# Understanding Dataset Difficulty with $\mathscr{V}$-Usable Information
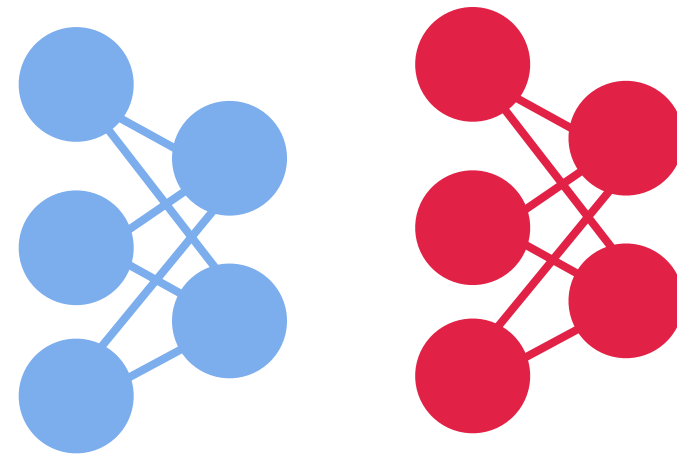
## ICML 2022

**Kawin Ethayarajh**

**Yejin Choi**
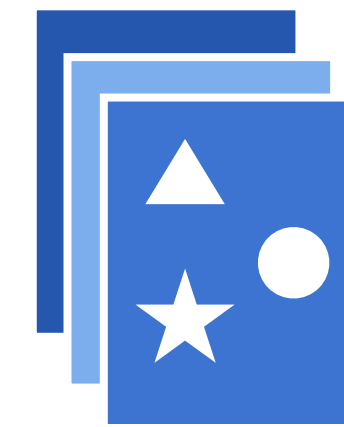
**Swabha Swayamdipta**

compare models $\mathscr{V}$

compare datasets $(X, Y)$

compare attributes $X_i$

compare instances $(x, y)$

compare slices $\{(x, y)\}_i$

compare models $\mathscr{V}$

compare datasets $(X, Y)$

compare attributes $X_i$

accuracy, F1
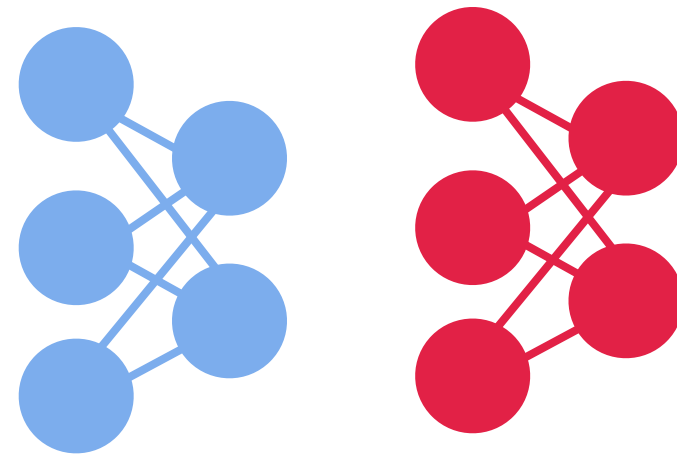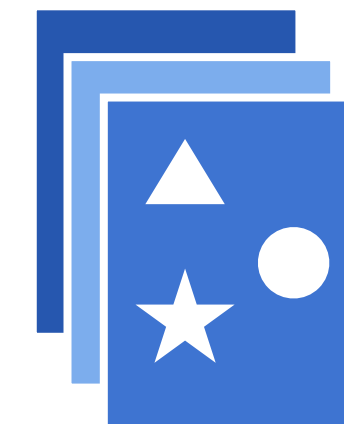
compare instances $(x, y)$

compare slices $\{(x, y)\}_i$

compare models $\mathscr{V}$

compare datasets $(X, Y)$

compare attributes $X_i$
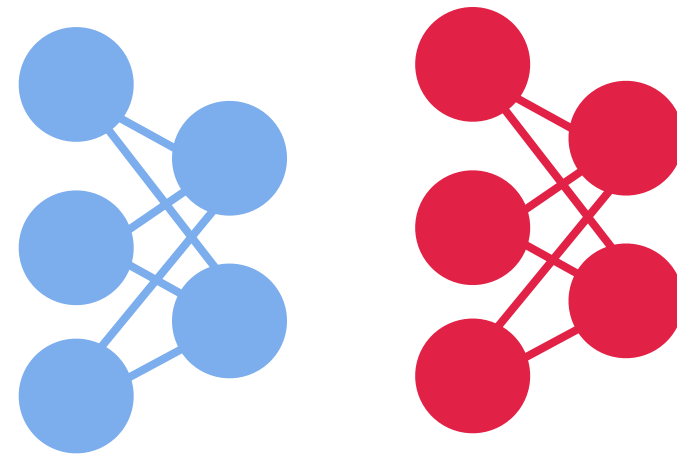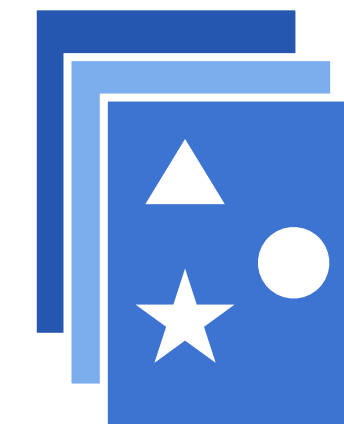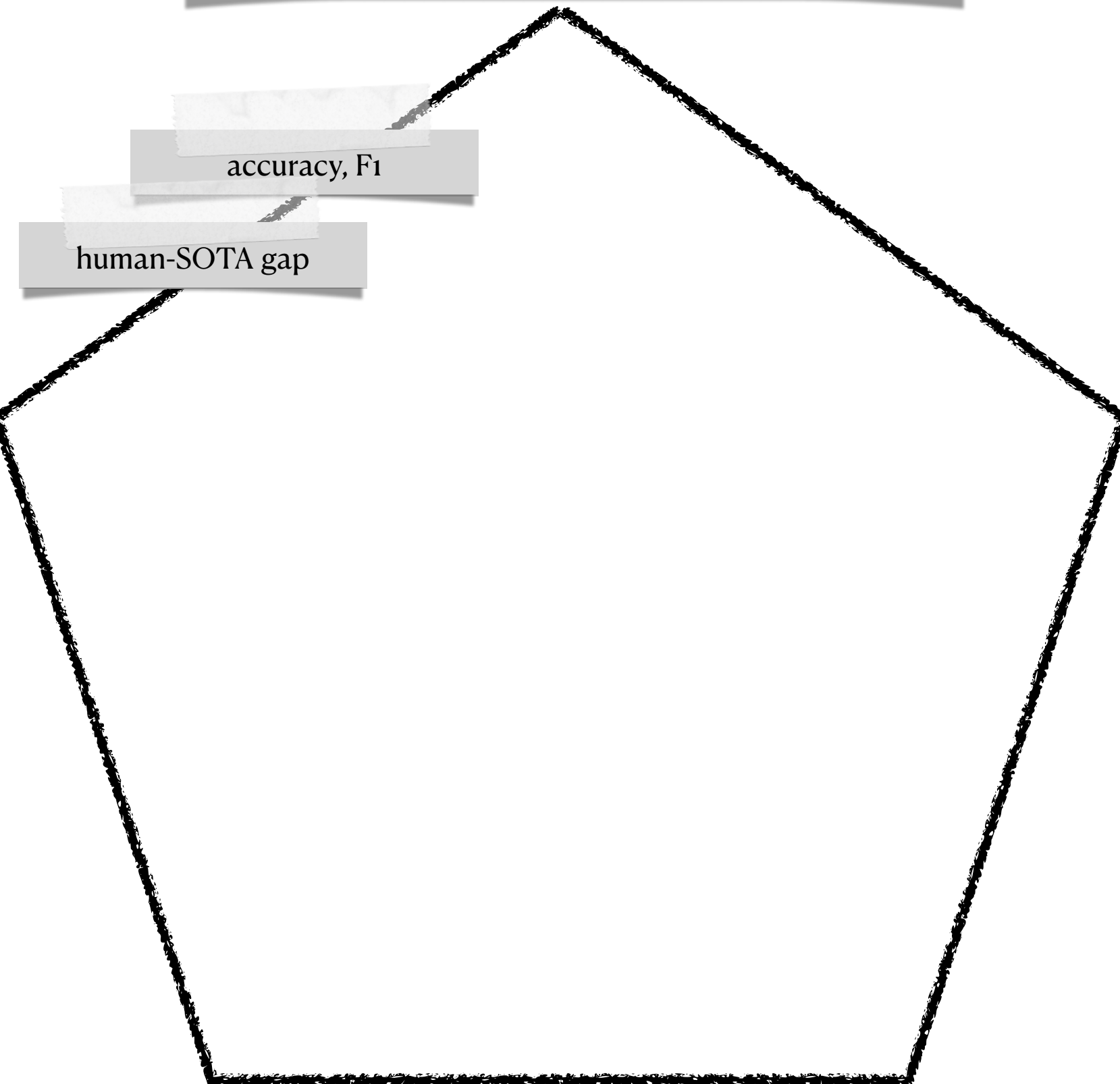
accuracy, F1

human-SOTA gap

compare instances $(x, y)$

compare slices $\{(x, y)\}_i$

2

compare models $\mathscr{V}$

compare datasets $(X, Y)$

compare attributes $X_i$

Dynascore (Ma et al., 2020)

accuracy, F1

(O'Connor & Andreas, 2021)

human-SOTA gap

DIME (Zhang et al., 2020)

IRT (Rodriguez et al., 2021)

MDL (Perez et al., 2021)

(Suguwara et al., 2018)

(Sen & Saffari, 2020)

(Rondeau & Hazen, 2018)

compare instances $(x, y)$

Cartography (Swayamdipta et al., 2020)

RHO-LOSS (Mindermann et al., 2022)

