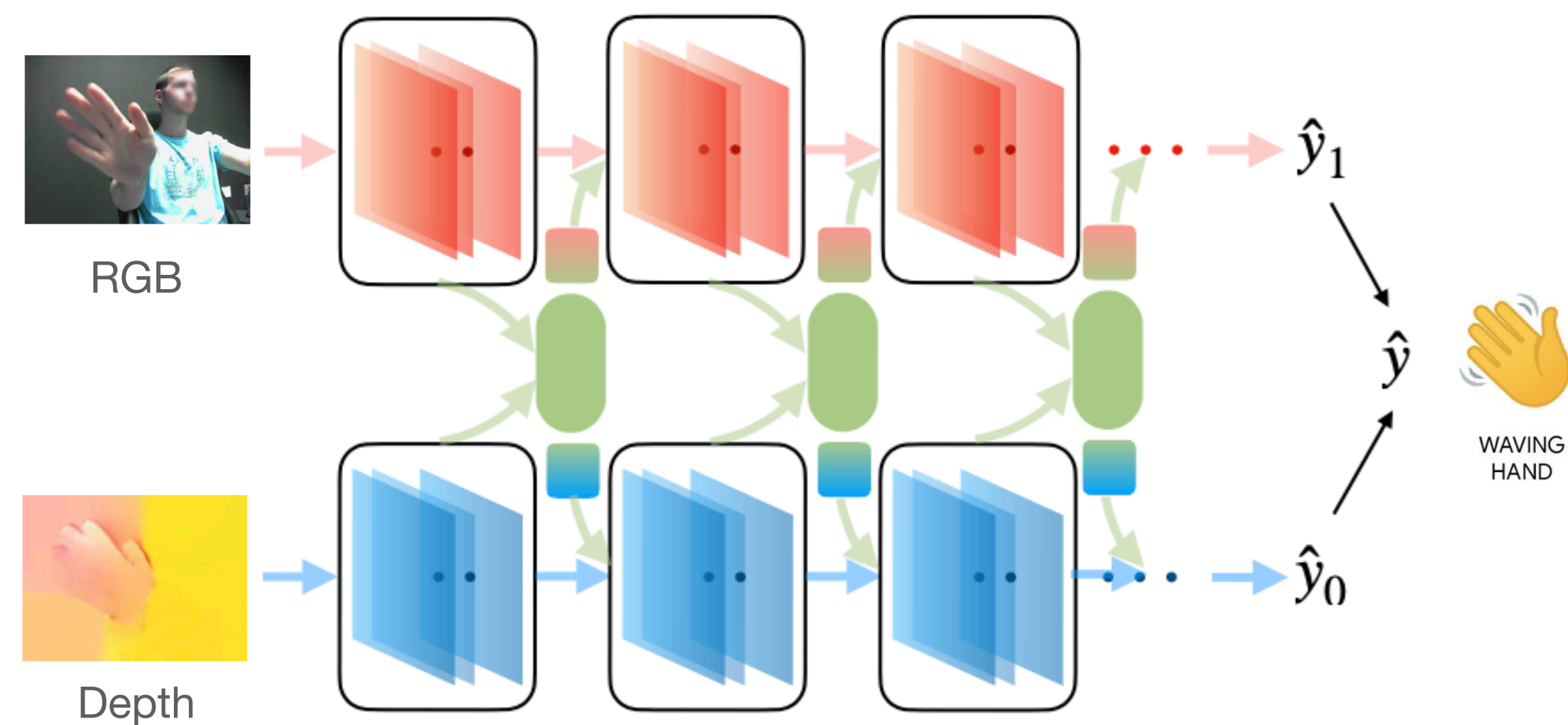
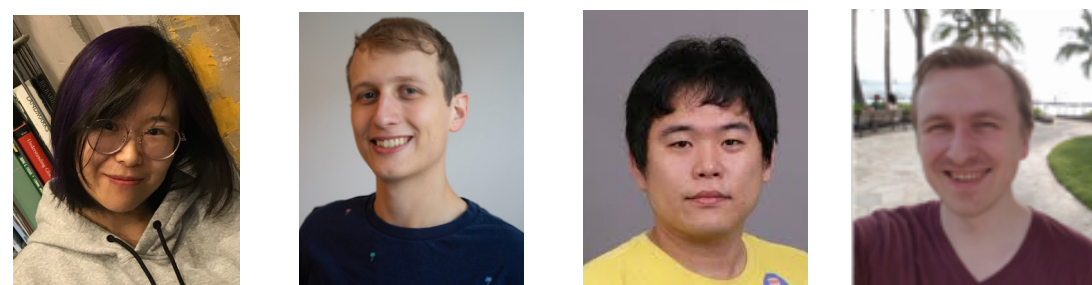


# Characterizing and Overcoming the Greedy Nature of Learning in Multi-Modal Deep Neural Networks

## *Dynamic Hand Gestures Recognition*

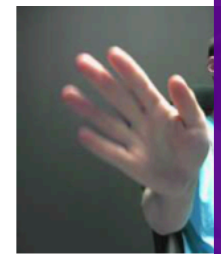


**Nan Wu**, Stanisław Jastrzębski, Kyunghyun Cho and Krzysztof J. Geras.

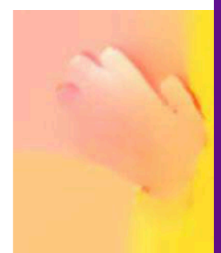


# Characterizing and Overcoming the Greedy Nature of Learning in Multi-Modal Deep Neural Networks

Dynamic Hand Gesture Recognition



RGB



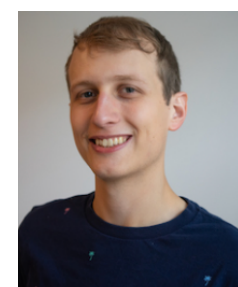
Depth

**Do multi-modal DNNs attend to all modalities?**

**Can we make multi-modal DNNs utilize all modalities?**

**Does better utilization of all modalities imply better generalization?**

**Nan Wu**, Stanisław Jastrzębski, Kyunghyun Cho and Krzysztof J. Geras.



# Metric 1: Conditional utilization rate

$$\hat{y}_0 = f_0(\mathbf{x}_{m_0}, \mathbf{x}_{m_1})$$

$$\hat{y}_1 = f_1(\mathbf{x}_{m_0}, \mathbf{x}_{m_1})$$

# Metric 1: Conditional utilization rate

$$\hat{y}_0 = f_0(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_0 = f'_0(\mathbf{x}_{m_0})$$

$$\hat{y}_1 = f_1(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_1 = f'_1(\mathbf{x}_{m_1}).$$

# Metric 1: Conditional utilization rate

$$\hat{y}_0 = f_0(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_0 = f'_0(\mathbf{x}_{m_0})$$

$$\hat{y}_1 = f_1(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_1 = f'_1(\mathbf{x}_{m_1}).$$

$$u(m_1|m_0) = \frac{A(f_0) - A(f'_0)}{A(f_0)}$$

$$u(m_0|m_1) = \frac{A(f_1) - A(f'_1)}{A(f_1)},$$

**The relative change in accuracy between the two models:**

- one using all modalities,
- the other using only one.

# Metric 1: Conditional utilization rate

$$\hat{y}_0 = f_0(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_0 = f'_0(\mathbf{x}_{m_0})$$

$$\hat{y}_1 = f_1(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_1 = f'_1(\mathbf{x}_{m_1}).$$

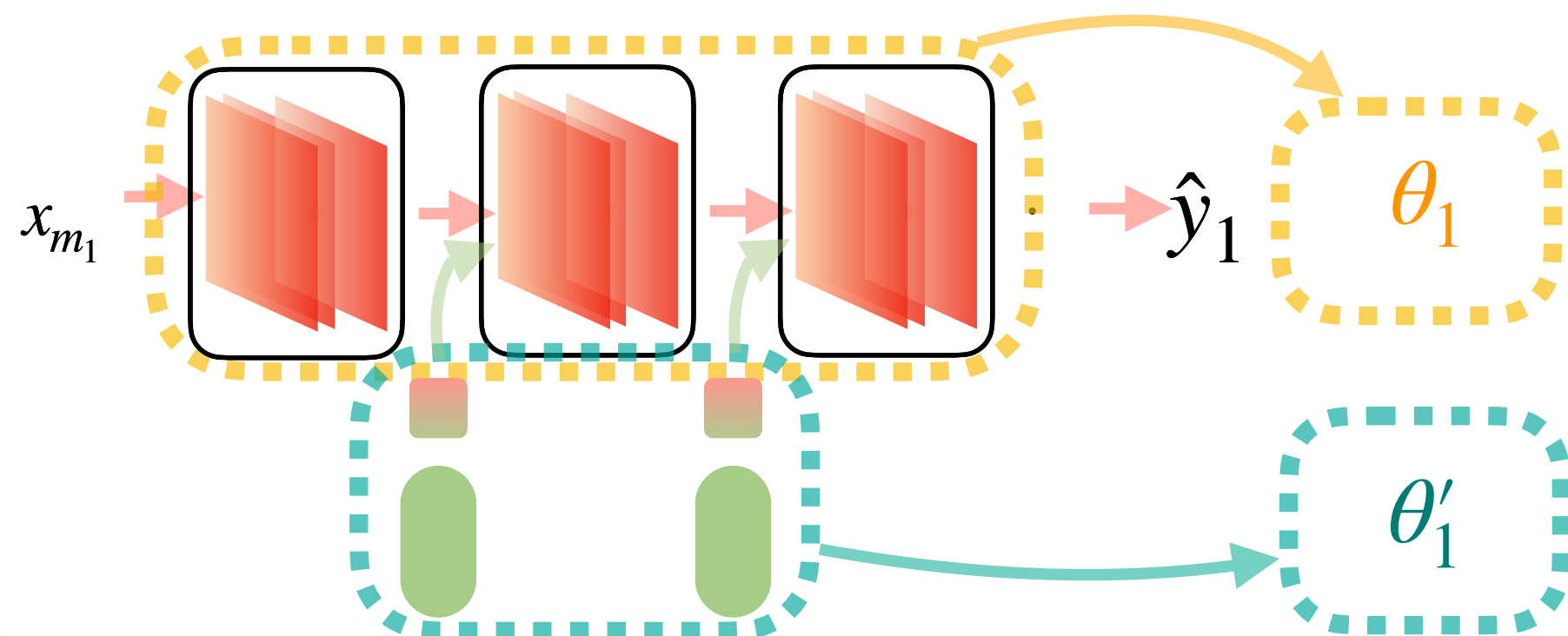
$$u(m_1|m_0) = \frac{A(f_0) - A(f'_0)}{A(f_0)}$$

$$u(m_0|m_1) = \frac{A(f_1) - A(f'_1)}{A(f_1)},$$

The relative change in accuracy between the two models:

- one using all modalities,
- the other using only one.

