



# Are Large Language Models Actually Good at Text Style Transfer?

Sourabrata Mukherjee<sup>1</sup>, Atul Kr. Ojha<sup>2</sup>, Ondrej Dusek<sup>1</sup>

<sup>1</sup>Charles University, Faculty of Mathematics and Physics, Prague, Czechia

<sup>2</sup>Insight SFI Research Centre for Data Analytics, DSI, University of Galway, Ireland



# Introduction

---

# Text Style Transfer (TST)



- Change style of given input text
- Preserve style-independent content
- **Style:**
  - demographic attributes (personality, gender)
  - sentiment
  - toxicity
  - politeness
  - ...



# Sentiment Transfer

- A sub-task of TST
- Positive to negative text & vice versa
- Keeps sentiment-independent content
- Example:

*The food is yummy.* ➔ *The food is tasteless.*

- **Uses:**
  - Marketing
  - Content Moderation
  - Communication improvement



# Detoxification

- A sub-task of TST
- Converts toxic text to clean text
- Without changing the intent of the text as much as possible
- Example:

*You're an idiot.* ➡ *You made a mistake.*

- **Uses:**
  - Hate/Toxic Comment Removal
  - Offensive Language Mitigation
  - Adjusting Political Extremism Language



# Our Work

# Overview

- Systematic LLM evaluation on TST
- Sentiment transfer & text detoxification
- English, Hindi & Bengali
- Zero-shot, few-shot prompting & parameter-efficient fine-tuning
- Automatic metrics, human evaluation & GPT-4-based evaluation
- Compare previous SOTA trained on dedicated datasets
- GPT-3.5 and some open LLMs show promising results but do not surpass previous SOTA
- Fine-tuning significantly improves open LLMs' performance (➔close to GPT-3.5 and SOTA)
- Dedicated datasets & tailored models still useful for TST



# LLMs

- Open LLMs & size variants

Model	Size Variants
BLOOM (BigScience Workshop, 2023)	560M, 1B, 3B, and 7B
BLOOMz (Muennighoff et al., 2023)	560M, 1B, 3B, and 7B
ChatGLM (Du et al., 2022)	6B
ChatGLM2 (Du et al., 2022)	6B
Falcon (Penedo et al., 2023; Almazrouei et al., 2023)	7B
Llama (Touvron et al., 2023a)	7B, 13B, and 30B
Llama-2 (Touvron et al., 2023b)	7B, and 13B
Llama-2-Chat (Touvron et al., 2023b)	7B, and 13B
Llama-3 (AI@Meta, 2024)	8B
Llama-3-Instruct (AI@Meta, 2024)	8B
Mistral-Instruct (Jiang et al., 2023)	7B
OPT (Zhang et al., 2022)	1.3B, 2.7B, 6.7B, 13B, and 30B
Zephyr (Tunstall et al., 2023)	7B

- + GPT-3.5 (*gpt-3.5-turbo*) via OpenAI API





# SOTA models

- Sentiment transfer: Mukherjee et al. (2024)\*
  - finetuned mBART
  - *Parallel* = single-language parallel data
  - *Joint* = data in multiple languages
- Text detoxification: Mukherjee et al. (2023)<sup>†</sup>
  - finetuned mBART
  - *Seq2seq + CLS\_OP* = multi-task w/text classif.
  - *KT* = knowledge transfer from sentiment

\*<https://aclanthology.org/2024.inlg-main.41>

<sup>†</sup><https://aclanthology.org/2023.icon-1.13>



# Datasets

- Sentiment
  - English, Hindi & Bengali
  - average score of pos  $\Rightarrow$  neg & neg  $\Rightarrow$  pos
- Detoxification
  - English & Hindi
  - toxic  $\Rightarrow$  clean
- All experiments – 1,000 examples:
  - 400 fine-tuning (if applicable)
  - 100 development
  - 500 testing



# Evaluation

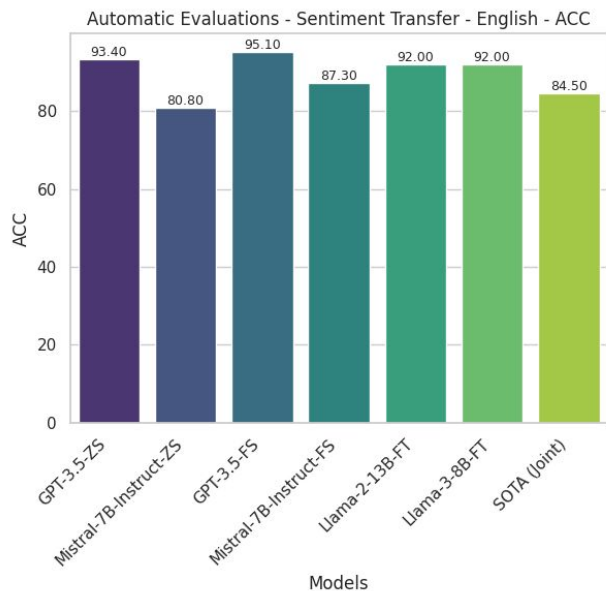
- Automatic metrics
  - sentiment / toxicity classifier
  - content preservation:  
BLEU & SBERT cosine similarity
- 50 sentiment outputs: GPT-4 & humans
  - 5-point Likert scales
  - style, content preservation, fluency

---

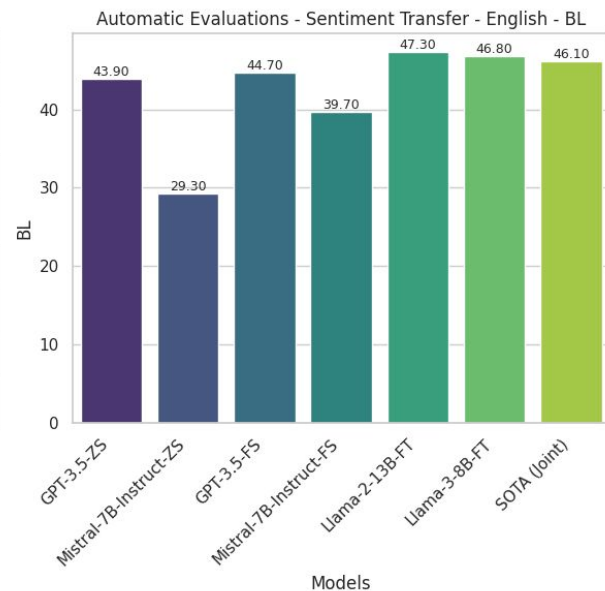
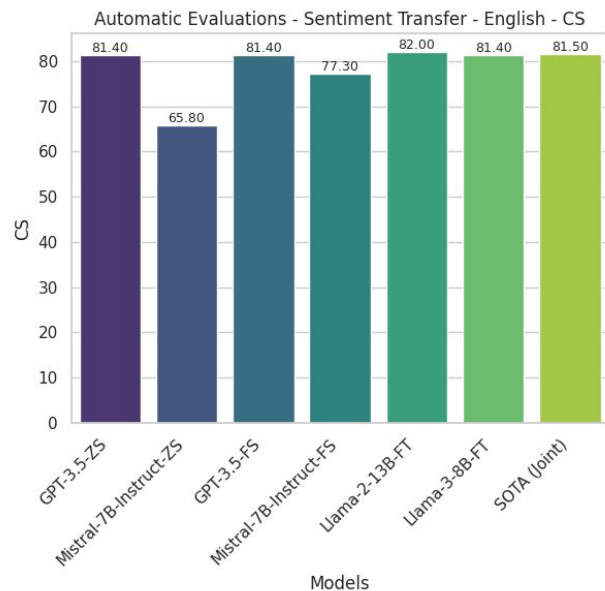
# Automatic Evaluation Results

# Sentiment Transfer: English

## Style Transfer Accuracy

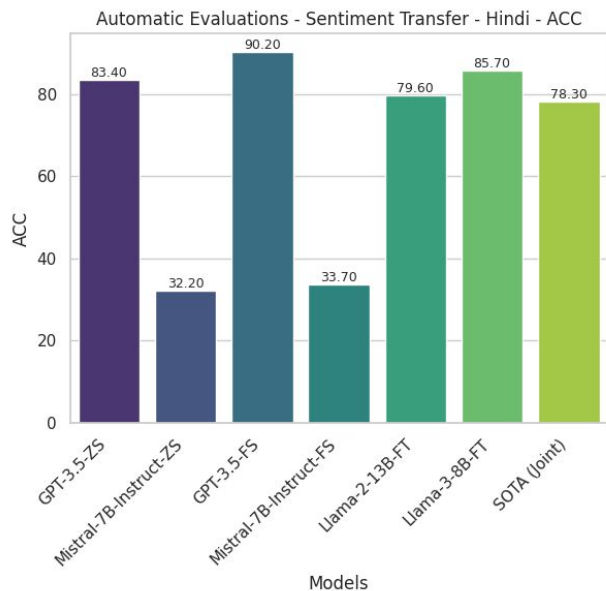


## Content Preservation

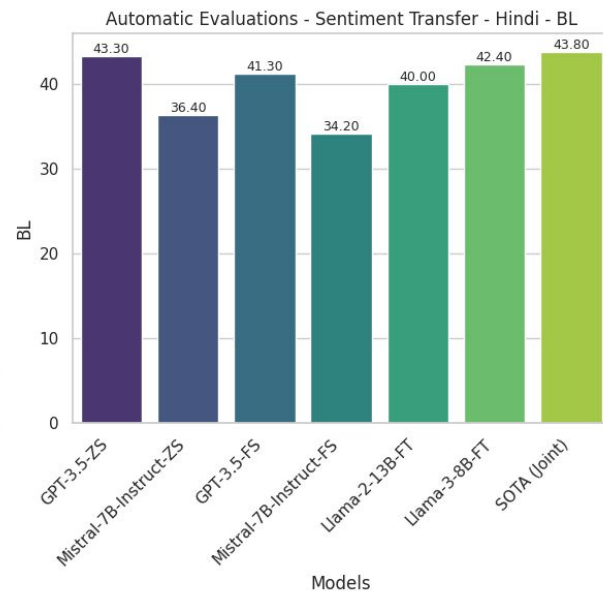
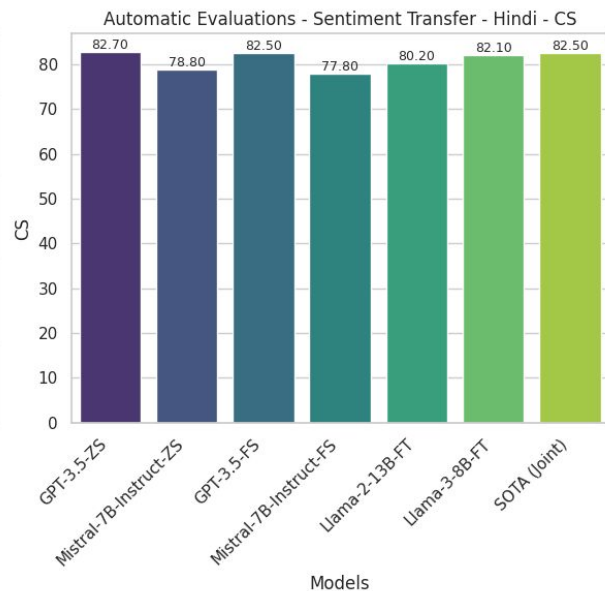


