

PERSONALIZED NLP



Przemysław Kazienko, Jan Kocoń
Department of Artificial Intelligence
Wrocław University of Science and Technology, Poland



AGENDA


1. Example and motivation
2. Subjective NLP tasks
3. Measuring diversity
4. Perspectives
5. Research on offensive content
6. Research on emotional dataset
7. Research on multiple tasks
8. Conclusions





1

MOTIVATION



"Your behaviour is inappropriate and your reaction is exaggerated. I am not sure if you should have administrator rights."

Wikipedia Detox Aggression

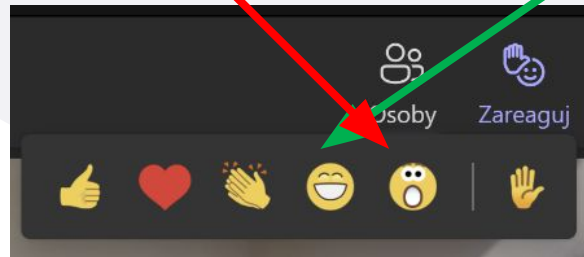
**Do you think, it is
aggressive or not?**



"Your behaviour is inappropriate and your reaction is exaggerated. I am not sure if you should have administrator rights."

Wikipedia Detox Aggression

Do you think, it is
aggressive or **not?**



MOTIVATION

COMMON GENERALIZED NLP

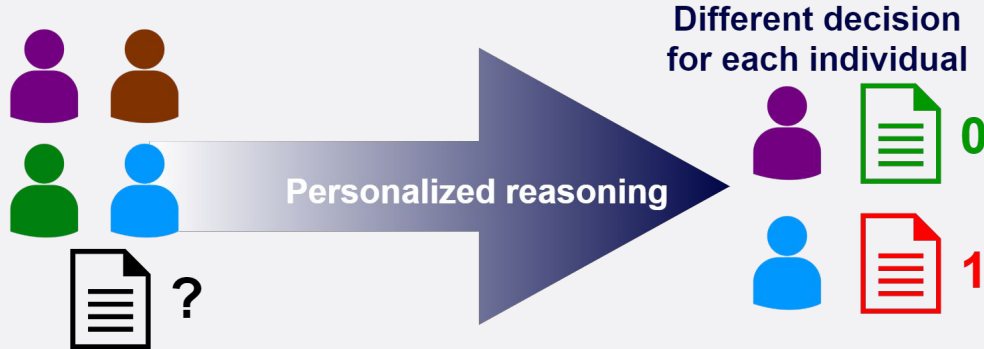


MOTIVATION

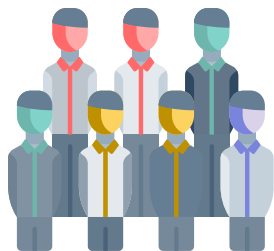
COMMON GENERALIZED NLP



OUR PERSONALIZED NLP



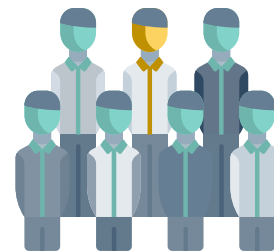
MOTIVATION



Representativeness

Hard to **acquire** data (annotations) from **all** social groups representing all diverse beliefs

"The people like me are not respected by the system"



Fairness

Common generalized solutions are **biased** toward the mainstream

"Since the system does not regard my individual beliefs, I do not trust in it"



2 SUBJECTIVE NLP TASKS

SUBJECTIVE NLP TASKS

1. **Reader** perspective: **perception** prediction

- a. **Emotions** (many models, multiple dimensions)
- b. **Offensive** content detection, incl. aggression, toxic, hate speech, cyberbullying, hostile, insulting
- c. **Humor**, funny
- d. Sarcasm and irony detection
- e. Antagonistic, provocative, trolling speech detection
- f. Counterspeech detection
- g. Hope, supportive speech detection
- h. Obscene language detection
- i. Dismissive, patronising, condescending
- j. Unfair generalisation
- k. Slur usage
- l. Unpalatable questions
- m. Persuasiveness
- n. Inflammatory text
- o. Subjective perception of **sentiment** polarization

2. **Author** perspective

- a. Sentiment analysis
- b. Content generation (e.g. style-based), summarization, adjustment

3. **Mixed**

- a. Conversations

The tasks often overlap



3

MEASURING DIVERSITY

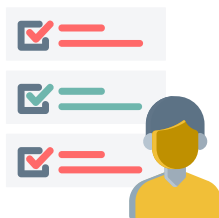
[Kan21, Mit21, Koc21b]

MEASURING DIVERSITY



Document-oriented

Document **Controversy**
(entropy-based) [Kan21]



Human-oriented

Human **Conformity**; general,
weighted, class-based [Kan21]

HB-measure – Human Bias
[Koc21b]; aggregated Z-score; for
emotions: PEB – **Personal
Emotional Bias** [Mit21]



Collection-oriented

Krippendorff's alfa [Koc21a]

WAVE kappa – Wroclaw
Annotators Variability Estimator;
Fleiss' kappa aggregated over
different no. of users [Koc21a]

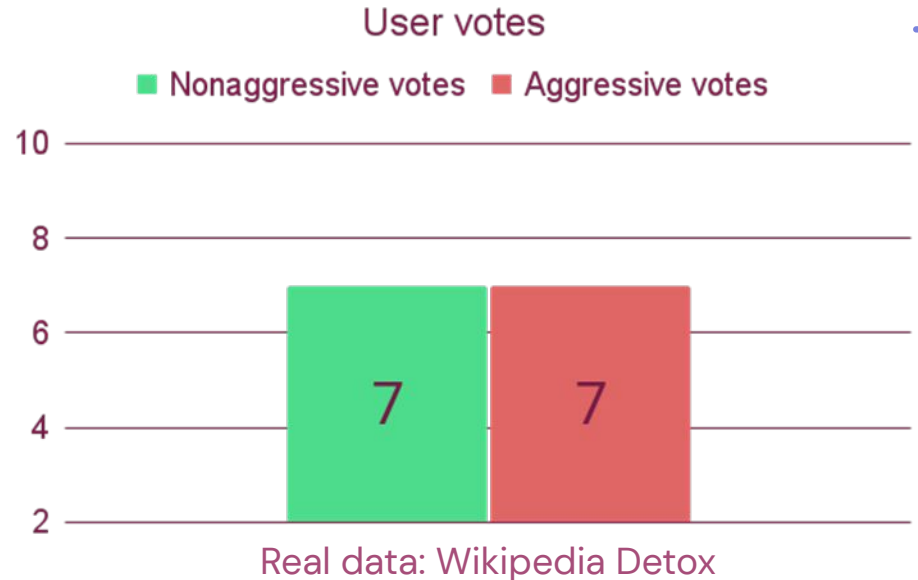
CONTROVERSY MEASURE

***“Your behaviour is inappropriate and your reaction is exaggerated.
I am not sure if you should have administrator rights.”***



CONTROVERSY = 1.0
(entropy-based)

$$\text{Contr}(d) = \begin{cases} 0, & \text{if } n_d^0 = n_d \vee n_d^1 = n_d \\ - \sum_{c=0,1} \frac{n_d^c}{n_d} \log_2 \left(\frac{n_d^c}{n_d} \right) & ; \end{cases}$$



CONTROVERSY MEASURE

inappropriate

**“Your behaviour is *terrible* and your reaction is exaggerated.
I am not sure if you should have administrator rights.”**

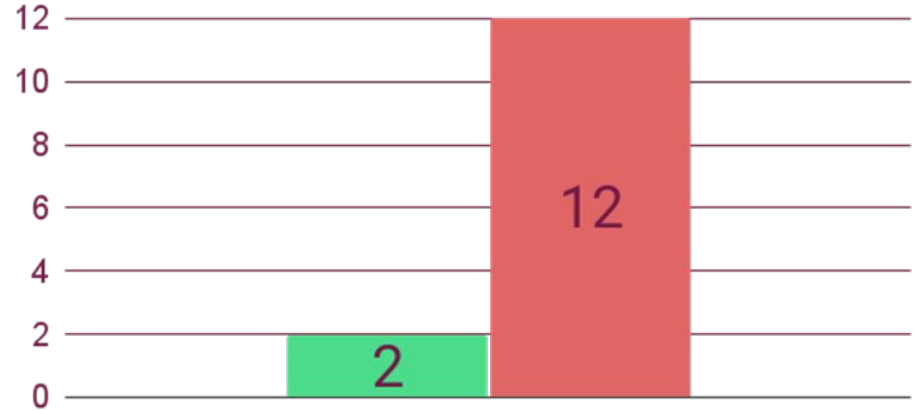


CONTROVERSY = 0.59 ↓
(entropy-based)

$$\text{Contr}(d) = \begin{cases} 0, & \text{if } n_d^0 = n_d \vee n_d^1 = n_d \\ - \sum_{c=0,1} \frac{n_d^c}{n_d} \log_2 \left(\frac{n_d^c}{n_d} \right) & ; \end{cases}$$

User votes

■ Nonaggressive votes ■ Aggressive votes



CONFORMITY MEASURE

***“Your behaviour is inappropriate and your reaction is exaggerated.
I am not sure if you should have administrator rights.”***

