

Bayesian Optimization for Mixture Optimization in Language Models?

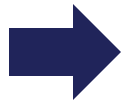
Tim Lindenau, Forschungspraxis
München, 16.10.2024



Motivation

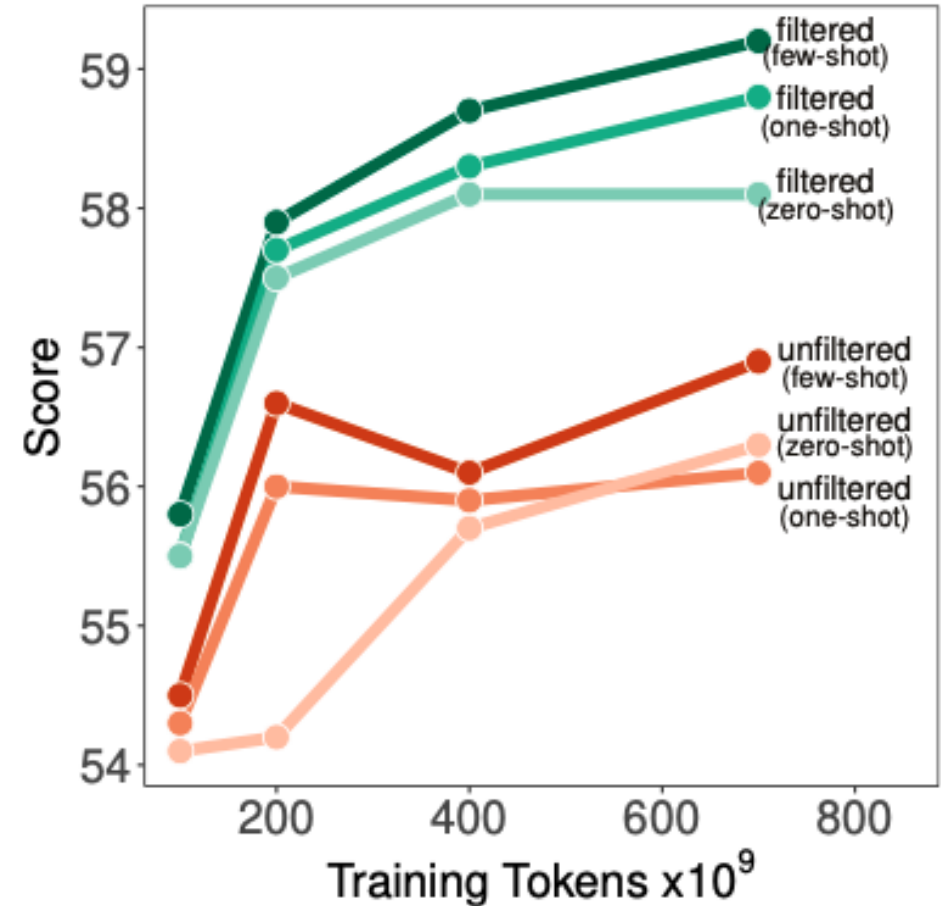
Motivation: The choice of data is crucial for LMs

- Massive **scaling of training data** has been key factor towards the current **LLM progress**
- **Unequal utility** of data. Even **simple filtering** can already achieve **big performance improvements**
- **Existing data curation** methods unprincipled, mostly around **heuristics** and **human intuition** about data utility



Principled data-curation required to unlock **faster training** and **better model performance**

Data filtering improves LM performance



Our Setting: Data curation by mixing different domains

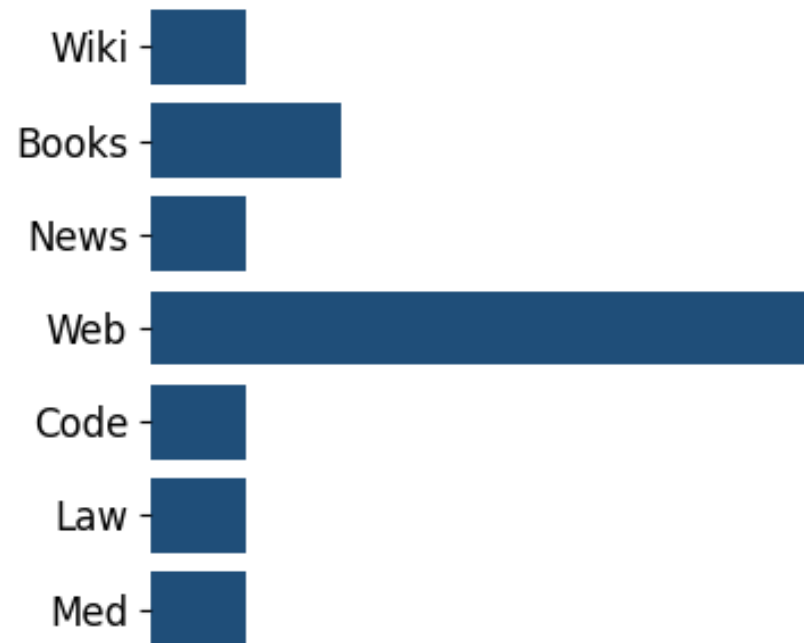
Problem Setup – Data Mixing:

- Training data from **k domains** (Wikipedia, Web, Code, etc.)
- For a **fixed budget**, how much of **each domain** to train on

Biggest Challenge – Runtime Overhead:

- **No** significant **compute overhead** in pre-training
- Mixture optimization on **small proxy model**, assuming scale invariance
- Focus on **sample efficiency** in training number of proxy models

Mixing proportions of different datasets for pre-training



Existing Mixture Optimization Approaches and their Pros & Cons

Name	Method	Pros & Cons
Online Mixing	<ul style="list-style-type: none">• Data Mixing while training the production grade model• Multi-Armed Bandit to model which domain to sample next token• Domain perplexity as reward function	<ul style="list-style-type: none">+ Negligible run-time overhead without pre-training proxy models- Missing interpretability into overall optimal mixture
DoReMi	<ul style="list-style-type: none">• Robust optimization minimizing worst case performance• Reference model is pre-trained on some mixture• Optimal mixture to minimize per-domain excess loss compared to reference	<ul style="list-style-type: none">+ Only two proxy models required+ Optimization criterion promises good generalizability- Identified mixture very sensitive to model architecture
Mixing Laws	<ul style="list-style-type: none">• Formulate parametrized law to model performance as function of mixture weights• Law-fitting based on limited number of evaluations• Optimal mixture predicted as minimizer of law	<ul style="list-style-type: none">+ Complexity reduction of search space bounded too law- Risk of modelling errors in law, (example: BiMix not modelling any domain interactions)

