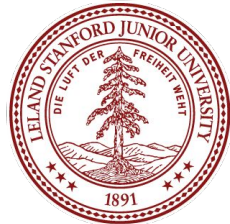


# Evaluating Self-Supervised Learning via Risk Decomposition

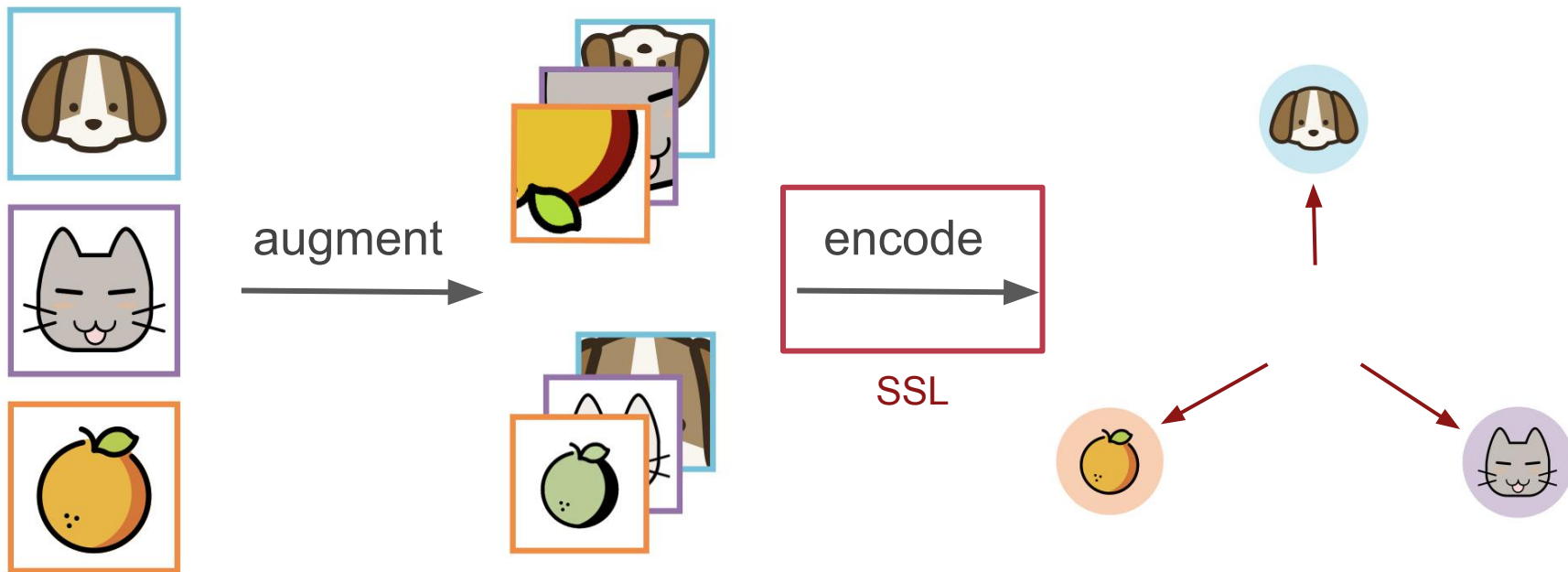
Yann Dubois, Tatsunori Hashimoto, Percy Liang

ICML 2023 Oral



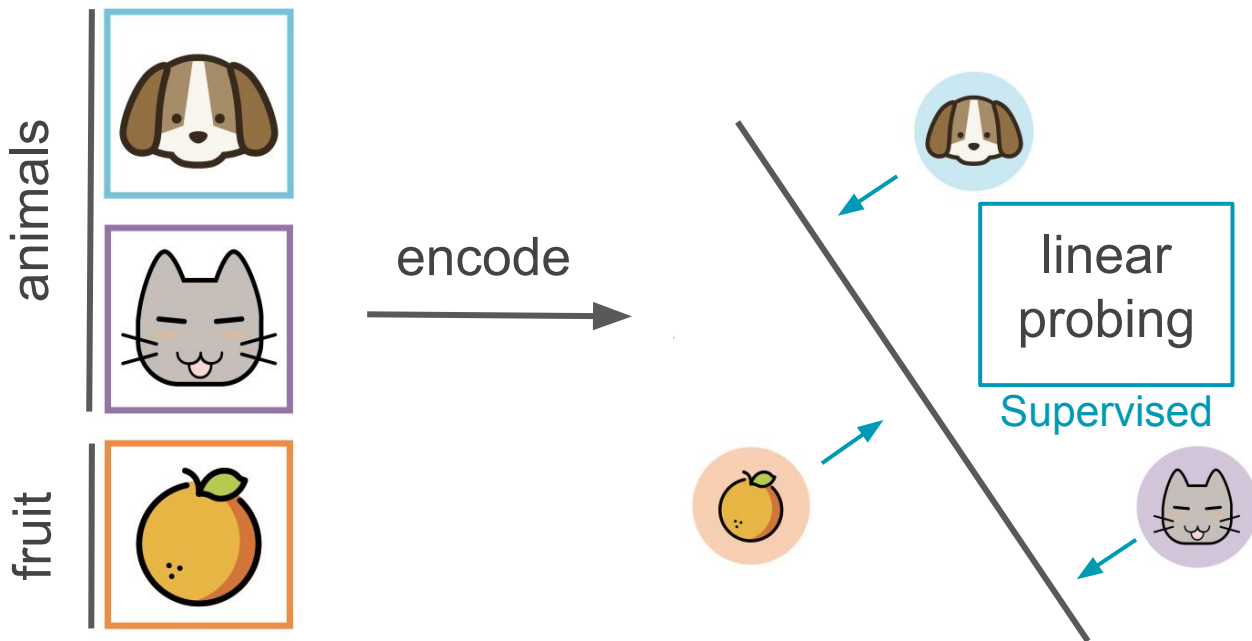
# Background: pretraining in SSL

- E.g. train on **unlabeled** ImageNet using **SSL**



# Background: linear probing in SSL

- E.g. train on **labeled** ImageNet using **supervised** learning



# Motivation: evaluating SSL

- There are many SSL pipelines that get proposed, with different:
  - SSL objectives
  - architectures
  - optimizers
  - pretraining data
- Evaluated on a single metric: linear probing on ImageNet.