# Sentiment Transfer: Hindi

## Style Transfer Accuracy

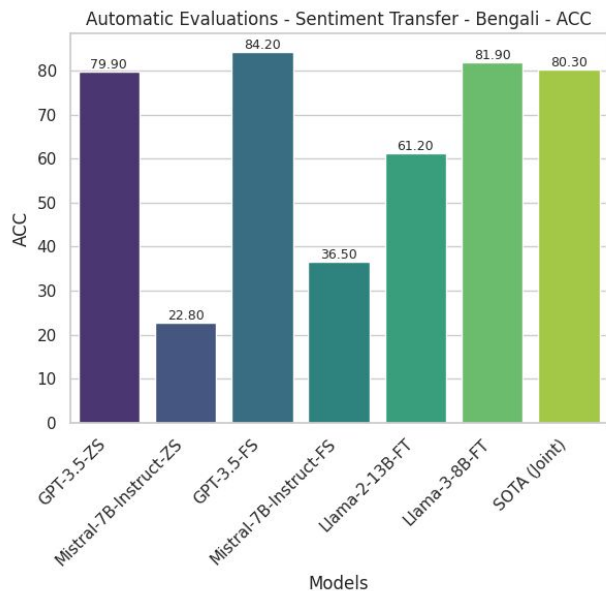


## Content Preservation

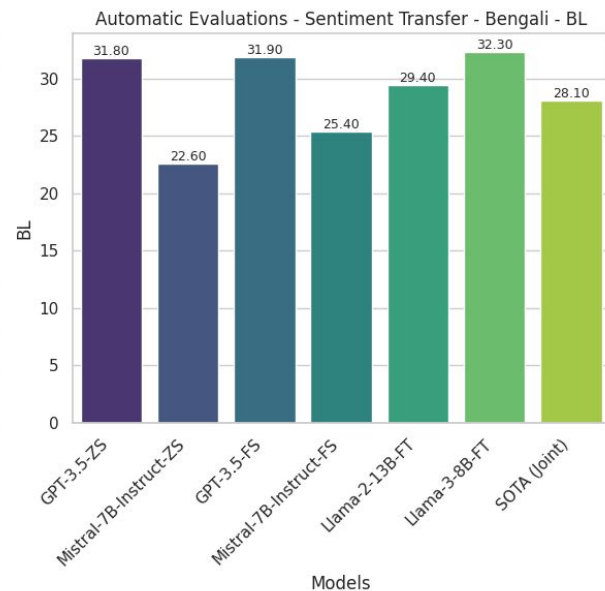
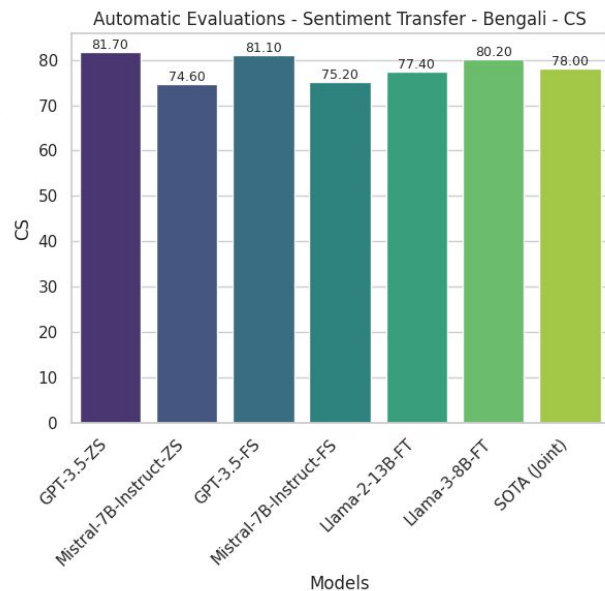


# Sentiment Transfer: Bengali

## Style Transfer Accuracy

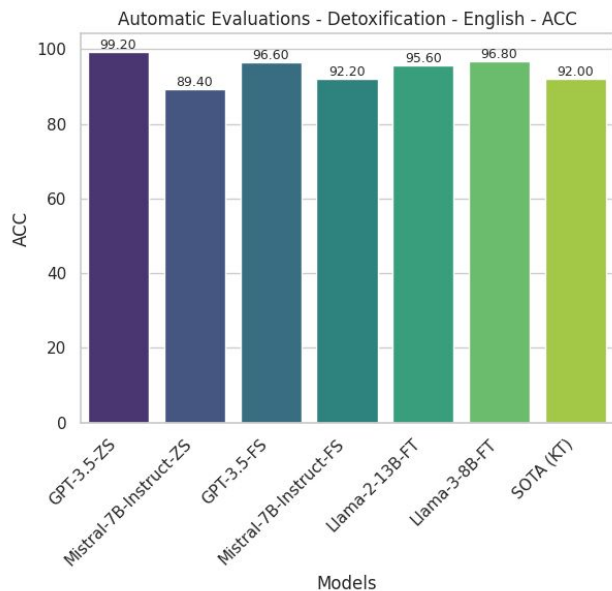


## Content Preservation

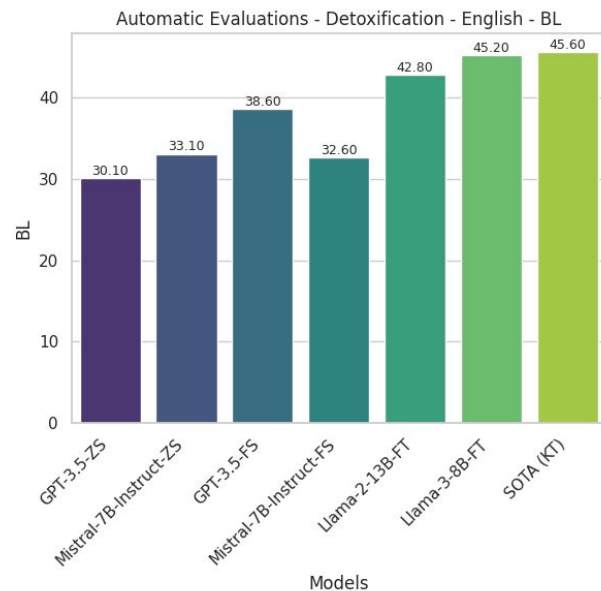
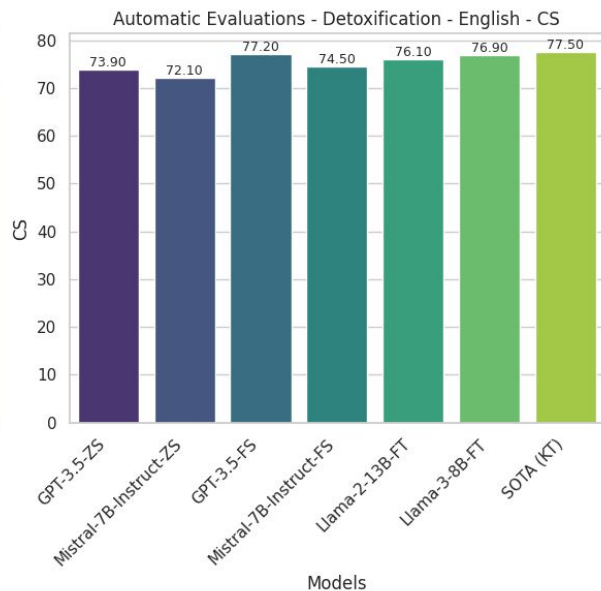


# Detoxification: English

## Style Transfer Accuracy



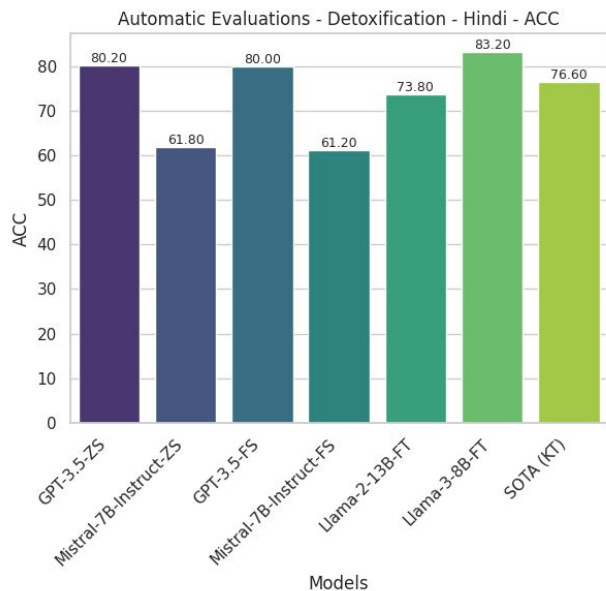
## Content Preservation



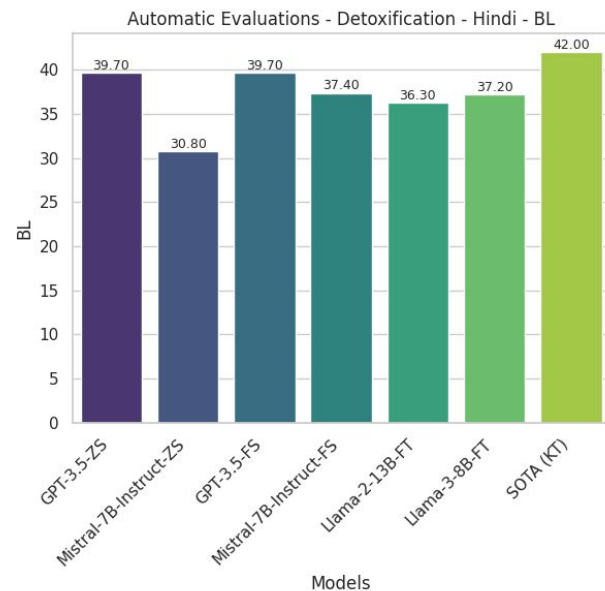
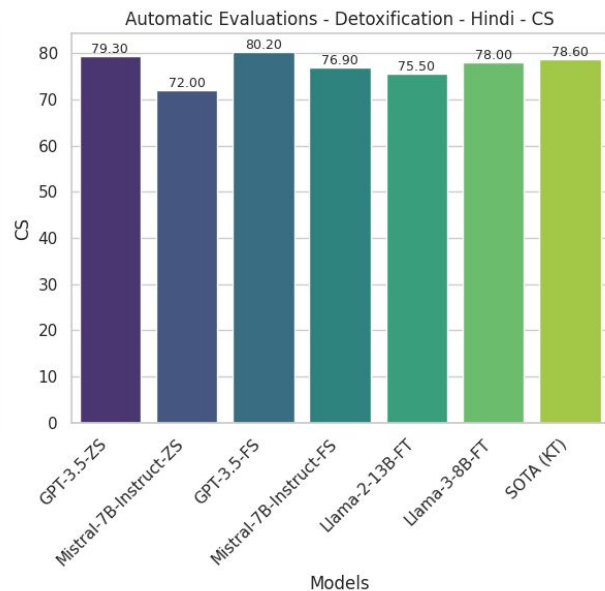