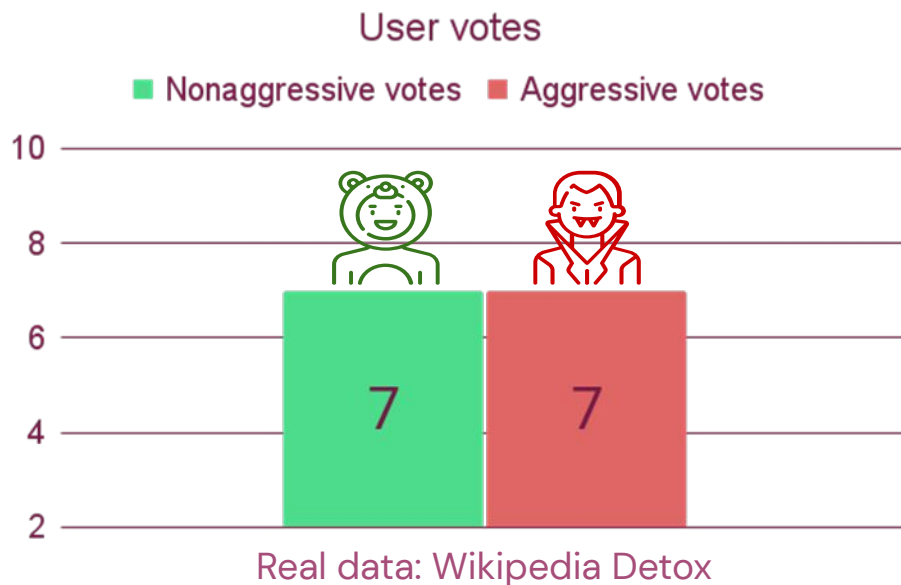


CONFORMITY = 0.50



CONFORMITY = 0.50

$$GConf(a, C) = \frac{\sum_{d \in A_a} \mathbb{1}_{\{l_d \in C \wedge l_d = l_{d,a}\}}}{\sum_{d \in A_a} \mathbb{1}_{\{l_d \in C\}}}$$



CONFORMITY MEASURE

*“Your behaviour is **terrible** and your reaction is exaggerated.
You **don't deserve** administrator rights.”*

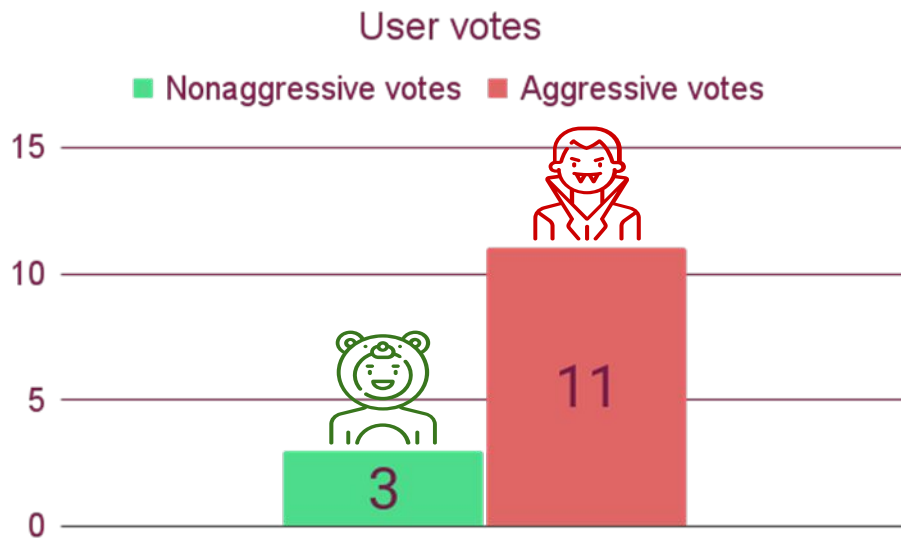


CONFORMITY = 0.21 = $\frac{3}{14}$



CONFORMITY = 0.79

$$GCon f(a, C) = \frac{\sum_{d \in A_a} \mathbb{1}\{l_d \in C \wedge l_d = l_{d,a}\}}{\sum_{d \in A_a} \mathbb{1}\{l_d \in C\}}$$





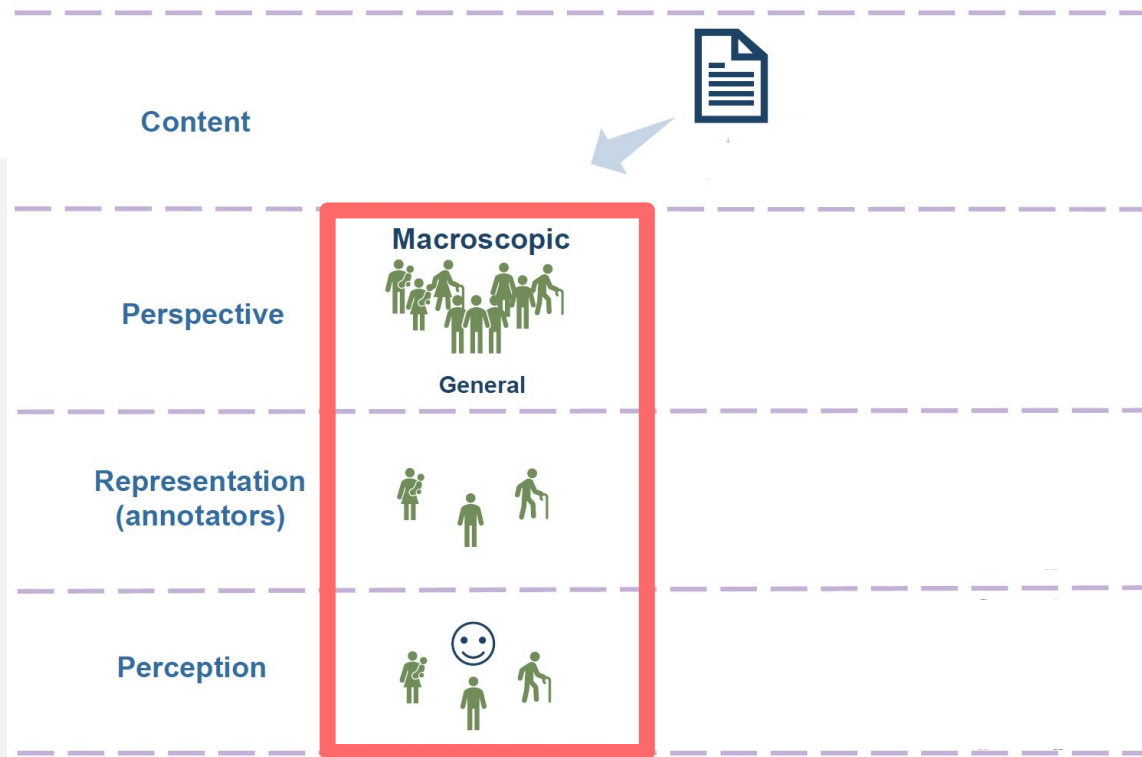
4

PERSPECTIVES

[Koc21a]



PERSPECTIVES: MACROSCOPIC



[Koc21a]




PERSPECTIVES: MACROSCOPIC



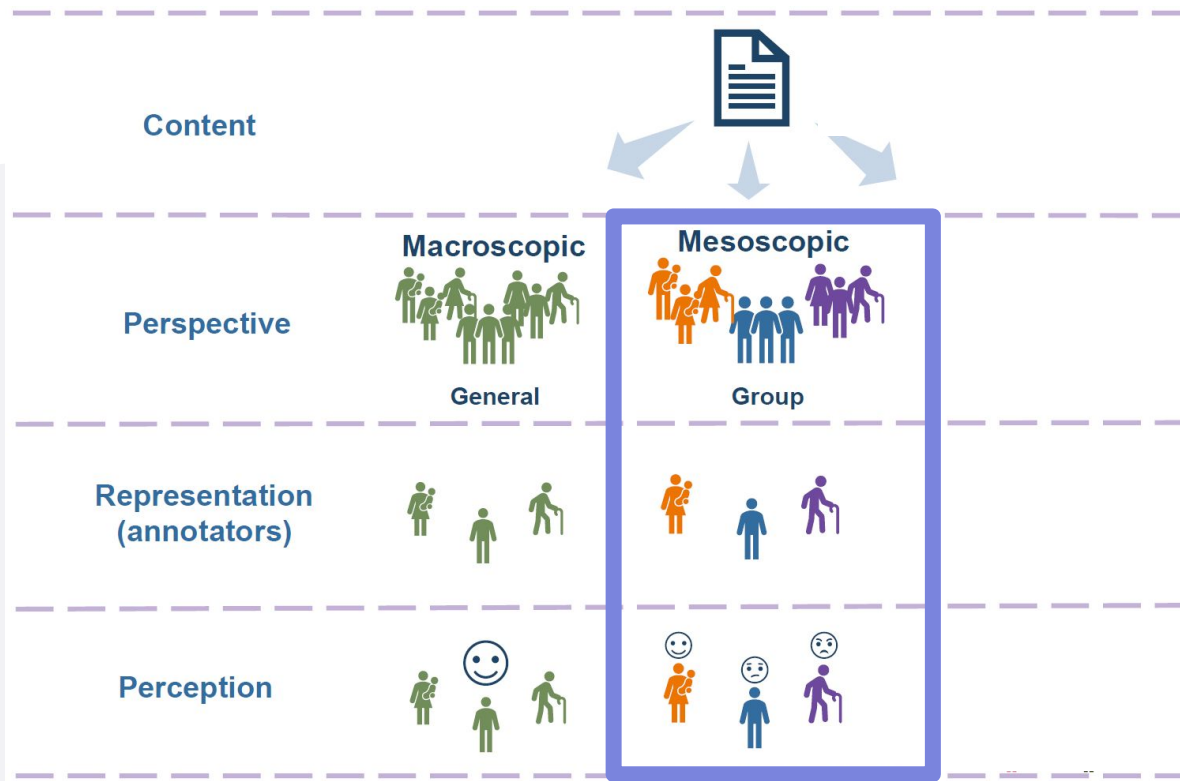
(general)



| Perspective profile | Statement | Information source | Annotation |
|--|---|--|---|
| <p>Society-based, global, general.</p> <p>Used in most research.</p> <p>Assumes the existence of common perception of the content</p> | <p><i>"People generally treat some content offensive/funny/sad/..."</i></p> | <p>(1) content (2) context of the content, e.g. source</p> | <p>Several trained/expert  annotators are able to express common perception (beliefs)</p> |



PERSPECTIVES: MESOSCOPIC



[Koc21a]



