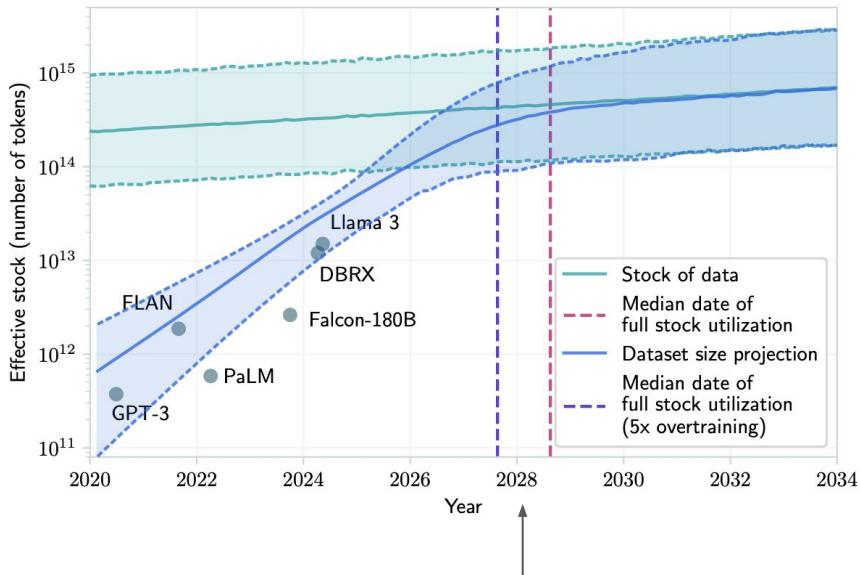


Data Attribution and Optimization for Data-Efficient AI

Xiaoqiang Lin

July 4th 2025 @ SJTU

Data Consumption VS Data Production

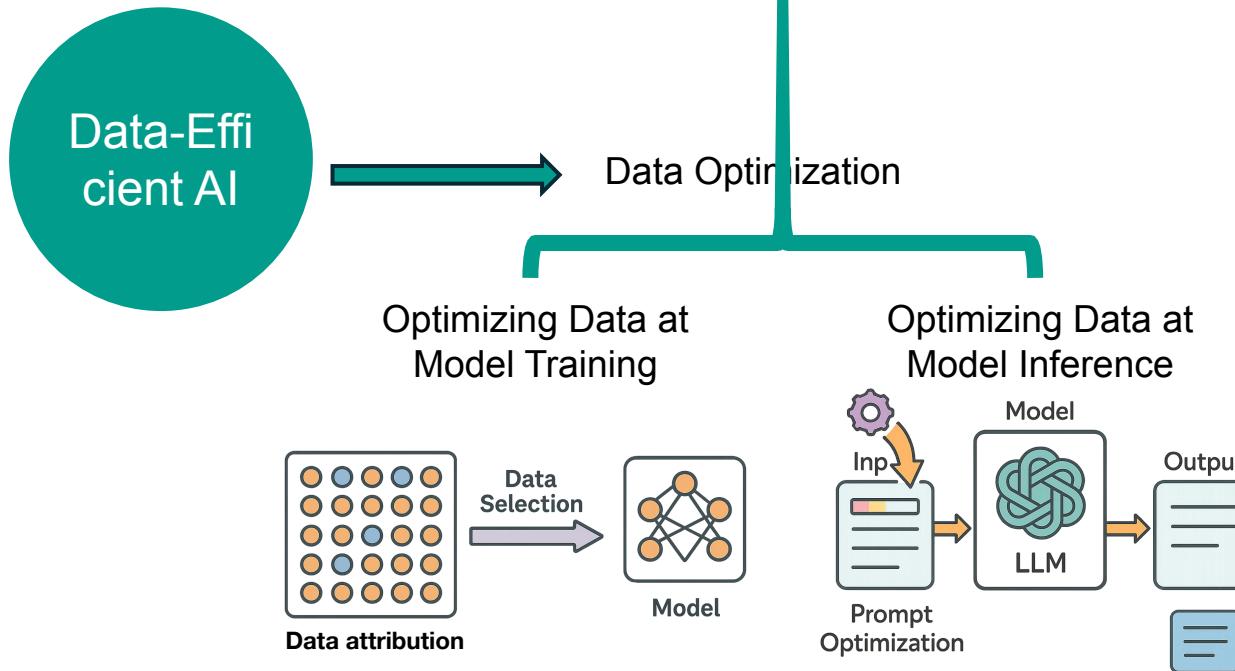


- Data consumption grows faster than production.
- Increase model does not improve performance proportionally.

We will run out of data stock here! [1]

[1] Villalobos, P., Ho, A., Sevilla, J., Besiroglu, T., Heim, L., & Hobbahn, M. (2024, July). Position: Will we run out of data? Limits of LLM scaling based on human-generated data. ICML 2024.

Data-Efficient AI



Agenda

Optimizing data at model training:

- FreeShap: Data Attribution for LLM Fine-tuning, ICML 2024
- NICE: Data Attribution for Non-differentiable Metrics, ICML 2025

Optimizing data at model inference:

- INSTINCT: Black-box Prompt Optimization, ICML 2024
- POHF: Prompt Optimization with Human Feedback, ICML workshop 2024

FreeShap

Helpful or Harmful Data? Fine-tuning-free Shapley Attribution for Explaining Language Model Predictions.

Jingtan Wang*, Xiaoqiang Lin*, Rui Qiao*, Chuan-Sheng Foo, Bryan Kian Hsiang Low.

ICML 2024.

Motivation

The Data Value / Data Contribution

The leave-one-out (LOO)
value of the -th data point:
(~influence function)

Marginal contribution of i for
the whole dataset

$$\xi_i = [U(N) - U(N \setminus \{i\})]$$

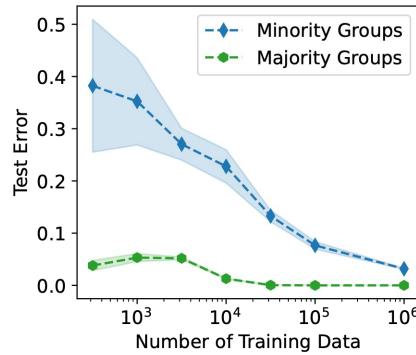
The Shapley value of
the i-th data point:

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} [(U(S \cup \{i\}) - U(S))]$$

Marginal contribution of i for
the subset

Preliminaries

- The harm of the mislabeled data is magnified when the dataset has a smaller size.

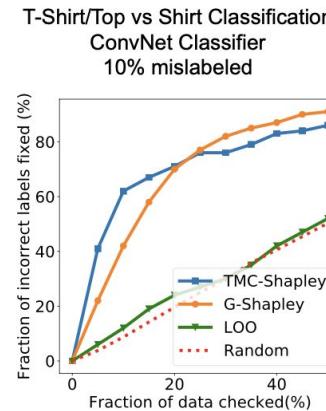
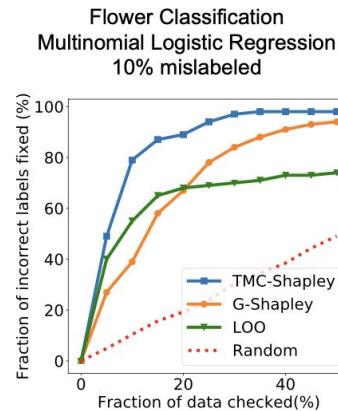
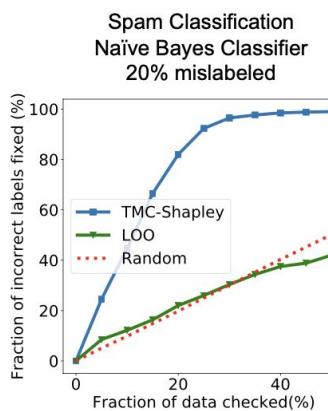


- The Shapley value is better at detecting mislabeled data because of the consideration of the smaller subsets.