Data Shapley (Ghorbani & Zhou, 2019)

compare slices $\{(x, y)\}_i$

compare models $\mathcal{V}$
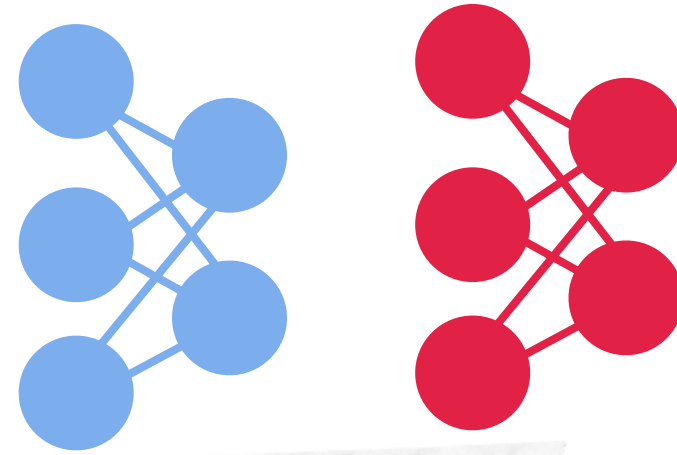
compare datasets $(X, Y)$

compare attributes $X_i$

Dynascore (Ma et al., 2020)

accuracy, F1

(O'Connor & Andreas, 2021)

human-SOTA gap

DIME (Zhang et al., 2020)

IRT (Rodriguez et al., 2021)

MDL (Perez et al., 2021)

$\mathcal{V}$-Usable Information

(Suguwara et al., 2018)

(Sen & Saffari, 2020)

(Rondeau & Hazen, 2018)

compare instances $(x, y)$

compare slices $\{(x, y)\}_i$

Cartography (Swayamdipta et al., 2020)

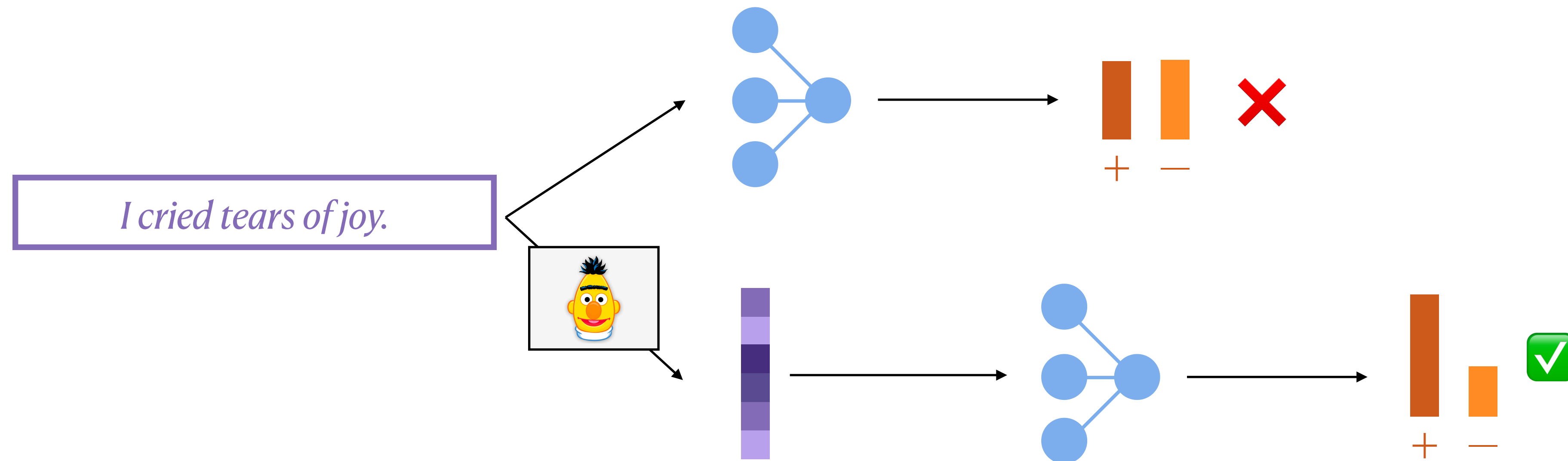RHO-LOSS (Mindermann et al., 2022)

Data Shapley (Ghorbani & Zhou, 2019)

3

# Transforming the input with $\tau$ can make information previously *unusable* by model family $\mathcal{V}$ now *usable*, despite $I(X; Y) \geq I(\tau(X); Y)$.



[ Xu et al., 2019 ]

**The predictive $\mathscr{V}$-information framework can be used to measure the amount of usable information $X$ contains about $Y$ w.r.t. $\mathscr{V}$.**

$$I_{\mathscr{V}}(X \to Y) = \underbrace{\inf_{f \in \mathscr{V}} \mathbb{E}[-\log_2 f[\varnothing](Y)]}_{\color{blue}{H_{\mathscr{V}}(Y)}} - \underbrace{\inf_{f \in \mathscr{V}} \mathbb{E}[-\log_2 f[X](Y)]}_{\color{red}{H_{\mathscr{V}}(Y|X)}}$$

[ Xu et al., 2019 ]

# The predictive $\mathscr{V}$-information framework can be used to measure the amount of usable information $X$ contains about $Y$ w.r.t. $\mathscr{V}$.

$$I_{\mathscr{V}}(X \to Y) = \underbrace{\inf_{f \in \mathscr{V}} \mathbb{E}[-\log_2 f[\varnothing](Y)]}_{\color{blue}{H_{\mathscr{V}}(Y)}} - \underbrace{\inf_{f \in \mathscr{V}} \mathbb{E}[-\log_2 f[X](Y)]}_{\color{red}{H_{\mathscr{V}}(Y|X)}}$$

train/finetune on null input $\varnothing$

[ Xu et al., 2019 ]

# The predictive $\mathcal{V}$-information framework can be used to measure the amount of usable information $X$ contains about $Y$ w.r.t. $\mathcal{V}$.

$$I_{\mathcal{V}}(X \to Y) = \underbrace{\inf_{f \in \mathcal{V}} \mathbb{E}[-\log_2 f[\varnothing](Y)]}_{\textcolor{blue}{H_{\mathcal{V}}(Y)}} - \underbrace{\inf_{f \in \mathcal{V}} \mathbb{E}[-\log_2 f[X](Y)]}_{\textcolor{red}{H_{\mathcal{V}}(Y|X)}}$$

train/finetune on <span style="color:blue">null input $\varnothing$</span>     train/finetune on <span style="color:red">actual input $X$</span>

[ Xu et al., 2019 ]

**The predictive $\mathscr{V}$-information framework can be used to measure the amount of usable information $X$ contains about $Y$ w.r.t. $\mathscr{V}$.**

$$I_{\mathscr{V}}(X \rightarrow \qquad X](Y)]$$

The lower the $\mathscr{V}$-usable information, the more difficult the dataset is for $\mathscr{V}$.

train/finetune on null input $\varnothing$      train/finetune on actual input $X$

[ Xu et al., 2019 ]

# SNLI

[Bowman et al., 2015]

natural language inference

PREMISE: Women enjoying a game of table tennis.

HYPOTHESIS: Women enjoying a game of ping pong.

● entailment
○ neutral
○ contradiction

# MultiNLI

[Williams et al., 2018]

natural language inference

PREMISE: The Old One always comforted Ca'daan, except today.

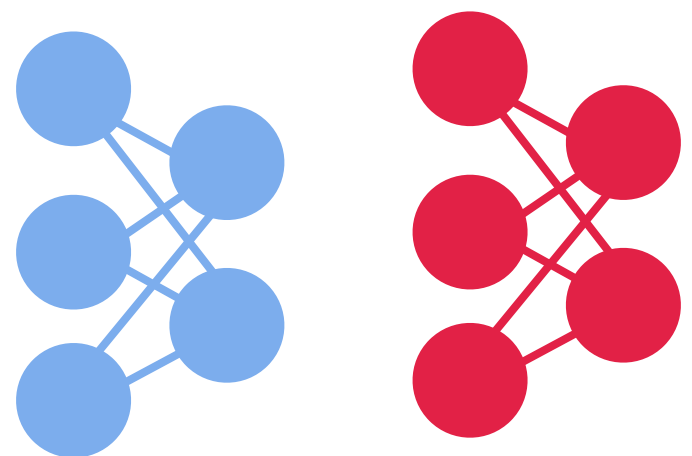HYPOTHESIS: Ca'daan knew the Old One very well.