Three key challenges to tackle for novel mixture optimization approach

Existing Mixture Optimization Approaches and their Pros & Cons

Name	Method	Pros & Cons
Online Mixing	<ul style="list-style-type: none">• Data Mixing while training the production grade model• Multi-Armed Bandit to model which domain to sample next token• Domain perplexity as reward function	<ul style="list-style-type: none">+ Negligible run-time overhead without pre-training proxy models- Missing interpretability into overall optimal mixture
DoReMi	<ul style="list-style-type: none">• Robust optimization minimizing worst case performance• Reference model is pre-trained on some mixture• Optimal mixture to minimize per-domain excess loss compared to reference	<ul style="list-style-type: none">+ Only two proxy models required+ Optimization criterion promises good generalizability- Identified mixture very sensitive to model architecture
Mixing Laws	<ul style="list-style-type: none">• Formulate parametrized law to model performance as function of mixture weights• Law-fitting based on limited number of evaluations• Optimal mixture predicted as minimizer of law	<ul style="list-style-type: none">+ Complexity reduction of search space bounded too law- Risk of modelling errors in law, (example: BiMix not modelling any domain interactions)

1

Compute Efficiency: Minimal overhead unlocked by high sample efficiency

2

Flexibility: Reduced dependency on pre-defined functional shapes to prevent modelling errors

3

Insight: Identified mixture must provide insight into the true quality/utility of each domain after training



**Method: Bayesian Optimization for Data
Mixing**

BO is a Sequential Optimization Process Towards Finding Function Optimums

- All **previous mixture weights** x_i and **performance value** y_i collected in dataset

$$D_t = \{(y_i, x_i)\}_{i=0}^t$$

- **Gaussian Process** provides **probability distribution** of performance at each sampling point

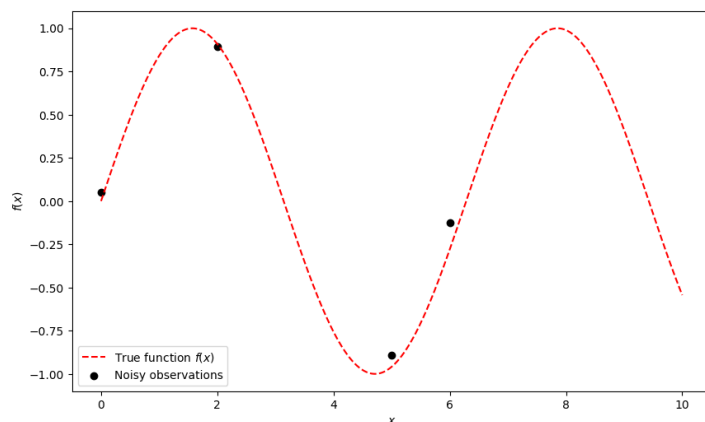
$$g(x, y | D) = \Pr(f(x) = y | D)$$

- **Acquisition function** measures utility of next-weight x_{t+1}
- **Tradeoff** between **exploitation** of previous well performing mixtures and **exploration** of unknown mixture regions
- **Maximizing acquisition** function gives **next mixture weight** in smart way

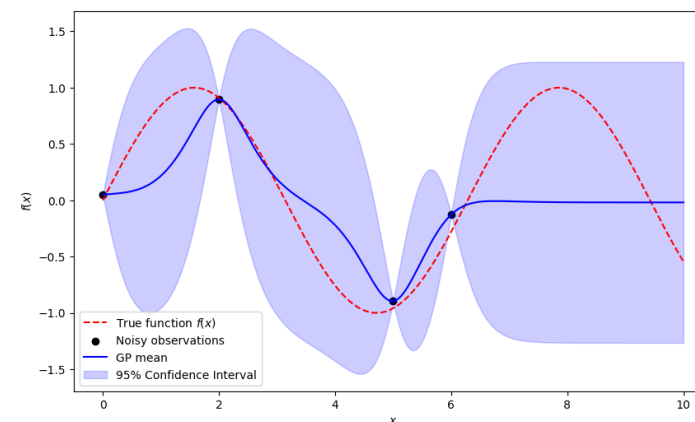
BO has potential for
high-sample efficiency,
without **systematic**
errors introduced by
mixing laws

BO Loop in Detail

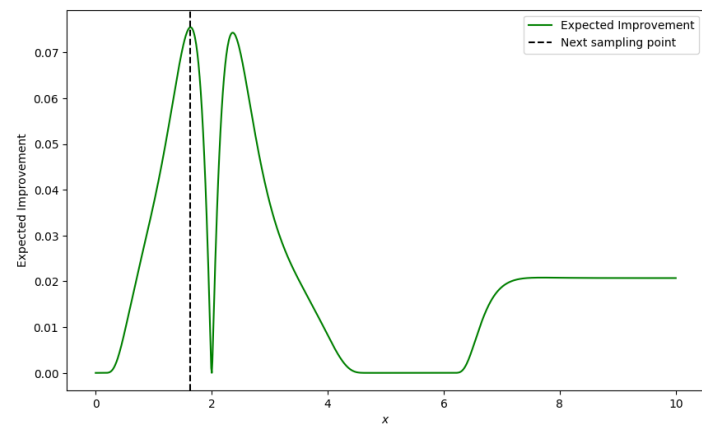
- 1 Existing function evaluations stored in dataset D_t