Motivation: Shapley Value is More Effective

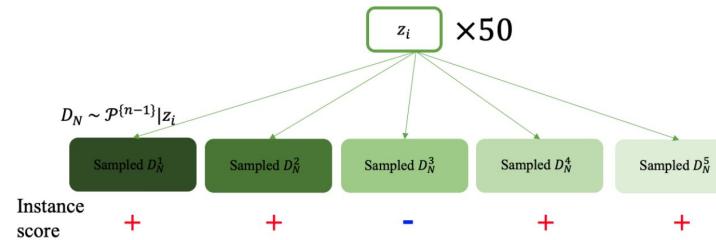
The Shapley Value

- The Shapley value is much better at identifying mislabeled data than LOO values



Methodology: Shapley Value is More Robust

- Place the same text sample into different training datasets and evaluate their instance scores



- Compared with LOO, Shapley value can output helpfulness/harmfulness consistently
 - Usefulness in one dataset is generalizable to other datasets

Text example	Label	Shap	Shapley score	LOO	LOO score
confident filmmaking and a pair of fascinating performances seriously, rent the Disney version	[pos] [neg]	5+/0- 5+/0-	0.45 ± 0.23 1.85 ± 0.72	3+/2- 3+/2-	0.23 ± 1.52 -0.70 ± 2.36
it's not going to be everyone's bag of popcorn, but it definitely gives you something to chew on' would fit chan like a \$ 99 bargain-basement special	[pos] [neg]	0+/5- 0+/5-	-1.91 ± 0.57 -0.36 ± 0.20	3+/2- 2+/3-	0.23 ± 1.34 0.00 ± 1.45

Challenge of SV in Large Models Training

Computational Scalability

$$\phi_i = \underbrace{\sum_{S \subseteq N \setminus \{i\}}}_{2^{|N|-1} \text{ number of calculations}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} \underbrace{(U(S \cup \{i\}) - U(S))}_{\text{requires fine-tuning}}$$

- The naive evaluation of the utility function requires fine-tuning a (large) pretrained model on the subset, then evaluate the model's performance.
- Repeated fine-tuning are extremely costly.

Solution

The Empirical “Kernel” Trick

$$\Psi = \begin{bmatrix} \nabla_{\theta} f(x_1)^\top \\ \nabla_{\theta} f(x_2)^\top \\ \nabla_{\theta} f(x_3)^\top \\ \vdots \\ \nabla_{\theta} f(x_n)^\top \end{bmatrix}$$

empirical Neural Tangent
Feature

$$K = \begin{pmatrix} \nabla_{\theta} f(x_1)^\top \nabla_{\theta} f(x_1) & \nabla_{\theta} f(x_1)^\top \nabla_{\theta} f(x_2) & \cdots \\ \nabla_{\theta} f(x_2)^\top \nabla_{\theta} f(x_1) & \nabla_{\theta} f(x_2)^\top \nabla_{\theta} f(x_2) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

empirical Neural Tangent
Kernel (eNTK)

Prior work: Kernel regression on the eNTK resembles fine-tuning.

FreeShap

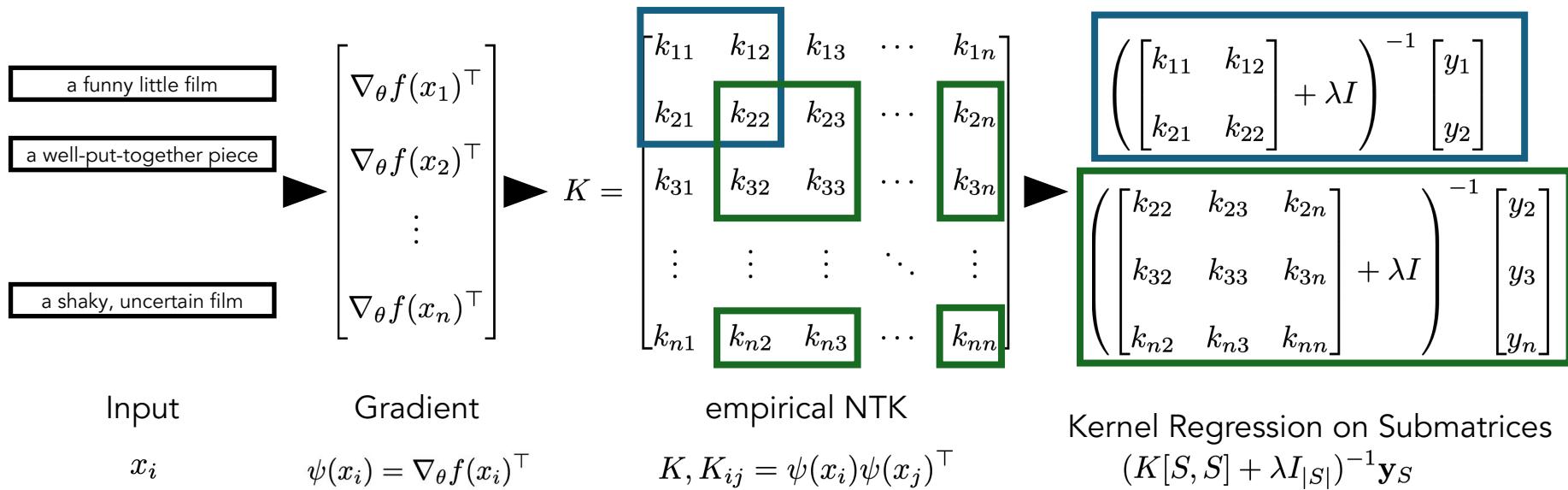
Fine-tuning-free Shapley Value

FreeShap amortizes the fine-tuning cost by pre-computing the eNTK matrix, then calculates the utility terms of the Shapley value using models obtained from kernel regressions.

