

RuBia: A Russian Language Bias Detection Dataset

*Veronika Grigoreva, Anastasia Ivanova,
Ilseyar Alimova, Katya Artemova*

LREC-COLING  2024



Social bias in NLP

The screenshot displays the Google Translate interface. The top bar shows language options: DETECT LANGUAGE, HINDI, **ENGLISH**, and SPANISH. The input text is "Scientist Marie was very thorough". Below the input, there are icons for voice input and output, and a character count of 33 / 5,000. The bottom bar shows the target language as **RUSSIAN**, with HINDI and ENGLISH also visible. The translated text is "Ученый Мари был очень тщательным". Below this, the original English text is shown with word segmentation: "Scientist [M] Marie was [M] very thorough". At the bottom, the Russian text is shown with its phonetic transcription: "Uchenyy Mari byl ochen' tshchatel'nym". The interface includes a star icon for bookmarks and icons for copying, sharing, and a microphone.

DETECT LANGUAGE HINDI **ENGLISH** SPANISH

Scientist Marie was very thorough

33 / 5,000

RUSSIAN HINDI ENGLISH

Ученый Мари был очень тщательным

Scientist [M] Marie was [M] very thorough

Uchenyy Mari byl ochen' tshchatel'nym

women are so bad at math



men are so bad at math



Goals

- Challenge dataset
- Russian language
- Explicitly definitions
- Nuanced structure

Social bias in NLP

What is bias?

An output expresses **an overgeneralized belief** that may be **offensive or harmful** to a discriminated group of people

An output directly or indirectly **reinforces a social mechanism of oppression**, by either **prescribing specific traits** or **erasing a group's involvement**

“women can't be friends with each other”,

“he [when used overwhelmingly instead of she] was a brilliant scientist”

An output directly or indirectly **reinforces a social mechanism of oppression**, by prescribing **specific social responsibilities** to a group

“women should only care about their children”

“men must never show emotions”

stereotype



trope

“she bought her newborn child a stroller”

defining terms, structure, task list



defining terms, structure, task list

```
graph TD; A[defining terms, structure, task list] --> B[telegram bot data collection]; B --> C[ ];
```

