# projektor | Performance Scaling via Optimal Transport: Enabling Data Selection from Partially Revealed Sources

**Feiyang Kang[1]\*, Hoang Anh Just[1]\*, Anit Kumar Sahu[2], Ruoxi Jia[1]**

[1]*Virginia Tech*  [2]*Amazon Alexa AI*

VIRGINIA TECH · amazon · REDS LAB RESPONSIBLE DATA SCIENCE · AI alexa

NEURAL INFORMATION PROCESSING SYSTEMS · NeurIPS | 2023

We approach practical data collection scenarios with multiple sources, where acquisition plans need to be made with only small samples. We propose a handy toolkit, **projektor**, that predicts model performance, projects it onto larger scales, and optimizes over predictions.
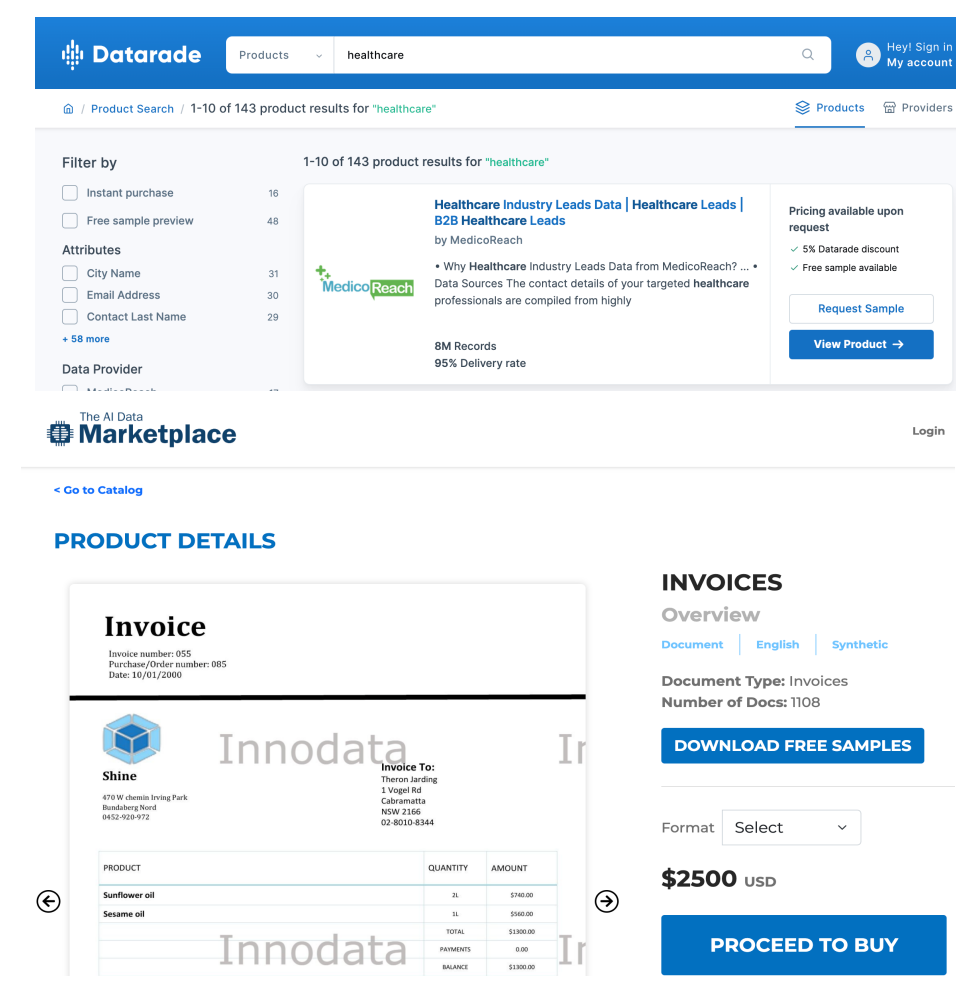
## What's the problem?

"data is the new oil"–the choice of training data is crucial for extracting the best performance out of a model.

