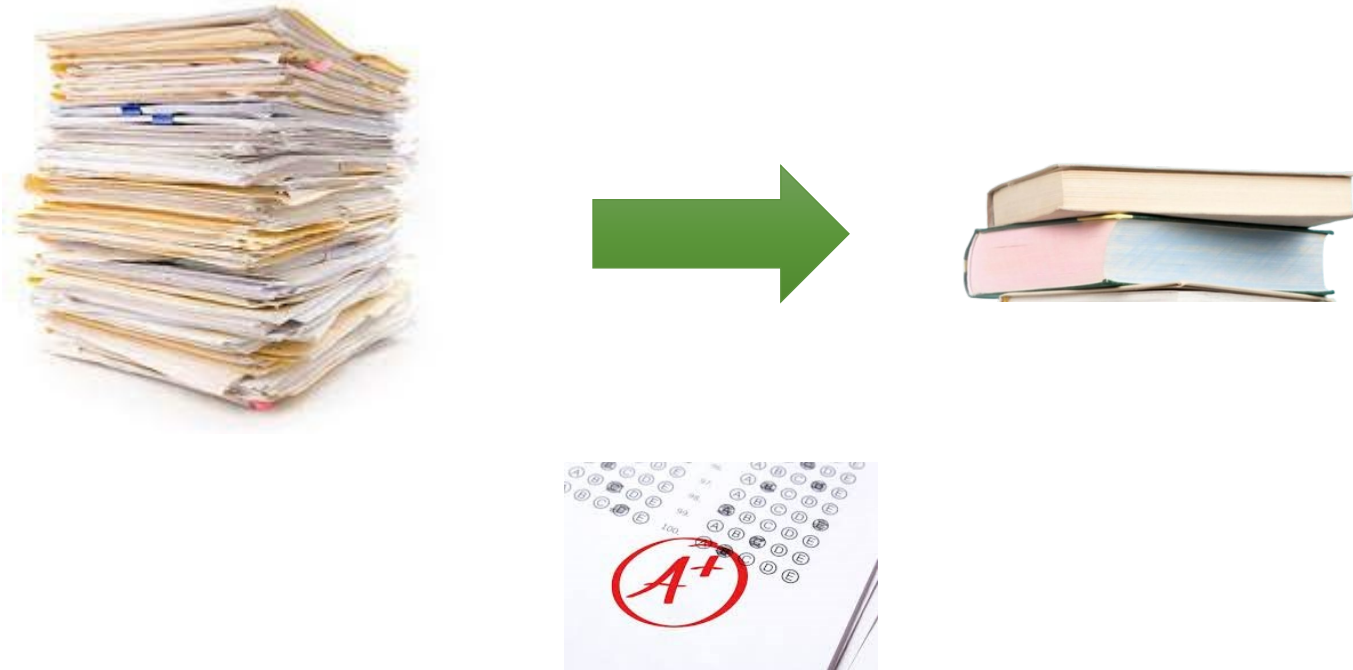# **Intuition**



- Less is More
- Small high quality textbook like data
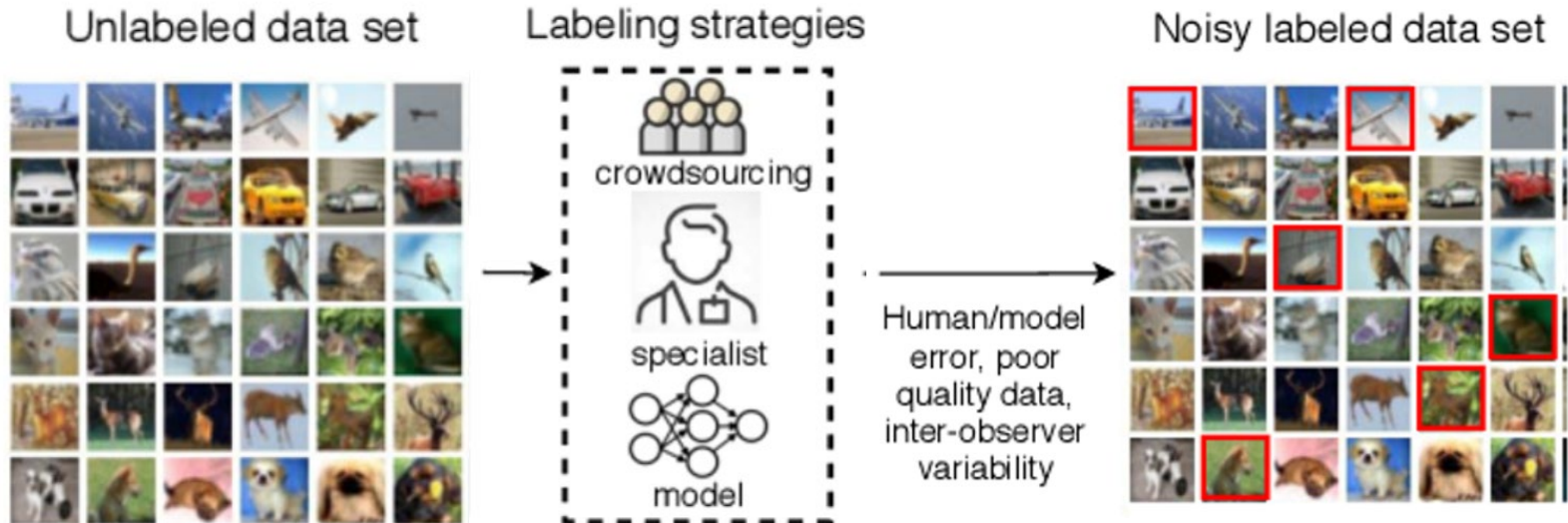
# Background

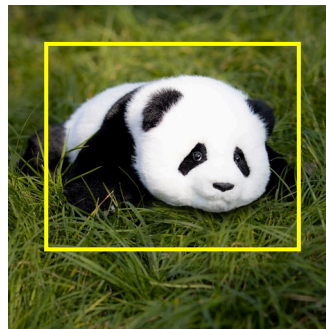## Data centric AI



Sachdeva, Noveen, et al. "Infinite recommendation networks: A data-centric approach." Neurips'22

# Motivation

## Data are noisy

2020, A Survey on Deep Learning with Noisy Labels: How to train your model when you cannot trust on the annotations?

# Motivation

## Big data cost



OpenAI, AI and compute, 2018

# Research Question

- Reduce data size while retaining the model's performance?
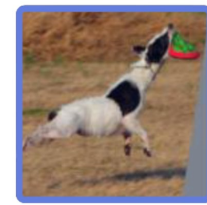


Data Distillation — << 50K distilled images — Train → Similar Performance

50K images — Train → Learning algorithm

Sachdeva, Noveen, and Julian McAuley. "Data distillation: A survey." TMLR'23

# Literature Review

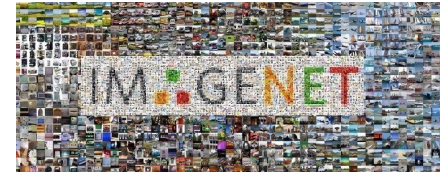

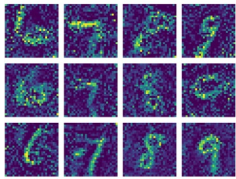a black and white dog is running through the field to catch something in its mouth

a black and white dog leaps to catch a frisbee in a field

MNIST



$$S_c \leftarrow S_c - \lambda \nabla_{S_c} D\left(\mathcal{A}\left(S_c\right), \mathcal{A}\left(T_c\right), \theta\right)$$

**2021~2022**
Different match objective function

**2023**
Large scale dataset, LLM, generative, multimodal

**2018 DD**

**2020 DC**



(b)