○ entailment
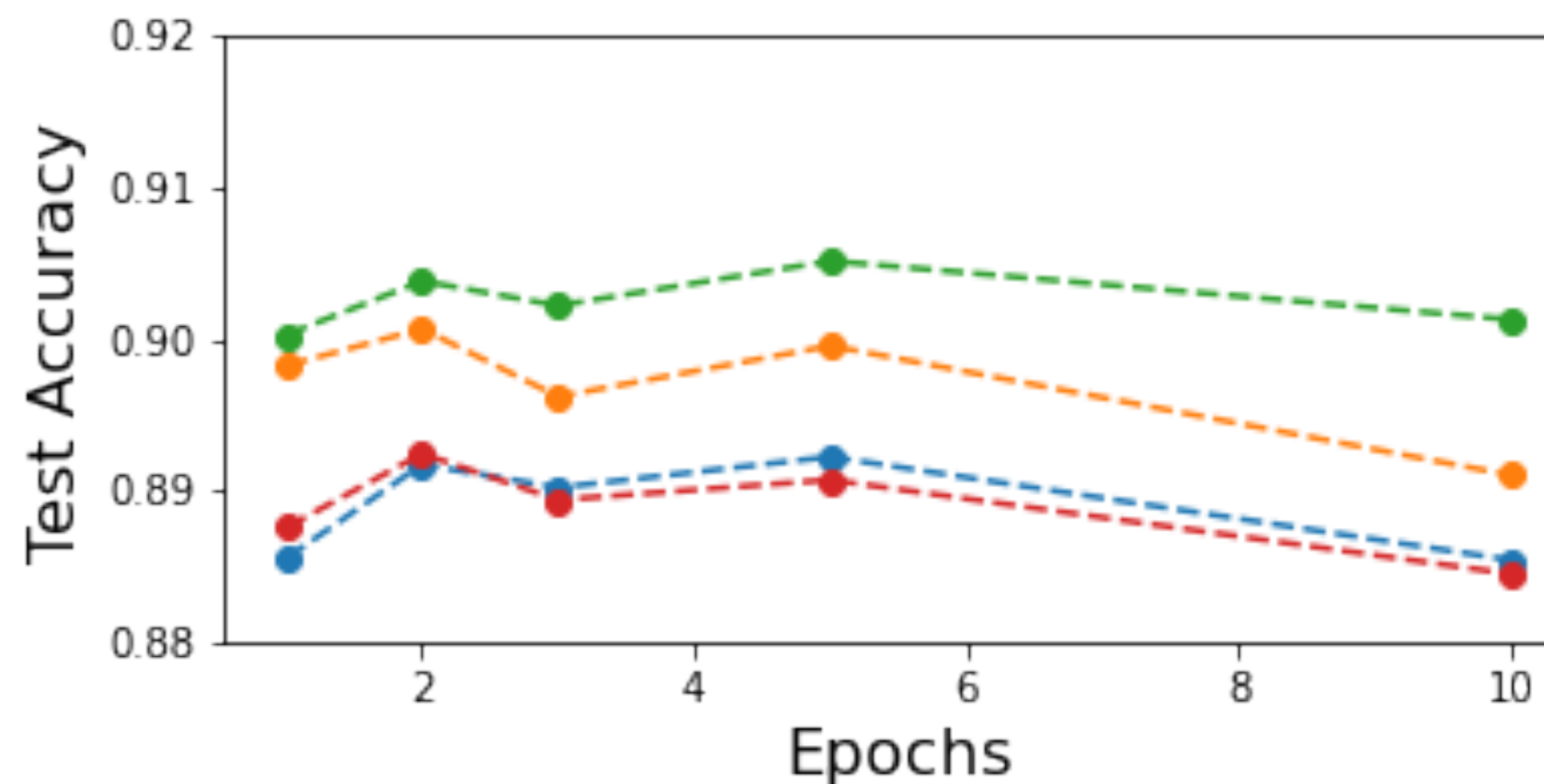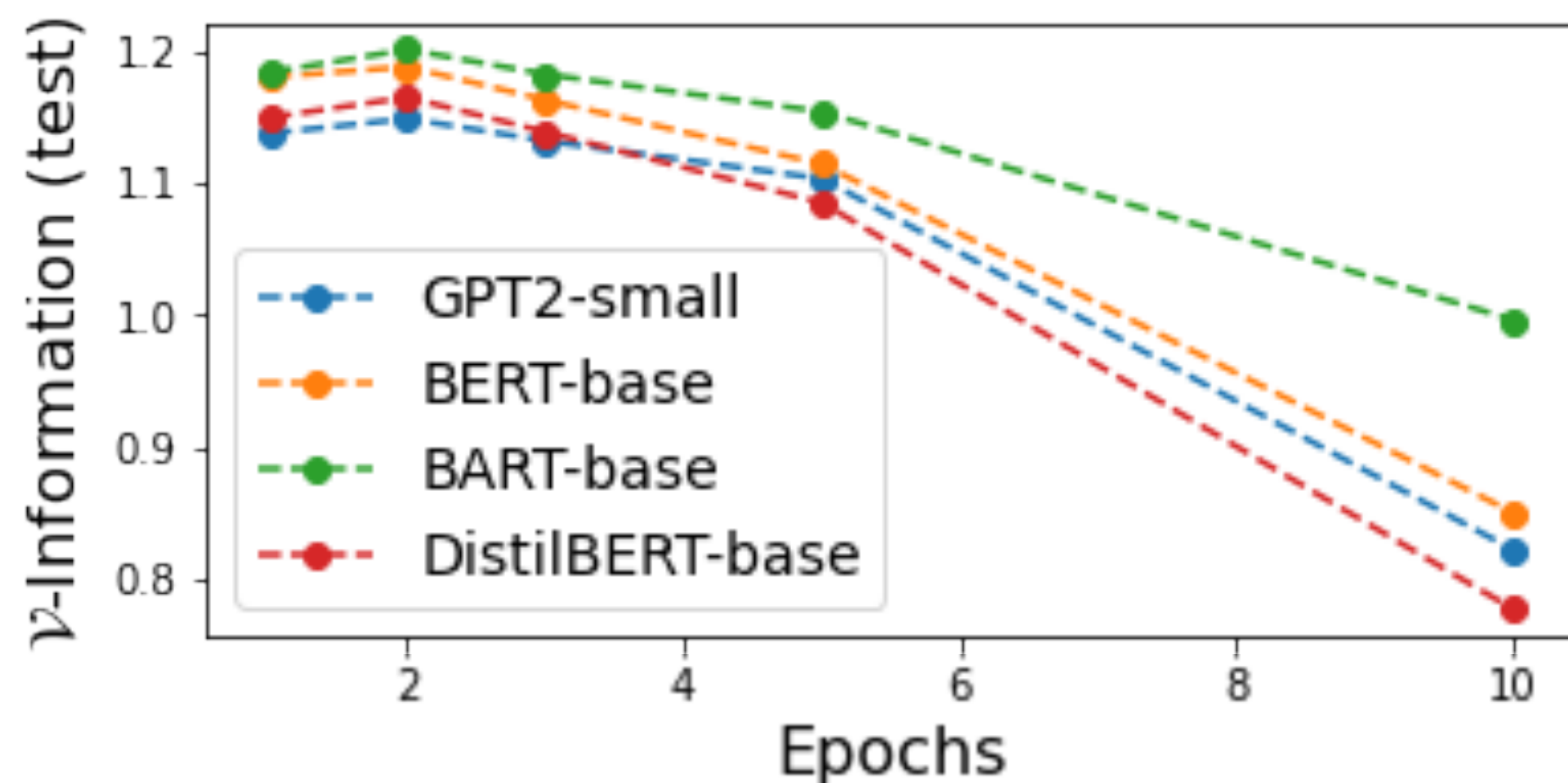● neutral
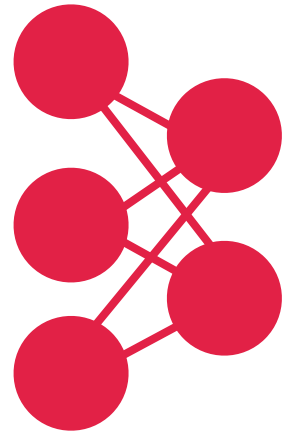○ contradiction

# CoLA

[Warstadt et al., 2018]

text classification
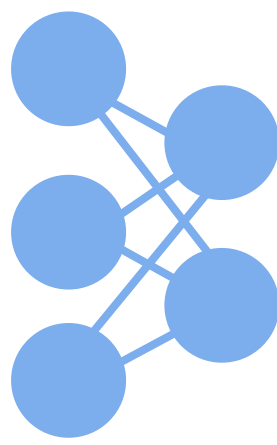
Wash you.