FreeShap

Fine-tuning-free Shapley Value

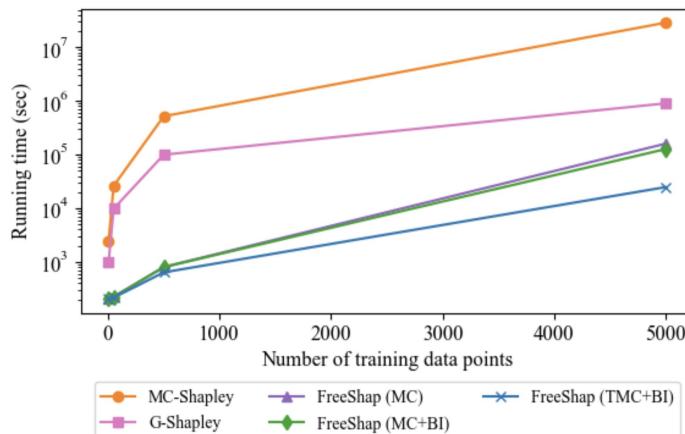
$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (U(S \cup \{i\}) - U(S))$$



FreeShap

Efficiency

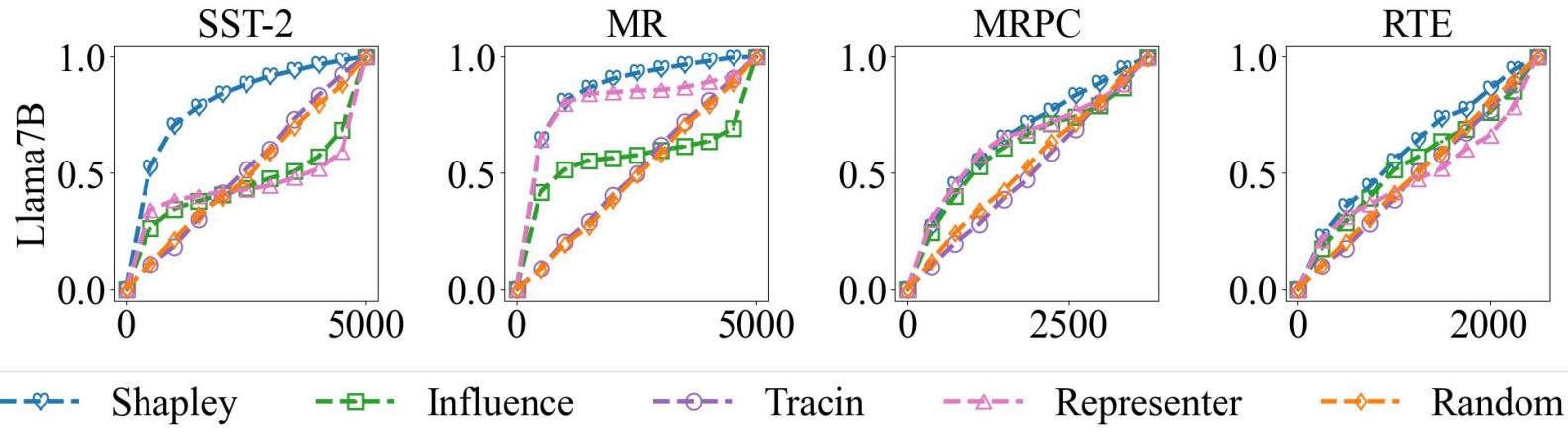
- FreeShap is significantly faster than other approximated Shapley value baselines.



Dataset	SST-2
Model	BERT

Data Curation: Wrong Label Detection

- Poison training set by flipping 10% of data, and then examine the poisoned data by reviewing data points in the order of their scores from lowest to highest.
- Shapley **excels in detecting mislabeled data** within datasets.



Data Selection

- Setting: A train set, a test set for calculating training point scores, and a *held-out* set evaluating the selected subsets.
- Sequentially add training data points with the highest scores.
- The higher the performance increase, the better the data curation approach is.**
- Shapley value is also effective for data curation when test distribution is unknown.

BERT					
	2%	4%	6%	8%	10%
Shapley	0.1951	0.2111	0.2148	0.2167	0.2223
Influence	0.0272	0.0647	0.0393	0.0647	0.0750
TracIn	0.1370	0.1904	0.2008	0.1548	0.1970
Representer	0.1182	0.1351	0.1388	0.0619	0.1238
Random	0.1970	0.1951	0.2073	0.2017	0.1529

Llama2					
	2%	4%	6%	8%	10%
Shapley	0.0816	0.1398	0.1754	0.2186	0.2458
Influence	-0.0038	-0.0038	-0.0038	-0.0038	-0.0038
TracIn	0.0638	0.1144	0.1304	0.1773	0.2092
Representer	-0.0019	0.0000	0.0019	0.0066	0.0056
Random	0.0638	0.1041	0.1529	0.1792	0.1951

Summary

- FreeShap provides an efficient and scalable approximation of the Shapley value.
- FreeShap demonstrates strong capability in mislabeled data detection, advancing data diagnostics, which can be used to increase model reliability.

NICE

NICE: Non-differentiable Evaluation Metric-based Data Selection for Instruction Tuning.

Jingtan Wang, Xiaoqiang Lin, Rui Qiao, Pang Wei Koh, Chuan-Sheng Foo, Bryan Kian Hsiang Low.

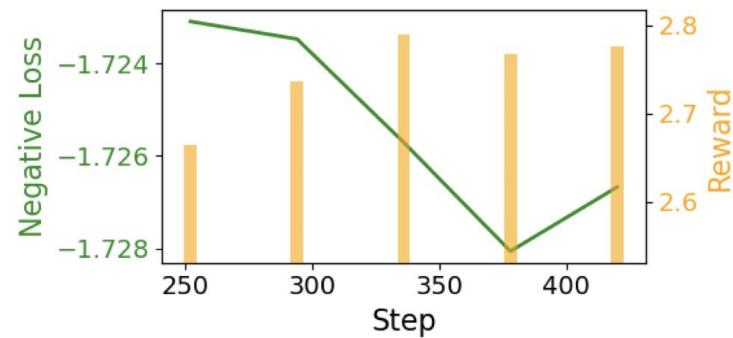
ICML 2025.

Motivation

Loss-based influence (TracIn, Influence function, etc.) estimates the effect of each training data on the validation loss via the gradient of the validation loss.

- Discrepancy Between Loss and Evaluation Metrics
- Reliance on Labeled Validation Data

$$\begin{aligned}\mathcal{I}_{\text{up},\text{loss}}(z, z_{\text{test}}) &\stackrel{\text{def}}{=} \frac{dL(z_{\text{test}}, \hat{\theta}_{\epsilon,z})}{d\epsilon} \Big|_{\epsilon=0} \\ &= \nabla_{\theta} L(z_{\text{test}}, \hat{\theta})^{\top} \frac{d\hat{\theta}_{\epsilon,z}}{d\epsilon} \Big|_{\epsilon=0} \\ &= -\nabla_{\theta} L(z_{\text{test}}, \hat{\theta})^{\top} H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta}).\end{aligned}$$



NICE: Non-differentiable evaluation metrics-based InfluenCe Estimation

Instead of using validation loss, we use $L_r(z_v; \theta) = \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta)}[-r(z_v, \hat{y}_v)]$,

Quantify the influence of including a training point z_i on the validation point z_v 's performance measured by a non-differentiable evaluation metric r

- Code pass rate, math correct rate
- Reward model

TracIn

$$Inf_{NICE} = \sum_{e=1}^E \bar{\eta}_e \left(\nabla_{\theta} L(z_i; \theta^e), \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^e)}[-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^e)) r(z_v, \hat{y}_v)] \right)$$

Training point's gradients w.r.t. loss
across checkpoints Validation point's **policy gradients** w.r.t. eval metrics

Influence function

$$Inf_{NICE_{IF}} = \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^E)}[-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^E)) r(z_v, \hat{y}_v)]^T H_{\theta^E}^{-1} \nabla_{\theta} L(z_i; \theta^E)$$

computed w.r.t final checkpoints Training point's gradients and Hessians w.r.t. loss
Validation point's **policy gradients** w.r.t. eval metrics

NICE: Non-differentiable evaluation metrics-based InfluenCe Estimation

Approximated via Monte-Carlo Sampling

$$Inf_{NICE} = \sum_{e=1}^E \bar{\eta}_e \left(\nabla_{\theta} L(z_i; \theta^e), \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^e)} [-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^e)) r(z_v, \hat{y}_v)] \right)$$

