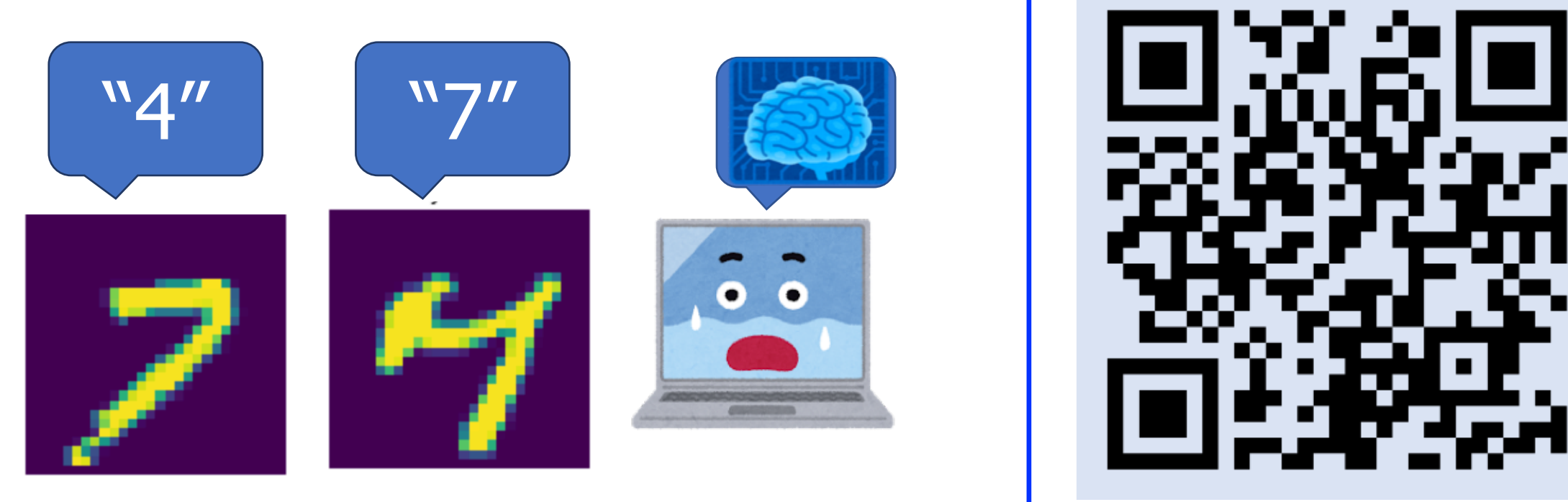


Data Cleansing for Models Trained with SGD

Satoshi Hara (Osaka Univ.), Atsushi Nitanda (Tokyo Univ./RIKEN AIP), Takanori Maehara (RIKEN AIP)

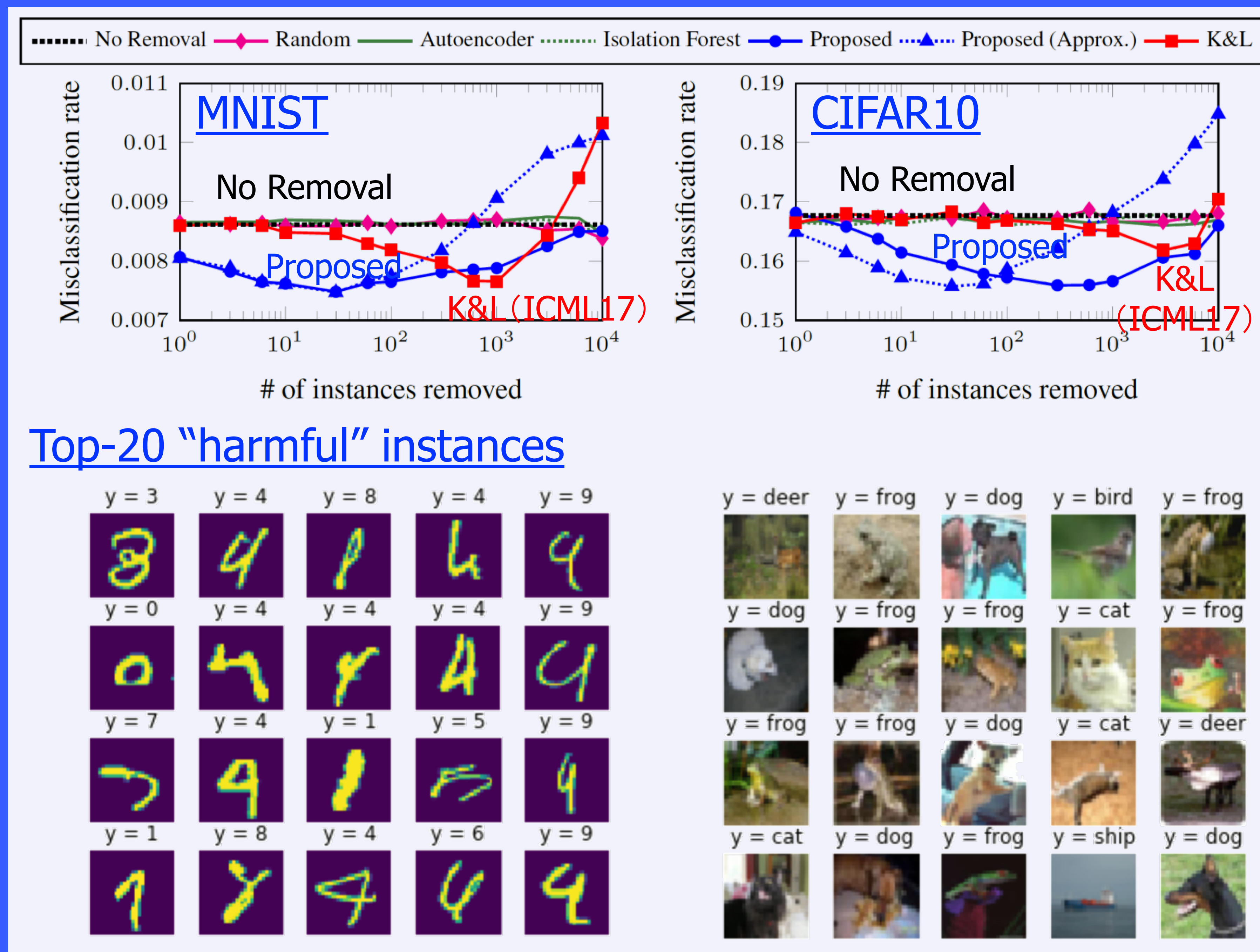
Data Cleansing

Remove "harmful" training instances, and improve the model's accuracy.



The Results on MNIST/CIFAR10

By removing "harmful" instances, the trained CNN became more accurate.



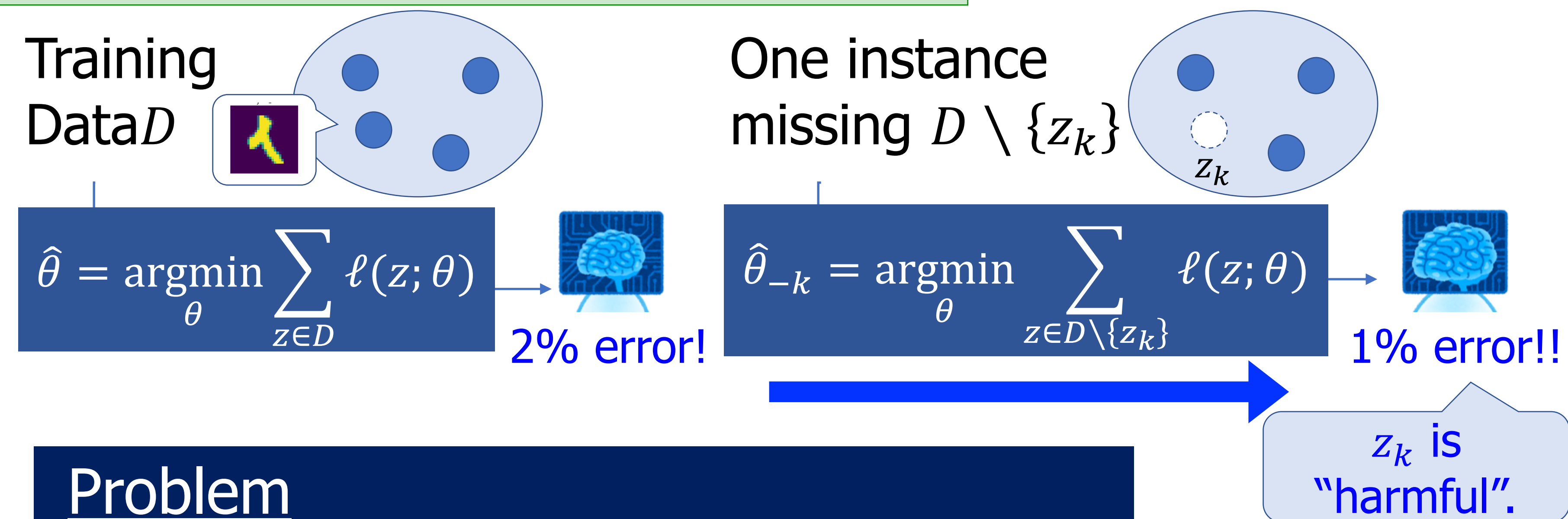
Algorithm

for $k = 1, 2, \dots, N$
 Compute $r_k = \langle u, \Delta\theta_{-k} \rangle$
 Identify "harmful" instance:
 $\hat{k} = \underset{k}{\operatorname{argmin}} r_k$
 Remove "harmful" instance:
 $D \leftarrow D \setminus \{z_{\hat{k}}\}$

estimated validation loss

The Proposed Estimator

Estimation of "harmful" instances



Problem

Find z_k that minimizes the validation loss.

$$\sum_{z \in D_V} (\ell(z; \hat{\theta}_{-k}) - \ell(z; \hat{\theta})) \approx \langle u, \hat{\theta}_{-k} - \hat{\theta} \rangle$$

- $u := \sum_{z \in D_V} \nabla_{\theta} \ell(z; \hat{\theta})$ for the val. set D_V .
- No model retraining.