Data is typically acquired from **various sources**, such as different organizations or vendors (e.g., data marketplace)
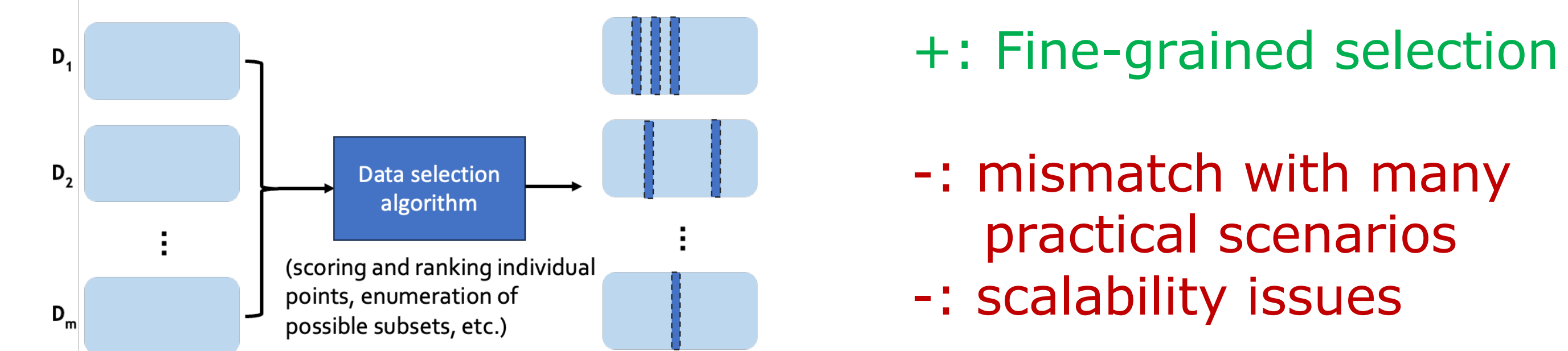
in practical scenarios, data providers often reveal **only a limited subset** of samples before an acquisition decision is made.

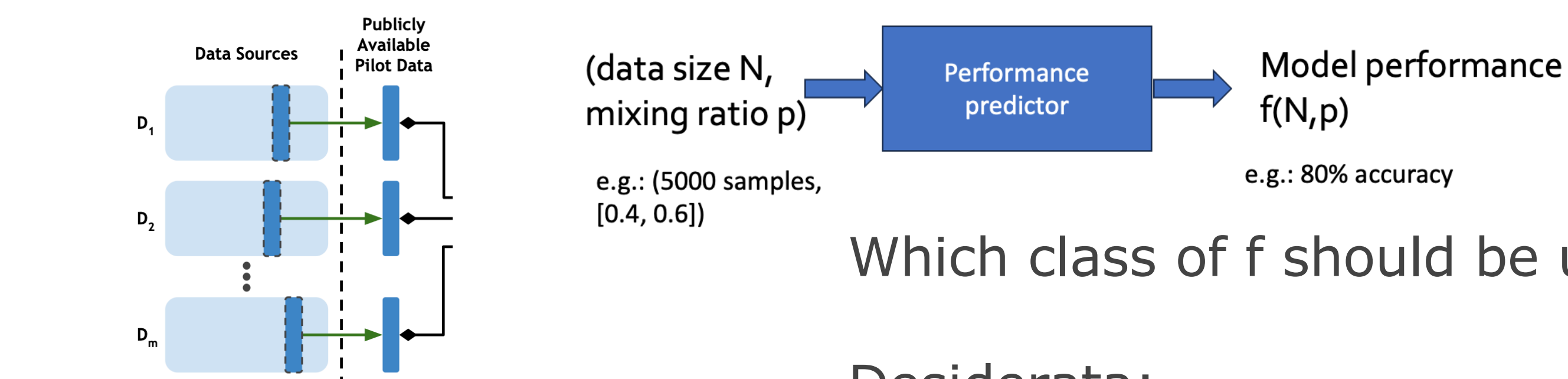**How to select and combine samples from these data sources?**

## Our Framework

strategic data selection in **partially observable** settings, where only limited samples of data sources (pilot datasets) are accessible.

→ The goal is to determine an **optimal allocation** of the selection budget to each source, based only on pilot datasets, such that the model trained on the mixture of collected data achieves the best result on given objectives.

## Limitations of Past Work

Most assumes **complete** access to potential data sources (coreset selection, active learning, data valuation, etc.)

+: Fine-grained selection

-: mismatch with many practical scenarios
-: scalability issues

(scoring and ranking individual points, enumeration of possible subsets, etc.)

