

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

Sadaf MD Halim<sup>†</sup> Chen Zhao<sup>‡</sup> Xintao Wu<sup>§</sup> Latifur Khan<sup>†</sup>  
Christan Earl Grant<sup>||</sup> Fariha Ishrat Rahman<sup>†</sup> Feng Chen<sup>†</sup>

<sup>†</sup>Department of Computer Science, The University of Texas at Dallas, Richardson, Texas

<sup>‡</sup>Department of Computer Science, Baylor University, Waco, Texas

<sup>§</sup>Department of Electrical Eng. & Comp. Sci., University of Arkansas, Fayetteville, Arkansas

<sup>||</sup>Dept. of Computer & Info. Sci. & Eng., University of Florida, Gainesville, Florida  
sxh190015@utdallas.edu

## Abstract

Large Language Models (LLMs) are powerful in-context learners, achieving strong performance with just a few high-quality demonstrations. However, fairness concerns arise in many in-context classification tasks, especially when predictions involve sensitive attributes. To address this, we propose **JUDGE**—a simple yet effective framework for selecting fair and representative demonstrations that improve group fairness in In-Context Learning. JUDGE constructs the demonstration set iteratively using a greedy approach, guided by a small, carefully selected *jury* set. Our method remains robust across varying LLM architectures and datasets, ensuring consistent fairness improvements. We evaluate JUDGE on *four* datasets using *four* LLMs, comparing it against *seven* baselines. Results show that JUDGE consistently improves fairness metrics without compromising accuracy.

## 1 Introduction

A key capability of Large Language Models (LLMs) is in-context learning (ICL) — the ability to learn from examples provided within a prompt, without requiring parameter updates (Brown et al., 2020; Dong et al., 2022). While research has advanced our understanding of ICL and techniques to enhance its effectiveness, a critical open question remains: how should we select fair and representative demonstration examples? This question becomes particularly critical in high-stakes domains where predictions directly impact human lives.

Consider a parole board using an LLM to assess recidivism risk. The model’s predictions are shaped by the examples it is shown—if those examples reflect historical biases or overlook key rehabilitation factors, the system may produce plausible-looking predictions that perpetuate or amplify existing disparities. In sensitive domains like criminal justice, healthcare, and hiring, the selection of demonstra-

tions directly influences both predictive reliability and equitable decision-making.

Existing demonstration selection strategies, with a few exceptions, largely focus on optimizing performance metrics such as accuracy (Peng et al., 2024; Wu et al., 2023). While these methods are effective for improving ICL performance, they often fail to account for fairness concerns. Parallel research has explored bias and fairness in LLM outputs (Gallegos et al., 2024) and their trustworthiness (Huang et al., 2024), but a key gap remains: how can we improve *group fairness* directly at the demonstration selection stage in in-context learning? Unlike prior work that dynamically selects demonstrations per test query (Wang et al., 2024), we explore an alternative: constructing a *single* demonstration set for an entire classification task.

In our work, we investigate several key questions. **Do different LLMs exhibit consistent fairness behavior across datasets?** We find significant variations in fairness outcomes across different LLMs, highlighting the need for adaptive approaches rather than one-size-fits-all solutions. **Do existing demonstration selection methods generalize across LLM architectures and datasets?** Our results reveal that it is challenging for some prior methods to maintain stable fairness improvements across different models due to inherent variability in LLM responses. **Can we design an effective, fairness-aware demonstration selection approach?** We propose a simple yet highly effective method, **JUDGE** (**JU**ry-based **D**emonstration Selection via **G**reedy **E**valuation) that leverages each LLM’s own predictions on a carefully curated set of *jury* examples to guide demonstration selection.

This paper makes several contributions. First, we provide a comprehensive analysis of existing approaches for fairness-aware demonstration selection across multiple datasets and architectures.

Second, we present **JUDGE**<sup>1</sup>, a consistent and efficient framework for extracting fair representative examples from large datasets for ICL. Third, we validate our approach through extensive empirical evaluation, showing significant fairness improvements without compromising accuracy across multiple fairness benchmarks. Finally, our systematic analysis demonstrates that the greedy construction approach is crucial for balancing fairness and accuracy, outperforming alternatives such as top-k selection and pooling-based methods. As LLMs continue to be deployed in increasingly sensitive domains, our work provides a practical framework for ensuring fairer outcomes while maintaining the efficiency that makes ICL attractive.

## 2 Preliminaries: Protected Groups, Attributes and Group Fairness

*Protected groups* are demographic subpopulations that should not face disparate treatment in model decisions. Let  $\mathcal{G}$  denote the set of protected groups, each defined by a *protected attribute* such as race, gender, or age. For any instance  $x$  in dataset  $\mathcal{D}$ , its protected group membership is given by  $g(x)$ . The population is partitioned into distinct protected groups  $\mathcal{G} = \{g_1, g_2, \dots, g_l\}$ , with each instance belonging to exactly one. For binary attributes (e.g., gender), this simplifies to  $\mathcal{G} = \{g_1, g_2\}$ .

*Group fairness* aims to ensure that a model’s behavior remains consistent across protected groups by enforcing that certain statistical measures are approximately equal across all protected groups, rather than focusing on individual-level fairness. We consider three popular metrics: **Demographic Parity Difference** ( $\Delta DP$ ), which measures the absolute difference in positive rates between protected groups (Padh et al., 2021); **Equalized Odds Difference** ( $\Delta EO$ ), which measures the absolute difference in true positive and false positive rates between groups (Hardt et al., 2016); and **Mutual Information** (MI) (Kamishima et al., 2012; Anahideh et al., 2022), which quantifies the mutual information between protected attributes and selection decisions (Details in Appendix A).

## 3 Proposed Approach

### 3.1 Problem Formulation: Fairest Prompt Search for In-Context Learning

Given a large language model  $M$  and an input  $x$ , ICL makes predictions by conditioning on a demon-

<sup>1</sup><https://github.com/smdh-hub/judge-icl>

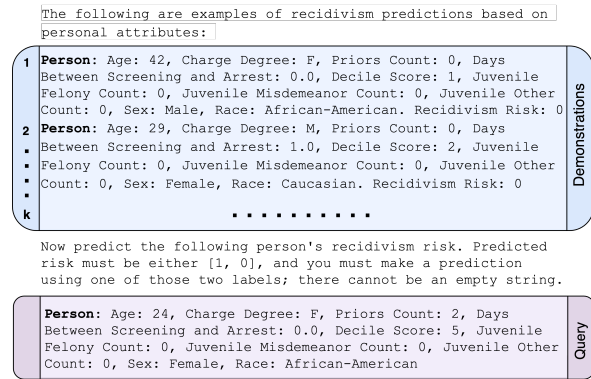


Figure 1: Example ICL Prompt on the COMPAS dataset

stration set  $\mathcal{S} = \{(x_1, y_1), \dots, (x_k, y_k)\}$ . The model processes these demonstrations along with the query input as:

$$\text{prompt}(\mathcal{S}, x) = [(x_1, y_1), \dots, (x_k, y_k), x] \quad (1)$$

Let  $\mathcal{X}$  denote the input space,  $\mathcal{Y}$  the label space, and  $\mathcal{L}$  the natural language space. The formatting function  $\phi$  maps  $k$  input-label pairs and the query to a natural language prompt as seen in Figure 1:

$$\phi: \underbrace{(\mathcal{X} \times \mathcal{Y})^k}_{k \text{ demonstration pairs}} \times \underbrace{\mathcal{X}}_{\text{query input}} \rightarrow \mathcal{L} \quad (2)$$

$$\phi(\text{prompt}(\mathcal{S}, x)) \in \mathcal{L} \quad (3)$$

The model then predicts:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} M(y | \phi(\text{prompt}(\mathcal{S}, x))) \quad (4)$$

where  $M(y | \phi(\text{prompt}(\mathcal{S}, x)))$  represents the model’s predicted probability distribution over the label space  $\mathcal{Y}$ , which we denote for brevity as:

$$\hat{y} = M(\mathcal{S}, x) \quad (5)$$

**The fundamental challenge** is selecting an effective and fair demonstration set  $\mathcal{S}$  from a large candidate pool. For a pool of size  $|\mathcal{D}|$  and desired demonstration set size  $k$ , there are  $\binom{|\mathcal{D}|}{k}$  possible combinations. For even modest values like  $|\mathcal{D}| = 1000$  and  $k = 5$ , this yields over 8 trillion possible demonstration sets, making exhaustive search intractable. Fairness constraints further complicate this selection.

**Our approach**, **JUDGE** addresses fair demonstration selection through a multi-step process as shown in Figure 2. Let:

**Train set  $\mathcal{D}_{train}$** : The pool of available examples, where each example  $x \in \mathcal{D}_{train}$  has associated features, a label  $y \in [0, 1]$ , and protected group membership  $g(x)$ .

**Jury set  $\mathcal{J}$ :** A small curated subset of examples extracted from  $\mathcal{D}_{train}$  that serves to evaluate the fairness and effectiveness of candidate demonstration sets.

**Candidate set  $\mathcal{D}_{candidate}$ :** The complement of the jury set with respect to the train set, defined as  $\mathcal{D}_{candidate} = \mathcal{D}_{train} \setminus \mathcal{J}$ , from which potential demonstrations can be selected.

**Reduced candidate set  $\mathcal{D}_{reduced} \subseteq \mathcal{D}_{candidate}$ :** A pruned subset of the candidate set, selected to maintain semantic diversity while reducing computational complexity. Demonstrations are selected from this subset.

**Protected groups  $\mathcal{G} = \{g_1, g_2, \dots, g_l\}$ :** The set of groups defined by protected attributes, where each example belongs to exactly one group. We consider a binary setting where  $\mathcal{G} = \{g_1, g_2\}$ .

**Selected set  $\mathcal{S} \subseteq \mathcal{D}_{reduced}$ :** The chosen subset of  $k$  examples that will serve as demonstrations, where  $k$  is typically small (e.g., 5-10) due to context length constraints.

Our *objective* is to find a demonstration set  $\mathcal{S}^*$  that optimizes both predictive accuracy ( $a$ ) and fairness ( $f$ ):

$$\mathcal{S}^* = \operatorname{argmax}_{\mathcal{S} \subseteq \mathcal{D}_{reduced}, |\mathcal{S}|=k} \operatorname{score}(\mathcal{S}, \mathcal{J}) \quad (6)$$

$$\operatorname{score}(\mathcal{S}, \mathcal{J}) = \omega \cdot f(\mathcal{S}, \mathcal{J}) + (1 - \omega) \cdot a(\mathcal{S}, \mathcal{J}) \quad (7)$$

The accuracy,  $a$  term measures the model’s predictive performance on the jury set:

$$a(\mathcal{S}, \mathcal{J}) = \frac{1}{|\mathcal{J}|} \sum_{(x,y) \in \mathcal{J}} \mathbb{I}[M(\mathcal{S}, x) = y] \quad (8)$$

For the fairness,  $f(\mathcal{S}, \mathcal{J})$  term, we use the widely used demographic parity difference (detailed in Appendix A) to assess the demonstration set’s fairness using the jury set:

$$f(\mathcal{S}, \mathcal{J}) = -|P(M(\mathcal{S}, x) = 1 \mid g(x) = g_1) - P(M(\mathcal{S}, x) = 1 \mid g(x) = g_2)| \quad (9)$$

Note that we negate the demographic parity difference since lower differences indicate better fairness, allowing both accuracy and fairness terms to be maximized in the same direction in Equation 6.

To summarise, JUDGE consists of three main steps. First, we construct a balanced and diverse jury set  $\mathcal{J}$  which evaluates candidate examples based on both fairness metrics and predictive performance. This jury set is drawn from the training

set and subsequently removed to form the candidate pool. Next, we prune the candidate pool to maximize semantic diversity and limit computational overhead. Finally, we employ a greedy selection algorithm that iteratively builds the demonstration set  $\mathcal{S}$  by adding, at each step, a demonstration from  $\mathcal{D}_{reduced}$  that maximizes the fairness-accuracy objective (Equation 7) over the jury set  $\mathcal{J}$ .

### 3.2 Jury Set Composition

The jury set  $\mathcal{J}$  is carefully constructed to ensure balanced representation across all protected groups and labels. We define all possible group-label pairs as  $\mathcal{C} = \{(g, y) : g \in \mathcal{G}, y \in \mathcal{Y}\}$ . For example, in a binary setting where  $g$  represents gender (Male, Female) and  $y$  represents income level ( $>50k$  as 1,  $\leq 50k$  as 0), we have  $\mathcal{C} = \{(\text{Female}, 0), (\text{Female}, 1), (\text{Male}, 0), (\text{Male}, 1)\}$

Each subset  $\mathcal{J}_{g,y}$  consists of  $m = |\mathcal{J}|/|\mathcal{C}|$  examples, selected to maximize semantic diversity.