# Motivation: evaluating SSL

- There are many SSL pipelines that get proposed, with different:
  - SSL objectives
  - architectures
  - optimizers
  - pretraining data
- Evaluated on a single metric: linear probing on ImageNet.
  - ✗ why is an SSL pipeline better?
  - ✗ when is an SSL pipeline better?
  - ✗ how to improve the SSL pipeline?

# Motivation: supervised risk decomposition

- **Supervised learning** monitor training/validation loss
  - underfitting  $\Rightarrow$  increase capacity
  - overfitting  $\Rightarrow$  regularize
  - small training loss: model would be better with large datasets
  - ...

# Motivation: supervised risk decomposition

- **Supervised learning** monitor training/validation loss
  - underfitting  $\Rightarrow$  increase capacity
  - overfitting  $\Rightarrow$  regularize
  - small training loss: model would be better with large datasets
  - ...

Risk = approximation error + generalization error

~training error

~training - validation error

# Motivation: supervised risk decomposition

- **Supervised learning** monitor training/validation loss
  - underfitting  $\Rightarrow$  increase capacity
  - overfitting  $\Rightarrow$  regularize
  - small training loss: model would be better with large datasets
  - ...
- **Self-supervised learning**
  - ?



# Motivation: supervised risk decomposition

- **Supervised learning** monitor training/validation loss
  - underfitting  $\Rightarrow$  increase capacity
  - overfitting  $\Rightarrow$  regularize
  - small training loss: model would be better with large datasets
  - ...
- **Self-supervised learning**
  - ?

Idea: generalize risk decomposition to SSL and estimate it

# Motivation: supervised risk decomposition

predictor's limitation

no constraints

linear

$\mathcal{F}$

$\infty$  data

$S$

finite data

0

approx.