Large training set

Update synthetic set

Matching Loss

CE Loss

Forward pass

Backpropagation

CE Loss

Small synthetic set

**2022**
Different data modality, application (privacy, med)



$(A, X, Y)$

Condense

$(A', X', Y')$

153,932 training nodes

154 training nodes

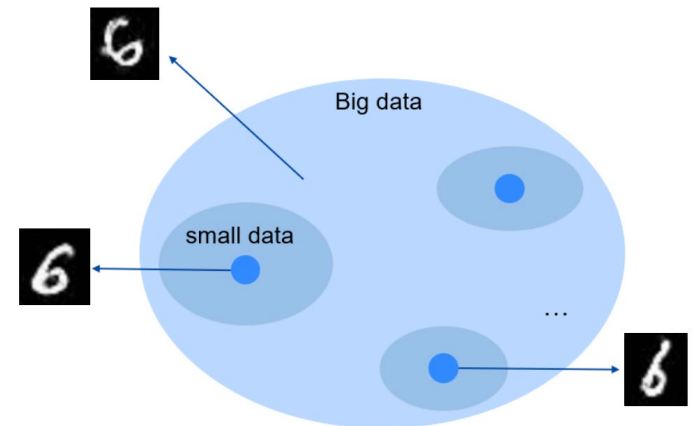Sachdeva, Noveen, and Julian McAuley. "Data distillation: A survey." TMLR'23

# Critical Literature Review 1

## Herustics data selection

### Select hard samples



### Cluster sampling
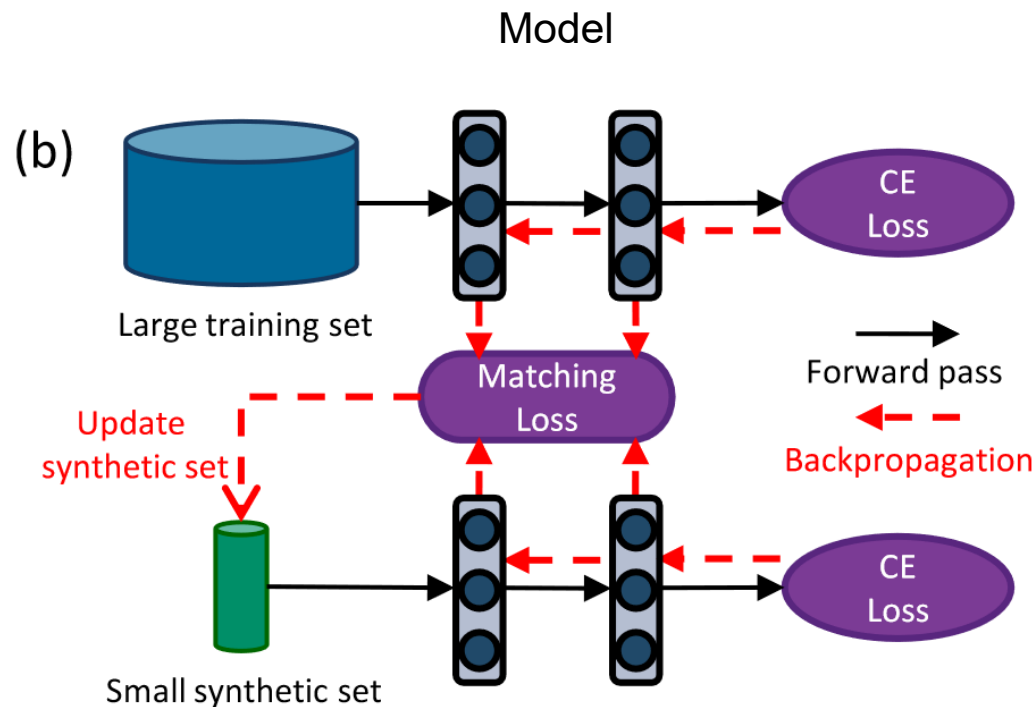


Fast, but may not be accurate

Jeong, Y et.al. (2023). Training data selection based on dataset distillation for rapid deployment in machine-learning workflows.