PERSPECTIVES: MESOSCOPIC

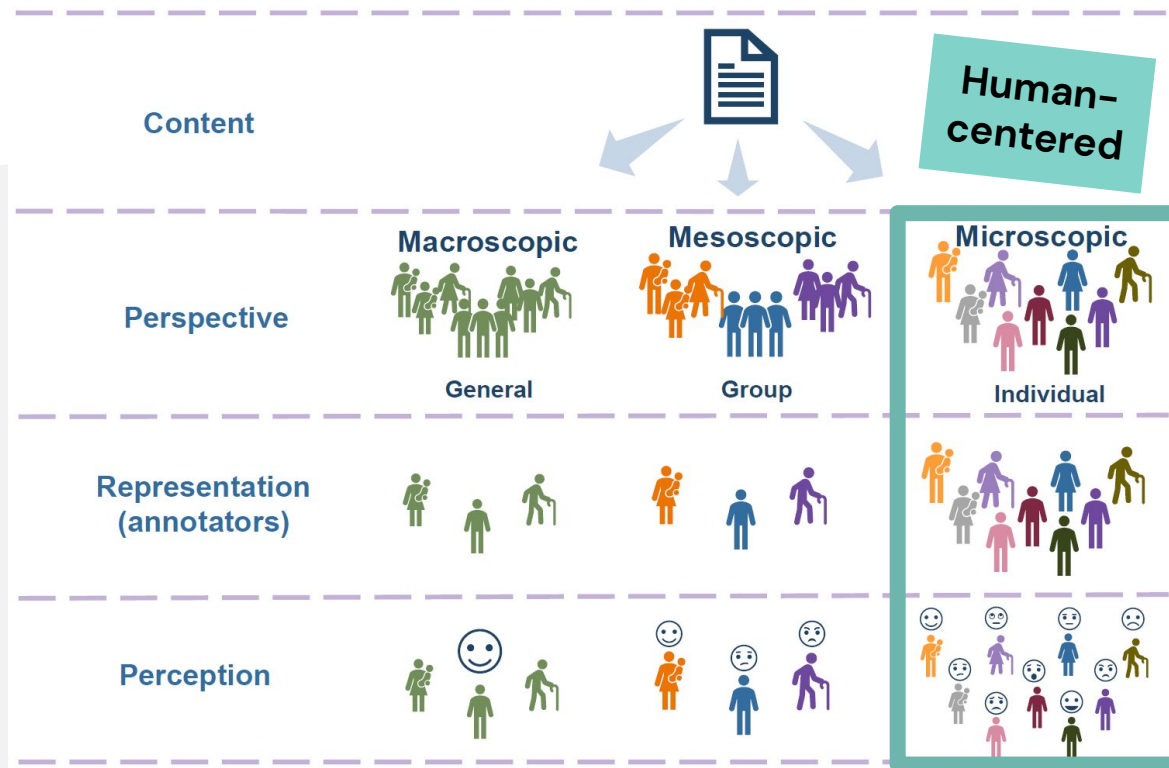
(group-based)



| Perspective profile | Statement | Information source | Annotation |
|---|--|--|--|
| Group-based, social or demographic groups. Perception is shared in social groups | <i>"There are some groups of people who perceive the content in the same way as offensive/funny/sad/..."</i> | (1) content (2) context of the content (3) group demographic profile , e.g. age (4) group context , e.g. culture, shared personality traits, religion | A lot of annotations per document ● are required. Annotator profiles need to be collected (surveys, behaviour) |



PERSPECTIVES: MICROSCOPIC



[Koc21a]



PERSPECTIVES: MICROSCOPIC (personalized)

• Human-centered

| Perspective profile | Statement | Information source | Annotation |
|---|---|---|---|
| Individual, fully personalized. Each individual may perceive content differently . | <i>"Perception of the content depends on a single human, i.e. on their individual and temporal concext"</i> | (1) content (2) context of the content (3) individual behaviour (4) individual demographics (5) individual social context (relationships with the author and the social group) (6) temporal affective state (mood, emotions) | An individual annotator beliefs need to be identified using surveys and/or previous annotations |

PERSONALIZED NLP: What we need?



Data about
human beliefs

Texts **earlier** annotated by a
given individual



Agreed, generalized
labels are useless

Usually obtained by
majority voting





5

RESEARCH ON OFFENSIVE CONTENT

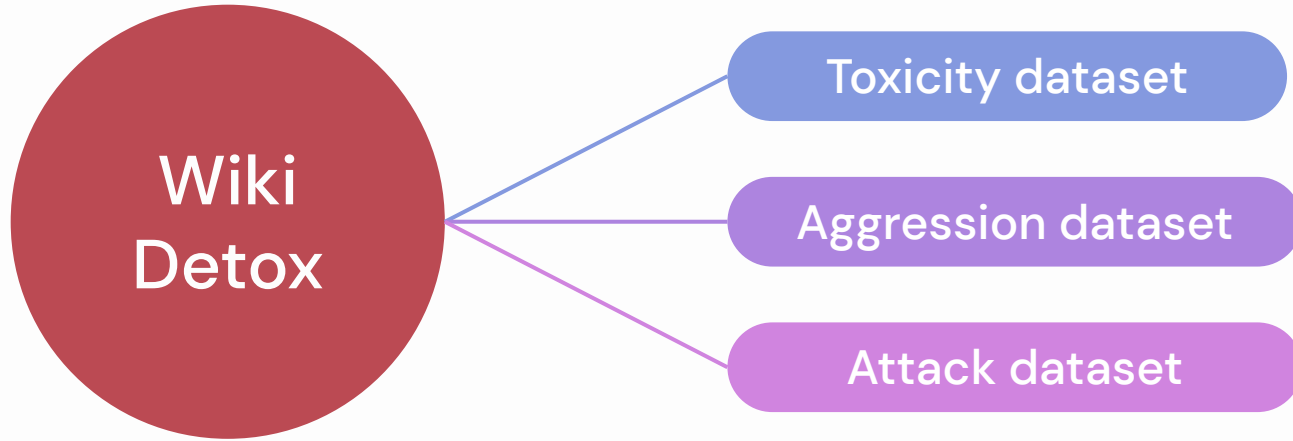
[Koc21a, Kan21, Koc21b]



5a

**OFFENSIVE CONTENT:
ANNOTATED DATA**

WIKI DETOX DATASETS (English)



Publicly available

WIKI: Toxicity



Classes

2

Texts

159,686

People

4,301

Annotations

1,598,289

Controversial Texts

40.5 %



WIKI: Aggression & Attack



Classes

2

Texts

115,864

People

4,053

2,190

Annotations

1,365,217

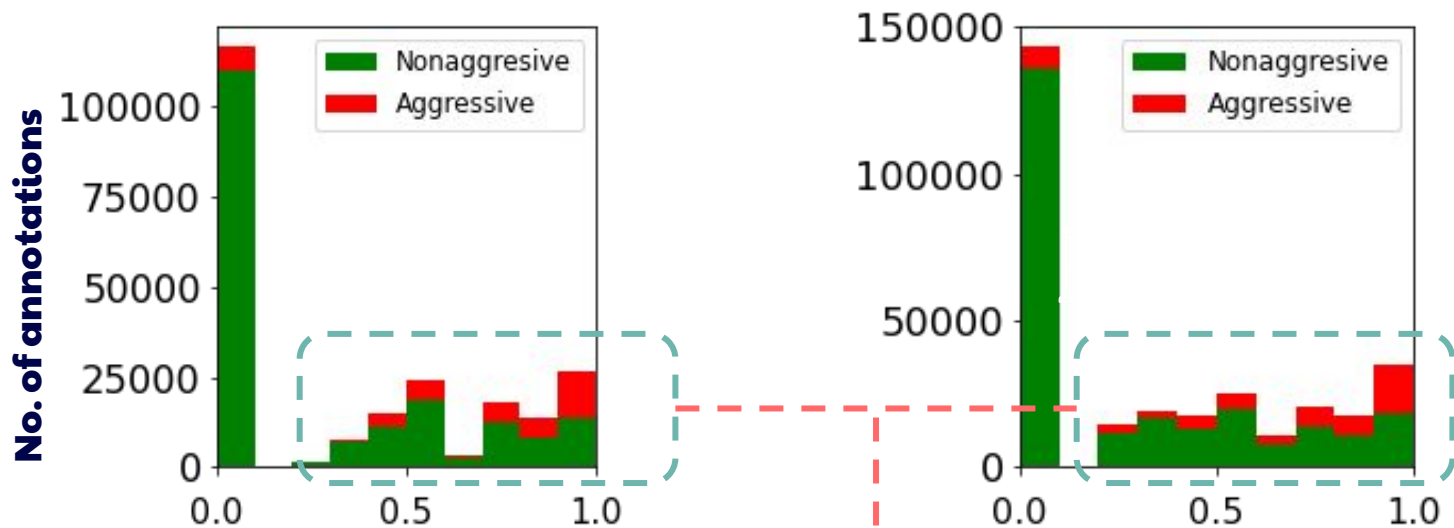
855,514

Controversial Texts

51.3% & 48%



WIKI: Aggressive



**Disagreement in ~50%
of annotations**



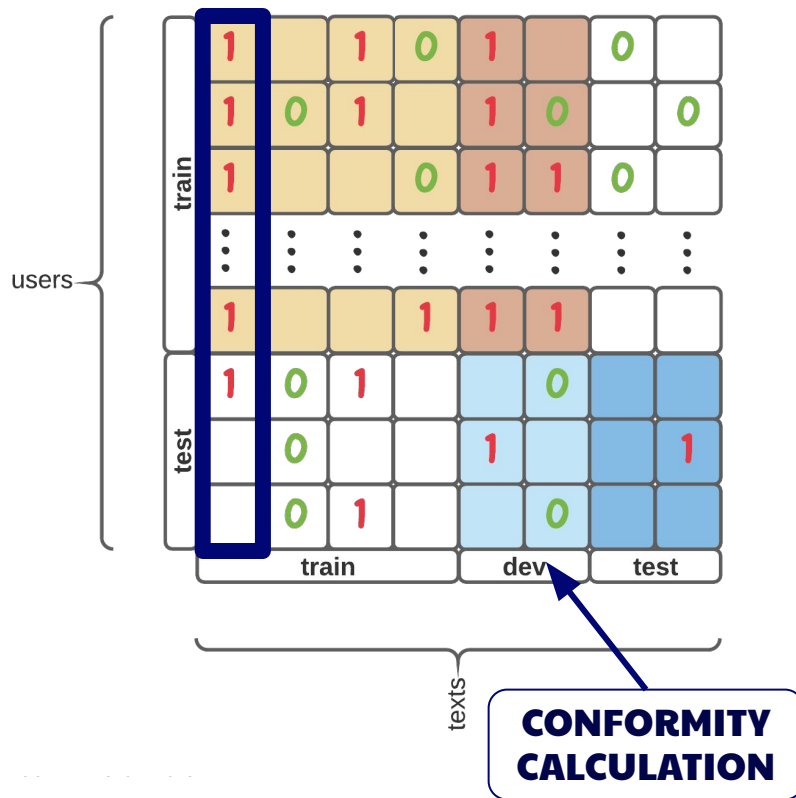
| | | | | | | | |
|-------|-------|---|---|-----|---|------|---|
| train | 1 | 0 | 1 | 0 | 1 | 0 | |
| | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| | 1 | | | 0 | 1 | 1 | 0 |
| | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| test | 1 | | | 1 | 1 | 1 | |
| | 1 | 0 | 1 | | 0 | | |
| | | 0 | 1 | | 1 | | 1 |
| | train | | | dev | | test | |

