

# Estimating Training Data Influence by Tracing Gradient Descent

Md Mahfuzur Rahman

Postdoctoral Research Associate

TReNDS Center

Georgia State University



---

# Estimating Training Data Influence by Tracing Gradient Descent

---

**Garima\***

Google

pruthi@google.com

**Frederick Liu\***

Google

frederickliu@google.com

**Satyen Kale**

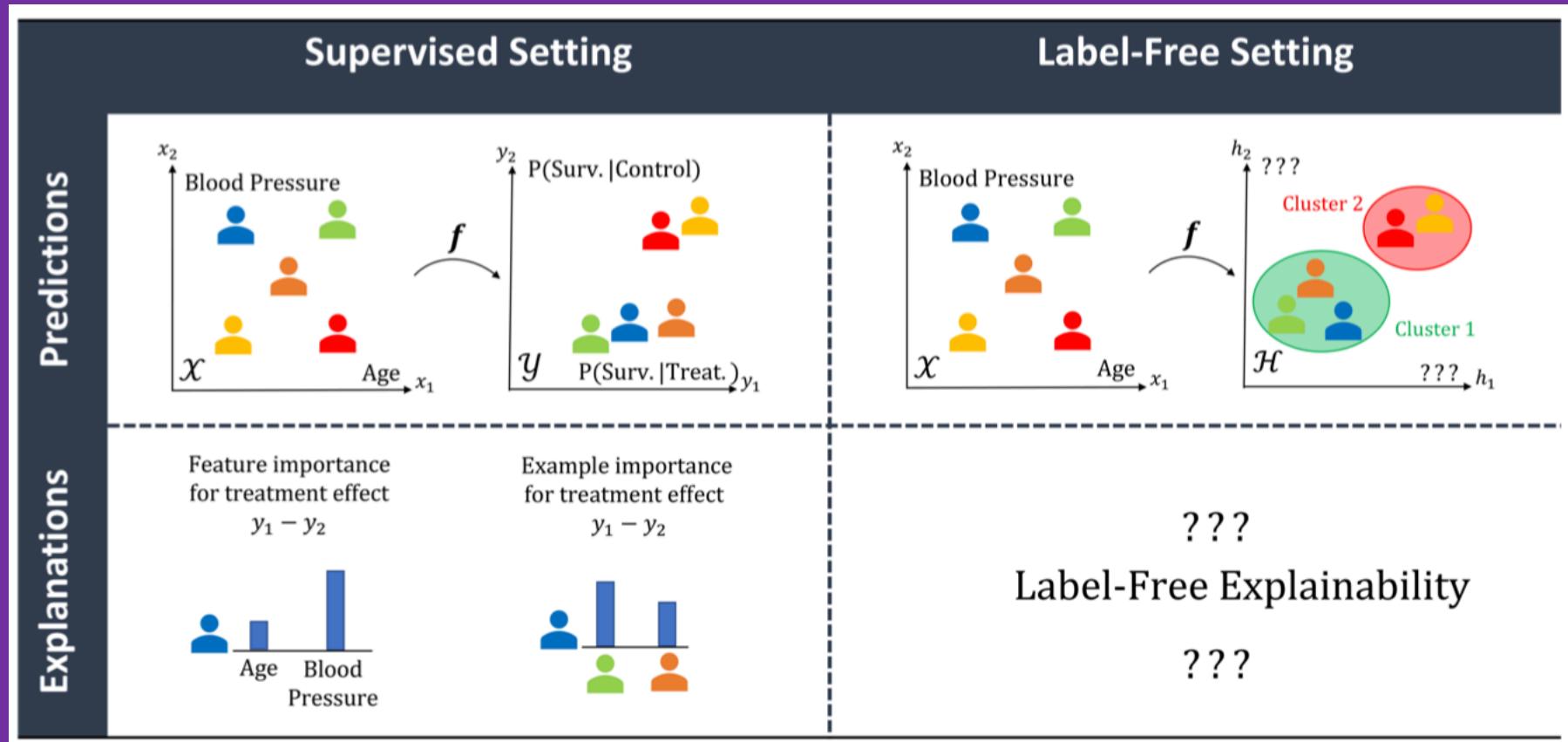
Google

satyenkale@google.com

**Mukund Sundararajan** †

Google

mukunds@google.com



Crabbé, Jonathan, and Mihaela van der Schaar. "Label-Free Explainability for Unsupervised Models." International Conference on Machine Learning. PMLR, 2022  
 Sutre, Elina Thibeau, et al. "How can data augmentation improve attribution maps for disease subtype explainability." Medical Imaging 2023: Image Processing. SPIE, 2023.

