



# ICLR

International Conference On  
Learning Representations

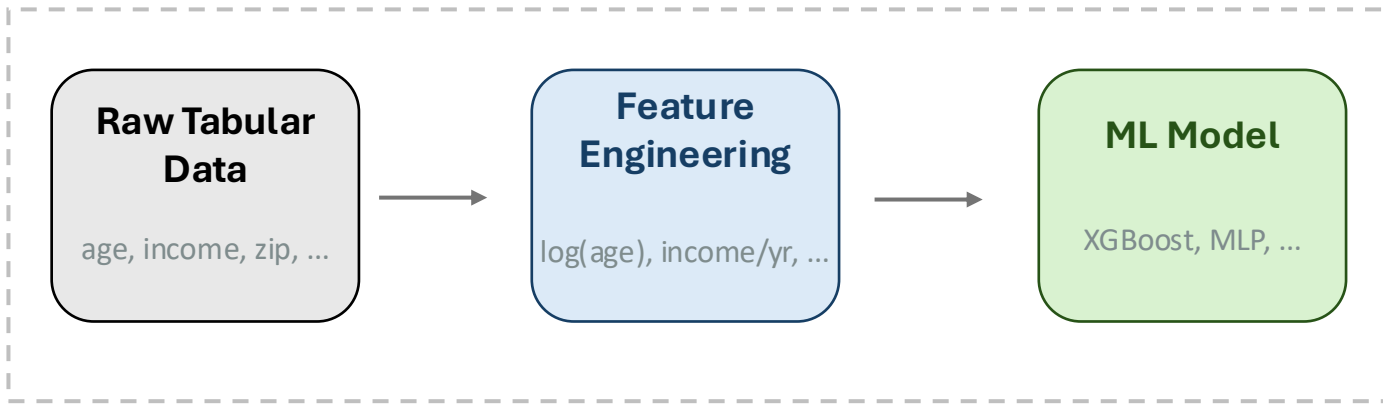
## Human-LLM Collaborative Feature Engineering for Tabular Learning

Zhuoyan Li<sup>1</sup>, Aditya Bansal<sup>2</sup>, Jinzhao Li<sup>1</sup>, Shishuang He<sup>3</sup>, Zhuoran Lu<sup>1</sup>, Mutian Zhang<sup>1</sup>,  
Yiwei Yang<sup>4</sup>, Qin Liu<sup>5</sup>, Swati Jain<sup>2</sup>, Ming Yin<sup>1</sup>, Yunyao Li<sup>2</sup>

PurdueUniversity<sup>1</sup>, Adobe<sup>2</sup>, UIUC<sup>3</sup>, University of Washington<sup>4</sup>, UC Davis<sup>5</sup>



# Feature Engineering for Tabular Data



LLMs can propose creative feature transformations by leveraging semantic understanding of the task and data.

E.g., CAAFE (Hollmann et al., 2023), OCTree (Nam et al., 2024)

## The Problem

Current methods assign the LLM both roles:

1. **Proposing** candidate features
2. **Selecting** which one to evaluate

The LLM acts as a **black-box optimizer**, relying entirely on internal heuristics for both generation and selection.

# Why is this a **problem**?

**1**

## **No Calibrated Uncertainty**

LLM cannot quantify how promising or uncertain a candidate is

**2**

## **Wasted Exploration**

Repeated low-yield operations without principled search strategy

**3**

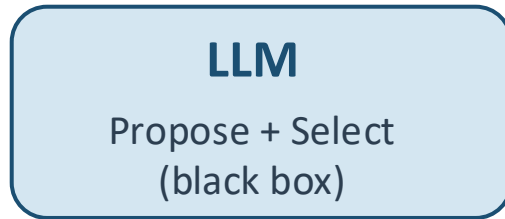
## **No Human Steerability**

No mechanism to incorporate expert domain knowledge

→ Suboptimal feature engineering, especially when iteration budgets are limited.