5b

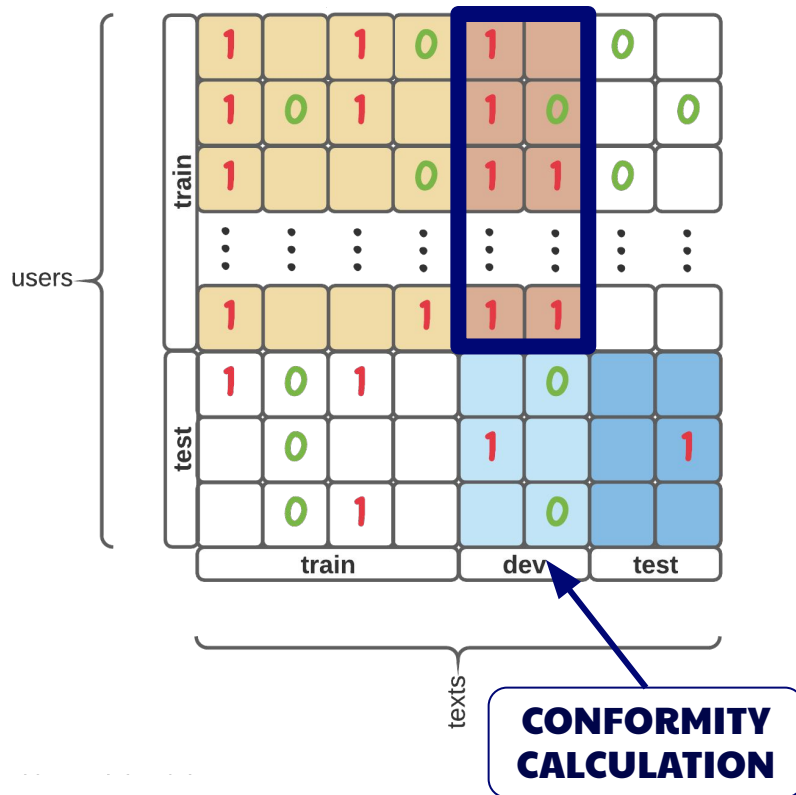
OFFENSIVE CONTENT: DATA SPLIT

Train-dev-test

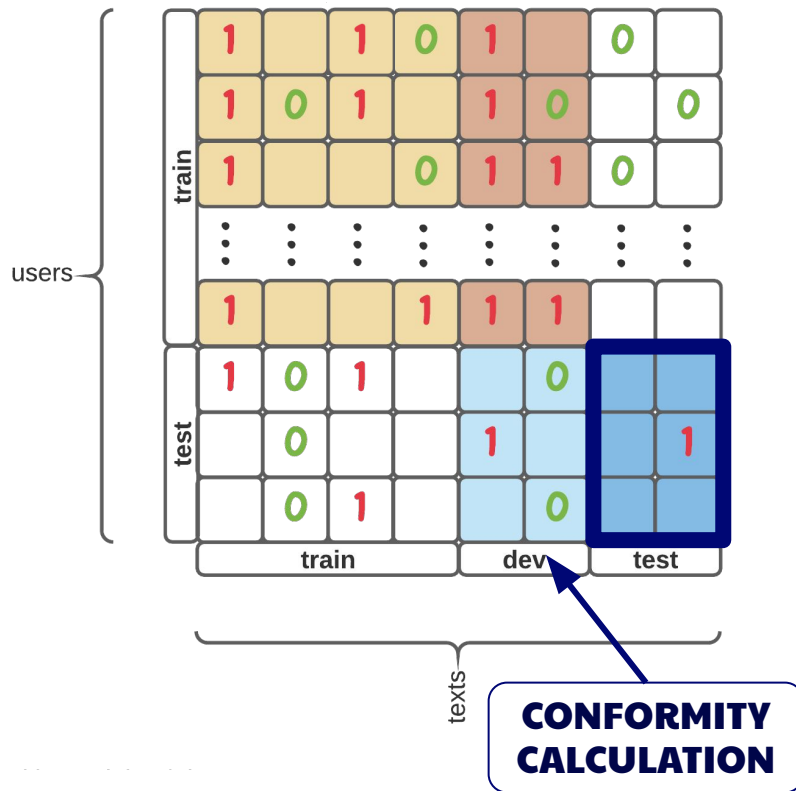
DATASET SPLIT: Wiki



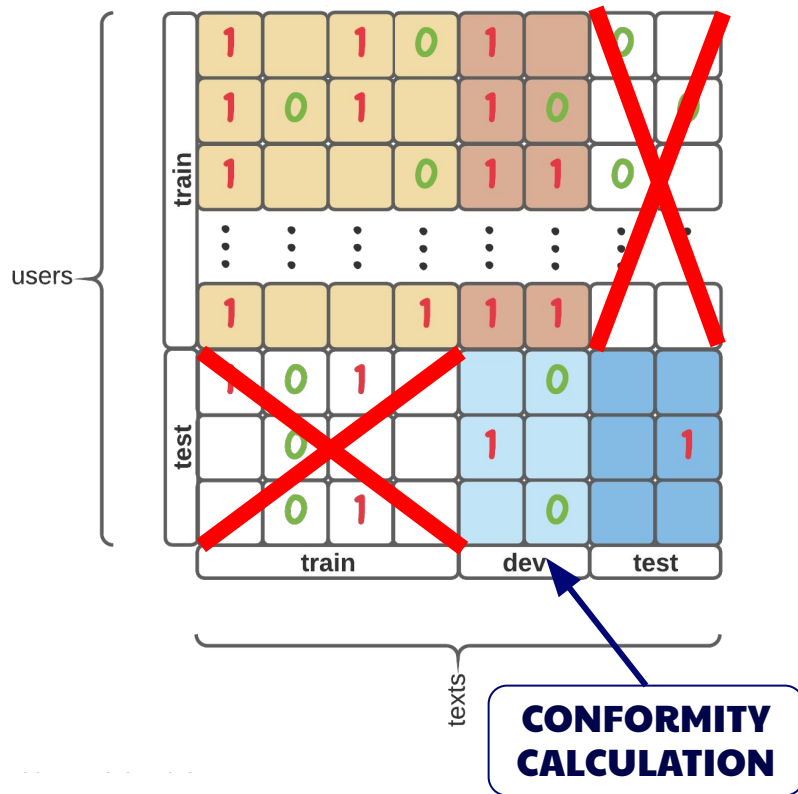
DATASET SPLIT: Wiki



DATASET SPLIT: Wiki



DATASET SPLIT: Wiki

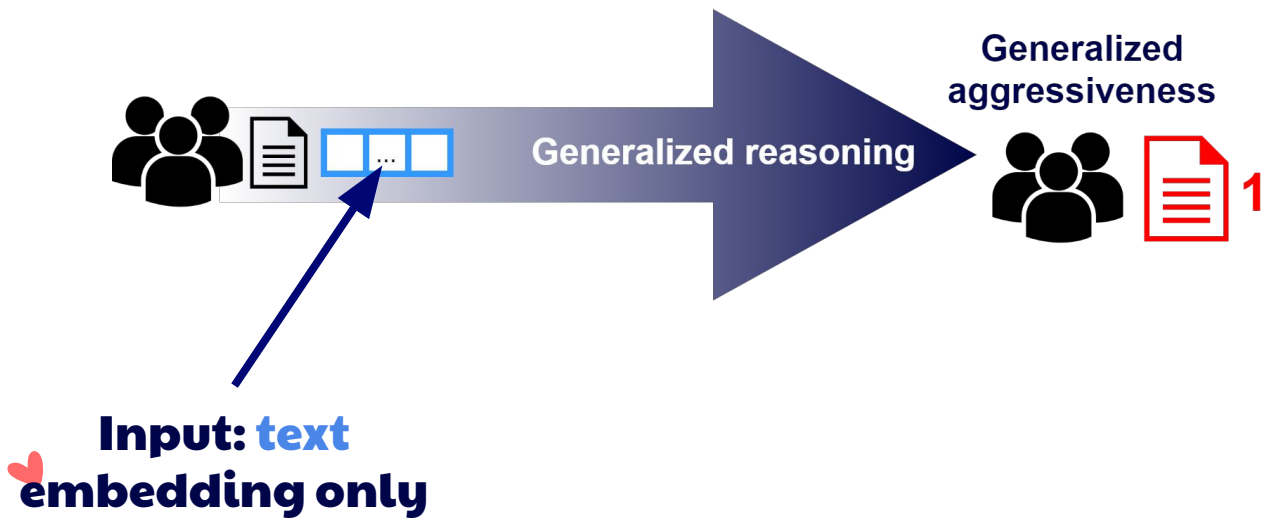




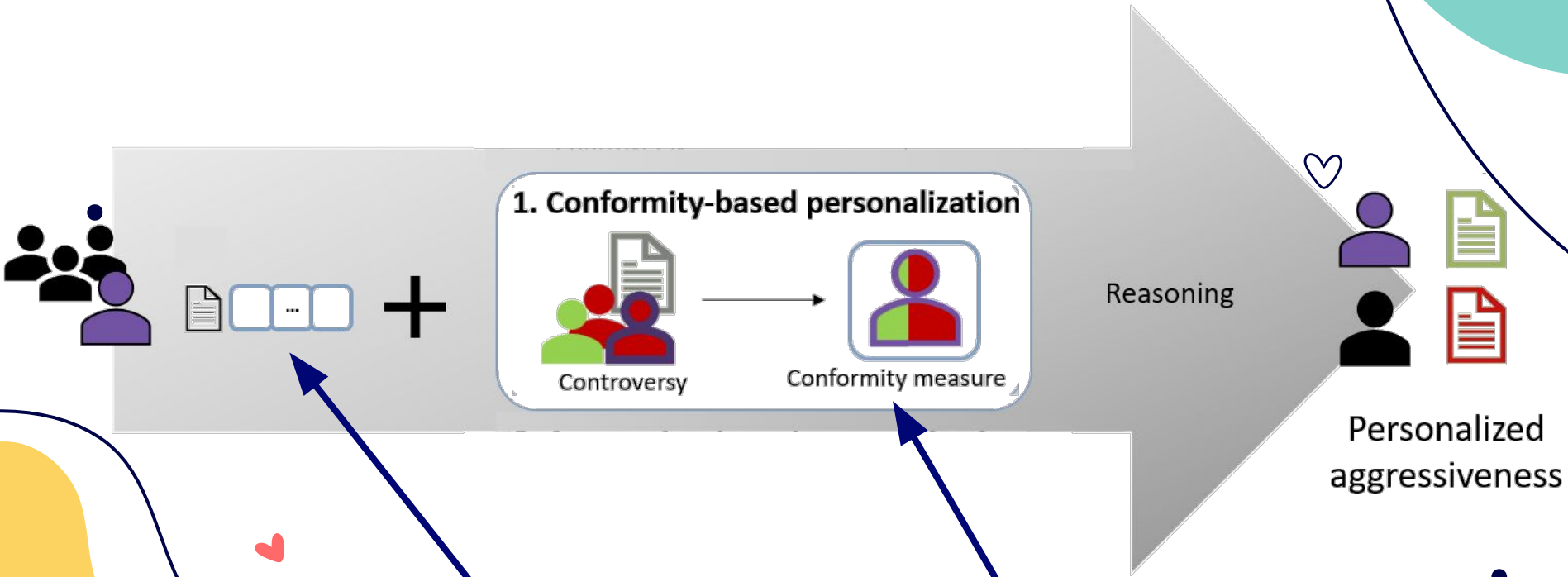
5c

**OFFENSIVE CONTENT:
METHODS**