For each example  $x$ , we compute an embedding  $e(x)$  using SentenceBERT (Reimers, 2019). We measure the semantic similarity between examples using cosine similarity:

$$\operatorname{sim}(x_i, x_j) = \frac{e(x_i) \cdot e(x_j)}{\|e(x_i)\| \|e(x_j)\|} \quad (10)$$

To construct a diverse subset  $\mathcal{J}_{g,y}$ , we iteratively select the next example  $x_{\text{next}}$  that minimizes its maximum similarity to the previously selected examples.

$$x_{\text{next}} = \operatorname{arg} \min_{x_i \notin \mathcal{J}_{g,y}} \max_{x_j \in \mathcal{J}_{g,y}} \operatorname{sim}(x_i, x_j) \quad (11)$$

Therefore, the subset  $\mathcal{J}_{g,y}$  is calculated as:

$$\mathcal{J}_{g,y} = \{x_1, \dots, x_m\} \text{ where } x_i = \operatorname{arg} \min_{x \in \mathcal{D}_{g,y} \setminus \{x_1, \dots, x_{i-1}\}} \max_{j < i} \operatorname{sim}(x, x_j) \quad (12)$$

where  $\mathcal{D}_{g,y}$  represents the subset of examples in  $\mathcal{D}_{train}$  with protected group  $g$  and label  $y$ . The final jury set is the union of these diverse subsets:

$$\mathcal{J} = \bigcup_{(g,y) \in \mathcal{C}} \mathcal{J}_{g,y} \quad (13)$$

### 3.3 Diversity-Based Candidate Pruning

To efficiently reduce the size of the candidate pool while preserving coverage across the semantic space, we employ a selection strategy based on semantic similarity.

We construct the reduced set  $\mathcal{D}_{reduced}$  iteratively by selecting examples that are maximally distinct from those already chosen, following Sec 3.2. Given a target size  $n$ , the selection is defined as:

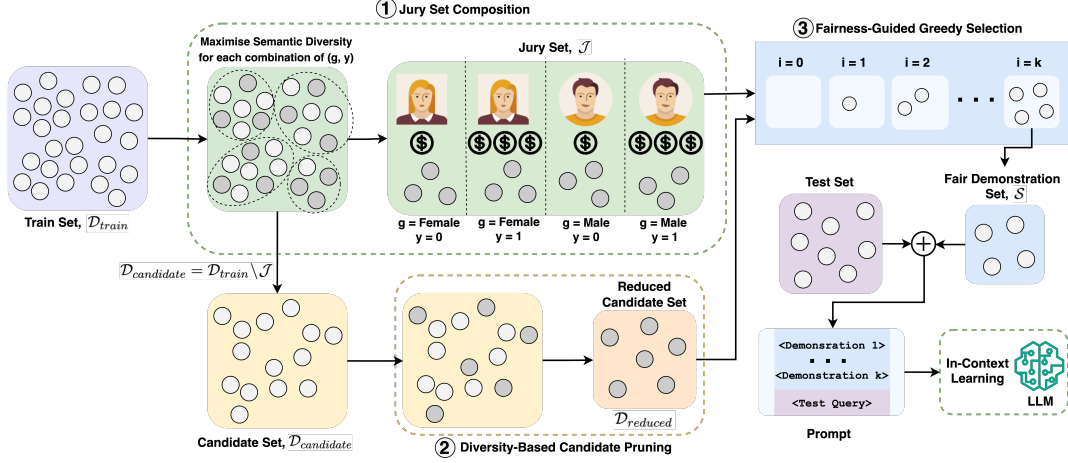


Figure 2: JUDGE consists of three main steps: (1) constructing a balanced and diverse jury set  $\mathcal{J}$  (2) pruning the candidate pool to reduce computational overhead, and (3) iteratively selecting demonstrations using a greedy algorithm that optimizes a weighted combination of fairness and accuracy scores over the jury set.

$$\mathcal{D}_{reduced} = \{x_1, \dots, x_n\} \text{ where}$$

$$x_i = \arg \min_{x \in \mathcal{D}_{candidate} \setminus \{x_1, \dots, x_{i-1}\}} \max_{j < i} \text{sim}(x, x_j) \quad (14)$$

This selection process ensures that the final subset  $\mathcal{D}_{reduced}$  preserves the semantic diversity of the original pool while being computationally tractable for subsequent operations.

### 3.4 Fairness-Guided Greedy Selection

The algorithm constructs the demonstration set  $\mathcal{S}$  iteratively using a greedy selection process, optimizing both fairness and accuracy over the jury set  $\mathcal{J}$ . At each iteration  $t$ , the example that maximizes the marginal improvement in the overall score is added to  $\mathcal{S}$ .

The process starts with an empty set,  $\mathcal{S}_0 = \emptyset$ . At  $t = 1$ , each example in  $\mathcal{D}_{reduced}$  is evaluated independently as the first demonstration, and the one yielding the highest fairness-accuracy score on the jury set is selected as  $x_1$ , and  $\mathcal{S}_1 = \mathcal{S}_0 \cup \{x_1\}$ . At  $t = 2$ , we evaluate each of the remaining candidates in  $\mathcal{D}_{reduced} \setminus \mathcal{S}_1$  in combination with  $x_1$ , forming two-example demonstration sets, selecting  $x_2$  that maximizes the score and  $\mathcal{S}_2 = \mathcal{S}_1 \cup \{x_2\}$ . This process continues until  $t = k$ .

Formally, starting with an empty set  $\mathcal{S}_0 = \emptyset$ , at each iteration  $t$  until  $|\mathcal{S}_t| = k$ , we select:

$$x_t = \operatorname{argmax}_{x \in \mathcal{D}_{reduced} \setminus \mathcal{S}_{t-1}} \text{score}(\mathcal{S}_{t-1} \cup \{x\}, \mathcal{J}) \quad (15)$$

where  $\mathcal{S}_t = \mathcal{S}_{t-1} \cup \{x_t\}$  and score is computed as defined in Section 3.1.

While this greedy approach does not guarantee finding the globally optimal demonstration set, it offers several advantages. First, it reduces search complexity from  $\binom{|\mathcal{D}_{reduced}|}{k}$  to  $O(k|\mathcal{D}_{reduced}|)$ , drastically reducing the search space. Second, it ensures interpretability, as each demonstration is chosen based on a clear improvement metric. Finally, by evaluating candidates based on their marginal contribution, it captures interaction effects, leading to a more effective and fair selection.

Our approach is detailed in Algorithm 1. The pseudocode for the helper function `DiverseSelect`, which is based on the description from Section 3.2, can be found in Algorithm 2 in the Appendix.

### 3.5 On the Greedy Selection Strategy

We employ a greedy algorithm for selecting the demonstration set, primarily due to its computational efficiency and strong empirical performance. While we do not guarantee that our objective function is submodular, the greedy approach is known to have provable approximation guarantees when applied to submodular objective functions (Nemhauser et al., 1978). Inspired by these theoretical underpinnings, and supported by empirical evidence showing that greedy algorithms often perform very well in practice across numerous subset optimization problems (Sener and Savarese, 2018; Wei et al., 2015), we adopted this strategy.

We considered alternative strategies such as beam search, but found them less suitable for our application. Beam search, while potentially yielding better sets, would increase the computational load by a factor of the beam width (Leblond et al.,

---

**Algorithm 1** JUDGE

---

**Require:** Training set  $\mathcal{D}_{train}$ , protected groups  $\mathcal{G}$ , labels  $\mathcal{Y}$ , desired size  $k$ , jury size per group  $m$ , candidate pool size  $n$ , trade-off  $\omega$

**Ensure:** Fair demonstration set  $\mathcal{S}$

```
1: // Step 1: Construct balanced jury set
2:  $\mathcal{C} \leftarrow \{(g, y) : g \in \mathcal{G}, y \in \mathcal{Y}\}$ 
3:  $\mathcal{J} \leftarrow \emptyset$ 
4: for  $(g, y) \in \mathcal{C}$  do
5:    $\mathcal{D}_{g,y} \leftarrow \{x \in \mathcal{D}_{train} : \mathbf{g}(x) = g \wedge \text{label}(x) = y\}$ 
6:    $\mathcal{J}_{g,y} \leftarrow \text{DiverseSelect}(\mathcal{D}_{g,y}, m)$ 
7:    $\mathcal{J} \leftarrow \mathcal{J} \cup \mathcal{J}_{g,y}$ 
8: end for
9: // Step 2: Prune candidate pool
10:  $\mathcal{D}_{reduced} \leftarrow \text{DiverseSelect}(\mathcal{D}_{train} \setminus \mathcal{J}, n)$ 
11: // Step 3: Greedy selection
12:  $\mathcal{S}_0 \leftarrow \emptyset$ 
13: for  $t \leftarrow 1$  to  $k$  do
14:    $x_t \leftarrow \text{None}, s_{\max} \leftarrow -\infty$ 
15:   for  $x \in \mathcal{D}_{reduced} \setminus \mathcal{S}_{t-1}$  do
16:      $\mathcal{S}_{\text{temp}} \leftarrow \mathcal{S}_{t-1} \cup \{x\}$ 
17:      $f \leftarrow f(\mathcal{S}_{\text{temp}}, \mathcal{J}), a \leftarrow a(\mathcal{S}_{\text{temp}}, \mathcal{J})$ 
18:      $s \leftarrow \omega \cdot f + (1 - \omega) \cdot a$ 
19:     if  $s > s_{\max}$  then
20:        $x_t \leftarrow x, s_{\max} \leftarrow s$ 
21:     end if
22:   end for
23:    $\mathcal{S}_t \leftarrow \mathcal{S}_{t-1} \cup \{x_t\}$ 
24: end for
25: return  $\mathcal{S}_k$ 
```

---

2021). This would involve evaluating multiple candidate sets concurrently with the jury set using LLM feedback, rendering the approach prohibitively expensive for users.

## 4 Complexity Analysis

The complexity is dominated by **LLM inference**. In JUDGE, each demonstration in  $\mathcal{D}_{reduced}$  is evaluated with every jury member to compute demographic parity and accuracy. Since this is repeated  $k$  times to build a  $k$ -sized set, the overall complexity is  $O(k \cdot |\mathcal{D}_{reduced}| \cdot |\mathcal{J}|)$ . Unlike our method, **exhaustive search** evaluates all possible subsets of size  $k$  from  $N$  demonstrations, i.e., a complexity of  $O(N^K)$  which is infeasible for large  $N$  and  $k$ . A detailed complexity comparison with other baselines can be found in [Appendix C](#).

## 5 Results

### 5.1 Datasets

We use *four* widely studied fairness datasets across different domains and protected attributes (details in [Appendix B.2](#)). **Adult Income** (Dua and Graff, 2019) to predict whether income exceeds \$50K (protected attribute: *gender*). **COMPAS** (Angwin et al., 2016) to predict recidivism risk (protected attribute: *race*). **Law School (LSAC)** (Wightman,

1998) to predict whether a student passes the bar (protected attribute: *race*). **ACS Income** (Ding et al., 2021) to predict whether income exceeds \$50K (protected attribute: *gender*).

### 5.2 Language Models

To assess the generalizability of JUDGE, we evaluate our approach using *four* open-source language models of varying parameters from different sources: Meta’s **LLaMA-3 8B** (Dubey et al., 2024), Mistral AI’s **Mistral 7B** (Jiang et al., 2023), Google’s **Gemma-2 9B** (Riviere et al., 2024), and Alibaba’s **Qwen-2.5 32B** (Hui et al., 2024).

### 5.3 Baselines

We compare our approach against *seven* baseline methods for demonstration selection. **Random** selects  $k$  demonstrations randomly from the training set. **Balanced** employs stratified random sampling to maintain equal representation across protected groups and label. **Counterfactual** (Li et al., 2023) selects from privileged groups and generates counterfactual examples by flipping sensitive attributes while preserving other features. **Instruct** (Atwood et al., 2024) guides the model toward fairness via explicit prompt instructions. **FCG** (Hu et al., 2024) uses clustering and evolutionary strategies to curate diverse, representative demonstrations while considering fairness metrics. **FairICL** (Bhaila et al., 2024) leverages latent concept variables to evaluate demonstration fairness and guide selection, learning fair concepts from training data to promote fairness while maintaining utility. **FADS** (Wang et al., 2024) implements a two-stage filtering approach (data and model bias mitigation) followed by similarity-based selection with balanced representation across groups and labels. Unlike adaptive methods which select demonstrations per test instance, **JUDGE** selects a *single fixed set* for all test examples. We evaluate our method against both *fixed* and *adaptive* approaches.

### 5.4 Experimental Setup

For each dataset-model combination, we conduct experiments with two demonstration set sizes:  $k = 5$  and  $k = 10$ , using 20% of the data for testing where standard splits are not provided. We evaluate performance using four metrics: Accuracy (Acc.), Demographic Parity Difference ( $\Delta DP$ ), Equalized Odds Difference ( $\Delta EO$ ), and Mutual Information (MI), as defined in Section 2. All results reported in Tables 1-8 show the mean of 3 reproduction runs.

For space constraints, results for  $k = 10$  (Tables 5-8), as well as important additional experiment details can be found in the Appendix.

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>
	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>
	JUDGE	<b>0.769</b> <sub>0.009</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>
	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>
	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>

### 5.4.1 Intrinsic Fairness Differences Among LLMs and Datasets

We note an interesting pattern across our results: *different LLMs report significantly different fairness metrics*. This is evident when examining the *Random* baseline. For instance, with  $k = 5$  on Adult (Table 1), Gemma-2 produces a  $\Delta DP$  score of 0.394, compared to LLaMA-3’s 0.185, more than twice the disparity in demographic parity. These variations persist across datasets, with Gemma-2 often exhibiting greater unfairness, e.g.,  $\Delta DP = 0.310$  on COMPAS, compared to Mistral’s 0.097 (Table 2), over three times the value.