# Our Approach: Decouple Proposal from Selection

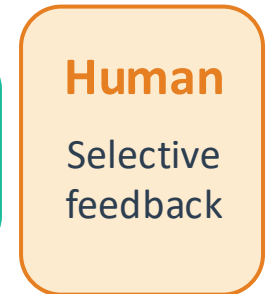
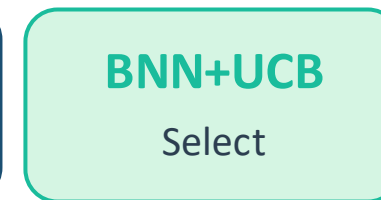
## Previous Approaches



- Single model handles both roles
- No uncertainty quantification
- No human interaction mechanism

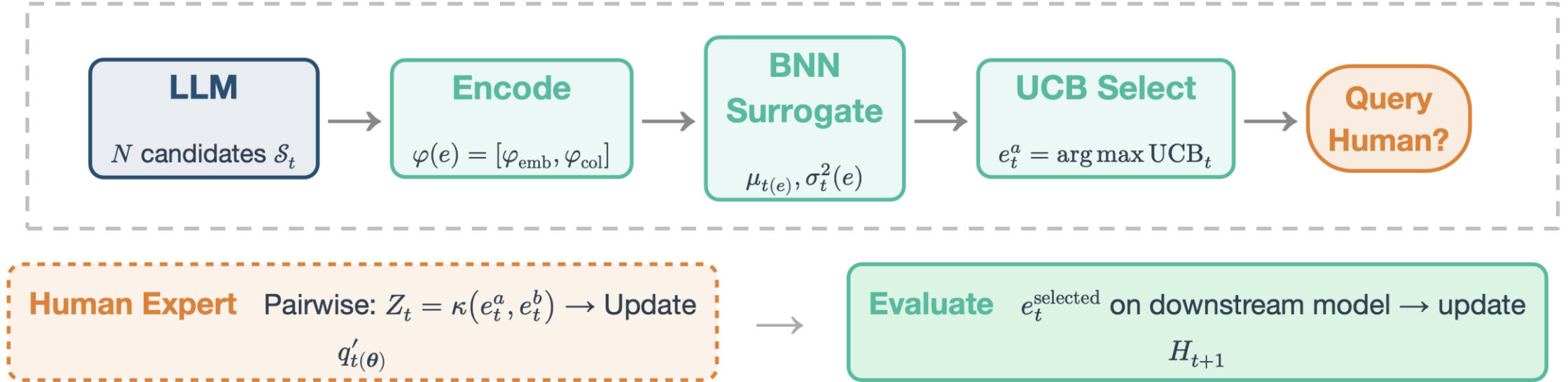


## Our Framework



- Principled Bayesian selection
- Calibrated uncertainty estimates
- Selective human preference queries

# Framework Overview



$$\text{UCB}_{t(e)} = \mu_{t(e)} + \sqrt{\beta_t} \cdot \sigma_{t(e)} \quad \text{Balances exploitation (high } \mu) \text{ with exploration (high } \sigma)$$

**Lemma 3.1:** With probability  $\geq 1 - \delta$ ,  $|g(e) - \mu_{t(e)}| \leq \sqrt{\beta_t} \cdot \sigma_{t(e)}$  for all  $e \in \mathcal{S}_t$

# Surrogate Model: Encoding & BNN

## Feature Operation Encoding

$$\varphi(e) = [\varphi_{\text{embedding}}(e), \varphi_{\text{column}}(e)]$$

$$\varphi_{\text{embedding}}(e)$$

Semantic text embedding via  
text-embedding-3-small

Captures *what* the operation  
does

$$\varphi_{\text{column}}(e)$$

Binary column-usage vector  
 $m \in \{0, 1\}^d$

Captures *which* columns are  
used

*Why both?* Semantic embeddings alone can't distinguish  $\log(\text{age})$  from  $\log(\text{income})$  — column mask resolves this.