GENERAL METHOD - BASELINE

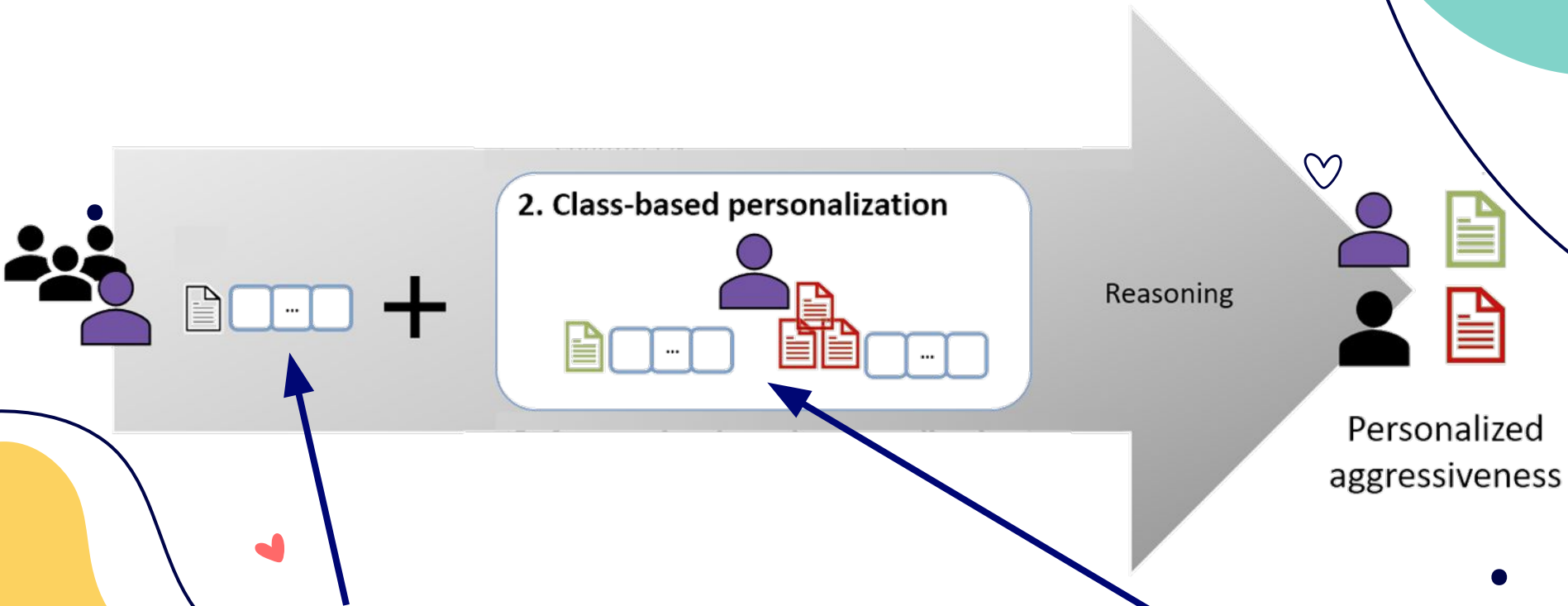


1. CONFORMITY-BASED PERSONALIZATION



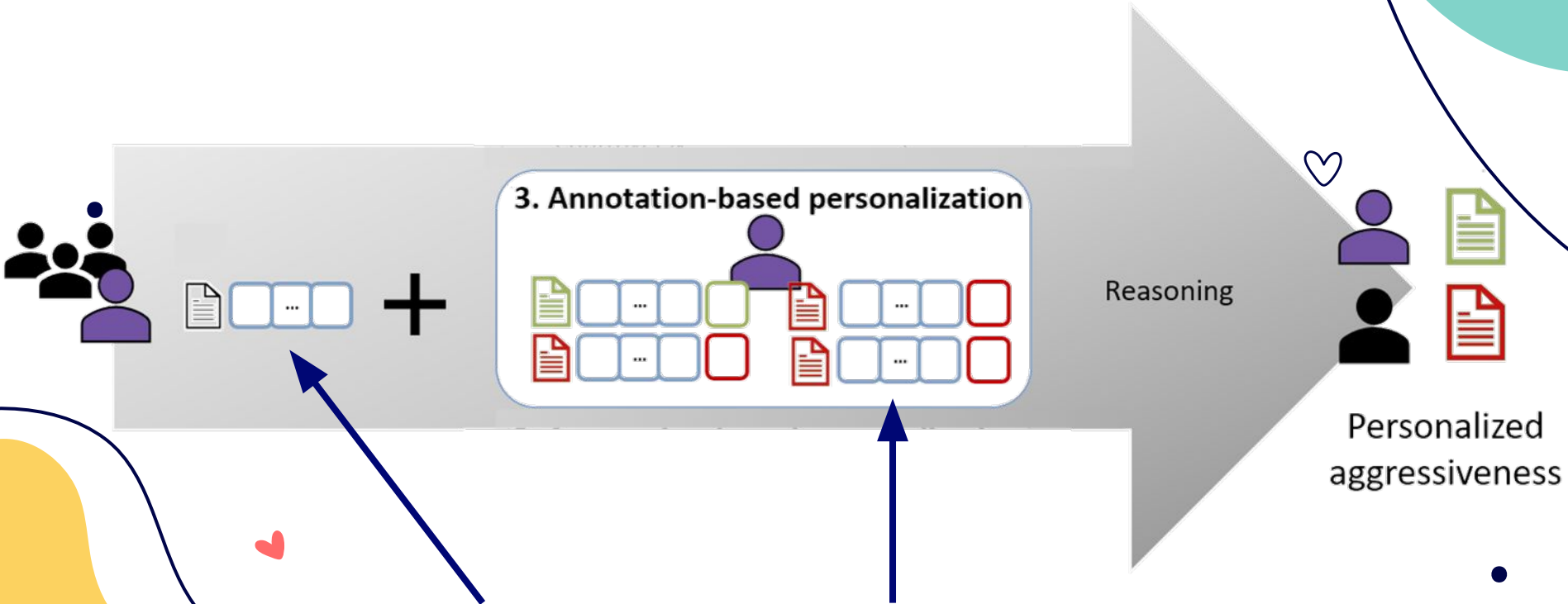
Input: text embedding + user conformity measures (6 features)

2. CLASS-BASED PERSONALIZATION

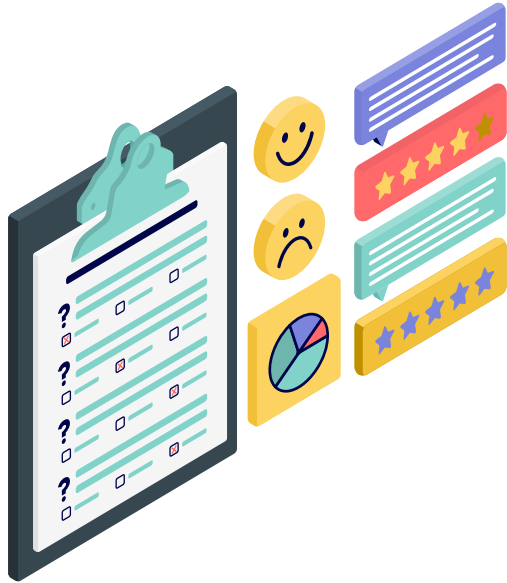


Input: text embedding + texts seen by user as **aggressive / **non-aggressive** (avg. of their embeddings)**