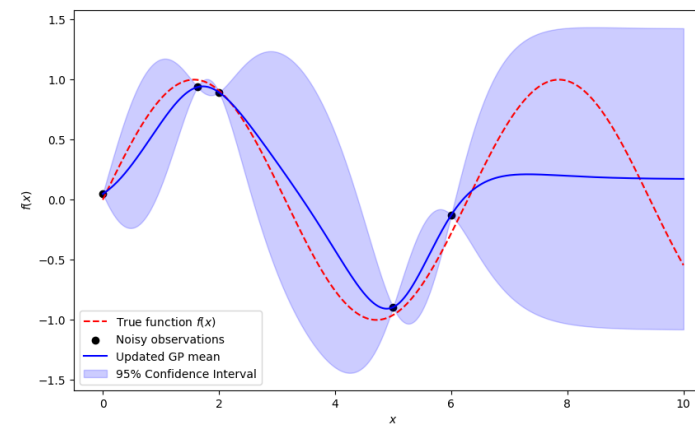
- 2 Gaussian Process builds model of the function of interest



- 3 Acquisition function measures utility of next sampling point



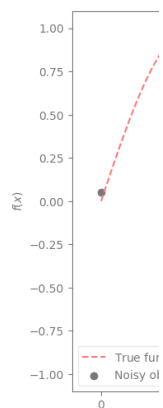
- 4 GP updated based on novel function evaluation



BO Loop in Detail

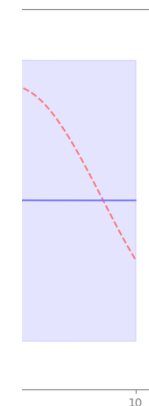
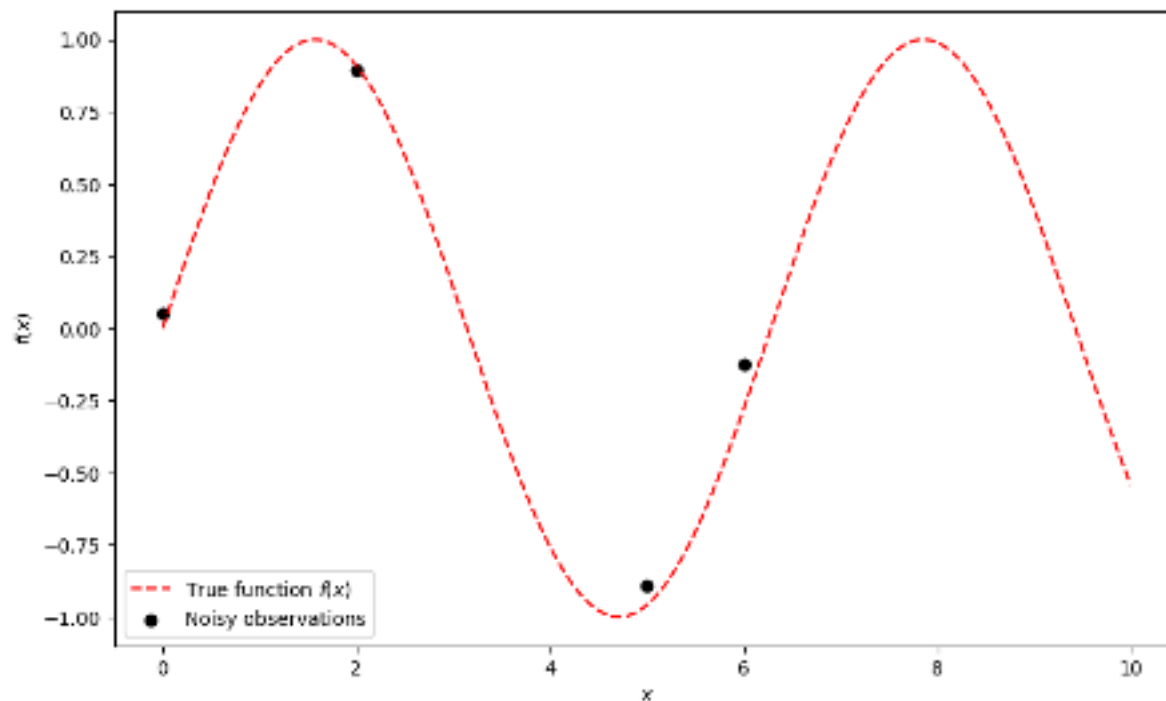
1 Existing function evaluations stored in dataset D_t

2 Gaussian Process builds model of the function

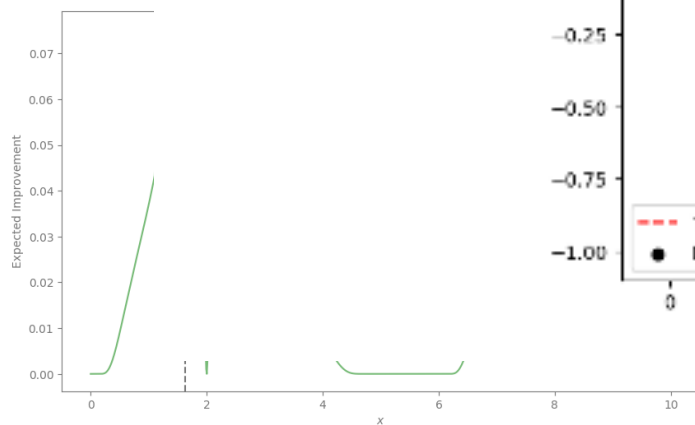


1

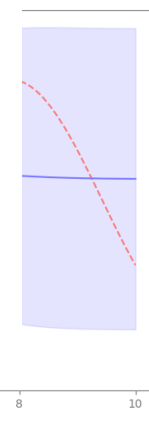
Existing function evaluations stored in dataset D_t



