

Data Selection via Optimal Control for Language Models

Yuxian Gu^{1,2}, Li Dong², Hongning Wang¹, Yaru Hao²,
Qingxiu Dong², Minlie Huang¹, Furu Wei²

¹The CoAI Group, Tsinghua University

²Microsoft Research



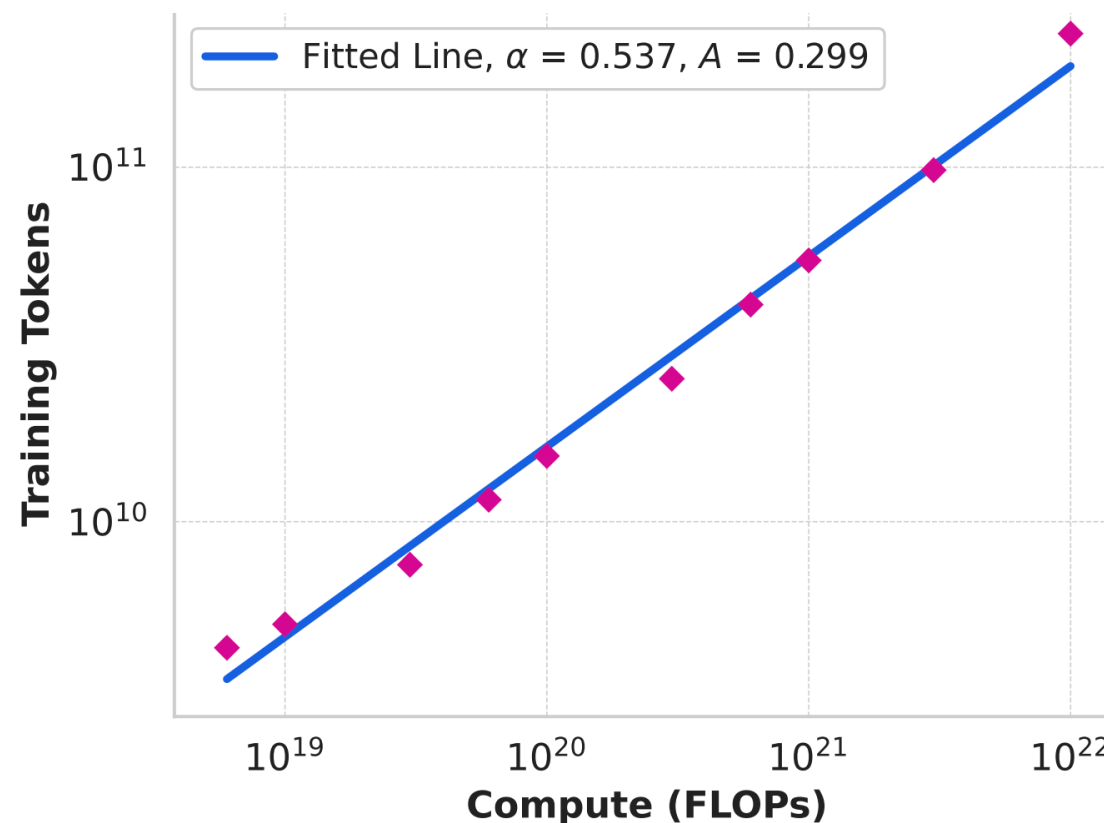
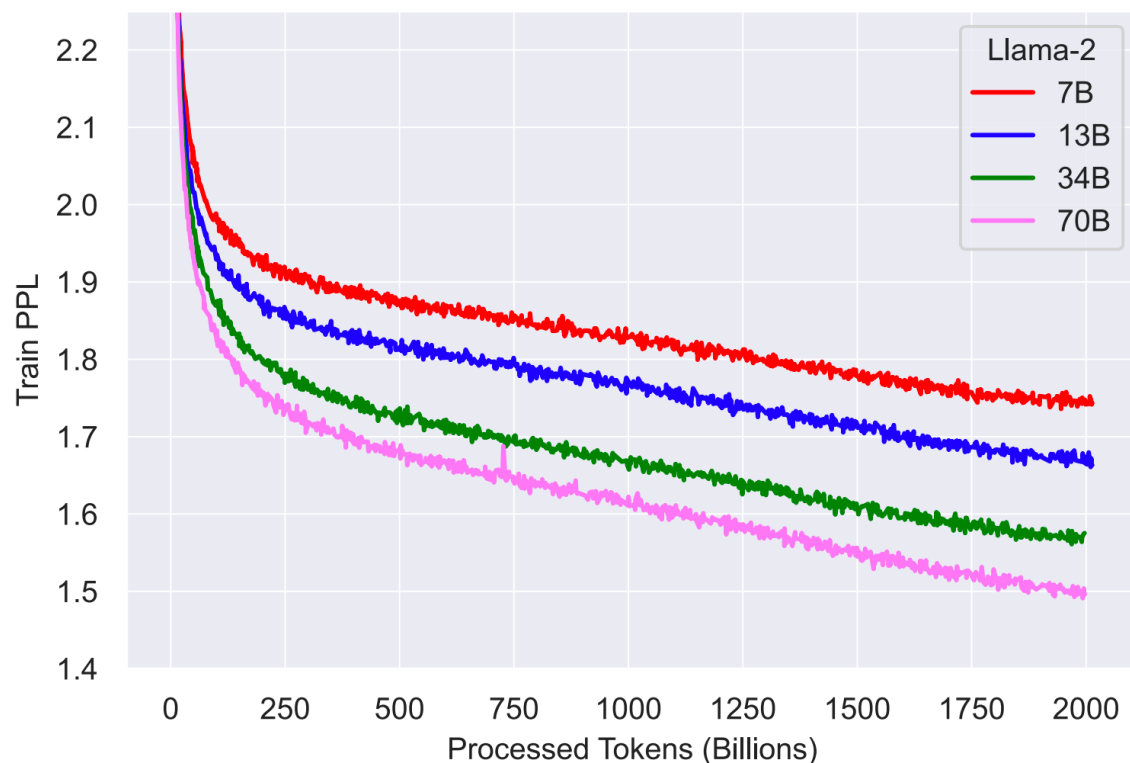
清华大学
Tsinghua University



Data challenges for pre-training LMs



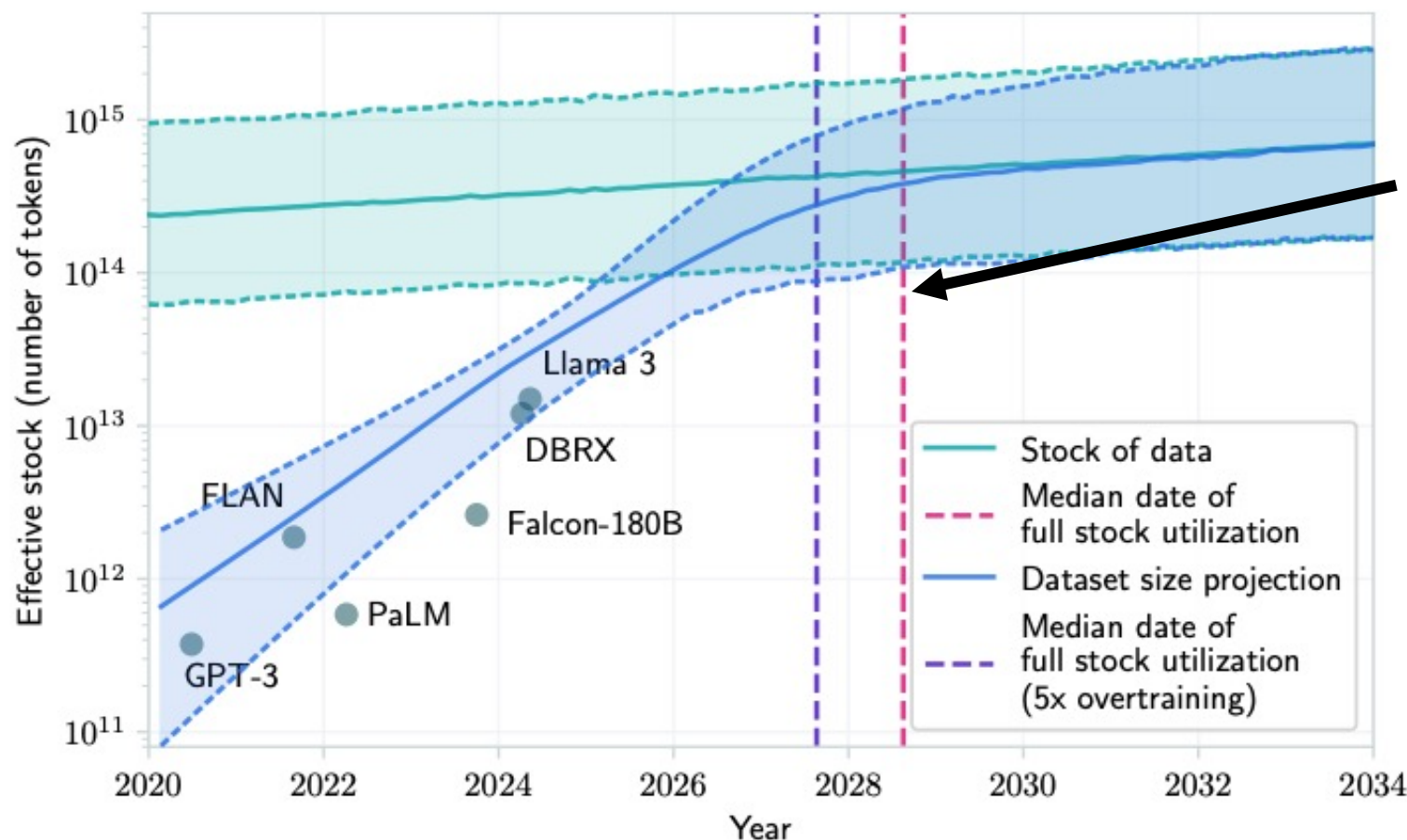
- Large amount of data makes pre-training quite **inefficient**.



Data challenges for pre-training LMs



- High-quality pre-training data is **running out**.

