

Step 0: Warm up w/ random data



Next round

Validation Data

Training Data

Proxy Model

Score Model

Predicts data scores

LLM

Step 3: Pretrain

Select data

Step 4: Evaluate

Downstream task

Step 1: Bilevel optimization for score model and proxy model

Step 2: Select data based on the score ranking