Input:

$$\boldsymbol{x} \in \mathbb{R}^d$$

Model:

$$F : \mathbb{R}^d \rightarrow \mathbb{R}^C$$

Class-specific logit:

$$F_c(\boldsymbol{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$$

Input

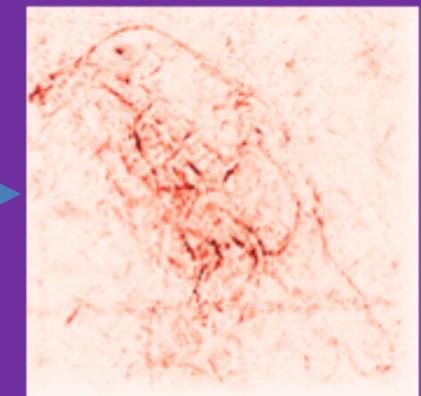


Explanation:

$$E : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Explanation  
 $E$

Model  
(usually  
Black-box)



Explanation

- *TracIn* - computes the influence of a training example on a prediction.



church



church



church



church



castle



castle



castle

test sample

proponents

opponents

The main idea:

- how does the loss on the test point change during the training process?

For scalable implementation:

- a. Use first-order gradient approximation.
- b. Saved checkpoints of standard training procedures
- c. Cherry-picking layers of a deep neural network.

- FTC decomposes the difference between a function at two points
- Uses the gradients along the path.
- The integration of the gradients between two points is equal to the difference.
- Approximation to integration is available.
- TracIn decomposes the difference between the losses of the test point

$$A(x + h) - A(x) = f(x).h + (\text{Excess})$$

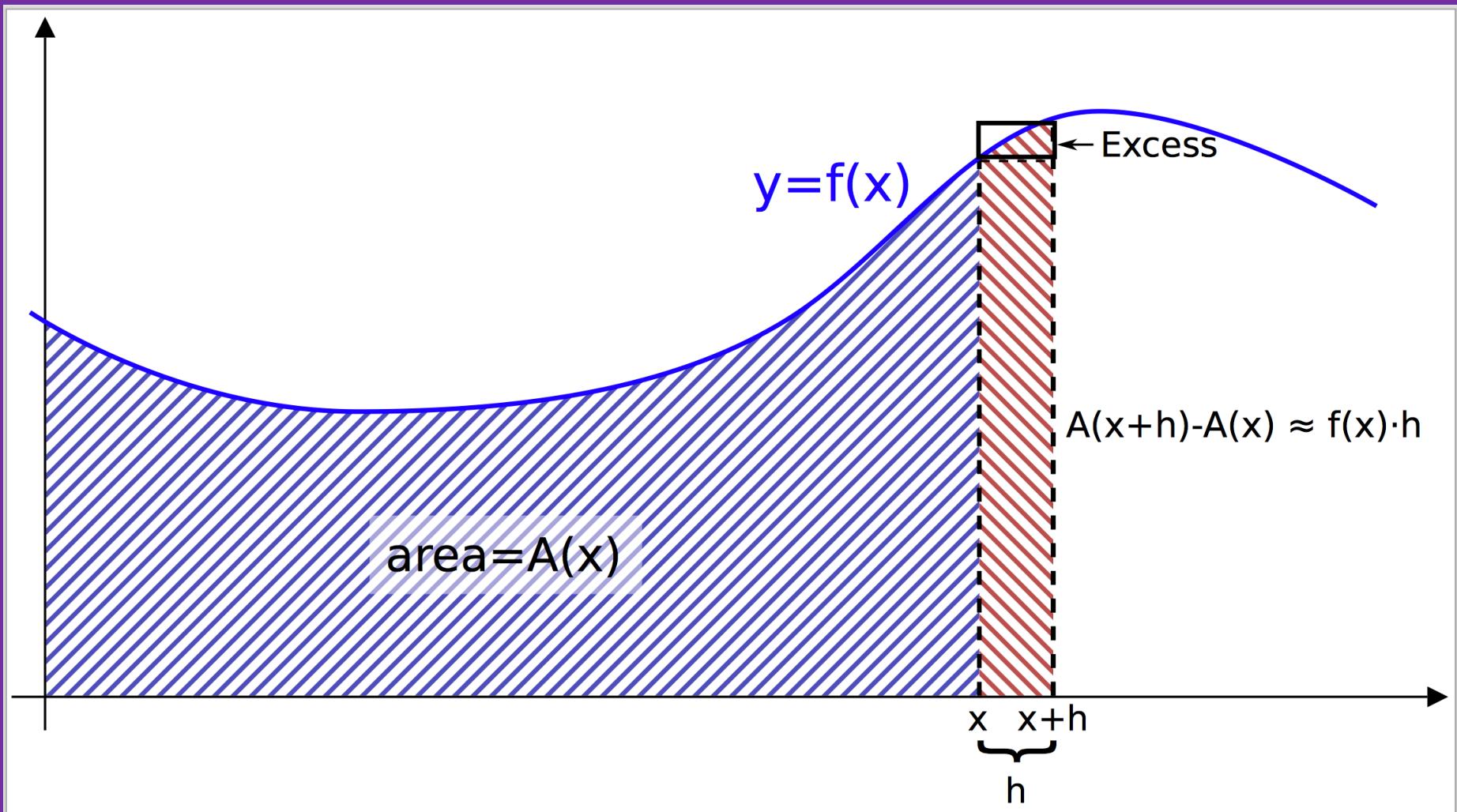


Image: [wikipedia](#)