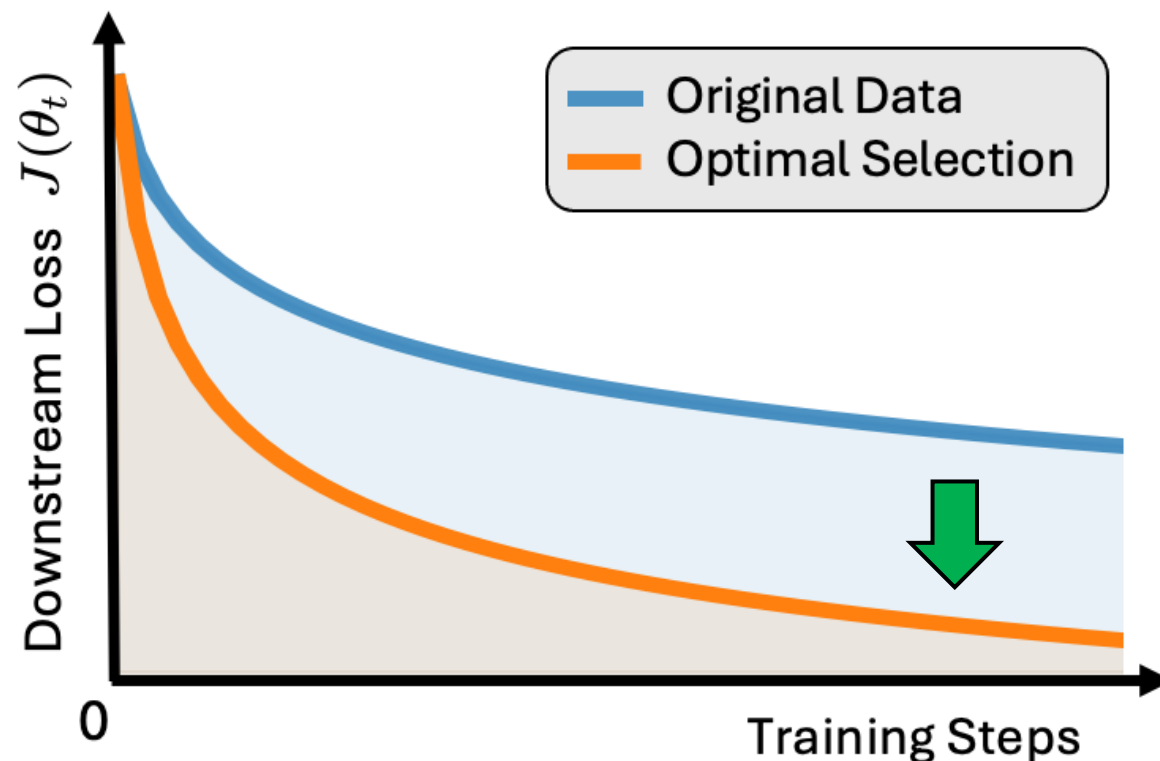
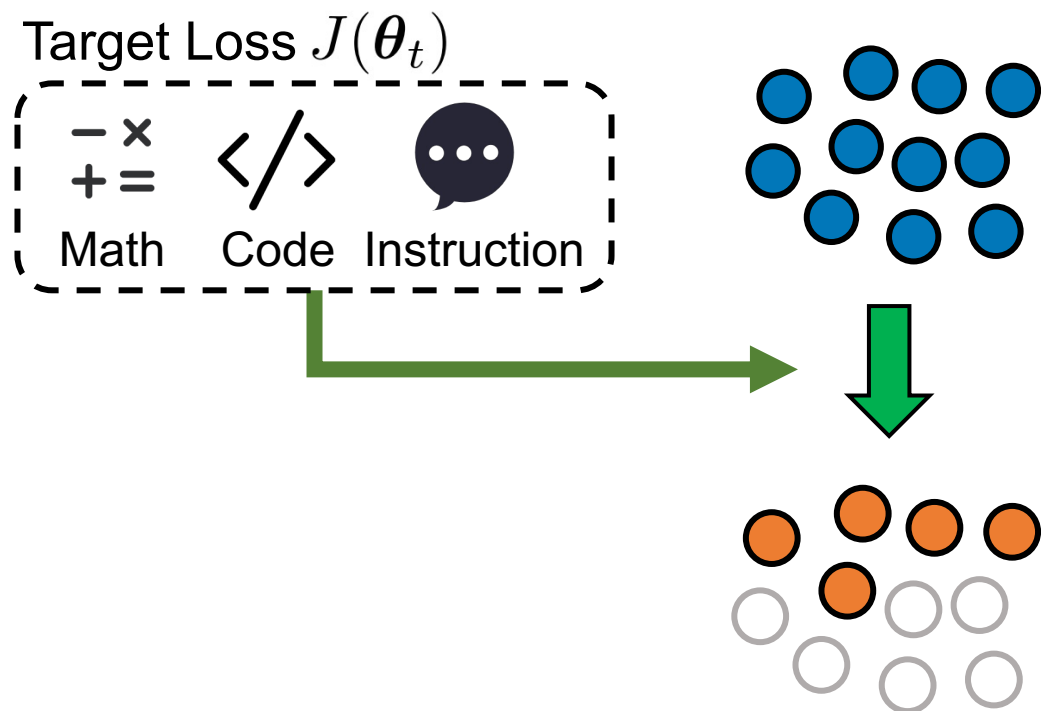


Models consume faster than humans produce.

Possible Solution: Data Selection?



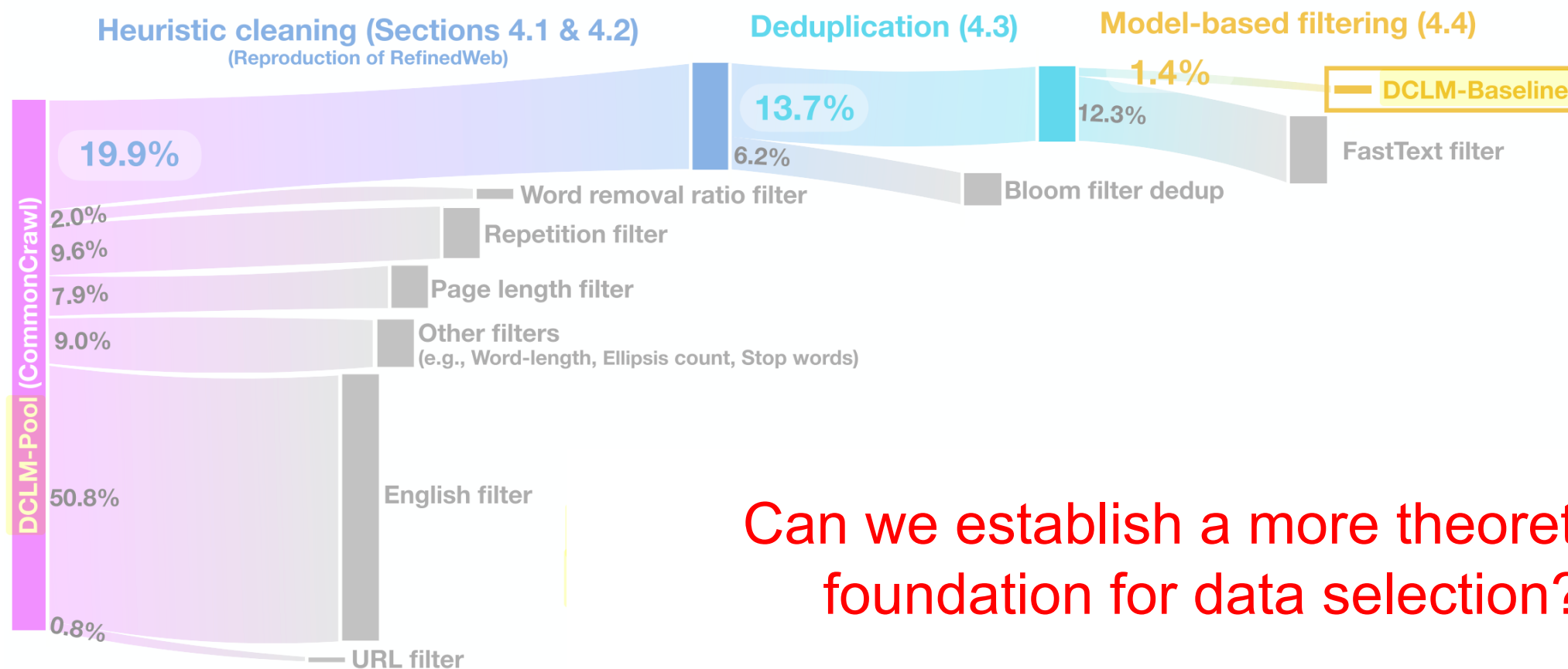
- ◉ Select a Pre-Training Corpus Subset for Better Target Performance
 - ◆ Target: Math, Code, High-Quality Instruction, etc..



Challenges for Data Selection



- Current data selection/filtering is **heuristic-based** and **tricky** task.



Can we establish a more theoretical foundation for data selection?

Overview



◉ Data challenges for pre-training LMs

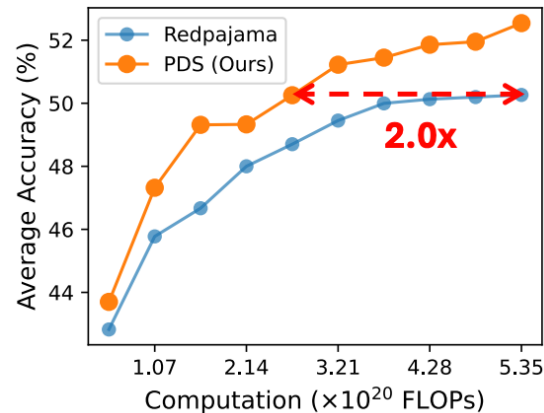
- ◆ Large amount of data makes pre-training quite **inefficient**.
- ◆ High-quality pre-training data is **running out**.
- ◆ Data selection/cleaning is a **heuristic-based** tricky task.

◉ **PDS**: Data Selection via Optimal Control for Pre-Training

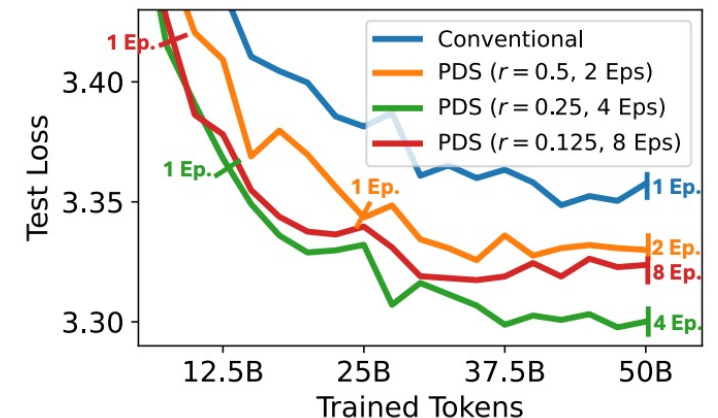
Theorem 2.1: PMP Conditions for Data Selection

$$\begin{aligned}\theta_{t+1}^* &= \theta_t^* - \eta \nabla L(\theta_t^*, \gamma^*), \quad \theta_0^* = \theta_0, \\ \lambda_t^* &= \lambda_{t+1}^* + \nabla J(\theta_t^*) - \eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*, \\ \gamma^* &= \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \lambda_{t+1}^{*\top} \nabla l(x_n, \theta_t^*) \right]\end{aligned}$$