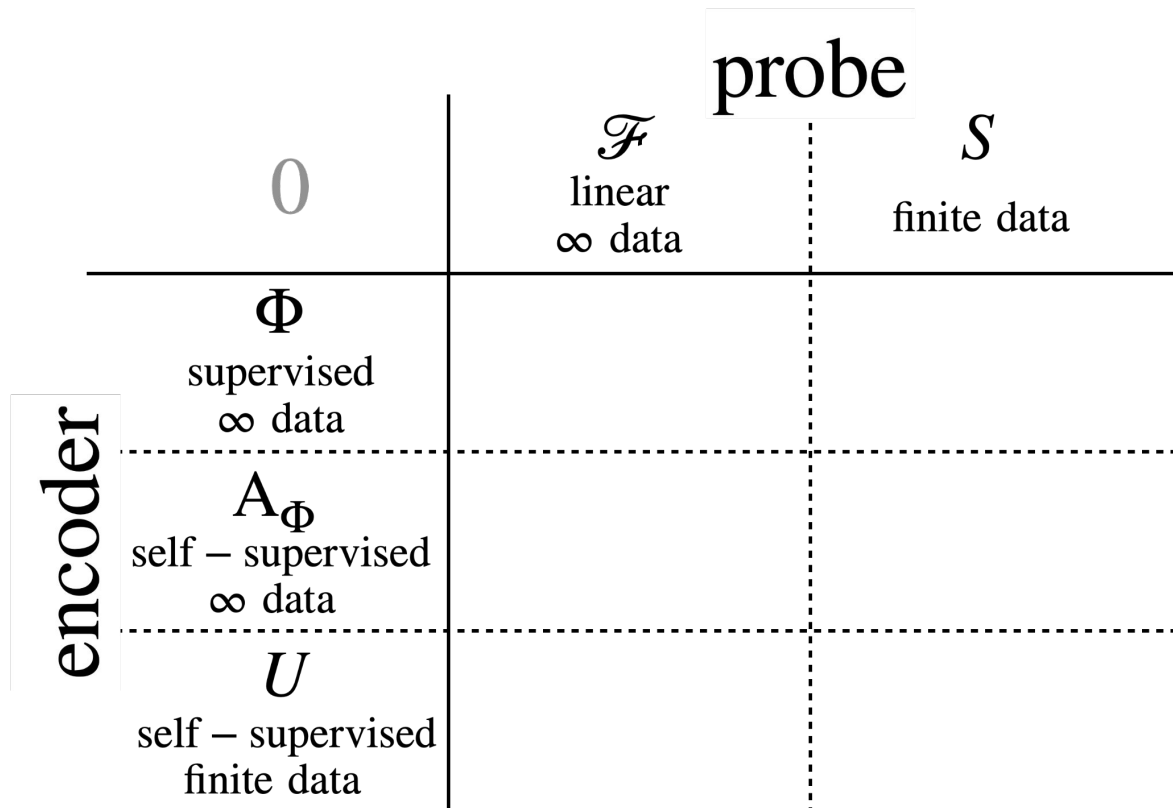
$R_{\mathcal{F}}$

estimation

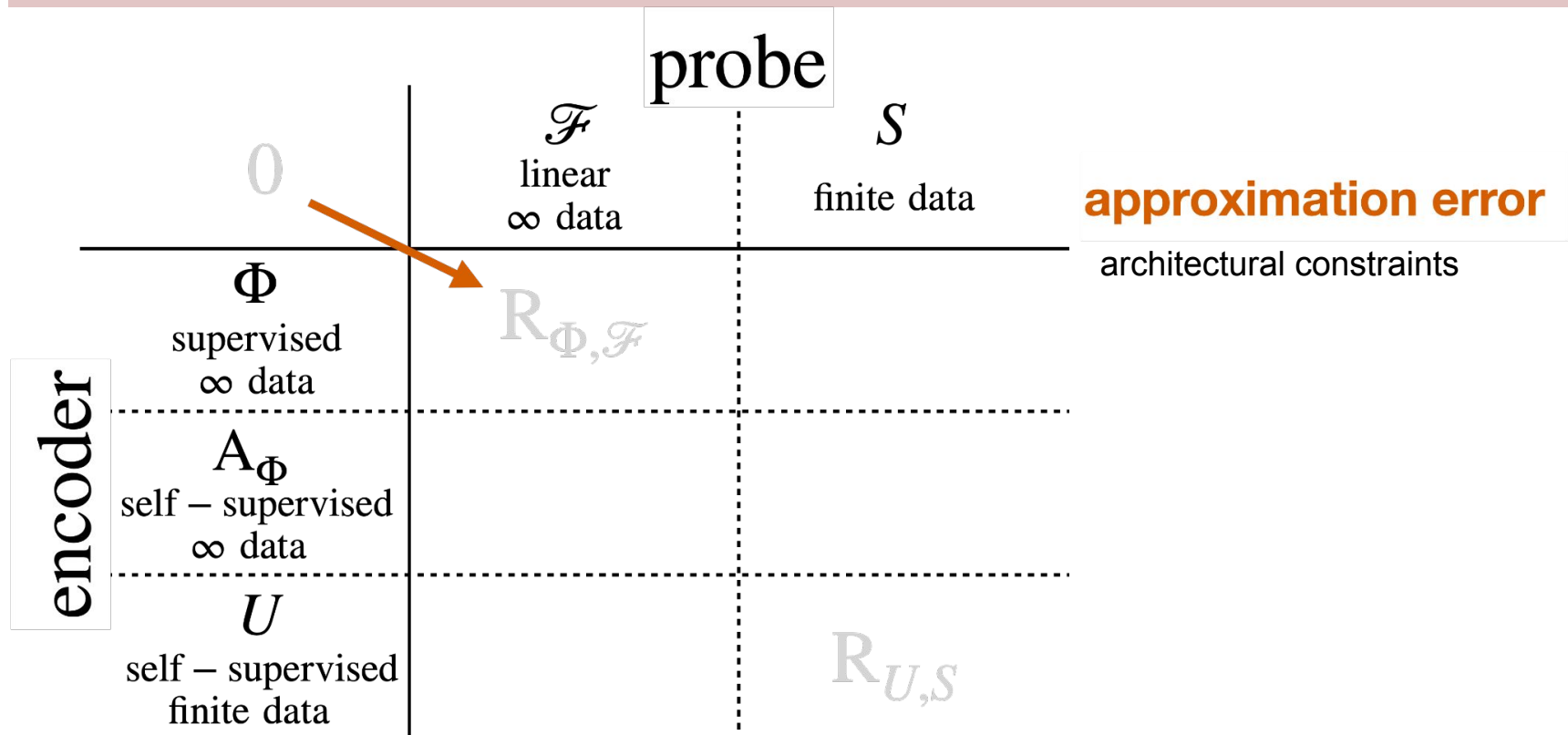
$R_S$

Total risk

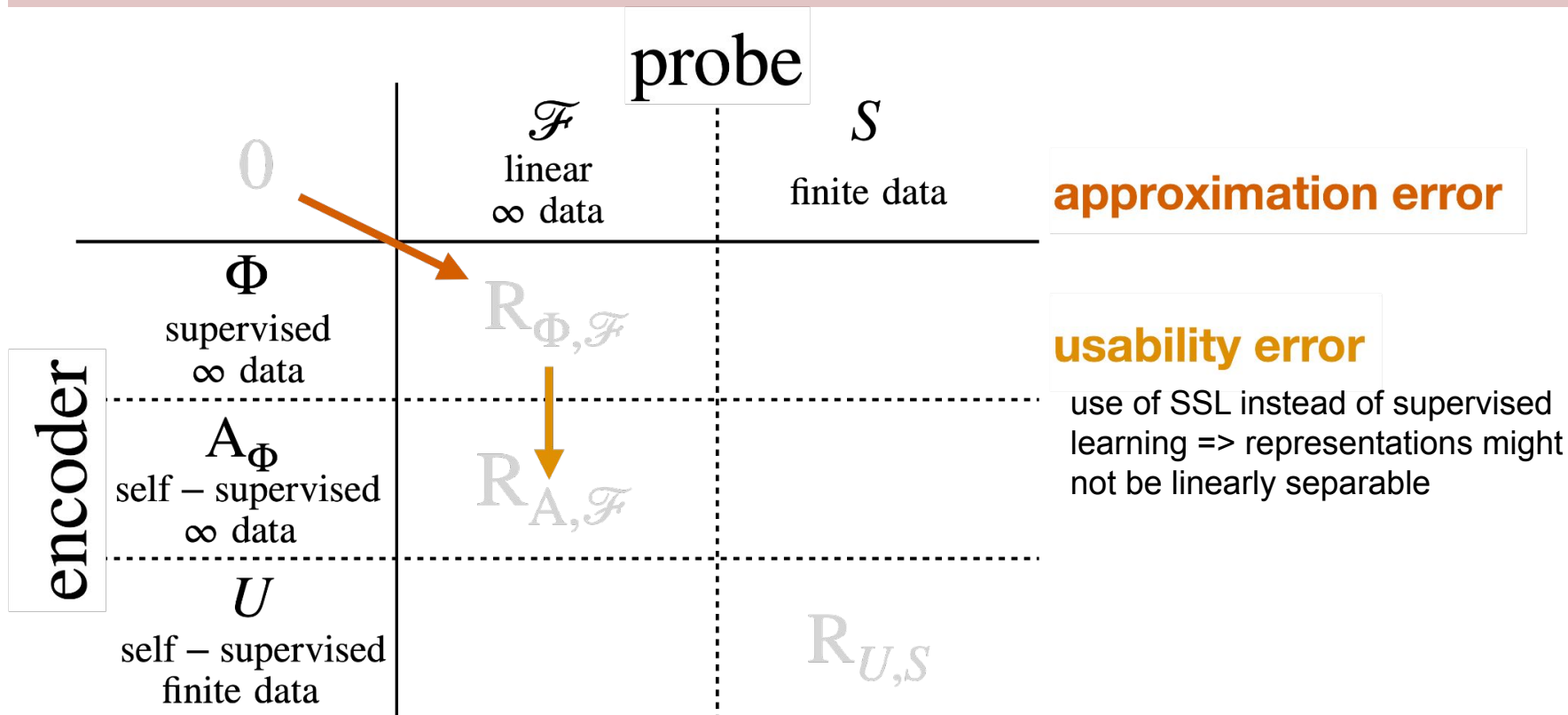
# Proposed: SSL risk decomposition



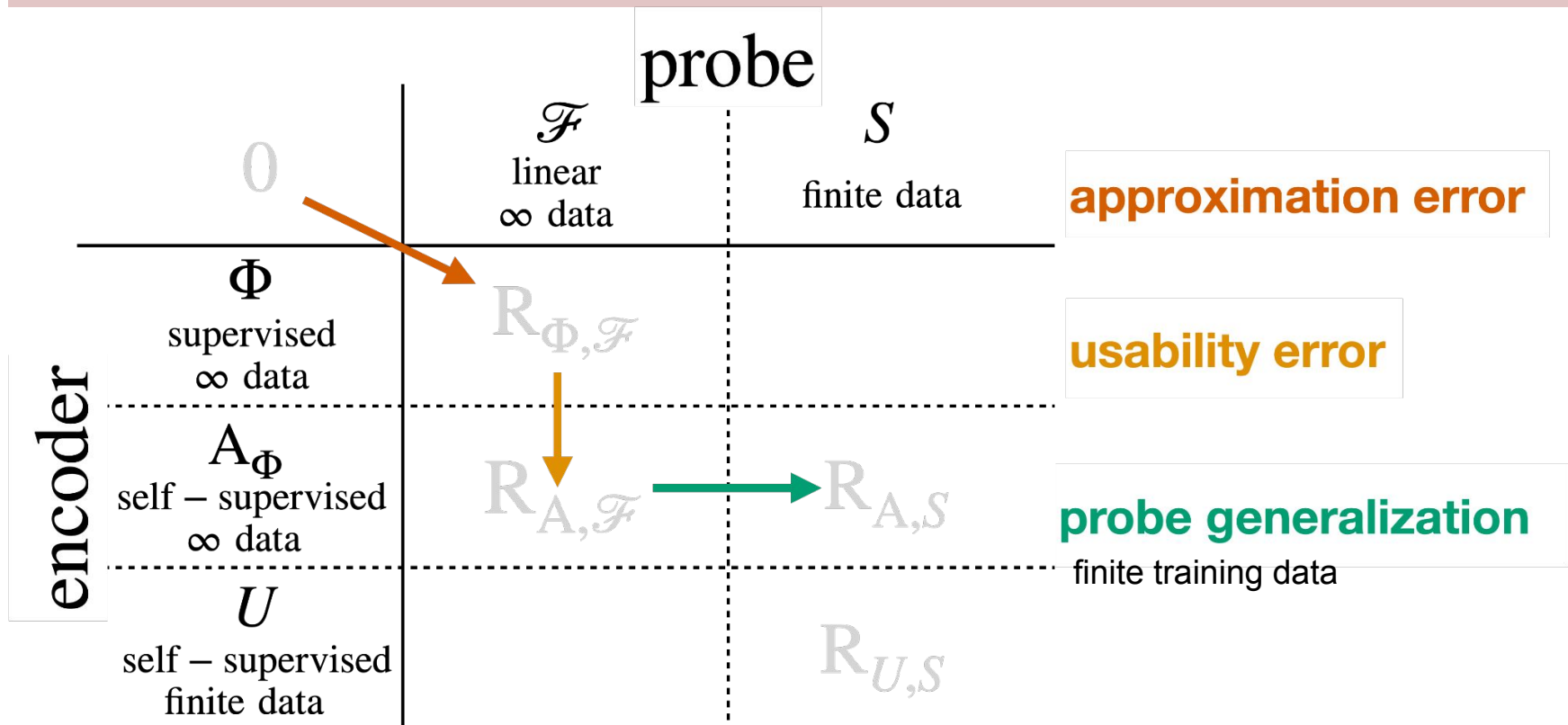
# Proposed: SSL risk decomposition



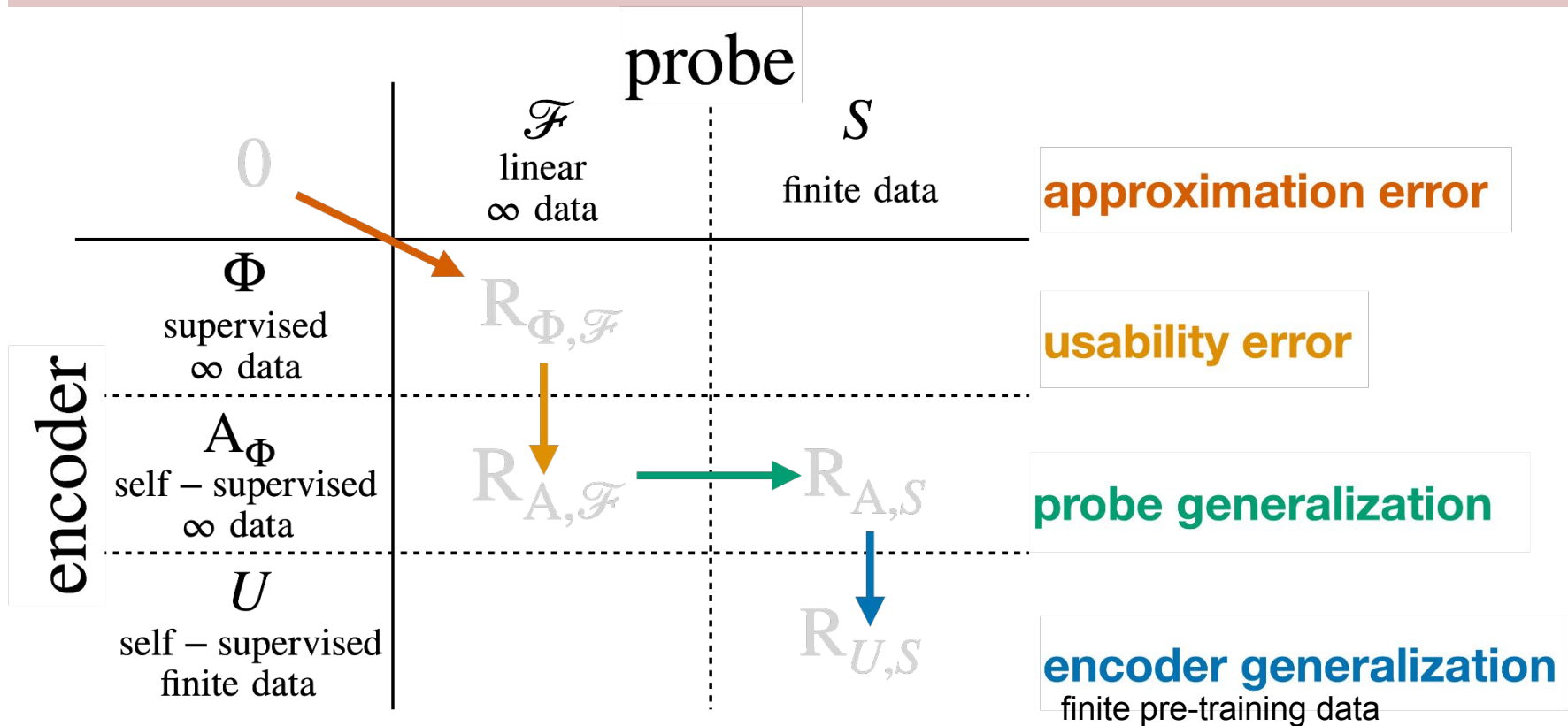
