# Efficient Model Reconstruction Leveraging Interpretability

Pasan Dissanayake
PhD Advisor: Sanghamitra Dutta
Northrop Grumman Tech Lead: Kerry Brown

UNIVERSITY OF MARYLAND 1856 — DEPARTMENT OF ELECTRICAL & COMPUTER ENGINEERING

NORTHROP GRUMMAN

NEURAL INFORMATION PROCESSING SYSTEMS — AISTATS

## MOTIVATION

### Ever-Increasing Model Size

Resource Constrained Environments:

*The evolution of GPT's number of parameters over time. (source: Medium)*

### Trust in High-Stakes Applications
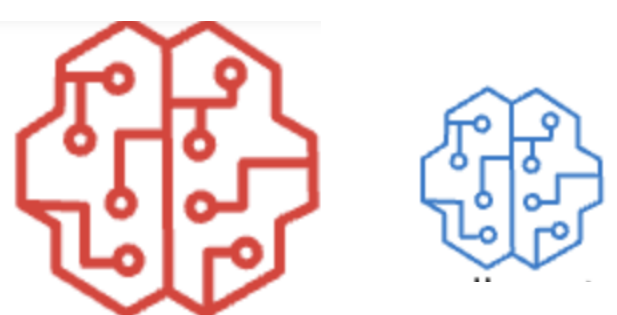
HIRING — FINANCE — HEALTHCARE — DEFENSE — AUTONOMY

**Can we build Efficient and Trustworthy AI by systematically leveraging Interpretability?**

---

## MODEL RECONSTRUCTION USING COUNTERFACTUAL EXPLANATIONS

**Dissanayake** & Dutta
**NeurIPS 2024**

### Model Compression
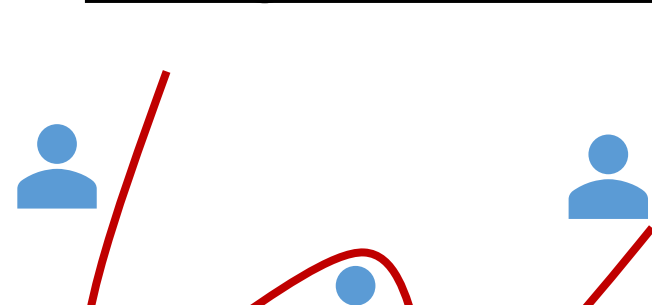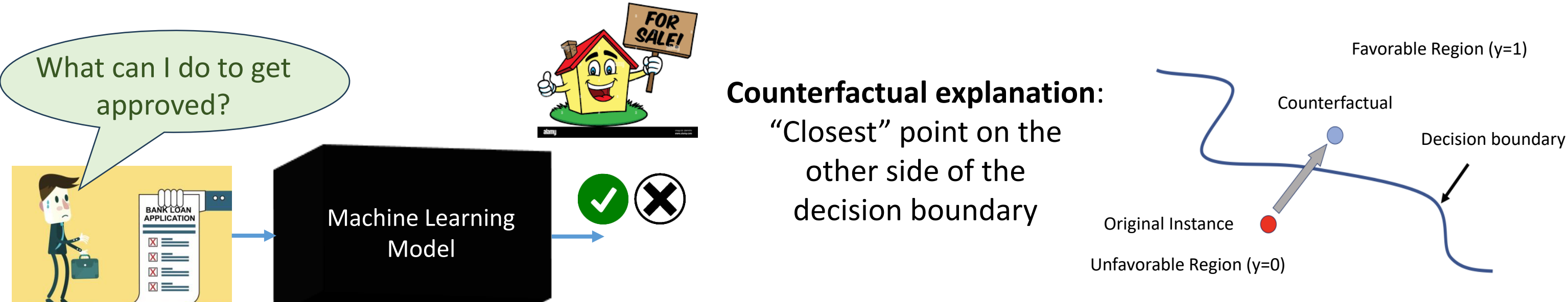Replicate A Large Model Using A Small Model

