

Training Data Attribution with Approximate Unrolling

Juhan Bae^{1,2}, Wu Lin², Jonathan Lorraine^{1,2,3}, Roger Grosse^{1,2,4,5}

¹University of Toronto, ²Vector Institute, ³NVIDIA, ⁴Anthropic, ⁵Schwartz Reisman Institute



VECTOR INSTITUTE

Training Data Attribution

- **Training data attribution (TDA)** techniques are motivated by understanding the relationship between training data and the properties of trained models.
- Many TDA methods aim to perform a **counterfactual prediction**, which estimates how a model's behavior would change if certain data points were removed from (or added to) the training dataset.

Implicit-differentiation-based TDA

- Implicit-differentiation-based TDA (e.g., **influence functions**) uses the Implicit Function Theorem to estimate the sensitivity of the optimal solution θ^* to downweighting a training data point \mathbf{z} :

$$\theta_{\text{removed}}^* \approx \theta^* + \frac{1}{N} \mathbf{H}_{\theta^*}^{-1} \nabla_{\theta} \mathcal{L}(\theta^*, \mathbf{z}).$$

- These methods **provide convenient estimation algorithms** that depend solely on the optimal model parameters rather than intermediate checkpoints throughout training.
- However, the classical formulation relies on assumptions such as **the uniqueness of and convergence to the optimal solution**.

Unrolling-based TDA

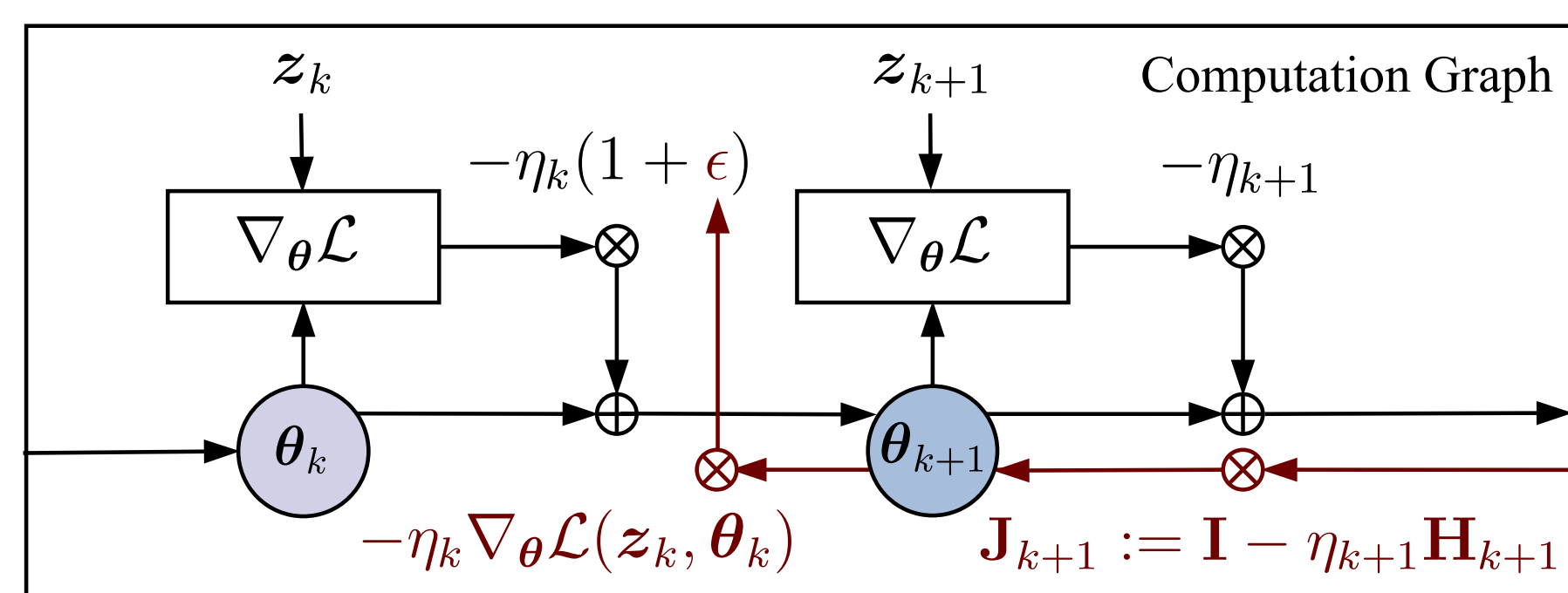
- Consider an update rule at iteration k with a data point weight ϵ :

$$\theta_{k+1} \leftarrow \theta_k - (1 + \epsilon) \eta_k \nabla_{\theta} \mathcal{L}(\theta_k, \mathbf{z}_k).$$