# Proposed: SSL risk decomposition



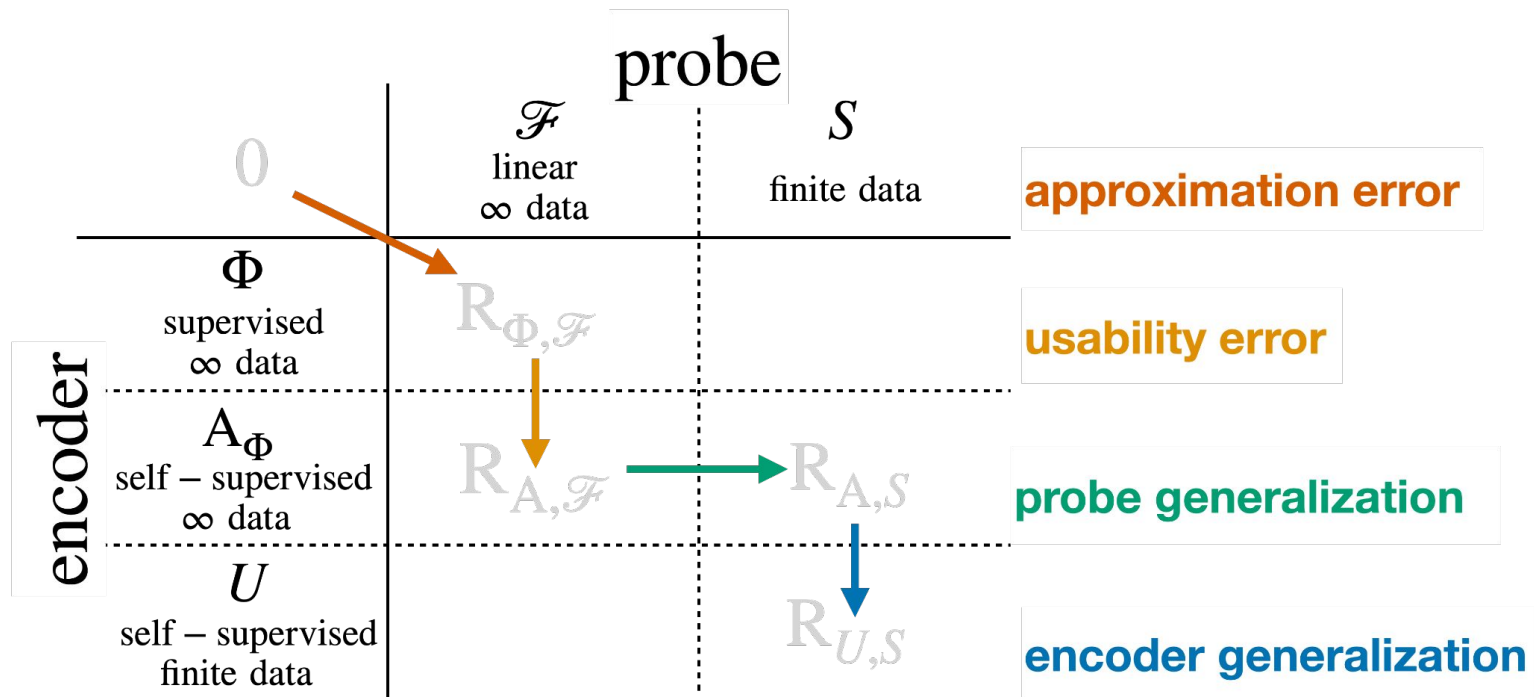
# Proposed: SSL risk decomposition



# Proposed: SSL risk decomposition



# Proposed: SSL risk decomposition



We provide efficient estimators for each component!



# Experiments: supervised risk decomposition

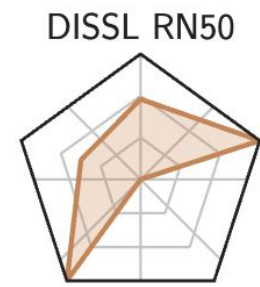
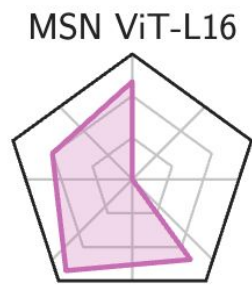