# Metric 2: Conditional learning speed



# Metric 1: Conditional utilization rate

$$\hat{y}_0 = f_0(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_0 = f'_0(\mathbf{x}_{m_0})$$

$$\hat{y}_1 = f_1(\mathbf{x}_{m_0}, \mathbf{x}_{m_1}) \quad \hat{y}'_1 = f'_1(\mathbf{x}_{m_1}).$$

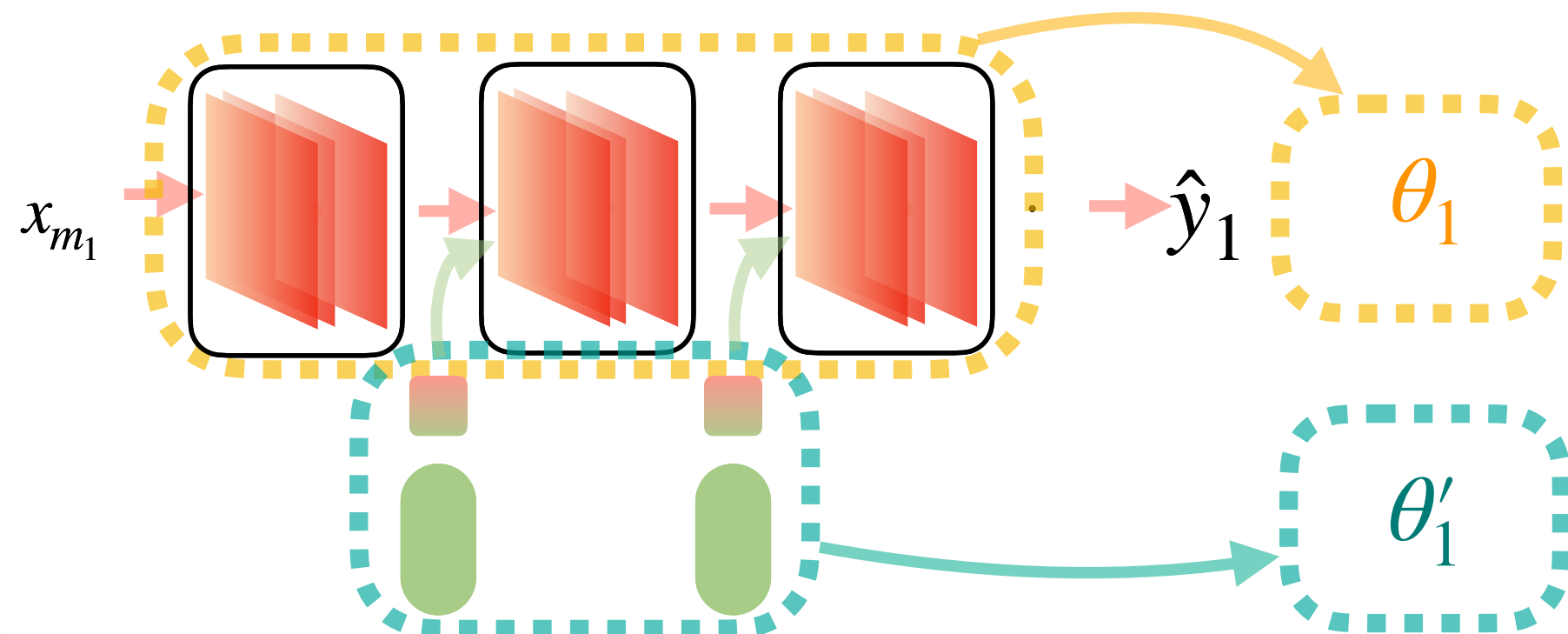
$$u(m_1|m_0) = \frac{A(f_0) - A(f'_0)}{A(f_0)}$$

$$u(m_0|m_1) = \frac{A(f_1) - A(f'_1)}{A(f_1)},$$

The relative change in accuracy between the two models:

- one using all modalities,
- the other using only one.

# Metric 2: Conditional learning speed



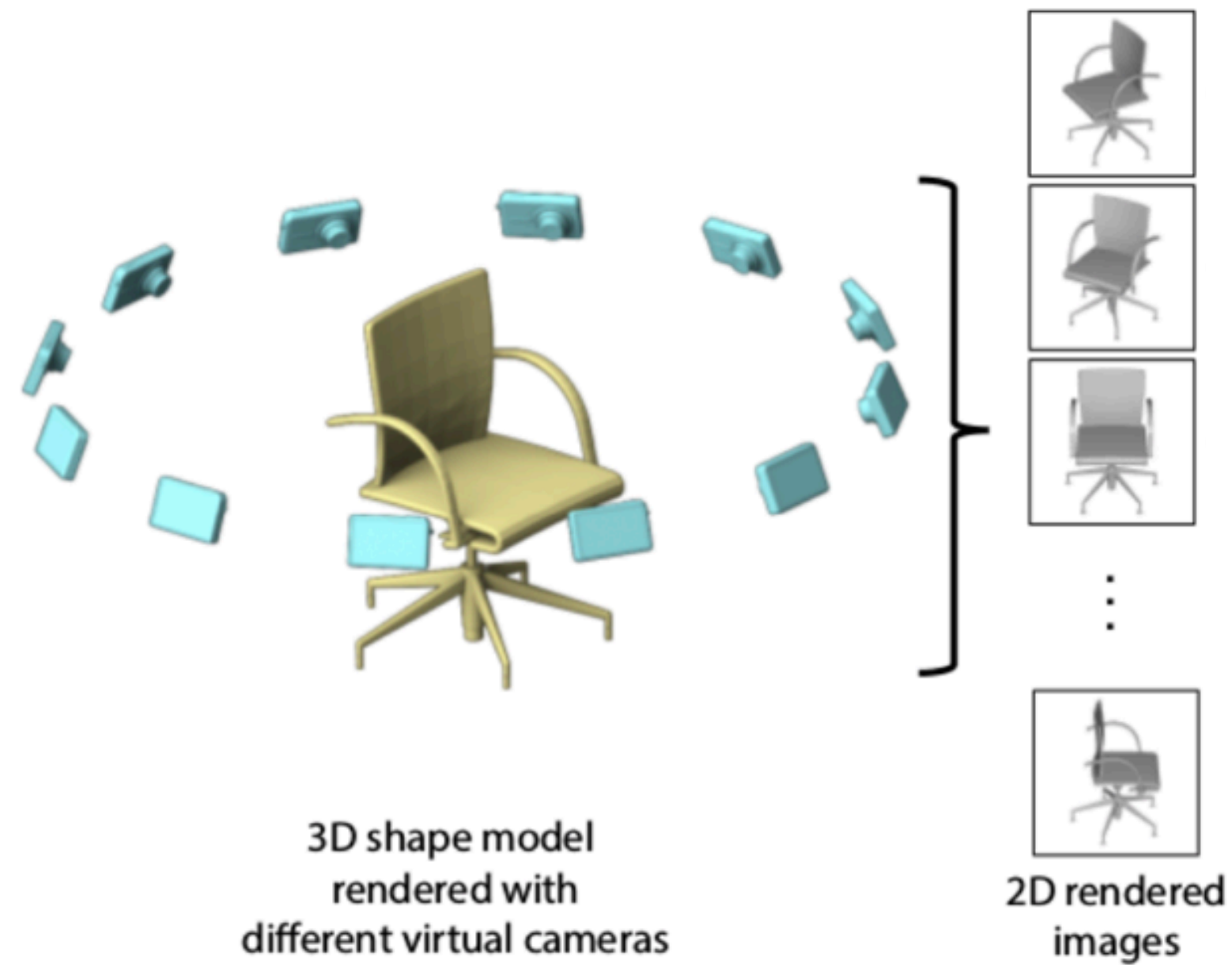
$$s(m_0|m_1; t) = \log \frac{\sum_{i=1}^t \mu(\theta'_1; i)}{\sum_{i=1}^t \mu(\theta_1; i)},$$