### Security
MLaaS
Understand Limits of Model Extraction & Stealing

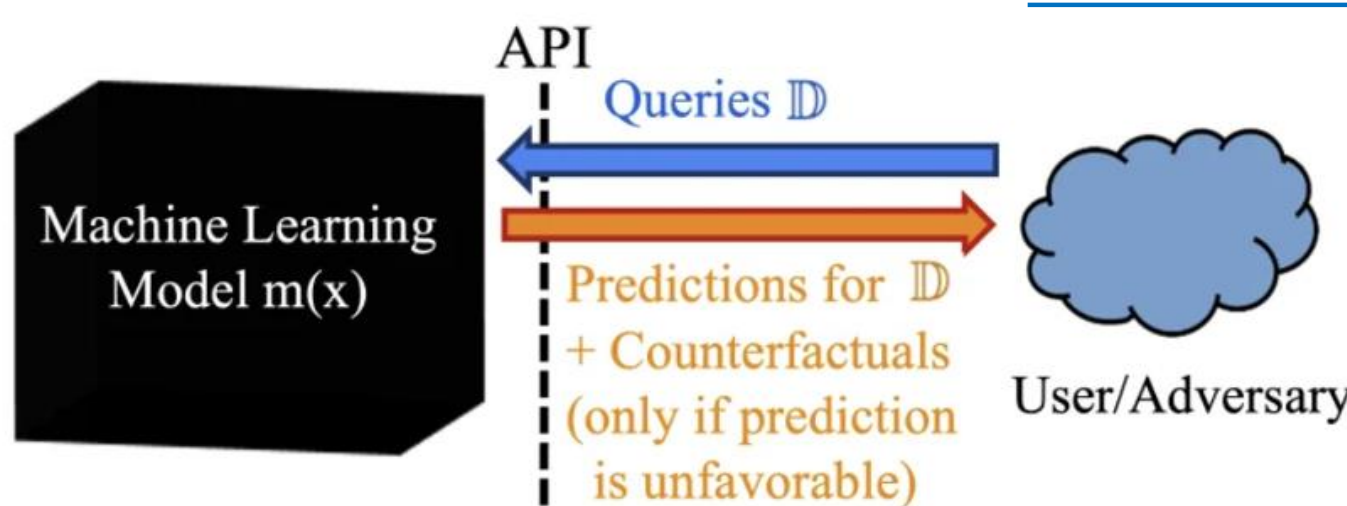### Global Explainability/Audit in High-Stakes ML
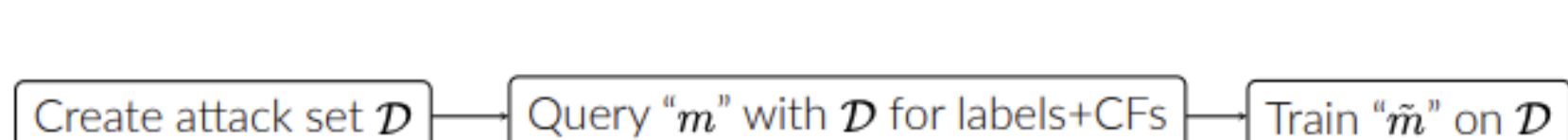Global Understanding From Local Crowdsourced Information

This Work: Model Reconstruction using an Interpretability technique called **Counterfactual Explanations**
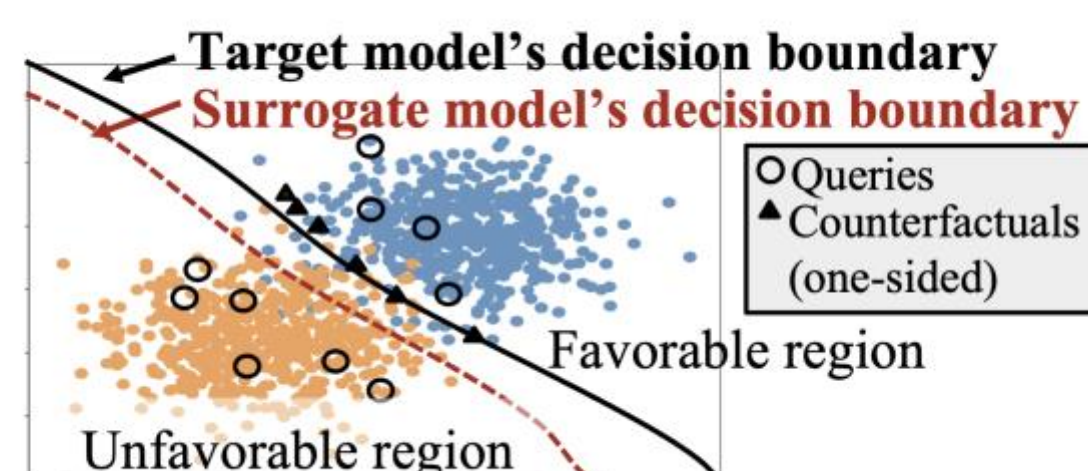
What can I do to get approved?

**Counterfactual explanation:** "Closest" point on the other side of the decision boundary

Favorable Region (y=1)
Counterfactual
Decision boundary
Original Instance
Unfavorable Region (y=0)

### Our Problem Setup

API
Queries $\mathbb{D}$
Machine Learning Model m(x)
Predictions for $\mathbb{D}$ + Counterfactuals (only if prediction is unfavorable)
User/Adversary

For each query, user knows the following:
Accepted (Predicted Label =1)
OR
Denied (Predicted Label =0) & Counterfactual (Closest Accepted Point)

How faithfully can one reconstruct a model using counterfactual explanations?

Create attack set $\mathcal{D}$ → Query "$m$" with $\mathcal{D}$ for labels+CFs → Train "$\hat{m}$" on $\mathcal{D}$

Target model's decision boundary
Surrogate model's decision boundary
Queries
Counterfactuals (one-sided)
Favorable region
Unfavorable region

Counterfactuals treated as ordinary labelled instances?
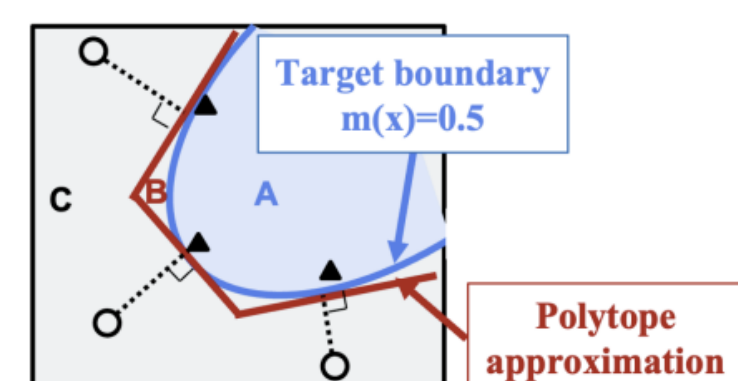**Boundary shift issue**

**Question:** Can we improve model reconstruction using counterfactuals specifically leveraging that the counterfactuals are quite close to the boundary?

**Main Contribution:** New Reconstruction Strategies & Fundamental Limits From Polytope Theory

### MAIN RESULTS

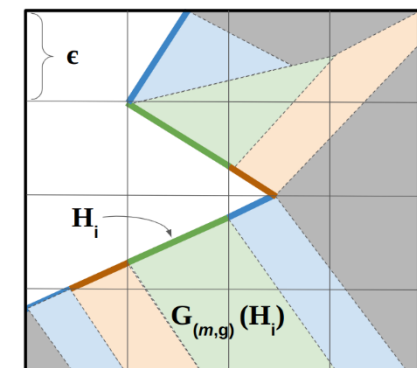**1. Convex Decision Boundaries and Closest Counterfactuals**

Theoretical guarantees on volume approximation using counterfactuals leveraging polytope theory

Target boundary m(x)=0.5
Polytope approximation

**2. ReLU Networks and Closest Counterfactuals**

$$\mathbb{P}[Reconstruction] \geq 1 - k(\epsilon)(1 - v^*(\epsilon))^n$$

Continuous Piece-Wise Linear (CPWL) Functions
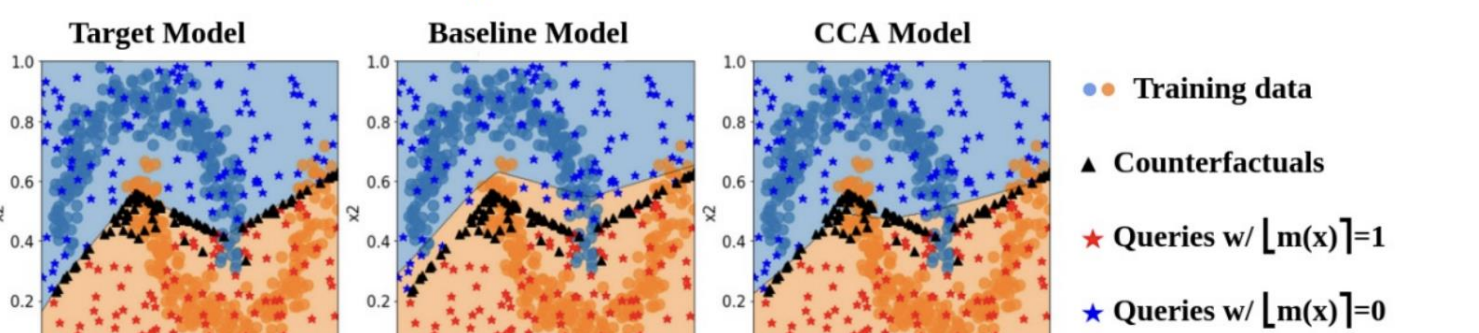
**3. Beyond Closest Counterfactuals**

**Theorem 3.11.** Suppose the target $m$ and surrogate $\hat{m}$ are locally Lipschitz (not necessarily ReLU) such that $m(w) = \hat{m}(w)$ for every counterfactual $w$. Assume the counterfactuals are well-spaced out and forms a $\delta$-cover over the decision boundary. Then $|\hat{m}(x) - m(x)| \leq (\gamma_m + \gamma_{\hat{m}})\delta$, over the target decision boundary.

Target boundary m(x)=0.5
Matching Points
Surrogate boundary

**4. Counterfactual Clamping Attack (CCA)**

$$L_k(\hat{m}(x), y_x) = \mathbb{1}[y_x = 0.5, \hat{m}(x) \leq k](L(\hat{m}(x), k) - h(k)) + \mathbb{1}[y_x \geq 0.5](L(\hat{m}(x), y_x)$$

neglect CFs that already have g(w) > k (for CFs) — for normal examples

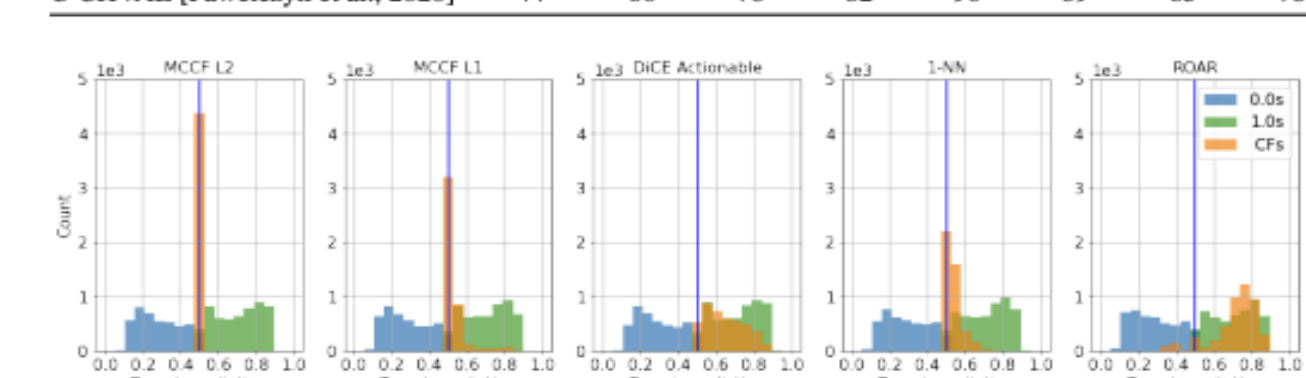Regular Dataset
Regular Dataset & Counterfactual
Clamping Loss Function
Counterfactuals

Target Model — Baseline Model — CCA Model
Training data
Counterfactuals
Queries w/ $\lfloor m(x)\rceil =1$
Queries w/ $\lfloor m(x)\rceil =0$

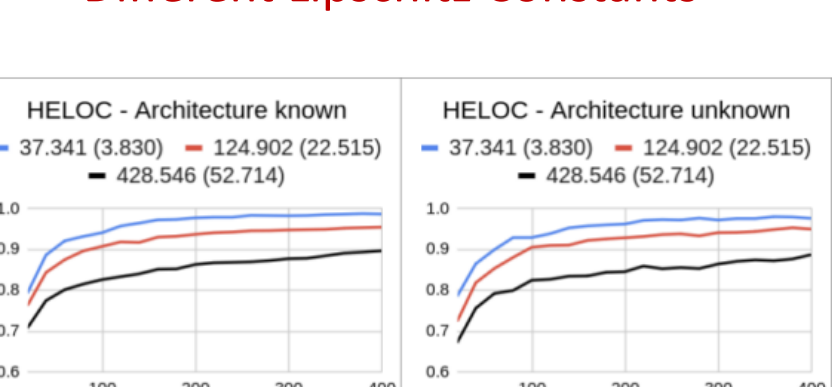| Dataset | Architecture known (model 0) Base. CCA (Dtest) | | | | Architecture unknown (model 1) Base. CCA (Dtest) | | | |
|---|---|---|---|---|---|---|---|---|
| | Base. | CCA | Base. | CCA | Base. | CCA | Base. | CCA |
| Adult In. | 91±3.2 | 94±3.2 | 84±3.2 | 91±3.2 | 91±4.5 | 94±3.2 | 84±4.2 | 90±3.2 |
| COMPAS | 93±2.2 | 96±2.0 | 94±1.7 | 96±2.2 | 91±8.9 | 96±3.2 | 94±2.0 | 94±8.9 |
| DCCC | 89±8.9 | 99±0.9 | 92±2.2 | 96±1.4 | 90±7.7 | 97±4.5 | 95±2.2 | 97±5.8 |
| HELOC | 91±4.7 | 96±2.2 | 92±2.8 | 94±2.4 | 90±7.4 | 95±5.5 | 91±3.3 | 93±3.2 |

**Other Counterfactual Generation Techniques**

Table 2: Fidelity achieved with different counterfactual generating methods in HELOC dataset. Target model has hidden layers with neurons (20, 30, 10). Surrogate model architecture is (10, 20).

| CF method | Fidelity over $\mathbb{D}_{test}$ n=100 Base. CCA | | n=100 Base. CCA | | Fidelity over $\mathbb{D}_{test}$ n=100 Base. CCA | | n=100 Base. CCA | |
|---|---|---|---|---|---|---|---|---|
| MCCF L2-norm | 91 | 95 | 93 | 96 | 91 | 93 | 93 | 95 |
| MCCF L1-norm | 93 | 95 | 94 | 96 | 93 | 92 | 91 | 93 |
| DiCE Actionable | 93 | 94 | 95 | 95 | 90 | 91 | 93 | 94 |
| 1-Nearest-Neighbor | 91 | 95 | 93 | 96 | 91 | 93 | 93 | 95 |
| ROAR [Upadhyay et al., 2021] | 87 | 91 | 92 | 92 | 87 | 85 | 92 | 92 |
| C-CHVAE [Pawelczyk et al., 2020] | 77 | 82 | 92±2.8 | | 79 | 89 | 83 | 78 |

**Different Lipschitz Constants** — **Different Model Architectures**

**CCA outperforms baselines:** A small set of curated data points are sufficient for model extraction with high fidelity!
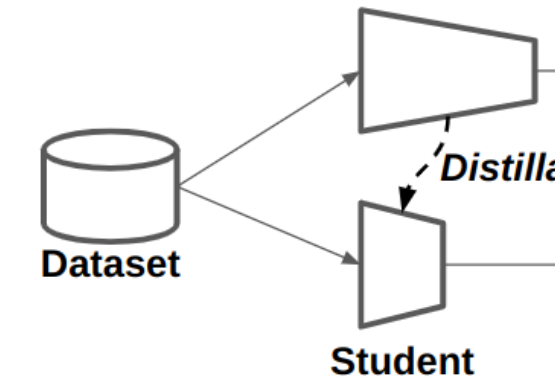
---

## KNOWLEDGE DISTILLATION USING PARTIAL INFORMATION DECOMPOSITION

**Dissanayake**, Hamman, Halder, Sucholutsky, Zhang, Dutta
**AISTATS 2025**

**Benefits of distillation:**
- Training is energy/data efficient
- Simpler student → runs on limited resources (e.g. edge devices)
- Student can be more interpretable → good for high-stakes applications

Teacher — Task 1 — Distillation — Task 2 — Student — Dataset

$$L(\eta_s) = \lambda_1 L_{ordinary}\left(Y, \hat{Y}(X)\right) + \lambda_2 L_{distill}(Y, S_{\eta_s}(X), T(X))$$

Propose using Redundant Information for Task-Aware Knowledge Distillation + Incorporate it into Optimization

**Teachers are not always helpful!!**

This Work: Explain and Quantify knowledge Distillation using **Partial Information Decomposition**

### Main Contributions:
- Formally show limits of existing distillation frameworks
- Quantify the **knowledge to distill** and the **transferred knowledge** using PID
- Provide a new technique of using redundant information as a regularizer
- Propose novel distillation framework – RID – with **alternating optimization**

Y = Task, T=Teacher, S=Student
$I(Y; T) = Uni(Y:T\backslash S) + Red(Y:T,S) \rightarrow$ constant
$I(Y; S) = Uni(Y:S\backslash T) + Red(Y:T,S) \rightarrow$ need to increase

### MAIN RESULTS

**Definition 3.1 (Knowledge to distill).** The knowledge to distill from T to S is defined as $Uni(Y:T\backslash S)$, the unique information about Y that is in T but not in S.

Total information in the teacher -- $I(Y;T)$ -- is constant. Therefore,

**Definition 3.2 (Transferred knowledge).** The transferred knowledge from T to S is defined as $Red(Y:T,S)$, the redundant information about Y between T and S.

I(Y:T) — I(Y:T,S) — I(Y:S)
Uni(Y:T|S) — Red(Y:T,S) — Uni(Y:S|T)
Syn(Y:T,S)

Exact computation of PID [Bertschinger et al.'14]:
**Definition 3.3 (Unique and redundant information).** Let P be the joint distribution of Y, T and S, and $\Delta$ be the set of all joint distributions over Y × T × S. Then,
- $Uni(Y:T\backslash S) = \min_{Q\in\Delta_P} I_Q(Y;T\mid S)$
- $Red(Y:T,S) = I(Y;T) - \min_{Q\in\Delta_P} I_Q(Y;T\mid S)$

**Theorem 4.1 (Transferred knowledge lower bound).** For three random variables Y, T and S,
$$Red_\cap(Y:T,S) \leq Red(Y:T,S)$$
where $Red_\cap(Y:T,S)$ and $Red(Y:T,S)$ are defined as per Definitions 4.1 and 3.3.

Computationally efficient definition of Redundant Info.:
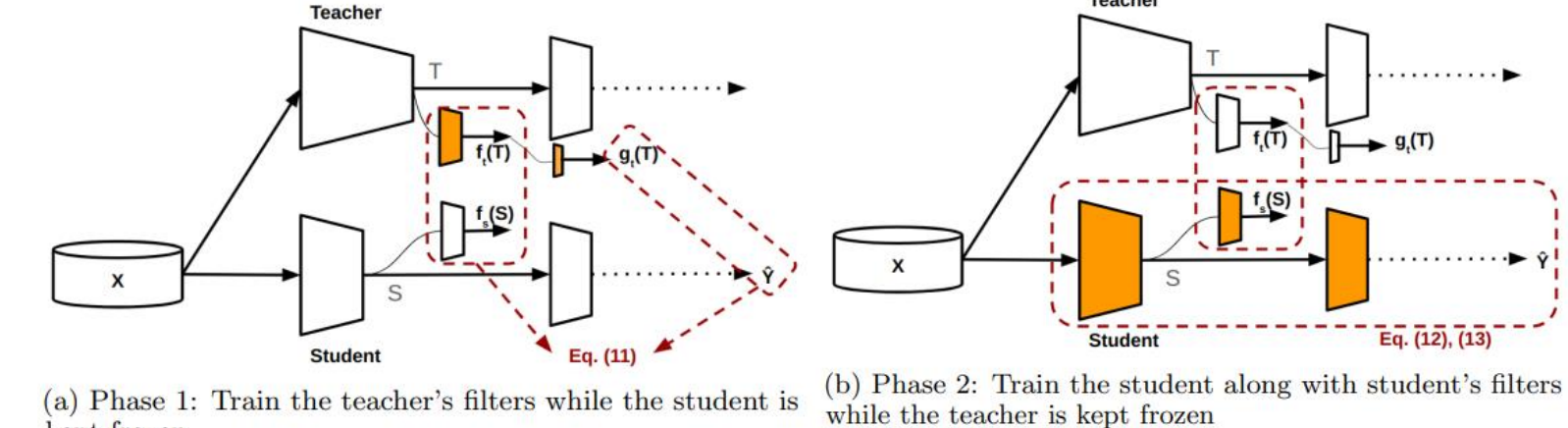**Definition 4.1 ($I_\alpha$ measure – Griffith & Ho, 2015).**
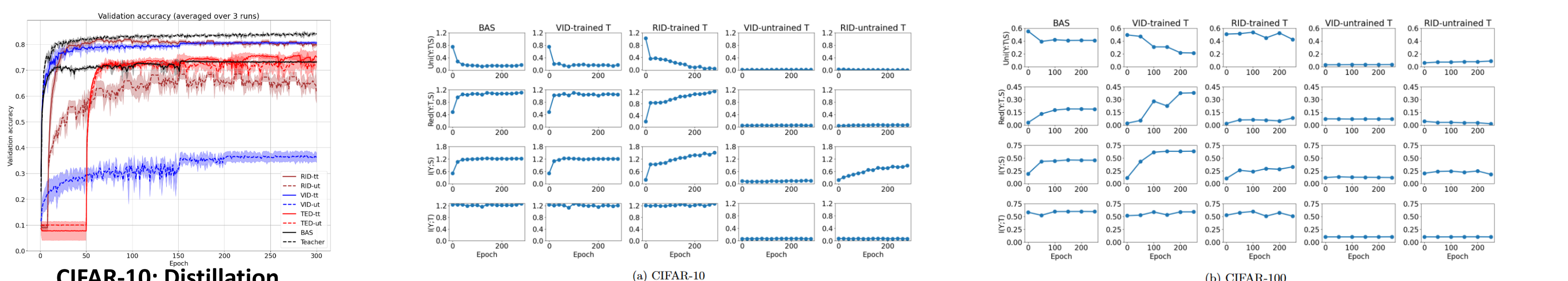$$Red_\cap(Y:T,S) = \max_{P(Q|Y)} I(Y:Q)$$
$$I(Y;Q\mid f_t(T)) = I(Y;Q\mid f_s(S)) = 0$$

**New bilevel optimization**

Set $Q = f_t(T)$ in Definition 4.1 which results in the optimization problem
$$\max_{\theta_t, \theta_s, \eta_s} I(Y; f_t(T; \theta_t)) \text{ subject to } I(Y; f_t(T; \theta_t)|f_S(S; \theta_s, \eta_s)) = 0$$
Minimize cross-entropy — Regression – minimize MSE

**Algorithm 1: Redundant Information Distillation**

Teacher — Student
(a) Phase 1: Train the teacher's filters while the student is kept frozen
(b) Phase 2: Train the student along with student's filters while the teacher is kept frozen

| Framework | Trained | Untrained |
|---|---|---|
| RID | 65% | 33% |
| VID | 70% | 11% |
| BAS | 36% | 36% |

ImageNet → CUB-200-2011: Transfer Learning

BAS — VID-trained T — RID-trained T
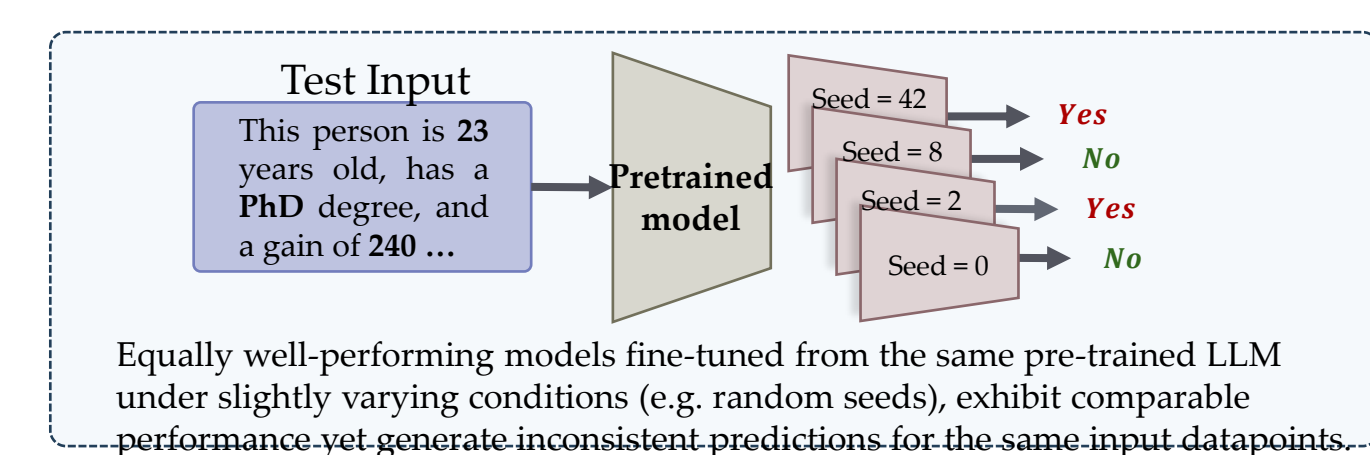VID-untrained T — RID-untrained T
(a) CIFAR-10 — (b) CIFAR-100

**CIFAR-10: Distillation**

Our strategy leads to more effective task-relevant distillation:
RID resists against nuisance or non-informative teachers, more robust to teacher instabilities

### OTHER SELECTED WORKS

**Quantifying Prediction Consistency Under Model Multiplicity in Tabular LLMs**
Hamman, Dissanayake, Mishra, Lecue, Dutta, **ICML 2025.**

Test Input: This person is 23 years old, has a PhD degree, and a gain of 240 ...
Pretrained model — Seed = 42 → Yes, Seed = 8 → No, Seed = 2 → Yes, Seed = 1 → No

Equally well-performing models fine-tuned from the same pre-trained LLM under slightly varying conditions (e.g. random seeds), exhibit comparable performance yet generate inconsistent predictions for the same input datapoints.

OURS: ~25 mins Corr: 0.88

**Few-Shot Knowledge Distillation of LLMs With Counterfactual Explanations**
Hamman, Dissanayake, Fu, Dutta. **In Review.**

Experiments on DeBERT-v3 and Qwen2.5 families and 6 benchmark datasets (NLP tasks)

**Achieves More With Less:** Improves accuracy in few-shot settings with as low as 8, 16, 32 samples

**References:**
[1] Y. Wang, H. Qian, and C. Miao. DualCF: Efficient model extraction attack from counterfactual explanations. In ACM FAccT 2022.
[2] C. Yadav, M. Moshkovitz, and K. Chaudhuri. Xaudit : A theoretical look at auditing with explanations. arXiv:2206.04740, 2023.
[3] N. Bertschinger, J. Rauh, E. Olbrich, J. Jost, N. Ay, Quantifying Unique Information. Entropy, 2014.
[4] V. Griffith and T. Ho. Quantifying redundant information in predicting a target random variable. Entropy, 2015.
[5] S. Ahn, S. X. Hu, A. Damianou, N. D. Lawrence, & Z. Dai, Variational information distillation for knowledge transfer. IEEE CVF 2019.
[6] P. P. Liang et al., Quantifying & modeling multimodal interactions: An information decomposition framework. NeurIPS 2023.