## Bayesian Neural Network Surrogate

### Why BNN over Gaussian Process?

- GP struggles with high-dimensional language-derived embeddings
- BNN provides greater scalability and expressiveness
- Naturally yields calibrated uncertainty via variational inference

### Variational Inference

Posterior:  $q_t(\theta) = \mathcal{N}(\theta; M_t, \Sigma_t)$

Predictions for candidate  $e$ :

$$\mu_t(e) = \mathbb{E}_{q_t}[\hat{g}(\varphi(e); \theta)] \quad (\text{expected utility})$$

$$\sigma_t^2(e) = \text{Var}_{q_t}[\hat{g}(\varphi(e); \theta)] \quad (\text{uncertainty})$$

Unlike a regular NN, BNN learns a distribution over weights → naturally provides calibrated uncertainty.

# Selective Human Preference Feedback

Human feedback elicited **only** when both conditions hold:

**C1 — Overlap**  $UCB_{t(e_t^b)} > LCB_{t(e_t^a)}$

Confidence intervals overlap — outcome uncertain

**C2 — Uncertainty  $\geq$  Cost**

$$\sqrt{\beta_t}(\sigma_{t(e_t^a)} + \sigma_{t(e_t^b)}) \geq \gamma_\kappa$$

Potential gain justifies cognitive cost

**Posterior update via probit model:**

$$\mathcal{P}(Z_t | \theta) = \Phi(\eta Z_t [\hat{g}(e_t^a) - \hat{g}(e_t^b)])$$

**When to ask the human?**

**No overlap → don't ask**

$e_t^a$ : | - [====] - |

$e_t^b$ : | - [====] - |

**Overlap → ask!**

$e_t^a$ : | - [=====] - |

$e_t^b$ : | - [=====] - |

*Overlap = feedback can change decision*

# Key Results: Classification (AUROC %)

Dataset	OpenFE	AutoGl.	CAAFE	OCTree	Ours (w/o)	Ours (w/)
<b>MLP</b>						
flight	93.3	92.6	92.9	94.8	<b>96.9</b>	<b>97.3</b>
loan	95.3	95.4	95.7	95.9	<b>96.0</b>	<b>96.1</b>
conversion	90.7	90.6	90.9	91.1	<b>92.6</b>	<b>92.9</b>
heart	92.2	92.3	92.6	93.1	<b>93.4</b>	<b>93.6</b>
<b>XGBoost</b>						
flight	95.7	95.4	95.2	96.4	<b>97.6</b>	<b>98.0</b>
conversion	91.2	91.9	92.1	92.4	<b>93.5</b>	<b>93.9</b>
heart	93.6	93.5	93.6	94.3	<b>95.1</b>	<b>94.8</b>

GPT-4o backbone. 18 datasets total, representative subset shown

**Avg Error Reduction vs. Best Baseline**

**MLP**

7.24% w/o Human

8.96% w/ Human

**XGBoost**

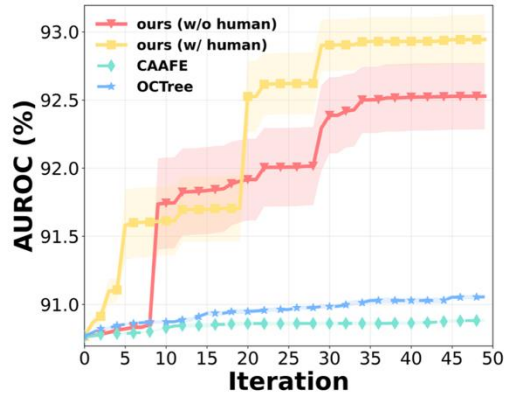
9.02% w/o Human

11.23% w/ Human

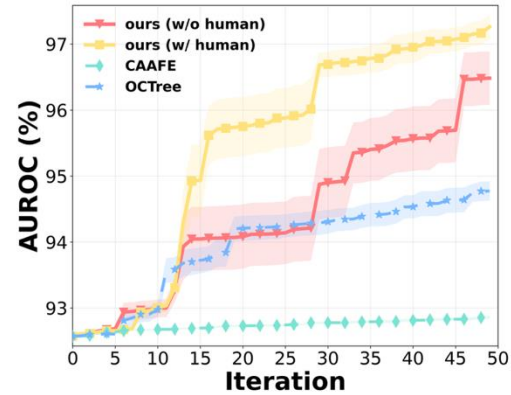
Consistent across 4 backbone LLMs (DeepSeek-v3, GPT-3.5, GPT-4o, GPT-5)