$$\mu(\theta; i) = \|\mathbf{G}\|_2^2 / \|\theta_{(i)}\|_2^2$$

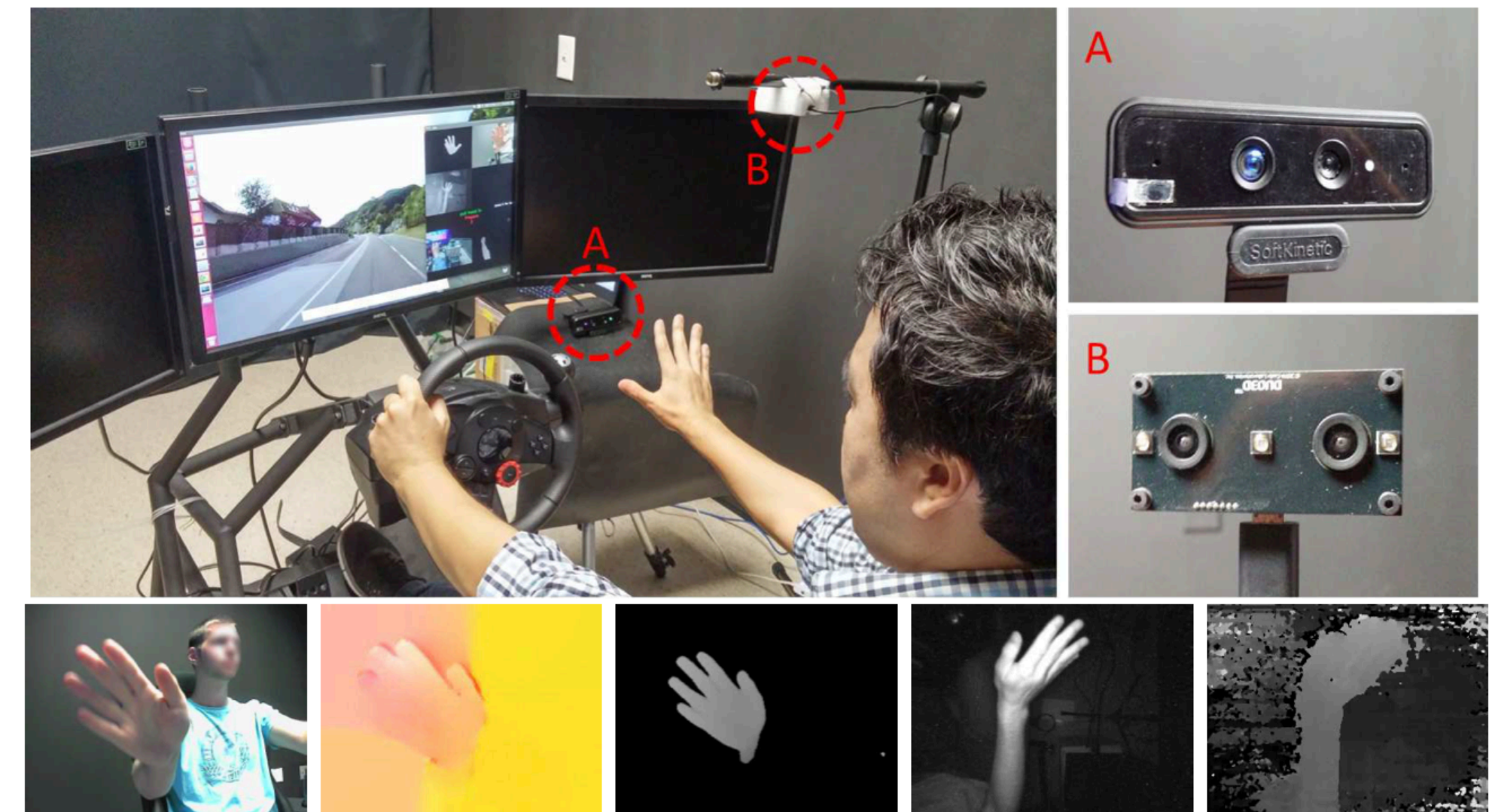
the log-ratio between the learning speed of

- the uni-modal branch and
- the corresponding fusion components.