3 Acquisition sampling point



Acquisition



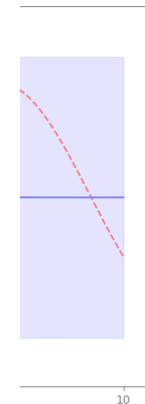
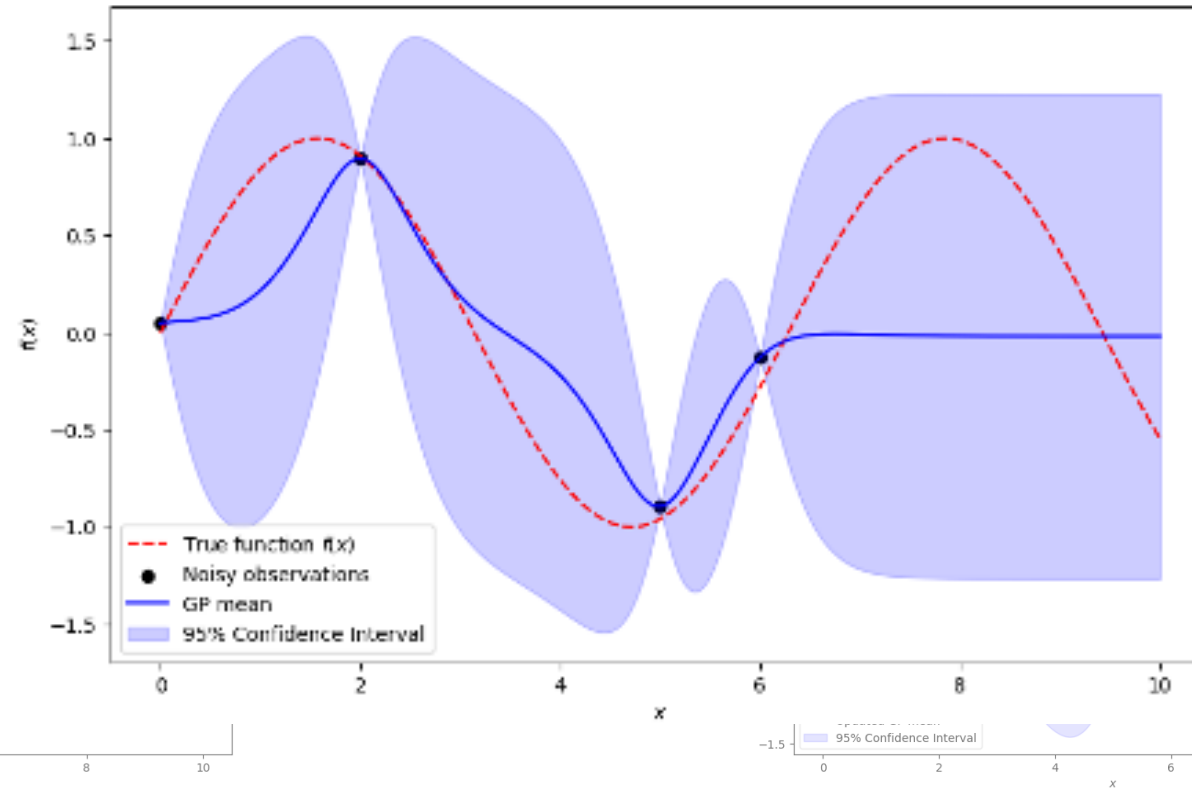
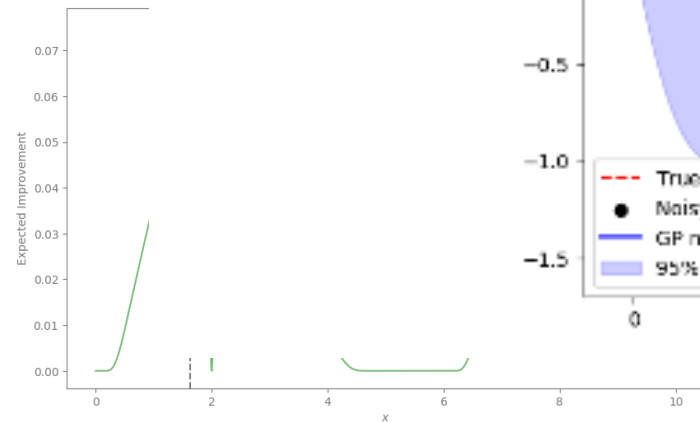
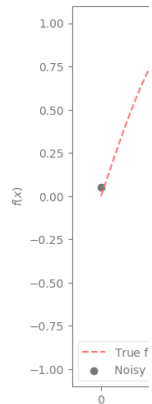
BO Loop in Detail

1 Existing function evaluations stored in dataset D_t

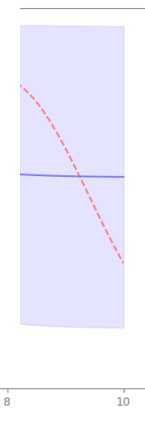
2 Gaussian Process builds model of the function of interest

2 Gaussian Process builds model of the function of interest

3 Acquisition sampling probability



ction



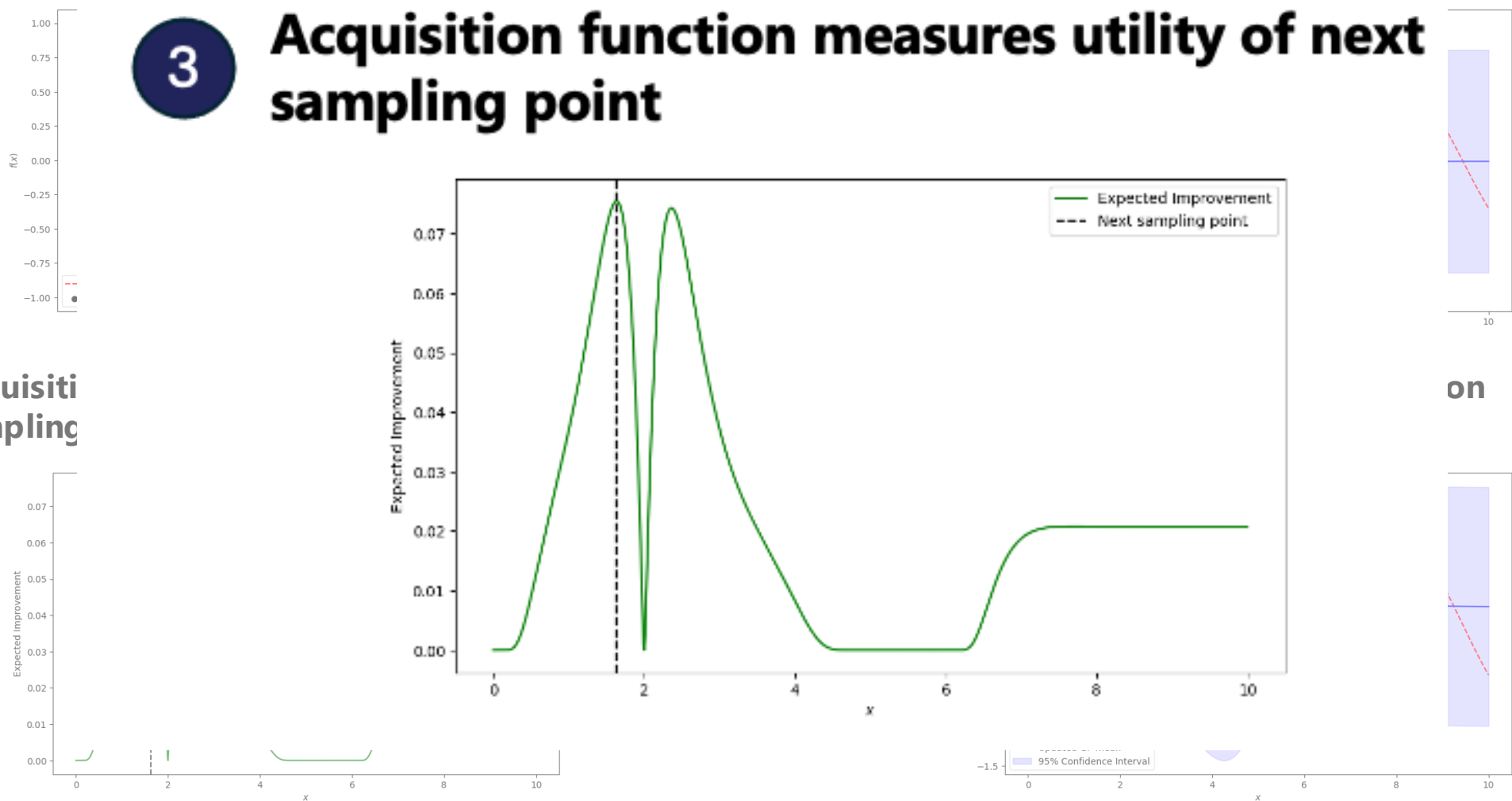
BO Loop in Detail

1 Existing function evaluations stored in dataset I

2 Gaussian Process builds model of the function $f(x)$

3 Acquisition function measures utility of next sampling point

3 Acquisition sampling



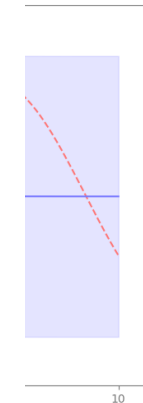
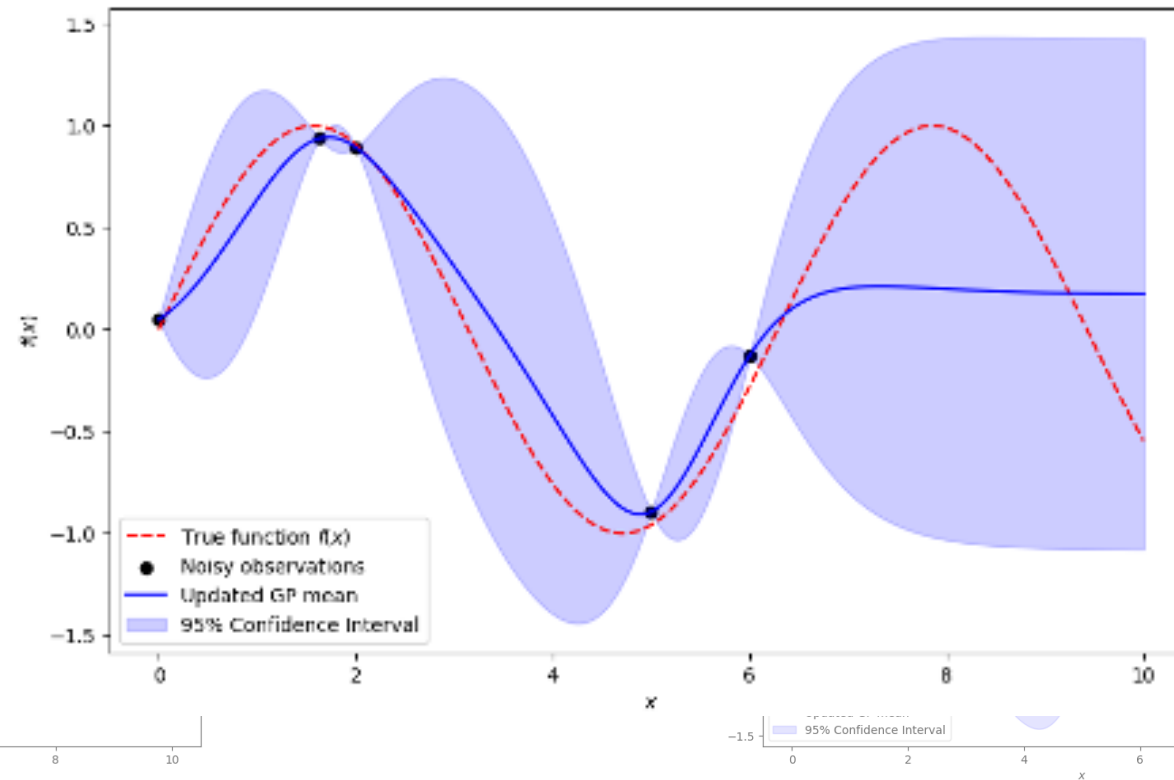
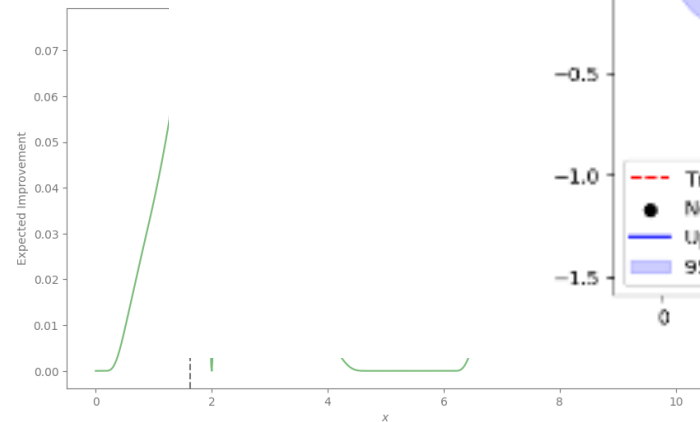
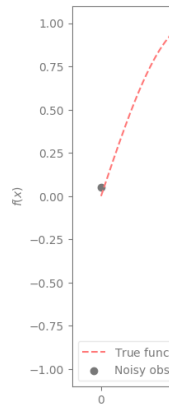
BO Loop in Detail

1 Existing function evaluations stored in dataset D_t

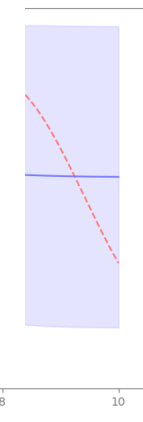
2 Gaussian Process builds model of the function

4 **GP updated based on novel function evaluation**

3 Acquisition function sampling point



ction





Experiments

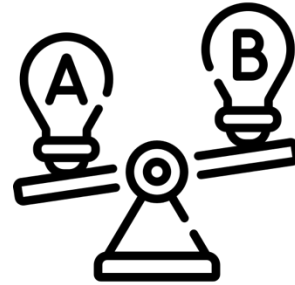
Experimental Setup

Data Sources



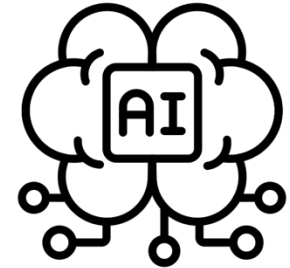
1. **Books:** Book Texts
2. **Common Crawl:** Web Data
3. **The Stack:** Code Data
4. **Pes20:** Academic Papers
5. **Reddit:** Social Media

Baselines



1. **Sobol:** Random Search
2. **Uniform:** 20% uniform weights
3. **LLaMa:** Weights reported
4. **BiMix:** SOTA Mixing Law
5. **MixLaws:** SOTA Mixing Law

Model & Metric

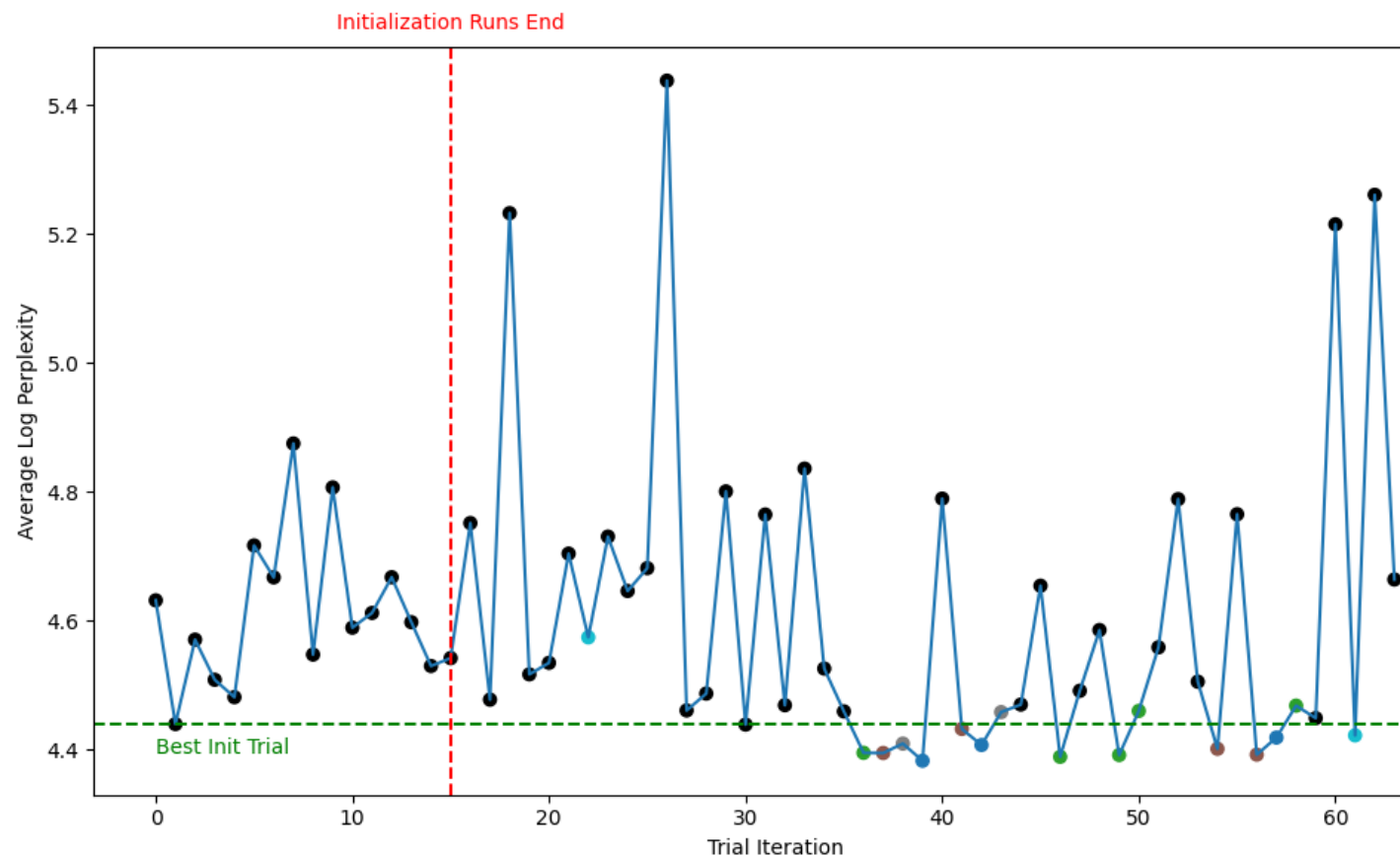


- **25M** parameters model size, trained **Chinchilla optimal**
- **64 trials** for each approach
- **Average perplexity** across all domains is **optimization criterion¹**

¹Each domain validated on held-out set, average compute across domains

Function Evaluations over time give insight into BO Process

BO performance over trial iterations¹



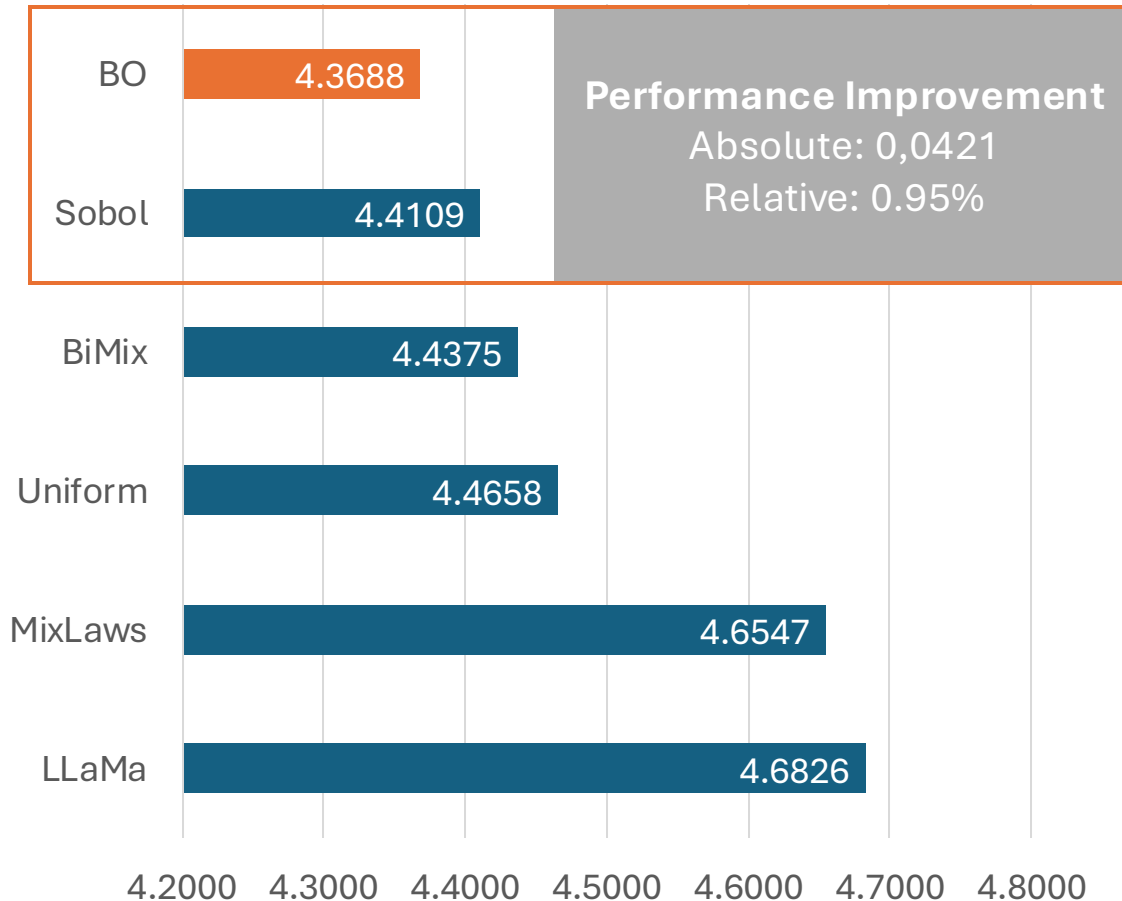
- BO achieves **significant improvement** over **initialization**
- **>35 trials required** before proposing good mixture
- **Exploitation:** many optimum share similar weight-cluster

¹Colors visualize top 5 clusters of similar weights

BO is effective and outperforms all other baselines

Performance Comparison against benchmarks

[Values in avg. log perplexity]

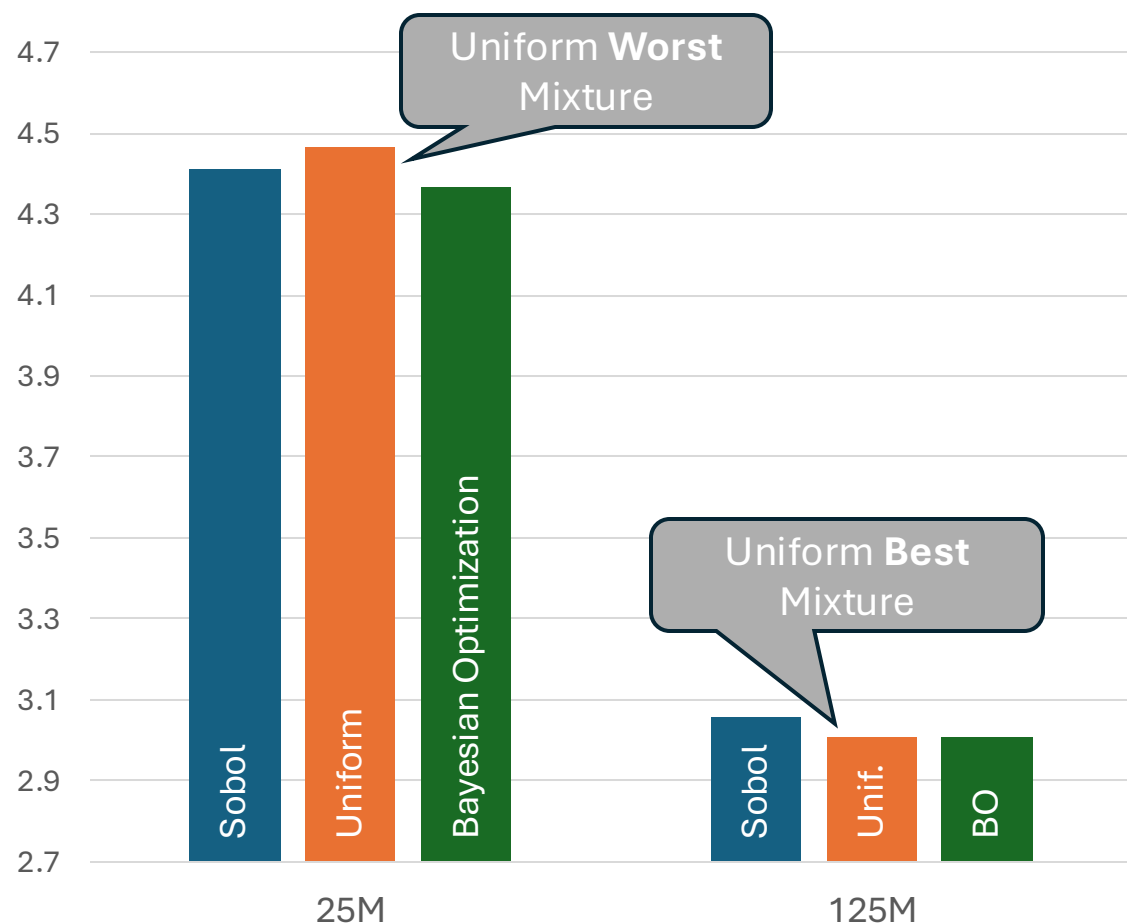


- **BO effective** and **identifies best-performing** mixture
- **$\Delta 0.04$ improvement significant** and larger than experiment variance of $\sigma < 0.02$
- **No other** approach capable to **outperform random search**
- **MixLaws** with major **struggle** to model **true function shape** leading to bad performance
- **LLaMa weights worst**. High focus on CC could be disadvantage in large model¹

¹According to Kang et al.

Missing Scale Invariance Hinders Real World Application of BO

Scaling up same mixture from 25M model to 125M model [Values in avg. log perplexity]

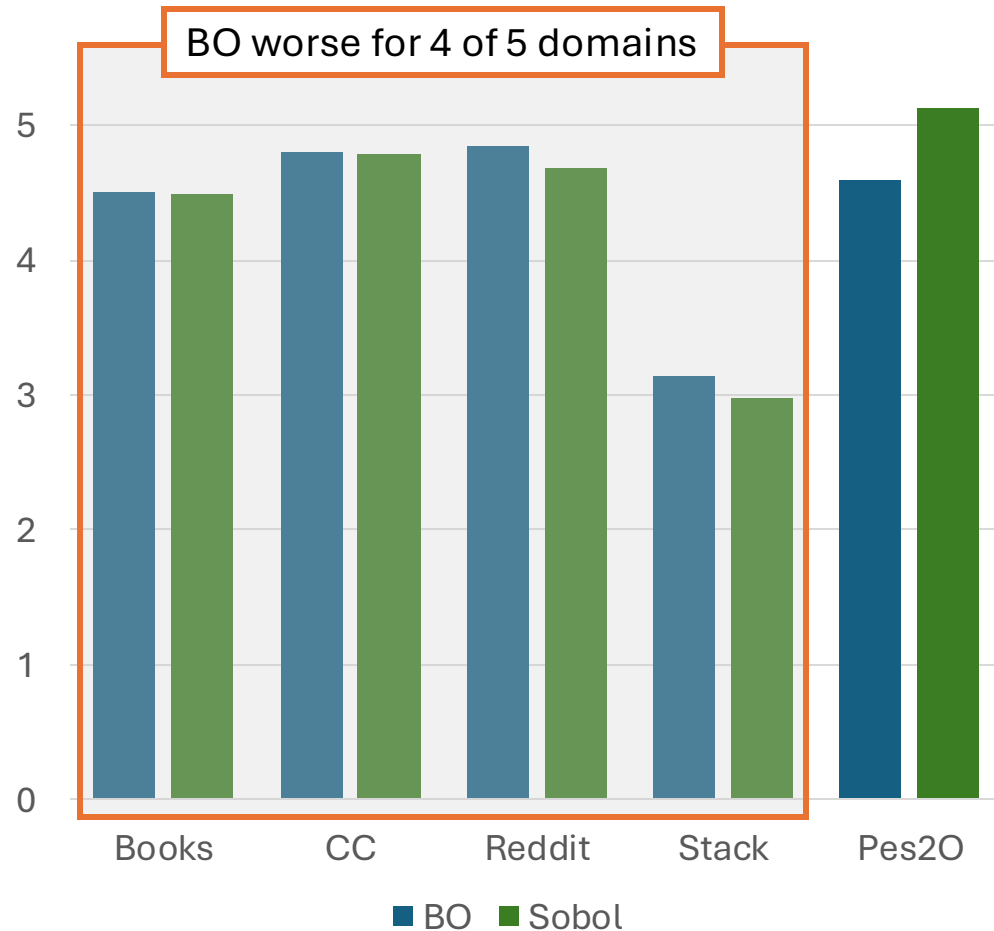


- **Assumption:** Optimal mixture identified at small proxy model translates to larger production model
- **Experiment:** Optimal 25M parameter mixture worse than non-optimal mixture when training 125M model
- When scaling up **uniform mixture** switches **from worst to best** performing mixture

➡ **Limited usability of our method.** More insight into scaling required to translate from proxy to production

Avg. Perplexity Criterion can Create Undesirable Optimization Incentives

Best BO mixture compared to best Sobol mixture per domain [Values in avg. log perplexity]



Goal of Data Selection:

- Create model that **performs well on training domains**
- Additional **generalization** to novel domains **desirable**

Limitations of Avg. Perplexity Criterion:

- Risk that model **trades of domains** to minimize overall loss
- Despite **better overall**, BO only outperforms Sobol on **one domain**



Conclusion

Conclusion: BO effective at identifying optimal mixtures at proxy model scale. Limitations currently prevent wide-scale adaptation

Result 1:

**BO effective for
optimizing mixtures
of one architecture**

- **BO only approach** capable of **outperforming random search**
- **Sample efficiency** expected from theory **validated** in practice
- **No magic:** Significant number of trials still required

Result 2:

**Major limitations
prevent adaptation
of our (and other)
methods in rea-
world**

- **Scale invariance** assumption **invalidated**. **Translating** insights from **proxy to production** model **difficult** in practice
- **Limitations of average performance criterion** identified. High **generalizability requires different optimization objective**
- **Not just our approach** suffers from limitations
→ **More research** required for **principle data mixtures**