# Detoxification: Hindi

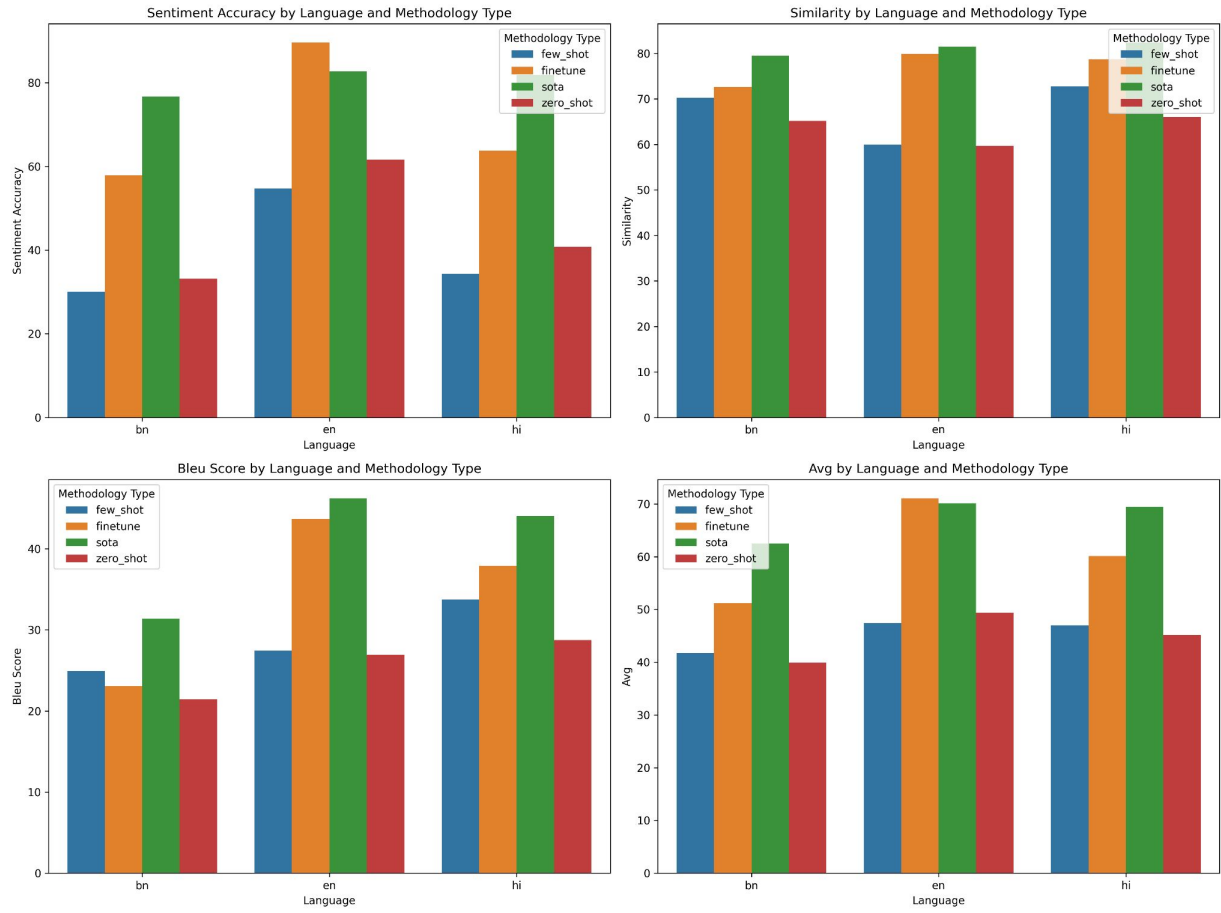
## Style Transfer Accuracy



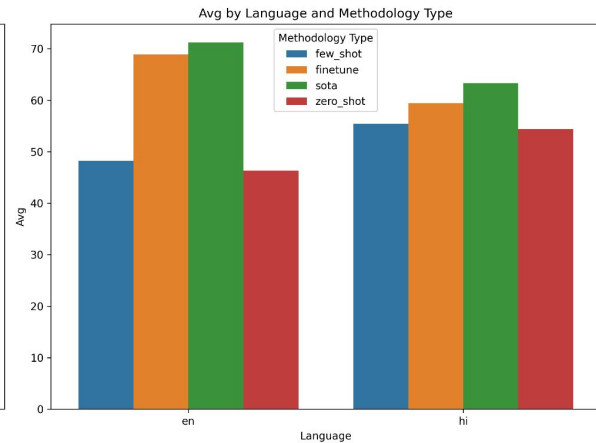
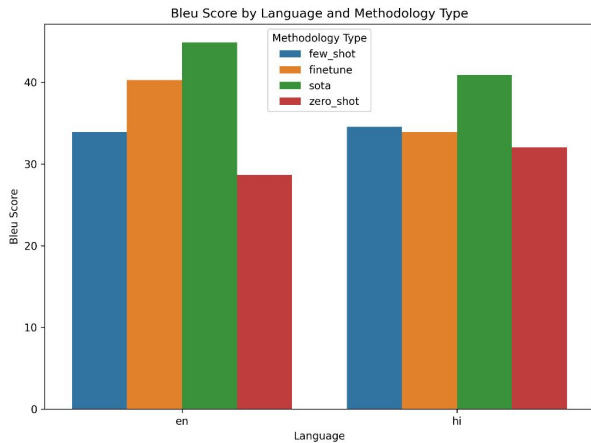
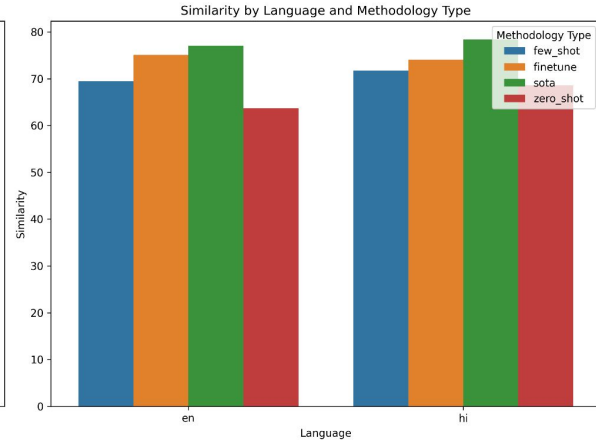
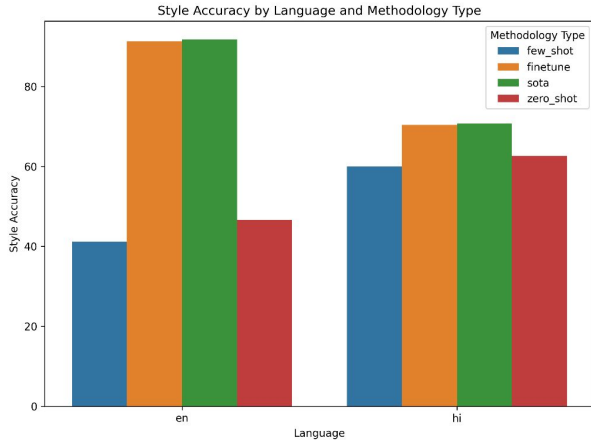
## Content Preservation