Assume:  $A$  represents loss function,  $f$  represents gradients

Training set:

$$\mathbb{S} = \{z_1, z_2, \dots, z_n \in \mathbb{Z}\}$$

$$\text{TracInIdeal}(z, z') = \sum_{t: z_t=z} \ell(w_t, z') - \ell(w_{t+1}, z')$$

- Assume SGD as an optimizer
- Consider the steps where training example  $z$  was used

**Lemma 3.1** Suppose the initial parameter vector before starting the training process is  $w_0$ , and the final parameter vector is  $w_T$ . Then  $\sum_{i=1}^n \text{TracInIdeal}(z_i, z') = \ell(w_0, z') - \ell(w_T, z')$

FOGA to loss function:

$$\ell(\mathbf{w}_{t+1}, \mathbf{z}') = \ell(\mathbf{w}_t, \mathbf{z}') + \nabla \ell(\mathbf{w}_t, \mathbf{z}') \cdot (\mathbf{w}_{t+1} - \mathbf{w}_t) + O(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2)$$

According to gradient update rule:

$$\mathbf{w}_{t+1} - \mathbf{w}_t = -\eta_t \nabla \ell(\mathbf{w}_t, \mathbf{z}_t)$$

Loss change now can be approximated as:

$$\ell(\mathbf{w}_t, \mathbf{z}') - \ell(\mathbf{w}_{t+1}, \mathbf{z}') \approx \eta_t \nabla \ell(\mathbf{w}_t, \mathbf{z}') \cdot \nabla \ell(\mathbf{w}_t, \mathbf{z}_t)$$

The influence of a training example:

$$\text{TracIn}(\mathbf{z}, \mathbf{z}') = \sum_{t: \mathbf{z}_t = \mathbf{z}} \eta_t \nabla \ell(\mathbf{w}_t, \mathbf{z}') \cdot \nabla \ell(\mathbf{w}_t, \mathbf{z})$$

Issue: not feasible, gradient is not available per training example

Summing up over all iterations t in which a particular training point z was chosen in  $B_t$

$$\text{TracIn}(z, z') = \frac{1}{b} \sum_{t: z \in B_t} \eta_t \nabla \ell(w_t, z') \cdot \nabla \ell(w_t, z)$$

- Access to the parameter vector at the specific iteration is not available.
- approximate it with the first checkpoint parameter vector after it.

K checkpoints:

$$\mathbf{w}_{t_1}, \mathbf{w}_{t_2}, \dots, \mathbf{w}_{t_k}$$

Influence based on checkpoints:

$$\text{TracInCP}(\mathbf{z}, \mathbf{z}') = \sum_{i=1}^k \eta_t \nabla \ell(\mathbf{w}_{t_i}, \mathbf{z}) \cdot \nabla \ell(\mathbf{w}_{t_i}, \mathbf{z}')$$

Derivation to another experimental setup could be different. But CP version would be quite same

**Remark 3.4 (Handling Variations of Training)** *In our derivation of TracIn we have assumed a certain form of training. In practice, there are likely to be differences in optimizers, learning rate schedules, the handling of minibatches etc. It should be possible to redo the derivation of TracIn to handle these differences. Also, we expect the practical form of TracInCP to remain the same across these variations.*

The same formulation could be used to study the influence of a hypothetical training point

**Remark 3.5 (Counterfactual Interpretation)** *An alternative interpretation of Equation I is that it approximates the influence of a training example on a test example had it been visited at each of the input checkpoints. Under this counterfactual interpretation, it is valid to study the influence of a point that is not part of the training data set, keeping in mind that such an example did not impact the training process or the checkpoints that arose as a consequence of the training process.*

Self-influence score:

the influence of a training point on its own loss, i.e., the training point  $z$  and the test point  $z'$  in  $\text{TracIn}$  are identical

- Proponents (positive value of influence) serve to reduce loss
- Opponents serve to increase loss.
- Incorrectly labelled examples are likely to be strong proponents for themselves.
- When sorted based on *self-influence*, an effective influence computation method would tend to rank mislabelled examples in the beginning.

Dataset and models:

ResNet-56 trained on the CIFAR-10.

