



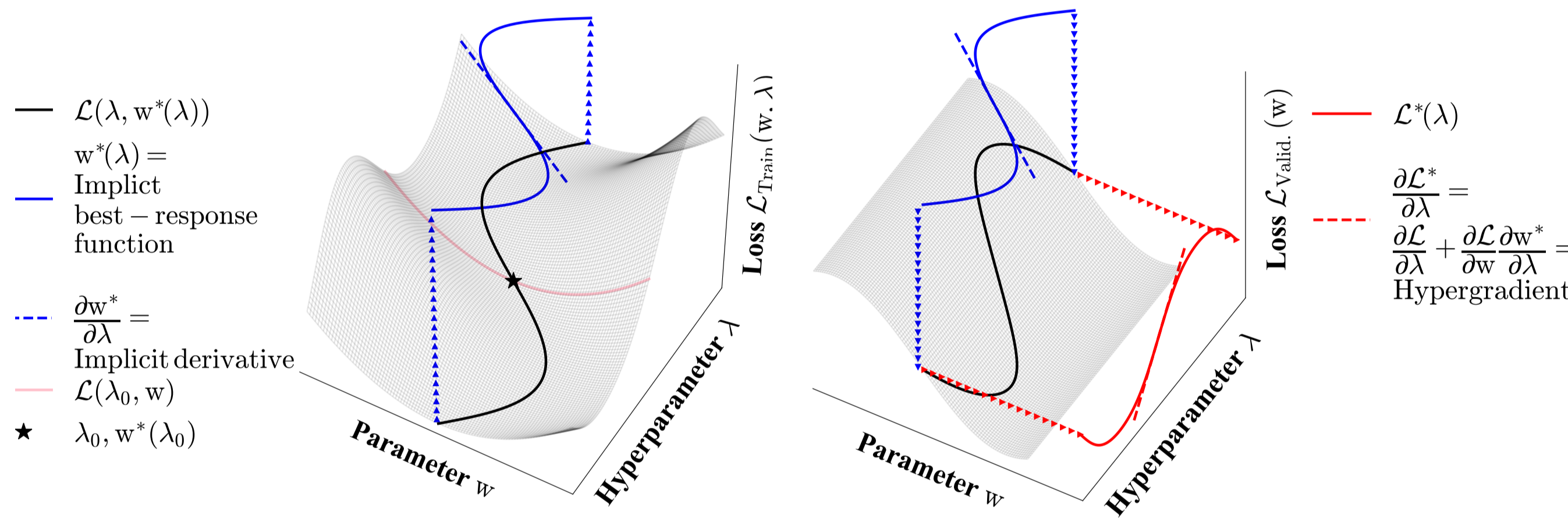
Motivation and Contributions

- Tuning hyperparameters can be critical for generalization, but **difficult when hyperparameters are high-dimensional**.
- If we can tune as many hyperparameters as weights, we have a new paradigm of learned, flexible regularization.
 - Learned data augmentation
 - Data distillation
 - Tuning millions of regularization hyperparameters
- We **scale** gradient-based hyperparameter optimization to high dimensions by combining **implicit differentiation with efficient & stable inverse-Hessian approximations**.

Hyperparameter Optimization is Nested Optimization

$$\lambda^* := \arg \min_{\lambda} \mathcal{L}_V^*(\lambda)$$

where $\mathcal{L}_V^*(\lambda) := \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))$ and $\mathbf{w}^*(\lambda) := \arg \min_{\mathbf{w}} \mathcal{L}_T(\lambda, \mathbf{w})$