# Methodology Averages: Sentiment Transfer

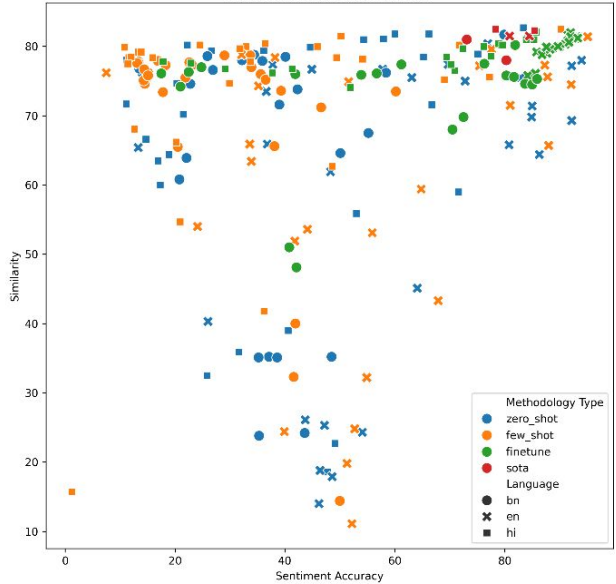


# Methodology Averages: Detoxification

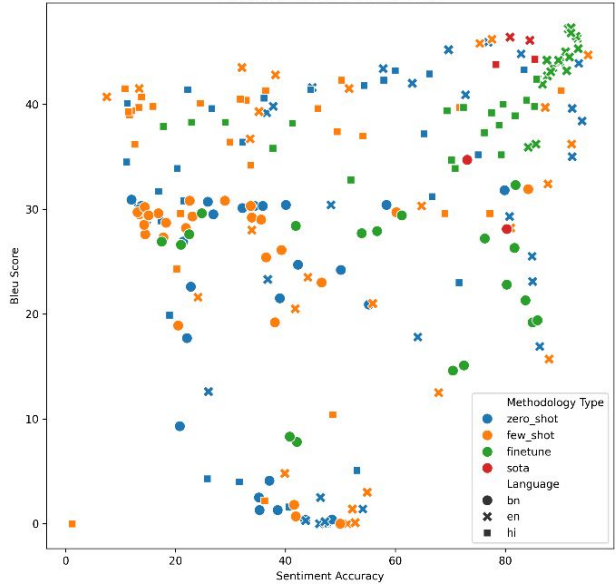


# Trade-offs: Sentiment Transfer

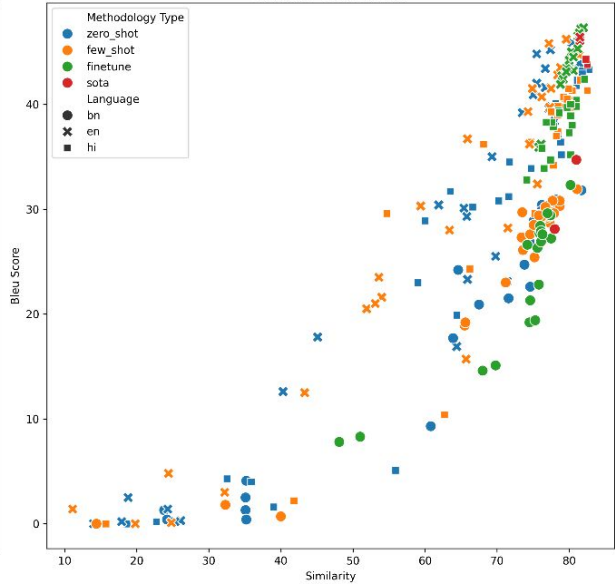
Sentiment Accuracy vs Similarity



Sentiment Accuracy vs Bleu Score

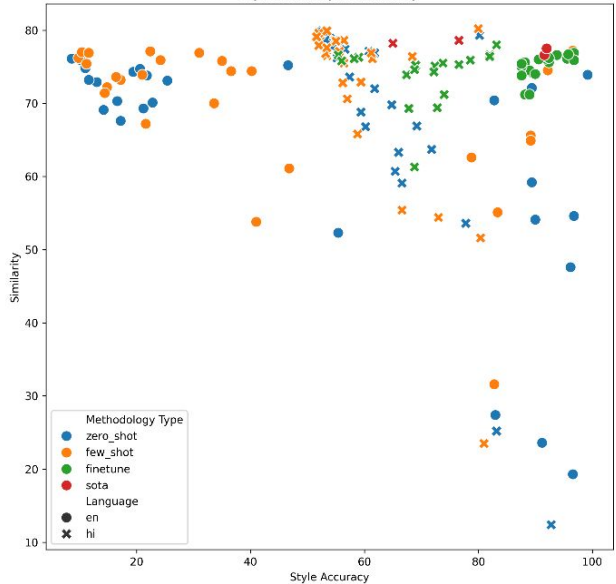


Similarity vs Bleu Score

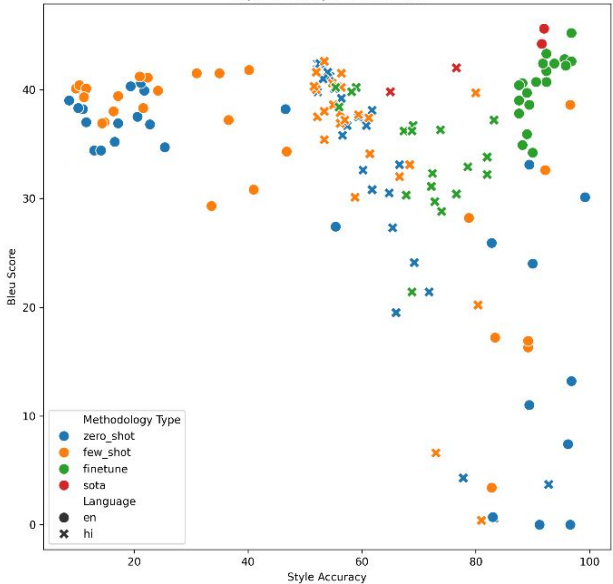


# Trade-offs: Detoxification

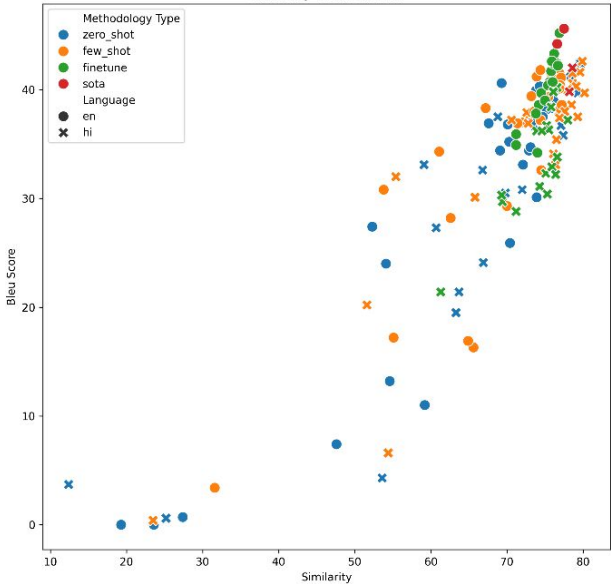
Style Accuracy vs Similarity



Style Accuracy vs Bleu Score



Similarity vs Bleu Score



# Automatic Evaluation

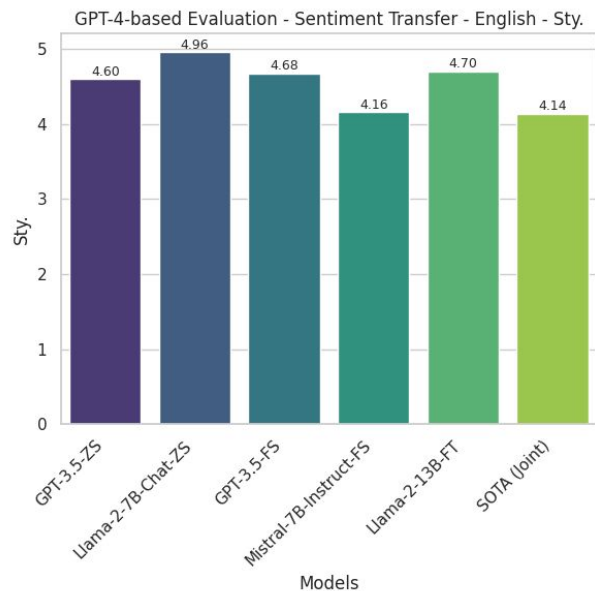
Models	Sentiment Transfer												Detoxification							
	English				Hindi				Bengali				English				Hindi			
	ACC	CS	BL	AVG	ACC	CS	BL	AVG	ACC	CS	BL	AVG	ACC	CS	BL	AVG	ACC	CS	BL	AVG
BLOOM-7B-ZS	37.8	77.4	39.8	51.6	26.6	79.4	39.6	48.6	34.4	78.8	30.3	47.8	8.6	76.1	39.0	41.2	52.2	79.1	39.8	57.0
BLOOMz-7B-ZS	26.0	40.3	12.6	26.3	31.6	35.9	4.0	23.9	35.2	35.1	2.5	24.2	14.2	69.1	34.4	39.2	64.8	69.8	30.5	55.0
ChatGLM2-6B-ZS	86.3	64.4	16.9	55.8	53.0	55.9	5.1	38.0	48.5	35.2	0.4	28.0	96.2	47.6	7.4	50.4	77.8	53.6	4.3	45.2
Falcon-7B-ZS	72.8	75.0	40.9	62.9	21.5	70.2	30.8	40.8	22.1	63.9	17.7	34.6	46.6	75.2	38.2	53.3	65.4	60.7	27.3	51.1
GPT-3.5-ZS	93.4	81.4	43.9	72.9	83.4	82.7	43.3	69.8	79.9	81.7	31.8	64.5	99.2	73.9	30.1	67.7	80.2	79.3	39.7	66.4
Llama-7B-ZS	36.8	65.9	23.3	42.0	22.2	80.2	41.4	47.9	12.0	78.2	30.9	40.4	11.6	73.2	37.0	40.6	52.6	79.7	42.4	58.2
Llama-2-7B-ZS	63.1	75.5	42.0	60.2	44.6	79.9	41.4	55.3	26.9	76.6	29.5	44.3	20.6	74.7	37.5	44.3	53.2	78.7	41.0	57.7
Llama-2-Chat-7B-ZS	94.0	78.0	38.4	70.1	65.2	78.5	37.2	60.3	39.0	71.6	21.5	44.0	82.8	70.4	25.9	59.7	61.8	76.9	38.1	58.9
Llama-3-8B-ZS	76.9	80.4	45.9	67.7	66.2	81.8	42.9	63.6	58.4	76.2	30.4	55.0	25.4	73.1	34.7	44.4	56.6	77.4	35.8	56.6
Llama-3-8B-Instruct-ZS	92.2	69.3	35.0	65.5	71.6	59.0	23.0	51.2	50.1	64.6	24.2	46.3	-	-	-	-	-	-	-	-
Mistral-7B-Instruct-ZS	80.8	65.8	29.3	58.6	32.2	78.8	36.4	49.1	22.8	74.6	22.6	40.0	89.4	72.1	33.1	64.9	61.8	72.0	30.8	54.9
OPT-6.7B-ZS	54.1	24.3	1.4	26.6	17.3	60.0	28.9	35.4	13.5	76.8	30.0	40.1	83.0	27.4	0.7	37.0	66.6	59.1	33.1	52.9
Zephyr-7B-ZS	85.0	71.4	23.1	59.8	66.7	71.6	31.2	56.5	55.2	67.5	20.9	47.9	96.8	54.6	13.2	54.9	71.8	63.7	21.4	52.3
BLOOM-7B-FS	32.1	78.8	43.5	51.5	24.5	80.2	40.1	48.3	16.9	77.9	29.6	41.5	22.4	77.1	41.1	46.9	52.0	79.6	41.6	57.7
BLOOMz-7B-FS	35.2	74.3	39.3	49.6	36.4	80.4	41.3	52.7	29.0	78.7	30.8	46.2	14.4	71.4	36.9	40.9	59.4	72.9	37.7	56.7
ChatGLM2-6B-FS	87.8	75.6	32.4	65.3	48.6	62.7	10.4	40.6	41.9	40.0	0.7	27.6	89.2	64.9	16.9	57.0	73.0	54.4	6.6	44.7
Falcon-7B-FS	77.6	79.6	46.2	67.8	15.9	78.4	39.8	44.7	17.8	73.4	27.3	39.5	24.2	75.9	39.9	46.7	56.4	75.5	40.2	57.3
GPT-3.5-FS	95.1	81.4	44.7	73.7	90.2	82.5	41.3	71.3	84.2	81.1	31.9	65.7	96.6	77.2	38.6	70.8	80.0	80.2	39.7	66.6
Llama-7B-FS	64.8	59.4	30.3	51.5	31.8	79.7	40.5	50.7	23.1	77.3	29.3	43.2	11.6	76.9	40.1	42.9	53.4	79.9	42.6	58.6
Llama-2-7B-FS	54.9	32.2	3.0	30.0	54.1	78.2	37.0	56.4	39.3	73.6	26.1	46.3	46.8	61.1	34.3	47.4	53.4	77.6	38.0	56.3
Llama-2-Chat-7B-FS	92.1	74.5	36.2	67.6	69.0	75.2	29.6	57.9	38.1	65.6	19.2	40.9	78.8	62.6	28.2	56.5	61.4	76.1	34.1	57.2
Llama-3-8B-FS	67.9	43.3	12.5	41.3	71.7	80.2	39.7	63.9	60.2	73.5	29.7	54.4	40.2	74.4	41.8	52.2	80.4	51.6	20.2	50.7
Llama-3-8B-Instruct-FS	52.2	11.1	1.4	21.6	1.2	15.7	0	5.6	50.0	14.4	0	21.5	-	-	-	-	-	-	-	-
Mistral-7B-Instruct-FS	87.3	77.3	39.7	68.1	33.7	77.8	34.2	48.6	36.5	75.2	25.4	45.7	92.2	74.5	32.6	66.5	61.2	76.9	37.4	58.5
OPT-6.7B-FS	33.9	63.4	28.0	41.8	11.4	77.5	39.3	42.7	15.1	75.8	29.4	40.1	11.2	75.4	39.3	42.0	57.0	70.6	37.2	54.9
BLOOM-7B-FT	91.2	80.6	43.2	71.7	83.9	81.0	40.4	68.4	81.7	75.6	26.3	61.2	92.4	75.8	41.7	70.0	82.0	76.6	33.8	64.1
BLOOMz-7B-FT	91.0	80.3	45.0	72.1	85.3	81.0	39.8	68.7	85.9	75.3	19.4	60.2	92.4	75.6	40.7	69.6	82.0	76.4	32.2	63.5
ChatGLM2-6B-FT	86.8	78.8	41.9	69.2	51.9	74.1	32.8	52.9	42.1	48.1	7.8	32.7	90.0	74.0	34.2	66.1	67.8	69.3	30.3	55.8
Falcon-7B-FT	88.3	79.6	43.1	70.3	37.7	76.2	35.8	49.9	40.8	51.0	8.3	33.4	87.6	73.8	37.8	66.4	68.8	61.3	21.4	50.5
Llama-7B-FT	91.5	81.6	47.2	73.4	69.4	78.5	39.4	62.4	41.9	76.0	28.4	48.8	91.8	76.1	42.4	70.1	67.4	73.9	36.2	59.2
Llama-2-7B-FT	92.9	81.2	46.5	73.5	77.5	78.6	39.2	65.1	56.7	76.1	27.9	53.6	92.4	76.2	43.3	70.6	68.8	74.6	36.2	59.9
Llama-2-13B-FT	92.0	82.0	47.3	73.8	79.6	80.2	40.0	66.6	61.2	77.4	29.4	56.0	95.6	76.1	42.8	71.5	73.8	75.5	36.3	61.9
Llama-3-8B-FT	92.0	81.4	46.8	73.4	85.7	82.1	42.4	70.1	81.9	80.2	32.3	64.8	96.8	76.9	45.2	73.0	83.2	78.0	37.2	66.1
OPT-6.7B-FT	91.7	80.6	44.5	72.3	29.1	76.8	38.3	48.1	22.5	76.3	27.6	42.1	95.8	76.7	42.2	71.6	58.2	76.1	39.8	58.0
SOTA (Joint)	84.5	81.5	46.1	70.7	78.3	82.5	43.8	68.2	80.3	78.0	28.1	62.1	-	-	-	-	-	-	-	-
SOTA (Parallel)	80.9	81.5	46.4	69.6	85.4	82.3	44.3	70.7	73.1	81.0	34.7	62.9	-	-	-	-	-	-	-	-
SOTA (CLS-OP)	-	-	-	-	-	-	-	-	-	-	-	-	91.6	76.6	44.2	70.8	65.0	78.2	39.8	61.0
SOTA (KT)	-	-	-	-	-	-	-	-	-	-	-	-	92.0	77.5	45.6	71.7	76.6	78.6	42.0	65.7