○ grammatical
● ungrammatical

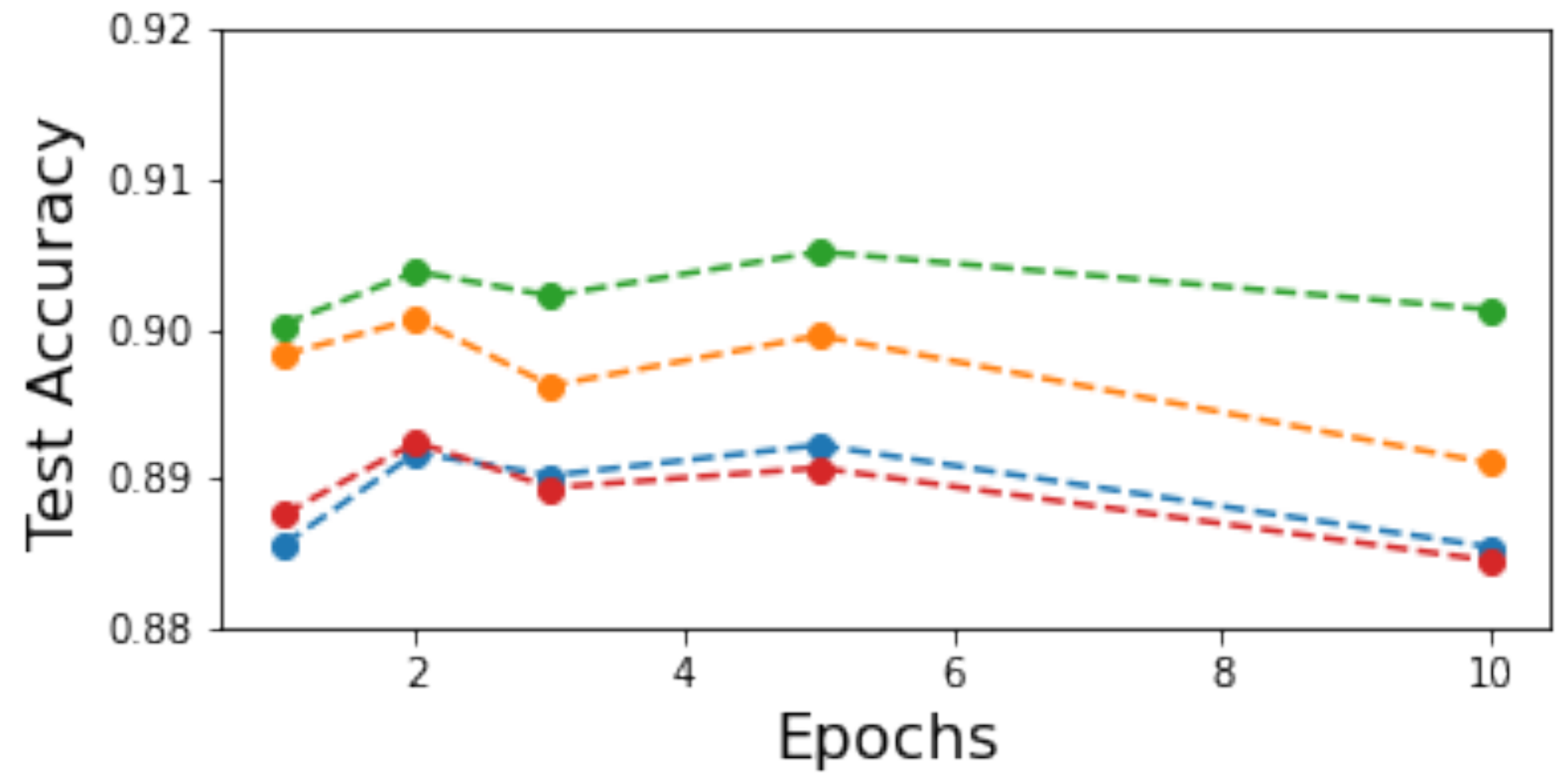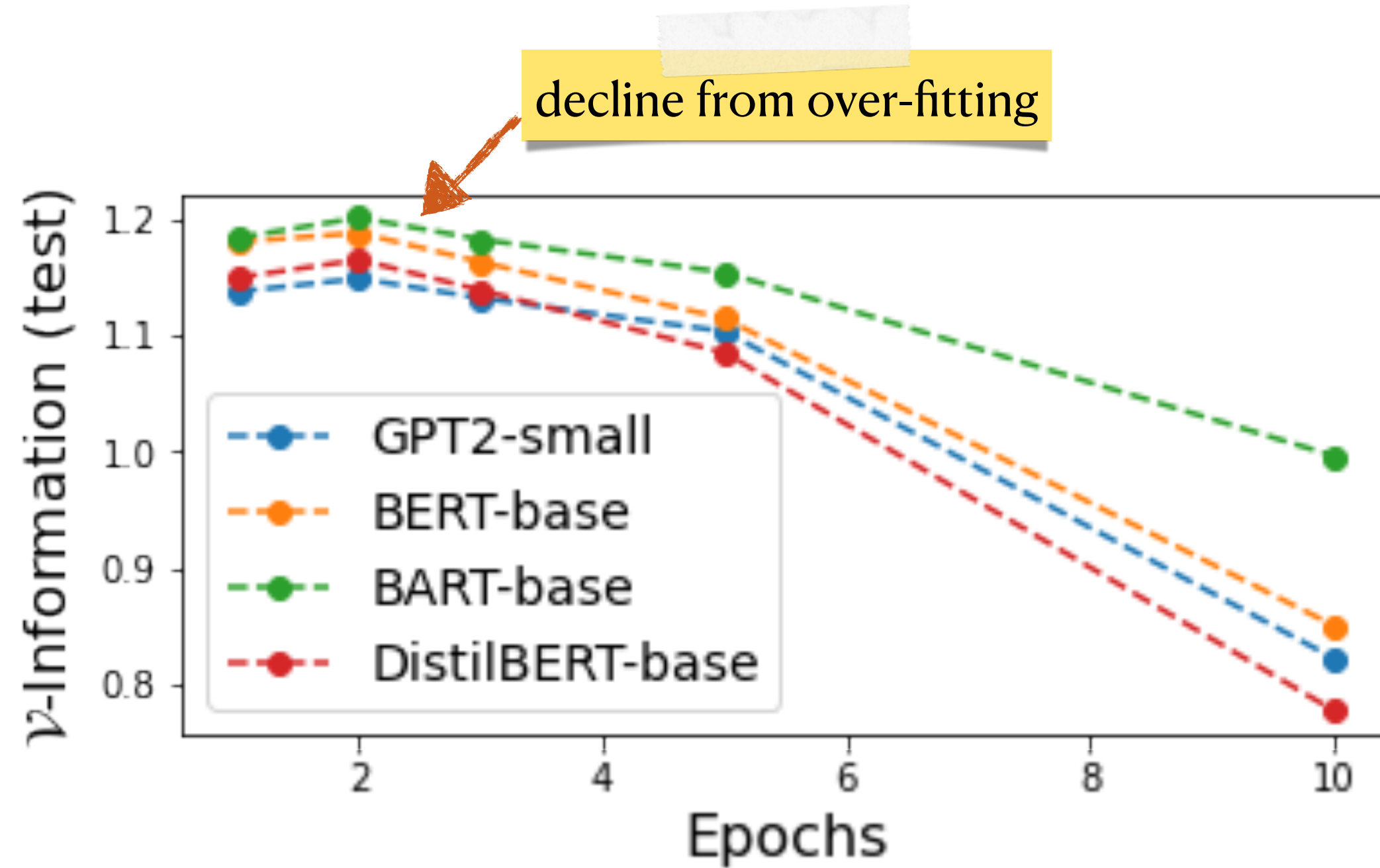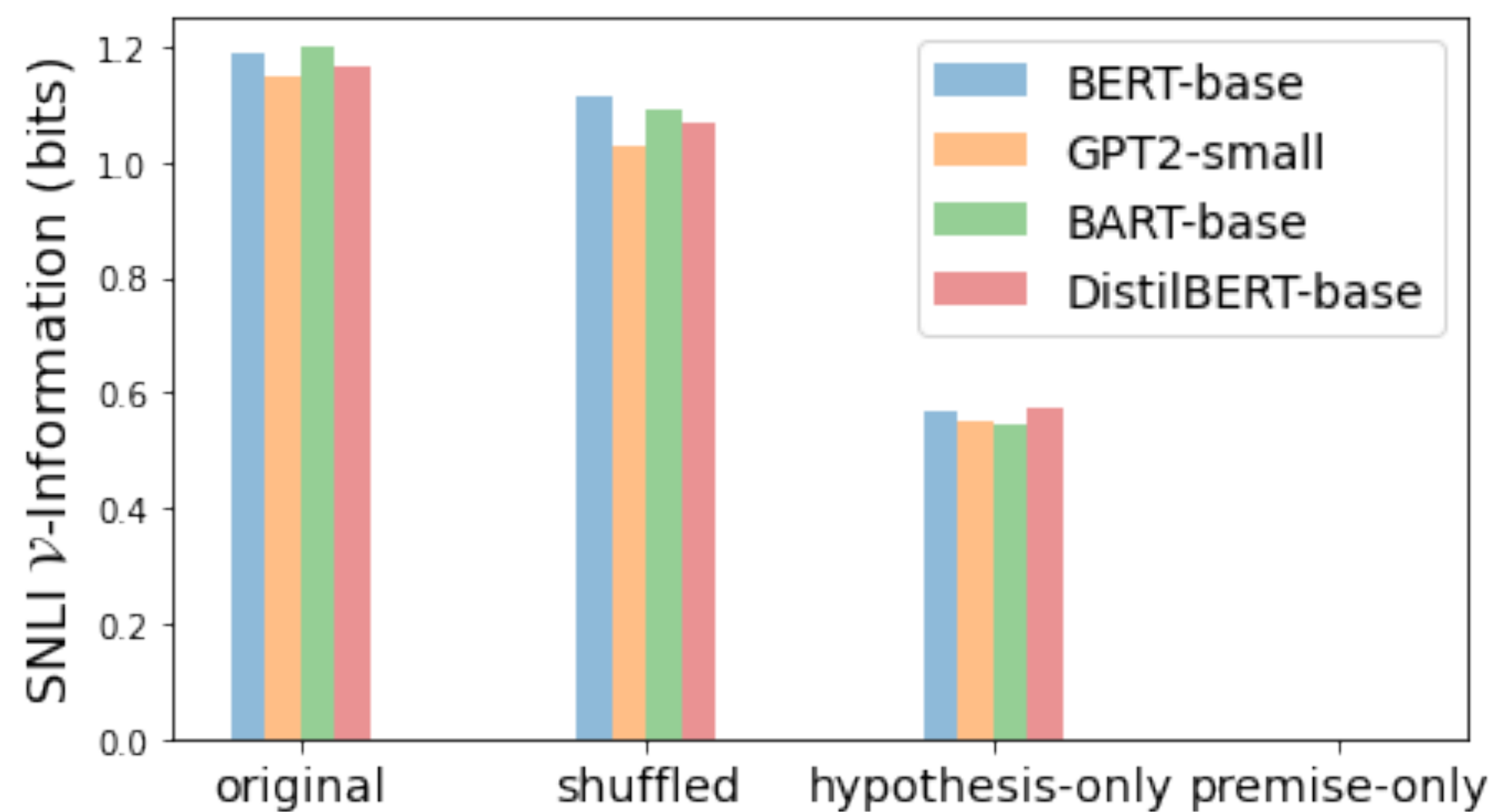**Compare different models $\mathscr{V}$ by computing $I_\mathscr{V}(X \to Y)$ for the same $(X, Y)$, shown here for SNLI.**

