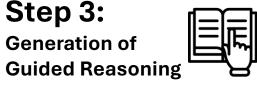
Step 2:

Shapley Value Extraction





Step 4:

Post-processing

Step 5: **Best Model** Selection



Step 6: Best Reasoning Selection





Step 1: **Dataset**

collection



Input:

 Toxic datasets: ParaDetox, APPDIA, Parallel Detoxification - Non-toxic dataset: Jigsaw

Unintended Bias

Process: Combine toxic and selected non-toxic samples

Outcome: Balanced toxic /

non-toxic dataset

Input:

- Models: HateBERT, Toxigen, Detoxtify, ToxDetect

- Balanced dataset

Process: Extract and aggregate Shapely values

Outcome: Enriched dataset with Shapley values

Input:

Reflect

- Models: Marco-o1, QwQ, OpenO1, Skywork, OpenAl o1 Enriched dataset **Process:** Generate reasoning using instruction + Shapley +

target paraphrase **Outcome:** Guided reasonings

Input:

and Filtering

Guided reasonings Process: Remove codeswitching (Google Translate), remove leakage and restating (Qwen 2.5) Outcome: Cleaned reasonings

Input:

Cleaned reasonings **Process:** Human A/B Evaluation, JudgeLLM A/B evaluation Outcome: Best performing models

Input:

dataset

Selection

Best-performing models' outputs **Process:** Evaluate using JudgeLLM and A/B tournament **Outcome:** Final few-shot

DetoxLLM)

Process: Instruction, n demonstrations, query generation, evaluation metrics (STA, TOX, J, FL,









Evaluation (ICL)

Balanced dataset, best-

Step 7:

Input:

performing models, Qwen 2.5 7B and LlaMA 3.1 8B, baselines (BART, PseudoParaDetox, CAPP, BERT-F1, BLEU, SIM)

Outcome: Comparative results of detoxtified outputs