Training point's gradients w.r.t. loss
Validation point's **policy gradients** w.r.t. eval metrics

across checkpoints

NICE: Non-differentiable evaluation metrics-based InfluenCe Estimation

$$Inf_{NICEIF} = \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^E)} \left[-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^E)) r(z_v, \hat{y}_v) \right]^T H_{\theta^E}^{-1} \nabla_{\theta} L(z_i; \theta^E)$$

computed w.r.t final checkpoints
Training point's gradients and Hessians w.r.t. loss

Validation point's **policy gradients** w.r.t. eval metrics

- NICE Preferred training points: Those whose gradient are most similar to the policy gradients of the validation performance measured by the reward function r , i.e., those training points can improve the validation performance more
 -  When r does not require the ground truth label y (reward model which only needs prompts x_y and generated response \hat{y}_v) → **NICE can output score without need of validation label!**

Generalization to Other Loss-based Influence Estimation Methods:

Method	AlpacaEval	TLDR	RLHF	HumanEval
IF (DataInf)	11.11	2.01	0.83	37.40
NICEIF	20.44	3.97	1.89	39.68

NICEIF consistently outperforms Vanilla Influence Function

Assisted Monte Carlo

Benefit of Monte-Carlo Sampling used in NICE: Utilizing multiple different responses, offering diverse guidance. The generated response can be better than the label response

Limitation: When the model is too weak, the MC samples may not contain high-quality responses with high rewards

Solution: Assisted Monte Carlo (AMC) uses a model that is better at the target task to assist generation

$$Inf_{NICE_{AMC}} = \sum_{e=1}^E \bar{\eta}_e \langle \nabla_{\theta} L(z_i; \theta^e), \mathbb{E}_{\hat{y}_v^* \sim g(y|x_v; \psi)} [-\nabla_{\theta} \log(f(\hat{y}_v^*|x_v; \theta^e)) r(z_v, \hat{y}_v^*)] \rangle$$

Assisted generated responses

Experimental Results

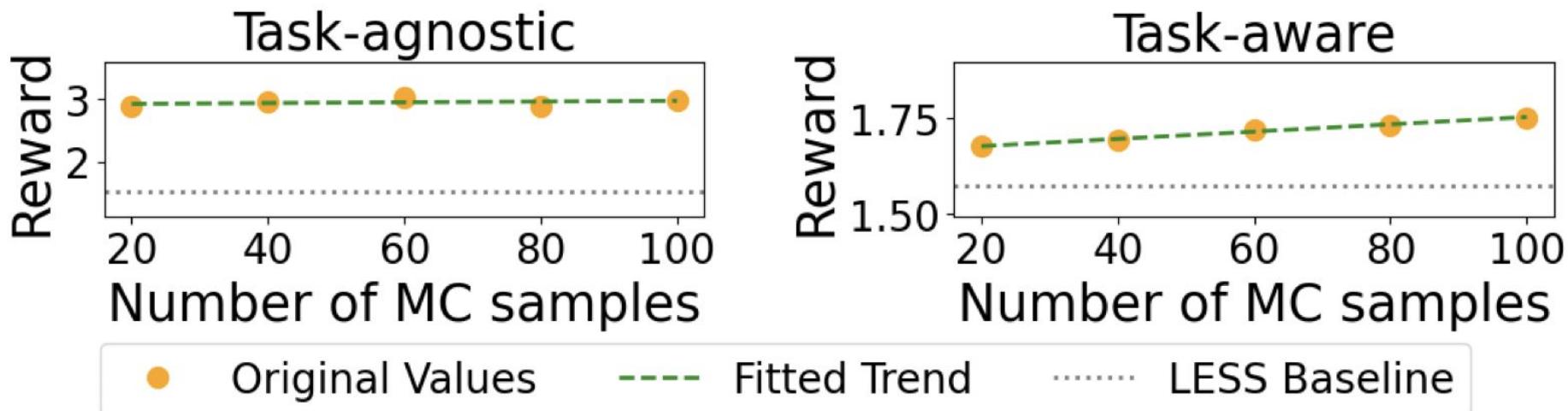
Major Findings:

- NICE outperforms the loss-based influence estimation.
- No labels? No problem! NICE outperform baselines that utilize the label response.
- Less is more: subset outperforms the full dataset.
- Assisted monte-carlo sampling can boost data selection when the size of training data is large (task-agnostic setting).

Model & Dataset		Full	Random	RDS	BM25	DSIR	TSDS	LESS	NICE	NICEAMC
<i>Task-agnostic</i>										
Llama2-7B	AlpacaEval	22.59	16.13 \pm 1.18	14.70	19.60	20.27	17.40 \pm 2.44	26.94 \pm 2.37	27.61 \pm 2.13	30.45 \pm 2.40
	TLDR	2.40	1.80 \pm 0.08	2.08	2.15	1.53	2.19 \pm 0.29	3.37 \pm 0.29	3.61 \pm 0.78	3.55 \pm 0.40
	RLHF	2.31	2.05 \pm 0.11	1.87	2.83	2.57	1.01 \pm 0.12	1.44 \pm 0.07	2.82 \pm 0.10	3.03 \pm 0.02
	HumanEval	47.44	44.30 \pm 2.36	45.29	46.19	42.22	43.68 \pm 1.82	47.50 \pm 1.57	48.59 \pm 2.08	45.10 \pm 2.84
Mistral-7B	AlpacaEval	33.77	24.99 \pm 4.28	21.70	28.47	29.31	35.84 \pm 0.53	41.09 \pm 1.56	41.43 \pm 3.00	47.40 \pm 2.94
	TLDR	2.79	3.06 \pm 0.24	2.90	2.41	3.48	3.28 \pm 0.41	4.40 \pm 0.12	4.80 \pm 0.12	4.59 \pm 0.20
	RLHF	2.56	2.13 \pm 0.04	1.78	2.88	2.94	1.83 \pm 0.15	1.70 \pm 0.09	3.10 \pm 0.06	3.42 \pm 0.05
	HumanEval	83.63	85.56 \pm 1.27	84.15	84.09	79.17	82.78 \pm 1.25	85.24 \pm 0.45	85.59 \pm 1.41	85.67 \pm 0.34
<i>Task-aware</i>										
Llama2-7B	RLHF	1.01	1.04 \pm 0.04	0.66	1.29	1.43	0.97 \pm 0.02	1.62 \pm 0.05	1.69 \pm 0.05	1.32 \pm 0.05
	HumanEval	51.27	51.91 \pm 1.61	54.74	52.23	53.10	49.85 \pm 3.17	52.67 \pm 0.71	55.09 \pm 1.66	50.67 \pm 1.24
Mistral-7B	RLHF	0.99	1.05 \pm 0.04	0.56	1.31	1.31	1.15 \pm 0.06	1.29 \pm 0.13	1.71 \pm 0.01	1.35 \pm 0.07
	HumanEval	84.27	83.34 \pm 2.54	86.75	84.81	79.91	85.51 \pm 1.28	85.26 \pm 1.13	87.35 \pm 1.03	84.18 \pm 1.63

Experimental Results

The Effect of the Number of Monte-Carlo Samples: Positive correlation between performance and generated MC samples



INSTINCT

Use Your INSTINCT: INSTRUCTION optimization for LLMs using Neural bandits Coupled with Transformers.