Good theoretical foundation



2x acceleration

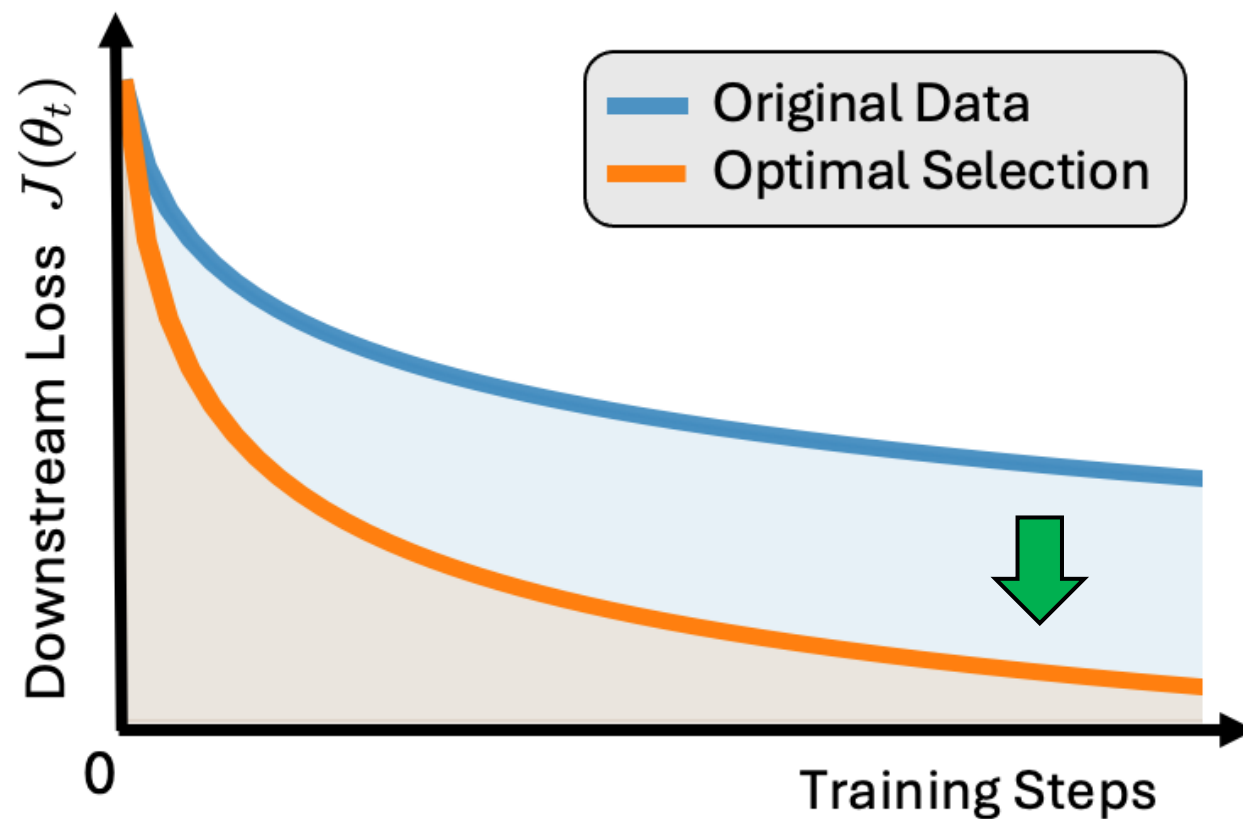
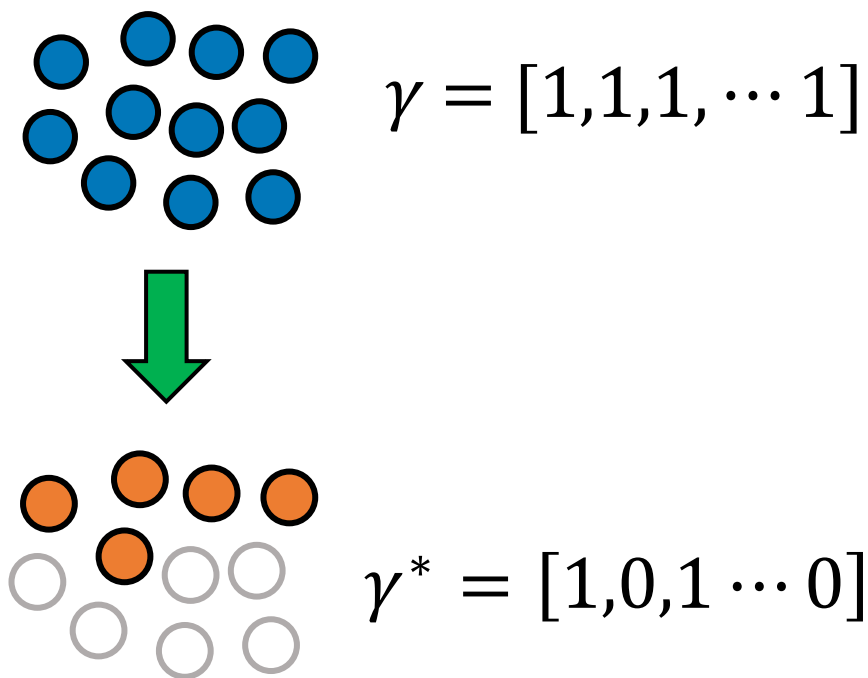


Improvement on limited data

Formulate Data Selection

- Optimize the data selection strategy for lower downstream loss

γ : indicates a sample is selected or not



Training on the Selected Data



- Loss: Treat the γ as the weights of the instance losses:

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \sum_{n=1}^{|\mathcal{D}|} \gamma_n l(x_n, \boldsymbol{\theta})$$

$\nearrow \gamma_n = 1: l(x_n, \theta)$ is selected

$\searrow \gamma_n = 0: l(x_n, \theta)$ is ignored

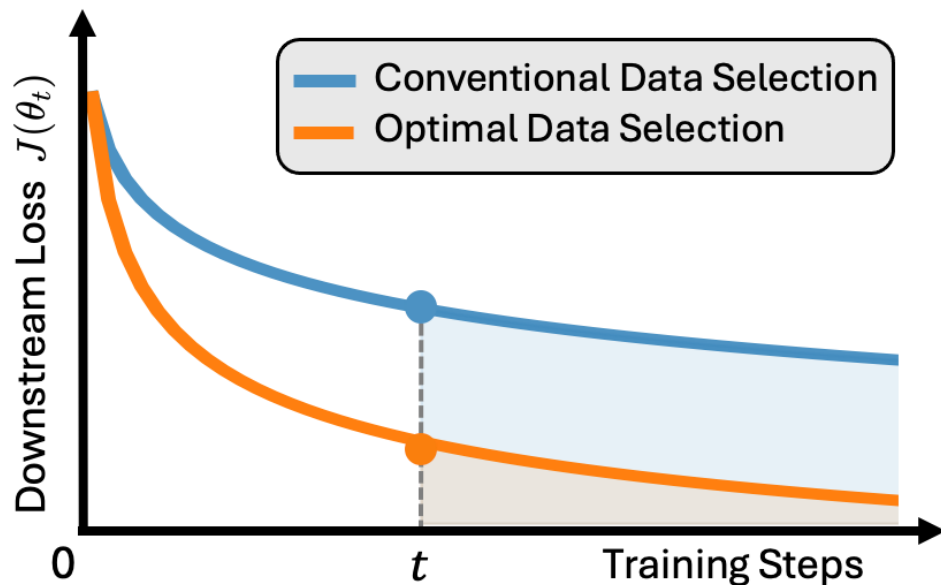
- The model is trained with:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla L(\boldsymbol{\theta}_t, \boldsymbol{\gamma})$$

Target to Optimize



- ◉ $J(\theta_t)$: downstream loss to minimize (like on math, code, etc.)
- ◉ Minimizing the Area Under the Loss curve (AUC)
 - ◆ AUC is directly related to the Scaling Law constants (see Appendix A in our paper)
 - ◆ Optimizing AUC is improving the Scaling Law!

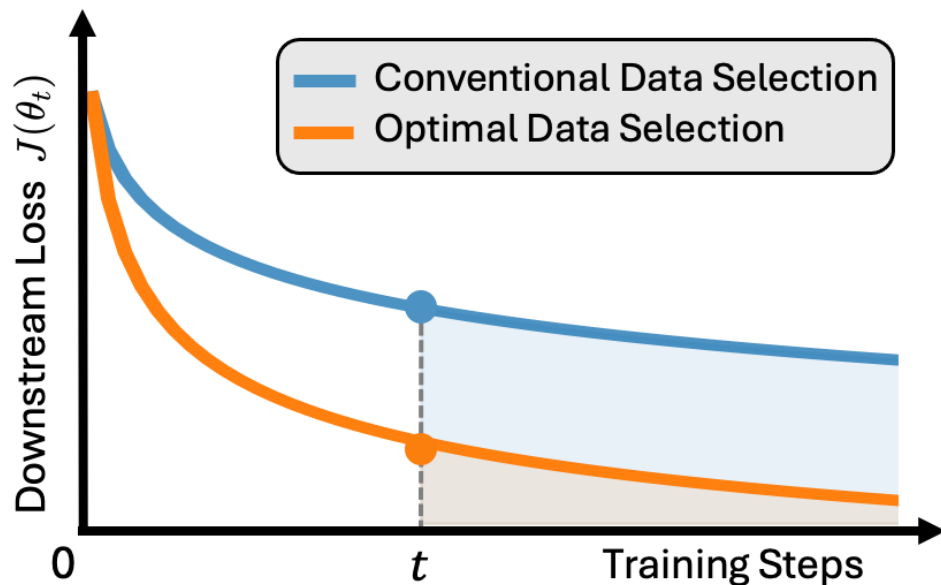


$$\min_{\gamma} \sum_{t=1}^T J(\theta_t),$$
$$\text{s.t. } \theta_{t+1} = \theta_t - \eta \nabla L(\theta_t, \gamma)$$

Target to Optimize



- ◉ $J(\theta_t)$: downstream loss to minimize (like on math, code, etc.)
- ◉ Minimizing the Area Under the Loss curve (AUC)
 - ◆ AUC is directly related to the Scaling Law constants (see Appendix A in our paper)
 - ◆ Optimizing AUC is improving the Scaling Law!