3. ANNOTATION-BASED PERSONALIZATION



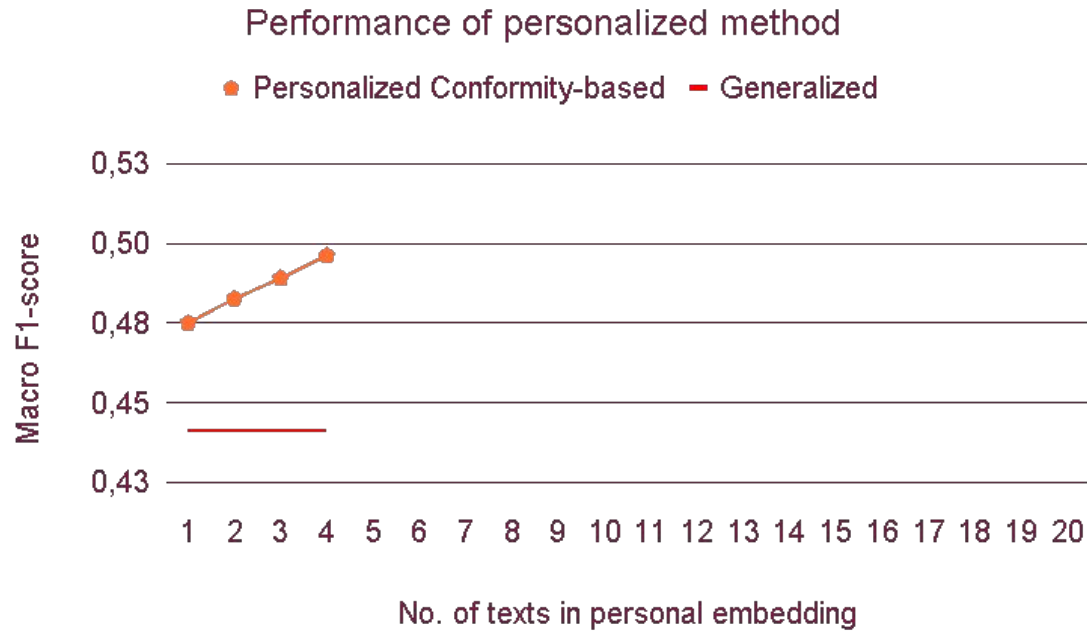
Input: text embedding + all texts prev. seen by the user with their annotations 1 - 0, raw embeddings



5d

**OFFENSIVE CONTENT:
RESULTS**

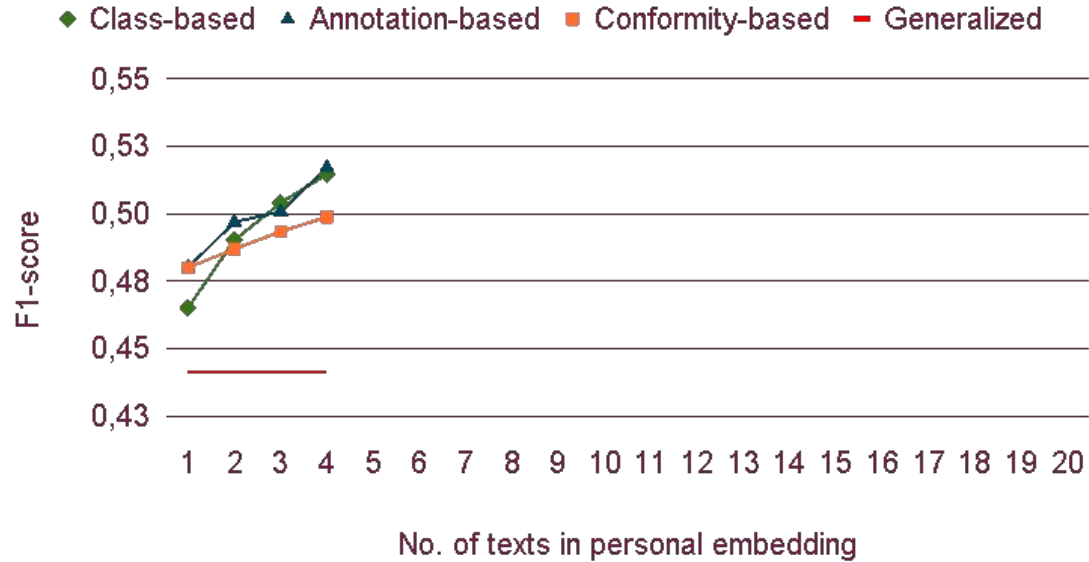
EVALUATION RESULTS



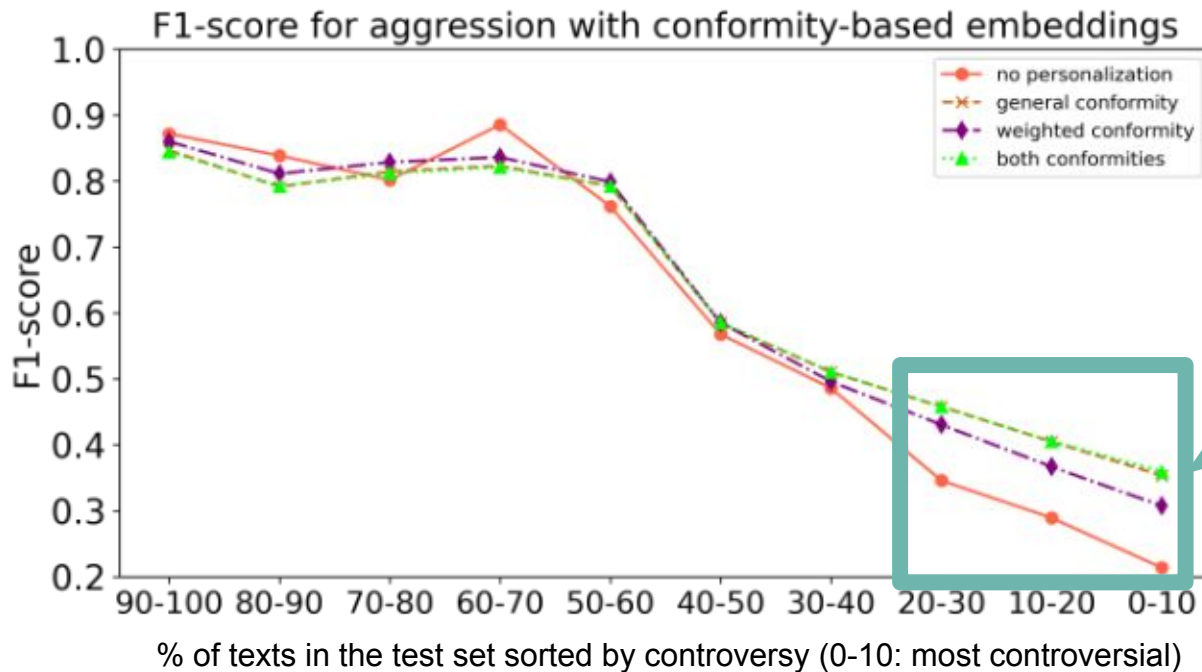
F1 for the *aggression* class only

EVALUATION RESULTS

Performance on aggression with most controversial scenario



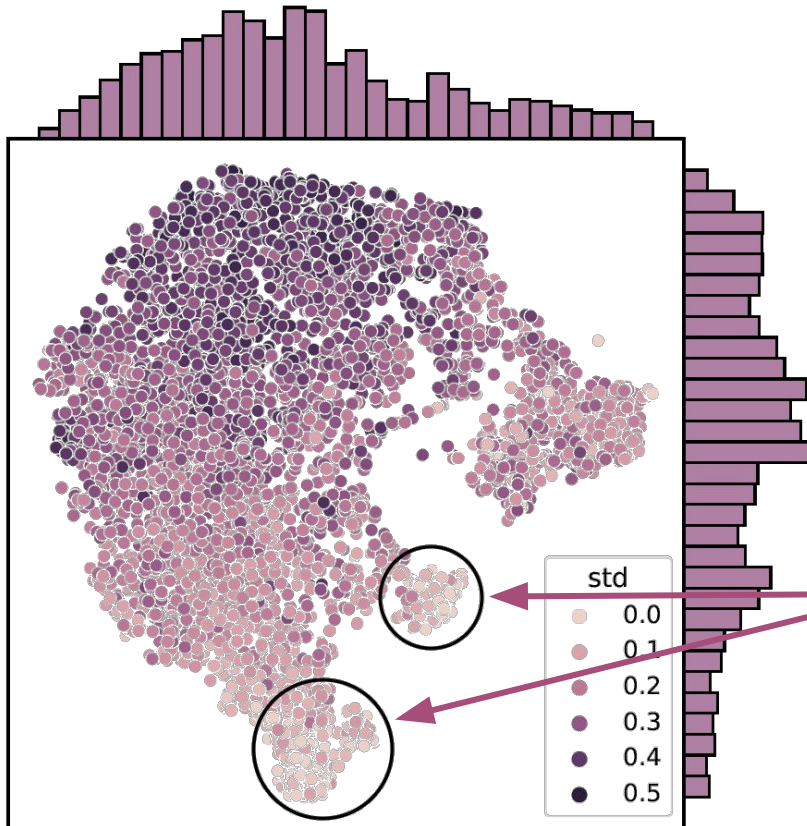
Where PNLP gains?