**Compare different models $\mathscr{V}$ by computing $I_{\mathscr{V}}(X \to Y)$ for the same $(X, Y)$, shown here for SNLI.**
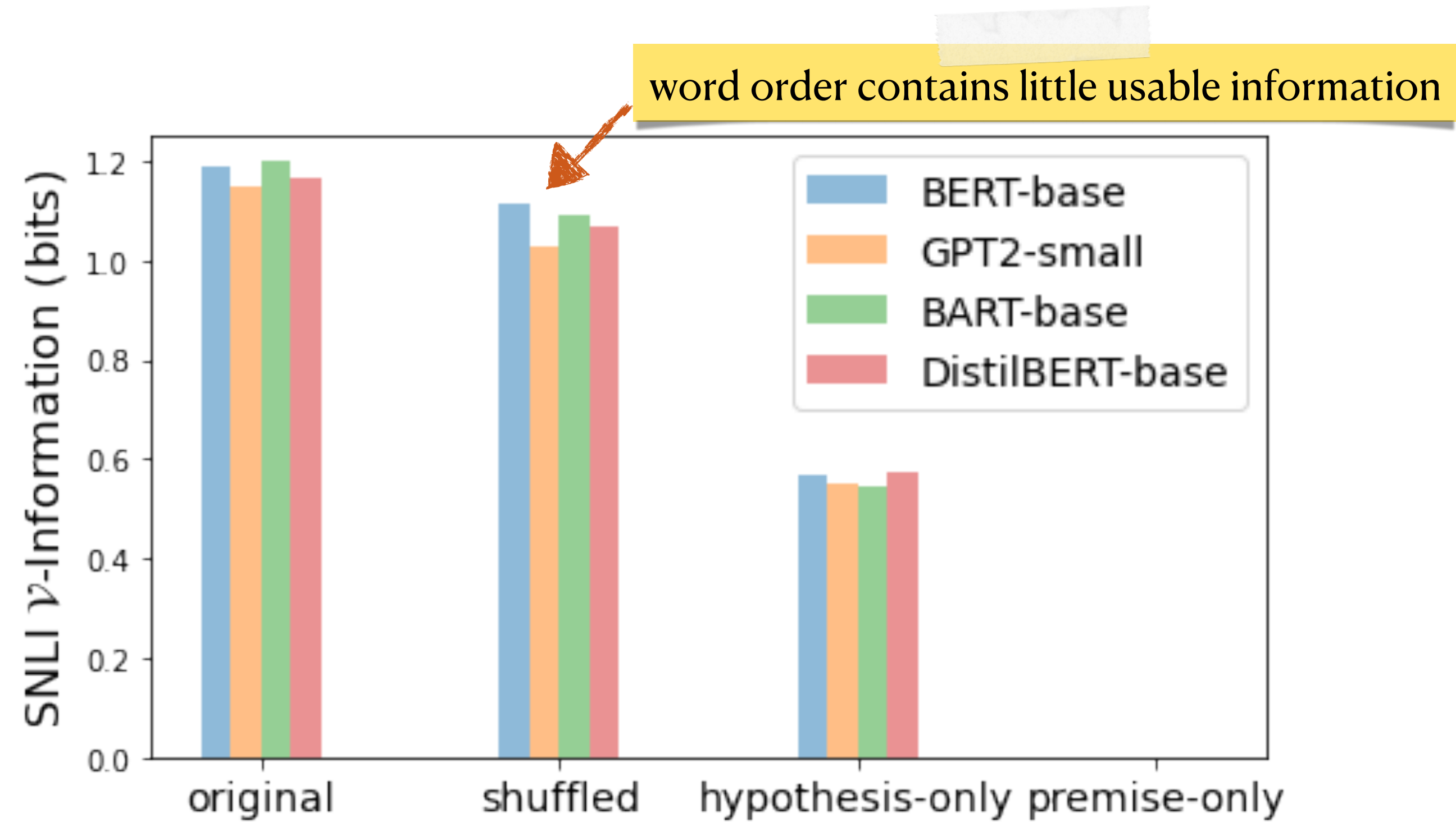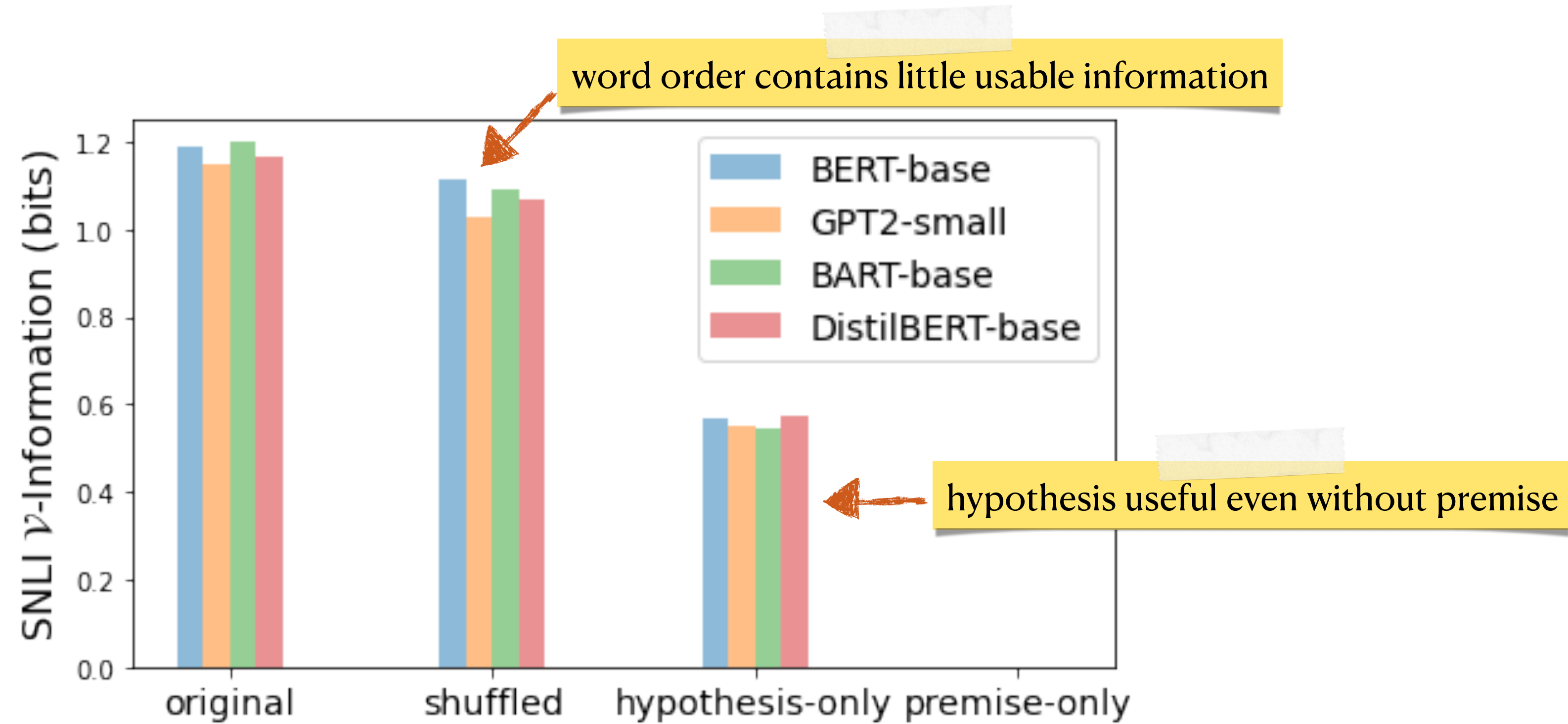


decline from over-fitting

**Compare different input attributes $X_i$ by computing $I_{\mathcal{V}}(X_i \rightarrow Y)$ for the same $Y, \mathcal{V}$.**

**Compare different input attributes $X_i$ by computing $I_{\mathscr{V}}(X_i \rightarrow Y)$ for the same $Y, \mathscr{V}$.**



word order contains little usable information

**Compare different input attributes $X_i$ by computing $I_{\mathscr{V}}(X_i \to Y)$ for the same $Y, \mathscr{V}$.**



word order contains little usable information

hypothesis useful even without premise

**We can measure instance-level difficulty (w.r.t. a distribution) with pointwise $\mathscr{V}$-information (PVI), the analogue of PMI.**

$$I_{\mathscr{V}}(X \to Y) = \mathbb{E}_{x,y \sim P(X,Y)}[\text{PVI}(x \to y)]$$

# We can measure instance-level difficulty (w.r.t. a distribution) with pointwise $\mathscr{V}$-information (PVI), the analogue of PMI.

$$I_{\mathscr{V}}(X \to Y) = \mathbb{E}_{x,y \sim P(X,Y)}[\text{PVI}(x \to y)]$$

$$I_{\mathscr{V}}(X \to Y) \in \mathbb{R}^{0+}; \ \text{PVI}(x \to y) \in \mathbb{R}$$

# We can measure instance-level difficulty (w.r.t. a distribution) with pointwise $\mathscr{V}$-information (PVI), the analogue of PMI.

$$I_{\mathscr{V}}(X \to Y) = \mathbb{E}_{x,y \sim P(X,Y)}[\text{PVI}(x \to y)]$$

$I_{\mathscr{V}}(X \to Y) \in \mathbb{R}^{0+}; \ \text{PVI}(x \to y) \in \mathbb{R}$

cross-epoch Pearson's $r \geq 0.747$

# We can measure instance-level difficulty (w.r.t. a distribution) with pointwise $\mathcal{V}$-information (PVI), the analogue of PMI.

$$I_{\mathcal{V}}(X \to Y) = \mathbb{E}_{x,y \sim P(X,Y)}[\text{PVI}(x \to y)]$$

$I_{\mathcal{V}}(X \to Y) \in \mathbb{R}^{0+}; \; \text{PVI}(x \to y) \in \mathbb{R}$

cross-epoch Pearson's $r \geq 0.747$

cross-seed Pearson's $r \geq 0.877$

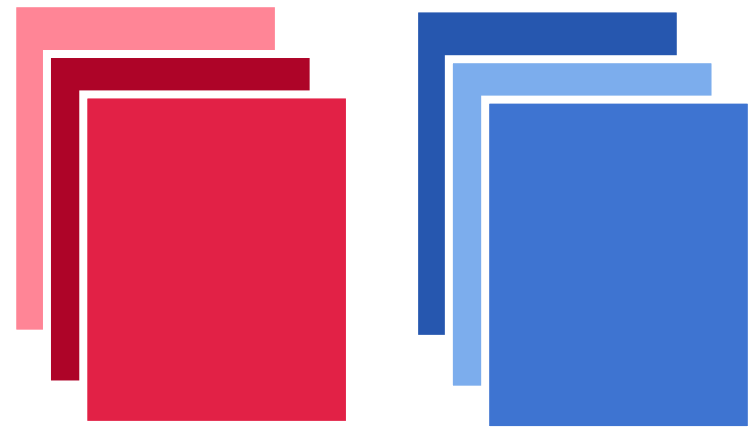**We can measure instance-level difficulty (w.r.t. a distribution) with pointwise $\mathscr{V}$-in**

The higher the PVI, the easier the instance is for $\mathscr{V}$ w.r.t. $P(X, Y)$.