$$\min_{\gamma} \sum_{t=1}^T J(\theta_t),$$

$$\text{s.t. } \theta_{t+1} = \theta_t - \eta \nabla L(\theta_t, \gamma)$$

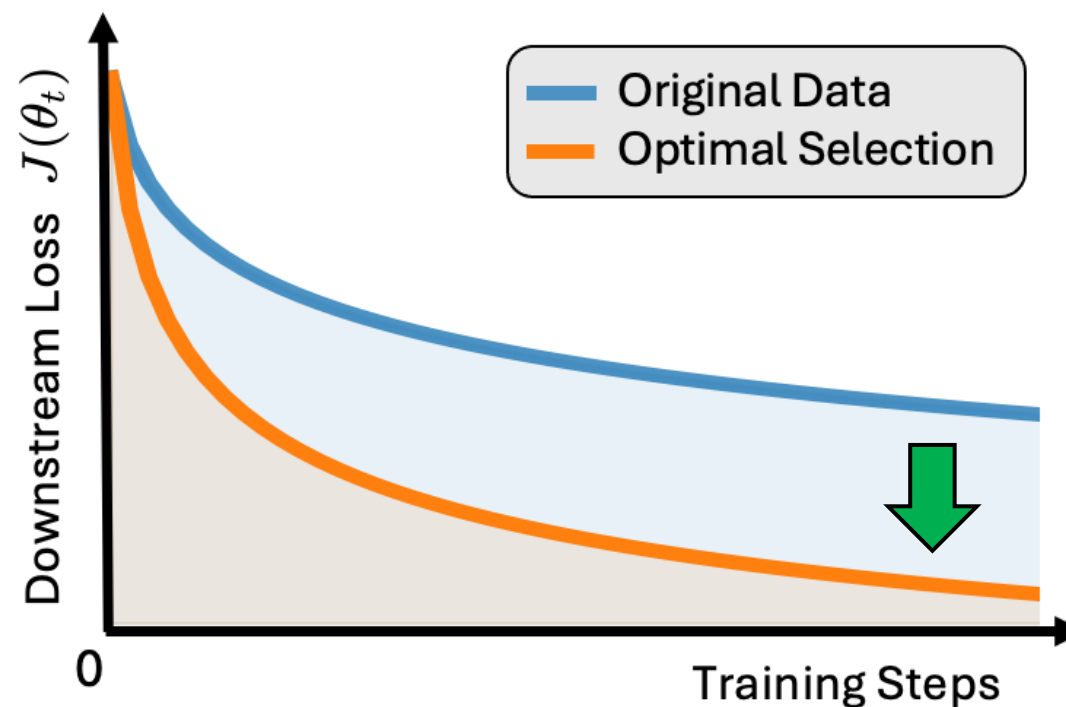
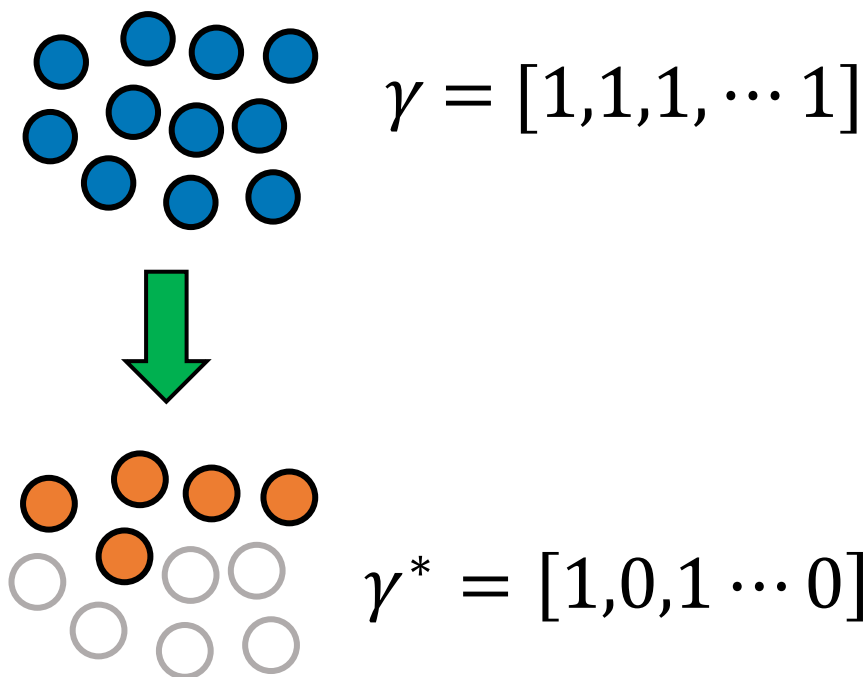
Hard to solve? Optimal Control!

Data Selection as a Control Problem



- Original Problem: Optimize the data selection strategy

γ : indicates a sample is selected or not



Data Selection as a Control Problem



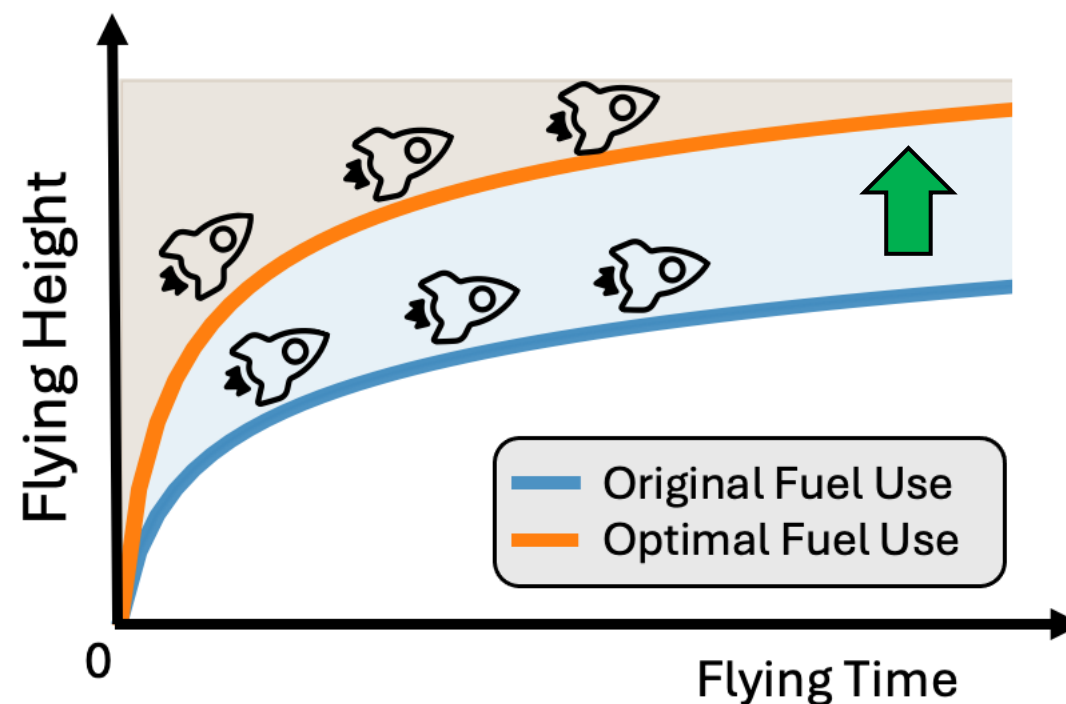
- Analogy: Optimizing fuel use when flying a rocket
 - ◆ Data is the “fuel” in pre-training language models



Original Fuel Use



Optimized Fuel Use

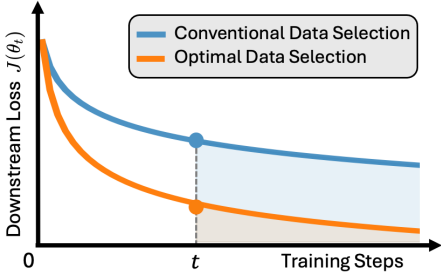



Mathematical Equivalence to Optimal Control



Data Selection for LMs

Fuel Use Optimization

		
Control Variable	Selection Strategy: γ	Fuel Consumption: u_t
Objective	Minimal AUC: $\min_{\gamma} \sum_{t=1}^T J(\theta_t),$	Maximize Distance: $\max_{u_t} x = \sum_{t=0}^T v_t \Delta t$
Constraints	Regularity: $\gamma \in U.$	Constant Total Fuel: $\sum_{t=0}^T u_t = U$
Dynamics	Gradient decent: $\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t, \gamma)$	Newton's Law: $\frac{v_{t+\Delta t} - v_t}{\Delta t} = -mg + F(u_t)$

Solving the Problem

- Pontryagin's **Maximum Principle** (PMP)

- ◆ Gives necessary conditions for the optimality of the problem

$$\begin{aligned} \min_{\gamma} \quad & \sum_{t=1}^T J(\boldsymbol{\theta}_t), \\ \text{s.t.} \quad & \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla L(\boldsymbol{\theta}_t, \gamma) \\ & L(\boldsymbol{\theta}, \gamma) = \sum_{n=1}^{|\mathcal{D}|} \gamma_n l(x_n, \boldsymbol{\theta}) \end{aligned}$$



Lev Pontryagin, 1908 - 1988

Conditions for Optimal Data Selection



Theorem 2.1 (PMP Conditions for Data Selection).

$$\left\{ \begin{array}{l} \text{\#1 } \boldsymbol{\theta}_{t+1}^* = \boldsymbol{\theta}_t^* - \eta \nabla L(\boldsymbol{\theta}_t^*, \gamma^*), \\ \text{\#2 } \boldsymbol{\lambda}_t^* = \boldsymbol{\lambda}_{t+1}^* + \nabla J(\boldsymbol{\theta}_t^*) - \eta \nabla^2 L(\boldsymbol{\theta}_t^*, \gamma^*) \boldsymbol{\lambda}_{t+1}^*, \\ \text{\#3 } \gamma^* = \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \boldsymbol{\lambda}_{t+1}^{*\top} \nabla l(x_n, \boldsymbol{\theta}_t^*) \right] \end{array} \right.$$

We can get the optimal data score γ^* here!

#1: Learning Condition

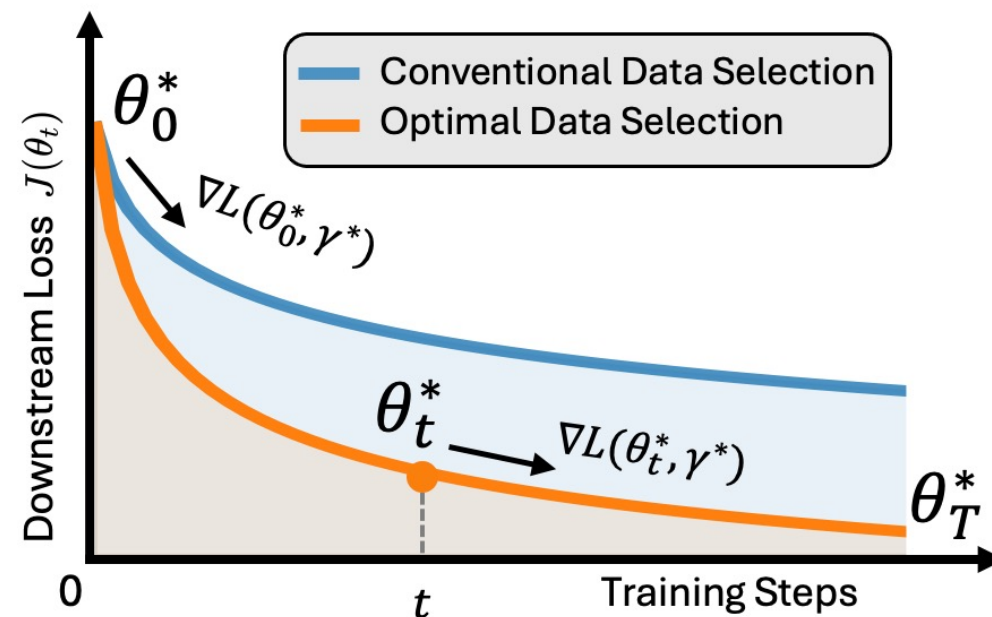


$$\theta_{t+1}^* = \theta_t^* - \eta \nabla L(\theta_t^*, \gamma^*)$$

Model Parameters

Optimal Data Selection Strategy

- Exactly the parameter updating policy of training LMs
- Constrains the θ_t^* to be reachable with GD under the optimal data selection



#2: Target Condition



$$\lambda_t^* = \lambda_{t+1}^* + \nabla J(\theta_t^*) - \eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*$$

