H}_{\mathbb{K}^c}}{n^2} \mathbf{H}_{\mathbb{K}^c}^{-2} + \mathbf{I}_{\mathbb{K}^c} \right)^{\frac{1}{2}} \mathbf{w}^* \right\|_{\mathbf{G}}, \\ \text{variance} &\gtrsim \frac{\sigma^2}{n} \left\langle \mathbf{G}, \mathbf{H}_{\mathbb{K}}^{-1} + \frac{n^2}{(\text{tr} \mathbf{H}_{\mathbb{K}^c})^2} \mathbf{H}_{\mathbb{K}^c} + \left(\frac{\text{tr} \mathbf{H}_{\mathbb{K}^c}}{n^2} \mathbf{H}_{\mathbb{K}^c}^{-2} + \mathbf{I}_{\mathbb{K}^c} \right) \cdot \left(\mathbf{G}_{\mathbb{J}^c}^{-1} + \frac{n^2}{(\text{tr} \mathbf{G}_{\mathbb{J}^c})^2} \mathbf{G}_{\mathbb{J}^c} \right) \right\rangle. \end{aligned}$$

[Extension to NTK regime] bias, variance satisfy similar results where

$$\mathbf{G} := \mathbb{E}_{\mathcal{D}^{(1)}} [\phi(\mathbf{x}^{(1)}) \phi(\mathbf{x}^{(1)})^{\top}], \mathbf{H} := \mathbb{E}_{\mathcal{D}^{(2)}} [\phi(\mathbf{x}^{(2)}) \phi(\mathbf{x}^{(2)})^{\top}].$$

Catastrophic Forgetting of OCL and ℓ_2 -RCL

For tasks solvable by Joint Learning,

to achieve $o(1)$ excess risk, for $\mathbf{G} = \text{diag}(\mu_i)_{i=1}^d, \mathbf{H} = \text{diag}(\lambda_i)_{i=1}^d$:

$$\cdot \text{OCL: } \sum_{i \in \mathbb{K}} \frac{\mu_i}{\lambda_i} + n^2 \sum_{i \in \mathbb{J} \cap \mathbb{K}^c} \mu_i \lambda_i = o(n)$$

$$\cdot \ell_2\text{-RCL: } \gamma = o(1), \sum_{i \in \mathbb{J} \cup \mathbb{K}} \left(\frac{\mu_i}{\lambda_i + \frac{1}{n} + \gamma} + \frac{\gamma}{\mu_i + \frac{1}{n}} \right) = o(n)$$

Large forgetting with dissimilar covariances \mathbf{G}, \mathbf{H} !

GRCL avoids Forgetting with Full Memory

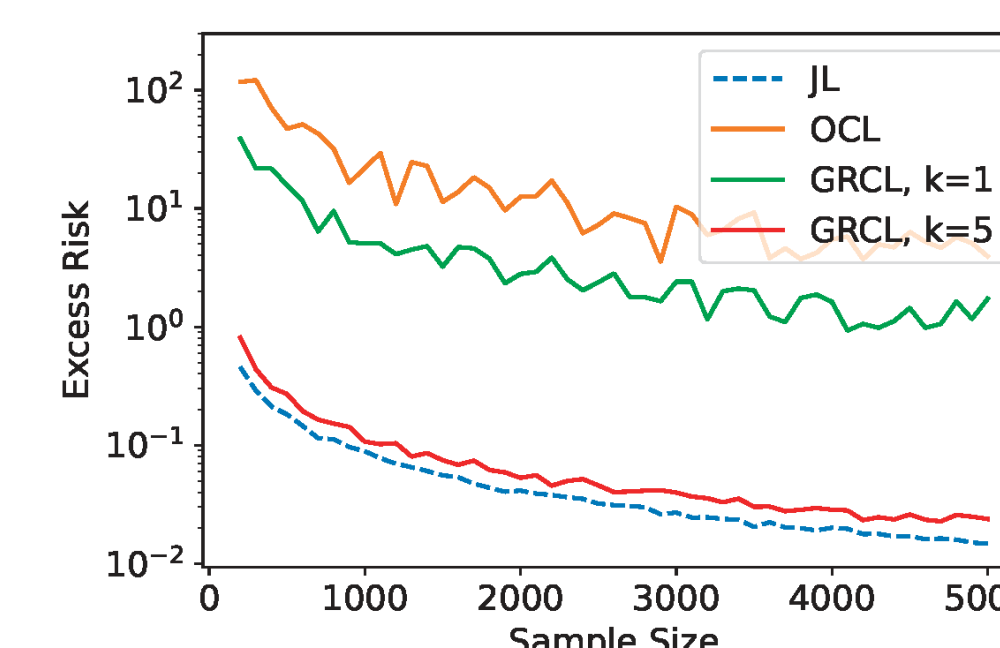
For regularization $\Sigma = \text{diag}(\gamma_i)_{i=1}^d$ aligned with \mathbf{G} ($\gamma_i = \mu_i$ for $\mu_i \geq 1/n$ and 0 otherwise),