# Observations: Imbalanced learning between modalities



(a) ModelNet40

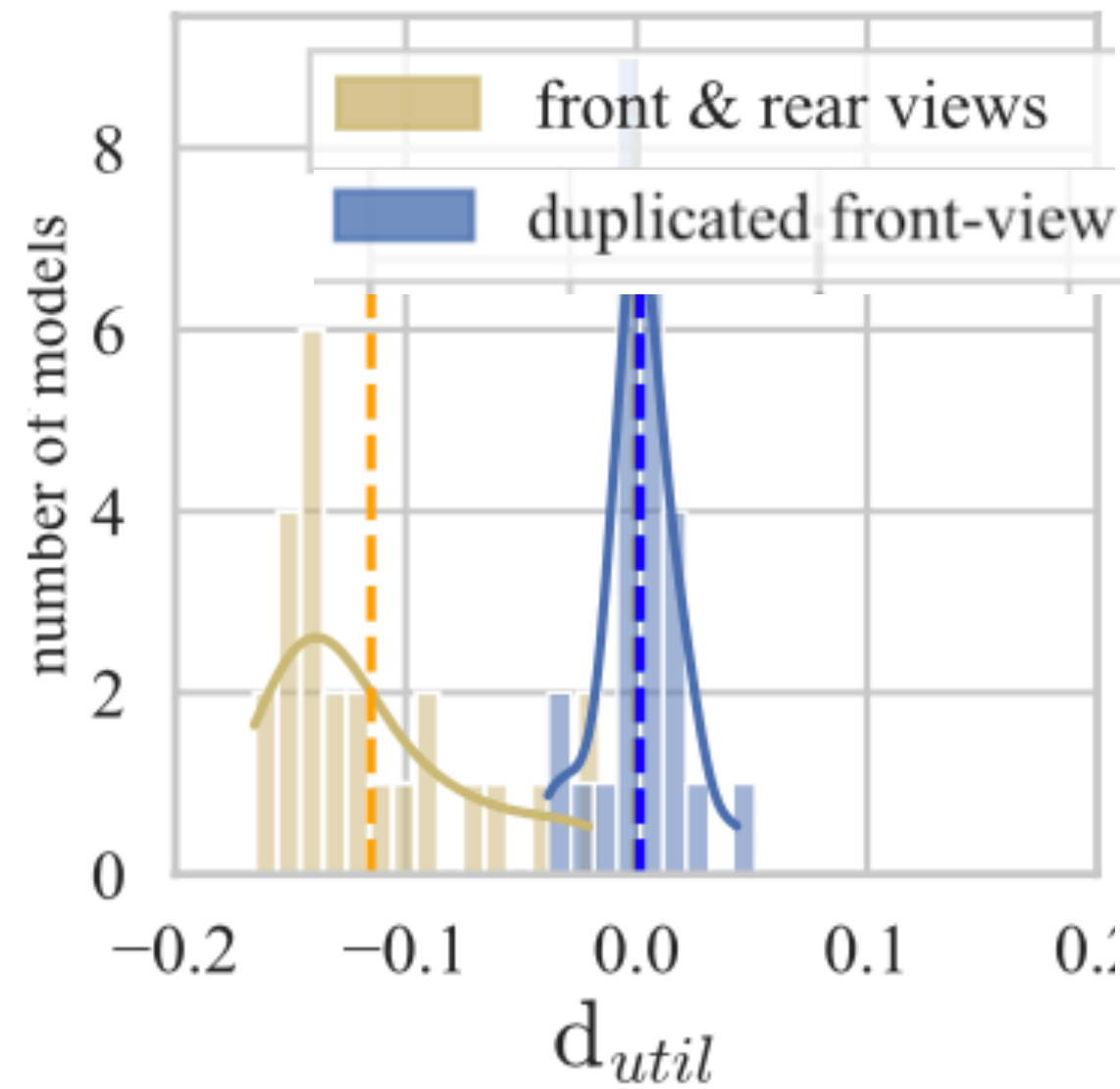


(b) NVGesture

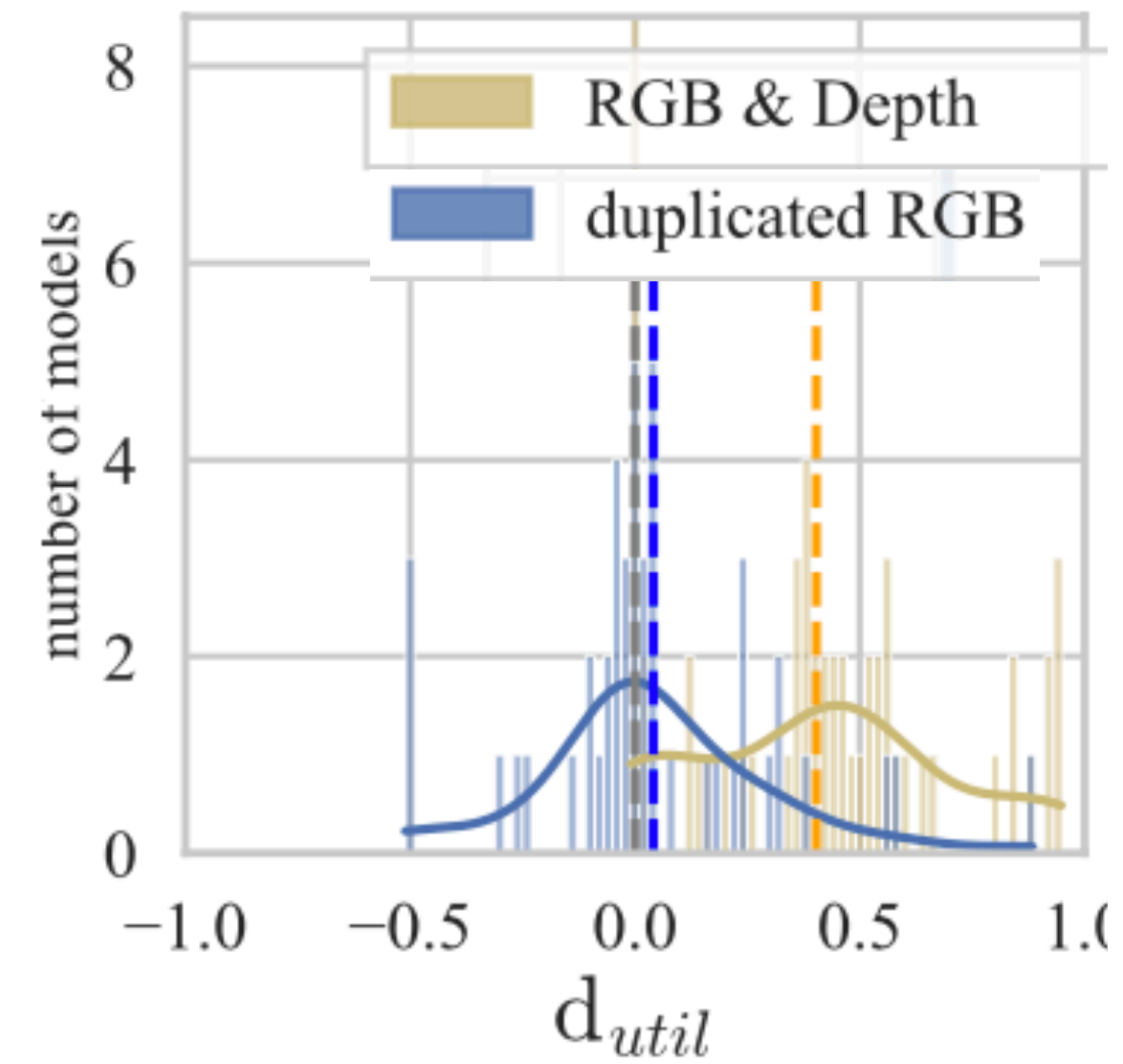
$$d_{\text{util}}(f) = \mathbf{u}(m_1|m_0) - \mathbf{u}(m_0|m_1)$$
$$d_{\text{speed}}(f; t) = \mathbf{s}(m_1|m_0; t) - \mathbf{s}(m_0|m_1; t)$$



# Observations: Imbalanced learning between modalities



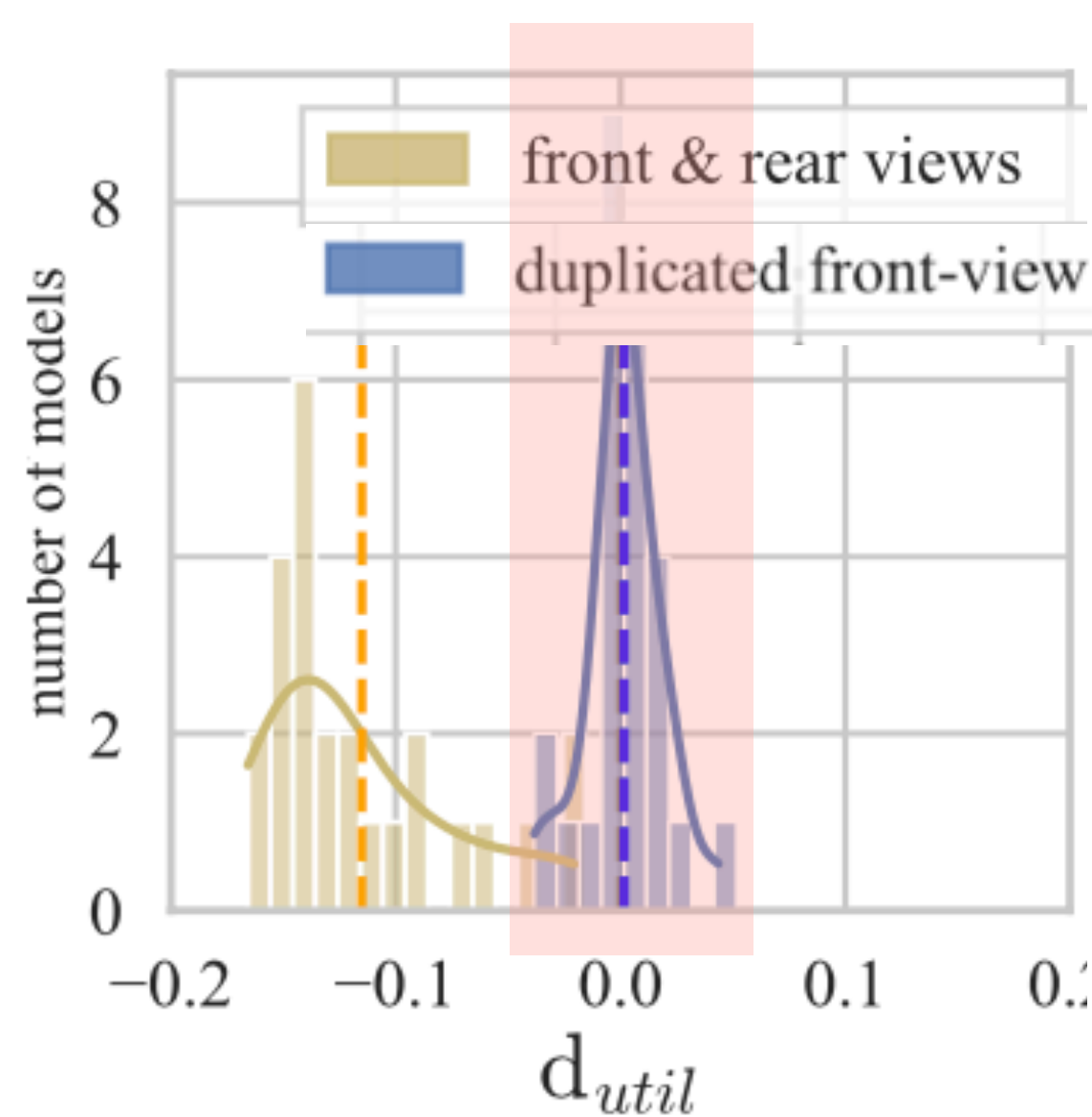
(a) ModelNet40



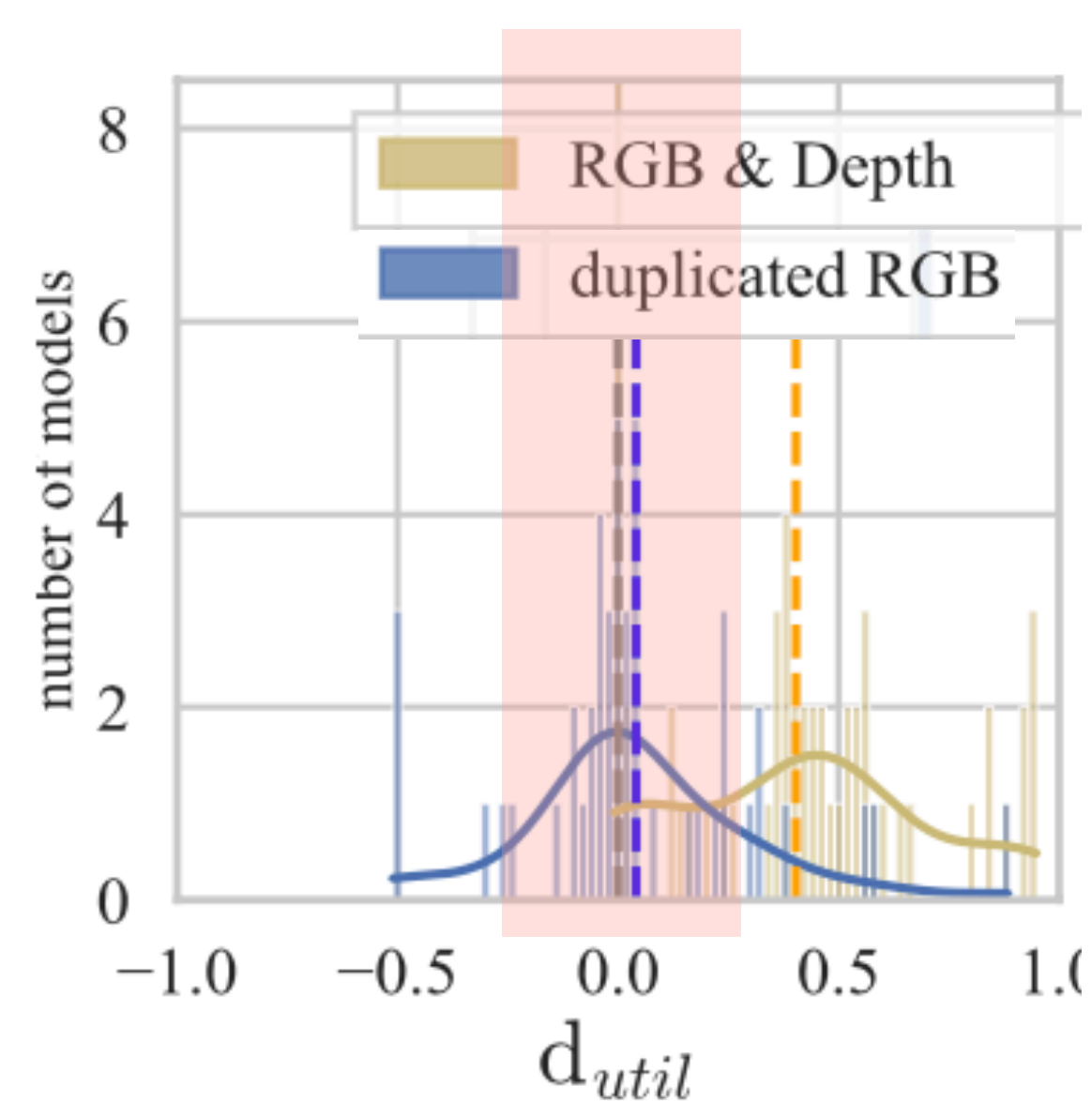
(b) NVGesture

$$d_{util}(f) = \mathbf{u}(m_1|m_0) - \mathbf{u}(m_0|m_1)$$
$$d_{speed}(f; t) = \mathbf{s}(m_1|m_0; t) - \mathbf{s}(m_0|m_1; t)$$