The model on the original dataset has 93.4% test accuracy

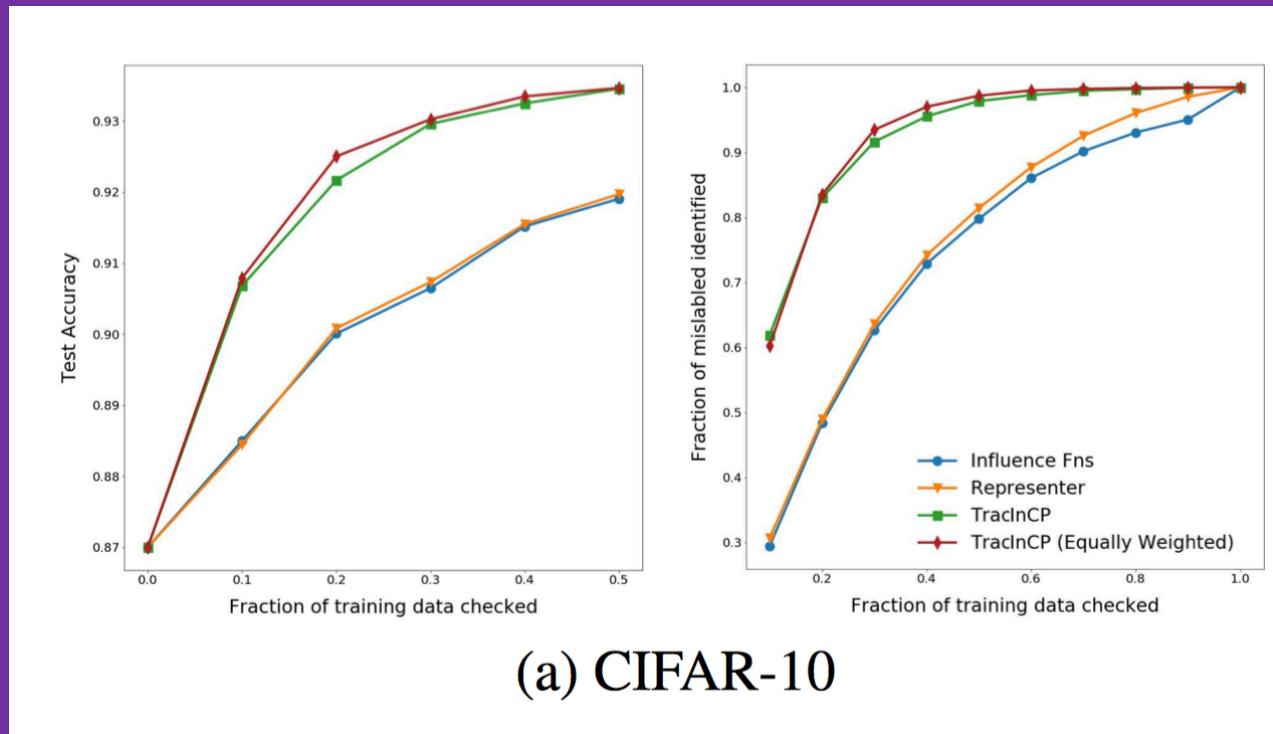
Setup for Mislabelled Example Identification:

- first train a model on correct data.
- Then, for 10% of the training data, labels were changed to the highest scoring incorrect labels.
- Test accuracy reduces from 93.4% to 87.0% (train accuracy is 99.6%).)
- Identify mislabelled examples

Metric:

the fraction of mislabelled data recovered for different prefixes of the rank order

- TracInCP with only the last layer.
- sample every 30 checkpoints starting from the 30th checkpoint
- Consider both “weighted” and “equally weighted” checkpoints.



Right: TracIn recovers more than 80% of the mislabelled data in 20% of the ranking, the other methods recover less than 50% at the same point

Left: fixing the mislabelled data found within a certain fraction of the training data, results in a larger improvement of test accuracy