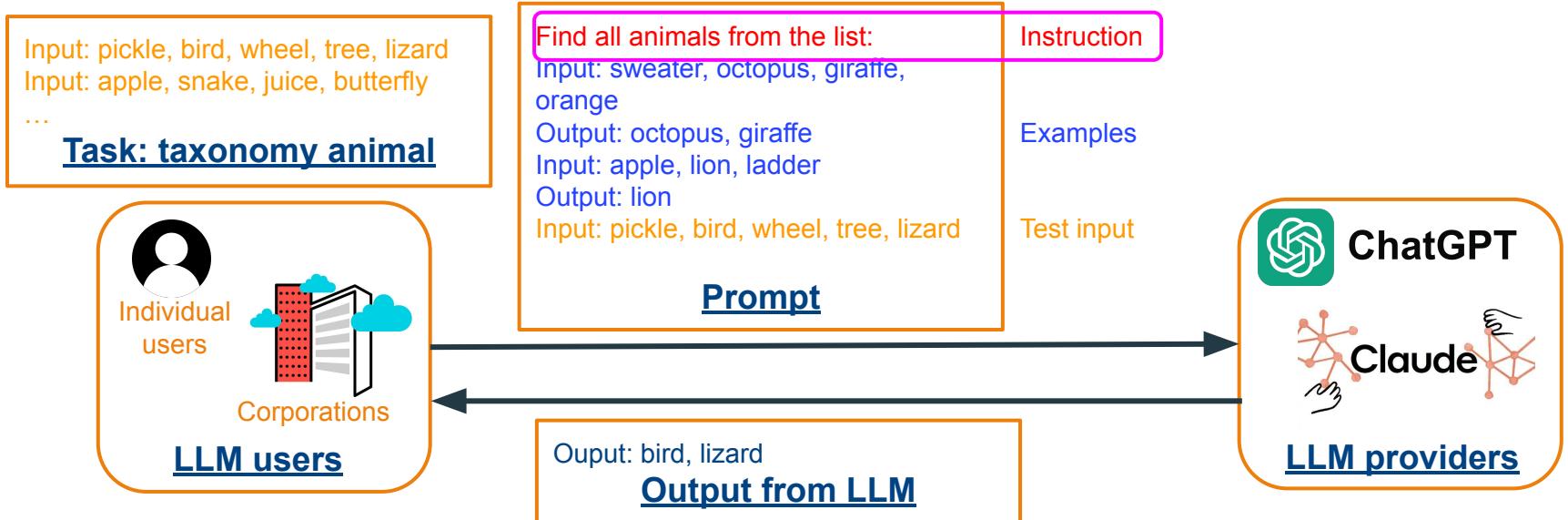
Xiaoqiang Lin*, Zhaoxuan Wu*, Zhongxiang Dai, Wenyang Hu, Yao Shu, See-Kiong Ng, Patrick Jaillet, Bryan Kian Hsiang Low.

ICML 2024.

Motivation

- Few-shot in-context learning

Good instruction is vital to the performance!



Motivation

- Human designed instruction can be **costly** and not good
- **Instruction optimization:** Automatically optimize the instructions to obtain the best performance of LLMs

In-context learning

Find all animals from the list:
Input: pickle, bird, wheel, tree, lizard

Instruction
Test input

Prompt

 ChatGPT

 Claude

LLM providers

Challenges

- Best performing LLMs are **black-box models**
 - ChatGPT, Claude
- Access to black-box LLMs is **costly**
 - API calls are expensive
 - A **query-efficient** approach is needed

Formulation: Instruction Optimization

- Objective:



$$\rho^* = \operatorname{argmax}_\rho h(\rho)$$

$$h(\rho) := \mathbb{E}_{(x,y) \in D_V} s(f(\rho, x), y)$$

Preliminary - Bayesian Optimization (BO)

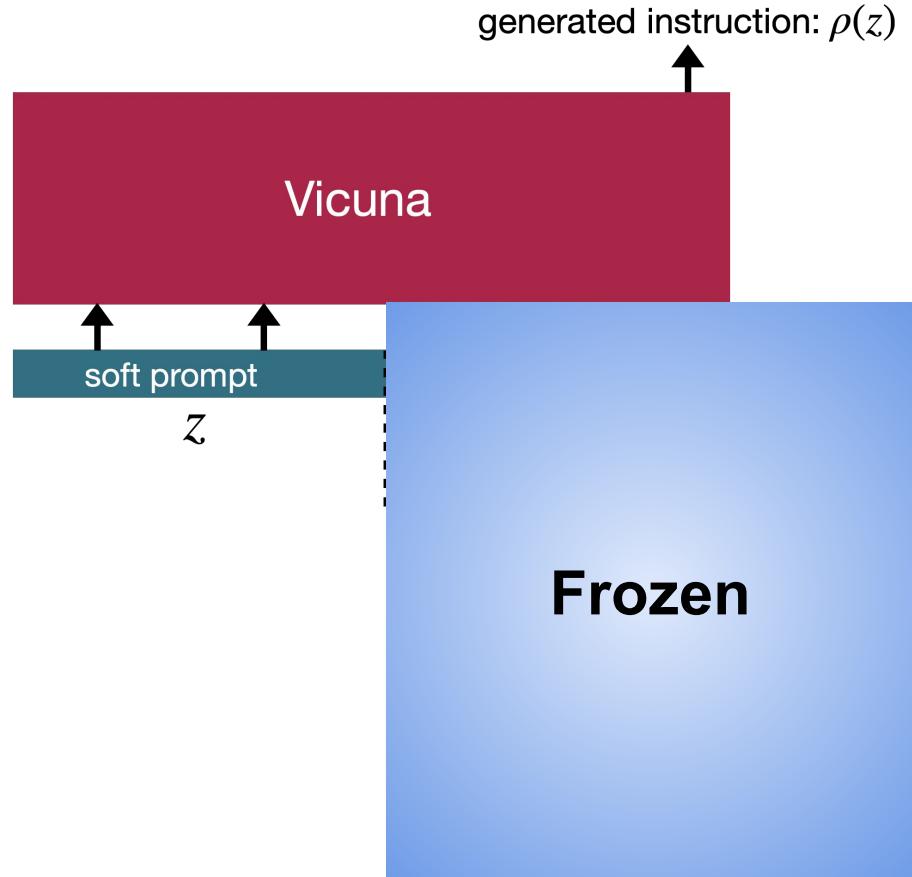
- Sequential black-box optimization: find $\rho^* = \operatorname{argmax}_\rho h(\rho)$
- To choose sequential queries intelligently: ρ_1, \dots, ρ_t
 - Uses a Gaussian process (GP) as a surrogate to model the objective function

Preliminary - Neural Bandits

- Problem with BO:
 - GP is not powerful enough to model the LLM performance.
 - Objective function $h(\theta)$ is not a simple function
- Solution: Use neural networks – neural bandits algorithm
 - Use the **neural networks (e.g., transformers)** as the **surrogate** model
 - Can model **highly complex functions**