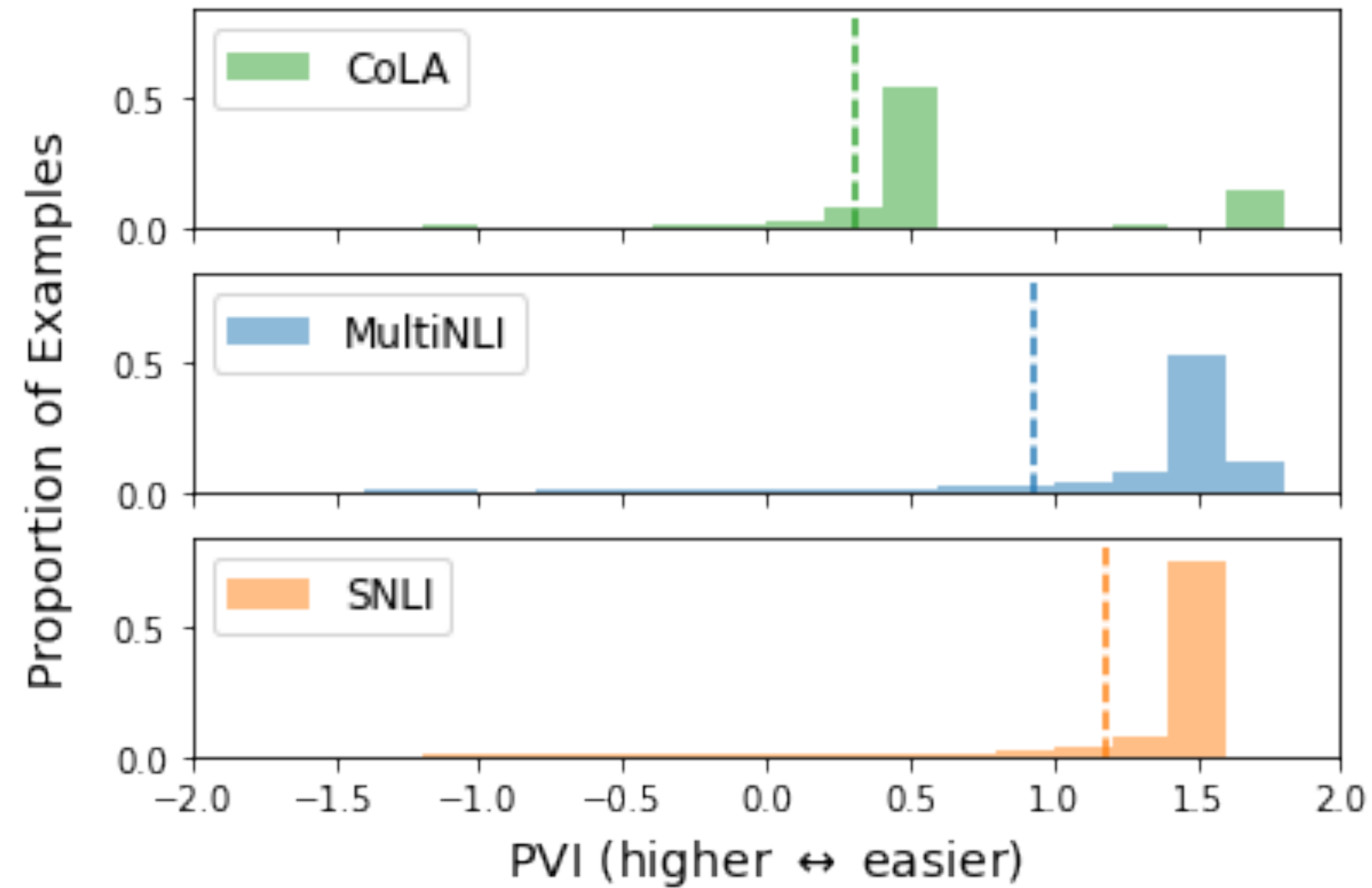$$I_{\mathscr{V}}(X \rightarrow Y) = \mathbb{E}_{x,y \sim P(X,Y)}[\text{PVI}(x \rightarrow y)]$$

$I_{\mathscr{V}}(X \rightarrow Y) \in \mathbb{R}^{0+}; \ \text{PVI}(x \rightarrow y) \in \mathbb{R}$
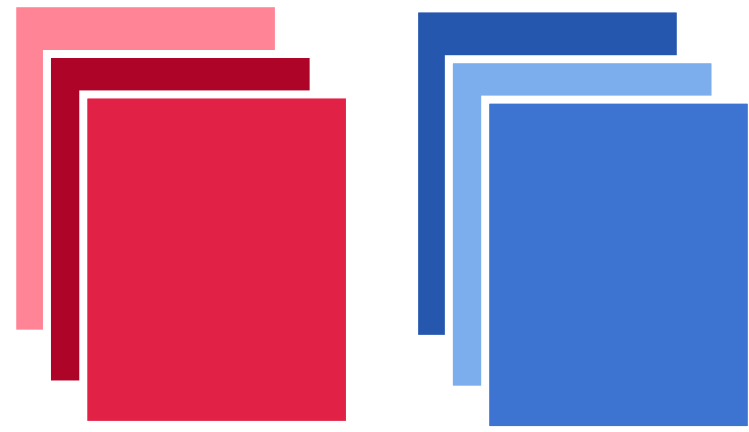
cross-epoch Pearson's $r \geq 0.747$
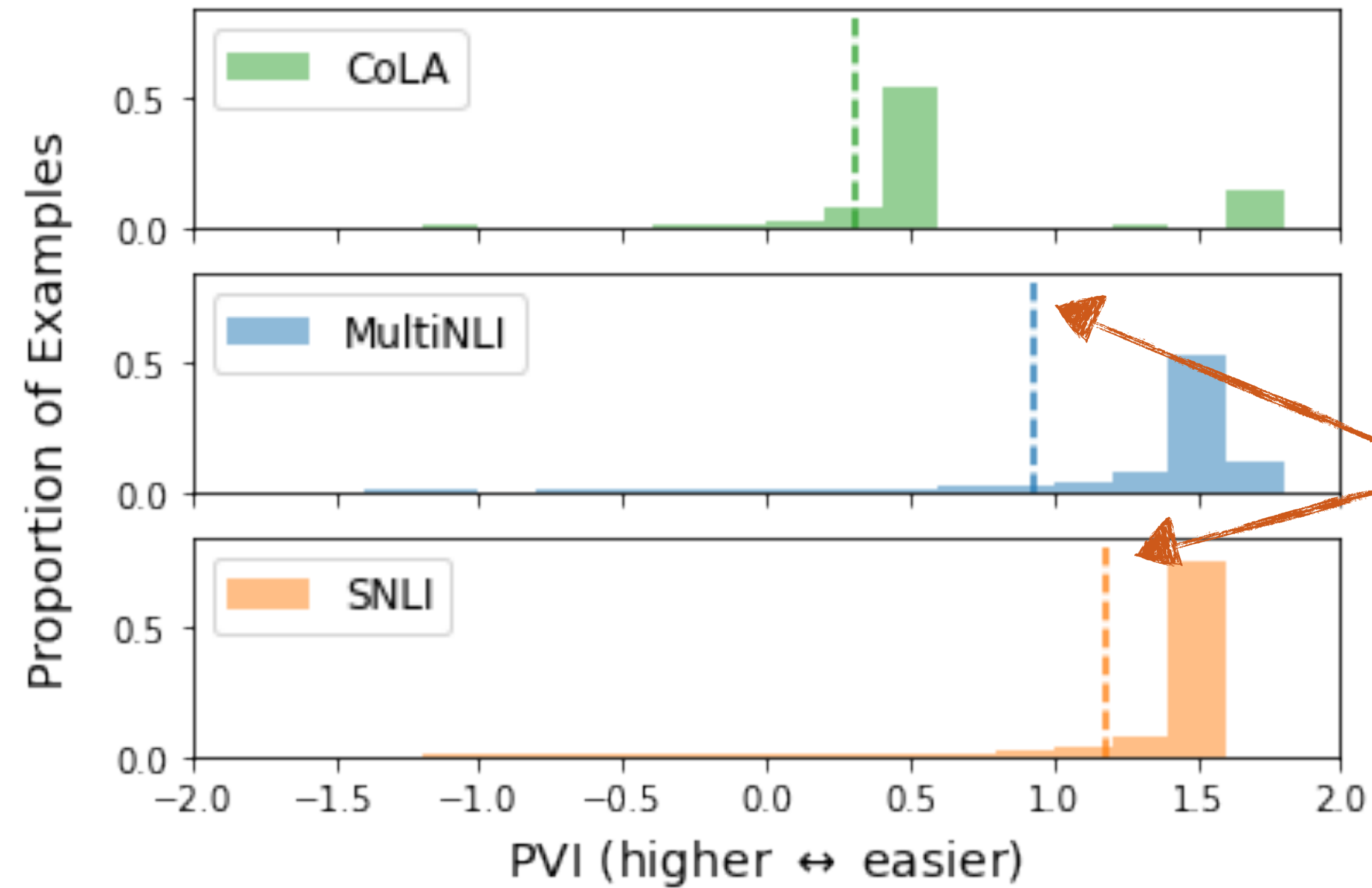
cross-seed Pearson's $r \geq 0.877$

**Compare different datasets $(X, Y)$ by estimating $I_\mathcal{V}(X \to Y)$ and PVI$(x \to y)$ for the same $\mathcal{V}$ across datasets.**

**Compare different datasets $(X, Y)$ by estimating $I_{\mathcal{V}}(X \to Y)$ and PVI$(x \to y)$ for the same $\mathcal{V}$ across datasets.**

# Compare **different instances** $(x, y)$ using $\text{PVI}(x \rightarrow y)$ for the same $\mathscr{V}, X, Y$, before and after transformations.

