

# Let The Jury Decide: Fair Demonstration Selection for In- Context Learning through Incremental Greedy Evaluation

Sadaf MD Halim, Chen Zhao, Xintao Wu, Latifur Khan, Christian Earl  
Grant, Fariha Ishrat Rahman, Feng Chen

# Fair In-Context Learning

Pretrained LLMs often encode demographic, societal, or linguistic preferences.

These issues can cause:

- Toxicity
- Stereotypical completions
- Irregularities in classification tasks

# In-Context Learning (ICL)

- ICL enables models to learn from a few examples at inference time.
- No fine-tuning required – examples are provided as part of the input prompt.

Used in:

- - Few-shot classification
  - - Question answering
  - - Summarization
- 
- Input Prompt = {Demo 1, Demo 2, ..., Demo k} + Test Query → Model Prediction

# Why Demonstration Selection Is Critical

- Demonstrations are the only supervision LLMs receive during inference.
- Poor selection can amplify irresponsible outputs and reduce accuracy.

Impacts of demo selection:

- Prediction accuracy
- Generalization across subgroups
- Unfair outcomes for certain groups

# Responsible In-Context Classification

Goal: Maximize metrics like Demographic Parity, Equalized Odds 

Challenges:

- - Select  $k$  (typically 5 or 10) from  $n$  demonstrations (where  $n$  is a large number)
- - Ensure responsible outputs
- - Retain predictive utility (accuracy, F1)

# JUDGE (**JU**ry-based **D**emonstration Selection via **G**reedy **E**valuation)

Our approach, JUDGE addresses demonstration selection through a multi-step process.

1. **Jury Set Selection:** Creating a set of examples for greedy evaluation.
2. **Candidate Pruning:** Reducing the pool of candidates to a much smaller pool.
3. **Iterative Greedy Selection:** Building the final demonstration set greedily by using performance on the jury set as a heuristic.

# The Jury Set, $\mathcal{J}$

A carefully constructed group of examples, providing a balanced representation across all combinations of groups and labels.

$$\mathcal{C} = \{(g, y) : g \in \mathcal{G}, y \in \mathcal{Y}\}$$

Each subset  $\mathcal{J}_{(g, y)}$  consists of  $|\mathcal{J}| / |\mathcal{C}|$  examples

Uses SentenceBERT embeddings and the cosine similarity measure:

$$\text{sim}(x_i, x_j) = \frac{e(x_i) \cdot e(x_j)}{\|e(x_i)\| \|e(x_j)\|}$$

# The Jury Set, $\mathcal{J}$

Each example is chosen such that it maximizes distance from existing examples.

$$\mathcal{J}_{g,y} = \{x_1, \dots, x_m\} \text{ where}$$
$$x_i = \arg \min_{x \in \mathcal{D}_{g,y} \setminus \{x_1, \dots, x_{i-1}\}} \max_{j < i} \text{sim}(x, x_j)$$

Finally, we have:

$$\mathcal{J} = \bigcup_{(g,y) \in \mathcal{C}} \mathcal{J}_{g,y}$$



# Candidate Pruning

We similarly prune the space of candidates (down to ~3%)

$$\mathcal{D}_{reduced} = \{x_1, \dots, x_n\} \text{ where}$$
$$x_i = \arg \min_{x \in \mathcal{D}_{candidate} \setminus \{x_1, \dots, x_{i-1}\}} \max_{j < i} \text{sim}(x, x_j)$$

# Objective and Score Functions

Our objective provides a balance between performance in terms of accuracy (denoted by  $a$ ) and metrics like Demographic Parity (denoted by  $f$ )

$$\mathcal{S}^* = \operatorname{argmax}_{\mathcal{S} \subseteq \mathcal{D}_{reduced}, |\mathcal{S}|=k} \operatorname{score}(\mathcal{S}, \mathcal{J})$$

$$\operatorname{score}(\mathcal{S}, \mathcal{J}) = \omega \cdot f(\mathcal{S}, \mathcal{J}) + (1 - \omega) \cdot a(\mathcal{S}, \mathcal{J})$$