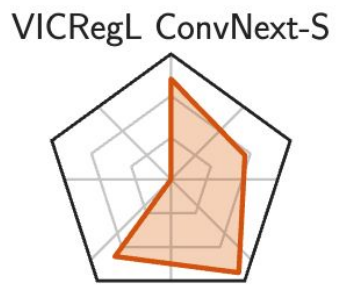
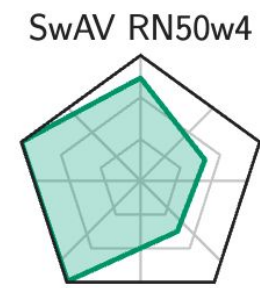
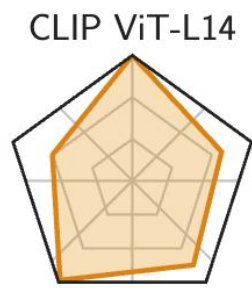
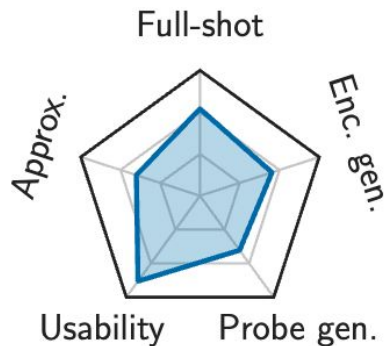
## Broad evaluation of SSL methods:

169 pretrained encoders, 28 objectives, 20 arch., 7 years

## Benchmark:

- linear probes on ImageNet (100%, 30-shot, 1%, 5-shot, 3-shot)
- estimate each error component