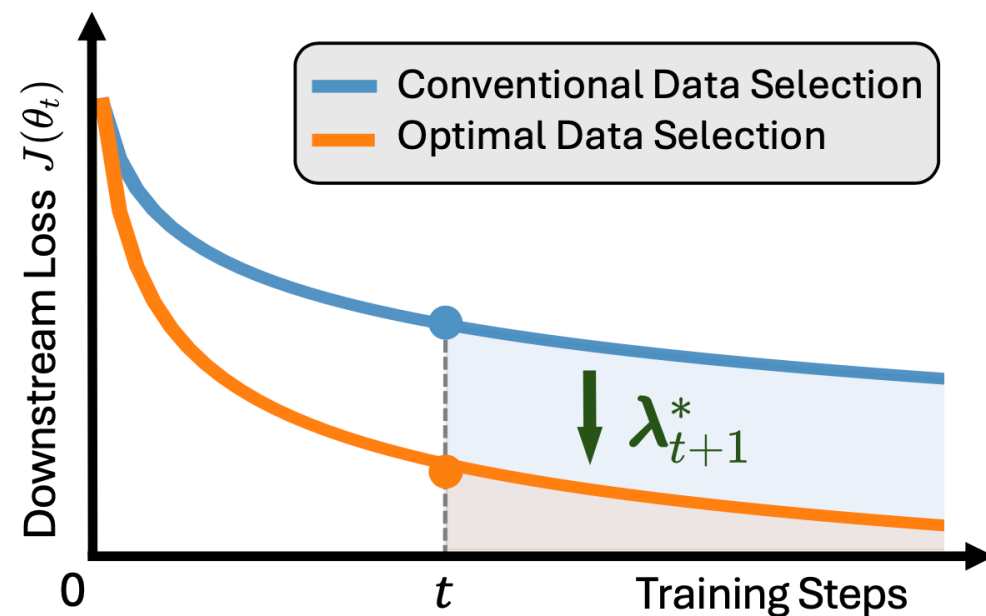
Ideal gradient = Target + Learning dynamics

- λ_t^* : the ideal gradient of the high-quality data points.



“Compass” for high-quality data.

- Ideal gradient includes information of Target Loss and Learning Dynamics.

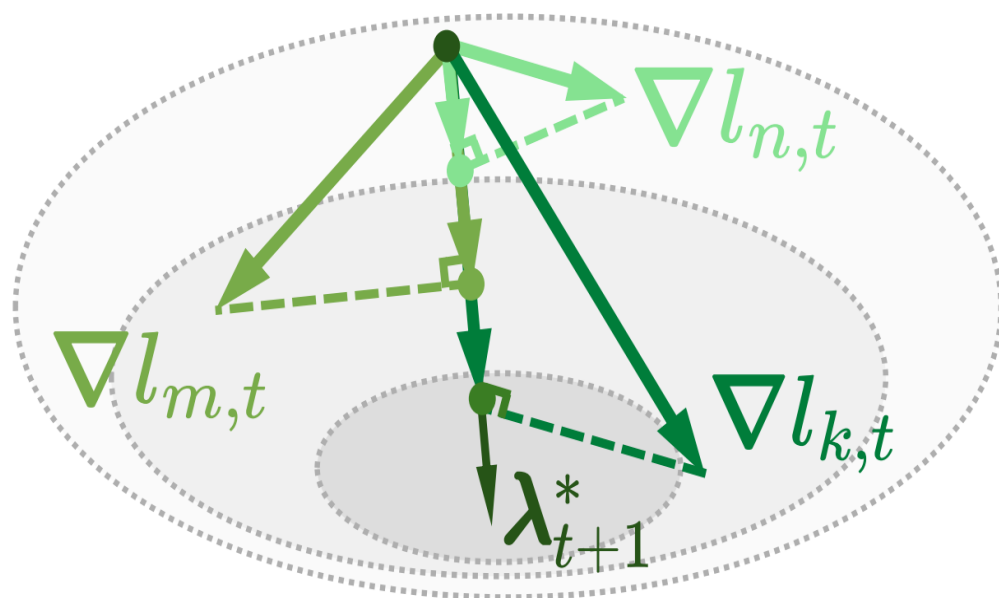


#3: Maximum Condition



$$\gamma^* = \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \lambda_{t+1}^{* \top} \nabla l(x_n, \theta_t^*) \right]$$

Gradient of each sample



$$\sum_t \lambda_{t+1}^{* \top} \nabla l_{n,t} < \sum_t \lambda_{t+1}^{* \top} \nabla l_{m,t} < \sum_t \lambda_{t+1}^{* \top} \nabla l_{k,t}$$

$$\Rightarrow \gamma_n < \gamma_m < \gamma_k$$

Examples with **closer gradients to λ_t**
should **have higher γ** .

Summing up



Theorem 2.1 (PMP Conditions for Data Selection).

#1 Learning Condition $\theta_{t+1}^* = \theta_t^* - \eta \nabla L(\theta_t^*, \gamma^*)$

#2 Target Condition $\lambda_t^* = \lambda_{t+1}^* + \nabla J(\theta_t^*) - \eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*$

#3 Maximum Condition $\gamma^* = \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \lambda_{t+1}^{*\top} \nabla l(x_n, \theta_t^*) \right]$

How to Solve (Theoretically)?



Use learning condition to forward compute θ_0 to θ_T

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

#2 Target Condition $\lambda_t^* = \lambda_{t+1}^* + \nabla J(\theta_t^*) - \eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*$

#3 Maximum Condition $\gamma^* = \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \lambda_{t+1}^{* \top} \nabla l(x_n, \theta_t^*) \right]$

How to Solve (Theoretically)?



Use Target Condition to reverse compute λ_T to λ_0

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

reverse Pass with #2: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

#3 Maximum Condition $\gamma^* = \arg \max_{\gamma} \sum_{n=1}^{|\mathcal{D}|} \gamma_n \left[\sum_{t=0}^{T-1} \lambda_{t+1}^{* \top} \nabla l(x_n, \theta_t^*) \right]$

How to Solve (Theoretically)?

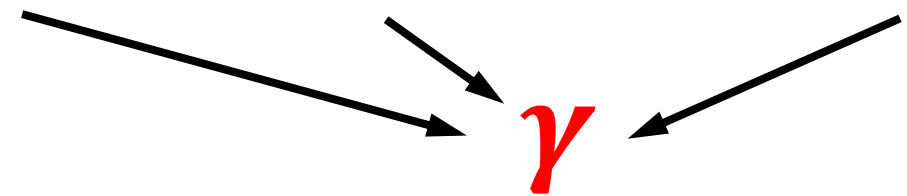


Use Maximum Condition to solve the final γ

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

reverse Pass with #2: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

maximum γ with #3:



How to Solve (Theoretically)?



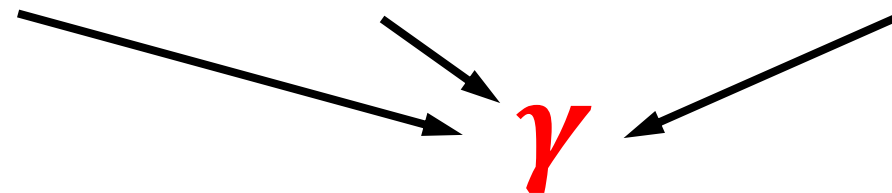
- Iteratively Solving γ until convergence (Algorithm 1)

while γ not converged, do

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

reverse Pass with #2: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

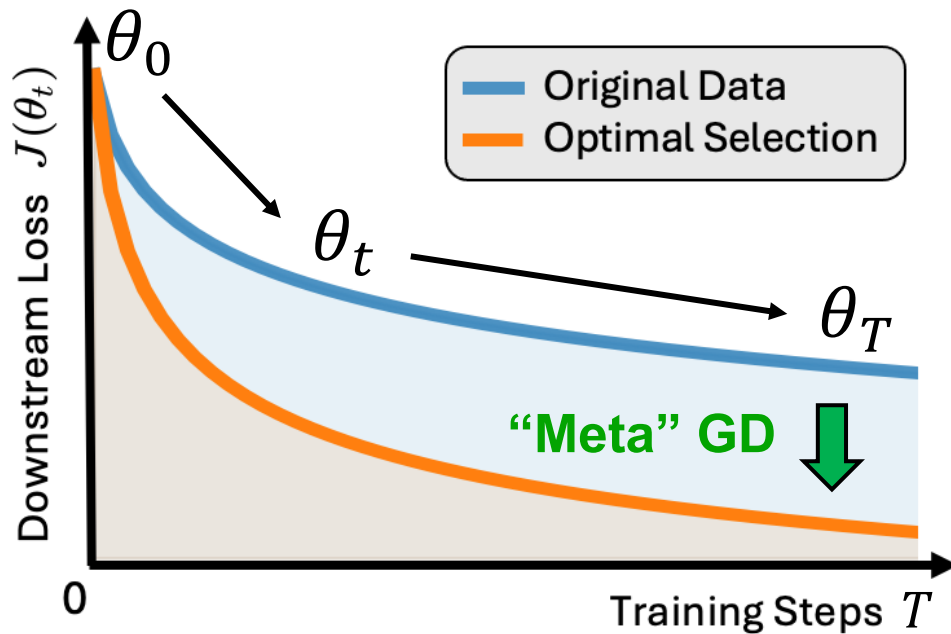
maximum γ with #3:



Equivalence to “Meta” Gradient Decent



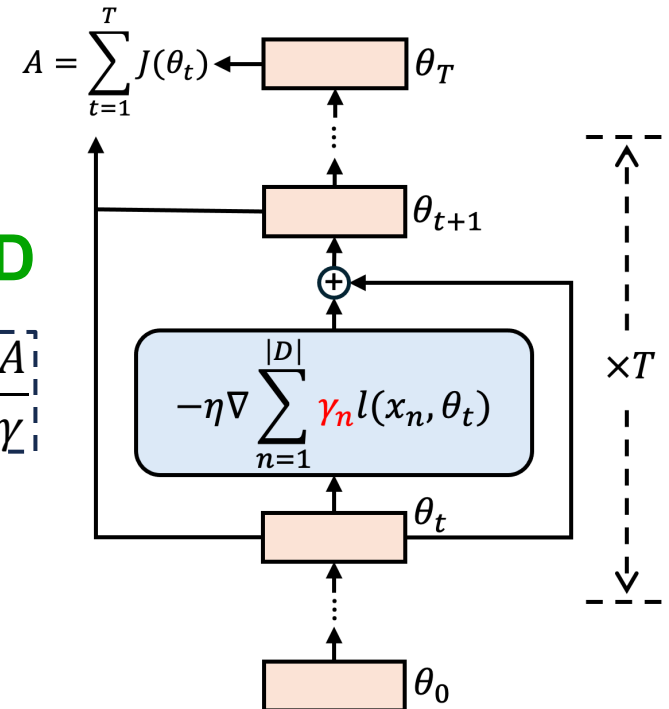
- Optimizing the whole training process with “Meta GD”
- A training process can be viewed as an NN (vertically)



Equivalent

“Meta” GD

$$\gamma \leftarrow \gamma - \alpha \frac{\partial A}{\partial \gamma}$$



Solve γ with “Meta” GD



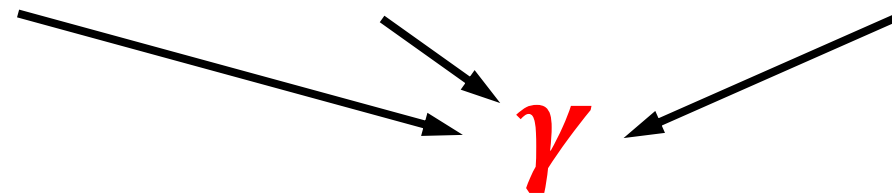
- Iteratively Solving γ until convergence (Algorithm 1)

while γ not converged, do

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

reverse Pass with #2: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

maximum γ with #3:



Solve γ with “Meta” GD



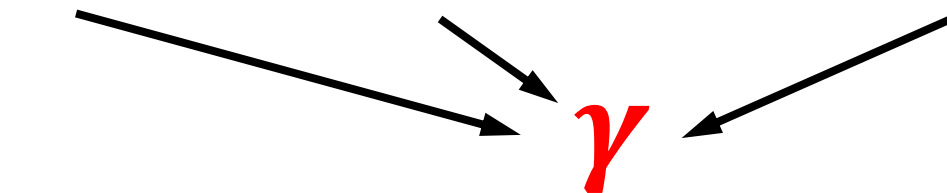
- Iteratively Solving γ until convergence (Algorithm 1)

while γ not converged, do

Forward pass of “Meta” GD: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

Backward pass of “Meta” GD: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

Step pass of “Meta” GD:



How to Solve (Theoretically)?