# Performance Trajectories Over Iterations

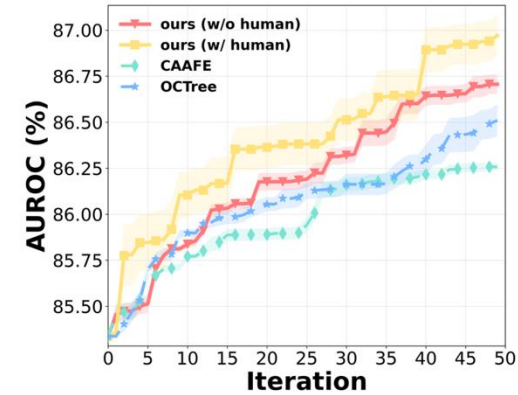
Our method (w/ and w/o human) steadily improves while baselines (CAAFE, OCTree) plateau. 50 iterations, MLP



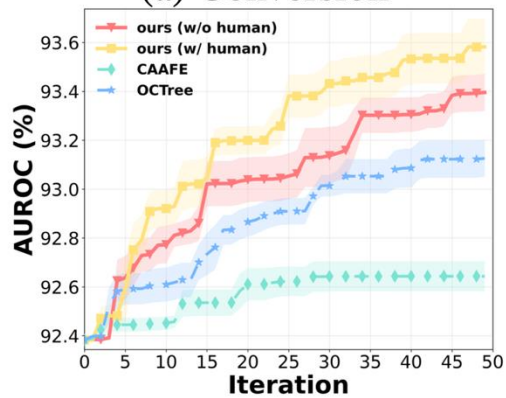
(a) Conversion



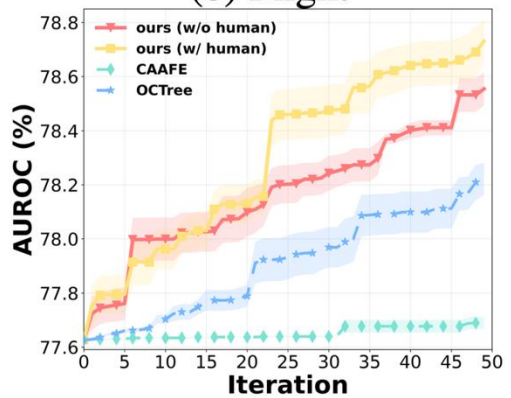
(b) Flight



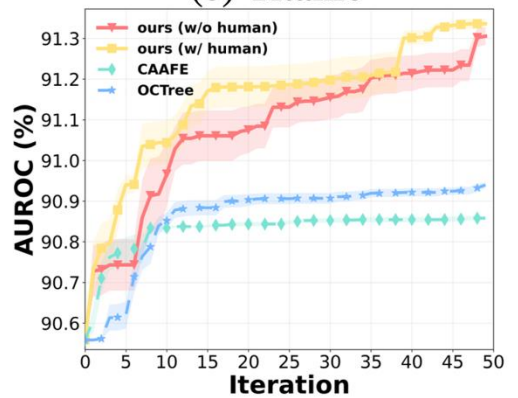
(c) Titanic



(d) Heart



(e) Wine



(f) Adult

# Computational Scalability

Our BNN surrogate and UCB computation add minimal overhead to the LLM-based pipeline.

**Runtime (s) vs. # Initial Features** (10K instances)

Features	LLM	Surrogate	UCB	Eval
10	1.82	0.17	0.006	1.79
100	1.82	0.19	0.005	1.78
1,000	1.82	0.20	0.009	8.4
10,000	1.82	0.57	0.018	23.4