---

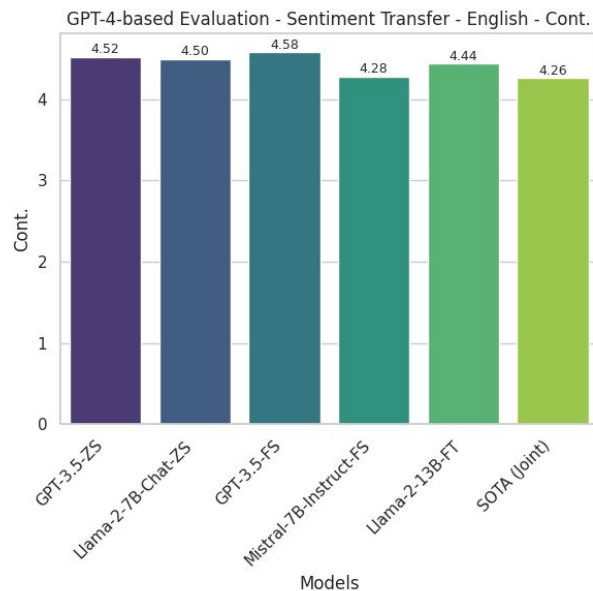
# GPT-4-based Evaluation Results

# Sentiment Transfer: English

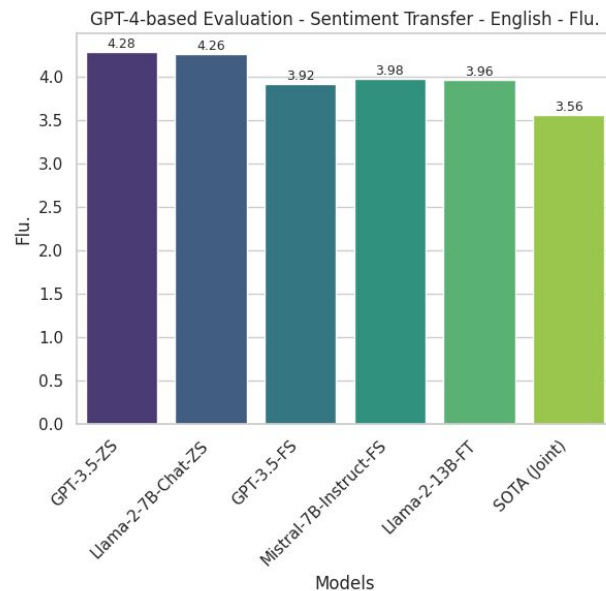
## Style Transfer Accuracy



## Content Preservation



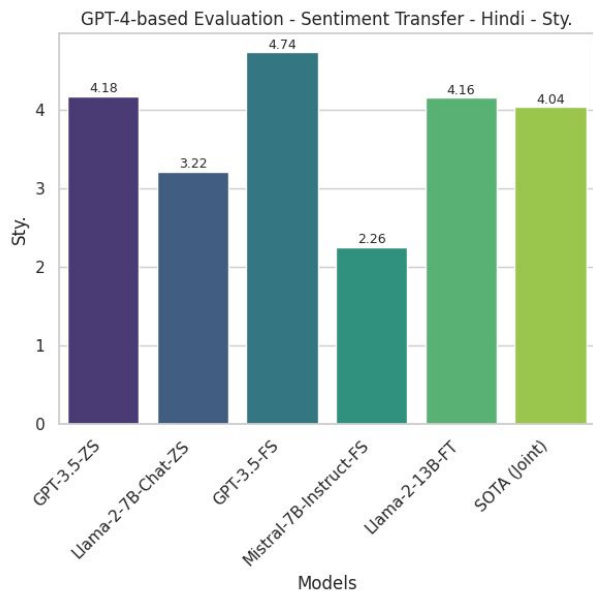
## Fluency



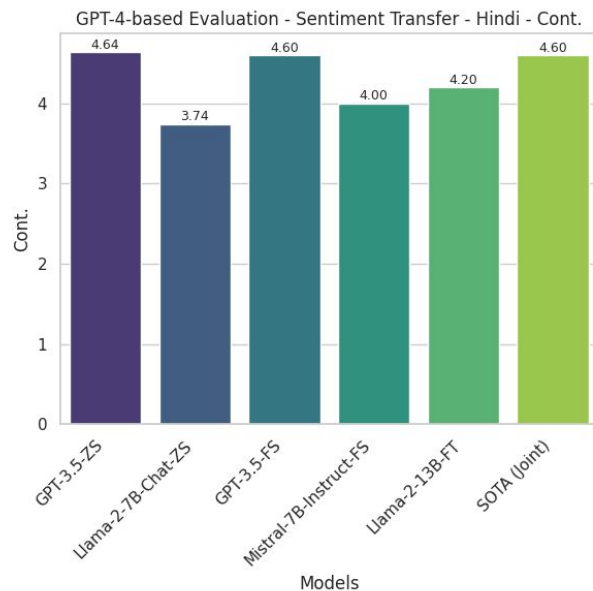


# Sentiment Transfer: Hindi

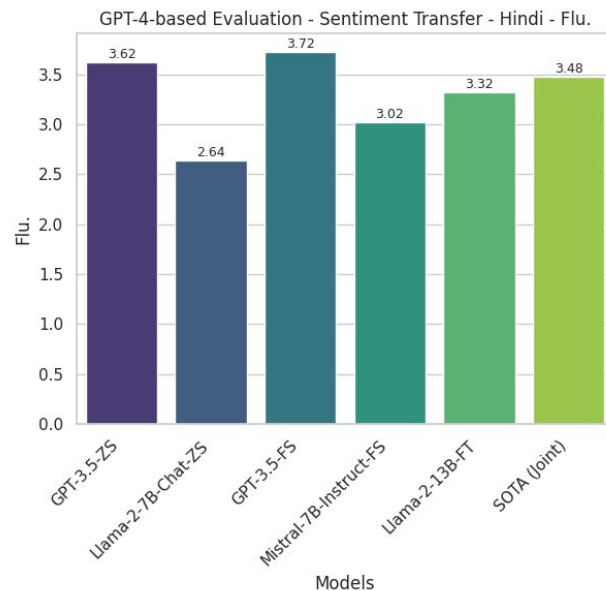
## Style Transfer Accuracy



## Content Preservation

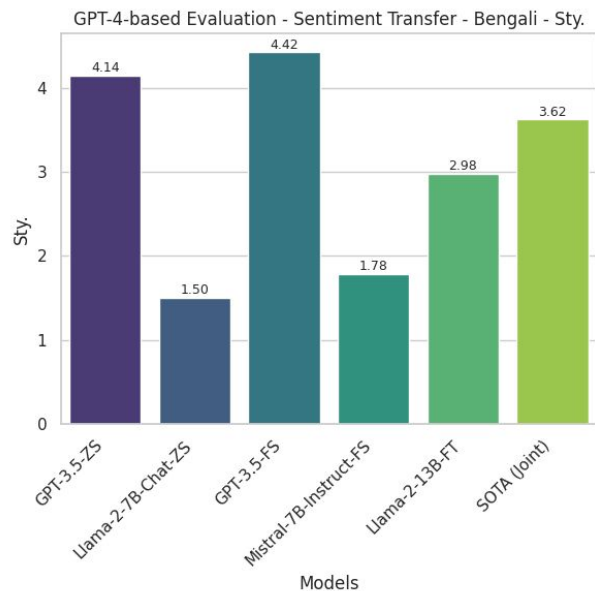


## Fluency

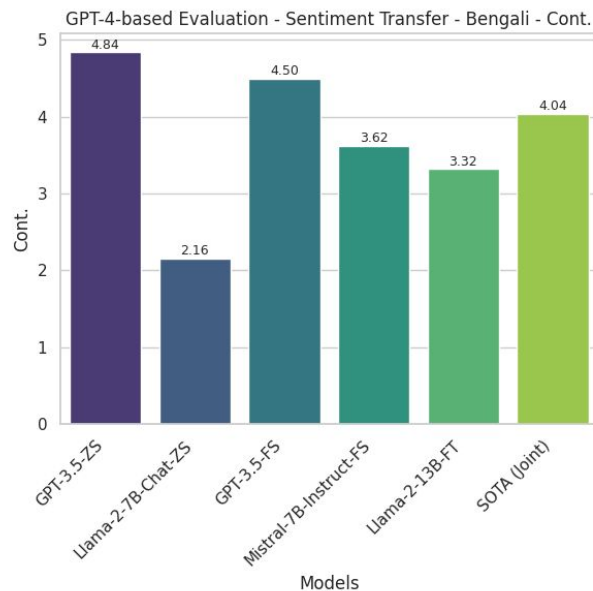


# Sentiment Transfer: Bengali

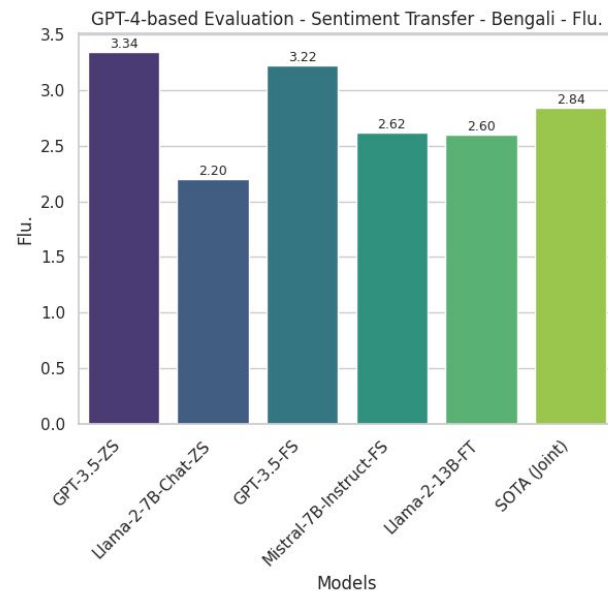
## Style Transfer Accuracy



## Content Preservation

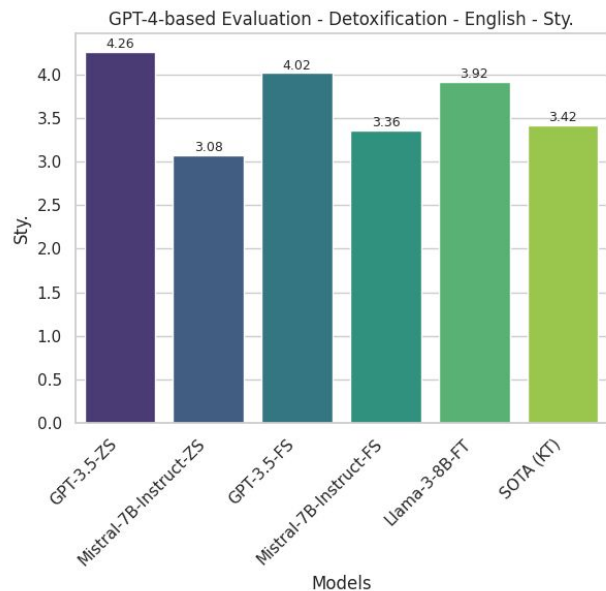


## Fluency

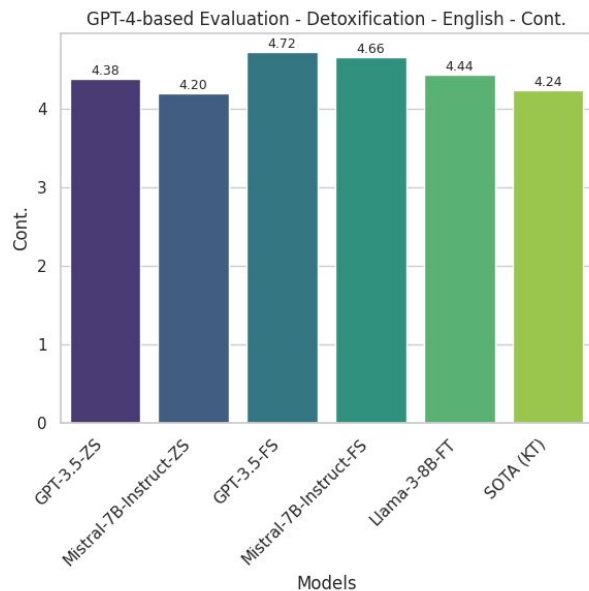


# Detoxification: English

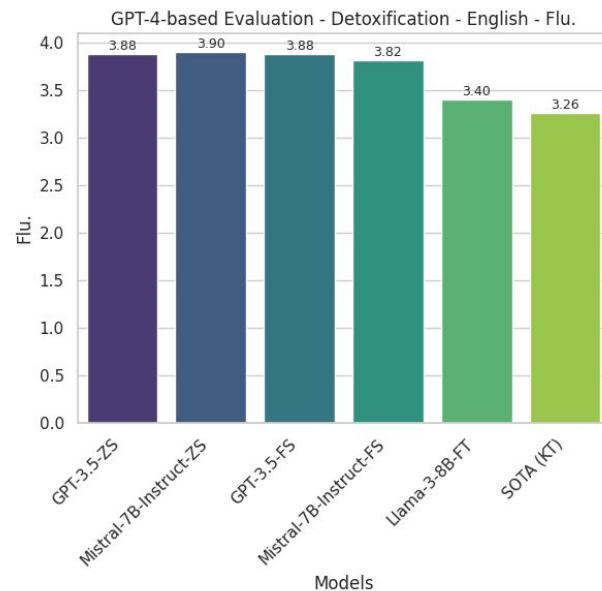