# Results: no model is uniformly better



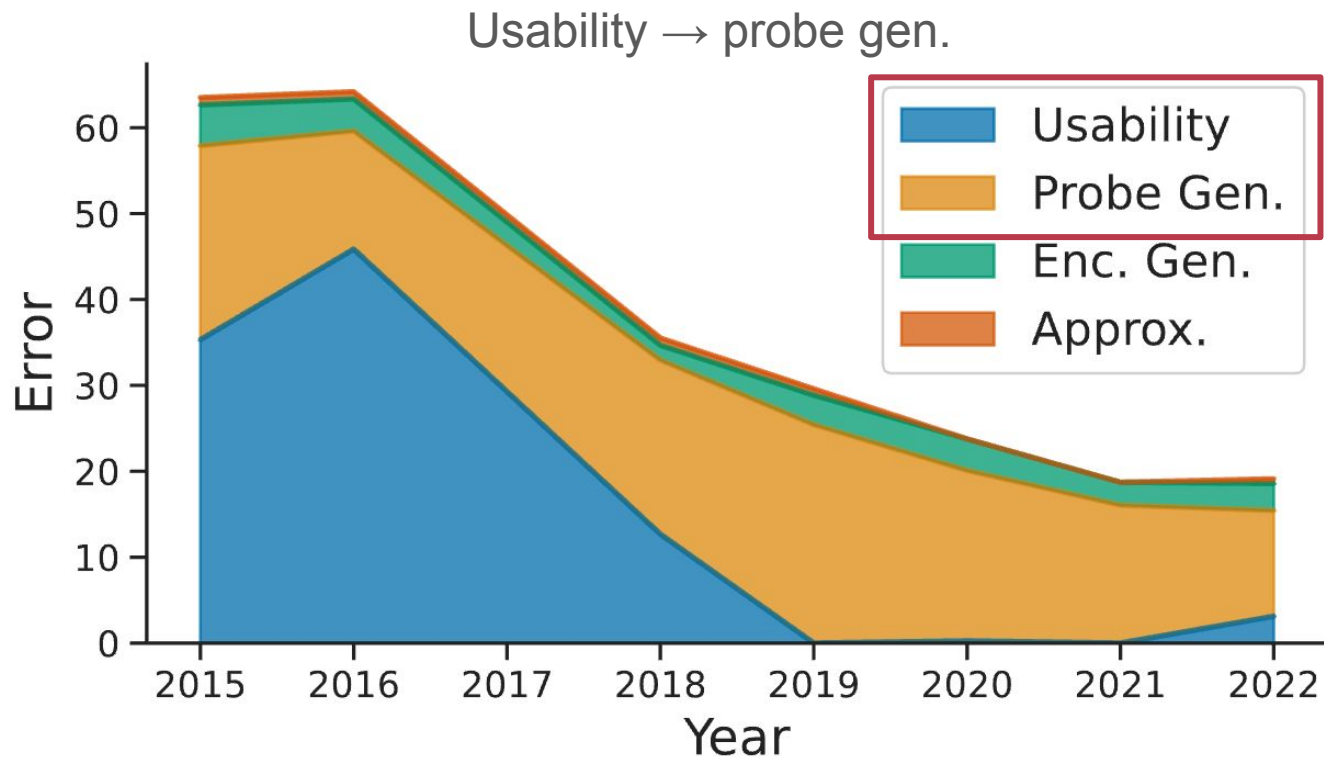
no model is uniformly better!

## Results: Full- vs Few-shot Tradeoff

Obj.	Arch.	Param.	ImageNet probe acc.		
			100%	1%	3-shot
MoCo-v3	RN50	24M	73.7	55.5	40.4
DINO	RN50	24M	74.2	52.9	35.9
SwAV	RN50w4	375M	76.2	56.2	36.9
VICRegL	CnvNxt-B	85M	74.8	64.3	56.3
MUGS	ViT-S16	22M	77.3	62.9	49.6
MSN	ViT-S16	22M	76.1	67.5	60.4
MSN	ViT-B4	86M	80.1	75.1	69.3
MUGS	ViT-L16	303M	80.9	74.0	68.5
MSN	ViT-L7	303M	79.9	74.9	69.8
CLIP	ViT-L14	304M	85.0	75.2	62.9
OpenCLIP	ViT-H14	632M	84.4	75.8	63.7

the best model in full-shot  
is always different than in  
few-shot

# Results: risk components over time



# Results: implication for SSL method design

	# dim. ↓	# views ↑	ViT	# param.↑	MLP proj.	generative SSL	# epoch ↑	Adam
Usability error	↑	↓		↓	↓	↑	↓	
Probe gen. error	↓	↓	↓		↓	↓	↓	↓
Full-shot error	↑	↓	↓	↓	↓	↑	↓	↓
3-shot error	↓	↓	↓	↓	↓	↓	↓	↓

## Results: dimensionality

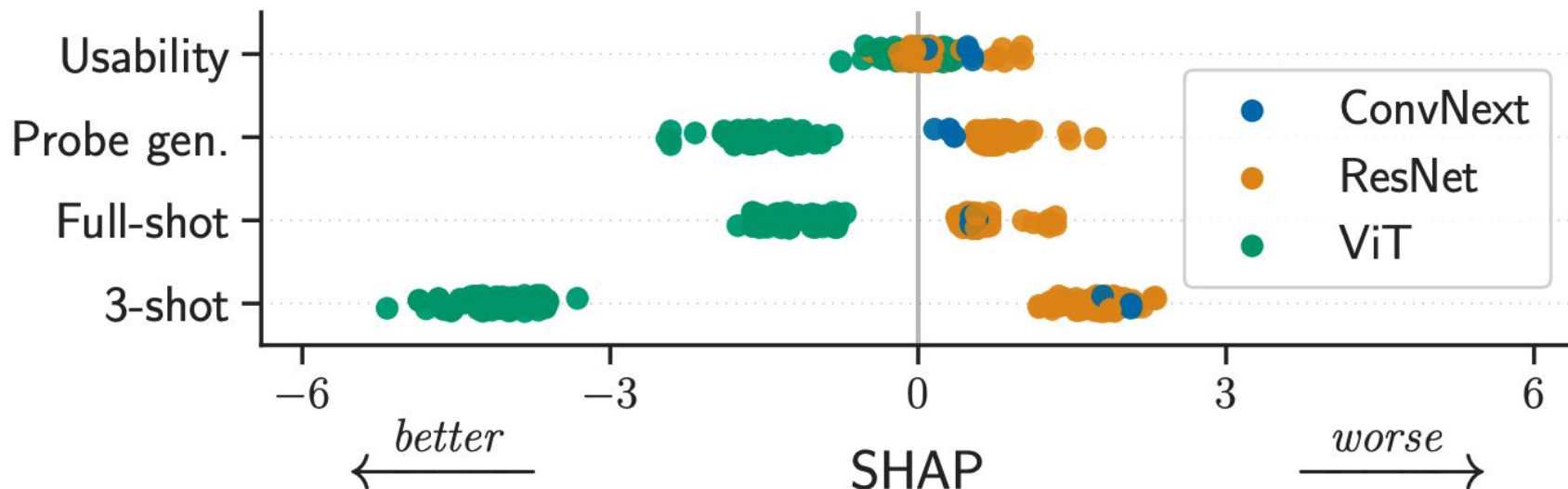
Some design choices (e.g. dimensionality) can control U-P tradeoff => full- vs few-shot