**Runtime (s) vs. # Instances** (100 features)

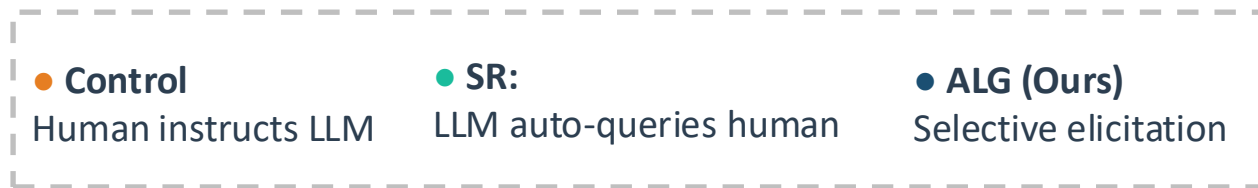
Samples	LLM	Surrogate	UCB	Eval
1,000	1.82	0.17	0.005	0.28
10,000	1.82	0.23	0.006	1.47
50,000	1.82	0.18	0.006	5.22
100,000	1.82	0.18	0.005	10.65

**Surrogate + UCB = only 2.2% of total runtime** (at 10K features). Our overhead operates at the feature-operation level and is independent of dataset size. Dominant cost is downstream model evaluation (shared across all methods).

# User Study: 31 ML Practitioners

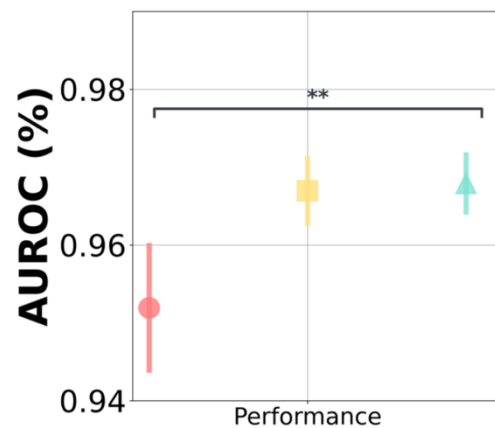
Task: Flight Satisfaction Prediction

3 conditions:

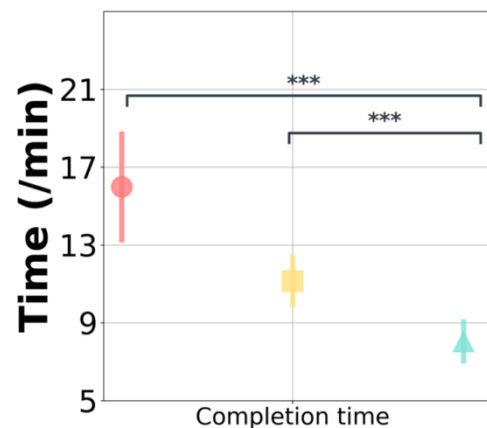


**ALG (Ours) achieves:**

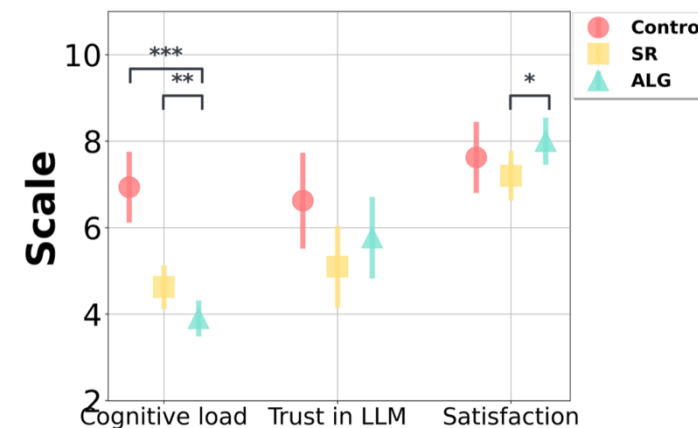
- ❖ Higher performance ( $p = 0.001$ )
- ❖ Faster completion ( $p < 0.001$ )
- ❖ Lower cognitive load ( $p < 0.001$ )
- ❖ Higher satisfaction ( $p = 0.072$ )



(a) Performance



(b) Completion time



(c) User experience

# Summary

1. **Decouple** LLM feature proposal from principled Bayesian selection with BNN surrogate
2. **Selective human preference** feedback with theoretical guarantees
3. Consistent improvement across **18 datasets** and **4 backbone** LLMs
4. Real **user study** validates better performance, faster completion, lower cognitive load

**Thank you!**