Perhaps less surprisingly, datasets themselves vary in fairness, with the same model reporting very different fairness metrics across different datasets. More interestingly, *certain baselines behave dramatically differently across models*. For example, *Instruct* achieves strong fairness on Law School with Mistral for both  $k = 5$  and 10, yet completely sacrifices fairness on Qwen-2.5B and Gemma-2, despite maintaining high accuracy (Table 4, 9).

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>
	JUDGE	<b>0.575</b> <sub>0.008</sub>	<b>0.167</b> <sub>0.006</sub>	<b>0.160</b> <sub>0.005</sub>	<b>0.014</b> <sub>0.002</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>0.003</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>
	JUDGE	<b>0.531</b> <sub>0.010</sub>	<b>0.091</b> <sub>0.005</sub>	<b>0.117</b> <sub>0.007</sub>	<b>0.015</b> <sub>0.009</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>
	JUDGE	<b>0.621</b> <sub>0.014</sub>	<b>0.307</b> <sub>0.011</sub>	<b>0.303</b> <sub>0.009</sub>	<b>0.053</b> <sub>0.009</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>
	JUDGE	<b>0.645</b> <sub>0.008</sub>	<b>0.224</b> <sub>0.006</sub>	<b>0.207</b> <sub>0.004</sub>	<b>0.025</b> <sub>0.003</sub>

One trend remains consistent: *methods behave similarly across demonstration sizes*, with performance staying stable across  $k = 5$  and  $k = 10$  for a given model, dataset, and method.

### 5.5 Performance Comparison

JUDGE consistently provides strong improvements across 32 settings (4 LLMs \* 4 Datasets \* 2 Demonstration set sizes), achieving the best performance in most cases and near-best results in the instances where it is not the top performer. We attribute this to its greedy approach, which iteratively selects demonstrations by maximizing their marginal contribution based on LLM feedback using a semantically diverse jury set. Given the high variability in data types and LLM architectures, we believe this step-by-step feedback is key to generalizability.

Notably, baselines that incorporate LLM feedback, like FCG, tend to perform better than those relying solely on heuristics, which often lack consistency—excelling in some cases but failing in others. For instance, Counterfactual selection significantly improves fairness over Random on Gemma-2 for Adult, but worsens on Qwen-2.5 for the same dataset (Table 1). Similarly, Instruct improves fairness over Random on LLaMA-3 for COMPAS (Ta-

Table 3: Results for ACS 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	0.693 <sub>0.009</sub>	0.122 <sub>0.008</sub>	0.106 <sub>0.009</sub>	0.008 <sub>0.001</sub>
	Balanced	0.689 <sub>0.016</sub>	0.089 <sub>0.009</sub>	0.070 <sub>0.008</sub>	0.004 <sub>0.000</sub>
	Cfact.	0.653 <sub>0.005</sub>	0.092 <sub>0.004</sub>	0.092 <sub>0.004</sub>	0.004 <sub>0.000</sub>
	Instruct	0.684 <sub>0.010</sub>	0.115 <sub>0.006</sub>	0.101 <sub>0.006</sub>	0.008 <sub>0.001</sub>
	FairICL	0.688 <sub>0.011</sub>	0.098 <sub>0.008</sub>	0.010 <sub>0.004</sub>	0.008 <sub>0.002</sub>
	FCG	0.759 <sub>0.010</sub>	0.066 <sub>0.005</sub>	0.071 <sub>0.004</sub>	0.002 <sub>0.006</sub>
	FADS	0.697 <sub>0.008</sub>	0.116 <sub>0.006</sub>	0.101 <sub>0.004</sub>	0.008 <sub>0.001</sub>
	<b>JUDGE</b>	<b>0.764</b> <sub>0.007</sub>	<b>0.045</b> <sub>0.002</sub>	<b>0.049</b> <sub>0.003</sub>	<b>0.001</b> <sub>0.000</sub>
MISTRAL-7B	Random	0.603 <sub>0.007</sub>	0.091 <sub>0.005</sub>	0.052 <sub>0.006</sub>	0.005 <sub>0.001</sub>
	Balanced	0.558 <sub>0.013</sub>	0.070 <sub>0.009</sub>	<b>0.032</b> <sub>0.009</sub>	0.003 <sub>0.000</sub>
	Cfact.	0.607 <sub>0.009</sub>	0.085 <sub>0.005</sub>	0.063 <sub>0.004</sub>	0.006 <sub>0.001</sub>
	Instruct	0.592 <sub>0.017</sub>	0.094 <sub>0.12</sub>	0.108 <sub>0.011</sub>	0.007 <sub>0.001</sub>
	FairICL	0.610 <sub>0.005</sub>	0.089 <sub>0.004</sub>	0.051 <sub>0.003</sub>	0.005 <sub>0.001</sub>
	FCG	0.648 <sub>0.008</sub>	0.051 <sub>0.007</sub>	0.069 <sub>0.004</sub>	0.007 <sub>0.001</sub>
	FADS	0.599 <sub>0.012</sub>	0.088 <sub>0.006</sub>	0.051 <sub>0.004</sub>	0.005 <sub>0.000</sub>
	<b>JUDGE</b>	<b>0.651</b> <sub>0.011</sub>	<b>0.031</b> <sub>0.005</sub>	0.036 <sub>0.006</sub>	<b>0.001</b> <sub>0.000</sub>
GEMMA-2-9B	Random	0.696 <sub>0.009</sub>	0.225 <sub>0.007</sub>	0.223 <sub>0.010</sub>	0.028 <sub>0.004</sub>
	Balanced	0.707 <sub>0.012</sub>	0.233 <sub>0.009</sub>	0.179 <sub>0.009</sub>	0.028 <sub>0.002</sub>
	Cfact.	0.690 <sub>0.010</sub>	0.227 <sub>0.008</sub>	0.227 <sub>0.009</sub>	0.027 <sub>0.003</sub>
	Instruct	0.690 <sub>0.016</sub>	0.263 <sub>0.010</sub>	0.278 <sub>0.010</sub>	0.039 <sub>0.005</sub>
	FairICL	0.691 <sub>0.014</sub>	0.211 <sub>0.007</sub>	0.218 <sub>0.012</sub>	0.027 <sub>0.004</sub>
	FCG	0.705 <sub>0.012</sub>	0.141 <sub>0.007</sub>	0.136 <sub>0.009</sub>	0.018 <sub>0.002</sub>
	FADS	<b>0.709</b> <sub>0.016</sub>	0.205 <sub>0.006</sub>	0.274 <sub>0.010</sub>	0.031 <sub>0.003</sub>
	<b>JUDGE</b>	0.704 <sub>0.010</sub>	<b>0.131</b> <sub>0.006</sub>	<b>0.124</b> <sub>0.006</sub>	<b>0.013</b> <sub>0.001</sub>
QWEN-2.5-32B	Random	0.727 <sub>0.014</sub>	0.101 <sub>0.008</sub>	0.059 <sub>0.010</sub>	0.005 <sub>0.001</sub>
	Balanced	0.728 <sub>0.012</sub>	0.076 <sub>0.008</sub>	<b>0.017</b> <sub>0.008</sub>	0.003 <sub>0.000</sub>
	Cfact.	0.731 <sub>0.005</sub>	0.087 <sub>0.003</sub>	0.032 <sub>0.003</sub>	0.004 <sub>0.001</sub>
	Instruct	0.735 <sub>0.014</sub>	0.191 <sub>0.009</sub>	0.125 <sub>0.010</sub>	0.018 <sub>0.002</sub>
	FairICL	0.724 <sub>0.011</sub>	0.091 <sub>0.006</sub>	0.076 <sub>0.007</sub>	0.005 <sub>0.001</sub>
	FCG	0.727 <sub>0.006</sub>	0.059 <sub>0.003</sub>	0.051 <sub>0.003</sub>	0.002 <sub>0.000</sub>
	FADS	0.729 <sub>0.003</sub>	0.097 <sub>0.002</sub>	0.046 <sub>0.004</sub>	0.005 <sub>0.001</sub>
	<b>JUDGE</b>	<b>0.739</b> <sub>0.010</sub>	<b>0.025</b> <sub>0.005</sub>	0.036 <sub>0.005</sub>	<b>0.001</b> <sub>0.000</sub>

ble 2) but significantly harms it on Adult using the same LLM. FADS, designed to mitigate both model and data bias, performs well in many cases with some exceptions on certain datasets. FairICL, which ranks demonstrations by learning latent concept variables using a local LLaMA model, can face challenges due to architectural differences between models. Overall, JUDGE remains very consistent across all settings, improving fairness across metrics while maintaining accuracy. Its LLM-driven, stepwise construction ensures robust, data- and model-agnostic performance, making it a reliable and consistent approach.

In summary, several key distinctions in its design contribute to JUDGE’s performance:

- JUDGE constructs a single, optimized demonstration set by explicitly exploring interaction effects between examples during selection, evaluating each candidate in combination with previously chosen ones. This approach captures how examples collectively influence fairness and accuracy for the overall task. In contrast, some methods select demonstration sets specific to each test example (e.g., based on criteria like similarity to the test instance),

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

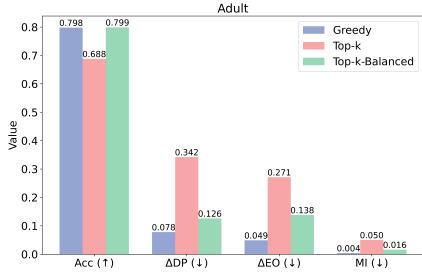
	Method	Acc. $\uparrow$	$\Delta DP \downarrow$	$\Delta EO \downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.895 <sub>0.012</sub>	0.299 <sub>0.009</sub>	0.493 <sub>0.015</sub>	0.054 <sub>0.004</sub>
	Balanced	0.663 <sub>0.016</sub>	0.406 <sub>0.008</sub>	0.377 <sub>0.006</sub>	0.047 <sub>0.005</sub>
	Cfact.	0.871 <sub>0.015</sub>	0.272 <sub>0.010</sub>	0.435 <sub>0.018</sub>	0.044 <sub>0.003</sub>
	Instruct	0.862 <sub>0.019</sub>	0.197 <sub>0.011</sub>	0.307 <sub>0.019</sub>	0.032 <sub>0.003</sub>
	FairICL	0.764 <sub>0.015</sub>	0.312 <sub>0.008</sub>	0.346 <sub>0.006</sub>	0.045 <sub>0.002</sub>
	FCG	0.909 <sub>0.016</sub>	0.082 <sub>0.016</sub>	0.178 <sub>0.020</sub>	0.019 <sub>0.003</sub>
	FADS	0.898 <sub>0.004</sub>	0.242 <sub>0.003</sub>	0.353 <sub>0.006</sub>	0.039 <sub>0.001</sub>
	<b>JUDGE</b>	<b>0.911</b> <sub>0.026</sub>	<b>0.057</b> <sub>0.027</sub>	<b>0.104</b> <sub>0.035</sub>	<b>0.005</b> <sub>0.001</sub>
MISTRAL-7B	Random	0.905 <sub>0.009</sub>	0.187 <sub>0.011</sub>	0.338 <sub>0.007</sub>	0.029 <sub>0.003</sub>
	Balanced	0.871 <sub>0.012</sub>	0.219 <sub>0.004</sub>	0.362 <sub>0.007</sub>	0.031 <sub>0.001</sub>
	Cfact.	0.904 <sub>0.012</sub>	0.200 <sub>0.011</sub>	0.418 <sub>0.006</sub>	0.029 <sub>0.001</sub>
	Instruct	0.913 <sub>0.010</sub>	<b>0.023</b> <sub>0.004</sub>	0.077 <sub>0.005</sub>	<b>0.004</b> <sub>0.000</sub>
	FairICL	0.902 <sub>0.010</sub>	0.173 <sub>0.006</sub>	0.311 <sub>0.009</sub>	0.026 <sub>0.003</sub>
	FCG	0.943 <sub>0.014</sub>	0.038 <sub>0.006</sub>	0.091 <sub>0.04</sub>	0.018 <sub>0.003</sub>
	FADS	0.934 <sub>0.006</sub>	0.103 <sub>0.004</sub>	0.227 <sub>0.003</sub>	0.019 <sub>0.002</sub>
	<b>JUDGE</b>	<b>0.946</b> <sub>0.013</sub>	0.027 <sub>0.003</sub>	<b>0.059</b> <sub>0.004</sub>	0.008 <sub>0.001</sub>
GEMMA-2-9B	Random	0.853 <sub>0.011</sub>	0.372 <sub>0.007</sub>	0.569 <sub>0.010</sub>	0.058 <sub>0.004</sub>
	Balanced	0.756 <sub>0.007</sub>	0.419 <sub>0.011</sub>	0.436 <sub>0.008</sub>	0.056 <sub>0.003</sub>
	Cfact.	0.747 <sub>0.007</sub>	0.366 <sub>0.004</sub>	0.358 <sub>0.006</sub>	0.042 <sub>0.003</sub>
	Instruct	<b>0.878</b> <sub>0.004</sub>	0.344 <sub>0.005</sub>	0.553 <sub>0.003</sub>	0.056 <sub>0.001</sub>
	FairICL	0.844 <sub>0.010</sub>	0.341 <sub>0.009</sub>	0.357 <sub>0.012</sub>	0.041 <sub>0.002</sub>
	FCG	0.845 <sub>0.013</sub>	0.258 <sub>0.009</sub>	0.267 <sub>0.011</sub>	0.029 <sub>0.003</sub>
	FADS	0.877 <sub>0.009</sub>	0.287 <sub>0.006</sub>	0.502 <sub>0.007</sub>	0.047 <sub>0.003</sub>
	<b>JUDGE</b>	0.855 <sub>0.013</sub>	<b>0.227</b> <sub>0.009</sub>	<b>0.214</b> <sub>0.008</sub>	<b>0.025</b> <sub>0.003</sub>
QWEN-2.5-32B	Random	0.865 <sub>0.007</sub>	0.327 <sub>0.005</sub>	0.414 <sub>0.004</sub>	0.052 <sub>0.002</sub>
	Balanced	0.831 <sub>0.011</sub>	0.392 <sub>0.005</sub>	0.493 <sub>0.006</sub>	0.061 <sub>0.002</sub>
	Cfact.	0.840 <sub>0.009</sub>	0.366 <sub>0.006</sub>	0.418 <sub>0.008</sub>	0.055 <sub>0.005</sub>
	Instruct	0.883 <sub>0.008</sub>	0.370 <sub>0.005</sub>	0.534 <sub>0.006</sub>	0.072 <sub>0.004</sub>
	FairICL	0.860 <sub>0.018</sub>	0.316 <sub>0.010</sub>	0.449 <sub>0.014</sub>	0.057 <sub>0.002</sub>
	FCG	0.862 <sub>0.016</sub>	0.248 <sub>0.013</sub>	0.293 <sub>0.012</sub>	0.035 <sub>0.003</sub>
	FADS	<b>0.889</b> <sub>0.021</sub>	0.238 <sub>0.012</sub>	0.419 <sub>0.016</sub>	0.044 <sub>0.009</sub>
	<b>JUDGE</b>	0.882 <sub>0.016</sub>	<b>0.214</b> <sub>0.013</sub>	<b>0.273</b> <sub>0.015</sub>	<b>0.027</b> <sub>0.001</sub>