$$\mathbb{E} \Delta(\mathbf{w}^{(2)}) \lesssim \mathbb{E} \Delta(\mathbf{w}_{\text{joint}}).$$

GRCL attain joint learning performance with sufficient memory.

[Extension to multi-task] $\mathbb{E} \Delta(\mathbf{w}^{(2)}) \lesssim \mathbb{E} \Delta(\mathbf{w}_{\text{joint}})$ holds when:

$$\gamma_i^{(t)} = \mu_i^{(t)} \text{ for } \mu_i^{(t)} \geq 1/n \text{ and } 0 \text{ otherwise.}$$



Memory-Statistics Trade-off

Catastrophic Forgetting in Low-Memory Settings

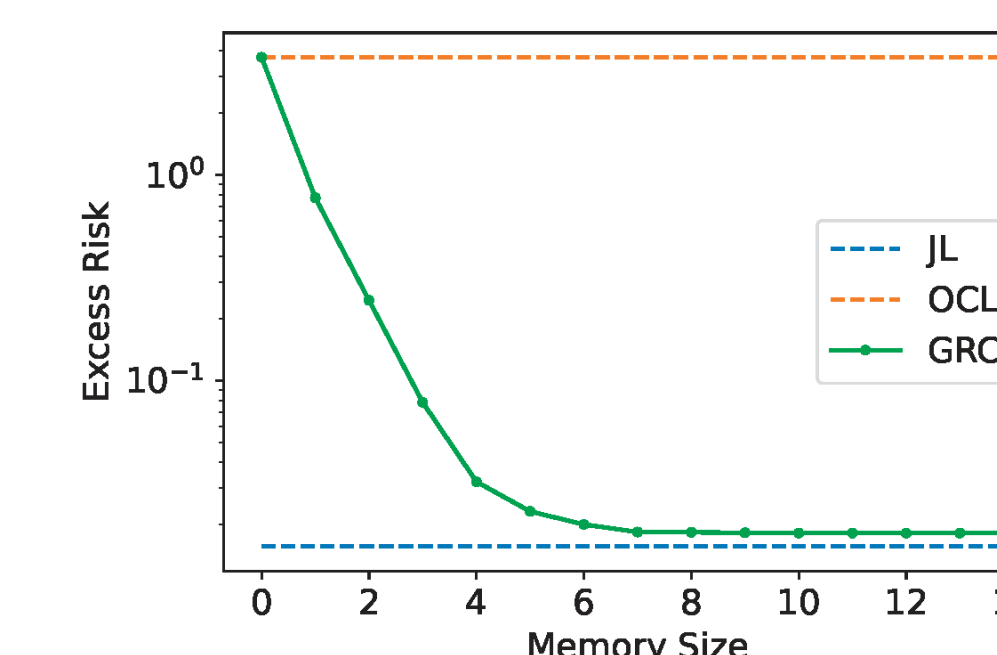
For any memory size $k \ll N$,

if $\mu_i = 1/(k+1)$ and $\lambda_i = 1/n$ for $1 \leq i \leq k+1$, then:

For any regularizer Σ of rank k , the GRCL excess risk is $\Omega(1)$.