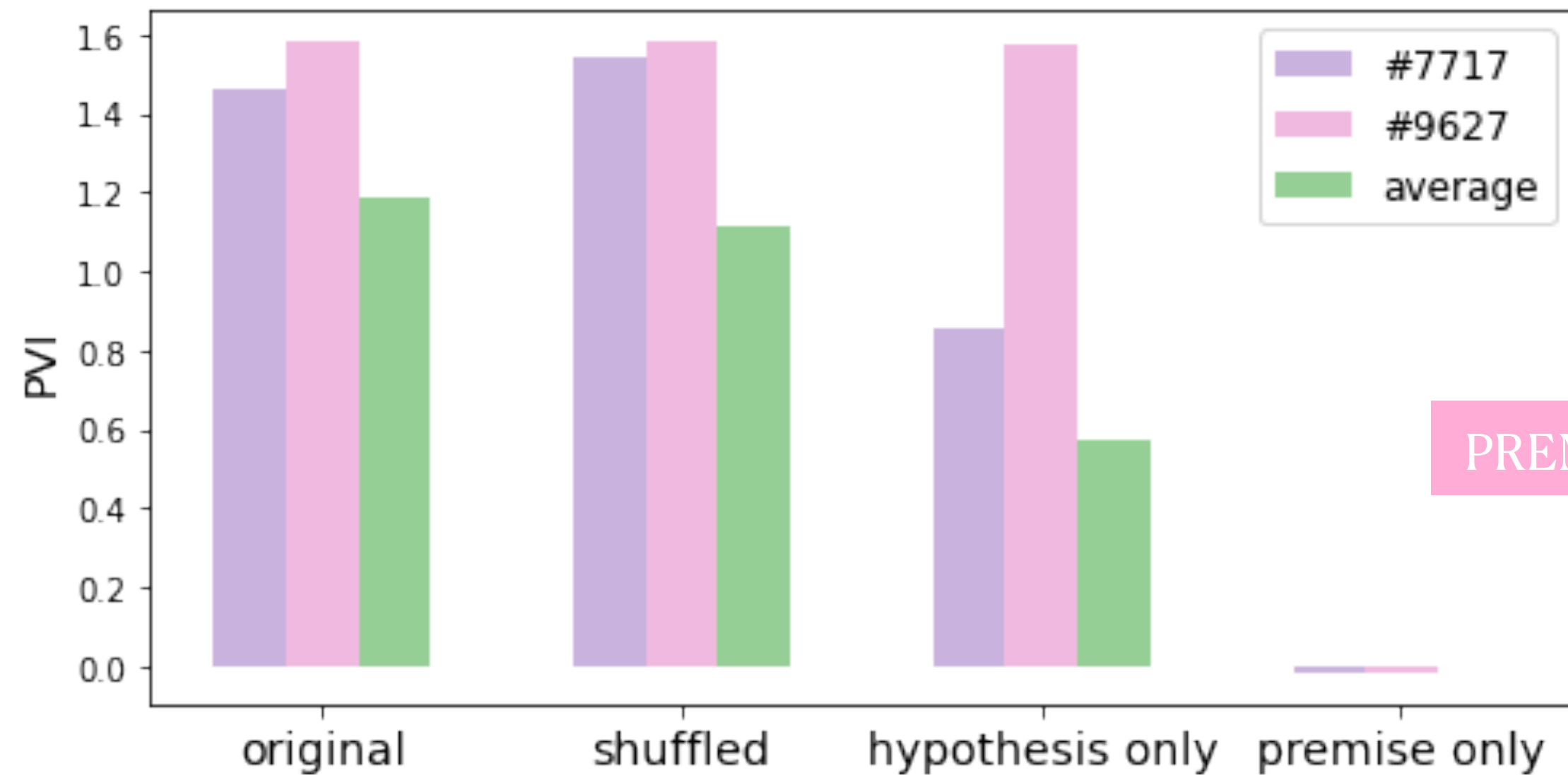PREMISE: Little kids play a game of running around a pole.

HYPOTHESIS: The kids are fighting outside.

PREMISE: A group of people watching a boy getting interviewed by a man.

HYPOTHESIS: A group of people are sleeping on Pluto.

**Compare different instances $(x, y)$ using PVI$(x \to y)$ for the same $\mathscr{V}, X, Y$, before and after transformations.**



PREMISE: Little kids play a game of running around a pole.

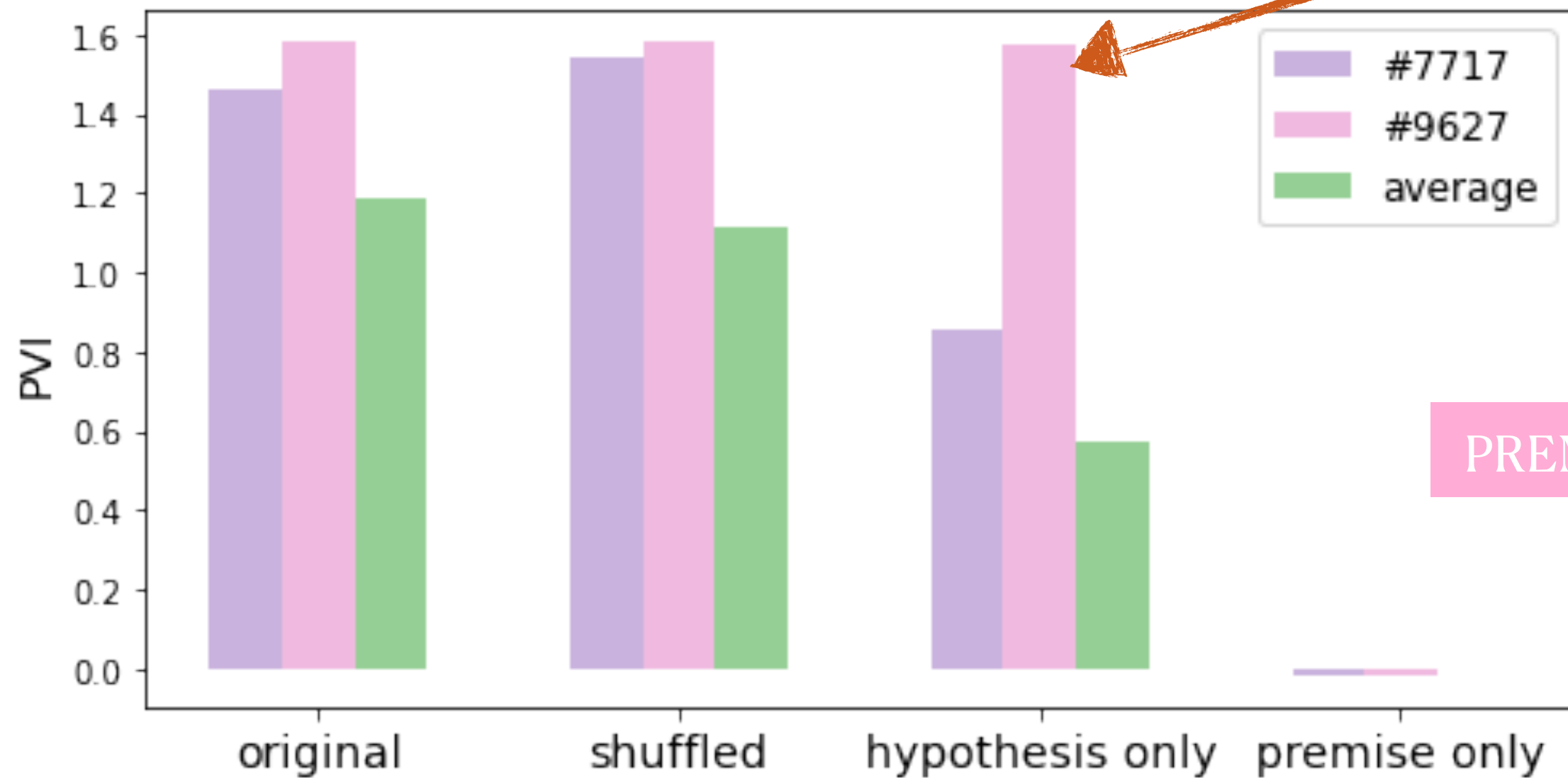HYPOTHESIS: The kids are fighting outside.

PREMISE: A group of people watching a boy getting interviewed by a man.

HYPOTHESIS: A group of people are sleeping on Pluto.

**Compare different instances** $(x, y)$ **using PVI**$(x \rightarrow y)$ **for the same** $\mathscr{V}, X, Y,$ **before and after transformations.**



hypothesis is what makes #9627 easier!

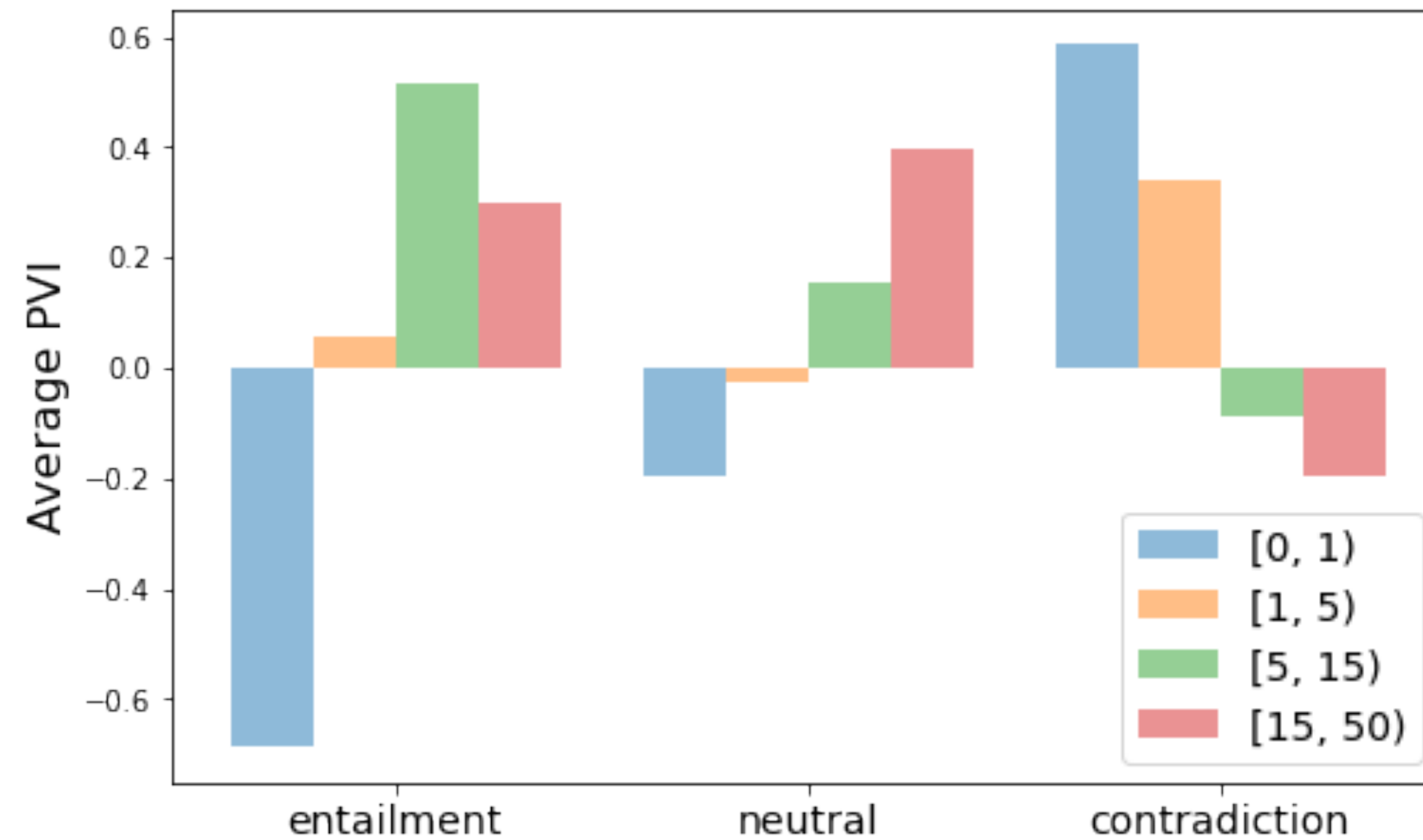PREMISE: Little kids play a game of running around a pole.

HYPOTHESIS: The kids are fighting outside.

PREMISE: A group of people watching a boy getting interviewed by a man.

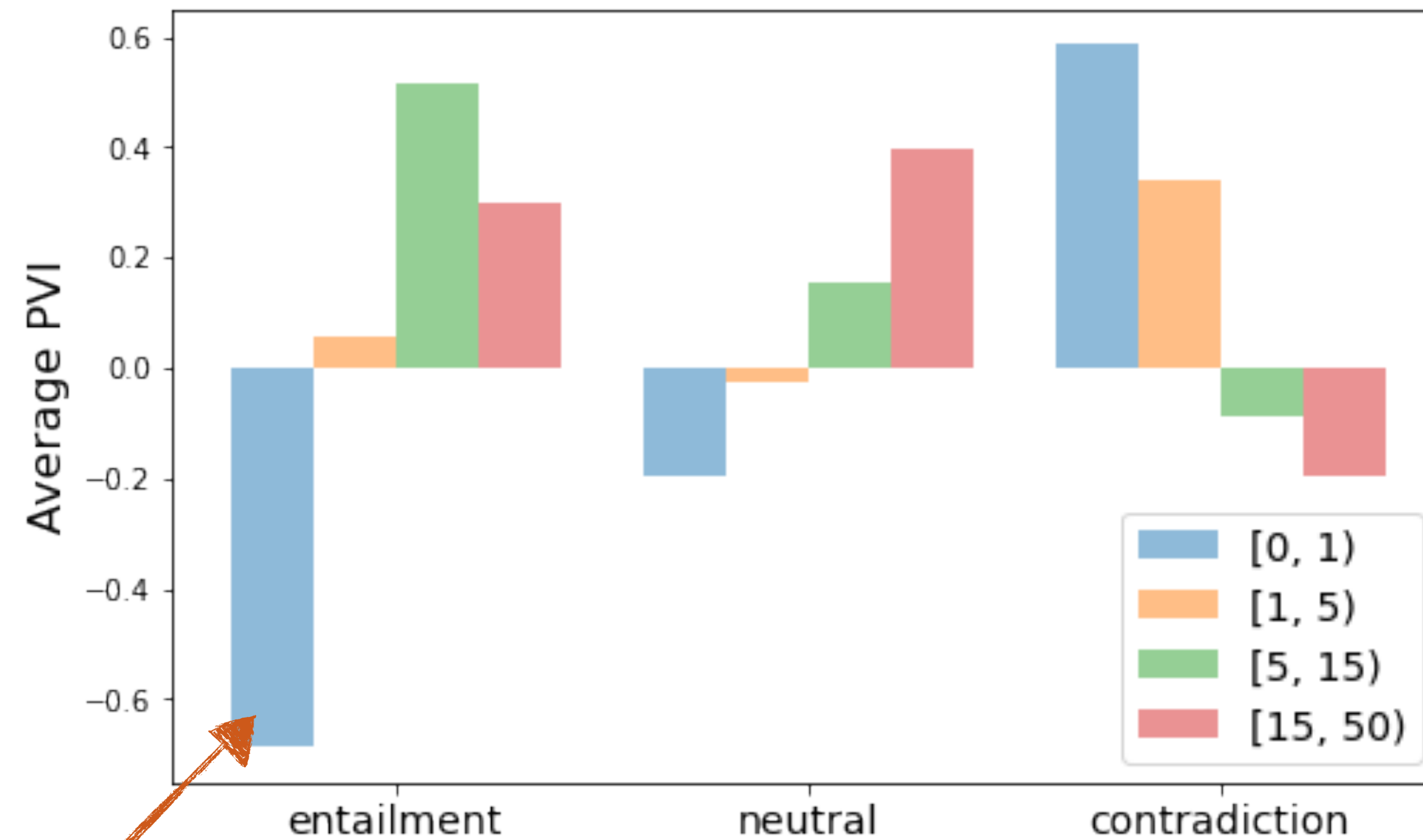HYPOTHESIS: A group of people are sleeping on Pluto.

**Compare different slices $\{(x, y)\}_i$ by estimating the average PVI$(x \rightarrow y)$ for each slice.**

**Compare different slices $\{(x, y)\}_i$ by estimating the average PVI$(x \rightarrow y)$ for each slice.**



what BERT finds hardest!

# Estimating the drop in $\mathscr{V}$-information after leaving out a token reveals token-level annotation artefacts.
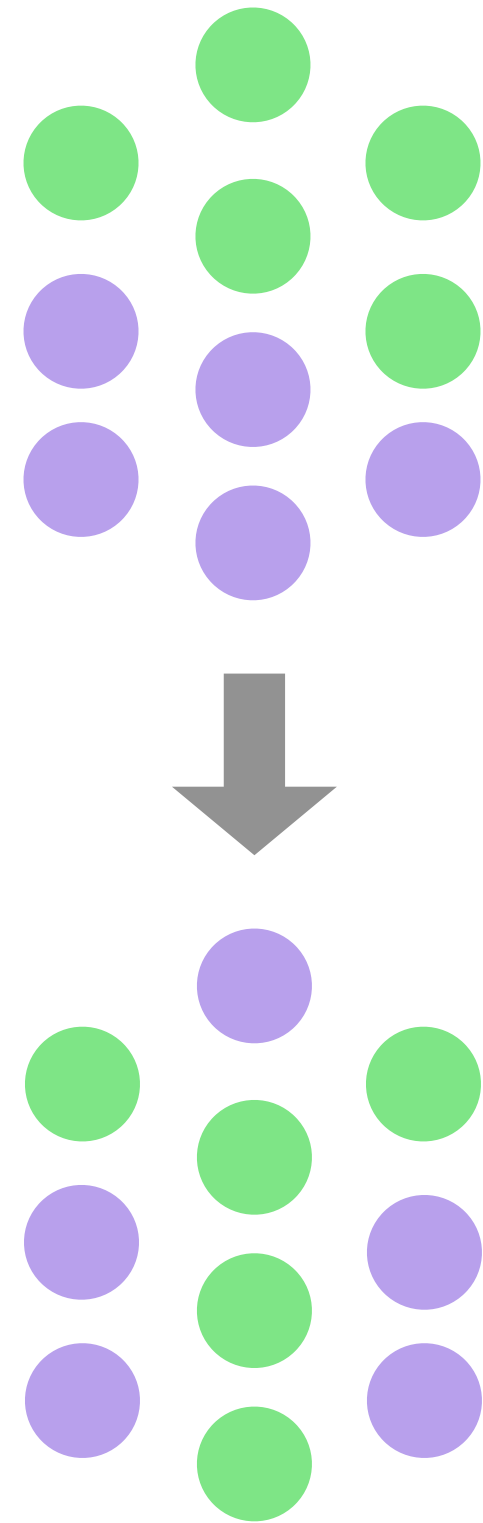
| Grammatical (CoLA) | |
|---|---|
| will | 0.267 |
| John | 0.168 |
| . | 0.006 |
| and | -0.039 |
| in | -0.050 |

| Ungrammatical (CoLA) | |
|---|---|
| book | 2.737 |
| is | 2.659 |
| was | 2.312 |
| of | 2.308 |
| in | 1.972 |

[ Gururangan et al., 2018 ]
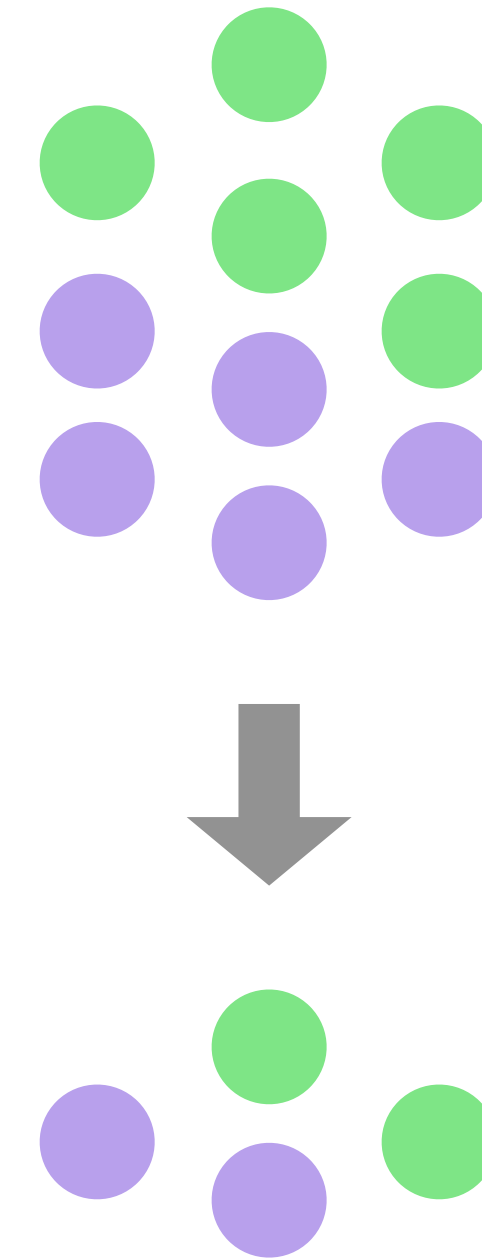
13

# Future Work

Making Tougher Datasets

Other Modalities

apple

time

Data Pruning

# Summary: A unified framework for interpreting datasets.



$I_{\mathscr{V}}(X \to Y)$

$\text{PVI}(x \to y)$

$I_{\mathscr{V}}(X_i \to Y)$

$I_{\mathscr{V}}(X \to Y)$

$\overline{\text{PVI}}(x \to y) \mid (x, y) \in S$