# Greedy Selection

1. We start with the empty set,  $S_0 = \emptyset$
2. Select the first example that maximizes score on the jury set.

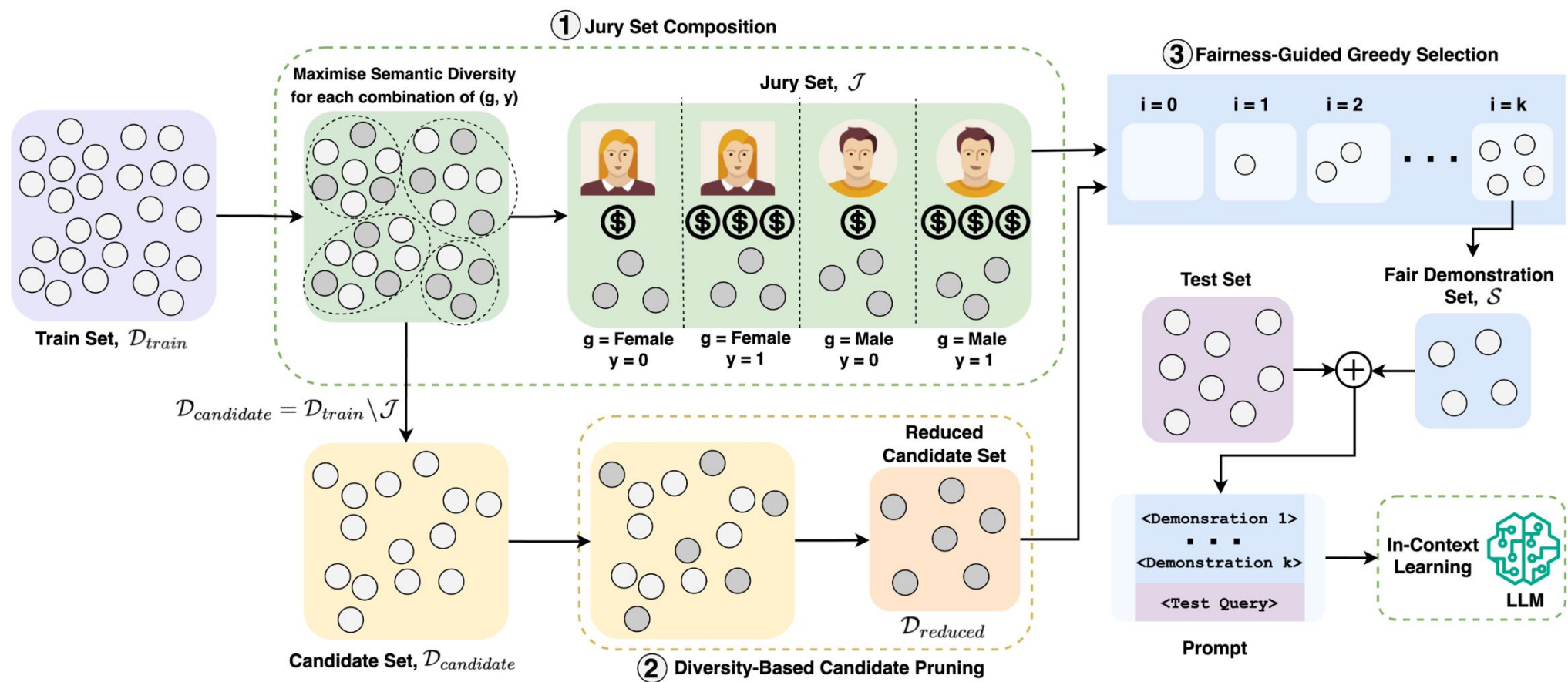
$$x_1 = \arg \max_{x \in \mathcal{D}_{\text{reduced}}} \text{score}(\{x\}, \mathcal{J})$$