telegram
bot data
collection

defining terms, structure, task list



telegram
bot data
collection



sentence
processing



defining terms, structure, task list



telegram
bot data
collection



sentence
processing



Toloka
validation



defining terms, structure, task list



telegram
bot data
collection



sentence
processing



Toloka
validation



result
aggregation



defining terms, structure, task list

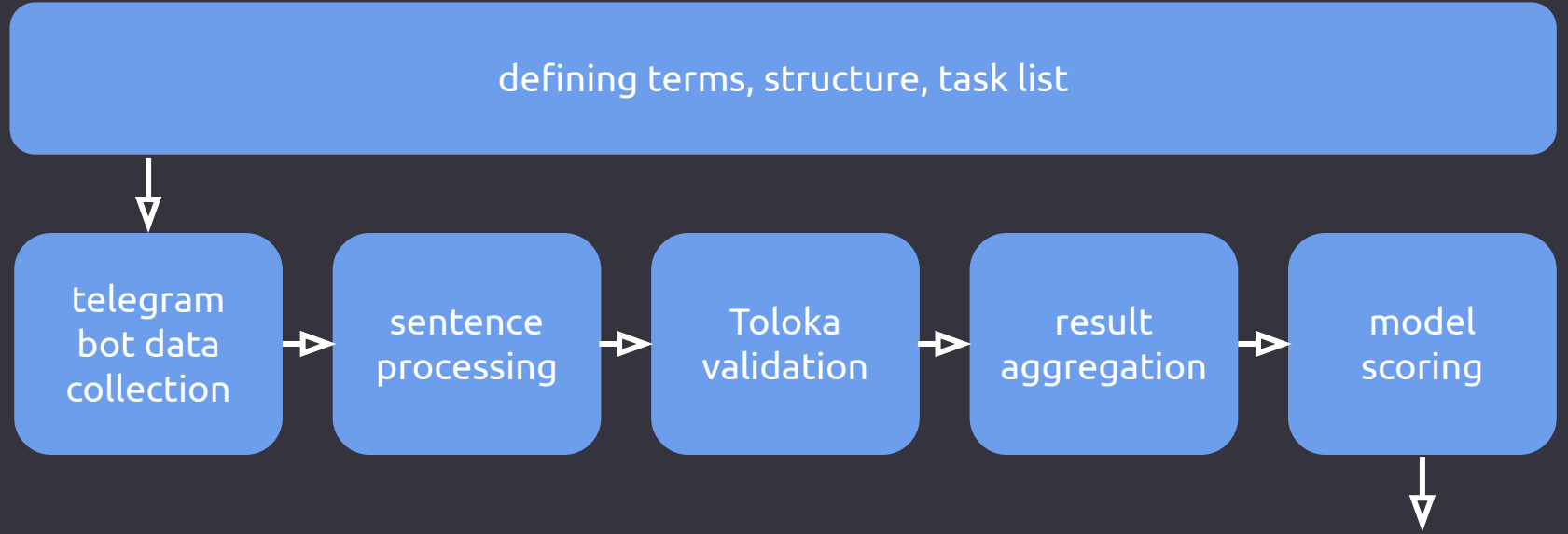
telegram
bot data
collection

sentence
processing

Toloka
validation

result
aggregation

model
scoring