Crucial gain
for most
controversial
texts

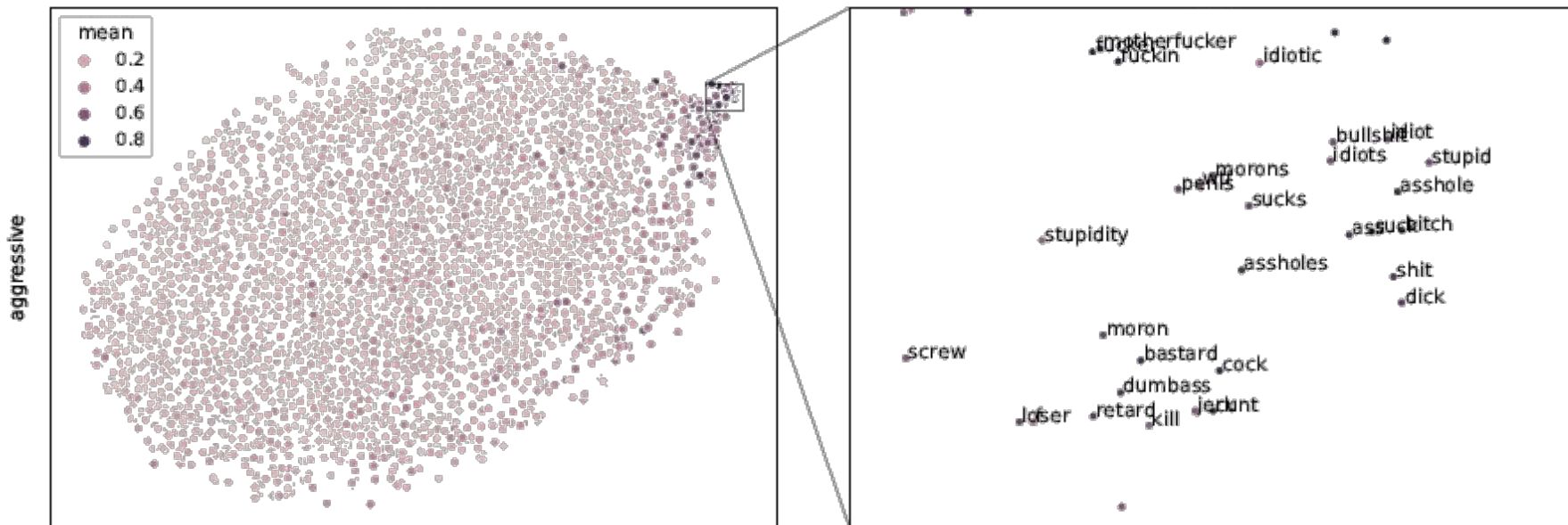


HUMAN EMBEDDINGS: Wiki Aggression

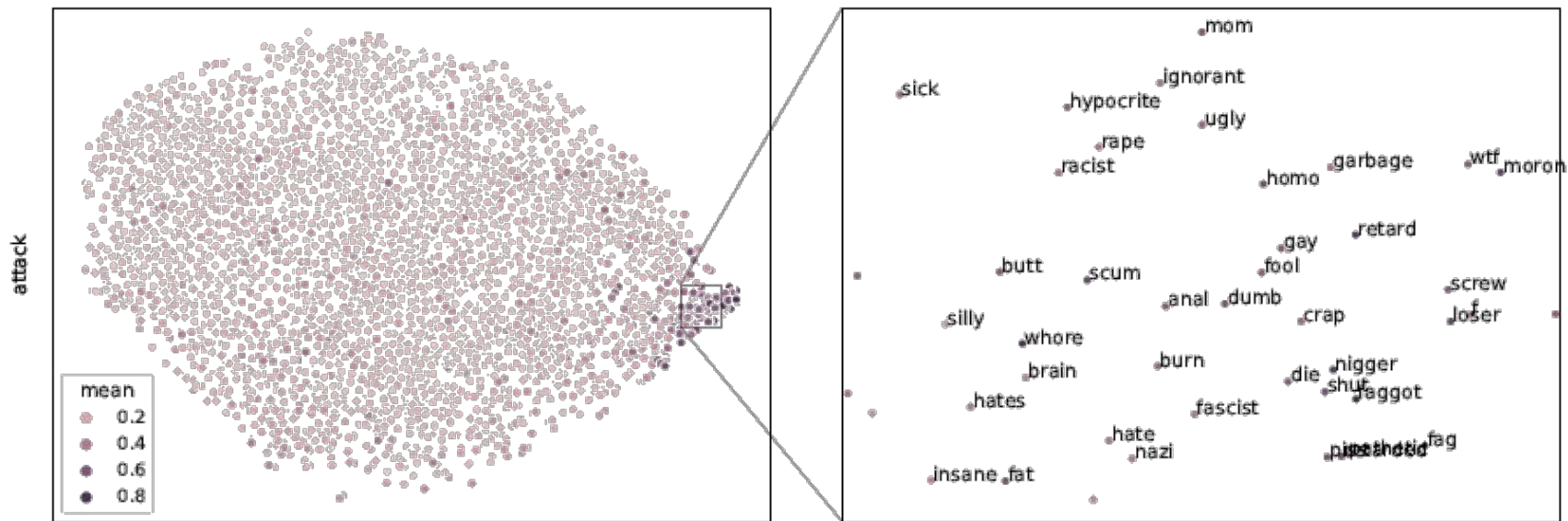


Low std. dev.
for some annotators
⇒ **not credible** ones?

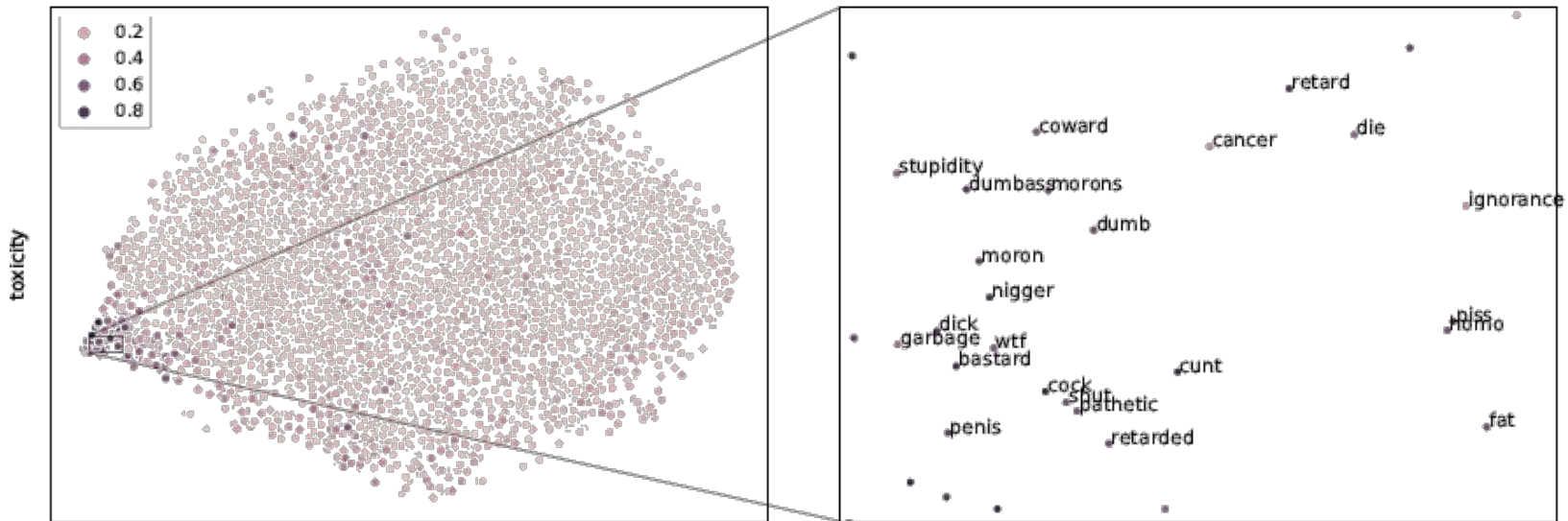
WORD EMBEDDINGS: Wiki Aggression

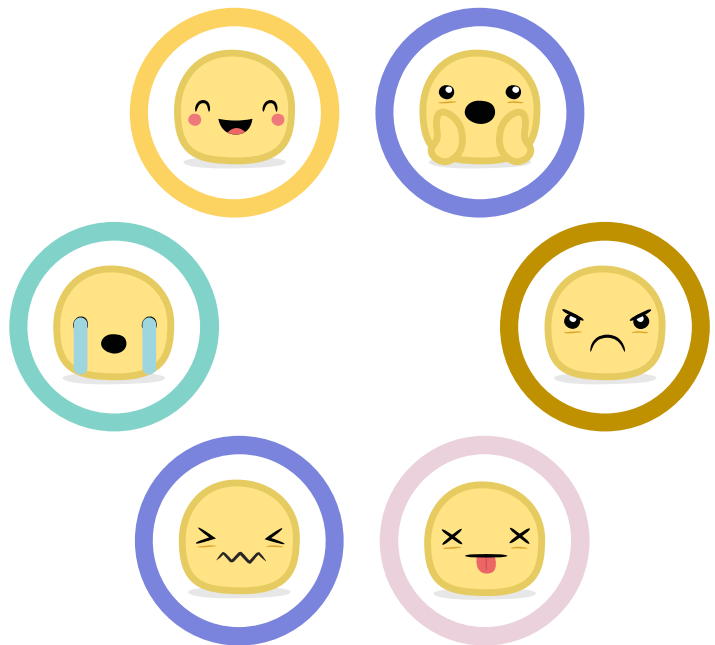


WORD EMBEDDINGS: Wiki Attack



WORD EMBEDDINGS: Wiki Toxicity



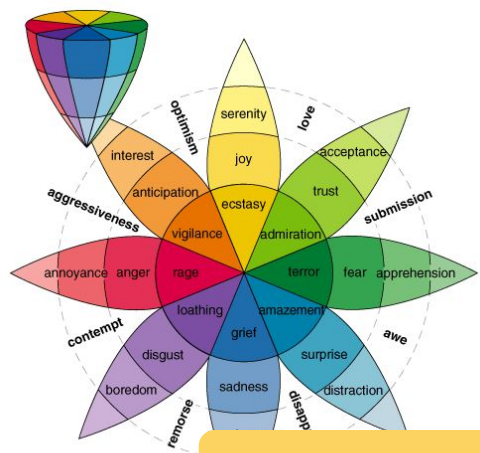


6

RESEARCH ON EMOTIONAL CONTENT PERCEPTION

ACL2021 – [Mit21]
ICDM2021 – [Koc21b]