which may not fully capture these broader interaction effects or the nuanced collective influence critical for robust performance. As we will show in Section 5.6, simply choosing the best examples (“Top-k”) or balanced top examples (“Top-k-Balanced”) is less effective than JUDGE’s incremental greedy selection, which adds demonstrations one at a time based on joint impact. This shows that demonstrations can be individually strong but sensitive to others that they are paired with.

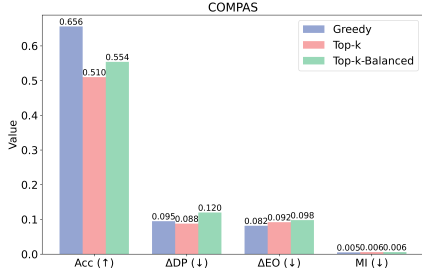
- Fairness metrics like  $\Delta DP$  and  $\Delta EO$  are group-based measures calculated across populations, not for individual examples. JUDGE optimizes directly for these group fairness metrics during selection using the jury set.
- JUDGE incorporates LLM feedback at all steps of the selection process, allowing it to adapt to each LLM’s specific characteristics. This explains its consistent performance across diverse LLMs.

## 5.6 Ablation: Greedy vs. Top-k Selection

To validate our greedy selection approach, we compare it against two alternatives: (1) Top-k, the top



(a) Adult dataset



(b) COMPAS dataset

Figure 3: Comparison of Greedy vs. Top-k alternatives

$k$  candidates that individually perform the best on the jury set, (2) Top-k-Balanced, which is a stratified selection that picks the top samples from each combination of protected group and label,  $(g, y)$ . Results show that greedy selection consistently outperforms both methods, highlighting the importance of marginal contribution of each example in building a fair and effective demonstration set. The effect of Top-k and Top-k-Balanced varies by dataset. As shown in Figure 3, on Adult (LLaMA-3-8B), Top-k exhibits a dramatic drop in fairness performance, while Top-k-Balanced fares better. On COMPAS, we see competitive fairness performance across variants, but upon closer inspection we observe that Top-k and Top-k-Balanced selection suffers large drops in accuracy. These findings underscore the inherent variability in ICL and reinforce the strength of the greedy approach, which incrementally selects candidates while considering their interactions with the existing set.

### 5.7 Impact of Jury Set Size

To assess the impact of jury size, we vary the number of examples per group-label combination  $m$  from 1 to 100, keeping all other parameters constant with  $k = 5$  demonstrations. Figure 4 (LLaMA-3 on Adult) shows results for accuracy and  $\Delta DP$  (full results in Appendix B.4, Figure 10). Results on *Adult* indicate that performance stabilizes as  $m$  increases, with diminishing returns beyond  $m > 25$  despite higher computational costs.

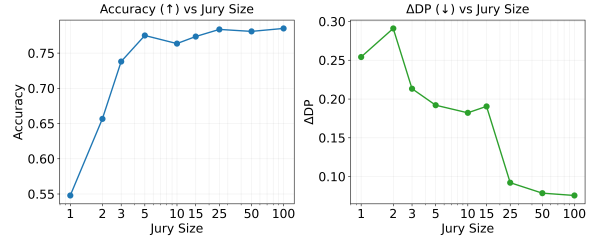
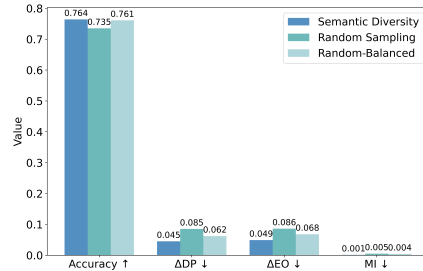
Figure 4: Accuracy and  $\Delta DP$  against the size of the jury set for Adult. Higher sizes show diminishing returns.

Figure 5: Comparison of diversity vs. other sampling techniques for the jury set on the ACS dataset.

Accuracy plateaus quickly, with  $m = 5$  or 10 being sufficient, while fairness improves up to  $m = 50$ .

### 5.8 Jury Set Diversity

To examine the impact of semantic diversity in jury set construction, we compare three methods: (1) Random Sampling, (2) Random-Balanced (random sampling after enforcing equal representation across protected group-label combinations), and (3) Semantic Diversity-based selection. We fix  $m = 25$  for this comparison. With jury sets constrained to be small for computational efficiency, Figure 5 shows that diversity-based selection outperforms both alternatives on the ACS dataset. A similar experiment on Adult is provided in Appendix B.5.

### 5.9 Sentence Embedding Models

We selected SentenceBERT for creating sentence embeddings. Like all pre-trained models, these embeddings may contain inherent biases. However, it is crucial to note that SentenceBERT is employed only for preprocessing, to select a diverse set of examples. The core fairness evaluation is conducted directly by the Large Language Model (LLM), not by the embedding model itself.

The choice of SentenceBERT was motivated by its proven general-purpose capabilities, which are well-suited for the diverse datasets used in our experiments. This ensures that our pipeline is broadly



applicable and maintains consistency across different domains. We recognize that for domain-specific applications, our method could be adapted to leverage specialized embeddings. For instance, LegalBERT (Chalkidis et al., 2020) could be employed for legal applications. To investigate the impact of the embedding model choice, we conducted a comparative analysis using LegalBERT on the COMPAS dataset. The results, obtained using LLAMA-3-8B with the number of demonstrations set to five, are presented in Table 5.

Table 5: Performance Comparison with SentenceBERT (SBERT) and LegalBERT (LBert) on COMPAS.

System	Acc. $\uparrow$	$\Delta$ DP $\downarrow$	$\Delta$ EO $\downarrow$	MI $\downarrow$
JUDGE (SBERT)	0.656	0.105	0.082	0.006
JUDGE (LBERT)	0.637	0.099	0.080	0.006

As shown in Table 5, the performance when using LegalBERT embeddings remains largely comparable to that of SentenceBERT. The system with LegalBERT exhibits a slight trade-off, with a marginally lower accuracy but improved fairness with  $\Delta$ DP and  $\Delta$ EO. This suggests that our method is robust to the choice of the underlying sentence embedding model, although domain-specific embeddings like LegalBERT may offer benefits or nuanced trade-offs in specialized contexts.

## 6 Related Work

**Demonstration Selection in ICL** The problem of selecting demonstrations for ICL has received significant attention. (Liu et al., 2022) showed that finding demonstrations which are semantically similar to the test data often shows promising results. Wu et al. (2023) addressed this challenge by establishing a select-then-rank framework where they first limit the search space of demonstrations and rank the remaining examples through heuristics. Peng et al. (2024) highlighted that both data and model factors contribute to variability in performance. Meanwhile, Ma et al. (2023) showed that predictive performance can be improved by selecting examples that minimize predictive bias. To address efficiency concerns, Yang et al. (2023) proposed a two-stage Determinantal Point Process (DPP) method to select a fixed, representative subset of demonstrations, improving efficiency while maintaining performance.

**Fair Demonstration Selection in ICL** The fairness of language models has received significant attention (Doan et al., 2024; Chu et al., 2024). Liu

et al. (2024) showed that LLMs exhibit significant bias in tabular classification. In ICL, fair demonstration selection is crucial. Hu et al. (2024) investigated how varying the composition of demonstrations affects fairness outcomes in ICL. The authors proposed a fairness-aware selection method that employs clustering and evolutionary strategies to curate a diverse and representative sample set from the training data. Meanwhile, Wang et al. (2024) introduced FADS, which addresses the challenge of fair demonstration selection by mitigating both model bias and bias in the data. Other approaches have explored leveraging counterfactual analysis. Bhaila et al. (2024) introduced a method that uses latent concept variables learned through counterfactual examples to evaluate the fairness of demonstrations. The idea of utilizing counterfactual examples is also presented by Li et al. (2023), which picks examples from the privileged group and flips the sensitive attribute to create new examples. Finally, Atwood et al. (2024) showed how prompting the model by explicitly asking it to be fair can also be effective. JUDGE adds to this growing body of literature focused on selecting fair and representative samples for ICL.

## 7 Conclusion

We propose JUDGE, a greedy framework for fair demonstration selection in ICL, guided by a jury set. Across four datasets and four LLM architectures, our method consistently improves fairness while maintaining accuracy, outperforming existing approaches. We further highlight the high variability of different methods across different datasets and language models, and establish the importance of considering demonstrations as a cohesive set rather than as individual examples to ensure fairness. As LLMs expand into critical applications, JUDGE offers a practical and robust solution for ensuring fairness in ICL.

## Limitations

This work investigates the problem of fairness aware demonstration selection for in-context learning. In order to do so, this work explores various open-source LLM architectures from Google, Meta, Mistral, and Alibaba. While these architectures have varied sizes ranging from 7B to 32B parameters, a key limitation in our work is that, due to hardware limitations we do not investigate the effect on truly massive models like LLAMA-

3-405B. Furthermore, financial constraints prevent us from using closed-source paid platforms like GPT-4o, given the large number of LLM queries required across our datasets, baselines, LLMs and demonstration sizes. Nonetheless, we believe we chose a diverse and representative set of highly performant open-source LLMs to make our study comprehensive. Furthermore, our study limits itself to exploring binary in-context classification as well as binary sensitive group settings. In the future, we plan to consider broader classification settings. Finally, in line with prior work, we aimed to conduct a comprehensive study across widely popular fairness datasets, which are typically tabular in nature and are thus serialized into a natural language prompt for the LLM. In the future, we hope to study other types of data in the context of fairness in large language models.

## Acknowledgments

The research reported herein was supported in part by NSF grant number 2147375 and NIST grant number 60NANB24D143, and the National Center for Transportation Cybersecurity and Resiliency (TraCR). Any opinions, findings, conclusions, and recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF, NIST or TraCR.