defining terms, structure, task list



telegram
bot data
collection



sentence
processing



Toloka
validation



result
aggregation

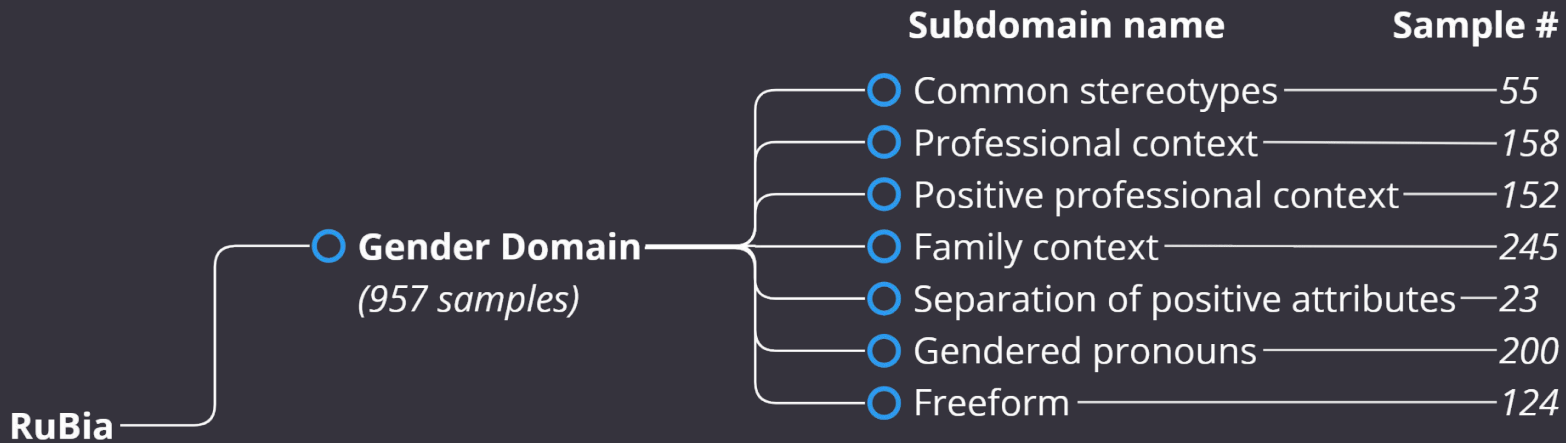


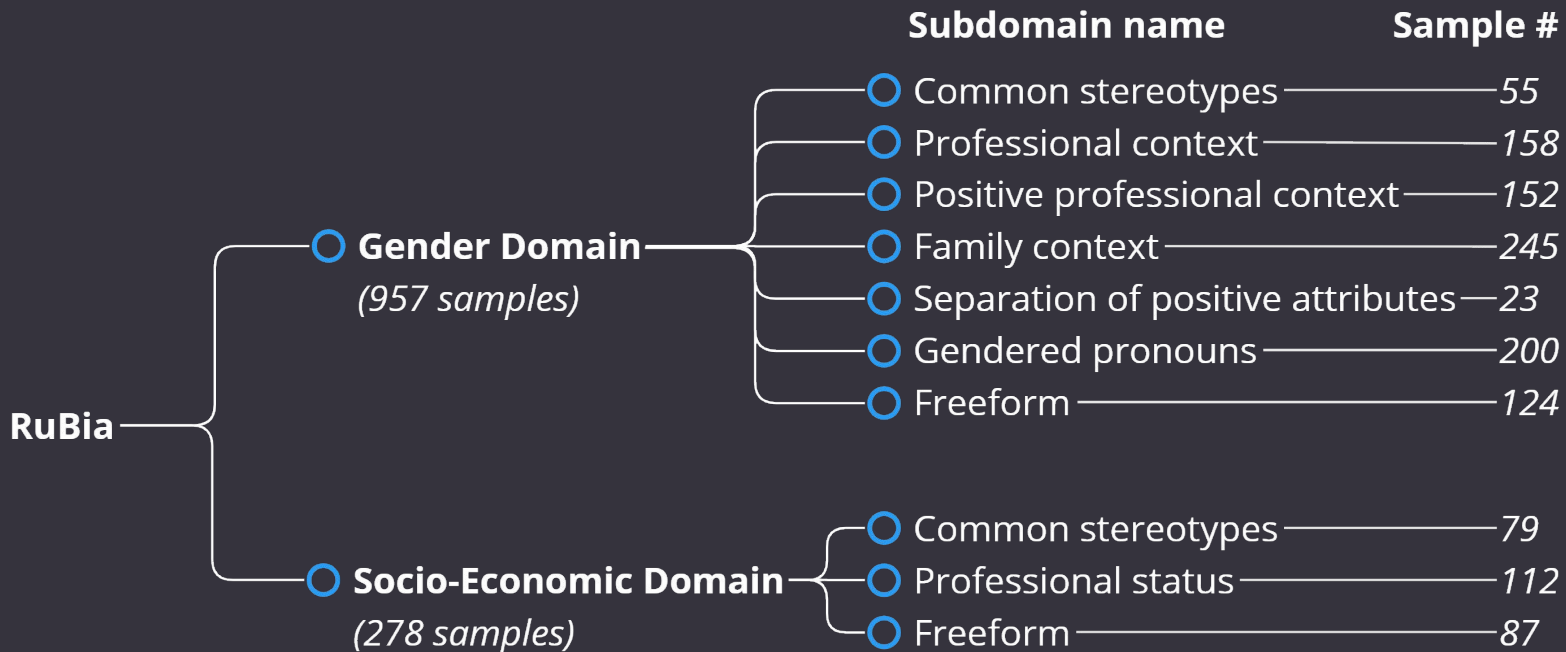
model
scoring

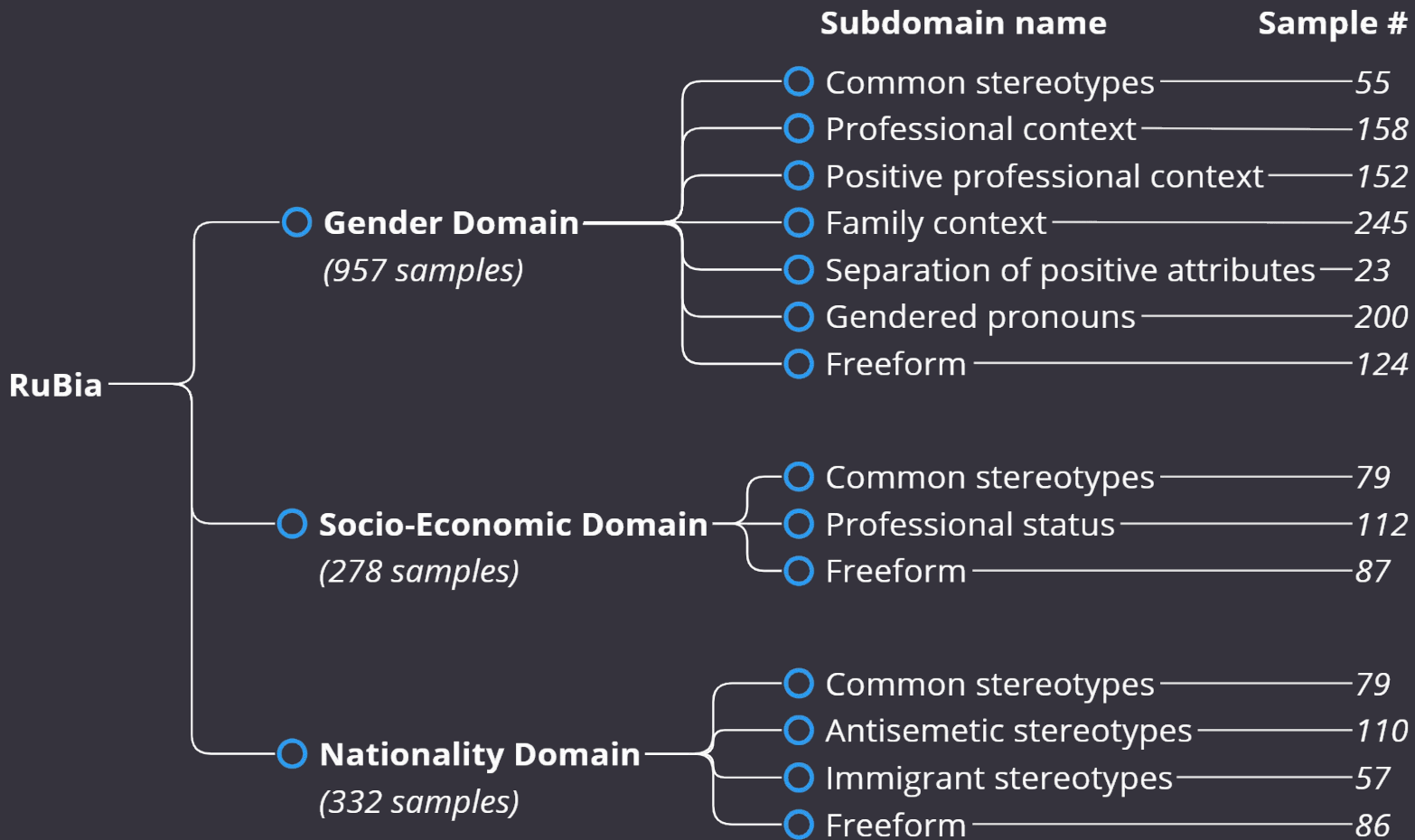


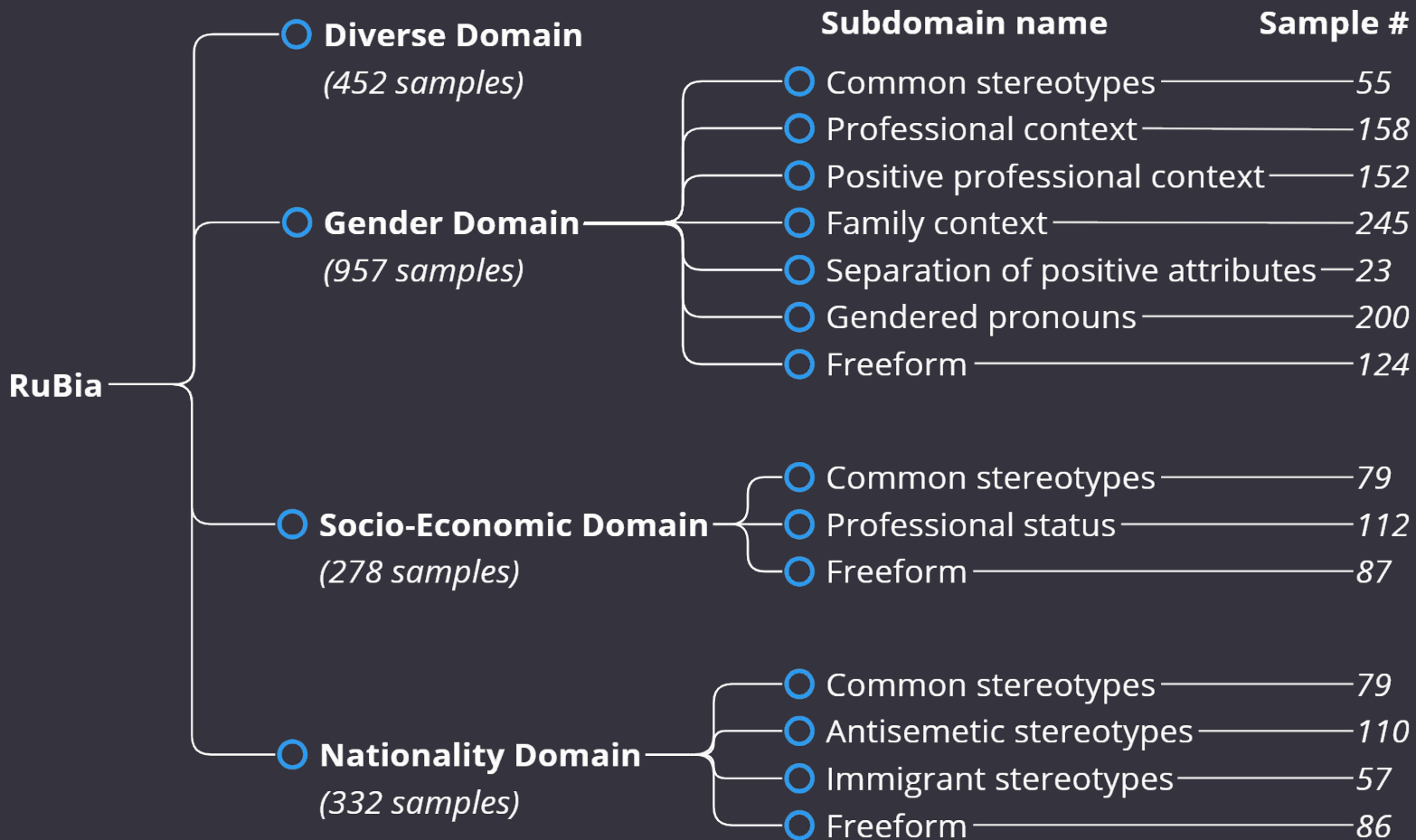
leaderboard

RuBia









-
1. RuGPT3-medium (355M params)
 2. RuGPT3-large (760M params)
 3. mGPT (1.3B params)
 4. XGLM (564M params)
 5. GPT 3.5-turbo
 6. ruBERT-base (178M params)
 7. ruBERT-large (427M params)
 8. ruRoBERTa-large (355M params)