Different checkpoints had focus on different classes

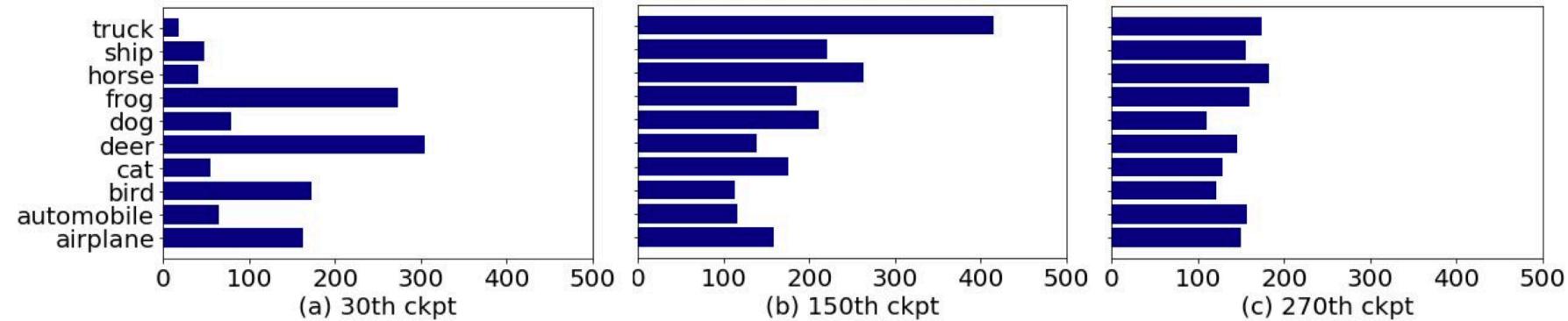
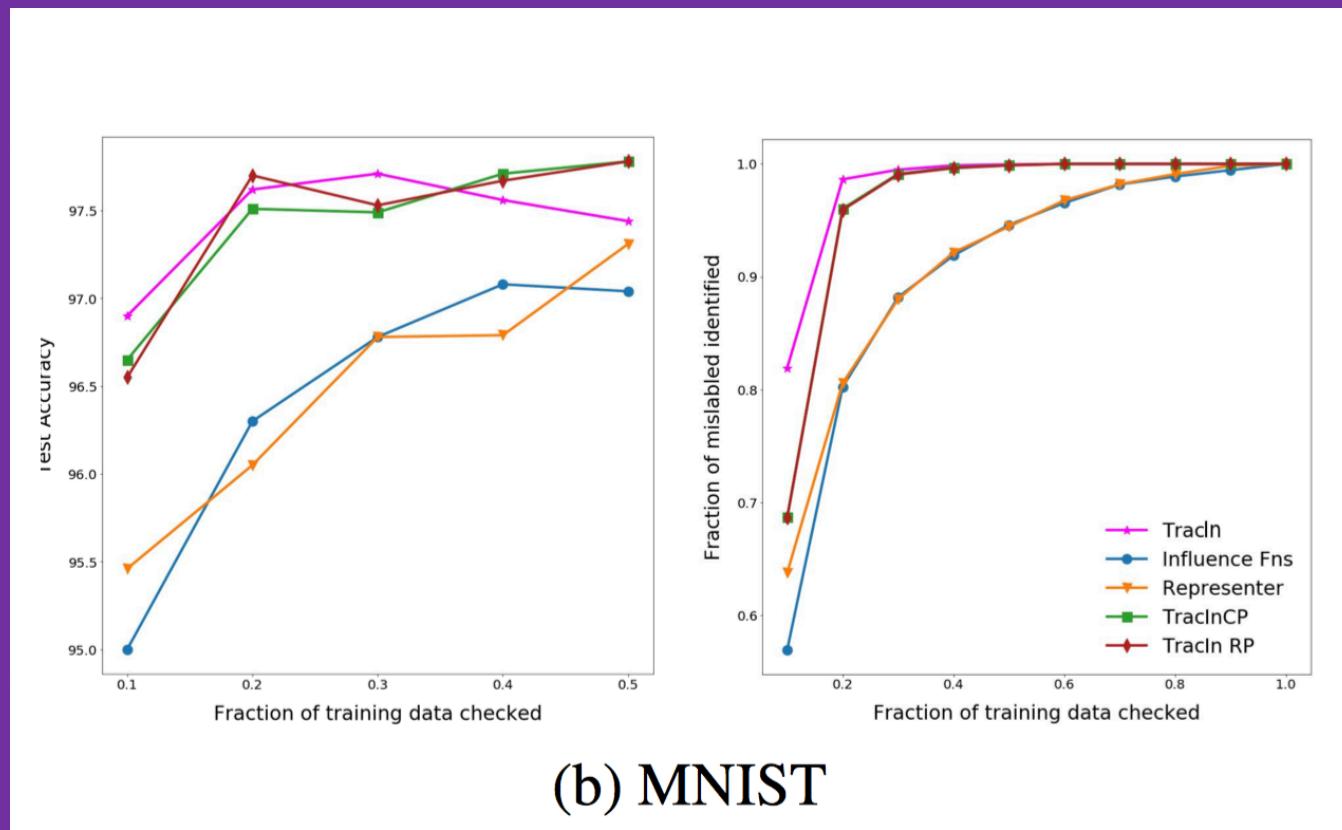
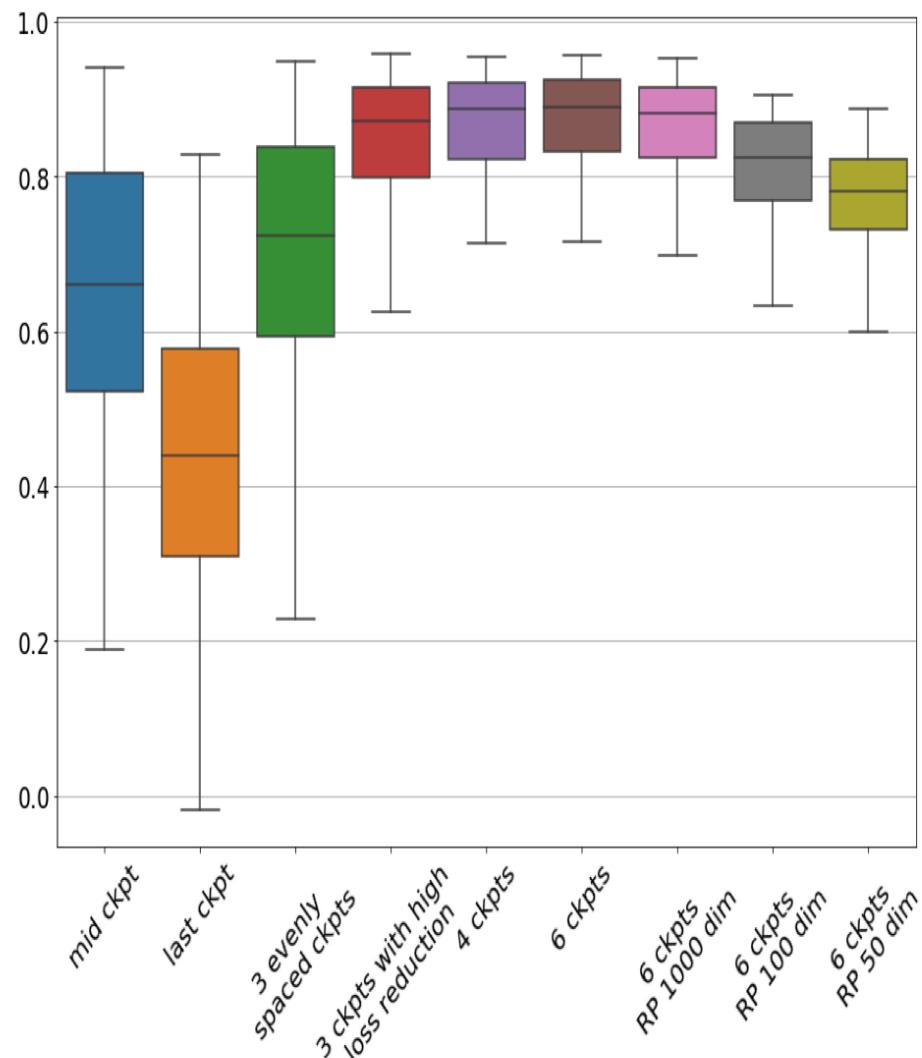