## References

- Hadis Anahideh, Abolfazl Asudeh, and Saravanan Thirumuruganathan. 2022. Fair active learning. *Expert Systems with Applications*, 199:116981.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias: There’s software used across the country to predict future criminals. and it’s biased against blacks. *ProPublica*, 23(2016):77–91.
- James Atwood, Preethi Lahoti, Ananth Balashankar, Flavien Prost, and Ahmad Beirami. 2024. Inducing group fairness in llm-based decisions. *arXiv preprint arXiv:2406.16738*.
- Karuna Bhaila, Minh-Hao Van, Kennedy Edemacu, Chen Zhao, Feng Chen, and Xintao Wu. 2024. Fair in-context learning via latent concept variables. *arXiv preprint arXiv:2411.02671*.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, et al. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. Legal-bert: The muppets straight out of law school. *arXiv preprint arXiv:2010.02559*.
- Zhibo Chu, Zichong Wang, and Wenbin Zhang. 2024. Fairness in large language models: A taxonomic survey. *ACM SIGKDD explorations newsletter*, 26(1):34–48.
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. 2021. Retiring adult: New datasets for fair machine learning. In *Advances in Neural Information Processing Systems*, volume 34, pages 6478–6490.
- Thang Viet Doan, Zichong Wang, Nhat Nguyen Minh Hoang, and Wenbin Zhang. 2024. [Fairness in large language models in three hours](#). In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM ’24*, page 5514–5517, New York, NY, USA. Association for Computing Machinery.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*.
- Dheeru Dua and Casey Graff. 2019. [Uci machine learning repository](#).
- Anirudh Dubey et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.
- Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29.
- Jingyu Hu, Weiru Liu, and Mengnan Du. 2024. Strategic demonstration selection for improved fairness in llm in-context learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7460–7475.
- Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, et al. 2024. [Position: TrustLLM: Trustworthiness in large language models](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 20166–20270. PMLR.
- B. Hui et al. 2024. Qwen2.5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- A. Q. Jiang et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

- Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-aware classifier with prejudice remover regularizer. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2012, Bristol, UK, September 24-28, 2012. Proceedings, Part II 23*, pages 35–50. Springer.
- Rémi Leblond, Jean-Baptiste Alayrac, Laurent Sifre, Miruna Pislariu, Lespiau Jean-Baptiste, Ioannis Antonoglou, Karen Simonyan, and Oriol Vinyals. 2021. Machine translation decoding beyond beam search. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8410–8434.
- Yunqi Li, Lanjing Zhang, and Yongfeng Zhang. 2023. Fairness of chatgpt. *arXiv preprint arXiv:2305.18569*.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
- Yanchen Liu, Srishti Gautam, Jiaqi Ma, and Himabindu Lakkaraju. 2024. [Confronting LLMs with traditional ML: Rethinking the fairness of large language models in tabular classifications.](#) In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3603–3620, Mexico City, Mexico. Association for Computational Linguistics.
- Huan Ma, Changqing Zhang, Yatao Bian, Lemao Liu, Zhirui Zhang, Peilin Zhao, Shu Zhang, Huazhu Fu, Qinghua Hu, and Bingzhe Wu. 2023. Fairness-guided few-shot prompting for large language models. *Advances in Neural Information Processing Systems*, 36:43136–43155.
- George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. 1978. An analysis of approximations for maximizing submodular set functions—i. *Mathematical programming*, 14:265–294.
- Kirtan Padh, Diego Antognini, Emma Lejal-Glaude, Boi Faltings, and Claudiu Musat. 2021. Addressing fairness in classification with a model-agnostic multi-objective algorithm. In *Uncertainty in artificial intelligence*, pages 600–609. PMLR.
- Keqin Peng, Liang Ding, Yancheng Yuan, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2024. [Revisiting demonstration selection strategies in in-context learning.](#) In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9090–9101, Bangkok, Thailand. Association for Computational Linguistics.
- N Reimers. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Morgane Riviere et al. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Ozan Sener and Silvio Savarese. 2018. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*.
- Song Wang, Peng Wang, Yushun Dong, Tong Zhou, Lu Cheng, Yangfeng Ji, and Jundong Li. 2024. On demonstration selection for improving fairness in language models. In *Workshop on Socially Responsible Language Modelling Research*.
- Kai Wei, Rishabh Iyer, and Jeff Bilmes. 2015. Submodularity in data subset selection and active learning. In *International conference on machine learning*, pages 1954–1963. PMLR.
- Linda F Wightman. 1998. Lsac national longitudinal bar passage study. Technical report, Law School Admission Council.
- Zhiyong Wu, Yaoxiang Wang, Jiacheng Ye, and Lingpeng Kong. 2023. [Self-adaptive in-context learning: An information compression perspective for in-context example selection and ordering.](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1423–1436, Toronto, Canada. Association for Computational Linguistics.
- Zhao Yang, Yuanzhe Zhang, Dianbo Sui, Cao Liu, Jun Zhao, and Kang Liu. 2023. [Representative demonstration selection for in-context learning with two-stage determinantal point process.](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5443–5456, Singapore. Association for Computational Linguistics.

## A Fairness Metrics Formulation

Here we provide detailed mathematical formulations of the fairness metrics used in our analysis. For all metrics, we take absolute values to ensure positive measures of disparity, where zero indicates perfect fairness and larger values indicate greater disparity.

### A.1 Demographic Parity Difference ( $\Delta DP$ )

The Demographic Parity Difference ( $\Delta DP$ ) measures the absolute difference in positive label rates between groups:

$$\Delta DP = |P(y = 1 | g(x) = g_1) - P(y = 1 | g(x) = g_2)| \quad (16)$$

where  $y = 1$  denotes a positive label. A  $\Delta DP$  of 0 indicates perfect demographic parity.

## A.2 Equalized Odds Difference ( $\Delta EO$ )

The Equalized Odds Difference ( $\Delta EO$ ) measures disparities in both true positive rates (TPR) and false positive rates (FPR) between groups:

$$\Delta EO = \max(|\text{TPR}_{g_1} - \text{TPR}_{g_2}|, |\text{FPR}_{g_1} - \text{FPR}_{g_2}|) \quad (17)$$

where

$$\text{TPR}_g = P(y = 1 \mid g(x) = g, y^* = 1) \quad (18)$$

$$\text{FPR}_g = P(y = 1 \mid g(x) = g, y^* = 0) \quad (19)$$

Here,  $y^*$  represents the true label, and  $y = 1$  represents the predicted positive label. A  $\Delta EO$  of 0 indicates perfect equalized odds.

## A.3 Mutual Information Fairness

The mutual information between protected group membership  $G$  and the positive label assignment is:

$$I(G; Y) = \sum_{g,y} P(g,y) \log \frac{P(g,y)}{P(g)P(y)} \quad (20)$$

where  $y$  denotes whether an instance receives a positive label. Lower mutual information indicates greater independence between the positive label assignment and protected group membership. This metric is naturally non-negative, with 0 indicating perfect independence.

## B Additional Experiment Details

### B.1 DiverseSelect

The pseudocode for maximizing semantic diversity in selection is shown in Algorithm 2.

---

**Algorithm 2** DiverseSelect: Diversity-Based Example Selection

---

**Require:** Initial pool  $\mathcal{D}$ , target size  $k$

**Ensure:** Diverse subset  $\mathcal{D}_{diverse}$

- 1: Compute  $S_{ij} = \text{sim}(x_i, x_j)$  for all  $x_i, x_j \in \mathcal{D}$
  - 2:  $\mathcal{D}_{diverse} \leftarrow x_r$  where  $x_r$  is randomly sampled from  $\mathcal{D}$
  - 3: **for**  $t \leftarrow 1$  to  $k - 1$  **do**
  - 4:   **for**  $x_i \in \mathcal{D} \setminus \mathcal{D}_{diverse}$  **do**
  - 5:      $s_i \leftarrow \max_{x_j \in \mathcal{D}_{diverse}} S_{ij}$
  - 6:   **end for**
  - 7:    $x_t \leftarrow \arg \min_{x_i \in \mathcal{D} \setminus \mathcal{D}_{diverse}} s_i$
  - 8:    $\mathcal{D}_{diverse} \leftarrow \mathcal{D}_{diverse} \cup x_t$
  - 9: **end for**
  - 10: **return**  $\mathcal{D}_{diverse}$
- 

## B.2 Dataset Details

**Adult Income** The UCI Adult dataset (Dua and Graff, 2019) contains demographic and employment information for 48,842 individuals. The task is to predict whether annual income exceeds \$50,000, with gender as the protected attribute. The prompt template for this dataset is shown in Figure 8.

**COMPAS** This dataset (Angwin et al., 2016) includes criminal history and demographic data for defendants. The classification task is predicting recidivism risk. We use a binarized race (Caucasian vs African-American) as the protected attribute. The prompt template for this dataset is shown in Figure 6.

**Law School** The LSAC dataset (Wightman, 1998) contains admissions data and academic performance for law school students. The model predicts whether a student passes the bar, with a binarized race (Caucasian vs Not-Caucasian) as the protected attribute. The prompt template for this dataset is shown in Figure 7.

**ACS Income** The ACS PUMS dataset (Ding et al., 2021) contains demographic and employment information from the American Community Survey. The task predicts if income exceeds \$50,000, using gender as the protected attribute. The ACS Income dataset in its original form contains over 1.66 million datapoints, which is far larger than all other datasets that we consider, combined. For LLM in-context classification, this becomes prohibitively expensive from a computation perspective. As a result, we randomly downsample ACS Income down to 48,842 samples, which is the same size as the closely related Adult Dataset. Both datasets track American income data, but ACS provides much newer information from 2018 instead of 1994 for Adult. The prompt template for this dataset is shown in Figure 9.

### B.3 Additional Results

This section presents results for all LLMs and all datasets with 10 demonstrations provided for In-Context Learning. These can be seen in Tables 5-8 for each of the four datasets.

### B.4 Full Results for the Effect of Jury Set Size

Here we provide the results for all metrics for our experiment in Section 5.7, which tests the effect of different jury set sizes on the Adult dataset using LLAMA-3-8B. This is shown in Figure 10. This

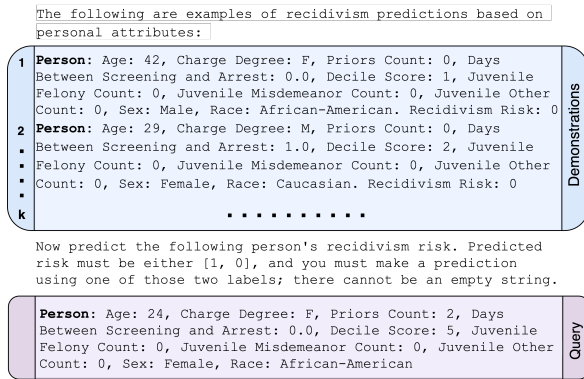


Figure 6: COMPAS Prompt Template

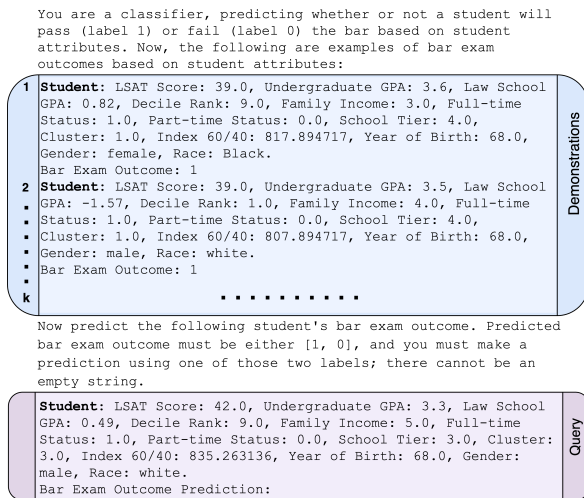


Figure 7: Law School Prompt Template

figure additionally shows the  $\Delta E O$  and MI metrics which show the same pattern as the  $\Delta D P$  metric.

### B.5 Additional Jury Set Diversity Experiment

We conduct the same jury set diversity experiment from Section 5.8 on the Adult dataset, comparing random selection, balanced-random selection, and our diversity-based approach. As with ACS, the jury set size is fixed at  $m = 25$  for this experiment, across all methods.

Figure 11 illustrates that diversity-based selection also outperforms other sampling strategies on the Adult dataset, reinforcing the importance of semantic diversity in jury construction.

### B.6 Sensitivity Analysis

To understand the impact of trade-off parameters in different methods, we conduct a sensitivity analysis by varying the fairness-accuracy balancing coefficients in JUDGE, FCG, and FairICL on the Adult dataset over LLAMA-3-8B. We select FCG and FairICL as baselines because their respective au-

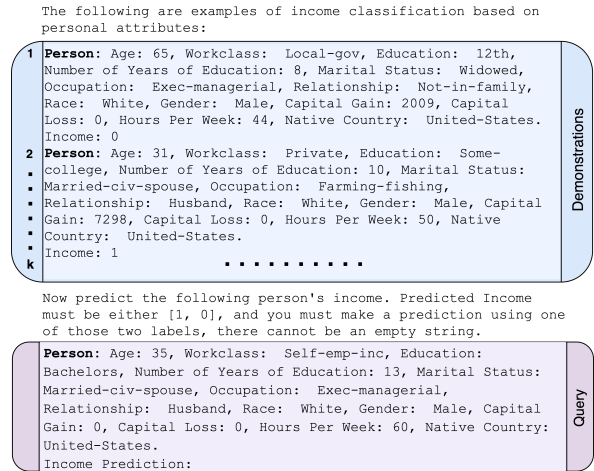


Figure 8: Adult Prompt Template

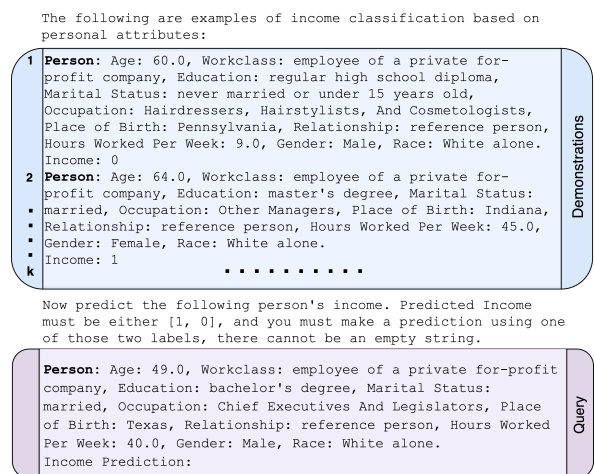


Figure 9: ACS Prompt Template

thors explicitly identify  $\alpha$  and  $\tilde{D}$  as key parameters that influence fairness, making them well-suited for comparison with JUDGE.