INSTINCT Algorithm

- Map a soft prompt \mathcal{Z} (a vector in continuous space) into instruction
 - Search in the continuous space $\rho(z)$



predicted score: $m(g(z); \theta)$

INSTINCT Algorithm

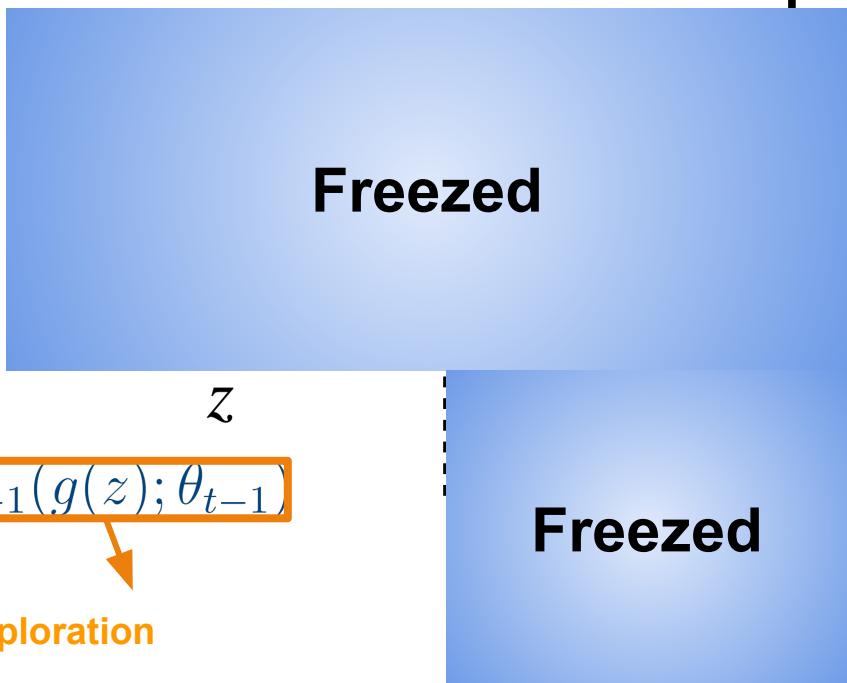
- Uses the whole Vicuna as surrogate model to leverage the expressive power of transformer: $m(g(z); \theta)$
- Acquisition function from NeuralUCB algorithm:

$$z_t = \operatorname{argmax}_{z \in Z} \text{NeuralUCB}_t(z)$$

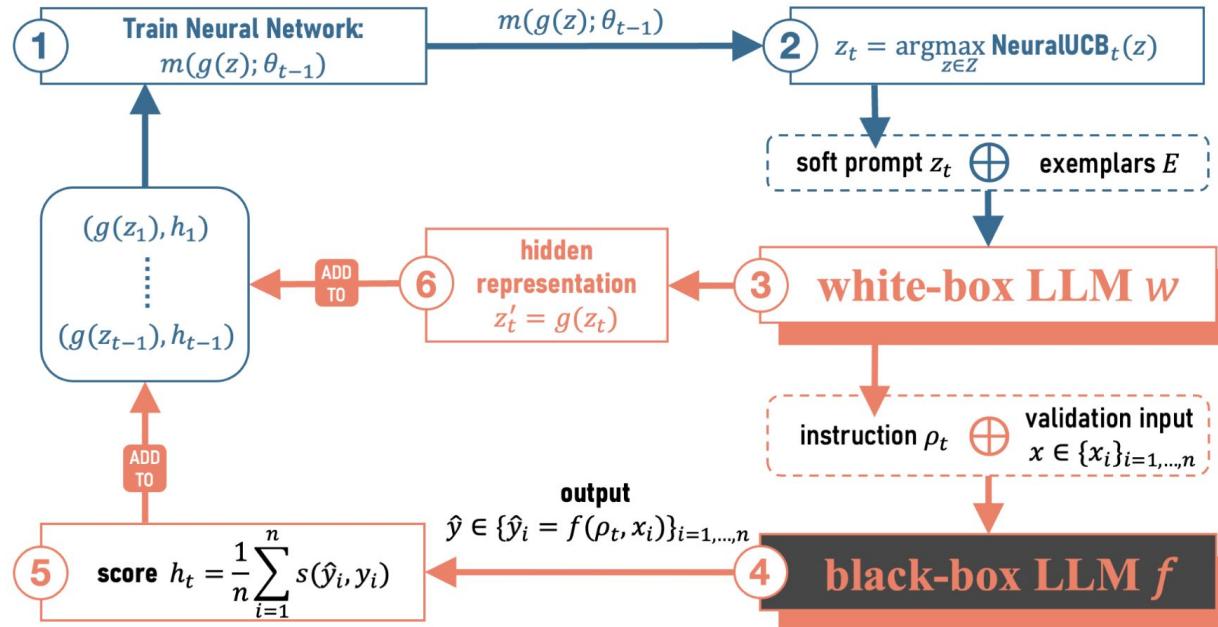
$$\text{NeuralUCB}_t(z) := m(g(z); \theta_{t-1}) + \nu_t \sigma_{t-1}(g(z); \theta_{t-1})$$

Exploitation

Exploration



INSTINCT Algorithm

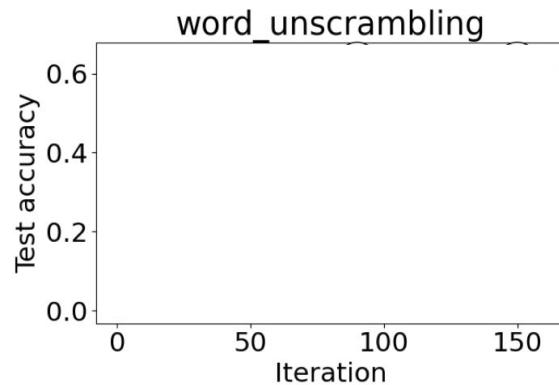


Experiments: Instruction Induction

- Given a task-specific training dataset, find the task-specific instruction that best describes the relationship between inputs and outputs
- Datasets
 - 30 instruction induction tasks curated by [1]

[1] Lichang Chen, Jiupei Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. InstructZero: Efficient instruction optimization for black-box large language models. arXiv preprint arXiv:2306.03082, 2023b.

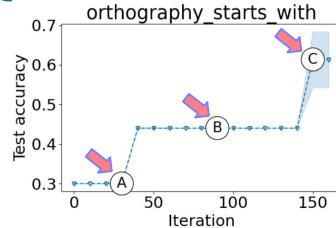
Experiments: Visualizing the Optimization Process



Task description: given a list of shuffled letters, rearrange the letters to form a meaningful word.

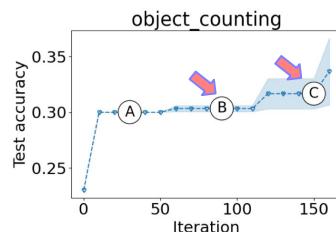
Iteration | Instruction

Experiments: Visualizing the Optimization Process



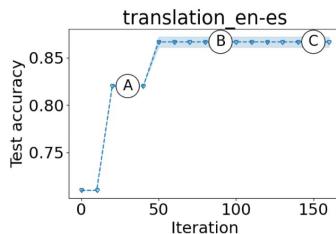
Task description: given a sentence and a letter, output the words that start with the letter in the sentence.

Iteration	Instruction
A	The instruction was to find a word that could be formed by rearranging the letters of the given word
B	The instruction was to find the word that the input corresponds to, and output it
C	The instruction was to output the word that starts with the letter that was inputted



Task description: given a sentence, output the number of objects in the sentence.