Without complete access of data, we cannot directly evaluate a plan

## Our Approach: projektor

**Key idea:** learn a performance predictor

(data size N, mixing ratio p) → Performance predictor → Model performance $f(N,p)$

e.g.: (5000 samples, [0.4, 0.6]) → e.g.: 80% accuracy
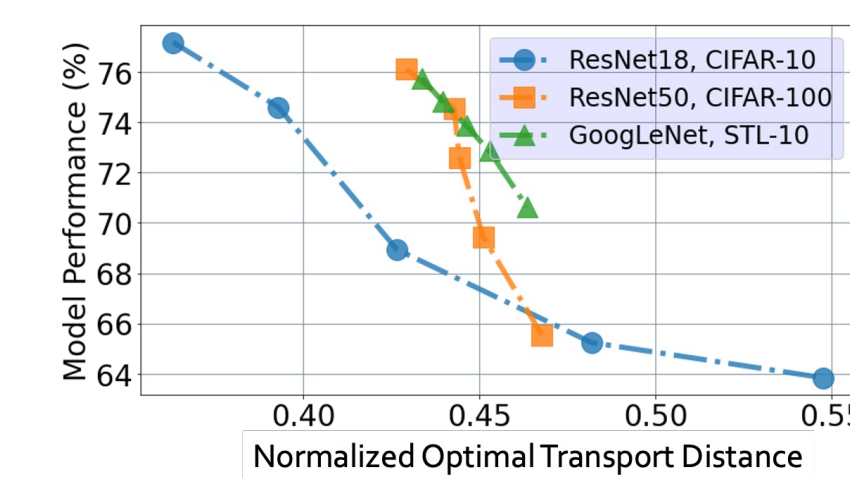
Which class of f should be used?

Desiderata:
- Accurate extrapolation to large N
- Easy to optimize over p

Intuition: The **more relevant** training data is to the validation, the **higher model performance**

Model Performance (%) vs Normalized Optimal Transport Distance
- ResNet18, CIFAR-10
- ResNet50, CIFAR-100
- GoogLeNet, STL-10

## Paper & Repository

**Paper.** openreview.net/pdf?id=quMBEd27x9

**Code.** github.com/ruoxi-jia-group/projektor

## Our Approach: projektor

Our approach: performance scaling via optimal transport

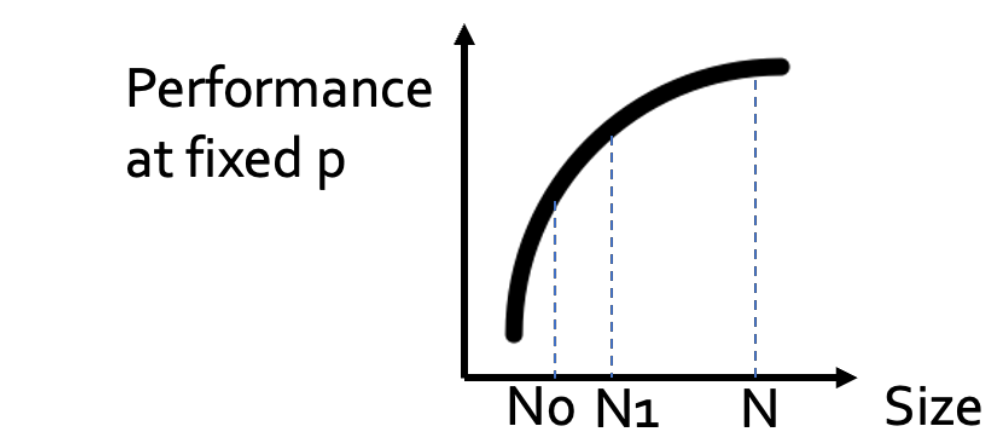Stage 1: Performance prediction **at any p but small N** — Validation set

$$\hat{\mathcal{L}}\left(\mathcal{A}(\mathcal{D}(N, \mathbf{p}), D^{\mathrm{val}}\right) = a_1 \cdot \mathrm{OT}\left(\mathcal{D}(N, \mathbf{p}), D^{\mathrm{val}}\right) + a_0$$

Performance of the model trained on D(N,p) and evaluated on D^val — Training set of size N and mixing ratio p

Stage 2: Parameter-free projection to **larger N**

Consider $\mathbb{E}_V[\mathcal{L}(\mathcal{A}(\mathcal{D}(N, \mathbf{p})); D^{\mathrm{val}})] = -\alpha(\mathbf{p})\log(N) + C(\mathbf{p})$