# Critical Literature Review 2

## Feedback optimization framework



Feedback

Select data

Accurate, but slow

Wang, T., et al. (2018). Dataset distillation. *arXiv preprint*

# Critical Literature Review 2

## Feedback optimization framework

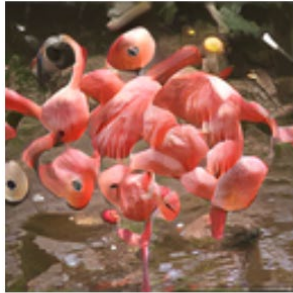Wang, T., et al. (2018). Dataset distillation. *arXiv preprint*

# Critical Literature Review 2

## Generative DD



ImageNet-Birds: Peacock, Flamingo, Macaw

ImageNet-Fruits: Pineapple, Banana, Strawberry
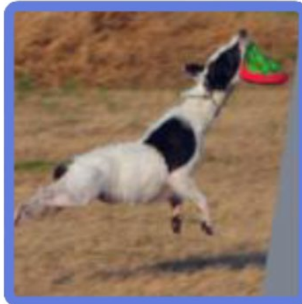
## DD



CIFAR-100: Apple — ImageNet: Banana, Camel

## Image-text



a black and white dog is running through the field to catch something in its mouth

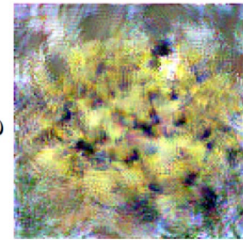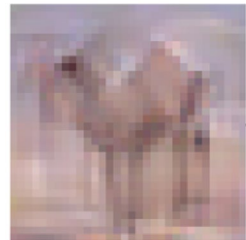a black and white dog leaps to catch a frisbee in a field

## Medical, privacy

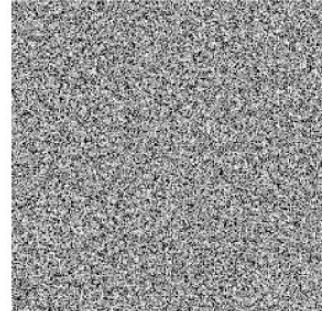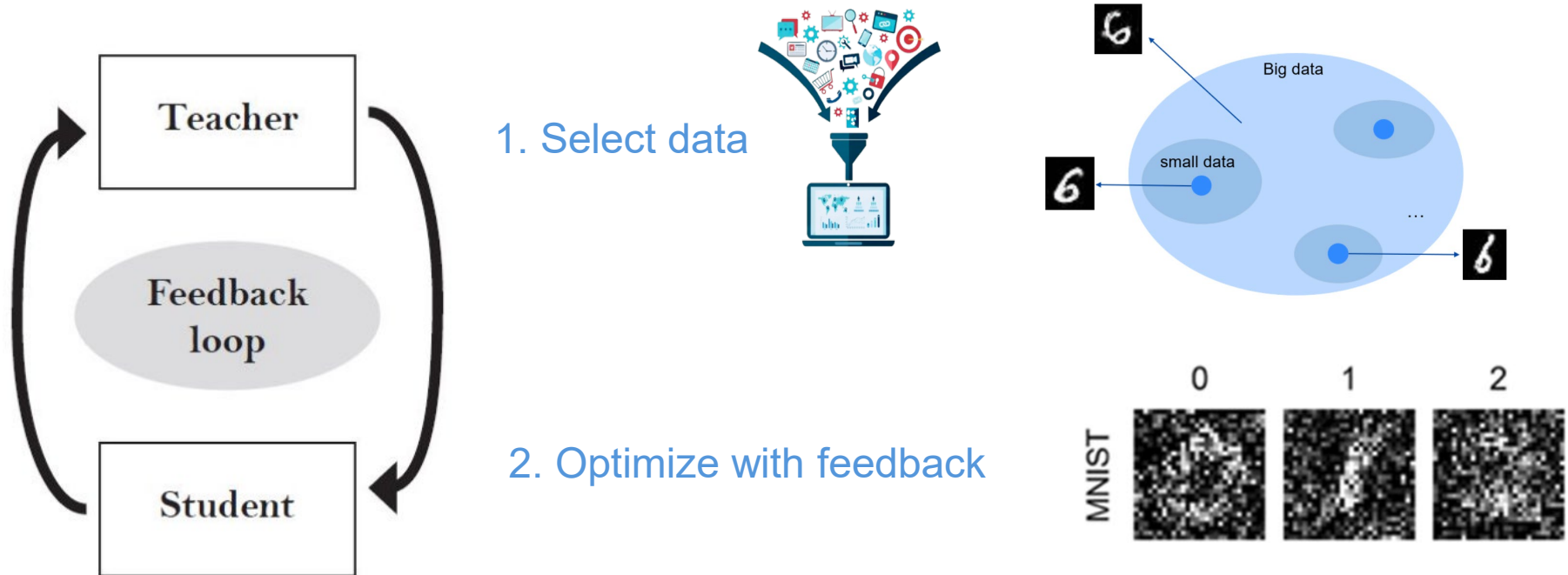Sachdeva, Noveen, and Julian McAuley. "Data distillation: A survey." TMLR'23