Iteration	Instruction
A	The instruction was to output the number of items that the speaker has, given the list of items that the speaker possesses
B	The instruction was to output the number of objects mentioned in the input
C	The instruction was to output the number of items the player has, but the player has entered the number of items instead



Task description: translate the words from English to Spanish.

Iteration	Instruction
A	The instruction was to translate the words from Spanish to English
B	The instruction was to translate the words from English to Spanish
C	The instruction was to translate the words from English to Spanish

Instructions found by
INSTINCT improves
over iterations

Experiments: Instruction Induction

Task	APE	InstructZero	INSTINCT (ours)
antonyms	0.6367(0.1416)	0.8267(0.0072)	0.8467(0.0027)
auto_categorization	0.2500(0.0094)	0.2567(0.0119)	0.2500(0.0330)
auto_debugging	0.2917(0.0340)	0.3750(0.0000)	0.2917(0.0340)
cause_and_effect	0.5733(0.0891)	0.8133(0.0109)	0.5867(0.0871)
common_concept	0.0691(0.0207)	0.0864(0.0398)	0.2129(0.0019)
diff	0.6733(0.2667)	0.6933(0.2224)	1.0000(0.0000)
informal_to_formal	0.5736(0.0026)	0.5310(0.0024)	0.5534(0.0000)
letters_list	1.0000(0.0000)	0.5900(0.1674)	1.0000(0.0000)
negation	0.7533(0.0109)	0.7767(0.0136)	0.8167(0.0027)
object_counting	0.3633(0.0191)	0.3600(0.0929)	0.3400(0.0698)
odd_one_out	0.6333(0.0144)	0.6133(0.0871)	0.7000(0.0163)
orthography_starts_with	0.4567(0.1477)	0.5067(0.0871)	0.6667(0.0272)
rhymes	0.1567(0.0640)	1.0000(0.0000)	1.0000(0.0000)
second_word_letter	0.7467(0.2028)	0.4333(0.1872)	0.1000(0.0411)
sentence_similarity	0.0000(0.0000)	0.0000(0.0000)	0.1400(0.0047)
sum	0.6733(0.2667)	1.0000(0.0000)	1.0000(0.0000)
synonyms	0.3600(0.0759)	0.2767(0.0925)	0.3067(0.0491)
taxonomy_animal	0.3467(0.2341)	0.7167(0.0838)	0.8567(0.0599)
word_sorting	0.3300(0.0374)	0.3100(0.1143)	0.5133(0.0027)
word_unscrambling	0.4400(0.1389)	0.5500(0.0170)	0.6333(0.0072)
# best-performing tasks	5	5	13
# second-best-performing tasks	5	10	5
average rank	2.25	2.0	1.45

Experiments: Instruction Induction (Summarization Task)

- INSTINCT also performs the best in another commonly used *SAMSum* benchmark dataset

Method	ROUGE-1	ROUGE-2	ROUGE-L
APE	0.32549	0.10308	0.30245
InstructZero	0.32595	0.10528	0.30061
INSTINCT	0.35580	0.13350	0.33600

Experiments: Improving Zero-shot CoT

- A well-known zero-shot instruction for chain-of-thought (CoT) reasoning form [1] is

“Let’s think step by step.”

- INSTINCT finds better ones:

Method	Dataset	Best Zero-Shot CoT Instruction	Score
Kojima et al. (2022) InstructZero INSTINCT (ours)	GSM8K	Let’s think step by step.	0.71797
	GSM8K	Let’s use the instruction to solve the problem.	0.74299
	GSM8K	Let’s think about it.	0.74526
Kojima et al. (2022) InstructZero INSTINCT (ours)	AQUA-RAT	Let’s think step by step.	0.52362
	AQUA-RAT	Let’s break down the problem.	0.54331
	AQUA-RAT	I have a new solution.	0.54724
Kojima et al. (2022) InstructZero INSTINCT (ours)	SVAMP	Let’s think step by step.	0.7625
	SVAMP	Let’s use the equation.	0.795
	SVAMP	Let’s use our brains.	0.81

[1] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In Proc. NeurIPS, 2022.

Conclusion

- We introduce the **INSTINCT** to optimize task-specific instructions for black-box LLMs
- **INSTINCT** achieves better performance due to the use of neural-bandits algorithm and the expressive power of the transformer.
- We demonstrate on multiple settings that **INSTINCT** achieves better performance with the same number of queries.

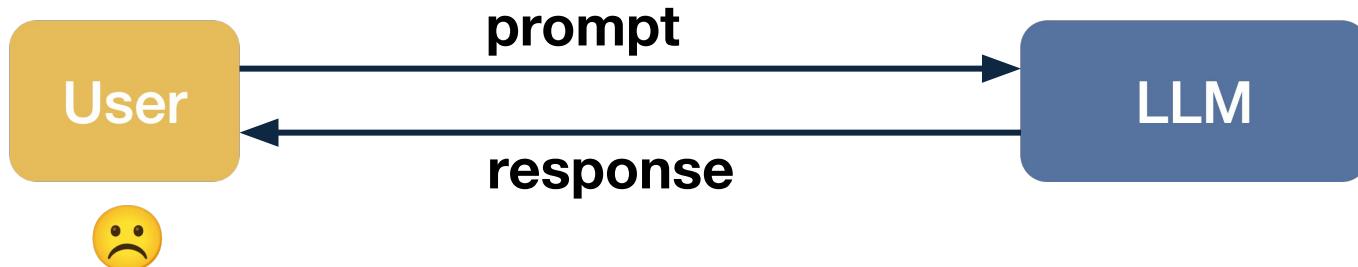
POHF

Prompt Optimization with Human Feedback.

Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-Kiong Ng, Patrick Jaillet, Bryan Kian Hsiang Low.

ICML 2024, Workshop on Models of Human Feedback for AI Alignment. Selected as Oral

Prompt Optimization



Prompt 1

Prompt 2

...

Prompt N



Response 1

Response 2

...

Response N

0.98

😊 Best prompt!

0.72

...