- Unrolling-based TDA methods estimate the effect of removing a data point \mathbf{z} on the final parameters θ_T by **backpropagating through the preceding optimization steps**:

$$\begin{aligned} \theta_T^{\text{removed}} &\approx \theta_T - \left. \frac{d\theta_T}{d\epsilon} \right|_{\epsilon=0} \\ &= \theta_T - \frac{\partial \theta_T}{\partial \theta_{T-1}} \cdots \frac{\partial \theta_{k+2}}{\partial \theta_{k+1}} \frac{\partial \theta_{k+1}}{\partial \epsilon} \Big|_{\epsilon=0} \\ &= \theta_T - (\mathbf{I} - \eta_{T-1} \mathbf{H}_{T-1}) \cdots (\mathbf{I} - \eta_{k+1} \mathbf{H}_{k+1}) (-\eta_k \nabla_{\theta} \mathcal{L}(\theta_k, \mathbf{z}_k)). \end{aligned}$$

- They do not rely on the **uniqueness of or convergence to the optimal solution** and can **incorporate details of training process**.
- However, they require storing **all intermediate variables** during the training process for backpropagation (e.g., parameter vectors for each optimization step).



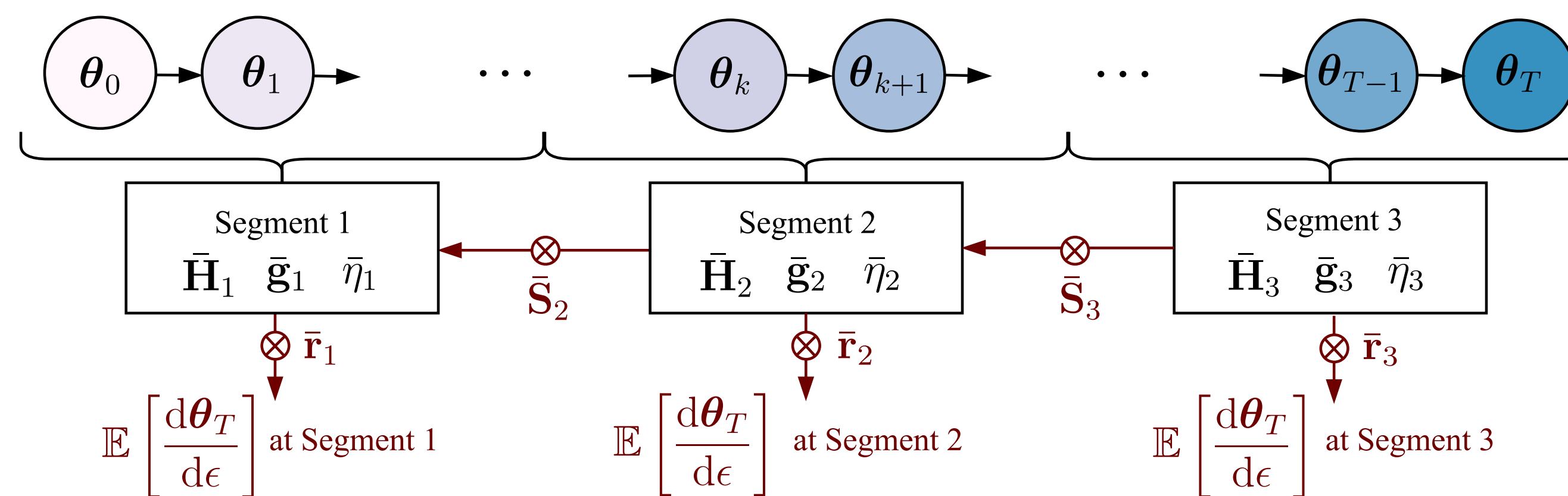
① → : Training
② ← : Gradient Accumulation

SOURCE (Our Proposed Algorithm)

TDA Strategy	Number of Checkpoints	Allows Non-Convergence	Supports Multi-Stage	Incorporates Optimizer
Implicit Differentiation	1	✗	✗	✗
Unrolling	T	✓	✓	✓
SOURCE (ours)	$C \ll T$	✓	✓	✓

- In this work, we connect implicit-differentiation-based and unrolling-based TDA approaches and introduce SOURCE that **enjoys the advantages of both methods**.
- SOURCE inherits three key advantages from unrolling-based methods:
 1. It enables TDA analysis for multi-stage training pipelines (e.g., foundational models and continual learning).
 2. It can incorporate algorithmic choices into the analysis (e.g., SGD vs. Adam).
 3. It maintains a close connection to counterfactual predictions even when implicit-differentiation assumptions fail (e.g., non-converged parameters).
- Unlike previous unrolling approaches, SOURCE achieves these benefits while **requiring only a small number of model checkpoints C** (e.g., $C = 5$) rather than storing the entire training trajectory.