# Research gap

- – Herustics methods are fast but not accurate
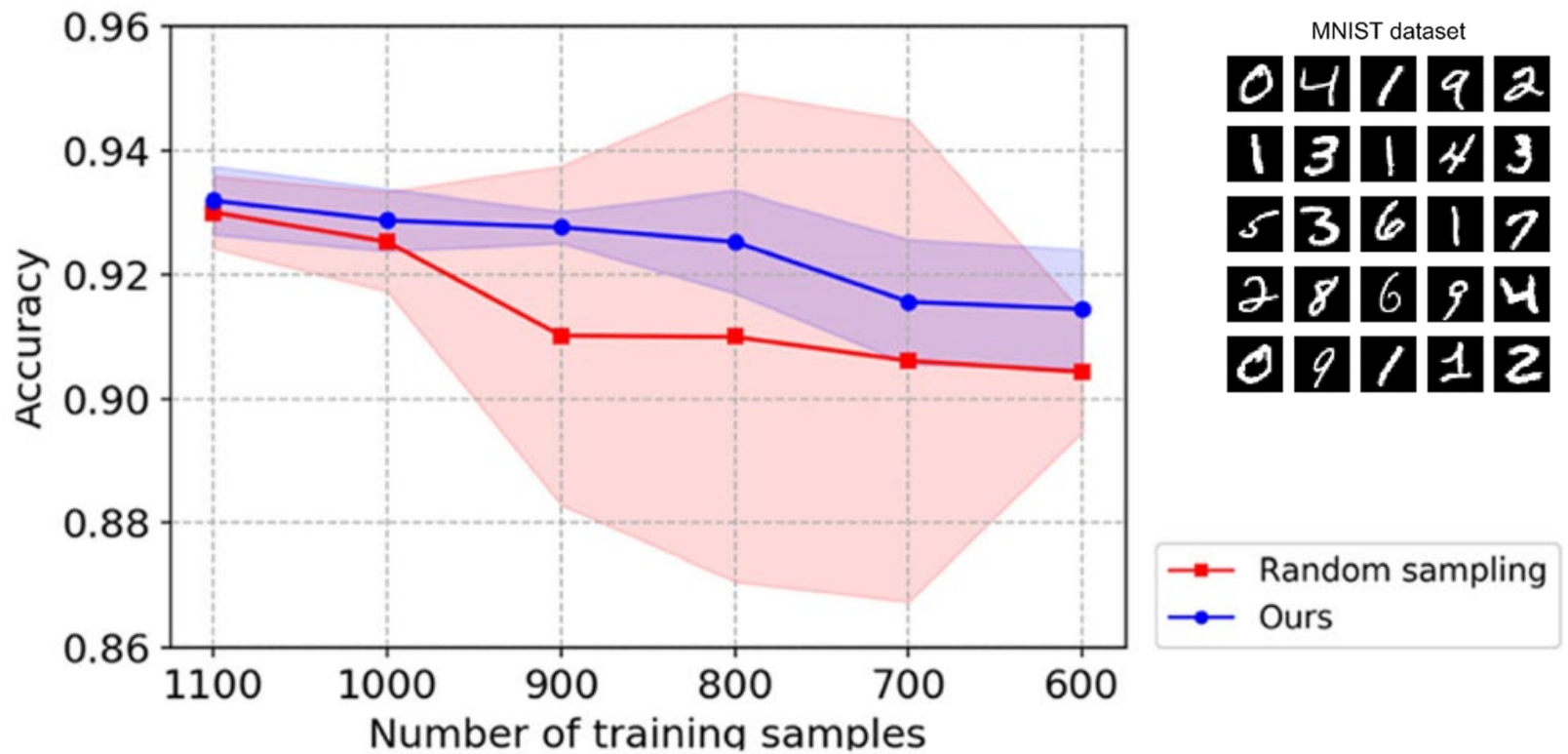- – Optimization methods are accurate but slow