Figure 8: Number of identified mislabelled examples by class for three checkpoints within the top 10% of ranking by self-influence. Different checkpoints highlight different labels—see Section 4.2.

- Model is small - 3 hidden layers and 240K parameters.
- 99.30% train, 97.55% test accuracy.
- 89.94% mislabelled train, 89.95% on test set
- Use each training step, and not just checkpoints.
- The whole model used, not just the last layer.



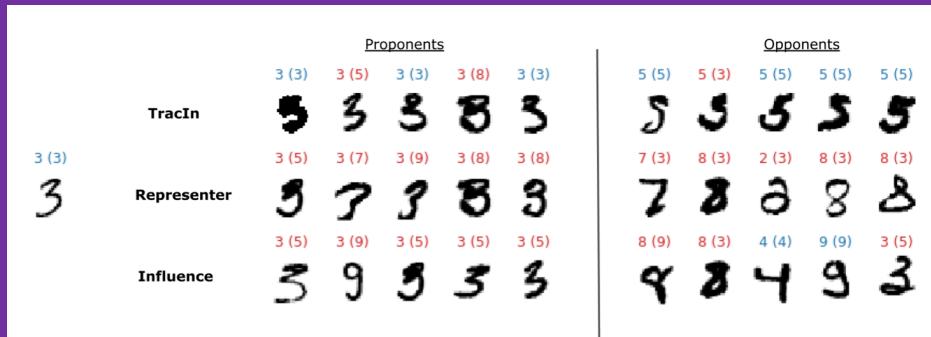
Note: approximate TracIn is able to recover mislabelled examples faster than TracInCP



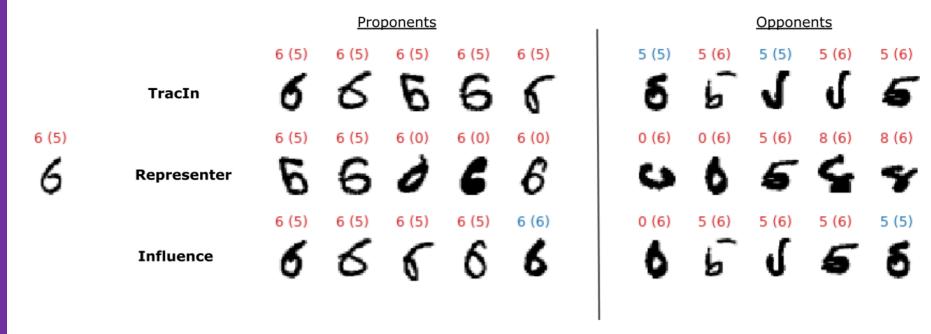
- correlation of the influence scores
- 100 test points
- TracInCP (different CPs) vs. approximation TracIn.
- Mid CP is more informative than last CP
- CPs with high loss reduction, are more informational
- Random Projection (RP) version saves computations.

RP formulation:

$$(\eta_t G \nabla \ell(\mathbf{w}_t, \mathbf{z}_t)) \cdot (G \nabla \ell(\mathbf{w}_t, \mathbf{z}'))$$



(a) Correctly classified 3.