Segmenting the Training Trajectory



- **Key Idea:** SOURCE partitions the training trajectory into one or more **segments** and **approximates the distributions of gradients and Hessians as stationary within each segment**.
- Given L segments, SOURCE approximates the expected total gradient over the data point ordering as:

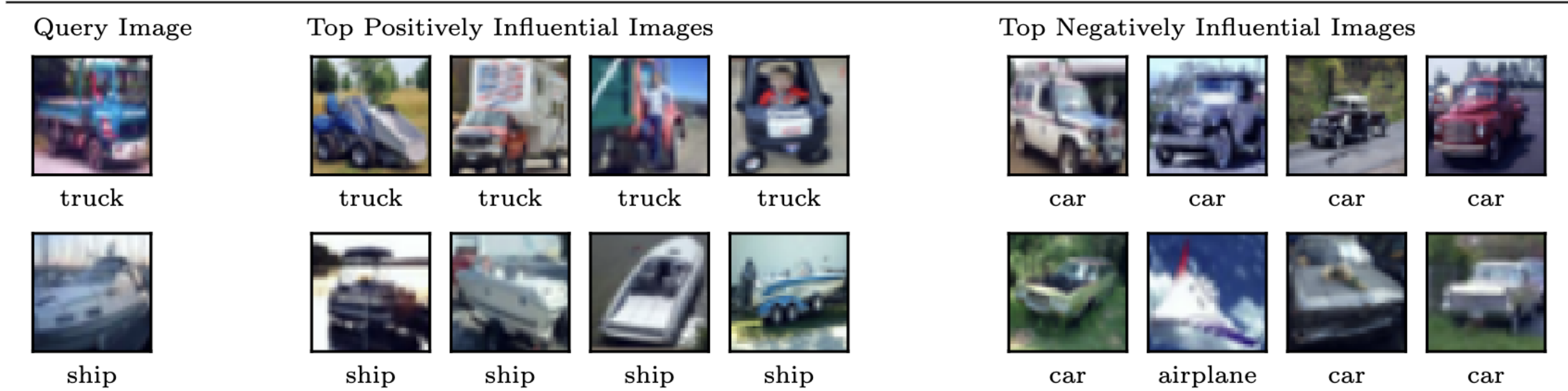
$$\mathbb{E} \left[\frac{d\theta_T}{d\epsilon} \right] \approx - \sum_{\ell=1}^L \left(\prod_{\ell'=L}^{\ell+1} \mathbb{E} [\mathbf{S}_{\ell'}] \right) \mathbb{E} [\mathbf{r}_{\ell}],$$

where these quantities are computed using segment-specific **averaged Hessians $\bar{\mathbf{H}}_{\ell}$** and **gradients $\bar{\mathbf{g}}_{\ell}$** . Note $\mathbb{E} [\mathbf{S}_{\ell}] := \exp(-\bar{\eta}_{\ell} K_{\ell} \bar{\mathbf{H}}_{\ell})$ and $\mathbb{E} [\mathbf{r}_{\ell}] := \frac{1}{N} (-\exp(-\bar{\eta}_{\ell} K_{\ell} \bar{\mathbf{H}}_{\ell})) \bar{\mathbf{H}}_{\ell}^{-1} \bar{\mathbf{g}}_{\ell}$, where K_{ℓ} is the total number of iterations performed in the specified segment.

