

An Exploration of Modern Text Detoxification Pipelines

A Modular Framework Study

Benjamin He & Kent Bourgoing

Berkeley MIDS W266

Why Detoxification, Not Just Blocking?

The Challenge

Online platforms still struggle with toxic language

Pure blocking = lost context, feels like censorship

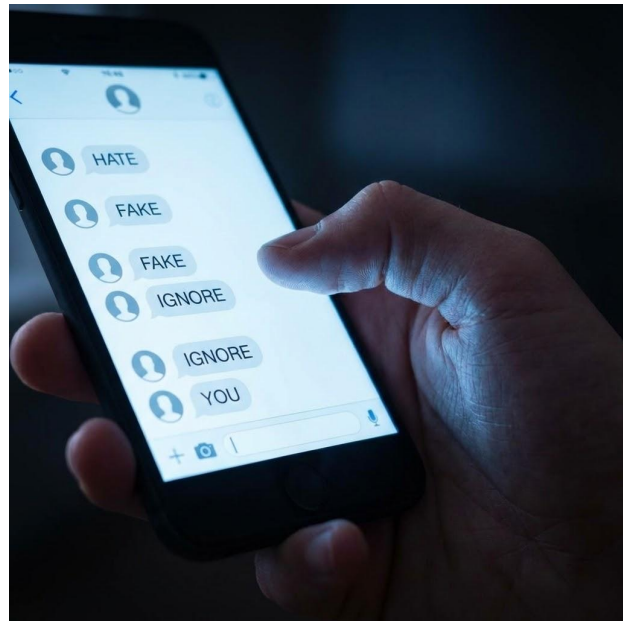
The Solution

Detoxification: rewrite toxic → non-toxic

Keep semantic meaning intact

Useful for:

- User-facing "gentler rewrite" suggestions
- Moderator tools and brand safety
- Pre/post-processing around LLMs



What Are We Trying to Learn?

Focus

Sentence-level detoxification (ParaDetox test set, 671 examples)

1. Masking

Compare masking strategies
DecompX-based vs LLM-based



2. Infilling

Compare infilling models
MaRCo vs Mistral-7B



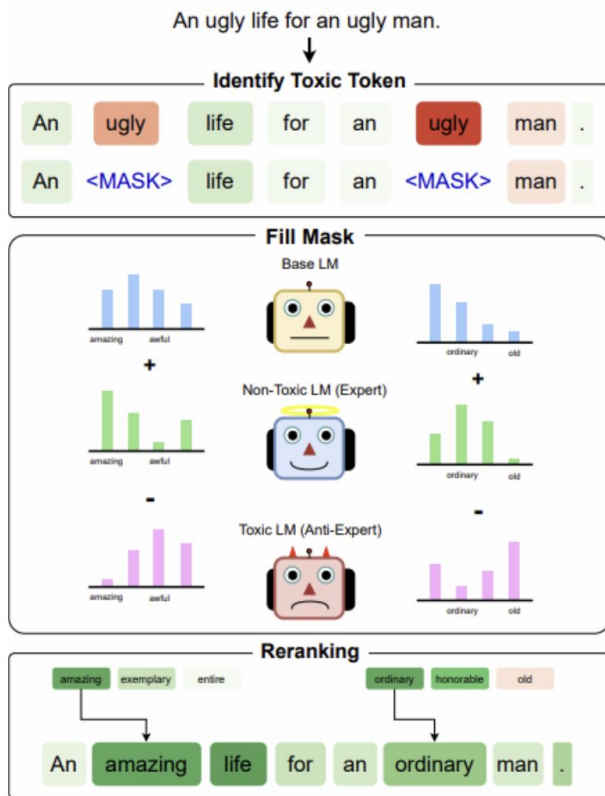
3. Reranking

Measure impact of reranking
DecompX vs Global reranker



Artifact: Reusable modular framework for detoxification pipelines

Modular Mask–Infill–Rerank Framework



1. Masker

Identify toxic spans → replace with <mask>

2. Infiller

Generate candidate rewrites

11 pipelines:

DecompX vs LLM masking ×

MaRCO vs LLM infilling ×

DecompX vs Global reranking

Models, Data, and Metrics

Base Model

T5-base fine-tuned on ParaDetox

Maskers

DecompX + RoBERTa (threshold 0.2)
Mistral-7B Instruct as LLM masker

Infillers

MaRCO (BART expert/anti-expert)
Mistral-7B as infiller

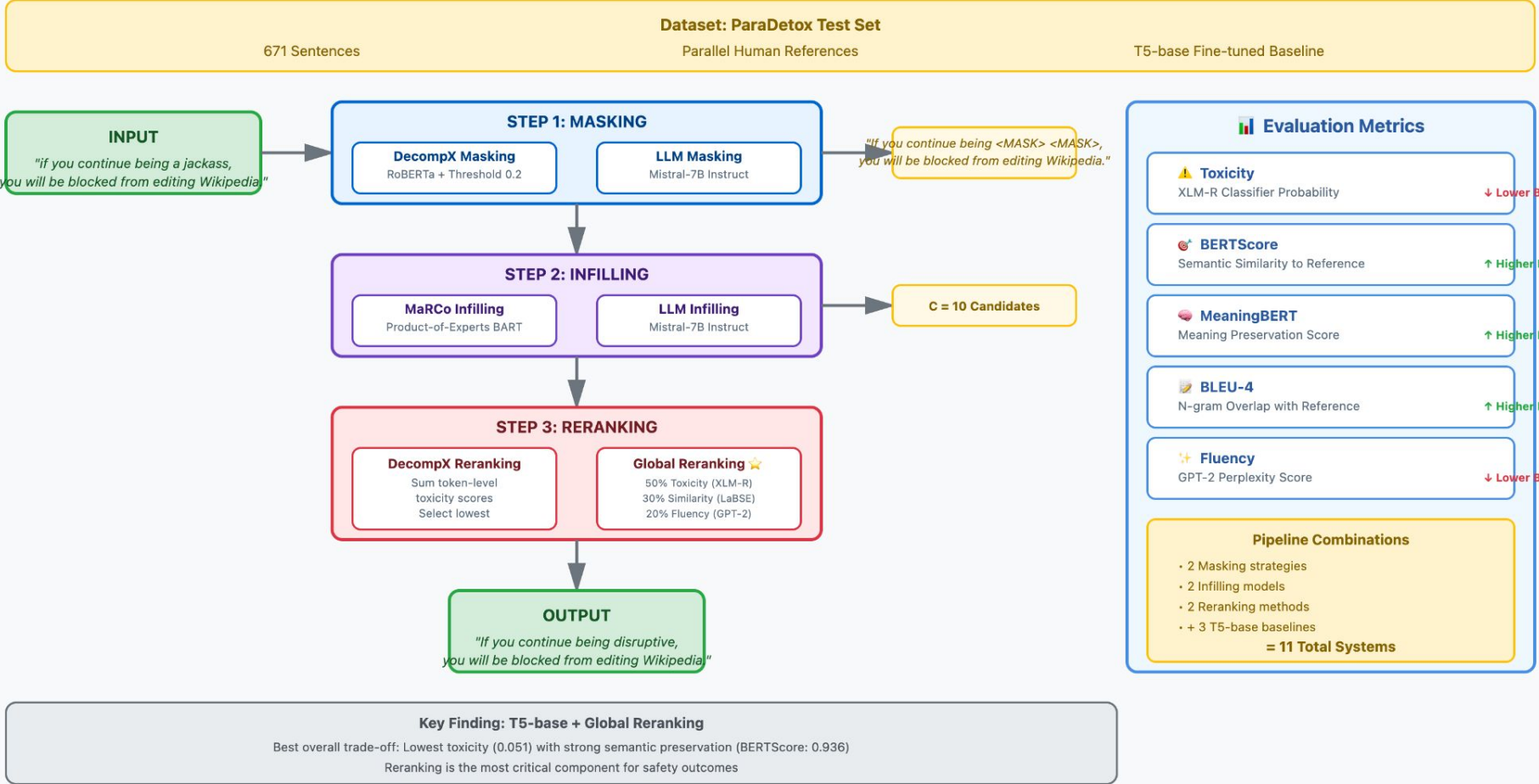
Rerankers

DecompX toxicity-sum
Global: Toxicity + Similarity + Fluency

Evaluation Metrics

- Toxicity (XLM-R)
- BERTScore
- MeaningBERT
- BLEU-4
- Perplexity (GPT-2)

Text Detoxification Pipeline: Models, Data, and Metrics



Text Detoxification Pipeline Results

Model	BERTScore	MeaningBERT	BLEU-4	Perplexity	Toxicity
T5-base	0.953	74.84	82.65	192.07	0.203
T5-base + DecompX Reranking	<u>0.947</u>	71.48	88.23	235.22	0.208
T5-base + Global Reranking	<u>0.936</u>	67.25	53.34	171.53	0.051
DecompX Masking + MaRCO Infilling + DecompX Reranking	0.944	<u>72.85</u>	68.99	136.08	0.132
DecompX Masking + MaRCO Infilling + Global Reranking	0.944	72.72	70.05	124.95	0.120
DecompX Masking + LLM Infilling + DecompX Reranking	0.938	66.16	<u>82.86</u>	200.29	0.171
DecompX Masking + LLM Infilling + Global Reranking	0.932	64.74	81.54	162.39	<u>0.103</u>
LLM Masking + MaRCO Infilling + DecompX Reranking	0.938	69.55	70.05	<u>90.65</u>	0.200
LLM Masking + MaRCO Infilling + Global Reranking	0.938	69.02	70.05	86.59	0.159
LLM Masking + LLM Infilling + DecompX Reranking	0.931	62.55	81.54	149.22	0.181
LLM Masking + LLM Infilling + Global Reranking	0.931	62.45	81.54	141.89	0.118

Main Quantitative



Best Overall: T5-base + Global Reranking

Lowest toxicity (0.051), slight drop in similarity, best safety–meaning trade-off

Model	BERTScore	MeaningBERT	BLEU-4	Perplexity	Toxicity
T5-base	0.953	74.84	82.65	192.07	0.203
T5 + DecompX Rerank	0.947	71.48	88.23	235.22	0.208
T5 + Global Rerank	0.936	67.25	53.34	171.53	0.051
DecompX + MaRCO + DecompX	0.944	72.85	68.99	136.08	0.132
DecompX + MaRCO + Global	0.944	72.72	70.05	124.95	0.120
DecompX + LLM + Global	0.932	64.74	81.54	162.39	0.103
LLM + LLM + Global	0.931	62.45	81.54	141.89	0.118

Masking Impact

DecompX → consistently lower toxicity,

Infilling Impact

LLM infilling generally safer than MaRCO

Key Insight

Global Reranker consistently improves safety for all generators/maskers

What Do the Outputs Actually Look Like?

T5-base without reranking

Keeps meaning but often adds new insults/profanity

T5-base + Global Reranking

Strong slurs and threats almost disappear

Residual: mild snark, odd paraphrases

MaRCo infilling

Fluent but problematic:

Can introduce severe slurs, graphic content, threats

LLM infilling

Safer templates:

"disrespectful person", "hurtful language"

Still some dehumanizing language, mild profanity

Trade-off: Perfect semantic overlap vs reduced toxicity

Conclusion, Limitations, Next Steps

Contributions from this work:

Modular detoxification framework (mask–infill–rerank)

Systematic comparison of 11 pipelines

Evidence that **global reranking is strong guardrail for safety**



Limitations

Single English benchmark (ParaDetox only)

toxic classifier may be biased

Moderate-size LLM (Mistral-7B) due to compute limits

Future Work

Learned rerankers optimizing toxicity + meaning

Stronger LLMs and more datasets

Better masking: DecompX + LLM judgments

Thank you! Questions?