# Observations: Imbalanced learning between modalities



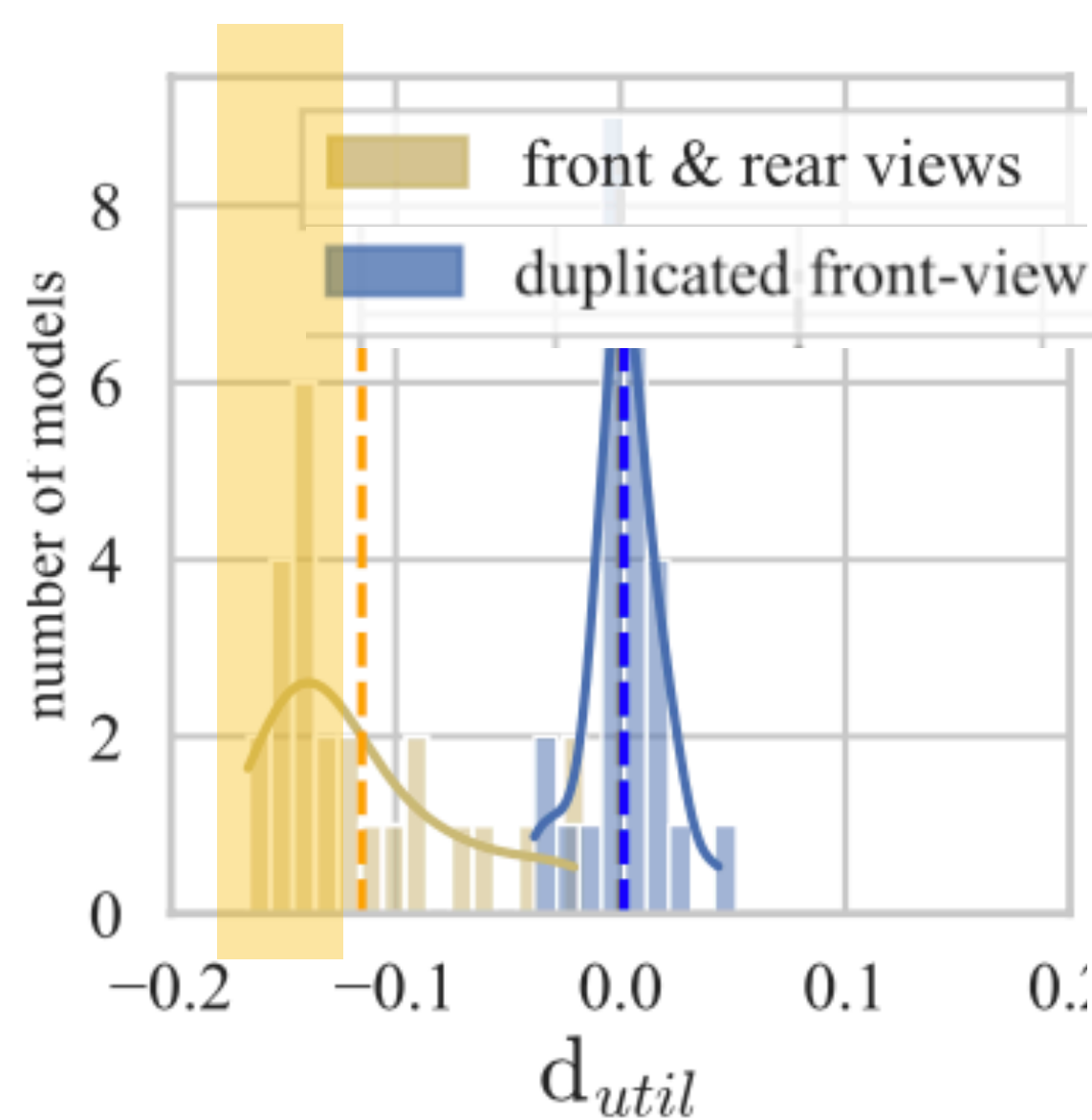
(a) ModelNet40



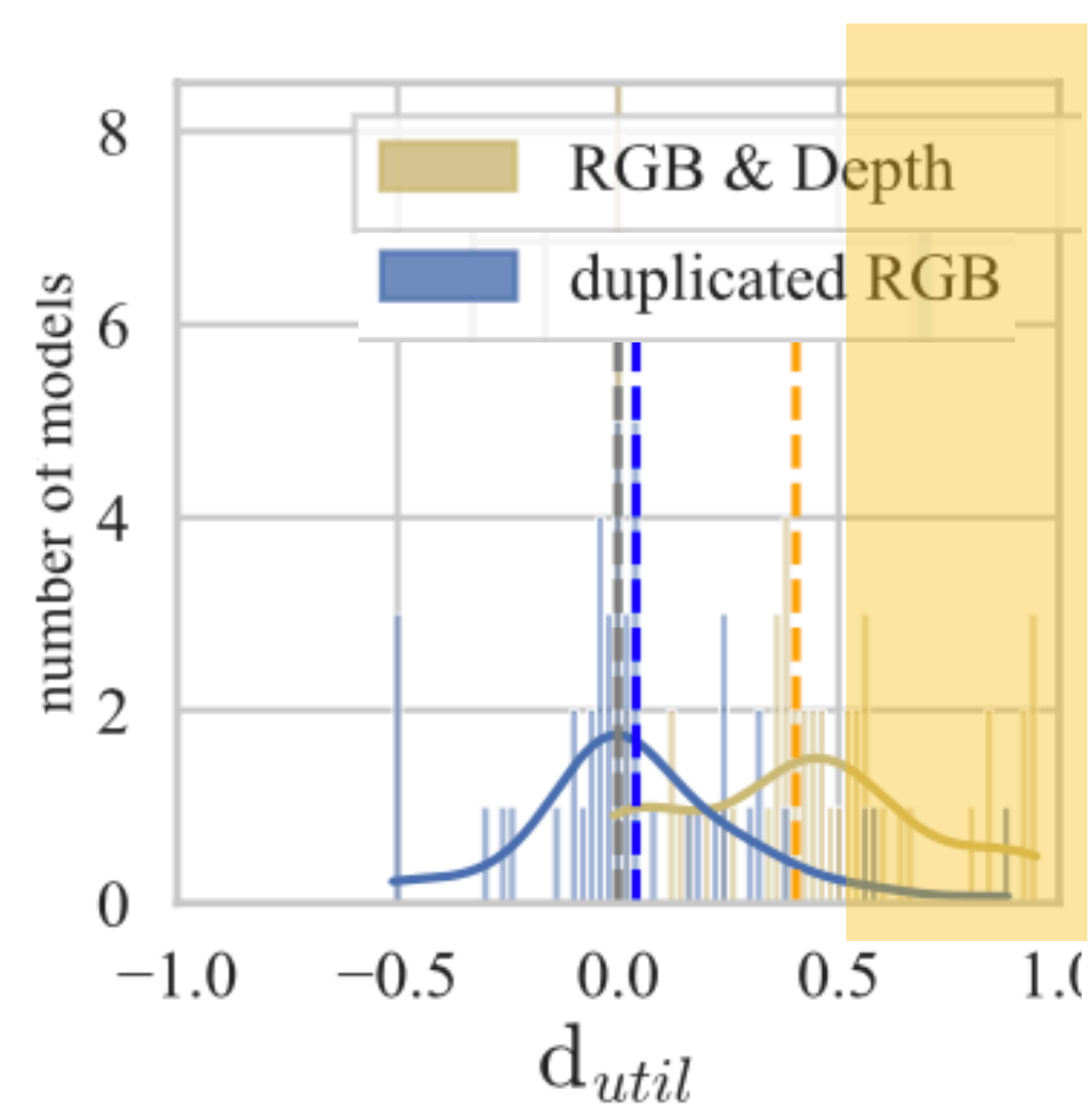
(b) NVGesture

$$d_{util}(f) = \mathbf{u}(m_1|m_0) - \mathbf{u}(m_0|m_1)$$
$$d_{speed}(f; t) = \mathbf{s}(m_1|m_0; t) - \mathbf{s}(m_0|m_1; t)$$