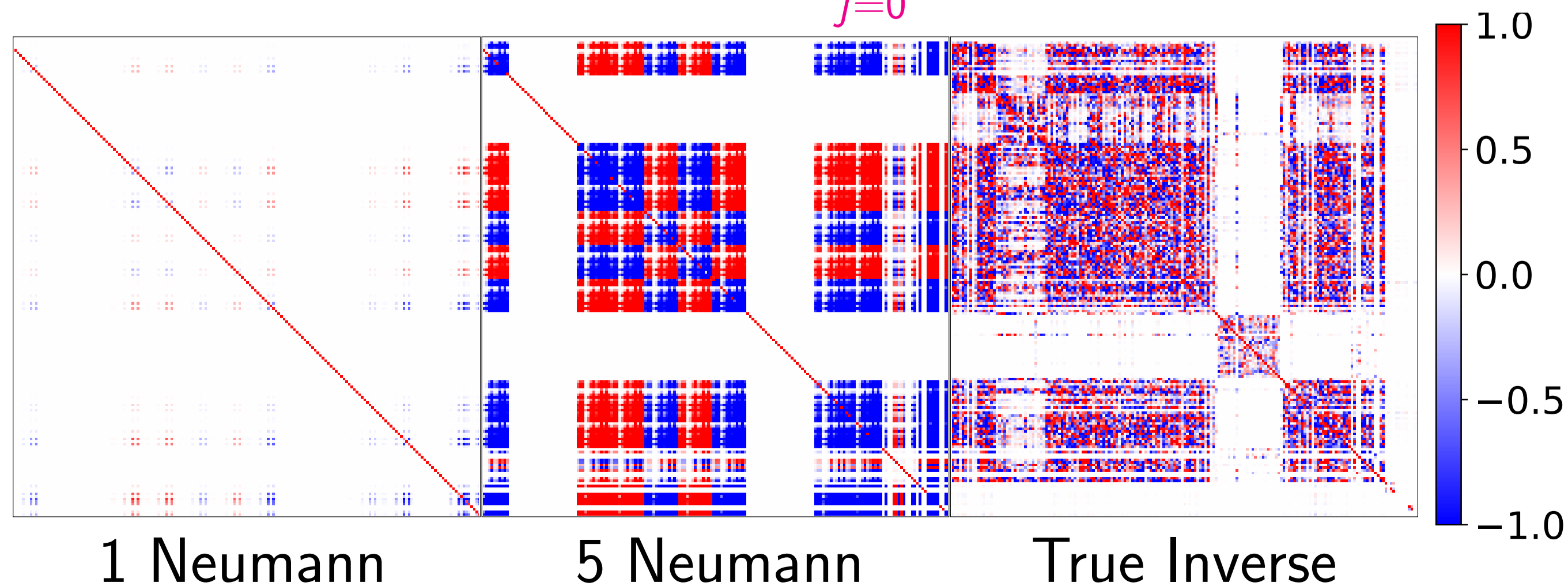
Hypergradients with Implicit Differentiation

$$\frac{\partial \mathcal{L}_V^*(\lambda)}{\partial \lambda} = \underbrace{\frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \lambda}}_{\text{hyperparam direct grad.}} + \underbrace{\frac{\partial \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda))}{\partial \mathbf{w}(\lambda)}}_{\text{parameter direct grad.}} \times \underbrace{\frac{\partial \mathbf{w}^*(\lambda)}{\partial \lambda}}_{\text{best-response Jacobian}}$$

$$\frac{\partial \mathbf{w}^*}{\partial \lambda} \Big|_{\lambda'} = - \underbrace{\left[\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1}}_{\text{training Hessian}} \times \underbrace{\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \lambda}}_{\text{training mixed partials}} \Big|_{\lambda', \mathbf{w}^*(\lambda')} \quad (\text{implicit differentiation})$$

Inverse-Hessian approximations via Neumann series

$$\left[\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1} = \lim_{i \rightarrow \infty} \sum_{j=0}^i \left[I - \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^j$$