## Style Transfer Accuracy



## Content Preservation

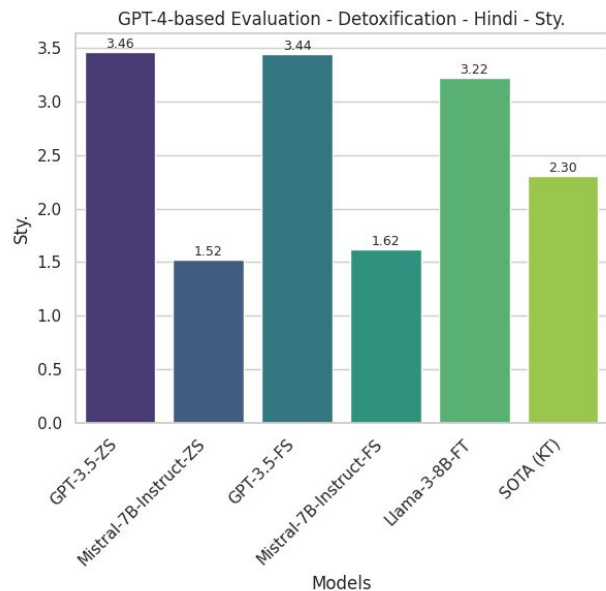


## Fluency

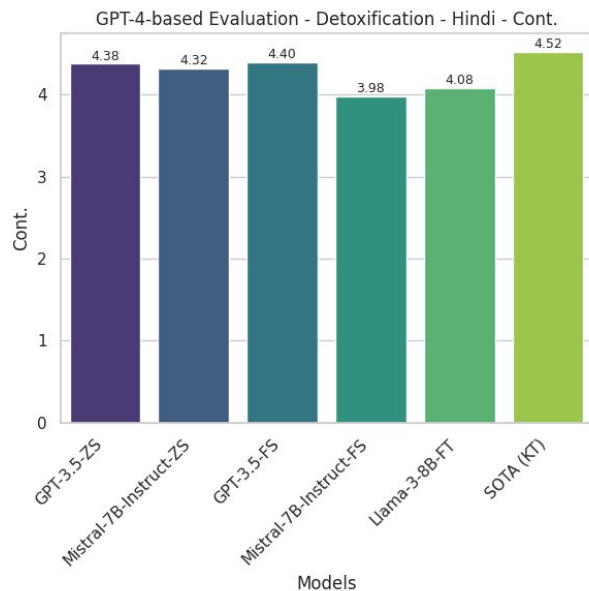


# Detoxification: Hindi

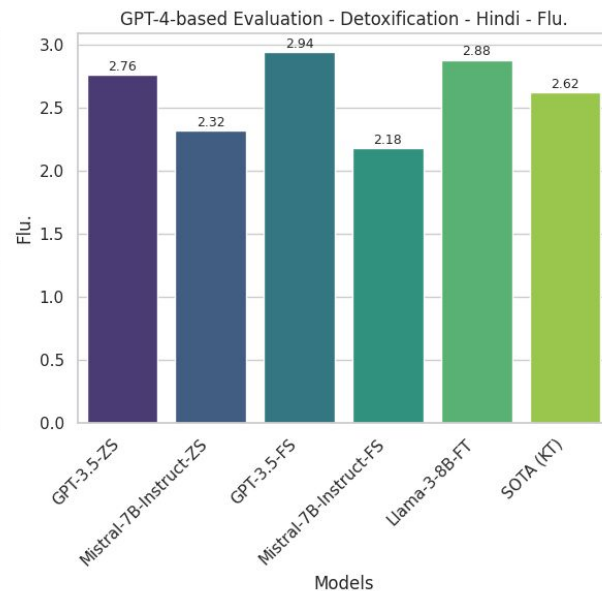
## Style Transfer Accuracy



## Content Preservation



## Fluency

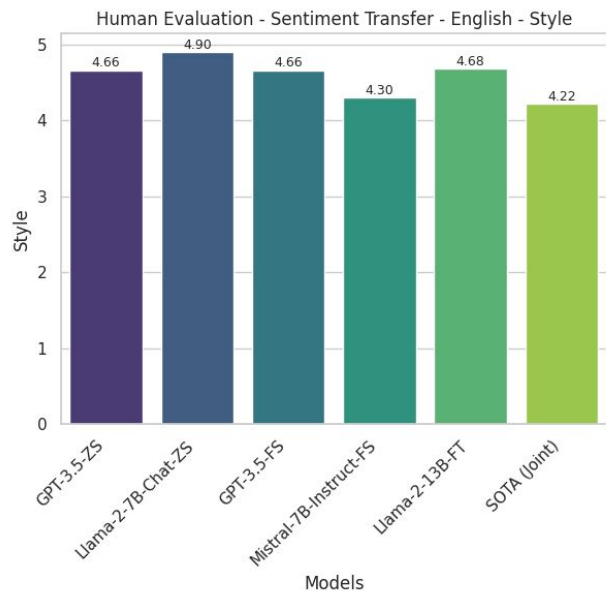




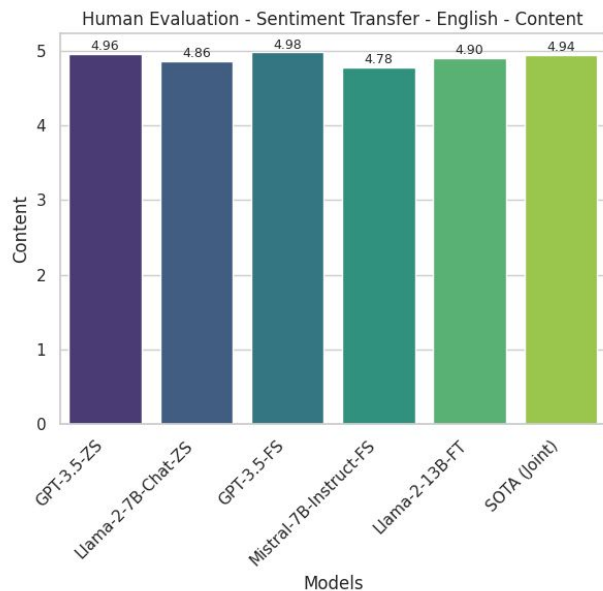
# Human Evaluation Results

# Sentiment Transfer: English

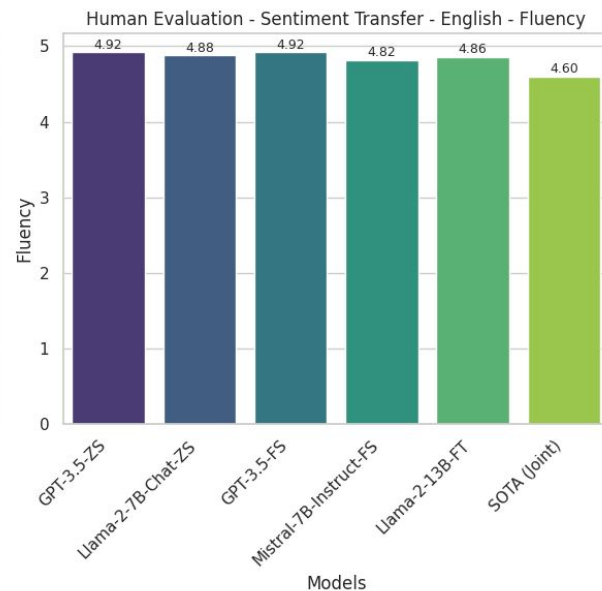
## Style Transfer Accuracy



## Content Preservation

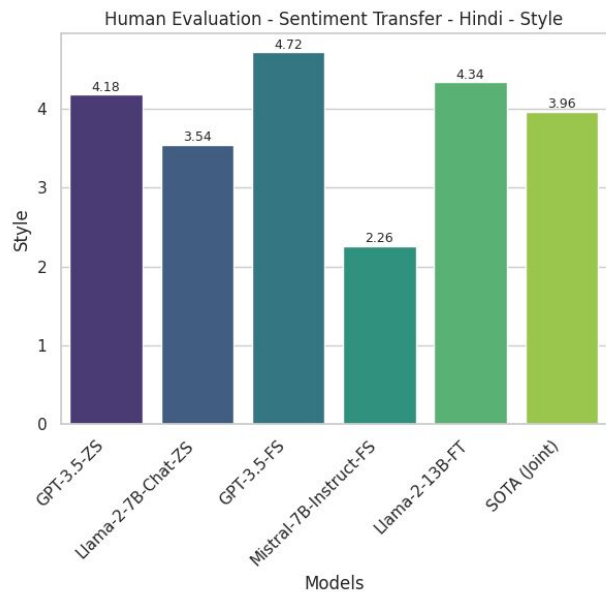


## Fluency

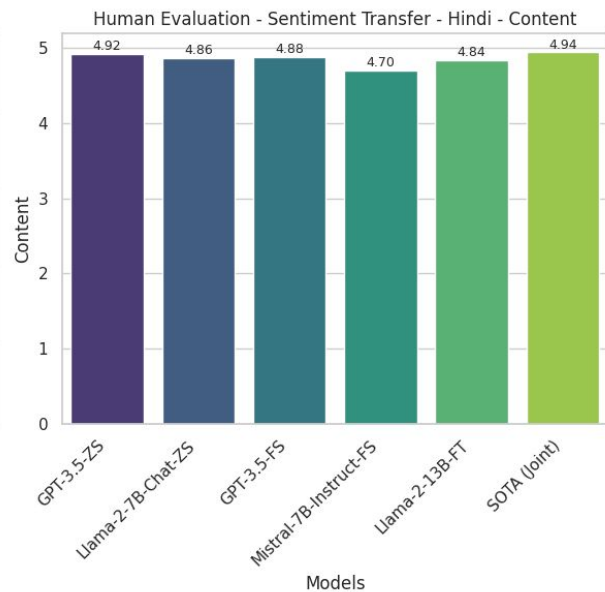


# Sentiment Transfer: Hindi

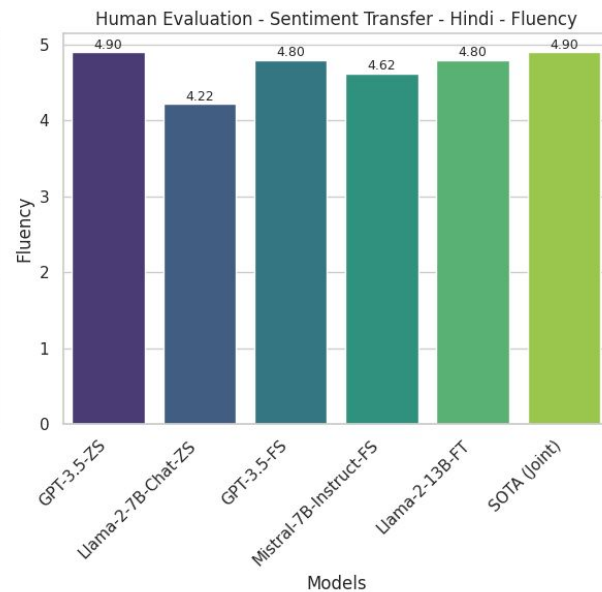
## Style Transfer Accuracy



## Content Preservation



## Fluency



# GPT-4-based Evaluation

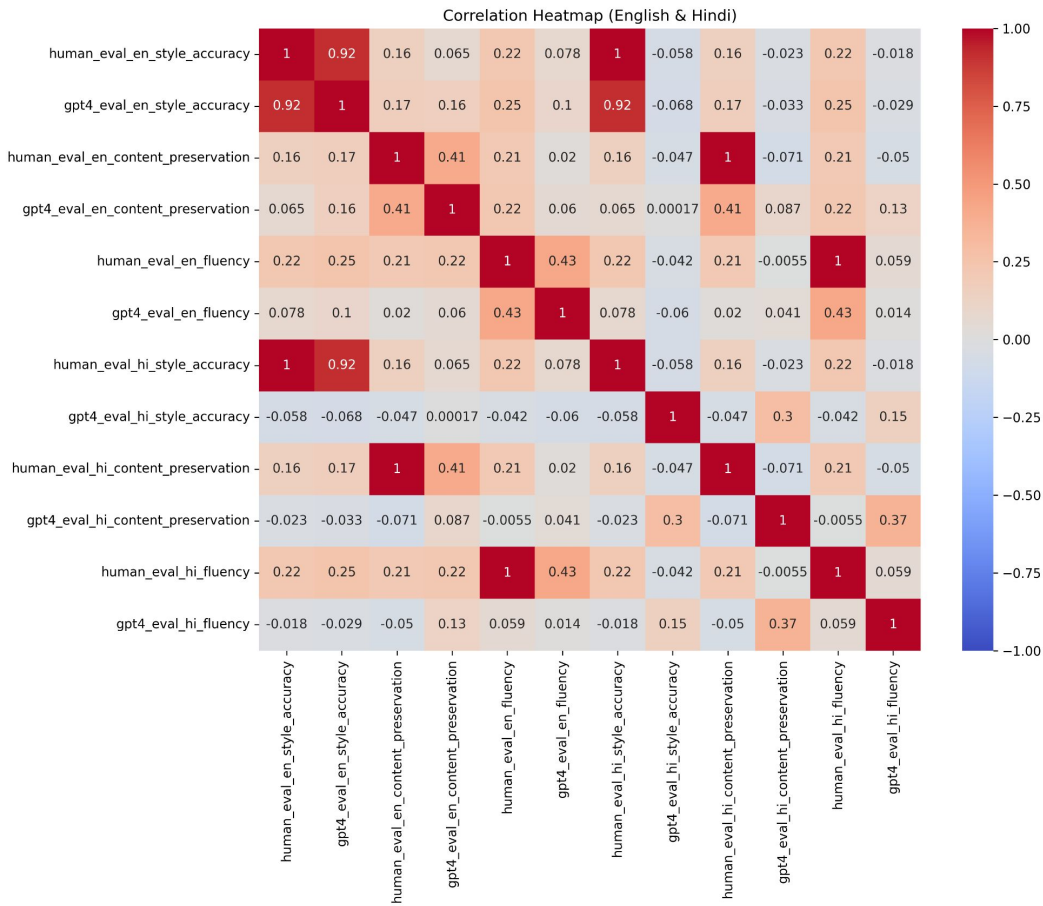