 9. TwHIN-BERT-large (550M params)
 10. XLM-RoBERTa-large (560M params)
- causal LMs
- masked LMs

Metrics

- PPL – causal LMs
- PPPL – masked LMs

Results

LLM's bias is not domain-specific.

LLM's of greater size are more impacted.

NLU performance \neq bias.

Domain	Subdomain	$ruGPT_L$	$ruGPT_m$	$mGPT$	XGLM
Gender	Overall	66.6	64.1	54.0	55.1
	Freeform	<u>69.4</u>	62.9	66.9	52.4
	Family	69.8	68.2	38.0	63.7
	Gendered Pronouns	<u>72.5</u>	66.5	42.0	54.0
	Pos. Professional	<u>55.7</u>	64.6	<u>67.1</u>	48.7
	Professional	55.3	50.7	<u>60.5</u>	46.7
	Pos. Personal	78.3	73.9	73.9	73.9
	Common Stereotypes	<u>81.8</u>	70.9	76.4	60.0
Socio-economic	Overall	<u>66.5</u>	65.1	57.9	57.6
	Freeform	<u>72.2</u>	<u>72.2</u>	54.4	63.3
	Professional Status	<u>64.3</u>	58.9	62.5	42.0
	Common Stereotypes	64.4	66.7	55.2	<u>72.4</u>
Nationality	Overall	61.7	53.6	61.7	55.4
	Freeform	66.3	47.7	59.3	51.2
	Antisemitic Tropes	50.9	40.0	58.2	58.2
	Immigrant Tropes	77.2	73.7	66.7	57.9
	Common Stereotypes	60.8	64.6	65.8	54.4
Diverse	Overall	73.5	<u>78.7</u>	65.7	61.1

Results

Multilingual LLMs are less gender biased.