### B.7 JUDGE: Sensitivity to $\omega$

JUDGE introduces  $\omega$  as a parameter that controls the trade-off between accuracy and fairness. The selection of demonstrations is influenced by  $\omega$ , where lower values prioritize fairness while higher values emphasize accuracy. We evaluate JUDGE at  $\omega \in \{0.4, 0.5, 0.6, 0.7, 0.8\}$ . We choose this set because we find in our experiments that  $\omega$  values that prioritize fairness slightly more than accuracy work well, improving fairness while also retaining predictive performance.

### B.8 FCG: Sensitivity to $\alpha$

FCG uses  $\alpha$  in the EvolScore function to balance accuracy and fairness. The original paper sets  $\alpha = 0.5$ , and we analyze values in

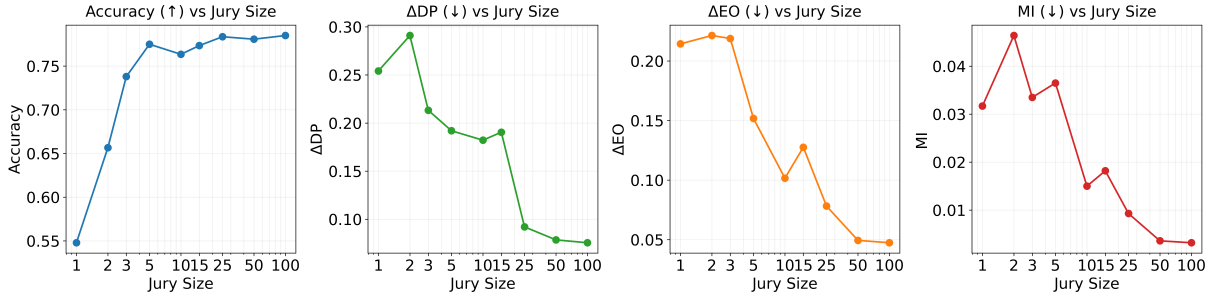


Figure 10: Comparing metrics against the size of the jury set for Adult. Higher sizes show diminishing returns.

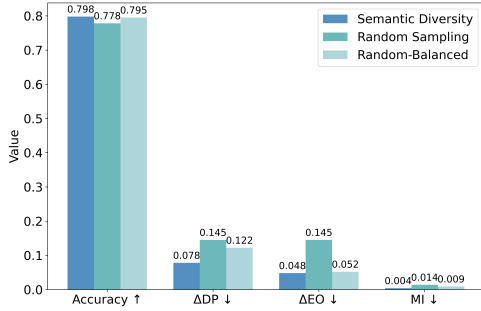


Figure 11: Comparison of diversity vs. other sampling techniques for the jury set on the Adult dataset.

{0.3, 0.4, 0.5, 0.6, 0.7} to assess how the fairness-accuracy trade-off shifts.

### B.9 FairICL: Sensitivity to $\tilde{D}$

FairICL introduces  $\tilde{D}$ , which represents the fraction of augmented data used for fairness-aware training. The original study evaluates FairICL at  $\tilde{D} \in \{0\%, 25\%, 50\%, 100\%\}$ , highlighting its influence on fairness. To provide a more fine-grained analysis, we add an additional evaluation at  $\tilde{D} = 75\%$ , resulting in the set  $\{0\%, 25\%, 50\%, 75\%, 100\%\}$ .

### B.10 Results: Accuracy vs. Fairness Trade-Off

Figure 12 presents a scatter plot where each method’s trade-off variations are shown along two axes: Accuracy (Y-axis) and  $\Delta DP$  (X-axis). Each point represents a model trained with a different trade-off parameter. Points closer to the top-left are preferred (high accuracy, low  $\Delta DP$ ).

We observe clear trade-offs for JUDGE and FCG, where higher accuracy comes at the cost of fairness and vice versa, guided by the weighting provided by  $\omega$  and  $\alpha$ . On FairICL, we find the relationship to be less strong, with most points clustered around a similar area.

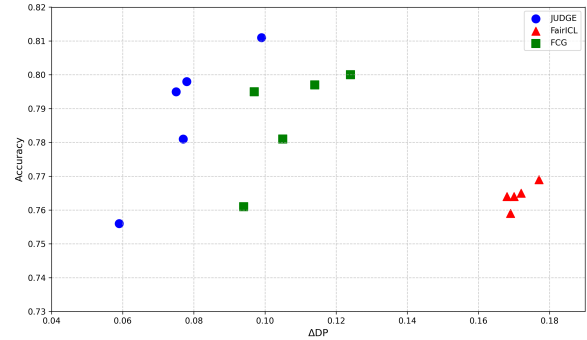


Figure 12: Scatter plot of Accuracy vs. Demographic Parity ( $\Delta DP$ ) for different trade-off parameter settings in JUDGE, FCG, and FairICL.

### B.11 Data Splits and Hyperparameters

For all datasets except Adult, we employ a consistent data splitting strategy:

- 20% for test set ( $\mathcal{D}_{\text{test}}$ )
- 70% for training set ( $\mathcal{D}_{\text{train}}$ )
- 10% for validation set ( $\mathcal{D}_{\text{validation}}$ )

For the Adult dataset, which provides a predefined train-test split, we maintain the original test set and split the training set into  $\mathcal{D}_{\text{train}}$  (90%) and  $\mathcal{D}_{\text{validation}}$  (10%).

It is important to note that this validation set is distinct from the jury set ( $\mathcal{J}$ ) used in our method. While the jury set is constructed from the training data to guide demonstration selection and is typically very small, the validation set is used exclusively for hyperparameter tuning.

To tune hyperparameters, we conduct a systematic grid search over two key hyperparameters:

1. The fairness-accuracy trade-off parameter  $\omega$  in the range  $[0.3, 0.9]$  with steps of 0.05
2. The number of examples per group-label combination  $m$  in the jury set, testing values:  $\{15, 20, 25, 35, 50\}$

Table 6: Results for Adult with 10 demonstrations. Each cell shows  $Mean_{S,D}$ .

	Method	Acc. $\uparrow$	$\Delta DP$ $\downarrow$	$\Delta EO$ $\downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.779 <sub>0.004</sub>	0.133 <sub>0.003</sub>	0.118 <sub>0.004</sub>	0.017 <sub>0.001</sub>
	Balanced	0.751 <sub>0.013</sub>	0.221 <sub>0.019</sub>	0.137 <sub>0.022</sub>	0.025 <sub>0.002</sub>
	Cfact.	0.776 <sub>0.018</sub>	0.144 <sub>0.014</sub>	0.142 <sub>0.015</sub>	0.015 <sub>0.003</sub>
	Instruct	0.781 <sub>0.022</sub>	0.252 <sub>0.017</sub>	0.289 <sub>0.019</sub>	0.046 <sub>0.004</sub>
	FairICL	0.777 <sub>0.015</sub>	0.146 <sub>0.012</sub>	0.164 <sub>0.135</sub>	0.014 <sub>0.003</sub>
	FCG	0.788 <sub>0.017</sub>	0.189 <sub>0.014</sub>	0.163 <sub>0.017</sub>	0.023 <sub>0.003</sub>
	FADS	0.772 <sub>0.013</sub>	0.161 <sub>0.011</sub>	0.098 <sub>0.005</sub>	0.020 <sub>0.003</sub>
	<b>JUDGE</b>	<b>0.794</b> <sub>0.011</sub>	<b>0.082</b> <sub>0.009</sub>	<b>0.092</b> <sub>0.008</sub>	<b>0.008</b> <sub>0.001</sub>
MISTRAL-7B	Random	0.755 <sub>0.014</sub>	0.209 <sub>0.012</sub>	0.262 <sub>0.009</sub>	0.023 <sub>0.004</sub>
	Balanced	0.585 <sub>0.009</sub>	0.220 <sub>0.011</sub>	0.170 <sub>0.008</sub>	0.023 <sub>0.003</sub>
	Cfact.	0.731 <sub>0.014</sub>	0.141 <sub>0.016</sub>	0.090 <sub>0.016</sub>	0.010 <sub>0.003</sub>
	Instruct	0.742 <sub>0.014</sub>	0.182 <sub>0.013</sub>	0.212 <sub>0.018</sub>	0.012 <sub>0.005</sub>
	FairICL	0.763 <sub>0.011</sub>	0.143 <sub>0.008</sub>	0.155 <sub>0.006</sub>	0.013 <sub>0.002</sub>
	FCG	0.758 <sub>0.022</sub>	0.122 <sub>0.013</sub>	0.083 <sub>0.016</sub>	0.012 <sub>0.003</sub>
	FADS	<b>0.775</b> <sub>0.012</sub>	0.192 <sub>0.009</sub>	0.244 <sub>0.013</sub>	0.022 <sub>0.003</sub>
	<b>JUDGE</b>	<b>0.771</b> <sub>0.010</sub>	<b>0.010</b> <sub>0.012</sub>	<b>0.058</b> <sub>0.010</sub>	<b>0.010</b> <sub>0.001</sub>
GEMMA-2-9B	Random	0.774 <sub>0.009</sub>	0.365 <sub>0.005</sub>	0.484 <sub>0.012</sub>	0.090 <sub>0.006</sub>
	Balanced	0.721 <sub>0.015</sub>	0.389 <sub>0.022</sub>	0.455 <sub>0.029</sub>	0.116 <sub>0.015</sub>
	Cfact.	0.762 <sub>0.020</sub>	0.276 <sub>0.017</sub>	0.383 <sub>0.019</sub>	0.076 <sub>0.013</sub>
	Instruct	0.764 <sub>0.013</sub>	0.408 <sub>0.013</sub>	0.523 <sub>0.011</sub>	0.107 <sub>0.009</sub>
	FairICL	0.763 <sub>0.013</sub>	0.301 <sub>0.021</sub>	0.323 <sub>0.024</sub>	0.072 <sub>0.011</sub>
	FCG	0.778 <sub>0.011</sub>	0.176 <sub>0.013</sub>	<b>0.179</b> <sub>0.014</sub>	0.053 <sub>0.003</sub>
	FADS	0.766 <sub>0.010</sub>	0.378 <sub>0.009</sub>	0.414 <sub>0.016</sub>	0.087 <sub>0.007</sub>
	<b>JUDGE</b>	<b>0.792</b> <sub>0.010</sub>	<b>0.173</b> <sub>0.015</sub>	0.197 <sub>0.019</sub>	<b>0.047</b> <sub>0.005</sub>
QWEN-2.5-32B	Random	0.741 <sub>0.016</sub>	0.210 <sub>0.015</sub>	0.129 <sub>0.004</sub>	0.021 <sub>0.003</sub>
	Balanced	0.728 <sub>0.017</sub>	0.223 <sub>0.019</sub>	0.152 <sub>0.013</sub>	0.026 <sub>0.004</sub>
	Cfact.	0.743 <sub>0.011</sub>	0.219 <sub>0.010</sub>	0.135 <sub>0.009</sub>	0.025 <sub>0.002</sub>
	Instruct	0.715 <sub>0.009</sub>	0.236 <sub>0.010</sub>	0.157 <sub>0.011</sub>	0.025 <sub>0.002</sub>
	FairICL	0.756 <sub>0.012</sub>	0.204 <sub>0.010</sub>	0.151 <sub>0.010</sub>	0.025 <sub>0.003</sub>
	FCG	<b>0.778</b> <sub>0.011</sub>	0.128 <sub>0.010</sub>	0.099 <sub>0.007</sub>	0.011 <sub>0.003</sub>
	FADS	0.706 <sub>0.010</sub>	0.206 <sub>0.007</sub>	0.132 <sub>0.005</sub>	0.022 <sub>0.004</sub>
	<b>JUDGE</b>	<b>0.775</b> <sub>0.009</sub>	<b>0.101</b> <sub>0.008</sub>	<b>0.078</b> <sub>0.004</sub>	<b>0.007</b> <sub>0.001</sub>

In Section 5.7, we demonstrated that jury sizes beyond  $m = 50$  yield diminishing returns, while very small values ( $m \in \{1, 2, 3, 5, 10\}$ ) show substantial performance gaps in fairness and accuracy. Based on these observations, we focus our parameter search on the more practical intermediate range. For jury set size selection, we do as follows: Starting from smaller values, we incrementally evaluate larger jury sizes until we observe diminishing returns in performance on the validation set. Specifically, if the relative improvement in both accuracy and fairness metrics between two consecutive jury sizes falls below 1%, we stop increasing the size. This process led to the selection of  $m = 25$  for Adult and COMPAS datasets, and  $m = 50$  for Law School and ACS datasets. The larger jury sizes for Law School and ACS datasets were chosen because these datasets exhibited continued performance improvements with larger jury sizes.

For the fairness-accuracy trade-off parameter  $\omega$ , we select the value that achieves the lowest  $\Delta DP$

Table 7: Results for COMPAS with 10 demonstrations. Each cell shows  $Mean_{S,D}$ .