1. Add the selected example to the demonstration set.  $S_1 = \{x_1\}$
2. At each step  $t$ , select the candidate that maximizes score when added to the current set.

$$x_t = \arg \max_{x \in \mathcal{D}_{\text{reduced}} - S_{t-1}} \text{score}(S_{t-1} \cup \{x\}, \mathcal{J})$$

1. Update the selected set with the newly chosen example.  $S_t = S_{t-1} \cup \{x_t\}$
2. Repeat the process until  $k$  examples are selected.  $|S_t| = k$

# Overview



# Datasets and Baselines

## Datasets:

Adult  
COMPAS  
Law School  
ACS-Income

## Baselines:

Random  
Balanced  
Counterfactual  
Instruct  
FairlCL  
FCG  
FADS

# Results

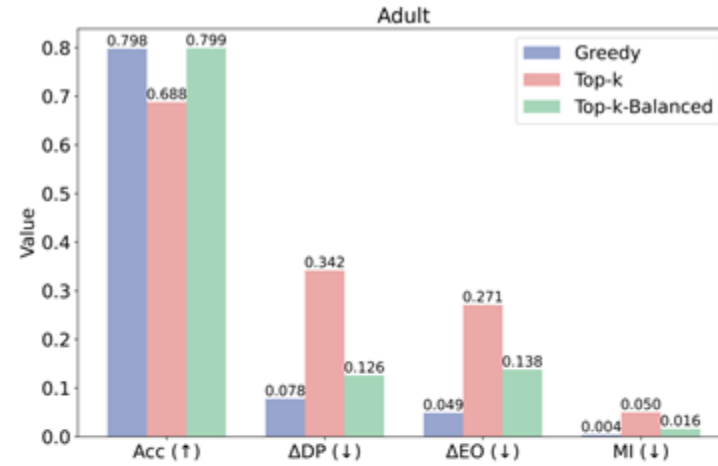
Table 1: Results for Adult with 5 demonstrations, across 4 LLMs. Each cell shows  $Mean_{S.D.}$ .