Models	Sentiment transfer									Detoxification					
	English			Hindi			Bengali			English			Hindi		
	Sty.	Cont.	Flu.	Sty.	Cont.	Flu.	Sty.	Cont.	Flu.	Sty.	Cont.	Flu.	Sty.	Cont.	Flu.
GPT-3.5-ZS	4.60	4.52	4.28	4.18	4.64	3.62	4.14	4.84	3.34	4.26	4.38	3.88	3.46	4.38	2.76
Llama-2-7B-Chat-ZS	4.96	4.50	4.26	3.22	3.74	2.64	1.50	2.16	2.20						
Mistral-7B-Instruct-ZS										3.08	4.20	3.90	1.52	4.32	2.32
GPT-3.5-FS	4.68	4.58	3.92	4.74	4.60	3.72	4.42	4.50	3.22	4.02	4.72	3.88	3.44	4.40	2.94
Mistral-7B-Instruct-FS	4.16	4.28	3.98	2.26	4.00	3.02	1.78	3.62	2.62	3.36	4.66	3.82	1.62	3.98	2.18
Llama-2-13B-FT	4.70	4.44	3.96	4.16	4.20	3.32	2.98	3.32	2.60						
Llama-3-8B-FT										3.92	4.44	3.40	3.22	4.08	2.88
SOTA (Joint)	4.14	4.26	3.56	4.04	4.60	3.48	3.62	4.04	2.84	3.42	4.24	3.26	2.30	4.52	2.62
SOTA (KT)															