Then $\hat{\mathcal{L}}(\mathcal{A}(\mathcal{D}(N, \mathbf{p}); D^{\mathrm{val}}) = \left(\log\frac{N_1}{N_0}\right)^{-1}\left[\log\frac{N}{N_0}\hat{\mathcal{L}}(\mathcal{A}(\mathcal{D}(N_1, \mathbf{p})); D^{\mathrm{val}}) - \log\frac{N}{N_1}\hat{\mathcal{L}}(\mathcal{A}(\mathcal{D}(N_0, \mathbf{p})); D^{\mathrm{val}})\right]$

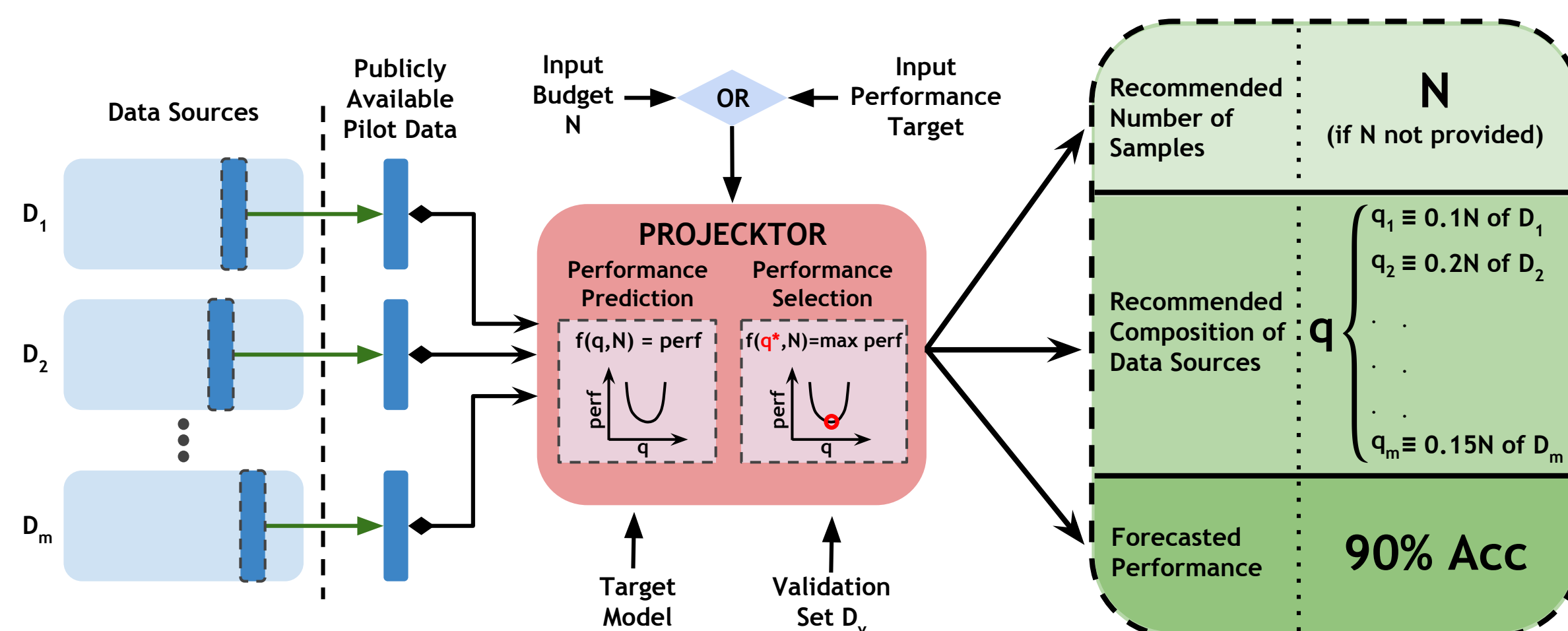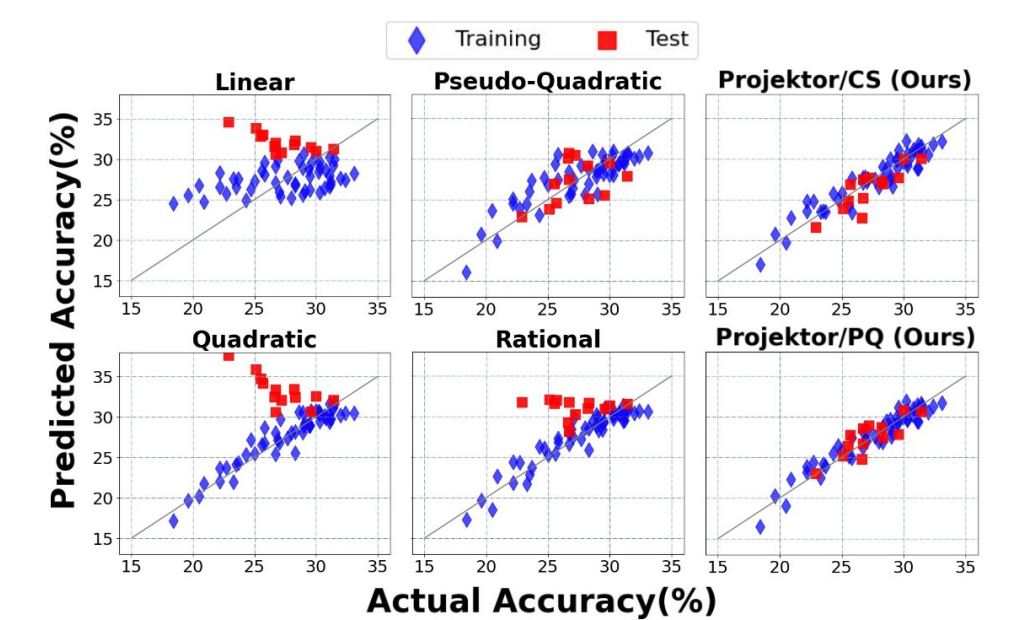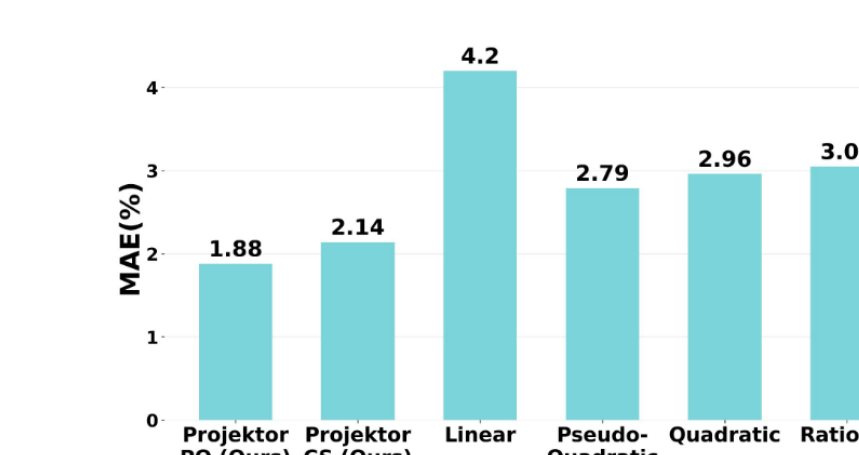Stage 3: Selection