Iteratively Solving γ (Algorithm 1)

while γ not converged, do

forward Pass with #1: $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$

reverse Pass with #2: $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$

maximum γ with #3:

The diagram illustrates the iterative solving of γ . It shows a forward pass (blue) and a reverse pass (green) within a dashed blue box. The forward pass is $\theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T$ and the reverse pass is $\lambda_0 \leftarrow \lambda_1 \leftarrow \dots \leftarrow \lambda_T$. Arrows from λ_0 , λ_1 , and λ_T point to a red γ . A red cloud contains the equation $\eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*$, with a red circle around θ_1 and a red arrow pointing from the cloud to the γ .

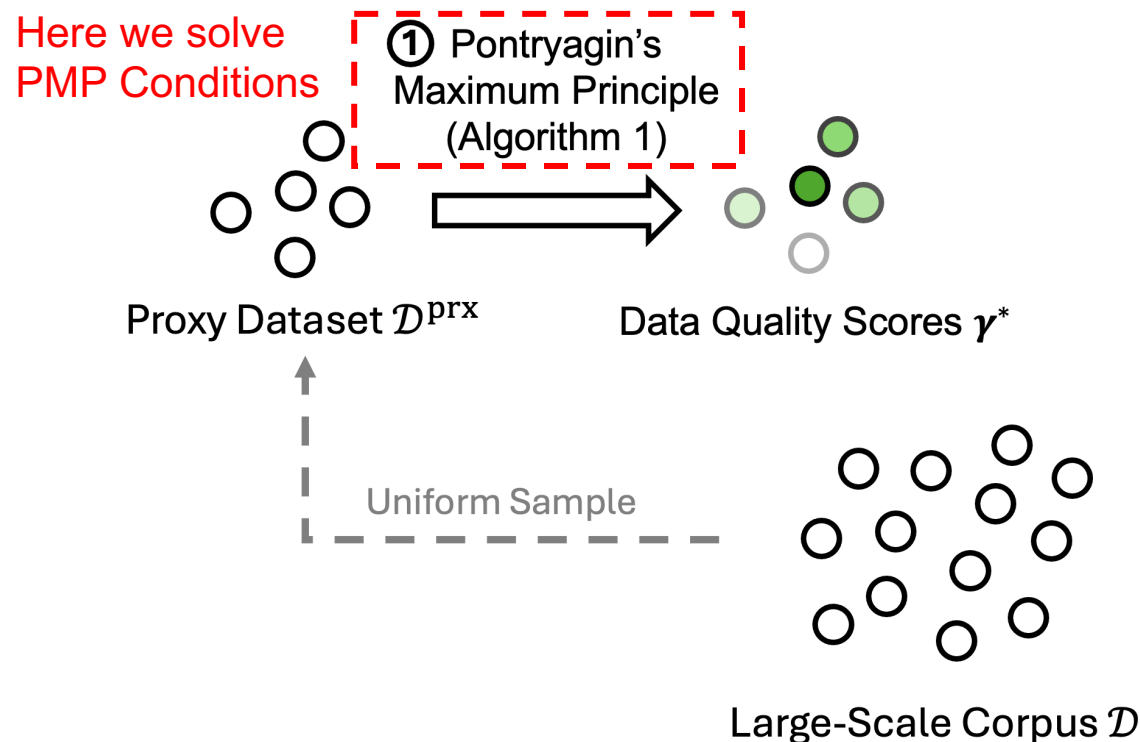
$$\eta \nabla^2 L(\theta_t^*, \gamma^*) \lambda_{t+1}^*$$

Forward and Reverse passes are computationally intensive!

Efficient Implementation: Proxy to Large



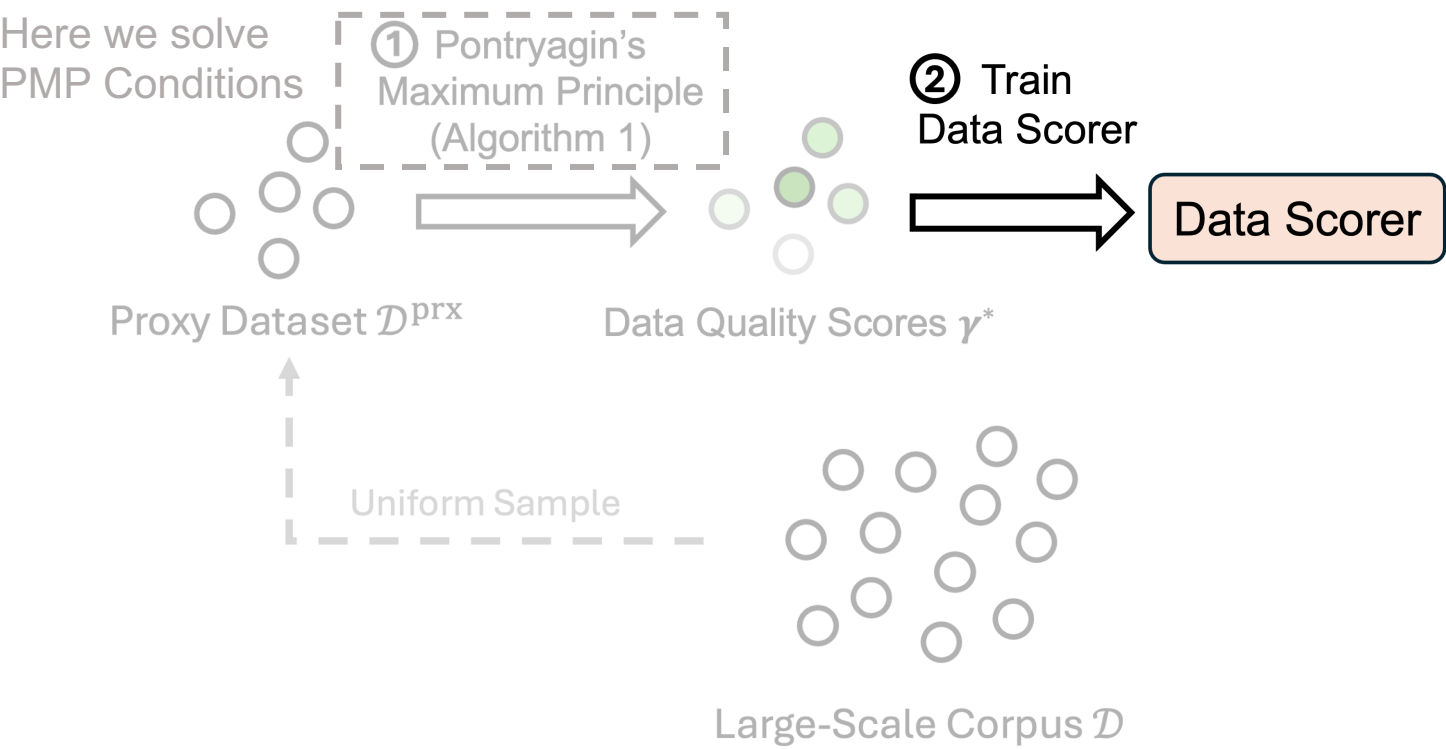
① Solve γ on a small model (e.g., 140M) and data (e.g., 160M tokens)



Efficient Implementation: Proxy to Large



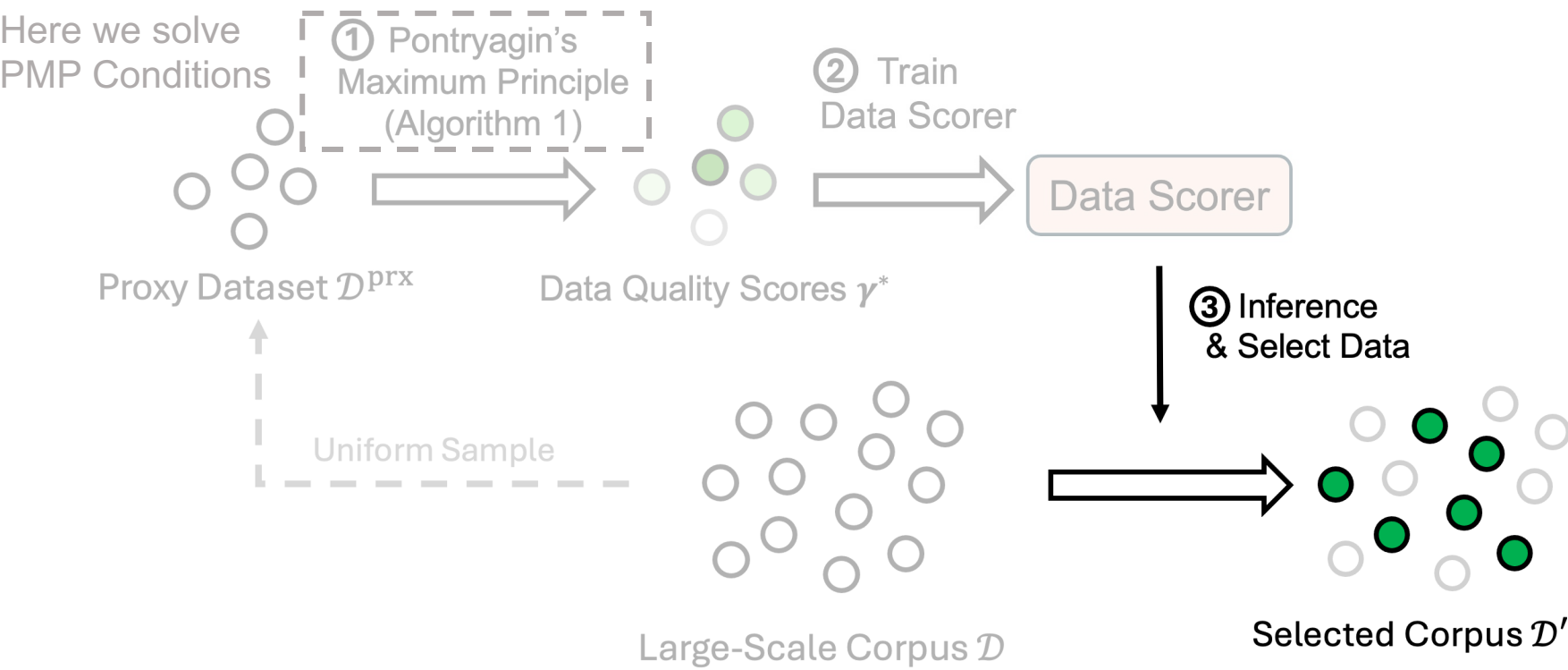
- ① Solve γ on a small model (e.g., 140M) and data (e.g., 160M tokens)
- ② Fit γ with a data scorer (e.g., a 140M LM with a regression head)