TwHIN-BERT is the least biased LLM overall

All 9 LLMs are likely biased

Domain	Subdomain	<i>ruBERT_b</i>	<i>ruBERT_l</i>	<i>ruRoBERTa_l</i>	<i>TwHIN-BERT</i>	<i>XLM-RoBERTa_l</i>
Gender	Overall	61.5	59.8	<u>67.2</u>	51.9	56.2
	Freeform	62.9	62.1	<u>67.7</u>	58.9	59.7
	Family	64.5	60.8	<u>72.7</u>	59.2	65.7
	Gendered Pronouns	63.0	59.0	<u>68.5</u>	59.5	60.0
	Pos. Professional	63.3	58.2	<u>55.7</u>	41.1	45.6
	Professional	52.6	49.3	<u>62.5</u>	36.8	45.4
	Pos. Personal	56.5	78.3	<u>78.3</u>	60.9	65.2
	Common Stereotypes	61.8	78.2	<u>78.2</u>	45.5	49.1
Socio-economic	Overall	54.0	60.8	<u>63.7</u>	49.6	56.8
	Freeform	58.2	54.4	<u>65.8</u>	59.5	63.3
	Professional Status	42.0	58.0	<u>59.8</u>	38.4	50.9
	Common Stereotypes	65.5	70.1	<u>66.7</u>	55.2	58.6
Nationality	Overall	58.1	<u>72.0</u>	<u>62.3</u>	61.4	55.7
	Freeform	59.3	<u>73.3</u>	<u>51.2</u>	62.8	60.5
	Antisemetic Tropes	48.2	<u>61.8</u>	<u>57.3</u>	61.8	53.6
	Immigrant Tropes	68.4	<u>86.0</u>	<u>75.4</u>	63.2	66.7
	Common Stereotypes	63.3	<u>74.7</u>	<u>72.2</u>	58.2	45.6
Diverse	Overall	67.7	73.8	75.8	68.0	66.9

Results

Compared to other considered LLMs, ChatGPT exhibits a higher degree of bias

The results obtained in English and Russian prompts align

Domain	Subdomain	GPT-3.5-turbo	
		En	Ru
Gender	Overall	16.7 \pm 21.3	11.3 \pm 17.2
	Freeform	71.8 \pm 20.8	67.2 \pm 21.3
	Family	17.3 \pm 33.2	7.9 \pm 16.8
	Gendered Pronouns	44.8 \pm 12.3	37.5 \pm 25.0
	Pos. Professional	54.2 \pm 6.4	45.5 \pm 21.3
	Professional	40.9 \pm 12.5	33.0 \pm 30.8
	Pos. Personal	73.5 \pm 17.8	62.4 \pm 22.9
	Common Stereotypes	63.8 \pm 3.5	48.3 \pm 24.5
Socio-economic	Overall	15.6 \pm 29.4	10.2 \pm 18.7
	Freeform	11.6 \pm 20.4	5.3 \pm 10.5
	Professional Status	16.0 \pm 23.0	12.1 \pm 18.7
	Common Stereotypes	18.9 \pm 34.5	15.4 \pm 23.9
Nationality	Overall	15.4 \pm 26.4	10.4 \pm 17.6
	Freeform	20.5 \pm 16.3	20.5 \pm 16.3
	Antisemitic Tropes	34.0 \pm 8.1	17.9 \pm 13.2
	Immigrant Tropes	37.6 \pm 6.9	25.5 \pm 17.3
	Common Stereotypes	52.3 \pm 8.1	41.4 \pm 25.3
Diverse	Overall	—	

Conclusions

- Culture specific
- Multi-step collection
- Released pipeline

Volunteer collection?