# Research Method

Heuristic selection (fast, not accurate) + Feedback (accurate, but slow)



1. Select data

2. Optimize with feedback

- Get student feedback on representative data not all

# Experiment

## Perform better w/ lower variance



MNIST dataset

# Experiment

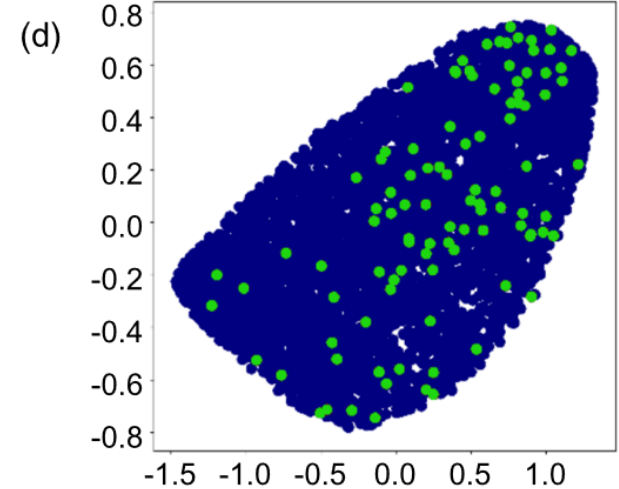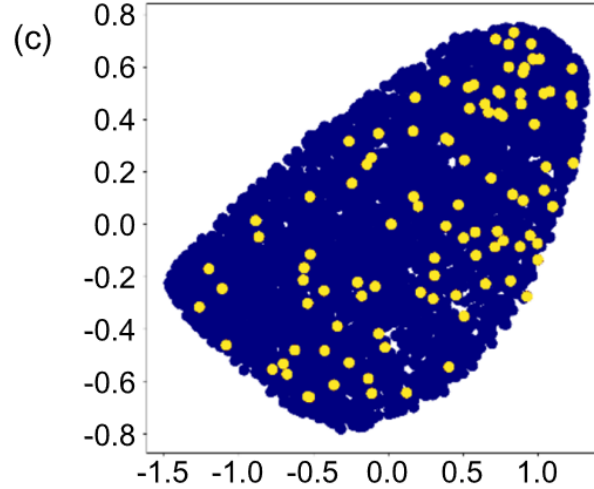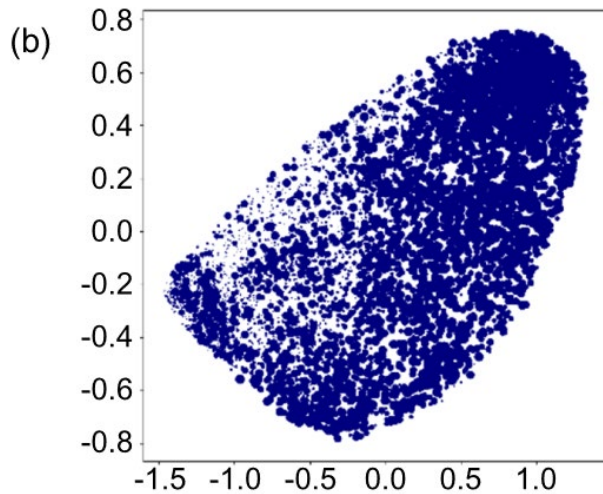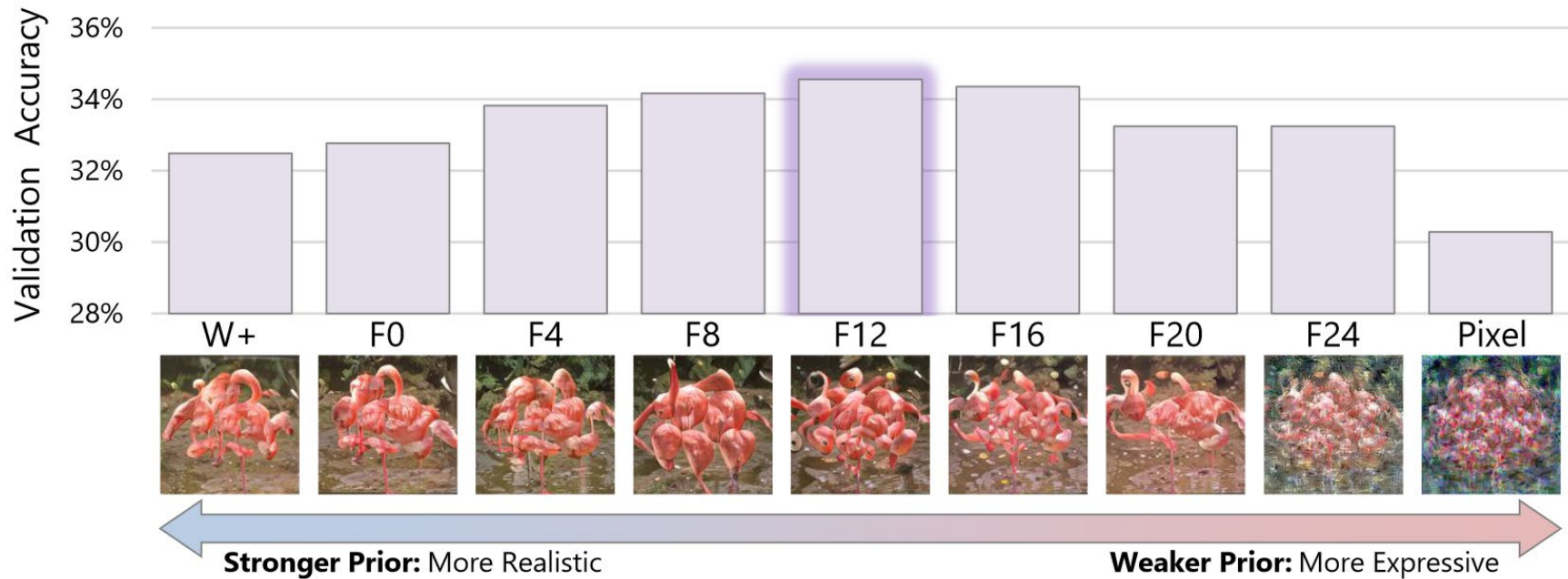b) Varying marker size w.r.t sample difficulty

c) Proposed method

d) Baselines

# Limitation



Trade-of between realism & expressiveness (performance)

Cazenavette, George, et al. "Generalizing dataset distillation via deep generative prior." CVPR'23.

# Conclusion

- Motivation
  - Big data cost
- RQ
  - Reduce big data to small while retaining performance
- Prev solution
  - Either slow or not accurate
- Proposed method
  - Heuristcs (fast) + feedback framework (accurate)

# Thank you for listening!

Any questions?



Small data



Original big data