Ours	Obj.	ViT	Dim.	100%	1%	3-shot
✗	MUGS	S16	1536	<b>77.3</b>	62.9	49.6
✓	MUGS	S16	384	77.0	<b>66.6</b>	<b>57.9</b>
✗	OpenCLIP	H14	1280	<b>84.4</b>	75.8	63.7
✓	OpenCLIP	H14	1024	84.3	<b>76.5</b>	<b>65.5</b>

by decreasing dimensionality we can greatly improve  
few shot performance without any retraining!

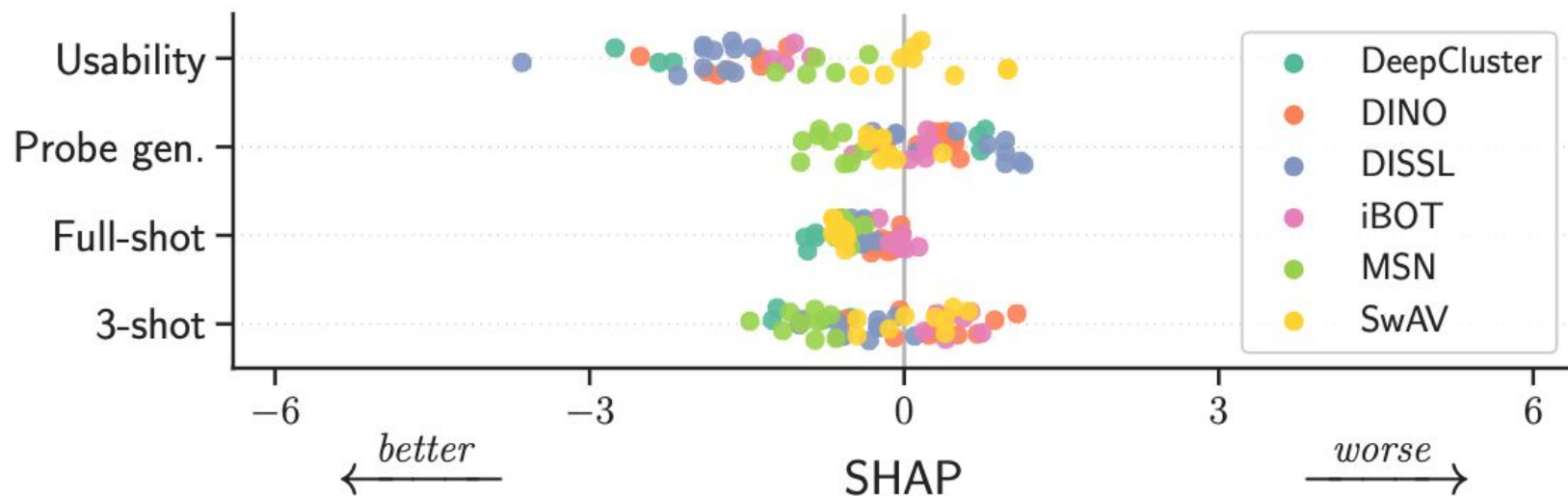
# Results: architecture

Other design choices (e.g. architecture) overcome the tradeoff => uniform improvement



# Results: exact objective

Other design choices (e.g. the exact objective in generative or contrastive) don't matter when controlling for confounders!





# Summary

- New risk decomposition for SSL with efficient estimators

$$\underbrace{R_{U,S} - R_*}_{\text{excess risk}} = \underbrace{R_{U,S} - R_{A,S}}_{\text{encoder generalization}} + \underbrace{R_{A,S} - R_{A,\mathcal{F}}}_{\text{probe generalization}} + \underbrace{R_{A,\mathcal{F}} - R_{\Phi,\mathcal{F}}}_{\text{representation usability}} + \underbrace{R_{\Phi,\mathcal{F}} - R_*}_{\text{approximation}}$$

---

## Algorithm 1 Estimating risk components in the standard SSL setting

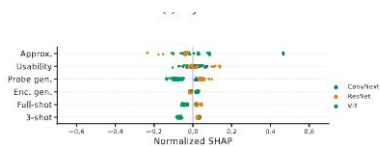
---

**Require:** Encoder family  $\Phi$ , probe family  $\mathcal{F}$ , training  $S_{tr}$  and testing  $S_{te}$  sets, SSL algorithm  $A_\Phi$ , evaluation loss  $\ell$ .