0.81

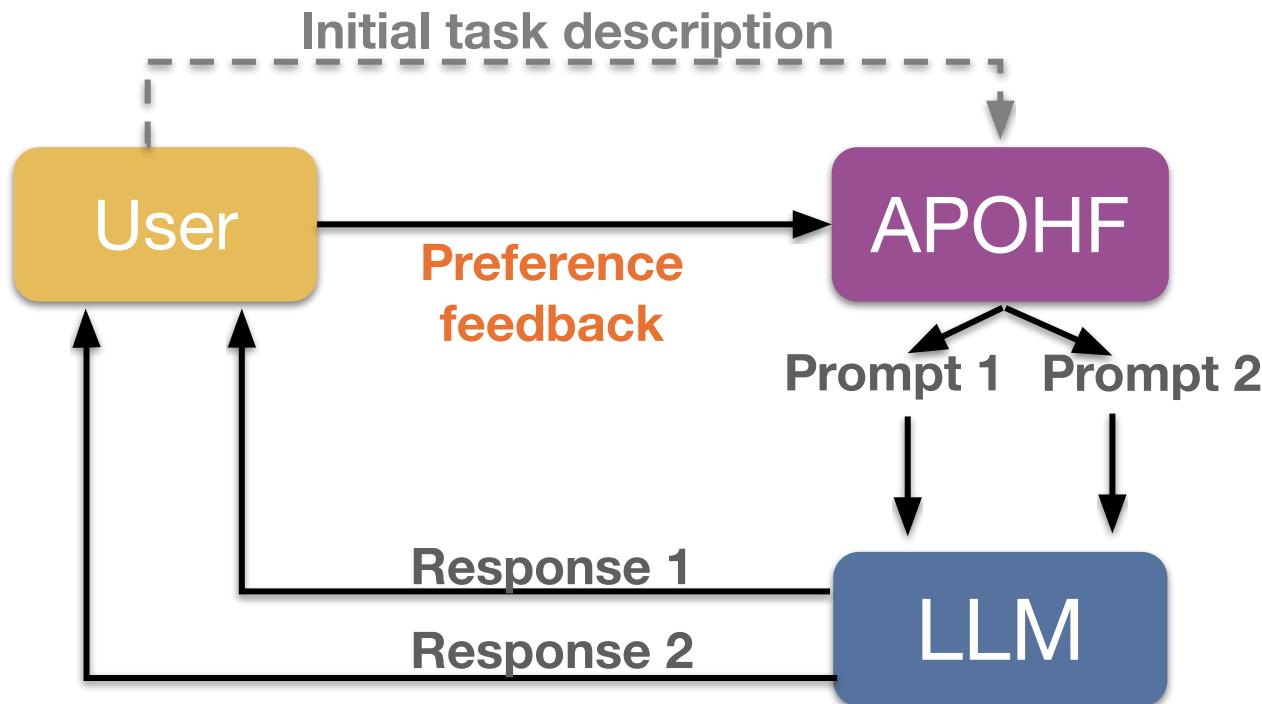
Scoring
method

[Chen et al. (2023);
Lin et al. (2024);
Yang et al. (2024)]

Prompt Optimization

- A scoring method may not be available or reliable
 - No validation dataset available
 - A scorer LLM may not be accurate
 - Human is not good at giving a score (Yue et al. 2012)
- Human is more reliable at providing preference feedback (Yue et al. 2012)
- Can we achieve prompt optimization using only human preference feedback?

Prompt Optimization with Human Feedback



Our algorithm - APOHF

- Using the neural network for latent score prediction
 - $h(x; \theta)$ mapping from prompt to latent score
- Preference feedback model - Bradley-Terry-Luce (BTL) model (Hunter et al. 2004)

$$P(x_1 > x_2) = \sigma(h(x_1; \theta) - h(x_2; \theta))$$

- Given the previous feedback $D_{t-1} = \{x_{s,1}, x_{s,2}, y_s\}_{s=1 \dots t-1}$, train the NN (h) by minimizing the following loss function:

$$\ell(\theta) = -\text{likelihood} \left(y, \sigma(h(x_1; \theta) - h(x_2; \theta)) \right) + \lambda \|\theta\|$$

Our algorithm - APOHF

- Selection of first prompt:

$$\mathbf{x}_{t,1} = \operatorname{argmax}_x \mathbf{h}(x; \boldsymbol{\theta}_t)$$

- Selection of second prompt:

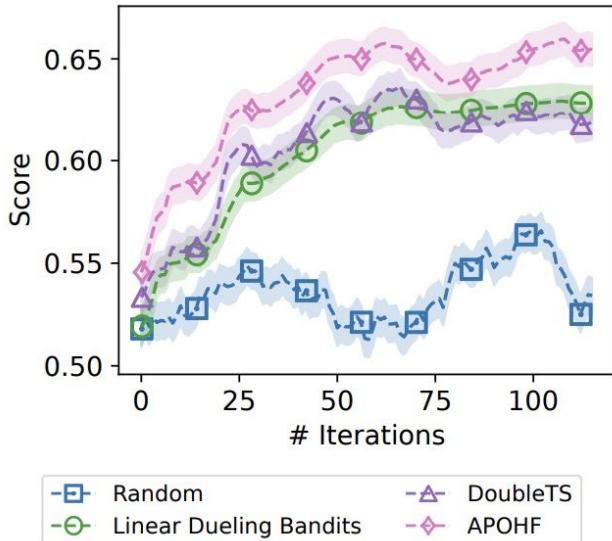
$$\mathbf{x}_{t,2} = \operatorname{argmax}_x \mathbf{h}(x; \boldsymbol{\theta}_t) + \nu \|\nabla \mathbf{h}(\mathbf{x}; \boldsymbol{\theta}_t) - \nabla \mathbf{h}(\mathbf{x}_{t,1}; \boldsymbol{\theta}_t)\|_{V_{t-1}^{-1}}$$

Exploitation:
Score
prediction

Exploration: Encourages $\mathbf{x}_{t,2}$ to
be different from the previously
selected prompts and $\mathbf{x}_{t,1}$

(Arun et al. (2025))

Experiment - Optimization of User Instructions



Task	Iter	Instruction	Score
antonyms	0	add the prefix "un-" to the given words to form their opposites.	0.45
	5	remove the "un-" prefix from each word.	0.45
	10	provide the opposite of the given words.	0.70

Our APOHF consistently and significantly outperforms the other methods.

Experiment - Prompt Optimization for Text-to-Image Generative Models

Ground Truth



Iteration 0



Iteration 10



Iteration 20



Iteration 30



Ground Truth



Iteration 0



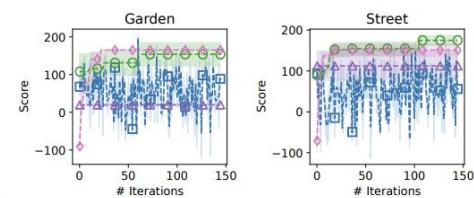
Iteration 5



Iteration 10



Iteration 15



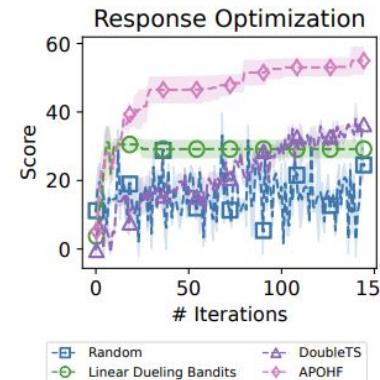
Our APOHF can efficiently produce images which better align with the image the user has in mind.

Experiment - Response Optimization with Human Feedback

The prompt is “Human: What is there to do in Atlantic City?”.

Iter	Response	Score
0	There are many activities to do in Atlantic City! You can visit the famous Atlantic City Boardwalk, try your luck at the casinos, enjoy the beach and water activities, shop at the outlets and specialty stores, dine at top-rated restaurants, enjoy live entertainment and nightlife, visit museums and historical sites, and attend events and festivals throughout the year.	81.96
10	There are plenty of things to do in Atlantic City! You can visit the famous boardwalk, try your luck at the casinos, enjoy a show or concert at one of the many entertainment venues, dine at a variety of restaurants, shop at the outlets and specialty stores, relax on the beach, or take a stroll through the historic district. There's something for everyone in Atlantic City!	91.46
20	There are many attractions and activities in Atlantic City, including: 1. Casinos: Atlantic City is known for its many casinos, where you can try your luck at slots, poker, blackjack, and more. 2. Boardwalk: ... 3. Beach: ... 4. Steel Pier: ... 5. Shows and concerts: ... 6. Nightlife: ... 7. Dining: ... 8. Shopping: ...	180.14

Our APOHF is able to further refine the response of an LLM to make it more preferable for human users.



Q & A

- Any questions?

THANK YOU