	Method	Acc. $\uparrow$	$\Delta DP \downarrow$	$\Delta EO \downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.772 <sub>0.008</sub>	0.185 <sub>0.004</sub>	0.191 <sub>0.006</sub>	0.023 <sub>0.002</sub>
	Balanced	0.706 <sub>0.015</sub>	0.216 <sub>0.011</sub>	0.146 <sub>0.014</sub>	0.022 <sub>0.001</sub>
	Cfact.	0.731 <sub>0.017</sub>	0.185 <sub>0.019</sub>	0.158 <sub>0.023</sub>	0.018 <sub>0.003</sub>
	Instruct	0.753 <sub>0.013</sub>	0.299 <sub>0.011</sub>	0.308 <sub>0.012</sub>	0.052 <sub>0.006</sub>
	FairICL	0.764 <sub>0.009</sub>	0.170 <sub>0.004</sub>	0.097 <sub>0.008</sub>	0.016 <sub>0.002</sub>
	FCG	0.795 <sub>0.011</sub>	0.097 <sub>0.009</sub>	0.157 <sub>0.006</sub>	0.011 <sub>0.001</sub>
	FADS	0.743 <sub>0.015</sub>	0.157 <sub>0.012</sub>	0.114 <sub>0.014</sub>	0.019 <sub>0.003</sub>
	JUDGE	<b>0.798</b> <sub>0.012</sub>	<b>0.078</b> <sub>0.011</sub>	<b>0.049</b> <sub>0.012</sub>	<b>0.004</b> <sub>0.001</sub>
MISTRAL-7B	Random	0.709 <sub>0.013</sub>	0.201 <sub>0.010</sub>	0.124 <sub>0.009</sub>	0.019 <sub>0.003</sub>
	Balanced	0.594 <sub>0.014</sub>	0.230 <sub>0.011</sub>	0.185 <sub>0.012</sub>	0.025 <sub>0.004</sub>
	Cfact.	0.722 <sub>0.011</sub>	0.143 <sub>0.008</sub>	0.193 <sub>0.013</sub>	0.011 <sub>0.003</sub>
	Instruct	0.729 <sub>0.021</sub>	0.162 <sub>0.019</sub>	0.171 <sub>0.023</sub>	0.015 <sub>0.004</sub>
	FairICL	0.761 <sub>0.006</sub>	0.151 <sub>0.011</sub>	0.159 <sub>0.007</sub>	0.012 <sub>0.002</sub>
	FCG	0.752 <sub>0.015</sub>	0.132 <sub>0.014</sub>	0.093 <sub>0.019</sub>	<b>0.006</b> <sub>0.001</sub>
	FADS	<b>0.769</b> <sub>0.009</sub>	0.180 <sub>0.008</sub>	0.129 <sub>0.005</sub>	0.021 <sub>0.002</sub>
	JUDGE	0.767 <sub>0.012</sub>	<b>0.101</b> <sub>0.009</sub>	<b>0.024</b> <sub>0.005</sub>	<b>0.006</b> <sub>0.001</sub>
GEMMA-2-9B	Random	0.754 <sub>0.006</sub>	0.394 <sub>0.008</sub>	0.423 <sub>0.013</sub>	0.091 <sub>0.005</sub>
	Balanced	0.701 <sub>0.014</sub>	0.482 <sub>0.023</sub>	0.413 <sub>0.026</sub>	0.113 <sub>0.021</sub>
	Cfact.	0.752 <sub>0.015</sub>	0.311 <sub>0.015</sub>	0.372 <sub>0.011</sub>	0.087 <sub>0.016</sub>
	Instruct	0.742 <sub>0.011</sub>	0.428 <sub>0.009</sub>	0.479 <sub>0.013</sub>	0.108 <sub>0.008</sub>
	FairICL	0.753 <sub>0.014</sub>	0.318 <sub>0.019</sub>	0.392 <sub>0.026</sub>	0.089 <sub>0.013</sub>
	FCG	0.755 <sub>0.017</sub>	0.233 <sub>0.025</sub>	0.192 <sub>0.018</sub>	<b>0.013</b> <sub>0.003</sub>
	FADS	0.759 <sub>0.013</sub>	0.353 <sub>0.011</sub>	0.387 <sub>0.016</sub>	0.072 <sub>0.006</sub>
	JUDGE	<b>0.769</b> <sub>0.012</sub>	<b>0.177</b> <sub>0.018</sub>	<b>0.101</b> <sub>0.009</sub>	0.018 <sub>0.003</sub>
QWEN-2.5-32B	Random	0.745 <sub>0.012</sub>	0.215 <sub>0.010</sub>	0.132 <sub>0.010</sub>	0.023 <sub>0.004</sub>
	Balanced	0.708 <sub>0.014</sub>	0.245 <sub>0.013</sub>	0.165 <sub>0.012</sub>	0.027 <sub>0.003</sub>
	Cfact.	0.748 <sub>0.014</sub>	0.225 <sub>0.014</sub>	0.143 <sub>0.011</sub>	0.025 <sub>0.003</sub>
	Instruct	0.733 <sub>0.007</sub>	0.239 <sub>0.013</sub>	0.161 <sub>0.009</sub>	0.026 <sub>0.005</sub>
	FairICL	0.743 <sub>0.009</sub>	0.192 <sub>0.012</sub>	0.147 <sub>0.015</sub>	0.027 <sub>0.009</sub>
	FCG	0.762 <sub>0.013</sub>	0.111 <sub>0.014</sub>	0.098 <sub>0.013</sub>	0.007 <sub>0.002</sub>
	FADS	0.712 <sub>0.009</sub>	0.220 <sub>0.007</sub>	0.141 <sub>0.006</sub>	0.023 <sub>0.003</sub>
	JUDGE	<b>0.771</b> <sub>0.008</sub>	<b>0.096</b> <sub>0.005</sub>	<b>0.062</b> <sub>0.004</sub>	<b>0.005</b> <sub>0.001</sub>