Calculating Hypergradients Efficiently

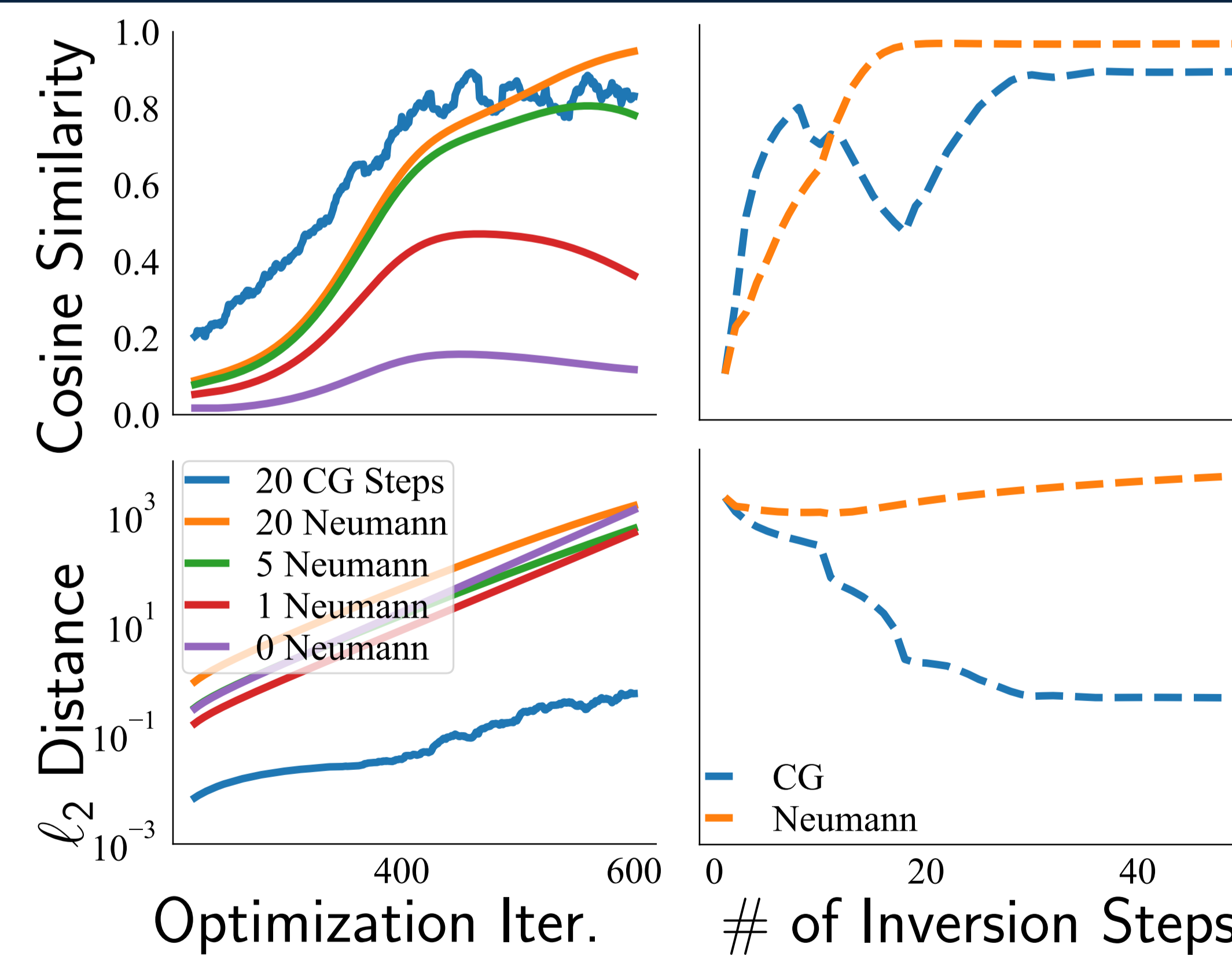
- while** not converged **do**
- for** $k = 1 \dots N$ **do**
- $\mathbf{w}' \leftarrow \frac{\partial \mathcal{L}_T}{\partial \mathbf{w}} \Big|_{\lambda', \mathbf{w}'}$
- $\lambda' \leftarrow \text{hypergradient}(\mathcal{L}_V, \mathcal{L}_T, \lambda', \mathbf{w}')$
- return** λ', \mathbf{w}'

$$\frac{\partial \mathcal{L}_V^*}{\partial \lambda} = \frac{\partial \mathcal{L}_V}{\partial \lambda} + \frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \underbrace{\left[\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1}}_{\text{vector-inverse Hessian product}} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \lambda}$$