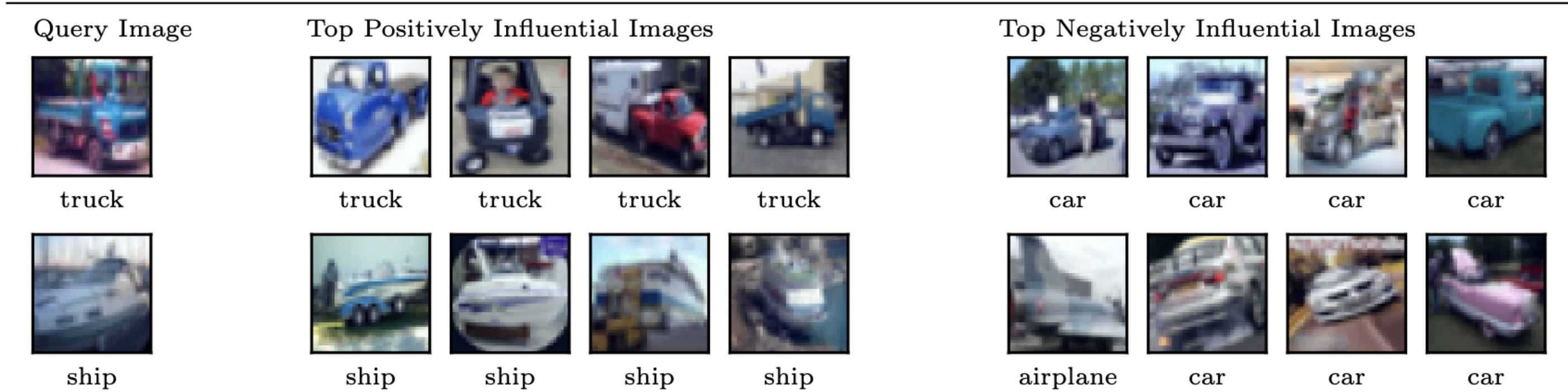
- For practical implementation, we use **EK-FAC** to efficiently approximate the Hessian, with both averaged Hessian and gradient estimates computed using **a set of checkpoints within each segment**.
- SOURCE is **C times more computationally expensive** than EK-FAC influence functions.

Qualitative Results

IF

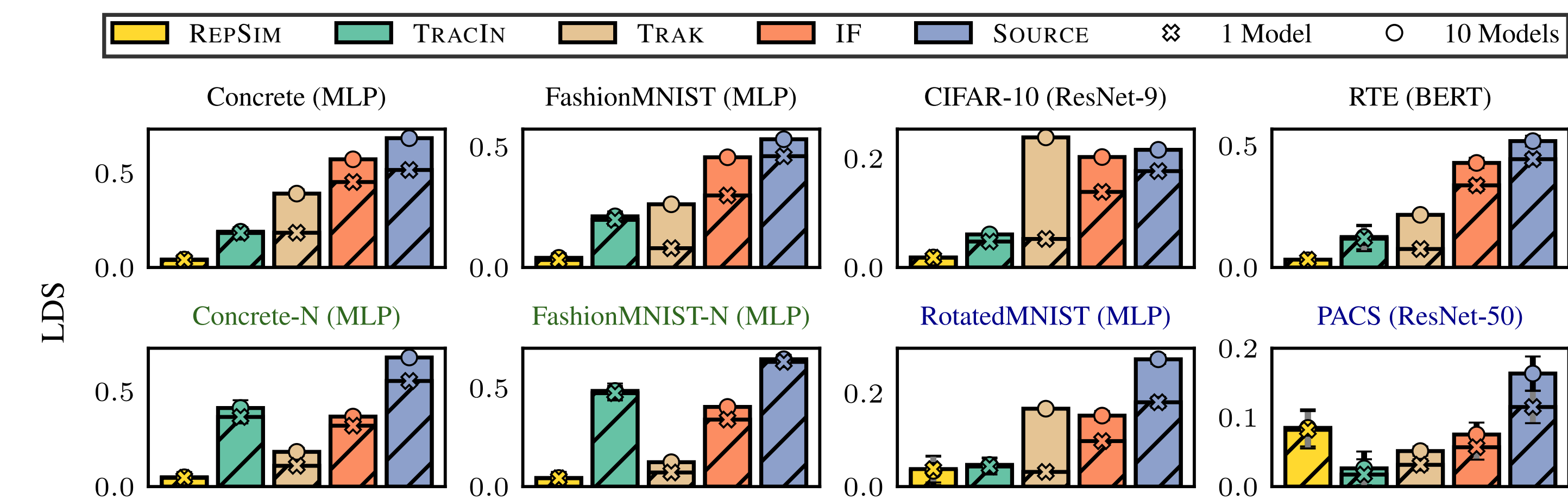


SOURCE



Linear Datamodeling Score (LDS)

- The LDS measures the Spearman correlation between the estimated quantities after retraining the model without a subset of data points and the predictions made by the TDA method.
- SOURCE especially performs strongly against other baseline techniques **on settings that pose challenges to implicit-differentiation-based approaches** (e.g., **non-converged models** and **models trained with multiple stages**).



Subset Removal Evaluation