(b) Incorrectly classified 6

(Predicted class in brackets)

- Proponents are consistently similar for TracInCP and Representer methods.
- Not necessarily examples from the same class are proponents always.
- TracInCP seems to pick pixel-wise similar opponents

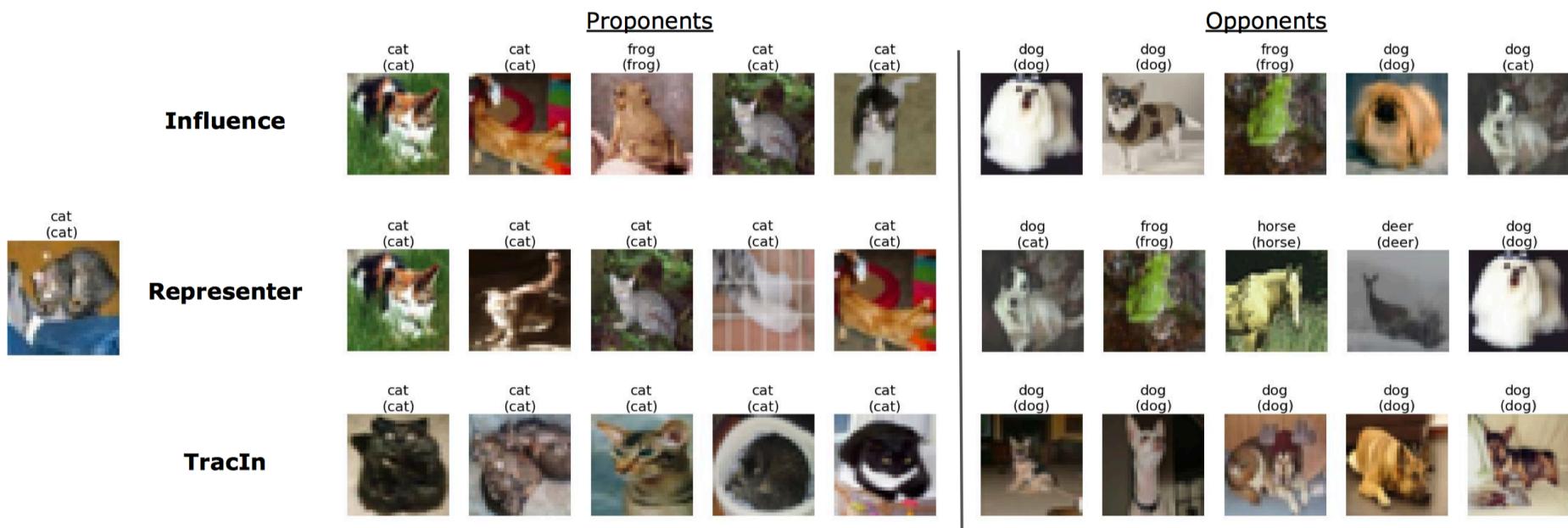


Figure 5: CIFAR-10 results: Proponents and opponents examples of a correctly classified cat for influence functions, representer point, and TracIn. (Predicted class in brackets)

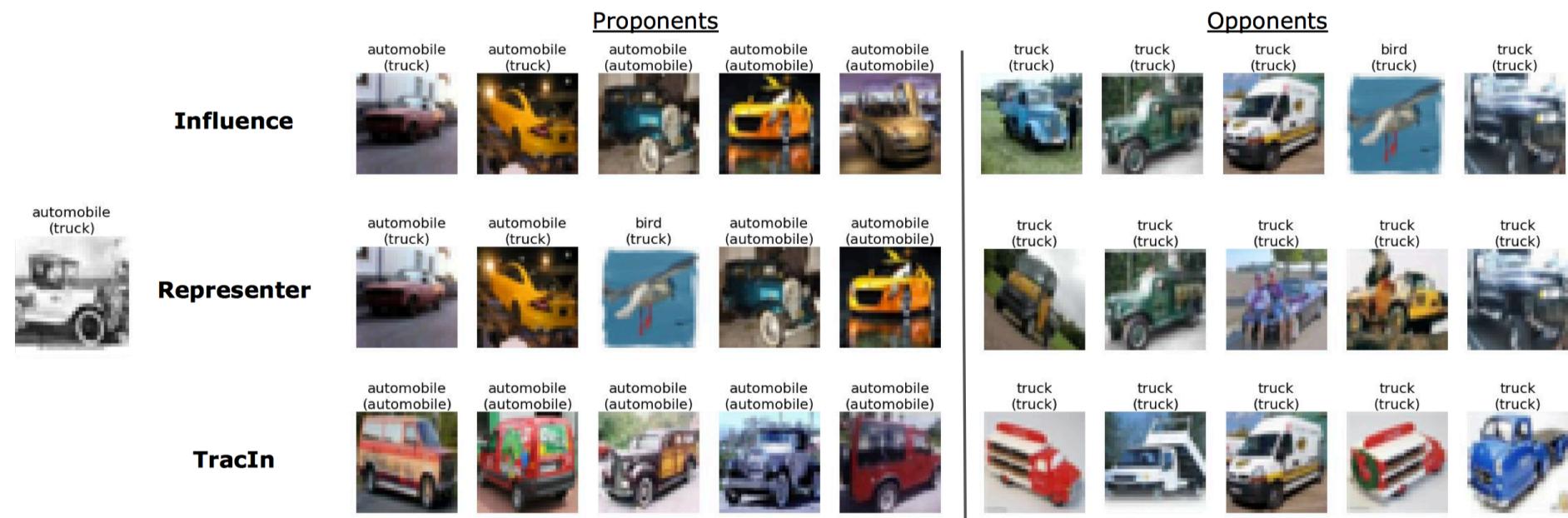


Figure 6: CIFAR-10 results: Proponents and opponents examples of an incorrectly classified automobile for influence functions, representer point, and TracIn. (Predicted class in brackets)