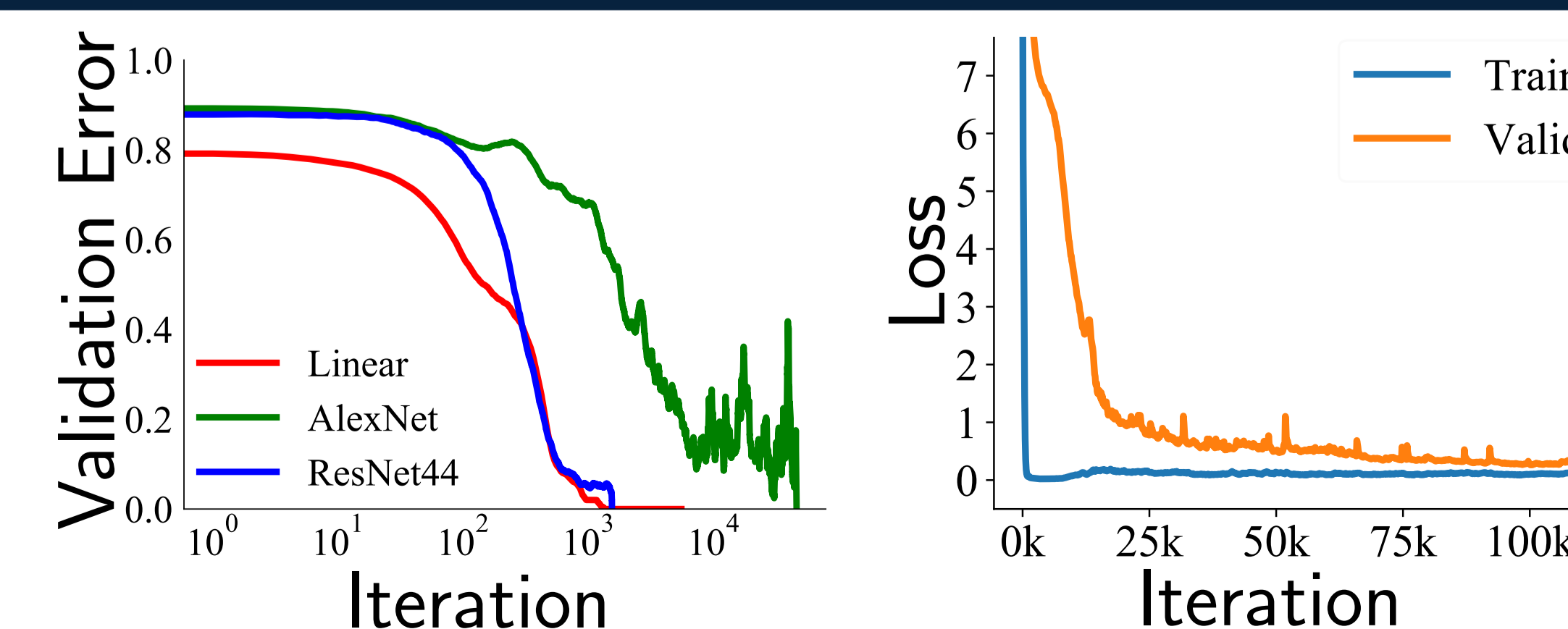
$$= \frac{\partial \mathcal{L}_V}{\partial \lambda} + \underbrace{\frac{\partial \mathcal{L}_V}{\partial \mathbf{w}} \times \left[\frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \mathbf{w}} \right]^{-1}}_{\text{vector-Jacobian product}} \frac{\partial^2 \mathcal{L}_T}{\partial \mathbf{w} \partial \lambda}$$

The entire hypergradient computation can be performed efficiently using **vector-Jacobian products**, given an inverse-Hessian-vector approximation.