# Observations: Imbalanced learning between modalities



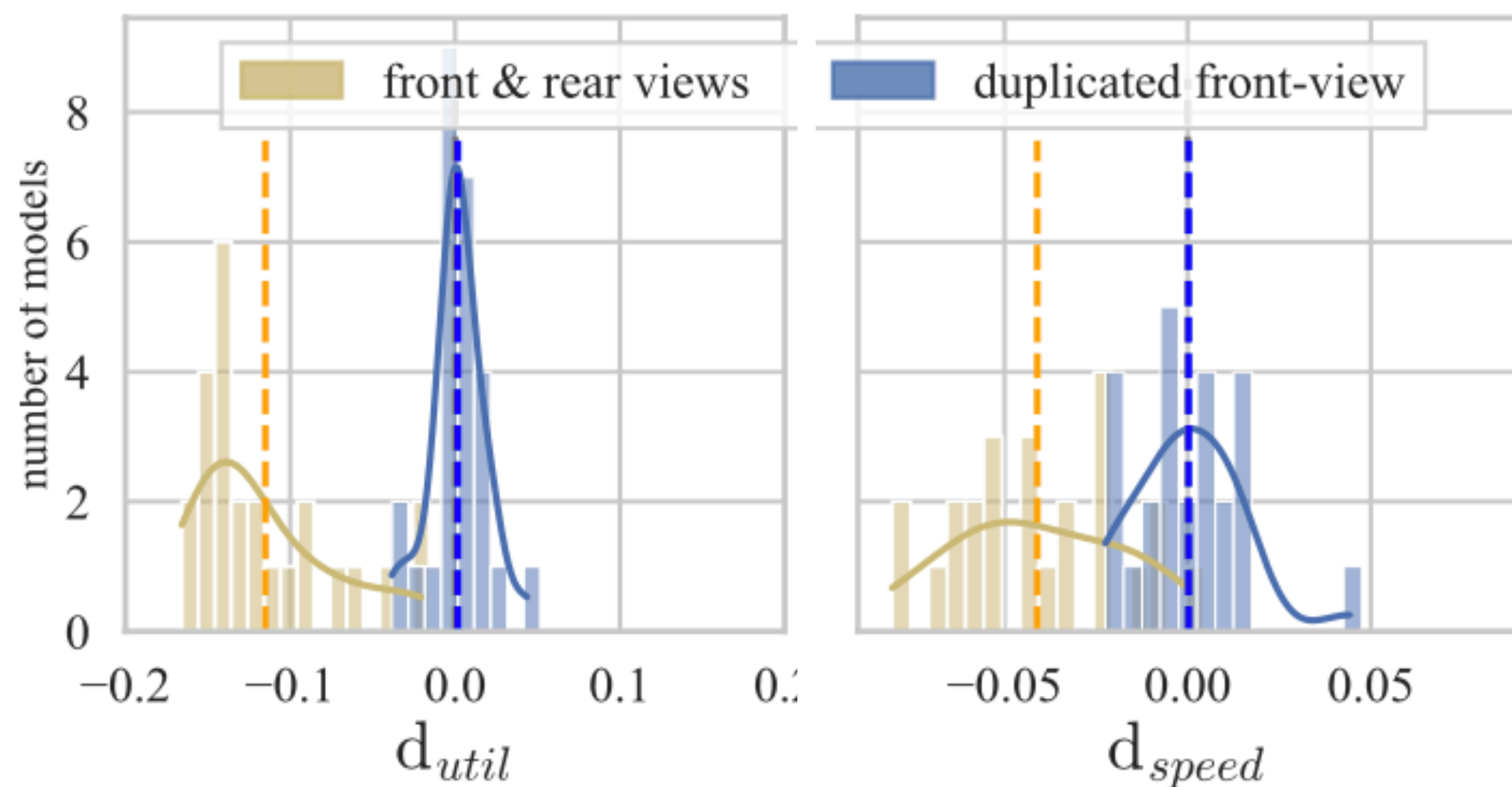
(a) ModelNet40



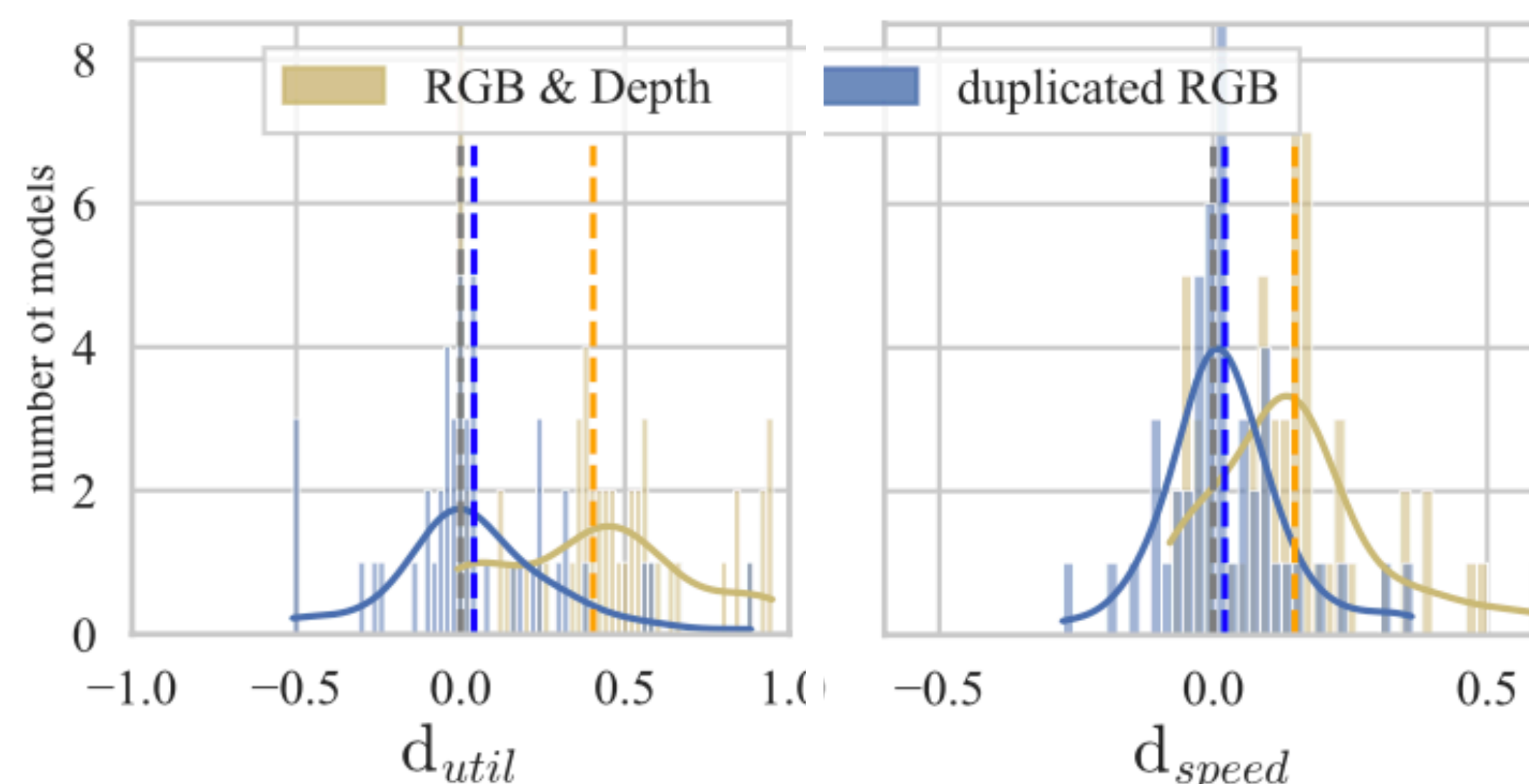
(b) NVGesture

$$d_{util}(f) = \mathbf{u}(m_1|m_0) - \mathbf{u}(m_0|m_1)$$
$$d_{speed}(f; t) = \mathbf{s}(m_1|m_0; t) - \mathbf{s}(m_0|m_1; t)$$

# Observations: Imbalanced learning between modalities



(a) ModelNet40



(b) NVGesture

$$d_{util}(f) = \mathbf{u}(m_1|m_0) - \mathbf{u}(m_0|m_1)$$
$$d_{speed}(f; t) = \mathbf{s}(m_1|m_0; t) - \mathbf{s}(m_0|m_1; t)$$

# Greedy learner hypothesis

A multi-modal learning process is greedy when it produces models that **rely on only one of the available modalities.**

# Greedy learner hypothesis

A multi-modal learning process is greedy when it produces models that **rely on only one of the available modalities.**

The modality that the multi-modal DNN primarily relies on is the modality that is the **fastest** to learn from.