TracIn was applied on the fully connected layer of ResNet-50 trained on ImageNet

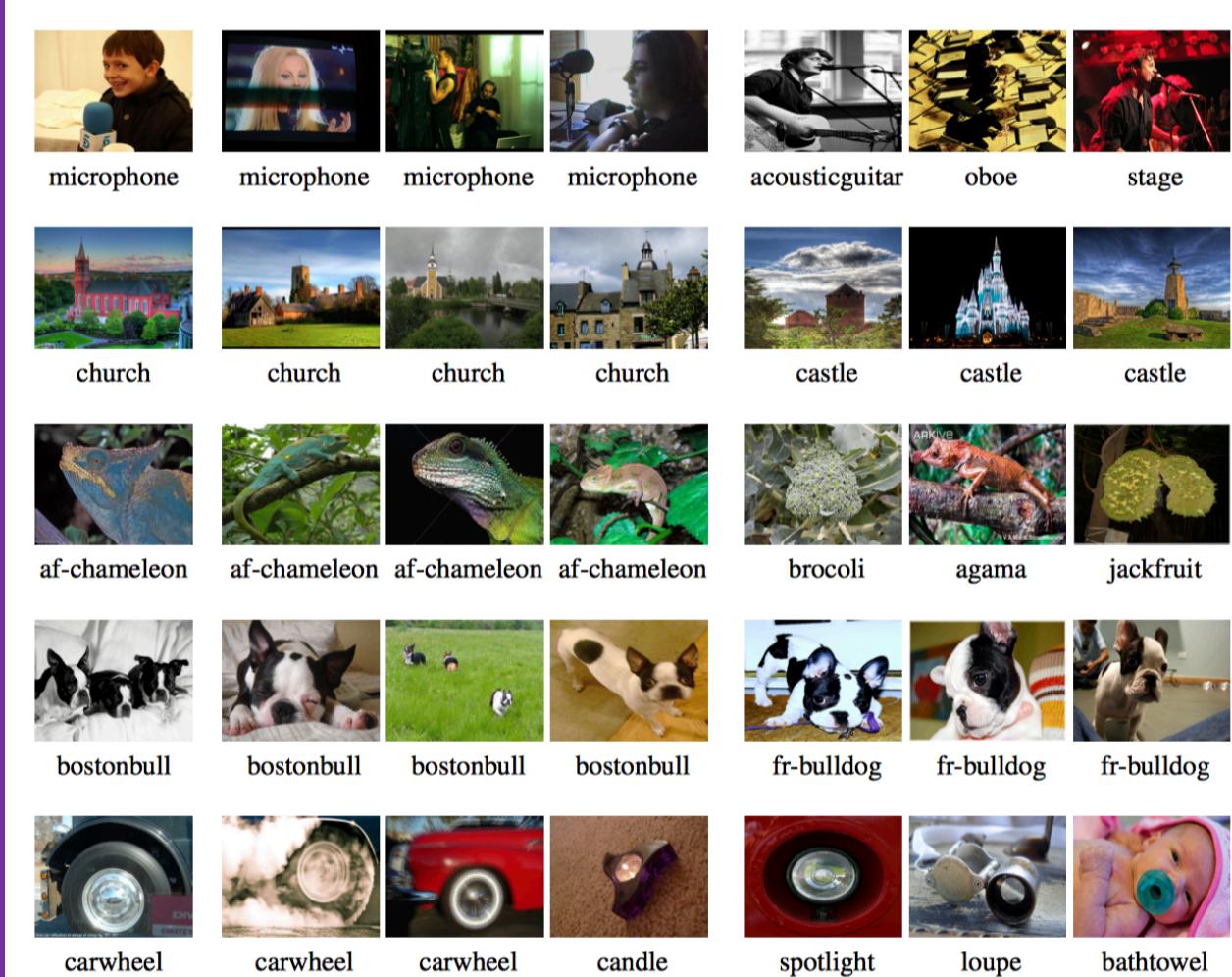


Figure 4: TracIn applied on Imagenet. Each row starts with the test example followed by three proponents and three opponents. The test image in the first row is classified as band-aid and is the only misclassified example. (af-chameleon: african-chameleon, fr-bulldog: french-bulldog)

- The method is simple
  - only requires concepts like gradients, checkpoints and loss functions.
- The method is general
  - applicable to any ML model using SGD or its variants,
  - agnostic of architecture, domain and task.
- The method is versatile
  - explain a single prediction
  - identify misslabelled examples, etc.

# Questions?



Contact: [mrahman21@gsu.edu](mailto:mrahman21@gsu.edu)



Influence functions:

what happens to the loss of a test point when you drop an individual training point and retrain?

TracIn:

how much change in loss of a test point happened between the start of training and the end of training, for a training example?