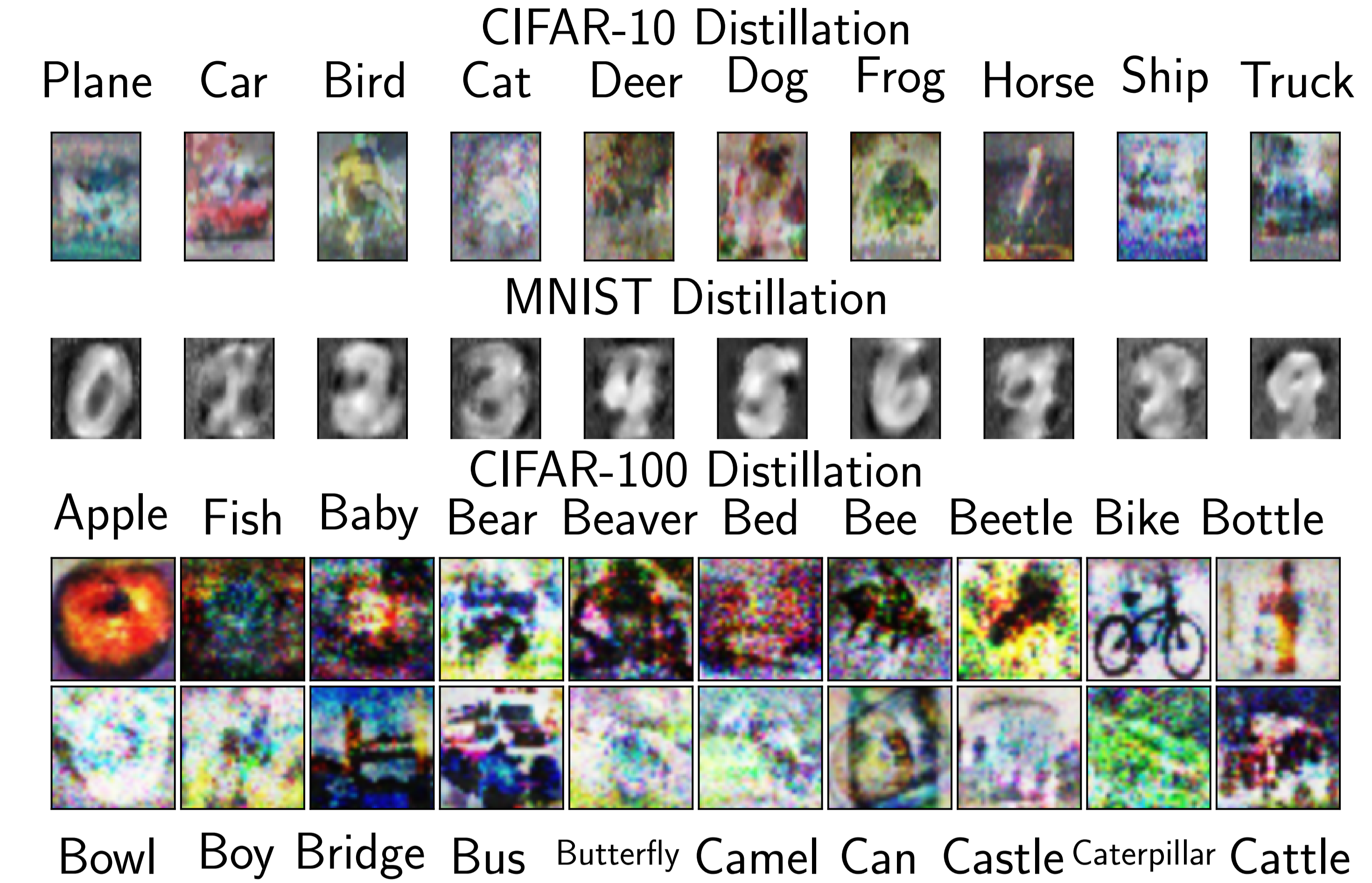
Comparing Inverse-Hessian Approximations



Overfitting the Validation Set



Data Distillation



We learn **one distilled image per class**, so after training a logistic regression classifier on the distillation, it generalizes to the rest of the data.

Learned Data Augmentation



We train a U-net for held-out data with weights as hyperparameters, which takes data & noise, and outputs augmented data.

Tuning Regularization Parameters

Method	Validation	Test	Time(s)
Grid Search	97.32	94.58	100k
Random Search	84.81	81.46	100k
Bayesian Opt.	72.13	69.29	100k
STN	70.30	67.68	25k
Ours	69.22	66.40	18.5k
Ours, Many	68.18	66.14	18.5k

For LSTMs on PTB, we tune the same 7 (and millions of) hyperparameters **faster, with comparable memory, and to a lower perplexity**.