# Greedy learner hypothesis

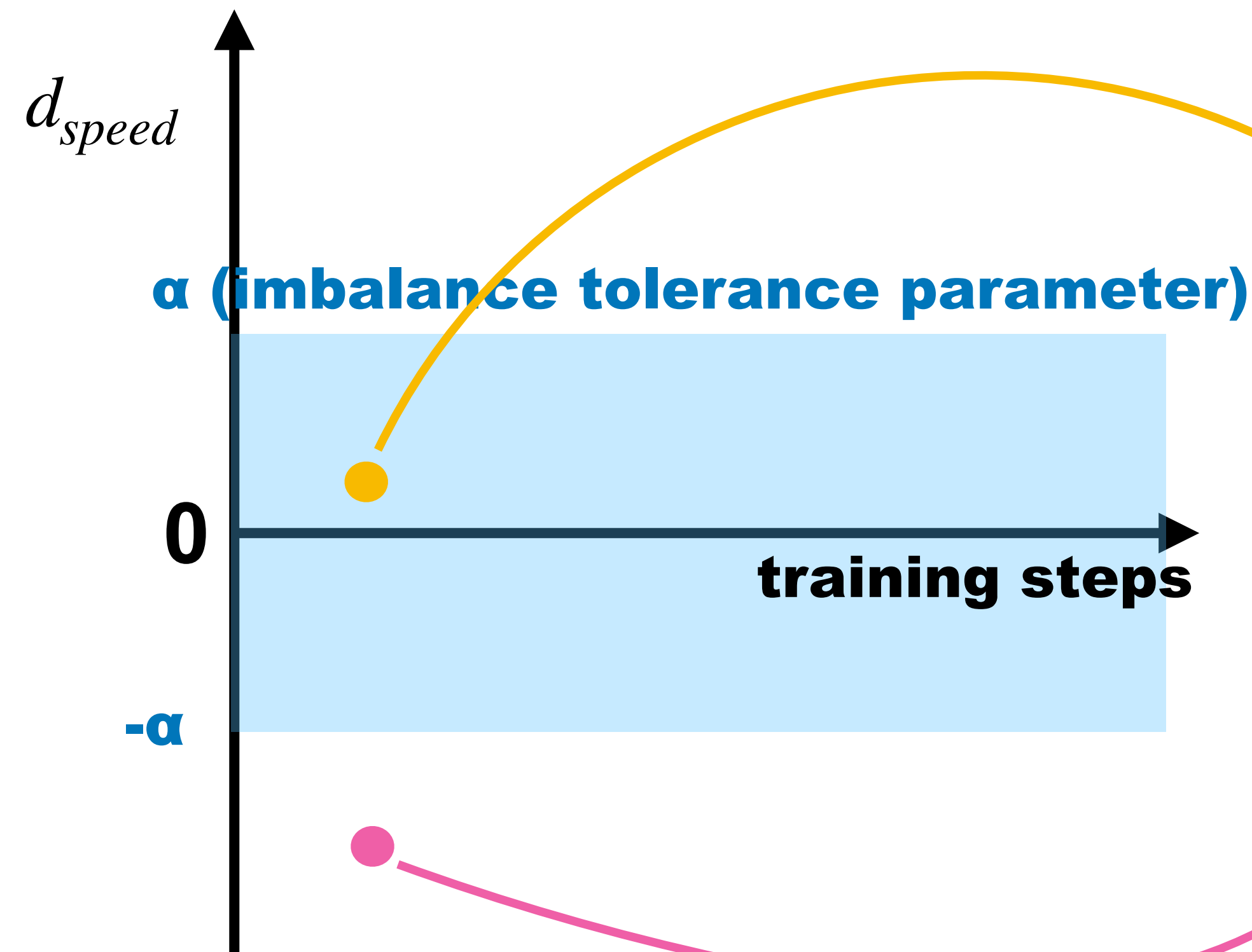
A multi-modal learning process is greedy when it produces models that **rely on only one of the available modalities.**

The modality that the multi-modal DNN primarily relies on is the modality that is the **fastest** to learn from.

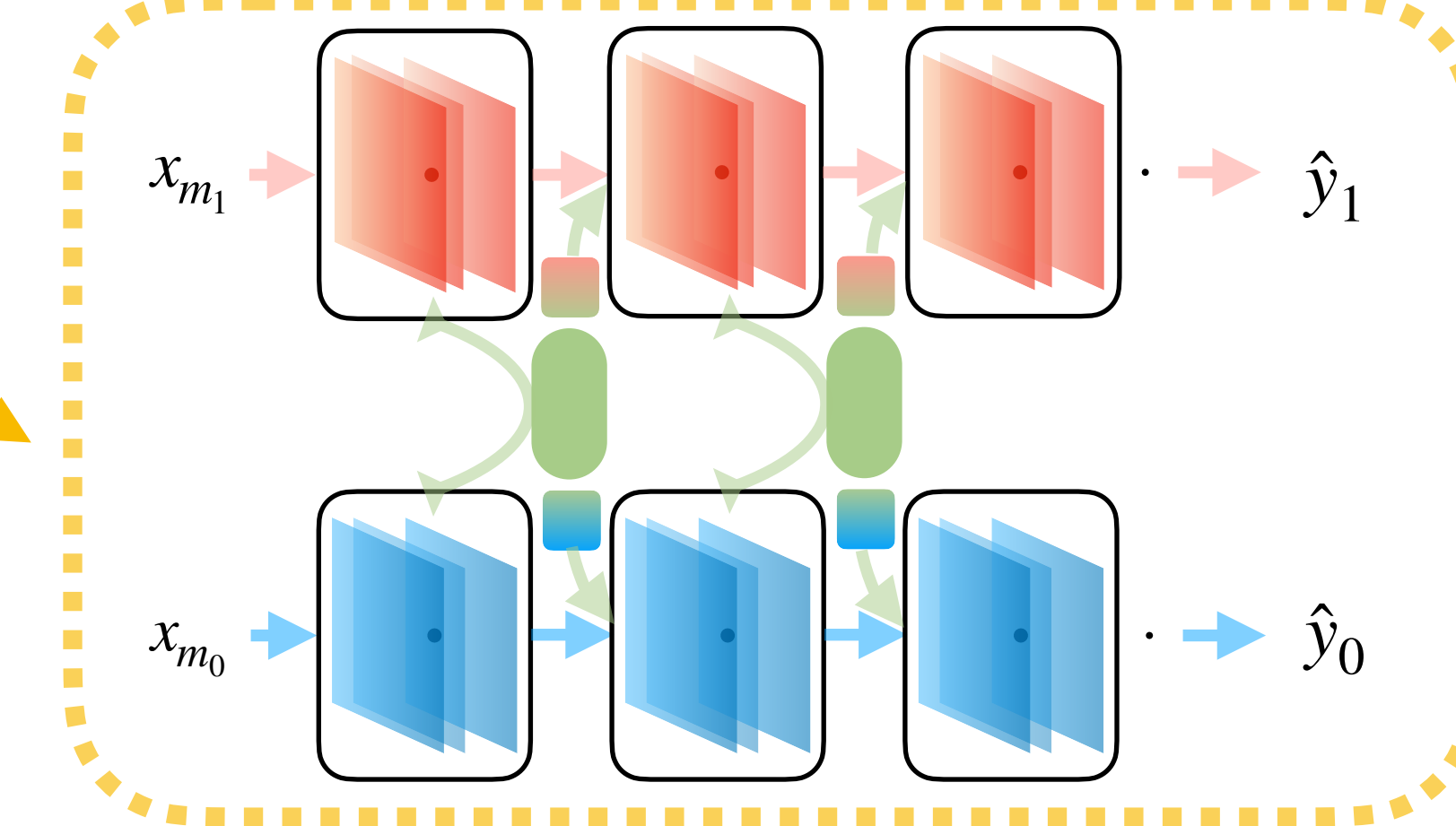
We **hypothesize that a multi-modal learning process, in which a multi-modal DNN is trained to minimize the sum of the modality-specific losses, is greedy.**

# Balanced multimodal learning

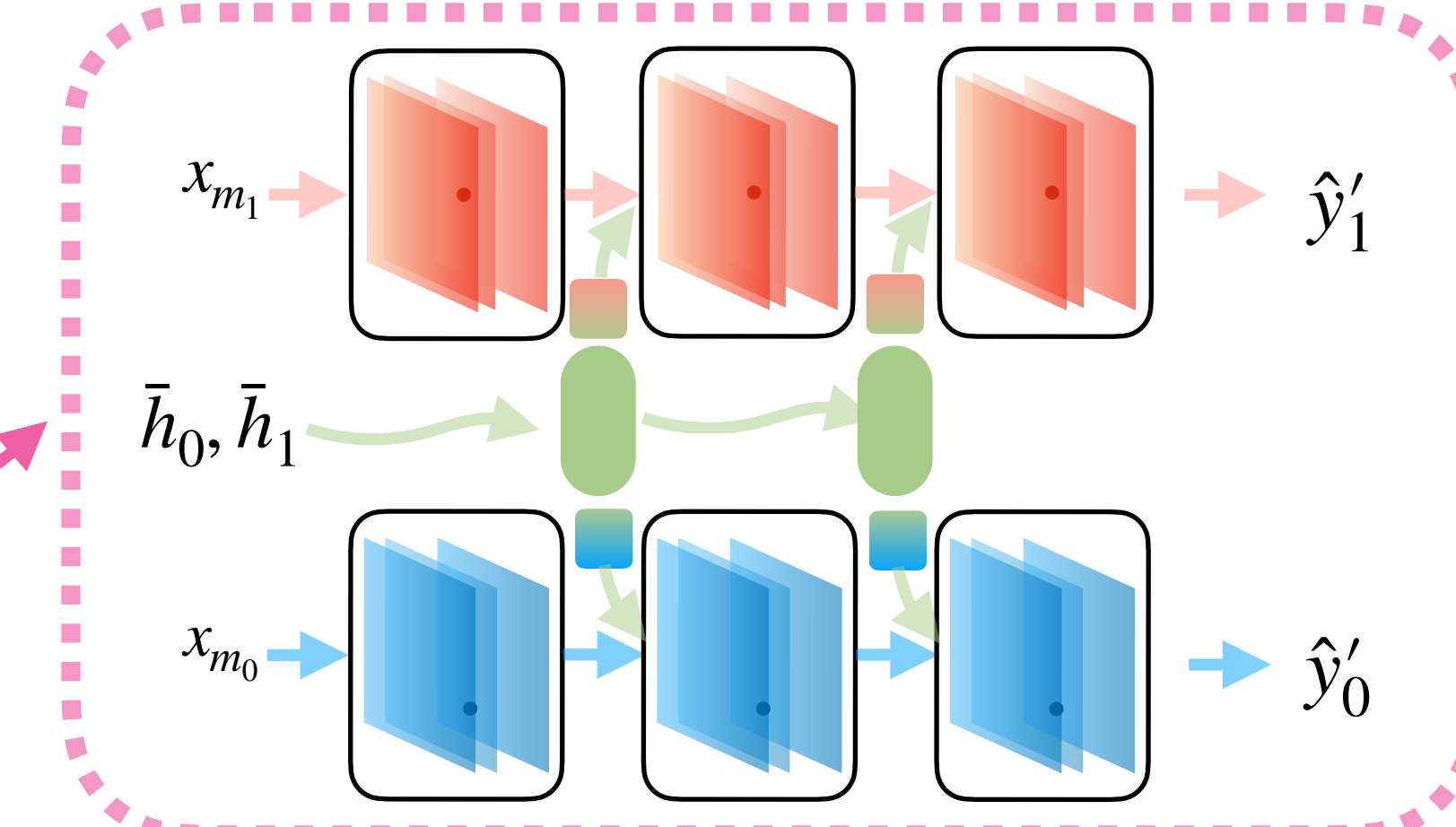
[Guided]



Regular steps



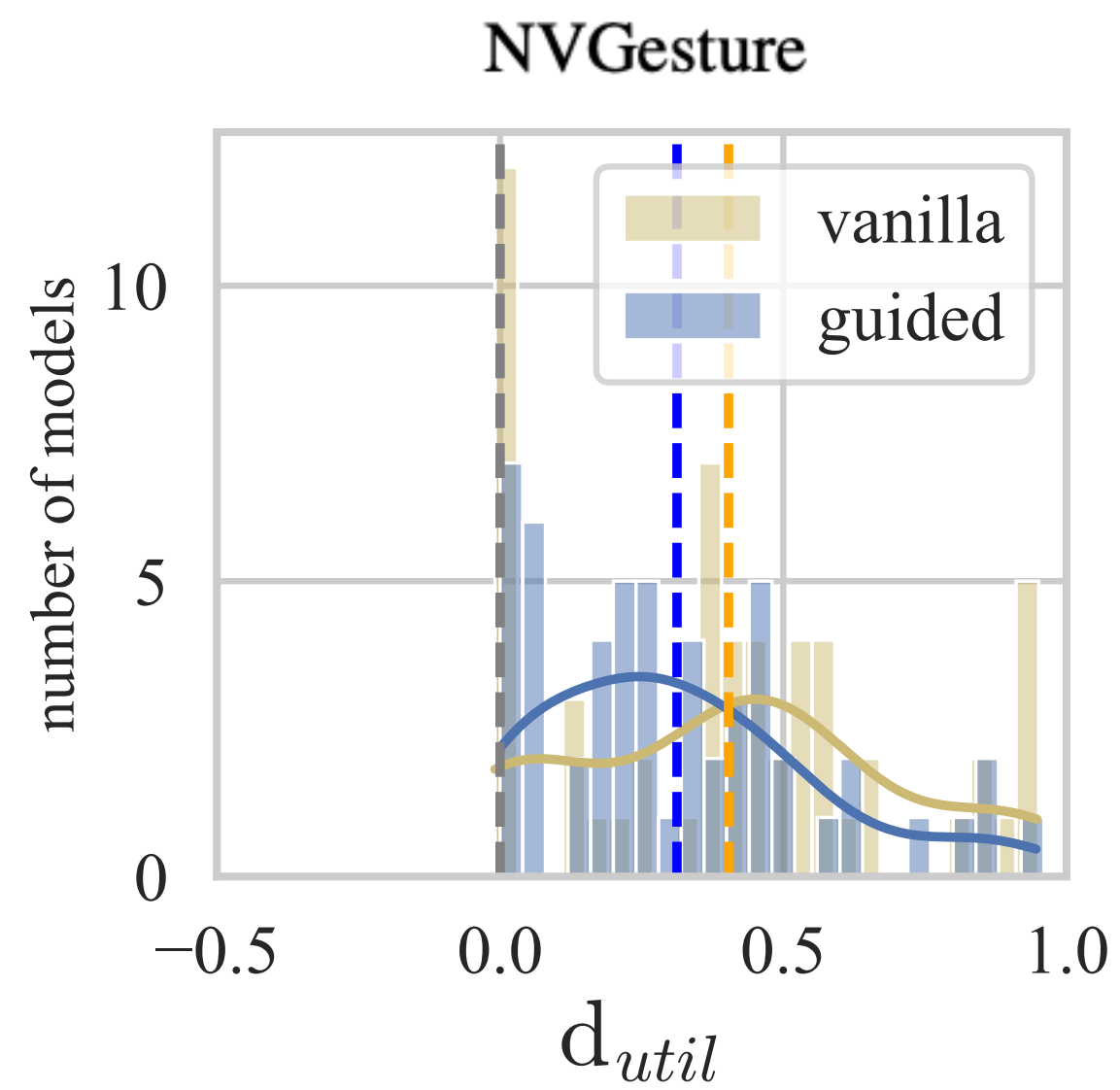
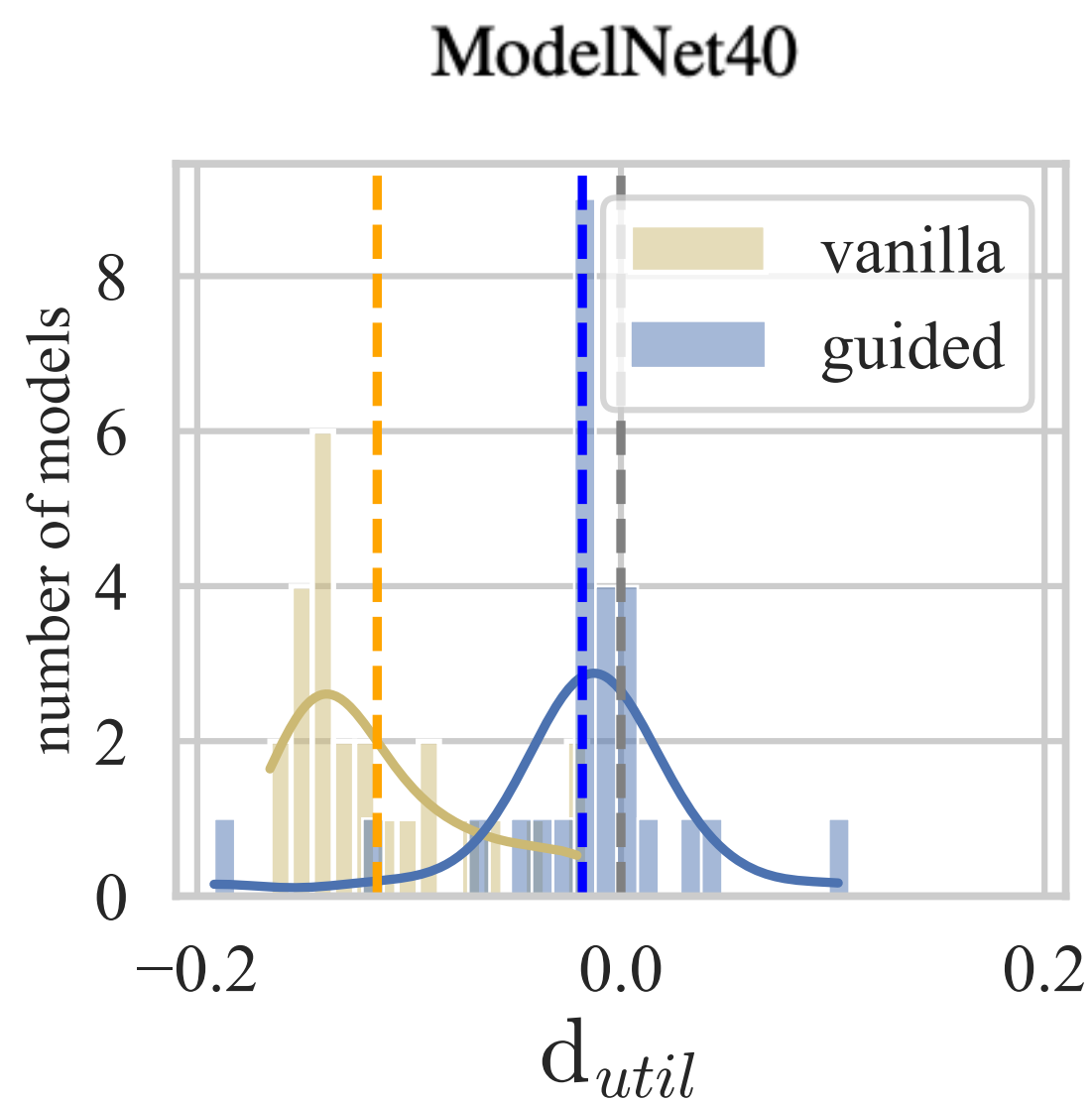
Re-balancing steps





# Results


- **Calibrating modality utilization**



- **Improving generalization**

	ModelNet40	NVGesture-scratch	NVGesture-pretrained
uni-modal (best)	89.34±0.39	77.59±0.55	78.98±2.02
multi-modal (vanilla)	90.09±0.58	79.81±1.14	83.20±0.21
+ RUBi (Cadene et al., 2019)	90.45±0.58	79.95±0.12	81.60±1.28
+ random (proposed)	91.36±0.10	79.88±0.90	82.64±0.84
+ guided (proposed)	<b>91.37±0.28</b>	<b>80.22±0.73</b>	<b>83.82±1.45</b>

# Thank you!

**Nan Wu** [email: [nan.wu@nyu.edu](mailto:nan.wu@nyu.edu); twitter:  **Nan Wu** @NanWu\_ ],  
**Stanisław Jastrzębski, Kyunghyun Cho and Krzysztof J. Geras.**  
**Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks. ICML, 2022.**

[https://github.com/nyukat/greedy\\_multimodal\\_learning](https://github.com/nyukat/greedy_multimodal_learning)