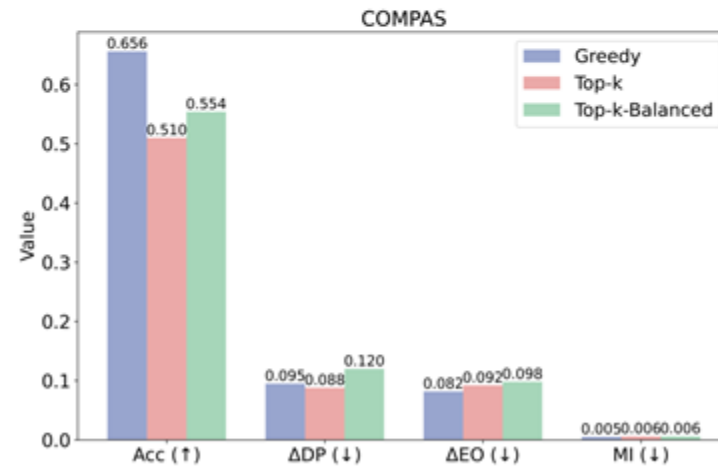
Table 2: Results for COMPAS with 5 demonstrations, across 4 LLMs. Each cell shows  $Mean_{S.D.}$ .

	Method	Acc. $\uparrow$	$\Delta DP \downarrow$	$\Delta EO \downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.617 <sub>0.011</sub>	0.209 <sub>0.009</sub>	0.199 <sub>0.008</sub>	0.021 <sub>0.003</sub>
	Balanced	0.620 <sub>0.012</sub>	0.235 <sub>0.011</sub>	0.218 <sub>0.013</sub>	0.027 <sub>0.002</sub>
	Cfact.	0.582 <sub>0.009</sub>	0.187 <sub>0.006</sub>	0.193 <sub>0.007</sub>	0.017 <sub>0.001</sub>
	Instruct	0.566 <sub>0.010</sub>	0.135 <sub>0.009</sub>	0.164 <sub>0.010</sub>	0.015 <sub>0.001</sub>
	FairICL	0.621 <sub>0.009</sub>	0.192 <sub>0.007</sub>	0.188 <sub>0.006</sub>	0.020 <sub>0.002</sub>
	FCG	0.614 <sub>0.007</sub>	0.182 <sub>0.005</sub>	0.197 <sub>0.005</sub>	0.019 <sub>0.001</sub>
	FADS	0.575 <sub>0.008</sub>	0.167 <sub>0.006</sub>	0.160 <sub>0.005</sub>	0.014 <sub>0.002</sub>
	JUDGE	<b>0.656</b> <sub>0.010</sub>	<b>0.105</b> <sub>0.008</sub>	<b>0.082</b> <sub>0.007</sub>	<b>0.006</b> <sub>0.001</sub>
MISTRAL-7B	Random	0.513 <sub>0.012</sub>	0.097 <sub>0.008</sub>	0.120 <sub>0.009</sub>	0.016 <sub>0.002</sub>
	Balanced	0.512 <sub>0.007</sub>	0.079 <sub>0.005</sub>	0.083 <sub>0.004</sub>	0.013 <sub>0.003</sub>
	Cfact.	0.487 <sub>0.010</sub>	0.059 <sub>0.009</sub>	<b>0.062</b> <sub>0.009</sub>	0.015 <sub>0.004</sub>
	Instruct	0.497 <sub>0.012</sub>	0.082 <sub>0.010</sub>	0.105 <sub>0.008</sub>	0.014 <sub>0.002</sub>
	FairICL	0.515 <sub>0.006</sub>	0.082 <sub>0.005</sub>	0.098 <sub>0.005</sub>	0.017 <sub>0.004</sub>
	FCG	0.489 <sub>0.009</sub>	0.074 <sub>0.004</sub>	0.108 <sub>0.006</sub>	0.013 <sub>0.003</sub>
	FADS	0.531 <sub>0.010</sub>	0.091 <sub>0.005</sub>	0.117 <sub>0.007</sub>	0.015 <sub>0.009</sub>