Efficient Implementation: Proxy to Large



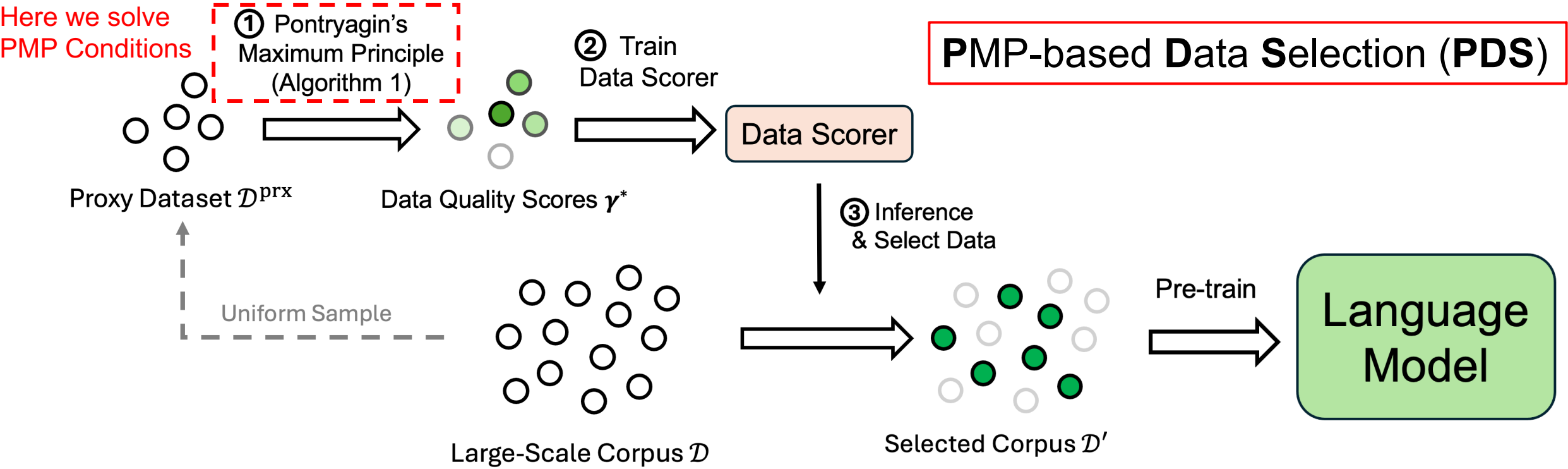
- ① Solve γ on a small model (e.g., 140M) and data (e.g., 160M tokens)
- ② Fit γ with a data scorer (e.g., a 140M LM with a regression head)
- ③ Infer γ on the whole dataset (e.g., 100B tokens)



Efficient Implementation: Proxy to Large



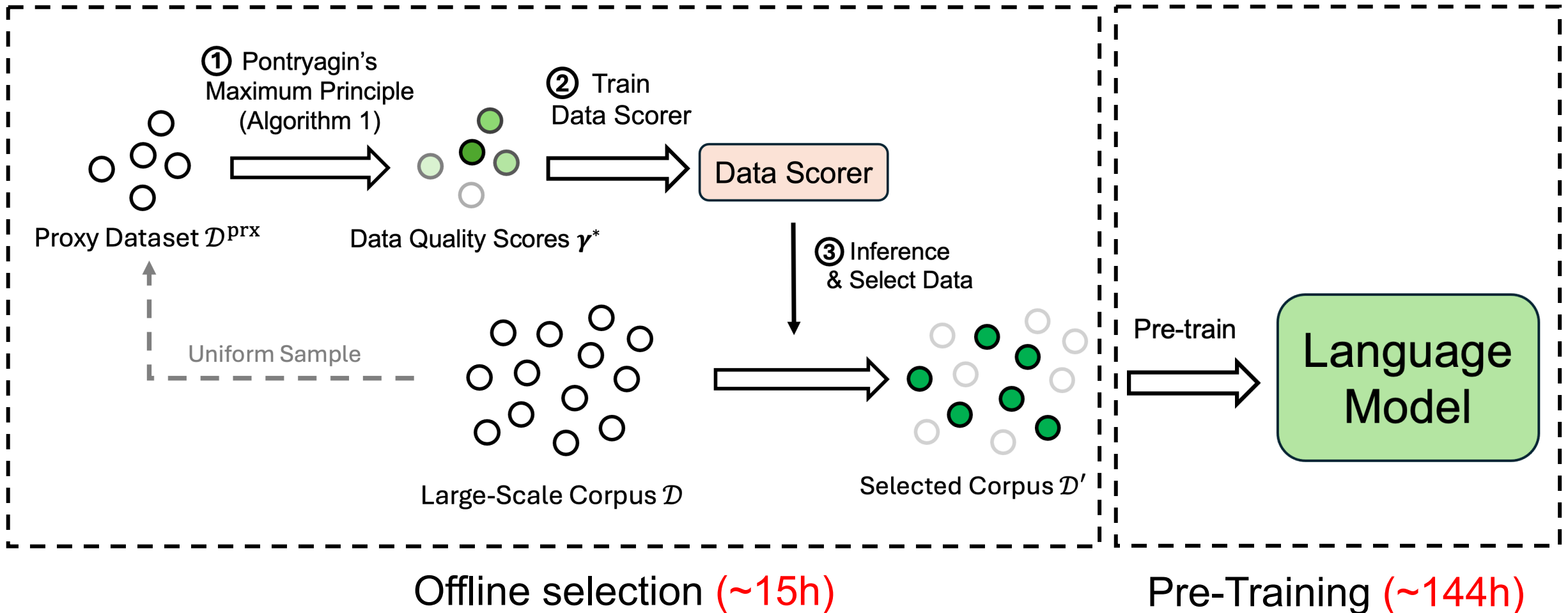
- ① Solve γ on a small model (e.g., 140M) and data (e.g., 160M tokens)
- ② Fit γ with a data scorer (e.g., a 140M LM with a regression head)
- ③ Infer γ on the whole dataset (e.g., 100B tokens)



Efficient Implementation: Proxy to Large



PMP-based Data Selection (PDS) is Efficient

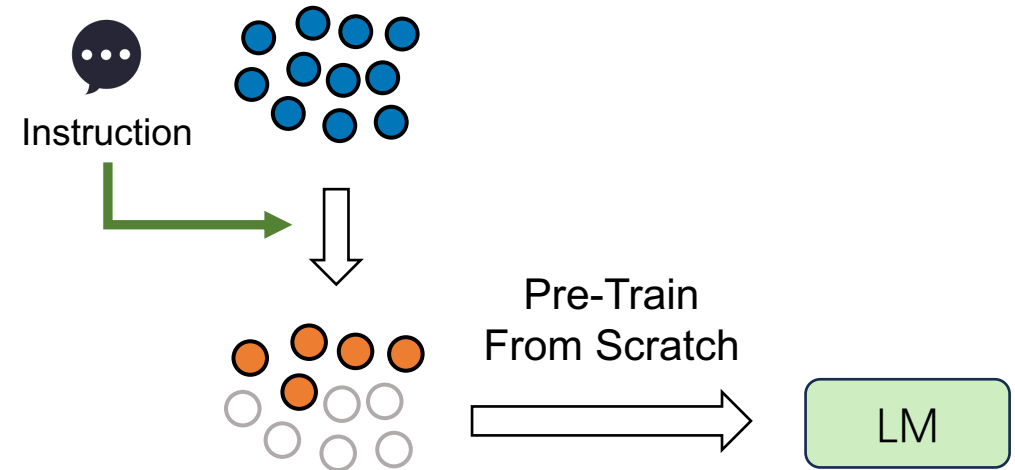


Experiment Setups



Training & Evaluation Setups

- ◆ Pre-training LMs from scratch
- ◆ Evaluate zero-shot performance



Data Setups

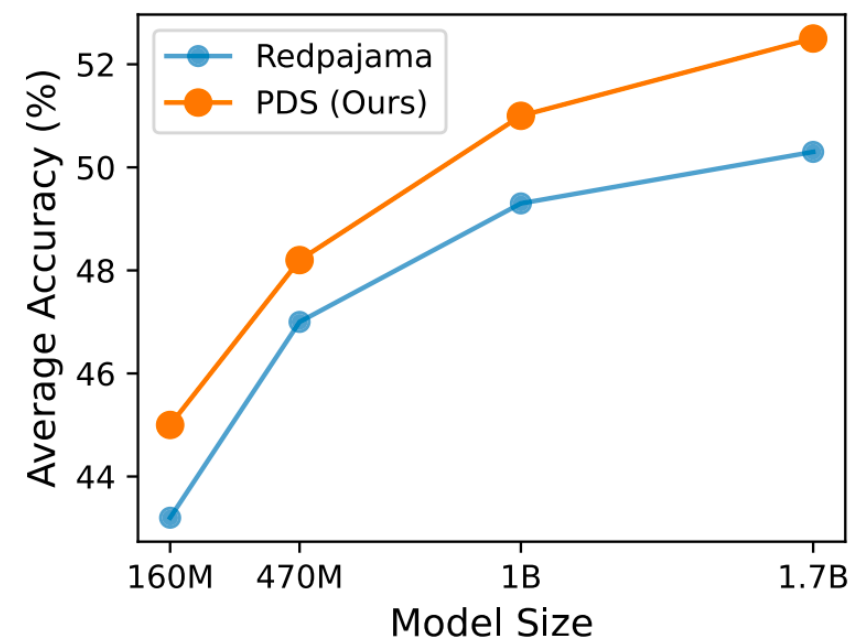
- ◆ Pre-Training data: CommonCrawl from Redpajama (100B tokens)
- ◆ Downstream loss $J(\theta)$: loss on LIMA (1k high-quality instruction-response pairs)
- ◆ Evaluation: 9 common NLP benchmarks: LAMBADA, Hellaswag, BoolQ, etc.