Performance at fixed p

$$\mathbf{p}^* = \arg\max_{\mathbf{p}} \hat{\mathcal{L}}(\mathcal{A}(\mathcal{D}(N_s, \mathbf{p})), D^{\mathrm{val}})$$

Optimization is solved via gradient descent

$$\mathbf{p}^{t+1} \leftarrow \mathbf{p}^t + d^t \cdot \frac{\partial \hat{\mathcal{L}}(\mathcal{A}(\mathcal{D}(N_s, \mathbf{p})), D^{\mathrm{val}})}{\partial \mathbf{p}}\Big|_{\mathbf{p}=\mathbf{p}^t}$$

No N0 N1 · N · Size

## Applications

Evaluation: Performance prediction at unseen mixing ratio (extrapolation on p)

| | Data Source 1 | Data Source 2 | Data Source 3 | Model Performance |
|---|---|---|---|---|
| Projektor/PQ (Ours) | 34 | 33 | 33 | 61% |
| Projektor/CS (Ours) | 35 | 32 | 33 | 60% |
| Linear | | | 100 | 46% |
| Pseudo-Quadratic | 35 | 42 | 23 | 57% |
| Quadratic | 28 | 36 | 36 | 58% |
| Rational | 29 | 37 | 35 | 58% |
| LOO | | 100 | | 46% |
| Shapley | | 100 | | 47% |
| Random | ? | ? | ? | 52% |

Data mixture ratio selection for ImageNet-100 selection for 50K budget from 10K samples.

Projected performance for the selected mixture ratio (from above) and the actual performance.

Performance projection from 1K samples to larger data scales (2-10K).