	Method	Acc. $\uparrow$	$\Delta DP$ $\downarrow$	$\Delta EO$ $\downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.603 <sub>0.013</sub>	0.224 <sub>0.010</sub>	0.221 <sub>0.009</sub>	0.025 <sub>0.002</sub>
	Balanced	0.605 <sub>0.007</sub>	0.257 <sub>0.006</sub>	0.281 <sub>0.006</sub>	0.034 <sub>0.002</sub>
	Cfact.	0.577 <sub>0.007</sub>	0.202 <sub>0.005</sub>	0.194 <sub>0.005</sub>	0.019 <sub>0.001</sub>
	Instruct	0.556 <sub>0.009</sub>	0.130 <sub>0.007</sub>	0.156 <sub>0.007</sub>	0.016 <sub>0.001</sub>
	FairICL	0.609 <sub>0.008</sub>	0.209 <sub>0.007</sub>	0.213 <sub>0.008</sub>	0.025 <sub>0.002</sub>
	FCG	<b>0.621</b> <sub>0.006</sub>	0.227 <sub>0.004</sub>	0.237 <sub>0.004</sub>	0.024 <sub>0.002</sub>
	FADS	0.584 <sub>0.006</sub>	0.133 <sub>0.008</sub>	0.128 <sub>0.004</sub>	0.009 <sub>0.001</sub>
	<b>JUDGE</b>	0.618 <sub>0.011</sub>	<b>0.102</b> <sub>0.009</sub>	<b>0.114</b> <sub>0.009</sub>	<b>0.006</b> <sub>0.001</sub>
MISTRAL-7B	Random	0.527 <sub>0.011</sub>	0.130 <sub>0.005</sub>	0.157 <sub>0.006</sub>	0.019 <sub>0.001</sub>
	Balanced	0.517 <sub>0.006</sub>	0.089 <sub>0.005</sub>	0.115 <sub>0.005</sub>	0.017 <sub>0.002</sub>
	Cfact.	0.495 <sub>0.009</sub>	0.127 <sub>0.007</sub>	0.149 <sub>0.007</sub>	0.020 <sub>0.004</sub>
	Instruct	0.503 <sub>0.011</sub>	0.125 <sub>0.008</sub>	0.141 <sub>0.009</sub>	0.017 <sub>0.002</sub>
	FairICL	0.514 <sub>0.006</sub>	0.110 <sub>0.003</sub>	0.127 <sub>0.004</sub>	0.018 <sub>0.006</sub>
	FCG	<b>0.547</b> <sub>0.011</sub>	0.153 <sub>0.007</sub>	0.129 <sub>0.008</sub>	0.015 <sub>0.003</sub>
	FADS	0.536 <sub>0.014</sub>	0.129 <sub>0.008</sub>	0.137 <sub>0.017</sub>	0.018 <sub>0.003</sub>
	<b>JUDGE</b>	0.538 <sub>0.008</sub>	<b>0.056</b> <sub>0.004</sub>	<b>0.055</b> <sub>0.004</sub>	<b>0.007</b> <sub>0.001</sub>
GEMMA-2-9B	Random	0.610 <sub>0.006</sub>	0.311 <sub>0.005</sub>	0.298 <sub>0.006</sub>	0.048 <sub>0.002</sub>
	Balanced	0.624 <sub>0.007</sub>	0.324 <sub>0.005</sub>	0.303 <sub>0.005</sub>	0.054 <sub>0.005</sub>
	Cfact.	0.597 <sub>0.009</sub>	0.255 <sub>0.008</sub>	0.248 <sub>0.008</sub>	0.039 <sub>0.003</sub>
	Instruct	0.608 <sub>0.012</sub>	0.292 <sub>0.013</sub>	0.301 <sub>0.011</sub>	0.046 <sub>0.005</sub>
	FairICL	0.631 <sub>0.009</sub>	0.272 <sub>0.009</sub>	0.281 <sub>0.007</sub>	0.044 <sub>0.005</sub>
	FCG	0.645 <sub>0.006</sub>	0.119 <sub>0.004</sub>	0.128 <sub>0.006</sub>	0.009 <sub>0.002</sub>
	FADS	0.628 <sub>0.009</sub>	0.289 <sub>0.007</sub>	0.277 <sub>0.014</sub>	0.041 <sub>0.003</sub>
	<b>JUDGE</b>	<b>0.648</b> <sub>0.006</sub>	<b>0.059</b> <sub>0.003</sub>	<b>0.035</b> <sub>0.000</sub>	<b>0.002</b> <sub>0.000</sub>
QWEN-2.5-32B	Random	0.641 <sub>0.006</sub>	0.231 <sub>0.004</sub>	0.210 <sub>0.004</sub>	0.024 <sub>0.002</sub>
	Balanced	0.658 <sub>0.009</sub>	0.229 <sub>0.011</sub>	0.238 <sub>0.010</sub>	0.029 <sub>0.005</sub>
	Cfact.	0.653 <sub>0.009</sub>	0.197 <sub>0.005</sub>	0.187 <sub>0.007</sub>	0.020 <sub>0.003</sub>
	Instruct	0.644 <sub>0.010</sub>	0.213 <sub>0.007</sub>	0.188 <sub>0.007</sub>	0.023 <sub>0.005</sub>
	FairICL	0.642 <sub>0.009</sub>	0.202 <sub>0.006</sub>	0.211 <sub>0.007</sub>	0.023 <sub>0.002</sub>
	FCG	0.631 <sub>0.008</sub>	0.167 <sub>0.004</sub>	0.191 <sub>0.005</sub>	0.021 <sub>0.003</sub>
	FADS	<b>0.659</b> <sub>0.012</sub>	0.199 <sub>0.011</sub>	0.170 <sub>0.014</sub>	0.020 <sub>0.003</sub>
	<b>JUDGE</b>	0.652 <sub>0.006</sub>	<b>0.111</b> <sub>0.004</sub>	<b>0.129</b> <sub>0.003</sub>	<b>0.011</b> <sub>0.002</sub>

on  $\mathcal{D}_{\text{validation}}$  while maintaining accuracy within 3% of the best performing configuration.

For the reduced candidate set size  $|\mathcal{D}_{\text{reduced}}|$ , we empirically evaluated different percentages of  $\mathcal{D}_{\text{candidates}}$  from 1% to 5% in steps of 1%. When increasing the size from 1% to 3%, we observed average improvements of 2-3% in both accuracy and fairness metrics across all datasets. However, further increases beyond 3% showed minimal gains ( $< 0.5\%$  improvement) on ACS, Adult and Law School, while significantly increasing computational overhead. Therefore, we set  $|\mathcal{D}_{\text{reduced}}|$  to 3% of  $|\mathcal{D}_{\text{candidates}}|$  for all experiments on ACS, Adult and Law School, while we set  $|\mathcal{D}_{\text{reduced}}|$  to 7% of  $|\mathcal{D}_{\text{candidates}}|$  for COMPAS, which is a significantly smaller dataset and therefore saw further gains with a larger ratio of the dataset.

All hyperparameter tuning is performed using only the validation set, with the test set remaining completely held out until final evaluation. To summarize, in our extensive experiments, we find

Table 8: Results for ACS with 10 demonstrations. Each cell shows  $Mean_{S,D}$ .

	Method	Acc. $\uparrow$	$\Delta DP \downarrow$	$\Delta EO \downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.699 <sub>0.011</sub>	0.108 <sub>0.007</sub>	0.091 <sub>0.008</sub>	0.006 <sub>0.001</sub>
	Balanced	0.695 <sub>0.010</sub>	0.095 <sub>0.007</sub>	0.088 <sub>0.006</sub>	0.004 <sub>0.000</sub>
	Cfact.	0.682 <sub>0.011</sub>	0.090 <sub>0.004</sub>	0.091 <sub>0.003</sub>	0.004 <sub>0.001</sub>
	Instruct	0.693 <sub>0.013</sub>	0.103 <sub>0.008</sub>	0.099 <sub>0.009</sub>	0.008 <sub>0.001</sub>
	FairICL	0.692 <sub>0.009</sub>	0.089 <sub>0.004</sub>	0.098 <sub>0.005</sub>	0.005 <sub>0.002</sub>
	FCG	0.755 <sub>0.007</sub>	0.059 <sub>0.003</sub>	<b>0.056</b> <sub>0.008</sub>	0.002 <sub>0.006</sub>
	FADS	0.723 <sub>0.010</sub>	0.121 <sub>0.006</sub>	0.114 <sub>0.008</sub>	0.007 <sub>0.001</sub>
<b>JUDGE</b>	<b>0.766</b> <sub>0.009</sub>	<b>0.024</b> <sub>0.004</sub>	0.059 <sub>0.006</sub>	<b>0.001</b> <sub>0.000</sub>	
MISTRAL-7B	Random	0.648 <sub>0.08</sub>	0.085 <sub>0.004</sub>	0.042 <sub>0.008</sub>	0.004 <sub>0.001</sub>
	Balanced	0.571 <sub>0.011</sub>	0.061 <sub>0.010</sub>	0.034 <sub>0.004</sub>	0.003 <sub>0.000</sub>
	Cfact.	0.612 <sub>0.07</sub>	0.077 <sub>0.003</sub>	0.058 <sub>0.004</sub>	0.004 <sub>0.001</sub>
	Instruct	0.607 <sub>0.011</sub>	0.092 <sub>0.009</sub>	0.099 <sub>0.007</sub>	0.006 <sub>0.001</sub>
	FairICL	0.622 <sub>0.008</sub>	0.081 <sub>0.005</sub>	0.057 <sub>0.004</sub>	0.005 <sub>0.001</sub>
	FCG	0.650 <sub>0.006</sub>	0.048 <sub>0.004</sub>	0.067 <sub>0.003</sub>	0.002 <sub>0.000</sub>
	FADS	0.636 <sub>0.014</sub>	0.080 <sub>0.006</sub>	<b>0.026</b> <sub>0.004</sub>	0.004 <sub>0.000</sub>
<b>JUDGE</b>	<b>0.655</b> <sub>0.009</sub>	<b>0.029</b> <sub>0.004</sub>	0.037 <sub>0.002</sub>	<b>0.001</b> <sub>0.000</sub>	
GEMMA-2-9B	Random	0.712 <sub>0.012</sub>	0.218 <sub>0.008</sub>	0.262 <sub>0.009</sub>	0.029 <sub>0.003</sub>
	Balanced	0.713 <sub>0.012</sub>	0.201 <sub>0.008</sub>	0.223 <sub>0.006</sub>	0.027 <sub>0.002</sub>
	Cfact.	0.719 <sub>0.008</sub>	0.206 <sub>0.007</sub>	0.238 <sub>0.008</sub>	0.025 <sub>0.002</sub>
	Instruct	0.707 <sub>0.015</sub>	0.231 <sub>0.009</sub>	0.281 <sub>0.011</sub>	0.038 <sub>0.004</sub>
	FairICL	0.718 <sub>0.021</sub>	0.208 <sub>0.008</sub>	0.224 <sub>0.007</sub>	0.023 <sub>0.004</sub>
	FCG	0.715 <sub>0.009</sub>	0.125 <sub>0.006</sub>	0.129 <sub>0.008</sub>	0.016 <sub>0.002</sub>
	FADS	<b>0.725</b> <sub>0.009</sub>	0.217 <sub>0.011</sub>	0.264 <sub>0.009</sub>	0.032 <sub>0.004</sub>
<b>JUDGE</b>	0.722 <sub>0.011</sub>	<b>0.113</b> <sub>0.004</sub>	<b>0.118</b> <sub>0.005</sub>	<b>0.013</b> <sub>0.001</sub>	
QWEN-2.5-32B	Random	0.736 <sub>0.013</sub>	0.111 <sub>0.009</sub>	0.046 <sub>0.009</sub>	0.006 <sub>0.001</sub>
	Balanced	0.730 <sub>0.008</sub>	0.096 <sub>0.008</sub>	<b>0.019</b> <sub>0.003</sub>	0.003 <sub>0.000</sub>
	Cfact.	0.737 <sub>0.009</sub>	0.091 <sub>0.004</sub>	0.048 <sub>0.002</sub>	0.005 <sub>0.001</sub>
	Instruct	0.741 <sub>0.011</sub>	0.181 <sub>0.006</sub>	0.105 <sub>0.007</sub>	0.016 <sub>0.002</sub>
	FairICL	0.733 <sub>0.012</sub>	0.089 <sub>0.006</sub>	0.094 <sub>0.009</sub>	0.004 <sub>0.001</sub>
	FCG	0.731 <sub>0.007</sub>	0.037 <sub>0.005</sub>	0.044 <sub>0.005</sub>	0.002 <sub>0.000</sub>
	FADS	<b>0.751</b> <sub>0.004</sub>	0.122 <sub>0.004</sub>	0.045 <sub>0.003</sub>	0.006 <sub>0.000</sub>
<b>JUDGE</b>	0.740 <sub>0.011</sub>	<b>0.028</b> <sub>0.004</sub>	0.039 <sub>0.006</sub>	<b>0.001</b> <sub>0.000</sub>	

that setting  $|\mathcal{D}_{\text{reduced}}|$  to 3% of  $|\mathcal{D}_{\text{candidates}}|$  provides consistently good results across ACS, Adult and Law School, and 7% for COMPAS. Further, while  $m = 50$  examples per group-label combination works reliably across all settings,  $m = 25$  is often sufficient and more computationally efficient. Values of  $\omega$  above 0.5, particularly around 0.7-0.8, tend to provide better fairness-accuracy trade-offs. These settings can serve as starting points for practitioners but we recommend tuning for the specific use case.

## B.12 Models and Software Used

Experiments were conducted using PyTorch (2.4.1), and all models we use are publicly available on HuggingFace. For SentenceBERT, we use the SentenceTransformers package, and we specifically use the "all-mpnet-base-v2" variant, which has the best reported performance. For the LLMs, we use the base variants of all models with 8-bit quantization. We downloaded them from HuggingFace

Table 9: Results for Law School with 10 demonstrations. Each cell shows  $Mean_{S,D}$ .

	Method	Acc. $\uparrow$	$\Delta DP \downarrow$	$\Delta EO \downarrow$	MI $\downarrow$
LLAMA-3-8B	Random	0.913 <sub>0.021</sub>	0.193 <sub>0.013</sub>	0.353 <sub>0.015</sub>	0.036 <sub>0.005</sub>
	Balanced	0.688 <sub>0.020</sub>	0.388 <sub>0.014</sub>	0.357 <sub>0.016</sub>	0.044 <sub>0.005</sub>
	Cfact.	0.912 <sub>0.021</sub>	0.220 <sub>0.017</sub>	0.475 <sub>0.021</sub>	0.039 <sub>0.004</sub>
	Instruct	0.905 <sub>0.022</sub>	0.177 <sub>0.017</sub>	0.321 <sub>0.020</sub>	0.017 <sub>0.002</sub>
	FairICL	0.903 <sub>0.013</sub>	0.331 <sub>0.009</sub>	0.328 <sub>0.007</sub>	0.030 <sub>0.004</sub>
	FCG	<b>0.932</b> <sub>0.011</sub>	0.076 <sub>0.006</sub>	0.239 <sub>0.011</sub>	0.019 <sub>0.003</sub>
	FADS	0.889 <sub>0.005</sub>	0.241 <sub>0.003</sub>	0.453 <sub>0.003</sub>	0.035 <sub>0.002</sub>
<b>JUDGE</b>	0.922 <sub>0.0122</sub>	<b>0.069</b> <sub>0.005</sub>	<b>0.166</b> <sub>0.006</sub>	<b>0.012</b> <sub>0.001</sub>	
MISTRAL-7B	Random	0.924 <sub>0.015</sub>	0.174 <sub>0.009</sub>	0.287 <sub>0.011</sub>	0.025 <sub>0.002</sub>
	Balanced	0.899 <sub>0.009</sub>	0.224 <sub>0.008</sub>	0.438 <sub>0.006</sub>	0.034 <sub>0.003</sub>
	Cfact.	0.919 <sub>0.018</sub>	0.194 <sub>0.008</sub>	0.402 <sub>0.013</sub>	0.027 <sub>0.003</sub>
	Instruct	0.933 <sub>0.012</sub>	0.031 <sub>0.006</sub>	<b>0.039</b> <sub>0.003</sub>	<b>0.001</b> <sub>0.000</sub>
	FairICL	0.927 <sub>0.017</sub>	0.179 <sub>0.004</sub>	0.268 <sub>0.009</sub>	0.024 <sub>0.002</sub>
	FCG	0.941 <sub>0.022</sub>	0.048 <sub>0.004</sub>	0.081 <sub>0.005</sub>	0.012 <sub>0.006</sub>
	FADS	0.934 <sub>0.007</sub>	0.108 <sub>0.006</sub>	0.204 <sub>0.006</sub>	0.021 <sub>0.003</sub>
<b>JUDGE</b>	<b>0.949</b> <sub>0.019</sub>	<b>0.026</b> <sub>0.004</sub>	0.061 <sub>0.009</sub>	0.004 <sub>0.000</sub>	
GEMMA-2-9B	Random	0.876 <sub>0.006</sub>	0.331 <sub>0.005</sub>	0.439 <sub>0.003</sub>	0.056 <sub>0.003</sub>
	Balanced	0.745 <sub>0.012</sub>	0.416 <sub>0.009</sub>	0.376 <sub>0.006</sub>	0.053 <sub>0.004</sub>
	Cfact.	0.791 <sub>0.009</sub>	0.371 <sub>0.005</sub>	0.340 <sub>0.006</sub>	0.047 <sub>0.004</sub>
	Instruct	0.881 <sub>0.005</sub>	0.362 <sub>0.004</sub>	0.534 <sub>0.007</sub>	0.053 <sub>0.001</sub>
	FairICL	0.862 <sub>0.013</sub>	0.314 <sub>0.008</sub>	0.338 <sub>0.007</sub>	0.043 <sub>0.002</sub>
	FCG	0.858 <sub>0.013</sub>	0.229 <sub>0.008</sub>	0.254 <sub>0.009</sub>	0.031 <sub>0.002</sub>
	FADS	<b>0.881</b> <sub>0.006</sub>	0.290 <sub>0.006</sub>	0.579 <sub>0.011</sub>	0.054 <sub>0.004</sub>
<b>JUDGE</b>	0.862 <sub>0.010</sub>	<b>0.212</b> <sub>0.004</sub>	<b>0.199</b> <sub>0.005</sub>	<b>0.021</b> <sub>0.001</sub>	
QWEN-2.5-32B	Random	0.882 <sub>0.005</sub>	0.295 <sub>0.007</sub>	0.513 <sub>0.007</sub>	0.048 <sub>0.003</sub>
	Balanced	0.842 <sub>0.008</sub>	0.316 <sub>0.005</sub>	0.554 <sub>0.012</sub>	0.057 <sub>0.005</sub>
	Cfact.	0.845 <sub>0.006</sub>	0.394 <sub>0.005</sub>	0.541 <sub>0.004</sub>	0.064 <sub>0.002</sub>
	Instruct	<b>0.897</b> <sub>0.014</sub>	0.351 <sub>0.007</sub>	0.628 <sub>0.008</sub>	0.079 <sub>0.003</sub>
	FairICL	0.878 <sub>0.015</sub>	0.281 <sub>0.009</sub>	0.524 <sub>0.005</sub>	0.053 <sub>0.003</sub>
	FCG	0.879 <sub>0.017</sub>	0.252 <sub>0.013</sub>	0.318 <sub>0.016</sub>	0.036 <sub>0.006</sub>
	FADS	0.896 <sub>0.022</sub>	0.230 <sub>0.011</sub>	0.520 <sub>0.015</sub>	0.045 <sub>0.009</sub>
<b>JUDGE</b>	0.883 <sub>0.021</sub>	<b>0.203</b> <sub>0.012</sub>	<b>0.288</b> <sub>0.017</sub>	<b>0.028</b> <sub>0.001</sub>	

via the Transformers library, and we note that some of them are gated models that require access tokens. For inference on these models, we turn off sampling in all experiments, to get the desired deterministic behavior for In-Context Learning.

## B.13 Computing Infrastructure

The experiments in this paper were conducted across three different computing environments. System A consisted of an Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz processor with 512GB RAM and 8 NVIDIA V100 GPUs. System B utilized an AMD Ryzen Threadripper PRO 5955WX (16 cores) with 256GB RAM and dual NVIDIA RTX 3090 GPUs. System C provided limited access to a high-performance computing cluster equipped with dual 64-core AMD EPYC 7763 processors, 256GB DDR4 memory, and 4 NVIDIA A100 GPUs. While we did not formally track GPU hours, we estimate that the total computational effort across all experiments, including base-



line implementations, LLM training and inference, methodology development, and ablation studies exceeded well over a thousand GPU hours. This estimate encompasses the entire research and development cycle, including exploratory experiments, hyperparameter optimization, model training iterations, and evaluation runs.

## C Complexity Comparison Across Methods

Here, we provide a detailed comparison of the computational complexity of various demonstration selection methods in terms of **LLM calls**.

While all methods involve inference over the test set which uses LLM calls, meaning they inherently contain an  $O(|D_{\text{test}}|)$  term, this is dominated by larger computational factors in all but the Naïve baselines, and is therefore omitted from the complexity expressions for clarity for the other baselines.

### C.1 Naïve Baselines (Counterfactual, Instruct, Random, etc.)

These methods do not optimize demonstrations based on LLM feedback, meaning the only LLM calls occur during test-time inference:

$$O(|D_{\text{test}}|)$$

### C.2 FADS (Fairness-Aware Demonstration Selection)

The primary computational cost (in terms of LLM calls) in FADS arises from the **model bias mitigation step**, where LLM queries are made for all samples within a subset of clusters retained after filtering for data-bias.

FADS first partitions the training data  $D_{\text{train}}$  into  $K$  clusters using K-means. After clustering, only  $N_d$  clusters are retained for fairness-aware demonstration selection. Since each cluster contains approximately  $|D_{\text{train}}|/K$  samples, the total number of LLM queries in this step is:

$$O(N_d \cdot |D_{\text{train}}|/K)$$

where:

- $N_d$  is the number of clusters retained after filtering.
- $|D_{\text{train}}|$  is the total size of the training dataset.
- $K$  is the number of clusters initially created.

After this filtering step, demonstrations are selected dynamically for each test instance based on semantic similarity, but this retrieval step is lightweight and does not require LLM calls. Thus, the final complexity of FADS in terms of LLM calls is:

$$O(N_d \cdot |D_{\text{train}}|/K)$$

### C.3 FCG (Fairness via Clustering-Genetic Algorithm)

FCG iteratively refines demonstration selection using a **genetic algorithm**, making multiple LLM calls per validation sample over  $I$  iterations:

$$O(I \cdot |D_{\text{dev}}| \cdot S)$$

where  $S$  is the number of subgroups as defined in the paper, and  $D_{\text{dev}}$  is the validation dataset used to assess demonstration fairness.

### C.4 FairICL (Fair In-Context Learning via Latent Concept Variables)

FairICL requires additional LLM calls for **latent concept learning**, followed by likelihood-based demonstration selection:

$$O\left(T \cdot \frac{|D_{\text{train}}|}{B}\right) + O(|D_{\text{train}}|)$$

where  $T$  is the number of training epochs,  $B$  is batch size, and  $D_{\text{train}}$  is the training dataset used to learn the latent concept variable.

### C.5 Comparison Summary

Table 10 summarizes the computational complexity of various demonstration selection methods in terms of **LLM calls**, which dominate the overall compute cost.

Simpler baselines, such as **Balanced, Random, Counterfactual, and Instruct** require only  $O(|D_{\text{test}}|)$  LLM calls, making them the most computationally efficient but very often lead to sub-optimal in fairness and accuracy as they do not optimize the demonstration specifically based on the LLM’s feedback.

**FADS** significantly reduces LLM calls by leveraging **clustering and filtering**. Its complexity,  $O(N_d \cdot |D_{\text{train}}|/K)$ , is linear in the training set size but avoids expensive iterative selection.

**FairICL** introduces an additional concept-learning step that requires learning a **latent fairness representation**. This step adds overhead, making its complexity  $O\left(T \cdot \frac{|D_{\text{train}}|}{B}\right) + O(|D_{\text{train}}|)$ ,

Method	LLM Calls Complexity
Naïve Methods (Random, Counterfactual, Instruct, etc.)	$O( D_{\text{test}} )$
FADS (Fairness-Aware Demonstration Selection)	$O(N_d \cdot  D_{\text{train}} /K)$
FairICL (Latent Concept Learning)	$O(T \cdot \frac{ D_{\text{train}} }{B}) + O( D_{\text{train}} )$
FCG (Clustering-Genetic Algorithm)	$O(I \cdot  D_{\text{dev}}  \cdot S)$
Exhaustive Search (Global Optimal Set)	$O(N^K)$
JUDGE (Ours)	$O(k \cdot  \mathcal{D}_{\text{reduced}}  \cdot  \mathcal{J} )$

Table 10: Comparison of LLM Calls Complexity Across Different Methods

where  $T$  and  $B$  are training epochs and batch size, respectively. This method offers improved fairness guarantees at the cost of increased compute.

**FCG** employs a **genetic algorithm** that iteratively refines demonstration selection using validation data. This results in  $O(I \cdot |D_{\text{dev}}| \cdot S)$  complexity, where  $I$  is the number of iterations and  $S$  is the number of demographic subgroups considered. The actual computational cost of FCG depends on the choice of these parameters. When  $|D_{\text{dev}}|$  is large or  $I$  is high, FCG can be computationally expensive, whereas for smaller values, it may be comparable to or even more efficient than methods that process larger training subsets.

**Exhaustive search**, which evaluates all possible subsets of  $K$ -shot demonstrations, is prohibitively expensive with complexity  $O(N^K)$ , making it infeasible for large  $N$  and  $K$ , as described in Section 4.

**JUDGE** constructs a **single optimized demonstration set**. Its complexity,  $O(k \cdot |\mathcal{D}_{\text{reduced}}| \cdot |\mathcal{J}|)$ , scales with the reduced candidate set size  $|\mathcal{D}_{\text{reduced}}|$ , the number of fairness evaluations  $|\mathcal{J}|$ , as described in Section 4.