Small validation loss \approx Small inner product

Our Contributions

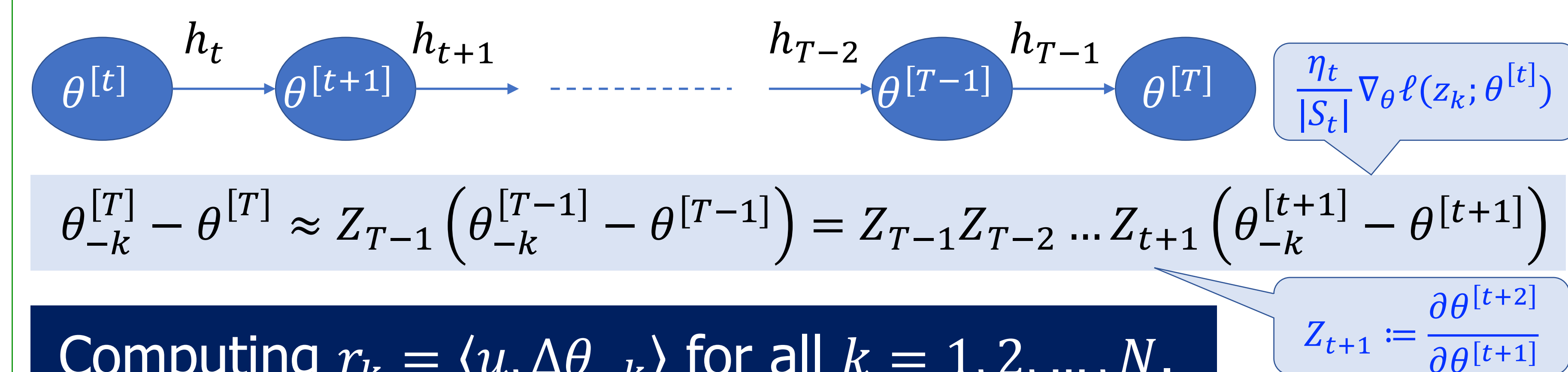
- To estimate $\hat{\theta}_{-k} - \hat{\theta}$, the current methods require "the optimal $\hat{\theta}$ " and "convex loss ℓ ".
 → Too restrictive for modern ML.
- Our estimator assumes that "the model is trained with SGD".
 → More appropriate for modern ML.

The Proposed Estimator

- $\theta^{[t+1]} \leftarrow \theta^{[t]} - \frac{\eta_t}{|S_t|} \sum_{z \in S_t} \nabla_{\theta} \ell(z; \theta^{[t]}) =: h_t(\theta^{[t]})$ Ordinary SGD
- $\theta_{-k}^{[t+1]} \leftarrow \theta_{-k}^{[t]} - \frac{\eta_t}{|S_t|} \sum_{z \in S_t \setminus \{z_k\}} \nabla_{\theta} \ell(z; \theta^{[t]})$ SGD with missing z_k

Our Estimator $\Delta\theta_{-k} \approx \theta_{-k}^{[T]} - \theta^{[T]}$

"Backprop" through SGD



Computing $r_k = \langle u, \Delta\theta_{-k} \rangle$ for all $k = 1, 2, \dots, N$.

for $t = T-1, T-2, \dots, 1$

$$r_j \leftarrow r_j + \left\langle u, \frac{\eta_t}{|S_t|} \nabla_{\theta} \ell(z_j; \theta^{[t]}) \right\rangle, \forall j \in S_t$$

$$u \leftarrow Z_t u$$

$O(|S_t|T)$ time in total

Theoretical Analysis

Convex Case

Our Estimator

$$\|(\theta_{-k}^{[T]} - \theta^{[T]}) - \Delta\theta_{-k}\| \leq \sqrt{2(h_k(\lambda)^2 + h_k(\Lambda)^2)}$$

Assumption

- The loss ℓ is strongly convex, smooth, and twice differentiable.
- $\exists \lambda, \Lambda > 0$ such that $\lambda I \preceq \nabla_{\theta}^2 \ell(z; \theta) \preceq \Lambda I$.

$$h_k(a) := \frac{\eta_{\pi(k)}}{|S_{\pi(k)}|} \prod_{t=\pi(k)+1}^{T-1} (1 - \eta_t a) \|\nabla_{\theta} \ell(z_k; \theta^{[\pi(k)]})\|$$

Non-Convex Case

Our Estimator

$$\|(\theta_{-k}^{[T]} - \theta^{[T]}) - \Delta\theta_{-k}\| \leq \frac{\gamma^2 T G^2 L}{\Lambda} \exp^{O(\gamma \Lambda \sqrt{T})}$$

Assumption

- The Hessian of the loss ℓ is Lipschitz.
- $\|\nabla_{\theta} \ell(z; \theta)\| \leq G, \nabla_{\theta}^2 \ell(z; \theta) \preceq \Lambda I$.
- $\eta = O(\gamma/\sqrt{T})$.