Performance Improvement



- ◉ Select 50B-token corpus from 125B-token corpus.
- ◉ Match the total training steps with the baselines (training computation)

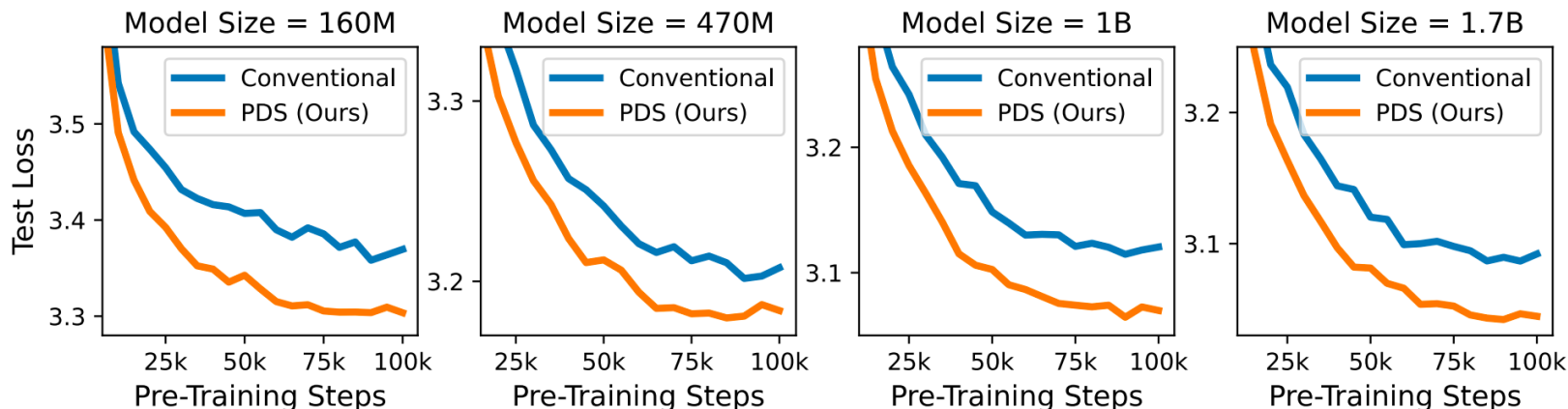
	HS	LAMB	Wino.	OBQA	ARC-e	ARC-c	PIQA	SciQ	BoolQ	Avg.
Model Size = 470M										
Conventional	36.7	41.4	52.4	30.4	44.8	25.2	61.0	70.6	60.4	47.0
RHO-Loss	36.6	42.4	53.0	29.4	43.7	25.2	60.4	72.8	59.8	47.0
DSIR	36.4	42.6	51.7	29.8	46.0	24.7	61.0	72.0	55.8	46.7
IF-Score	36.6	41.8	53.4	29.6	44.7	25.1	60.8	68.8	58.7	46.6
PDS	37.9	44.6	52.3	29.8	46.5	25.8	61.8	73.8	61.4	48.2
Model Size = 1B										
Conventional	39.9	47.6	52.4	30.6	49.3	26.4	63.1	73.7	60.9	49.3
RHO-Loss	39.8	47.0	53.0	30.8	48.0	26.4	62.9	71.1	61.0	48.9
DSIR	40.8	47.8	53.0	31.2	49.8	26.8	62.7	76.6	58.0	49.6
IF-Score	39.4	47.0	52.6	28.6	49.4	26.4	63.5	74.0	60.5	49.0
PDS	42.1	48.8	54.0	33.4	51.3	28.0	64.1	78.5	58.7	51.0



Extrapolation to 400B models on 15T tokens



○ Fitting the model loss curves with the Scaling Law



○ Extrapolating to 400B models on 10T tokens

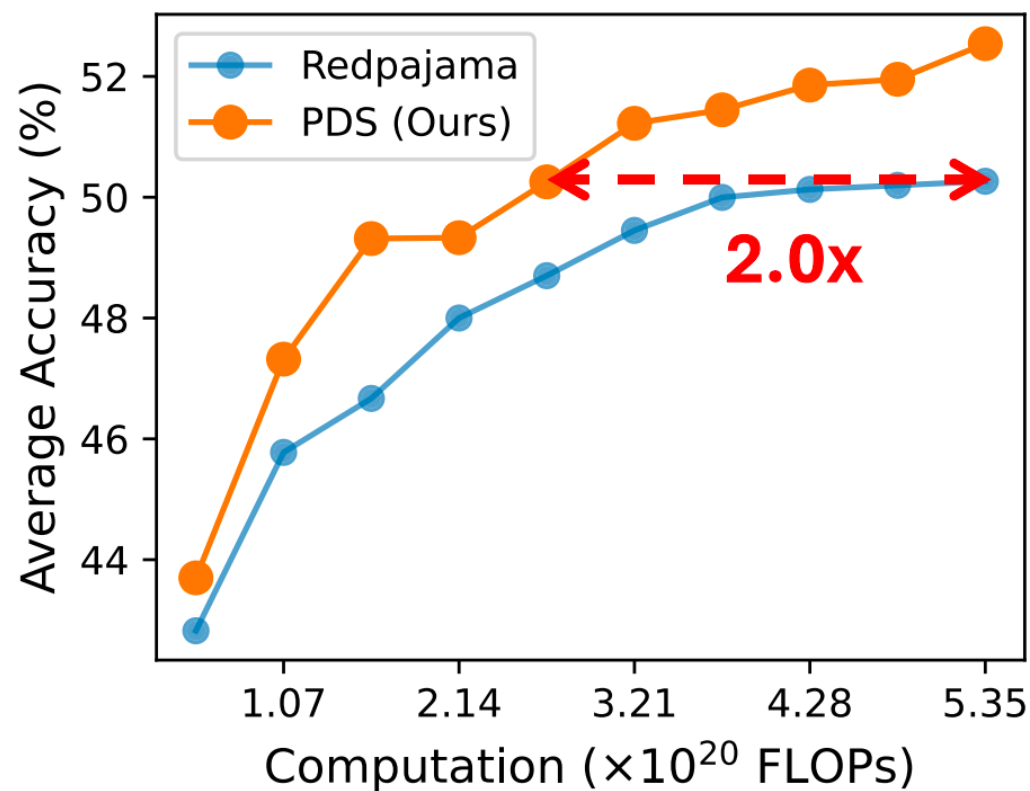
$$L(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta}$$

	N	D	Conventional	PDS
GPT-3 [12]	175B	300B	2.882	2.872
Llama [72]	6.7B	1.0T	2.942	2.896
Llama 2 [73]	70B	2.0T	2.877	2.855
Llama 3.1 [21]	405B	15T	2.851	2.838

Computation Saving



- 2.0x acceleration on 1.7B models
- PDS is efficient and offline
 - Select data once for all models



		FLOPs ($\times 10^{20}$)	Actual Time
PDS	Proxy γ -solver	0.49	15.2 Hours
	Data Scorer	0.063	1.50 Hours
	Data Selection	0.0	10.2 Minutes
Pre-Training		5.1	144 Hours

Data Utilization Improvement



Performance improvement with limit data (50B tokens)

1 Ep.

Pre-Training (w/o Data Selection)

1 Ep. 2 Ep.

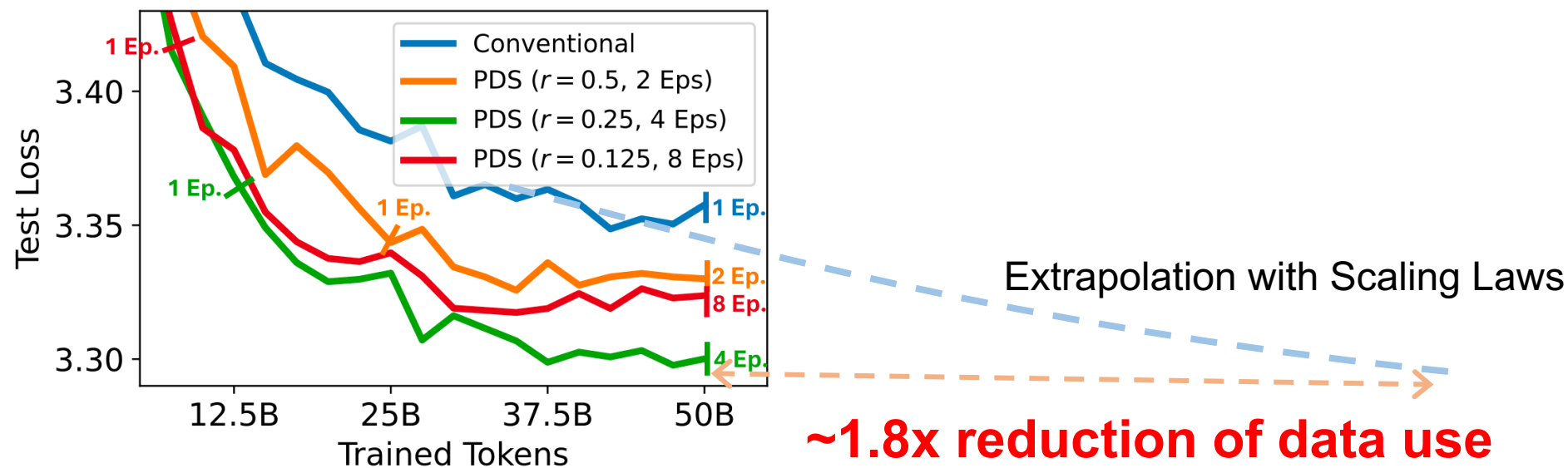
Select 50% data, train for 2 epochs

1 Ep. 2 Ep. 3 Ep. 4 Ep.

Select 25% data, train for 4 epochs

1 Ep. 2 Ep. 3 Ep. 4 Ep. 5 Ep. 6 Ep. 7 Ep. 8 Ep.

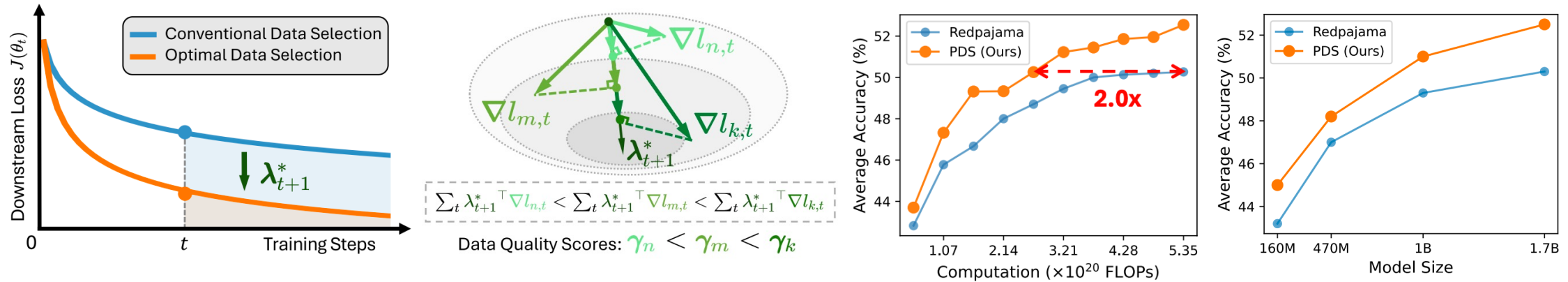
Select 12.5% data, train for 8 epochs



Conclusion



- ◉ A novel perspective for Data selection: Optimal Control problem



- ◆ Good theoretical guarantees ✓
- ◆ Efficient Implementation ✓
- ◆ Sound empirical results ✓

A **rigorous, theory-driven alternative** to the ad-hoc practices that currently dominate LM pre-training



Thanks!

 Paper: <https://arxiv.org/abs/2410.07064>

 GitHub: https://github.com/microsoft/LMOps/tree/main/data_selection

 HuggingFace: <https://huggingface.co/Data-Selection>

Paper:



Code:



HF:

