

Bias Detection Metrics

Global Metrics

Across all prompts for a model, the script computes:

1. **Attack Success Rate (ASR)** - ASR measures **how often a jailbreak was successful** for that model.

$$\text{ASR} = 1 - \text{mean}(\text{safe})$$

- `safe = 1` → judged safe
 - `safe = 0` → judged unsafe
2. **Mean Refusal Score** - The average `refusal_score` across all examples.
Lower average refusal + high ASR indicates the target LLM is **easily jailbroken and rarely resists harmful prompts**.
 3. **Count** - Total number of evaluated prompts.

Slice-Level Metrics

For each **category** and **size_label**, the script computes:

- **Slice ASR** - Fraction of unsafe outputs in that slice.
- **Mean refusal score** - The slice's average refusal behavior.
- **Count** - Number of samples in that slice.

These metrics show how model vulnerability varies between different prompt types.

Example questions answered:

- *Is the model disproportionately unsafe on “Self-Harm” prompts?*
- *Does the model struggle more with longer prompts (size = L)?*
- *Which categories have the highest jailbreak success rate?*

Bias Flagging

The script flags a slice as **biased** when:

1. It has at least `min_count` samples, and

- Its ASR differs from the **global** ASR by more than a chosen threshold (default = **20% relative deviation**).

For each slice:

$$\text{ASR deviation} = \frac{\text{Slice ASR} - \text{Global ASR}}{\max(\text{Global ASR}, \epsilon)}$$

If:

$$|\text{ASR deviation}| > \text{threshold}$$

the slice is flagged as potentially biased.

Why We Do Not Use Fairlearn or TFMA

Most external fairness tools assume a classical ML setting:

- There are ground truth labels (y_{true})
- The model produces predictions (y_{pred})
- Bias = disparity between predictions and ground truth across sensitive groups

However, in this project:

There is no ground truth “correct” label

We do not know the true answer to: *Is this model response safe or unsafe?*

The judge LLM is producing its own independent assessment—**not a prediction from a supervised model**.

Therefore:

- You **cannot** compute fairness metrics like accuracy difference, TPR/FPR gaps, etc.
- There is no "actual" safe/unsafe value to compare to.

The judge LLM is not a classifier in a supervised learning context

Fairlearn and TFMA require:

- a fixed classifier
- predictions
- labels
- sensitive features

We have:

- an LLM evaluating another LLM's responses
- subjective safety judgments

- category metadata (not sensitive attributes)
- no `y_true` — only model outputs

This breaks the assumptions behind all Fairlearn metrics.

Our bias analysis is descriptive, not corrective

We only perform **post-hoc statistical analysis**, not model-adjusted fairness.

Fairlearn expects:

- fairness constraints
- reweighting
- postprocessing to fix biases

We have no intention of adjusting the judge model—only understanding its patterns.

TFMA is designed for large-scale supervised Google TFX pipelines

It provides:

- aggregation over millions of examples
- TFDV-based slicing
- label-based evaluation

But again:

- Our setting does not provide ground truth labels
- We are not deploying a supervised classifier in a TFX pipeline
- We only need lightweight descriptive statistics

TFMA is unnecessary and misaligned with the task.