	JUDGE	<b>0.541</b> <sub>0.007</sub>	<b>0.055</b> <sub>0.004</sub>	0.075 <sub>0.004</sub>	<b>0.002</b> <sub>0.000</sub>
GEMMA-2-9B	Random	0.615 <sub>0.008</sub>	0.310 <sub>0.005</sub>	0.314 <sub>0.006</sub>	0.049 <sub>0.003</sub>
	Balanced	0.601 <sub>0.009</sub>	0.359 <sub>0.006</sub>	0.348 <sub>0.005</sub>	0.067 <sub>0.004</sub>
	Cfact.	0.604 <sub>0.007</sub>	0.261 <sub>0.004</sub>	0.272 <sub>0.005</sub>	0.044 <sub>0.005</sub>
	Instruct	0.609 <sub>0.011</sub>	0.291 <sub>0.009</sub>	0.309 <sub>0.012</sub>	0.047 <sub>0.006</sub>
	FairICL	0.622 <sub>0.010</sub>	0.265 <sub>0.011</sub>	0.282 <sub>0.012</sub>	0.040 <sub>0.005</sub>
	FCG	0.648 <sub>0.007</sub>	0.099 <sub>0.003</sub>	0.091 <sub>0.005</sub>	0.008 <sub>0.003</sub>
	FADS	0.621 <sub>0.014</sub>	0.307 <sub>0.011</sub>	0.303 <sub>0.09</sub>	0.053 <sub>0.009</sub>
	JUDGE	<b>0.665</b> <sub>0.006</sub>	<b>0.062</b> <sub>0.002</sub>	<b>0.039</b> <sub>0.003</sub>	<b>0.002</b> <sub>0.000</sub>
QWEN-2.5-32B	Random	0.637 <sub>0.007</sub>	0.242 <sub>0.005</sub>	0.221 <sub>0.006</sub>	0.029 <sub>0.003</sub>
	Balanced	<b>0.652</b> <sub>0.008</sub>	0.248 <sub>0.007</sub>	0.240 <sub>0.011</sub>	0.031 <sub>0.005</sub>
	Cfact.	0.611 <sub>0.008</sub>	0.244 <sub>0.006</sub>	0.228 <sub>0.006</sub>	0.031 <sub>0.004</sub>
	Instruct	0.633 <sub>0.006</sub>	0.234 <sub>0.003</sub>	0.214 <sub>0.004</sub>	0.026 <sub>0.002</sub>
	FairICL	0.639 <sub>0.008</sub>	0.211 <sub>0.005</sub>	0.218 <sub>0.005</sub>	0.025 <sub>0.003</sub>
	FCG	0.623 <sub>0.006</sub>	0.149 <sub>0.004</sub>	0.144 <sub>0.003</sub>	0.018 <sub>0.003</sub>
	FADS	0.645 <sub>0.008</sub>	0.224 <sub>0.006</sub>	0.207 <sub>0.004</sub>	0.025 <sub>0.003</sub>
	JUDGE	0.649 <sub>0.004</sub>	<b>0.138</b> <sub>0.005</sub>	<b>0.134</b> <sub>0.003</sub>	<b>0.010</b> <sub>0.001</sub>

# Greedy vs Top-K



(a) Adult dataset



(b) COMPAS dataset

# Effect of Jury Set Size