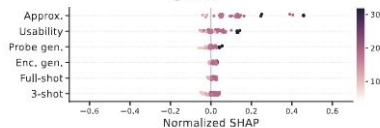
- 1: **function** RISK( $\mathcal{F}$ ,  $\mathcal{D}_{tr}$ ,  $\mathcal{D}_{te}$ ) ▷ Risk minimization
  - 2:    $\hat{f} \leftarrow \inf_{f \in \mathcal{F}} \sum_{(x,y) \in \mathcal{D}_{tr}} \ell(y, f(x))$  ▷ Test risk
  - 3:   **return**  $\frac{1}{|\mathcal{D}_{te}|} \sum_{(x,y) \in \mathcal{D}_{te}} \ell(y, \hat{f}(x))$
  - 4:  $\hat{R}_{\Phi,\mathcal{F}} \leftarrow \text{RISK}(\Phi \circ \mathcal{F}, S_{tr}, S_{tr})$  ▷ Supervised train performance
  - 5:  $\phi \leftarrow A_\Phi(\Phi, S_{tr})$  ▷ Pretrain SSL encoder
  - 6:  $S_{tr}^\phi \leftarrow [(\phi(x), y) \text{ for } x, y \text{ in } S_{tr}]$  ▷ Featurize data
  - 7:  $S_{te}^\phi \leftarrow [(\phi(x), y) \text{ for } x, y \text{ in } S_{te}]$
  - 8:  $S_{sub}^\phi \leftarrow \text{subset}(S_{tr}^\phi, n = \text{len}(S_{te}^\phi))$
  - 9:  $\hat{R}_{A,\mathcal{F}} \leftarrow \text{RISK}(\mathcal{F}, S_{tr}^\phi, S_{tr}^\phi)$  ▷ Risk without generalization
  - 10:  $\hat{R}_{A,S} \leftarrow \text{RISK}(\mathcal{F}, S_{tr}^\phi \setminus S_{sub}^\phi, S_{sub}^\phi)$  ▷ Risk with only probe gen.
  - 11:  $\hat{R}_{U,S} \leftarrow \text{RISK}(\mathcal{F}, S_{tr}^\phi, S_{te}^\phi)$  ▷ Risk with enc. and probe gen.
  - 12:  $\text{approx\_error} \leftarrow \hat{R}_{\Phi,\mathcal{F}}$
  - 13:  $\text{usability\_error} \leftarrow \hat{R}_{A,\mathcal{F}} - \hat{R}_{\Phi,\mathcal{F}}$
  - 14:  $\text{probe\_gen} \leftarrow \hat{R}_{A,S} - \hat{R}_{A,\mathcal{F}}$
  - 15:  $\text{encoder\_gen} \leftarrow \hat{R}_{U,S} - \hat{R}_{A,S}$
  - 16: **return** approx\_error, usability\_error, probe\_gen, encoder\_gen
-

# Summary

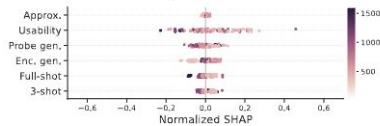
- New risk decomposition for SSL with efficient estimators
- Meta-analysis of 169 models and 30 design choices



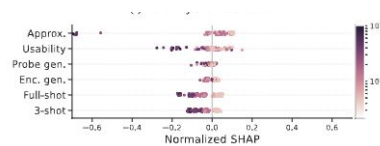
(g) Family



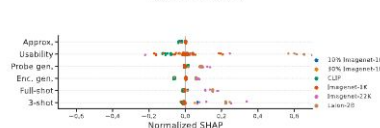
(i) Patch Size



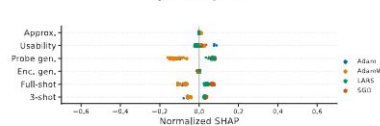
(k) Epochs



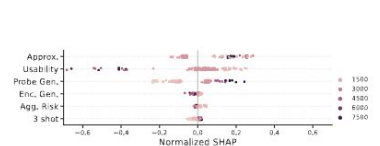
(h) Num. Parameters



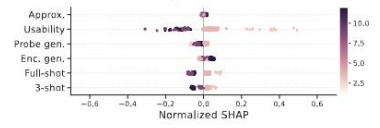
(j) Pretraining Data



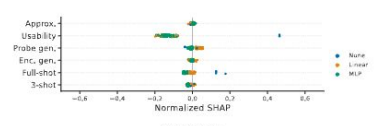
(l) Optimizer



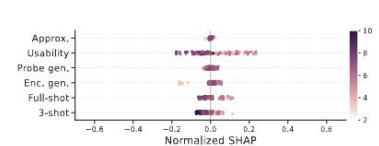
(a) Z Dim.



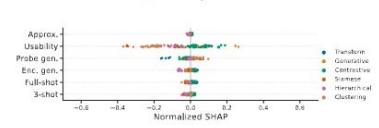
(c) Num. Views



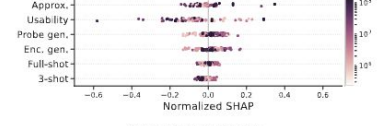
(e) Proj. Arch.



(b) Num. Augmentations



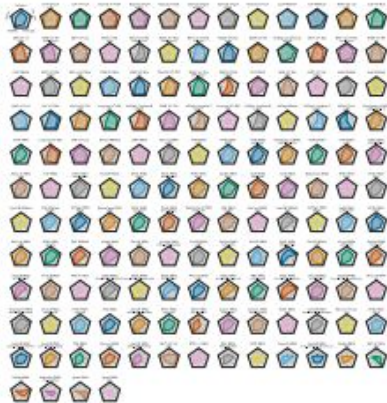
(d) SSL Mode



(f) Num. Projection Parameters

# Summary

- New risk decomposition for SSL with efficient estimators
- Meta-analysis of 169 models and 30 design choices
- Many more results in the paper!
  - Thorough analysis of each design choice
  - Large scale evaluation of SSL with different metrics



Model	Risk	Metric	Risk Decomposition					Information				
			Model	Risk	Metric	Model	Risk	Metric				
...	...	...	...	...	...	...	...	...	...	...	...	...

Model	Risk	Metric	Risk Decomposition					Information				
			Model	Risk	Metric	Model	Risk	Metric				
...	...	...	...	...	...	...	...	...	...	...	...	...

Model	Risk	Metric	Risk Decomposition					Information				
			Model	Risk	Metric	Model	Risk	Metric				
...	...	...	...	...	...	...	...	...	...	...	...	...

# Summary

- New risk decomposition for SSL with efficient estimators
- Meta-analysis of 169 models and 30 design choices
- Many more results in the paper!
- Torch Hub API & [code](#) to access any models or metadata in one line



[SSL-Risk-Decomposition](#)

```
import torch

# loads the desired pretrained model and preprocessing pipeline
name = "dino_rn50" # example
model, preprocessor = torch.hub.load('YannDubs/SSL-Risk-Decomposition:main', name, trust_repo=True)

# gets all available models
available_names = torch.hub.list('YannDubs/SSL-Risk-Decomposition:main')

# gets all results and hyperparameters as a dataframe
results_df = torch.hub.load('YannDubs/SSL-Risk-Decomposition:main', "results_df")
```