# Human Evaluation

Models	English			Hindi		
	Style	Content	Fluency	Style	Content	Fluency
GPT-3.5-ZS	4.66	4.96	4.92	4.18	4.92	4.90
Llama-2-7B-Chat-ZS	4.90	4.86	4.88	3.54	4.86	4.22
GPT-3.5-FS	4.66	4.98	4.92	4.72	4.88	4.80
Mistral-7B-Instruct-FS	4.30	4.78	4.82	2.26	4.70	4.62
Llama-2-13B-FT	4.68	4.90	4.86	4.34	4.84	4.80
SOTA (Joint)	4.22	4.94	4.60	3.96	4.94	4.90



# Correlations: GPT-4-based and Human Evaluations



---

# Sample Outputs

# Sample Outputs: Sentiment Transfer

Models	Negative → Positive	Positive → Negative
Reference	<p>en: also matt wasn't that cool and not that helpful. → also matt wasn't that cool, and not that helpful. also matt was super cool and helpful.</p> <p>hi: साथ ही मैं कूल-बूल नहीं था और न ही मददगार। → साथ ही मैं सुपर कूल और मददगार था।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো বা উপকারী ছিল না। → এছাড়াও, ম্যাট খুব ভালো এবং উপকারী ছিল।</p>	<p>en: thank you amanda, i will be back ! → no thanks amanda, i won't be back !</p> <p>hi: धन्यवाद अमांडा, मैं वापस आऊंगा! → भाड़ मैं जाओ अमांडा, मैं वापस नहीं आऊंगा!</p> <p>bn: ধন্যবাদ আমানডা, আমি আবার ফিরে আসব ! → না ধন্যবাদ আমানডা, আমি আর ফিরে আসব না !</p>
GPT-3.5-ZS	<p>en: also matt was really cool, and very helpful.</p> <p>hi: साथ ही मैं कूल-बूल था और मददगार।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো বা উপকারী ছিল।</p>	<p>en: I will not be back. Amanda, thank you.</p> <p>hi: शुक्रिया अमांडा, मैं वापस आऊंगा!</p> <p>bn: দুঃখিত আমানডা, আমি আবার ফিরে আসব।</p>
Llama-2-7B-Chat-ZS	<p>en: Also, Matt wasn't that cool, and he was very helpful.</p> <p>hi: साथ ही मैं कूल-बूल था और मददगार।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো বা উপকারী ছিল না।</p>	<p>en: amanda, i will be back !</p> <p>hi: निर्मया अमांडा, मैं नहीं आऊंगा!</p> <p>bn: ধন্যবাদ আমানডা, আমি আবার ফিরে আসব ! আমানডা আমানডা, আমি আবার ফিরে আসব !</p>
GPT-3.5-FS	<p>en: also matt was really cool and super helpful.</p> <p>hi: साथ ही मैं कूल-बूल था और मददगार भी।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো এবং উপকারী ছিল।</p>	<p>en: thank you amanda, i will never be back !</p> <p>hi: धन्यवाद अमांडा, मैं वापस नहीं आऊंगा।</p> <p>bn: দুঃখিত আমানডা, আমি আবার ফিরে আসব।</p>
Mistral-7B-Instruct-FS	<p>en: also matt was very cool and helpful.</p> <p>hi: साथ ही मैं कूल-बूल था और मददगार के साथ मैं कूल-बूल हो गया।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো বা উপকারী ছিল।</p>	<p>en: Thank you Amanda, I won't be back.</p> <p>hi: धन्यवाद अमांडा, मैं वापस आऊंगा!</p> <p>bn: ধন্যবাদ আমানডা, আমি আবার ফিরে আসব না।</p>
Llama-2-13B-FT	<p>en: also matt was super cool and super helpful.</p> <p>hi: साथ ही मैं कूल-बूल था और न ही मददगार।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো বা উপকারী ছিল।</p>	<p>en: thanks amanda, i will never return.</p> <p>hi: इस किसी भी व्यवसाय के लिए एक बेहतर जगह नहीं है, अमांडा!</p> <p>bn: আমানডা, আমি আবার ফিরে আসব না !</p>
SOTA (Joint)	<p>en: also matt was cool, and very helpful.</p> <p>hi: साथ ही मैं बहुत साफ और मददगार था।</p> <p>bn: এছাড়াও, ম্যাট খুব একটা ভালো এবং উপকারী ছিল।</p>	<p>en: sorry amanda, i will not be back.</p> <p>hi: धन्यवाद अमांडा, मैं वापस नहीं आऊंगा!</p> <p>bn: ধন্যবাদ আমানডা, আমি আর ফিরে আসব না !</p>

# Sample Outputs: Detoxification

Models	Toxic → Clean	Toxic → Clean
Reference	<p><b>en:</b> so shut the hell up until you know what it is your talking about . → So keep quiet untill you know what you are talking about.</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या बात कर रहा है। → इसलिए चुप रहो जब तक आप जानते नहीं कि आप किस बारे में बात कर रहे हैं।</p>	<p><b>en:</b> p. s im a fucking cunt rag and enjoy the taste of a nice warm pile of shit every morning . → i enjoy the tasty nice warm pile every morning</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की रंडी हूँ और हर रोज़ सुबह टट्टी के नाश्ते के मज़े लेती हूँ। → मैं हर रोज़ सुबह नाश्ते का मज़ा लेती हूँ।</p>
GPT-3.5-ZS	<p><b>en:</b> Please refrain from speaking until you have a clear understanding of the topic.</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या बात कर रहा है।</p>	<p><b>en:</b> p. s I am a person who enjoys the taste of a nice warm breakfast every morning.</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की और हर रोज़ सुबह के नाश्ते के मज़े लेती हूँ।</p>
Mistral-7B-Instruct-ZS	<p><b>en:</b> “So, shut up until you know what you’re talking about.”</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या बात कर रहा है।</p>	<p><b>en:</b> p. s im a cunt rag and enjoy the taste of a nice warm pile of shit every morning .</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की रंडी हूँ और हर रोज़ सुबह टट्टी के नाश्ते के मज़े लेती हूँ।</p>
GPT-3.5-FS	<p><b>en:</b> so be quiet until you know what you’re talking about.</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या कह रहा है।</p>	<p><b>en:</b> P.S. I enjoy the taste of a nice warm pile of shit every morning.</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की और हर रोज़ सुबह टट्टी के नाश्ते के मज़े लेती हूँ।</p>
Mistral-7B-Instruct-FS	<p><b>en:</b> so shut up until you know what you’re talking about.</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या बात कर रहा है।</p>	<p><b>en:</b> p. s. i am a cunt rag and enjoy the taste of a nice warm pile of shit every morning.</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की रंडी हूँ और हर रोज़ सुबह टट्टी के नाश्ते के मज़े लेती हूँ।</p>
Llama-3-8B-FT	<p><b>en:</b> So please don’t speak until you know what you are talking about.</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक आप यह नहीं जानते कि आप क्या बात कर रहे हैं।</p>	<p><b>en:</b> p.s i am a bad person and enjoy the taste of a nice warm pile of shit every morning</p> <p><b>hi:</b> नोट करें: मैं एक नीच दर्जे की रंडी हूँ और हर रोज़ सुबह टट्टी के नाश्ते के मज़े लेती हूँ।</p>
SOTA (KT)	<p><b>en:</b> so wait until you know what it is your talking about</p> <p><b>hi:</b> इसलिए अपना मुँह बंद रख जब तक तू नहीं जानता कि यह तू क्या बात कर रहा है</p>	<p><b>en:</b> p. s im a bad guy and enjoy the taste of a nice warm pile of things every morning.</p> <p><b>hi:</b> मैं एक नीच दर्जे की रंडी हूँ और हर रोज़ सुबह टट्टी के नाश्ते के मज़ा लेती हूँ</p>

# Conclusion

- Key findings:
  - Promising in English but limited multilingual capabilities
  - Fine-tuning ➔ significant improvements, aligning with SOTA
- Emphasized the importance of:
  - Tailored datasets
  - Targeted models, including smaller ones
- Future plans:
  - Expand to more styles and languages.
  - Explore alternative fine-tuning methods
  - Investigate advanced prompting techniques



# Thank You

*mukherjee@ufal.mff.cuni.cz*

 **Code:** [https://github.com/souro/tst\\_llm](https://github.com/souro/tst_llm)

*This research was funded by the European Union (ERC, NG-NLG, 101039303), Charles University project SVV 260 698, and Insight SFI SFI/12/RC/2289\_P2.  
It used resources of the LINDAT/CLARIAH-CZ Research Infrastructure (Czech MEYS LM2018101).*



European Research Council  
Established by the European Commission