No k -dimensional regularization in GRCL can mitigate forgetting across $k+1$ orthogonal directions.

Tempered Forgetting in Moderate Memory Settings

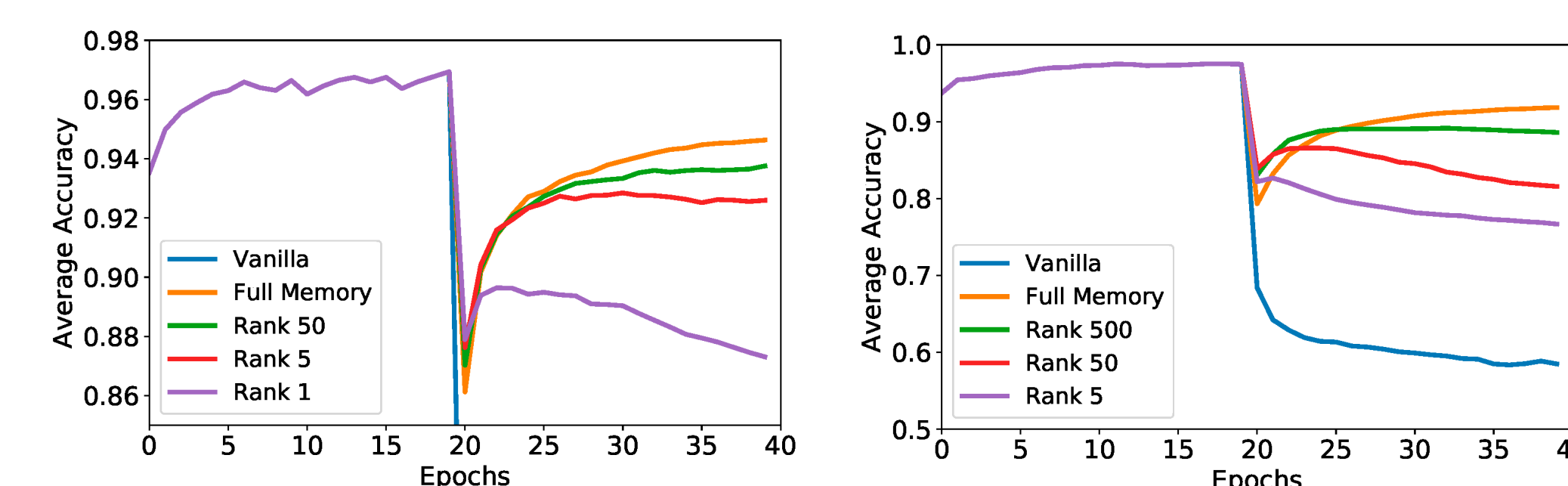


As the memory size k of GRCL increases, GRCL's excess risk decreases, approaching joint learning performance at a threshold.

GRCL with an intermediate memory constraint can partially alleviate catastrophic forgetting.

Neural Network Verification

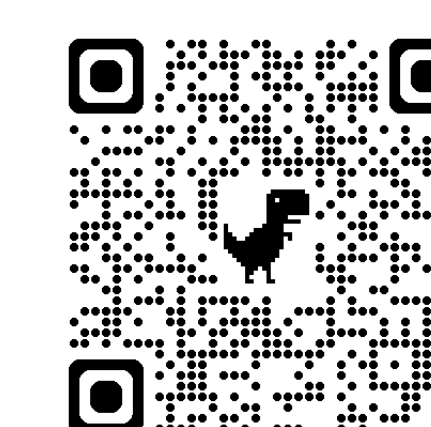
GRCL via PCA on two-task Permuted and Rotated MNIST:



(a) Permuted MNIST

(b) Rotated MNIST

Practical Algorithms



GRCL via PCA costly on NN.

Solution: GRCL via CountSketch!

← See our previous work

This work on ArXiv:

