

---

# ICLEF: In-Context Learning with Expert Feedback for Explainable Style Transfer

---



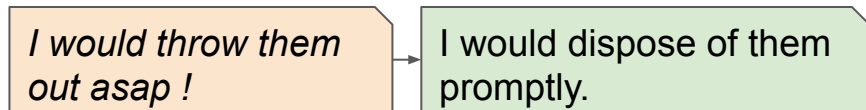
Arkadiy Saakyan  
Columbia University



Smaranda Muresan  
Columbia University

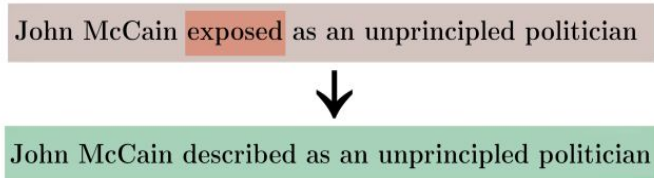
# Attribute style transfer

Current:



Formality

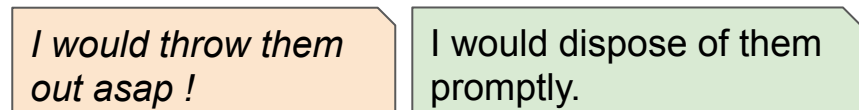
Dear Sir or Madam, May I introduce the GYAFC Dataset: Corpus, Benchmarks and Metrics for Formality Style Transfer  
Sudha Rao, Joel Tetreault



Bias

Automatically Neutralizing Subjective Bias in Text  
Reid Pryzant, Richard Diehl Martinez, Nathan Dass, Sadao Kurohashi, Dan Jurafsky, Diyi Yang

What we want:



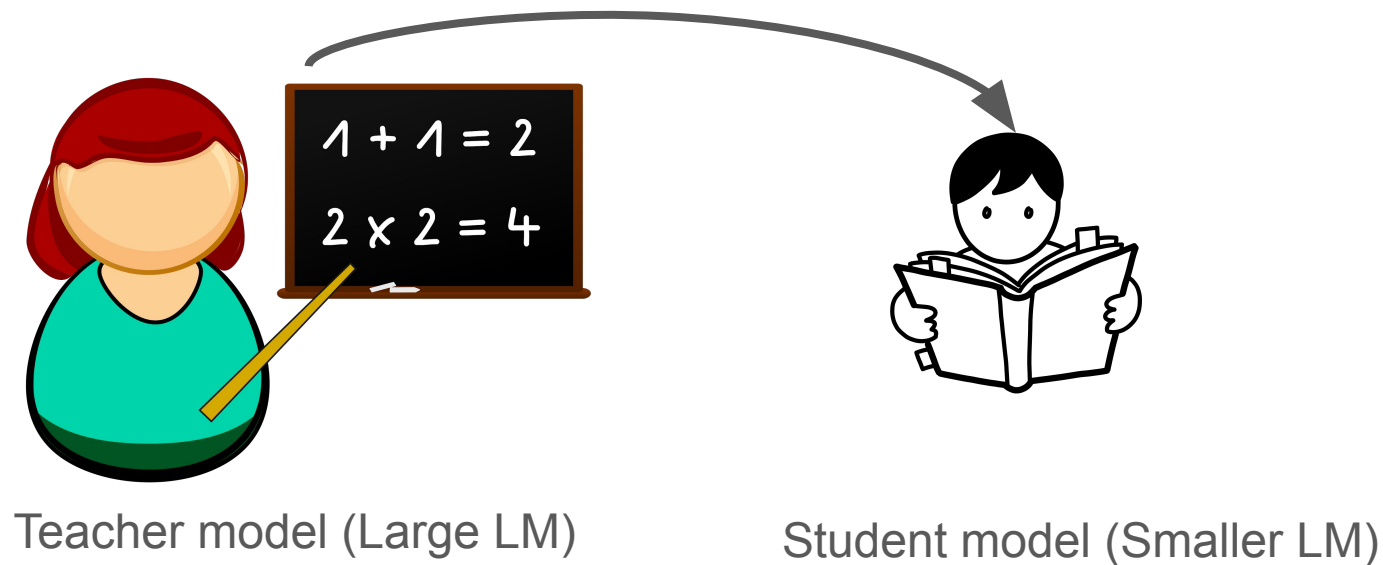
**Informal Attributes:**  
textese ("asap"),  
colloquialism ("throw out"),  
exclamation mark



**Bias Attributes:**  
Epistemological ("Exposed" is a factive verb that presupposes the truth of its complement that McCain is unprincipled)

# Model distillation framework

- Teacher model generates explanations
- Generations from teacher model are used to fine-tune smaller models



# Key Challenges

1. Necessity to rely on large models
2. Explanations may not be accurate
3. Scarcity of expert feedback

## Proposed Solution:

- Human-AI collaboration approach to *model distillation (ICLEF)* to create such datasets:
  - 1) distill natural language explanations from large LMs
  - 2) leverage in-context learning and self-critique abilities of LLMs
- Advantage: we can use this dataset to fine-tune ***smaller models for explainable style transfer***

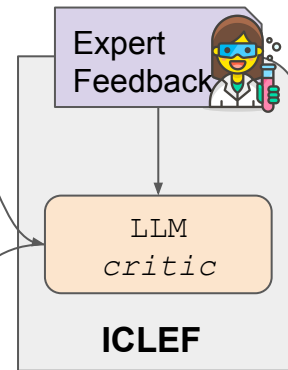
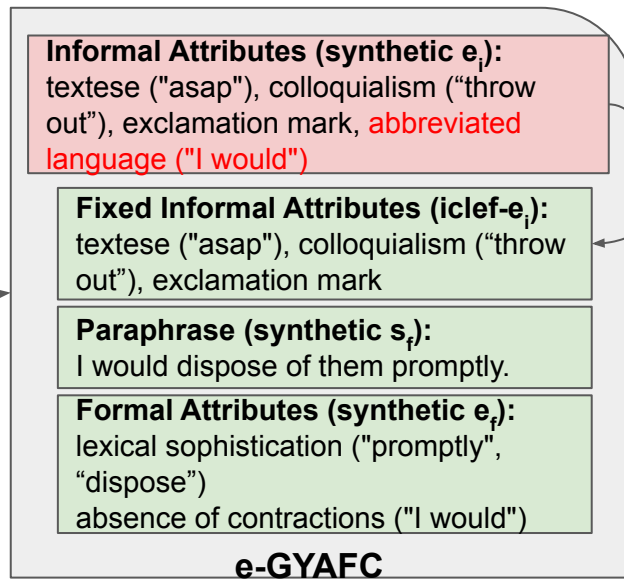
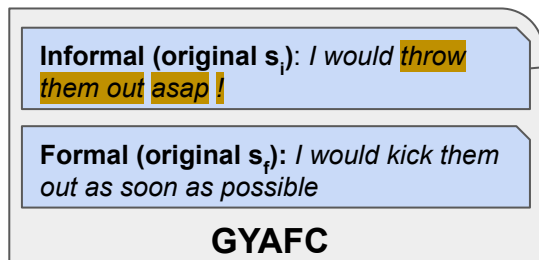
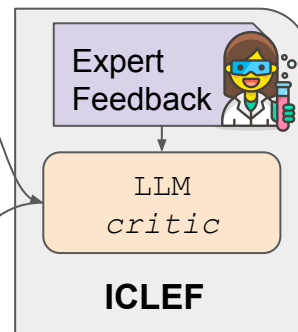
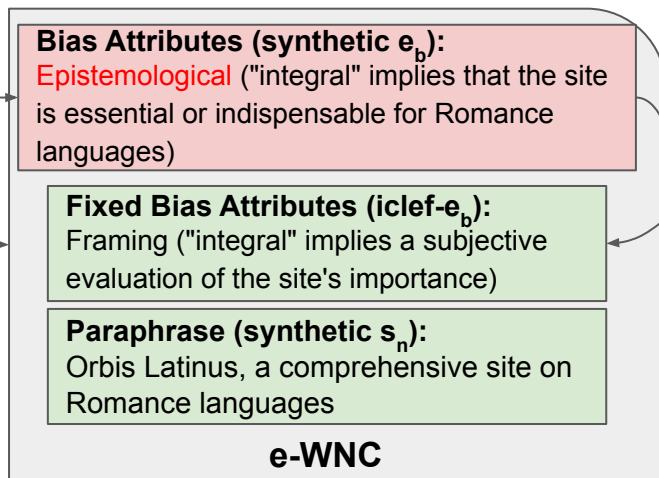
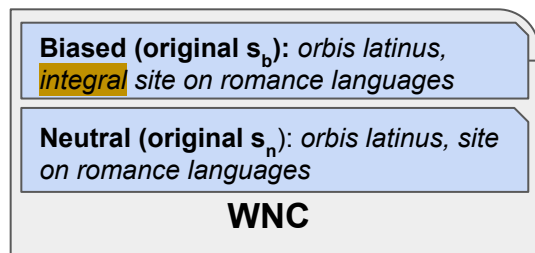
# Explainability

- Generating and Distilling natural language explanations:
  - [PINTO: Faithful Language Reasoning Using Prompt-Generated Rationales](#) (Wang et al., 2023)
  - [Large Language Models Are Reasoning Teachers](#) (Ho et al., 2022)
  - [Teaching Small Language Models to Reason](#) (Magister et al., 2023)
- Self-critique
  - [Self-Refine: Iterative Refinement with Self-Feedback](#) (Madaan et al., 2023)
  - [Constitutional AI: Harmlessness from AI Feedback](#) (Bai et al., 2022)
  - [Self-critiquing models for assisting human evaluators](#) (Saunders et al., 2022)
  - [Training Language Models with Language Feedback at Scale](#) (Scheurer et al., 2023)

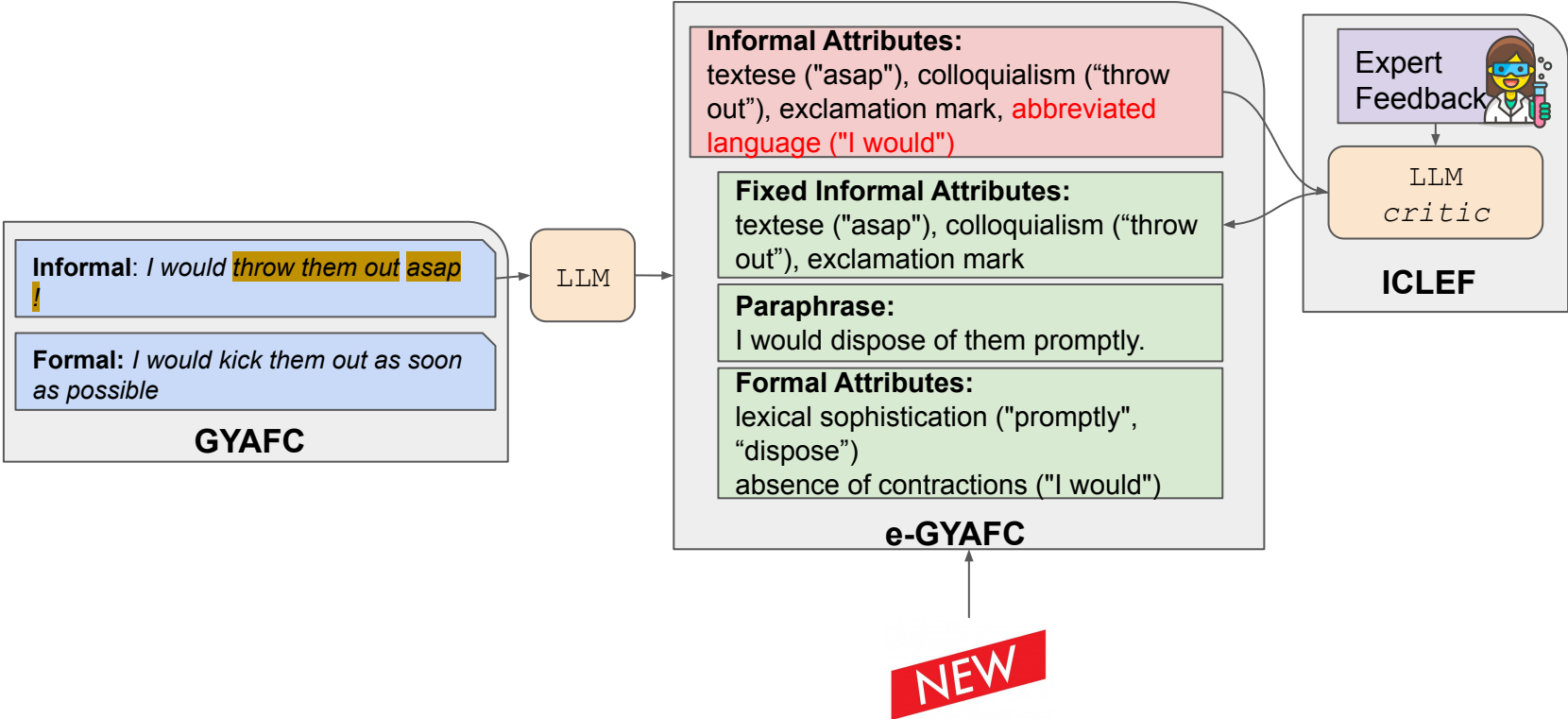
# Learning from human feedback

- Generated explanations might still have mistakes. How to leverage feedback to improve quality?
  - [RLHF](#) (Ouyang, 2022)
  - [Chain-of-Hindsight](#) (Liu et al., 2023),
  - [Sequence Likelihood Calibration](#) (Zhao et al., 2023),
  - [Direct Preference Optimization](#) (Rafailov et al., 2023)
- All work with large preference corpora from crowdworkers, but what to do with scarce expert feedback?

Taxonomy from Pryzant, et. al 2019

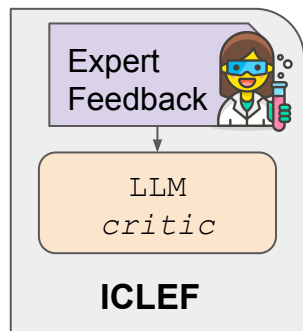


# Example for formality style transfer





# Incorporating expert feedback



Few-shot prompt with demonstrations of expert annotator corrections:

e-GYAFC:

...

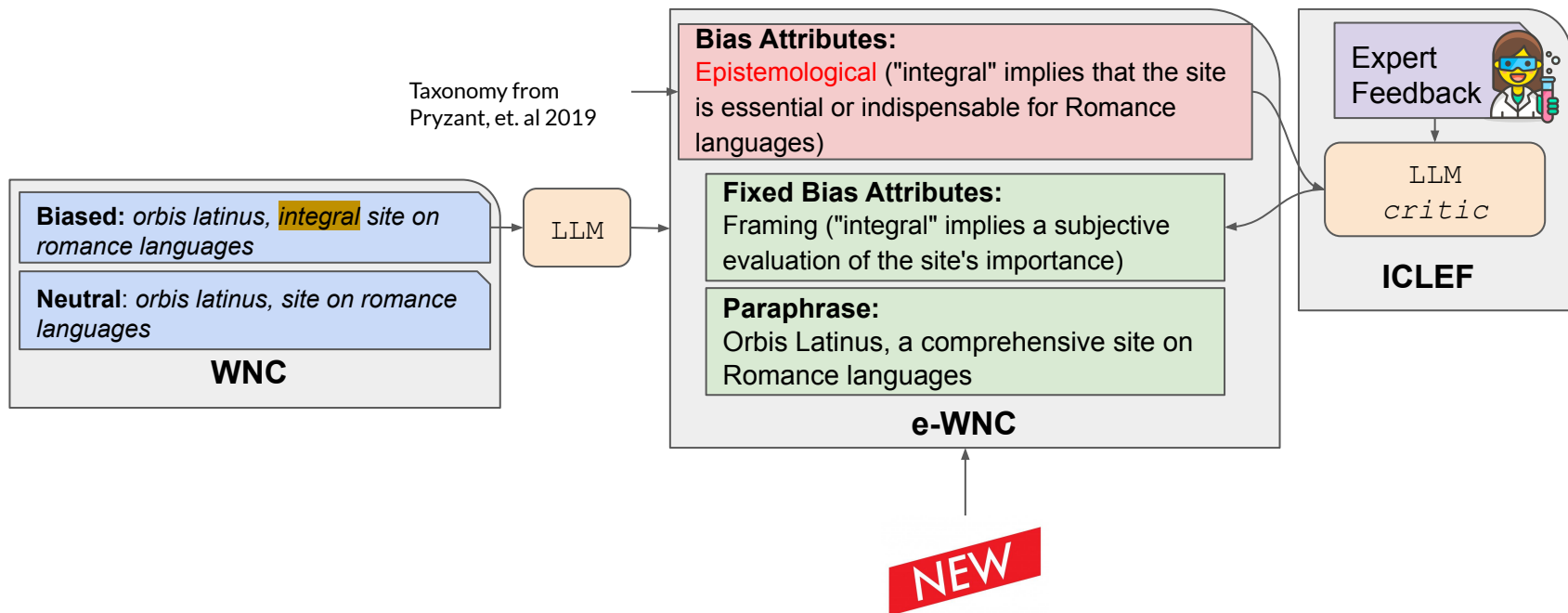
Informal Sentence: Look, If you really like this person, just tell her.

Informal Attributes: colloquialism ("just tell her"), contraction ("If you"), simple sentence structure

Attributes Listed Incorrectly: contraction ("If you" is not a contraction)

...

# Similarly for *biased* -> *unbiased* style transfer



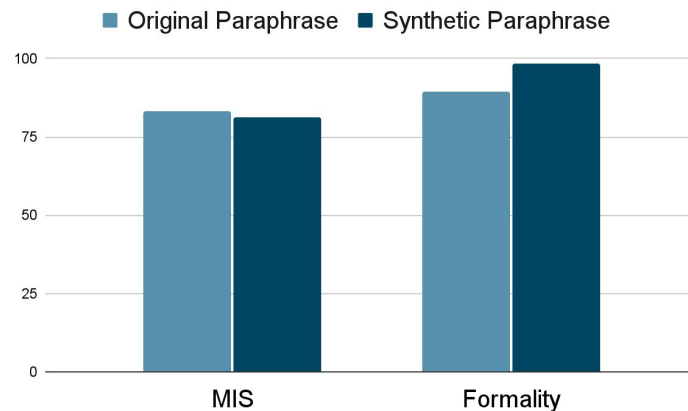
## Examples of ICLEF editing:

<b>Informal</b> ( $s_i$ )	<b>Gen. expl.</b> (synthetic $e_i$ )	<b>ICLEF expl.</b> ( $iclef-e_i$ )
hopefully you aren't too old or you are screwed.	<b>informal greeting ("hopefully")</b> , slang ("screwed"), contraction ("aren't")	slang ("screwed"), contraction ("aren't")
more info, we are both in our very late twenties.	[...], <b>omission of prepositions ("in our very late twenties")</b>	abbreviation ("info"), colloquialism ("very late twenties")
<b>Biased</b> ( $s_b$ )	<b>Gen. expl.</b> (synthetic $e_b$ )	<b>ICLEF expl.</b> ( $iclef-e_b$ )
[...] a play on the title of the popular mtv series, "unplugged".	<b>Epistemological ("popular" implies that the MTV series is universally well-liked)</b>	Framing ("popular" is a subjective term that implies the MTV series is widely liked)
[...] kendal, cbe (born 25 september 1946) is an english actress known in the united kingdom [...].	<b>Demographic ("actress" implies that the person is female and could perpetuate gender stereotypes or assumptions)</b>	This sentence does not contain bias.
claims for the existence of paranormal psychic abilities such as clairvoyance are highly controversial.	<b>This sentence does not contain bias.</b>	Epistemological ("highly controversial" implies that the existence of paranormal psychic abilities is widely disputed)

Table 1: Qualitative comparison of dataset instances before and after application of ICLEF.

# Dataset Quality

- Automatic Evaluation
  - **MIS**: semantic accuracy of the paraphrase
  - **Formality**: style accuracy



- Human Evaluation
  - **Acceptability**: does the explanation list correct attributes and does the paraphrase have the same meaning?
  - **Preference**: is the iclef-explanation preferred to the original one? Is the synthetic paraphrase preferred to the original one?

	Acceptability	Preference
<b>Informal explanation</b>	87	90
<b>Formal paraphrase</b>	77	77
<b>Formal explanation</b>	98	N/A

# Evaluation on the Explainable Style Transfer Task

## Teacher models:

- GPT models

## Student models:

- LLaMA-7B
- Alpaca-7B

## Other instruction-tuned models

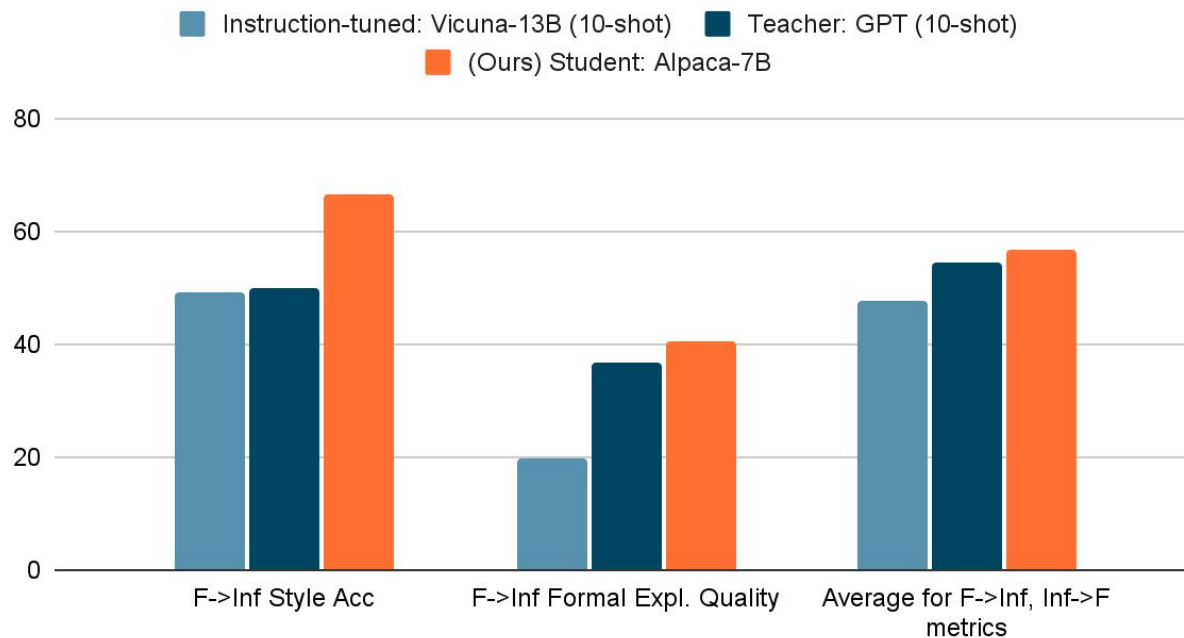
- Vicuna-13B
- More baselines in the paper

## Settings:

- One-shot
- Few-shot with ICLEF-improved data
  - 5-shot
  - 10-shot

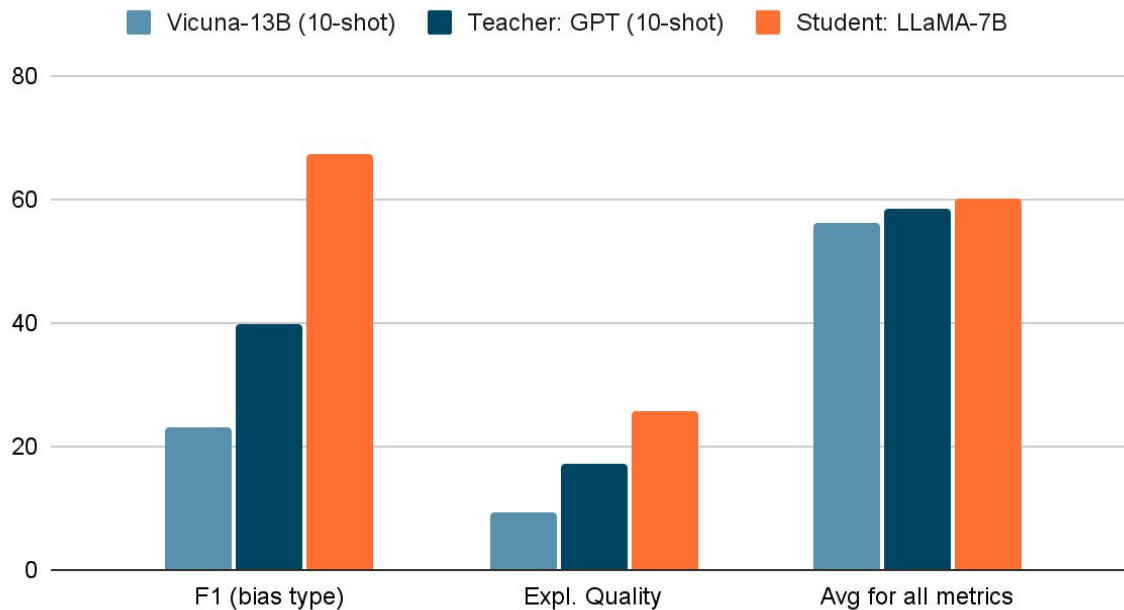
# Results: Formality Explainable Style Transfer

## Automatic Evaluation



# Results: Bias Style Transfer

## Automatic Evaluation (Explainable Bias Style Transfer)



# Human Evaluation

## Formal (Input)

If it is feasible, allow love to prevail.

## Teacher (1-shot)

Formal Attributes:

conditional clause ("if it is feasible"), elevated vocabulary ("allow love to prevail")

Informal Paraphrase:

**If it's possible, just go with love.**

Informal Attributes:

simple language ("if it's possible"), imperative form ("just go"), contraction ("it's")

## Student (Ours):

Formal Attributes:

precise language ("feasible"), use of conditional clause ("if it is feasible"), lexical sophistication ("prevail")

Informal Paraphrase:

**let love win if it's possible.**

Informal Attributes:

contraction ("it's"), casual tone ("let love win")

## Expert Annotator Preferences

Vicuna-13B (1-shot)

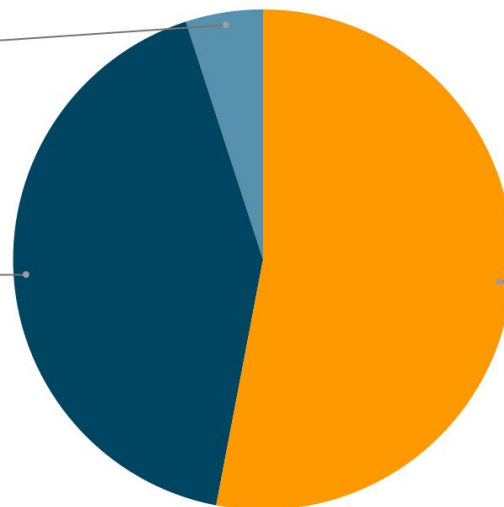
5.0%

Teacher (1-shot)

42.0%

Student (Ours)

53.0%

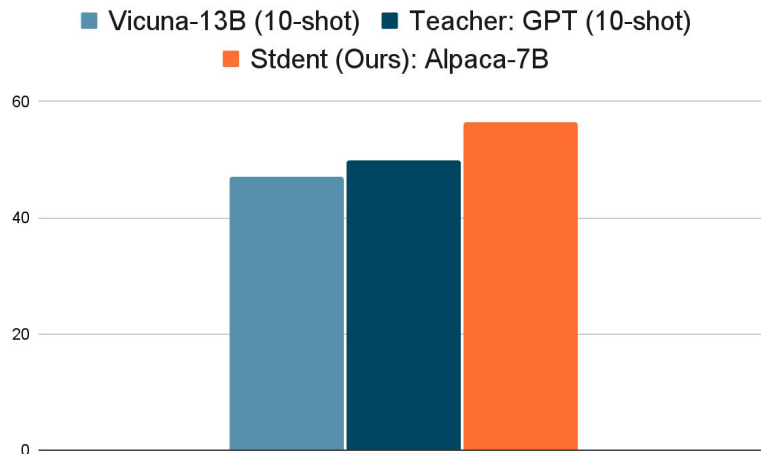




# Extrinsic Evaluation: Authorship Verification

- **Task:** decide if two texts belong to the same author
- **Approach:**
  - Apply explainable style transfer model to extract informality attributes
  - Use % of overlapping attributes as a score

Attribute	Evidence
Colloquialism	<i>“assumed they all started off low!?”</i> , <i>“typing it out”</i>
Textese	<i>“xx”</i>
Informal Tone	<i>“hoping to borrow a couple of charging leads”</i>



# Bonus: Interpretable Adversarial Attacks on AI-generated text detection

- How well can informal paraphrase be detected as AI-generated?
  - GPT-F: Formal (AI-generated)
  - GPT-Inf: ChatGPT informal paraphrase
  - Our model

GPTZero Classification score as AI-generated