EMOTIONAL DATA (in Polish)



Emotions

Texts

People

10 values

7,004

8,853

Annotations

Controversial Texts

3,774,338

NOT publicly available

100 %

EMOTIONAL TEXTS: example

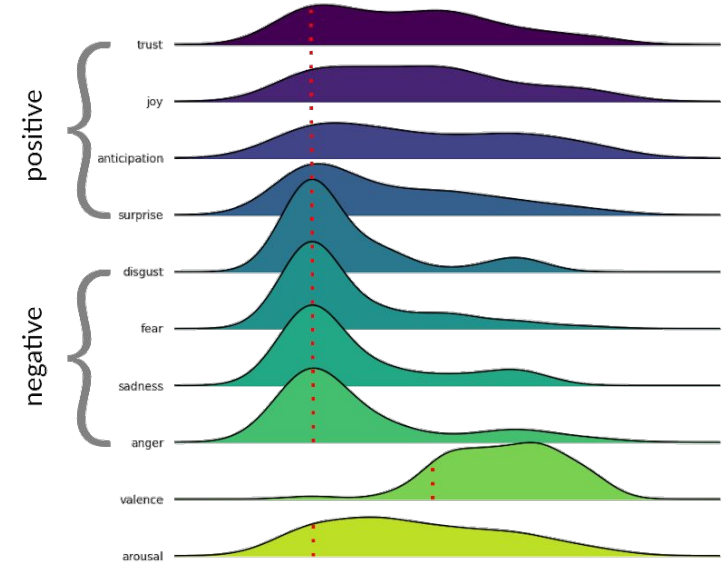
Example opinion

A modern, clean, well-maintained closed housing estate. Tastefully furnished apartments with full equipment. Great swimming pools, playground for children, exercise room - two treadmills and some other equipment, sauna. In fact, the car park is constantly full, we parked in front of the estate's gate. I do not recommend parking in prohibited places, because the security first sticker on the glass sticker, which is said to be hard to take off and then call the police. 10 minutes walk to the sea. Nearby a few places with home-made lunches, a little further on a grocery store. To the promenade on foot about half an hour.

Example annotation

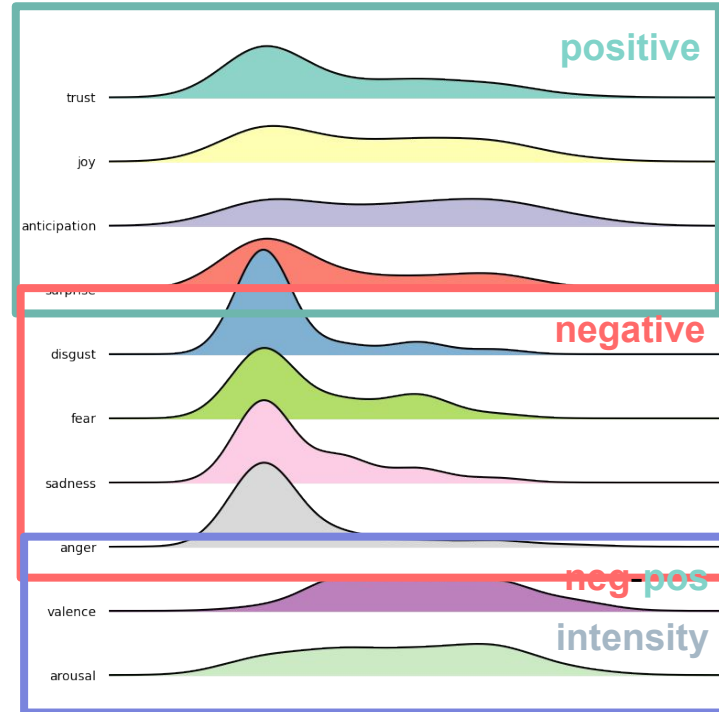


All annotations



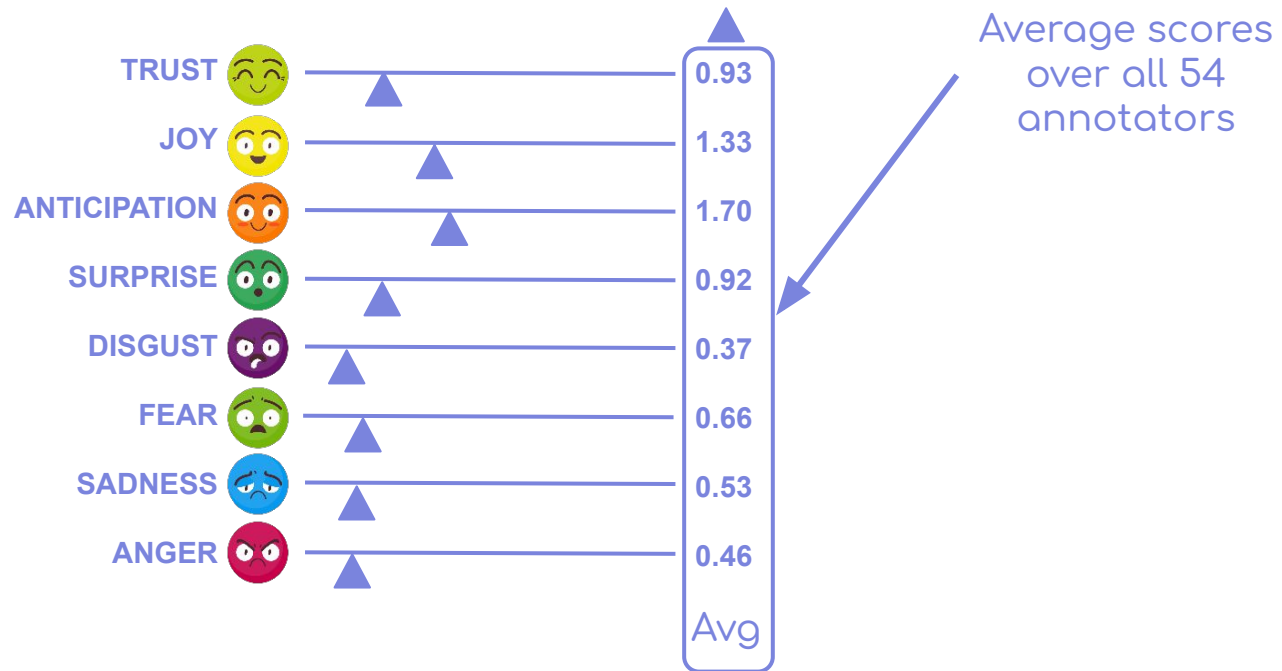
Example opinion

“She closed an unsuccessful chapter in her life and decided to start all over again.”



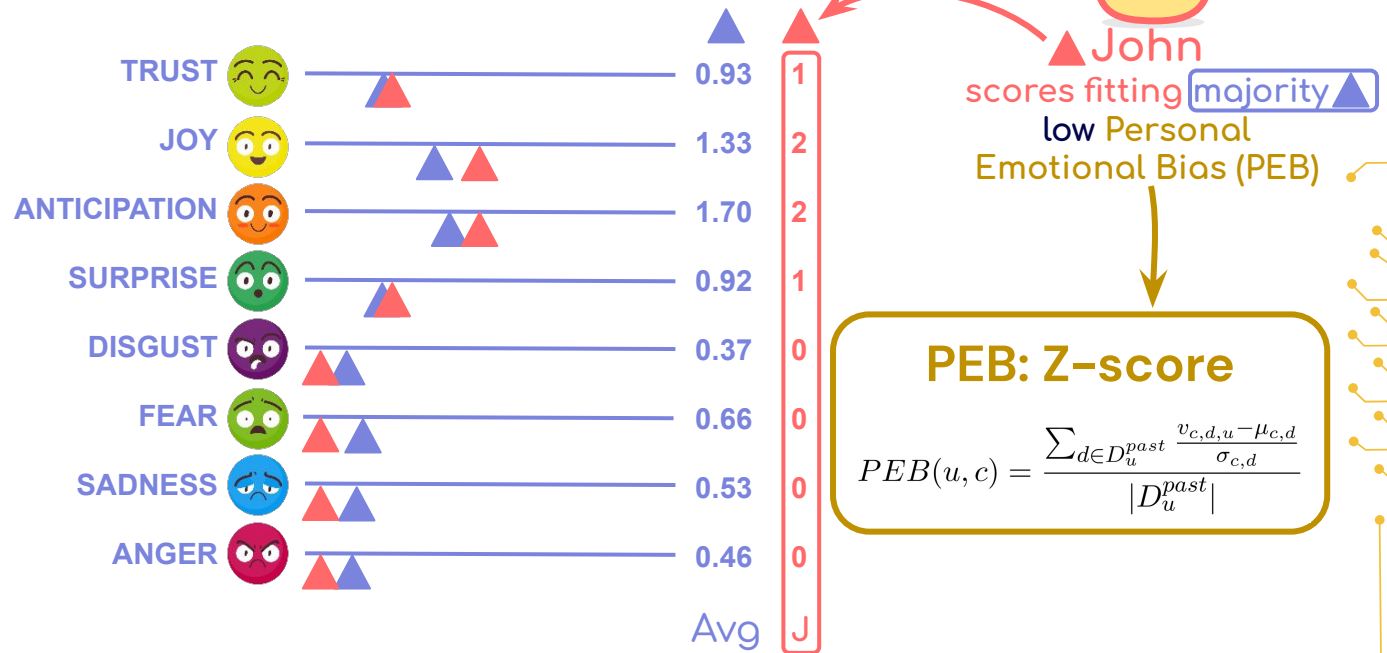
Different answers

“She closed an unsuccessful chapter in her life and decided to start all over again.”



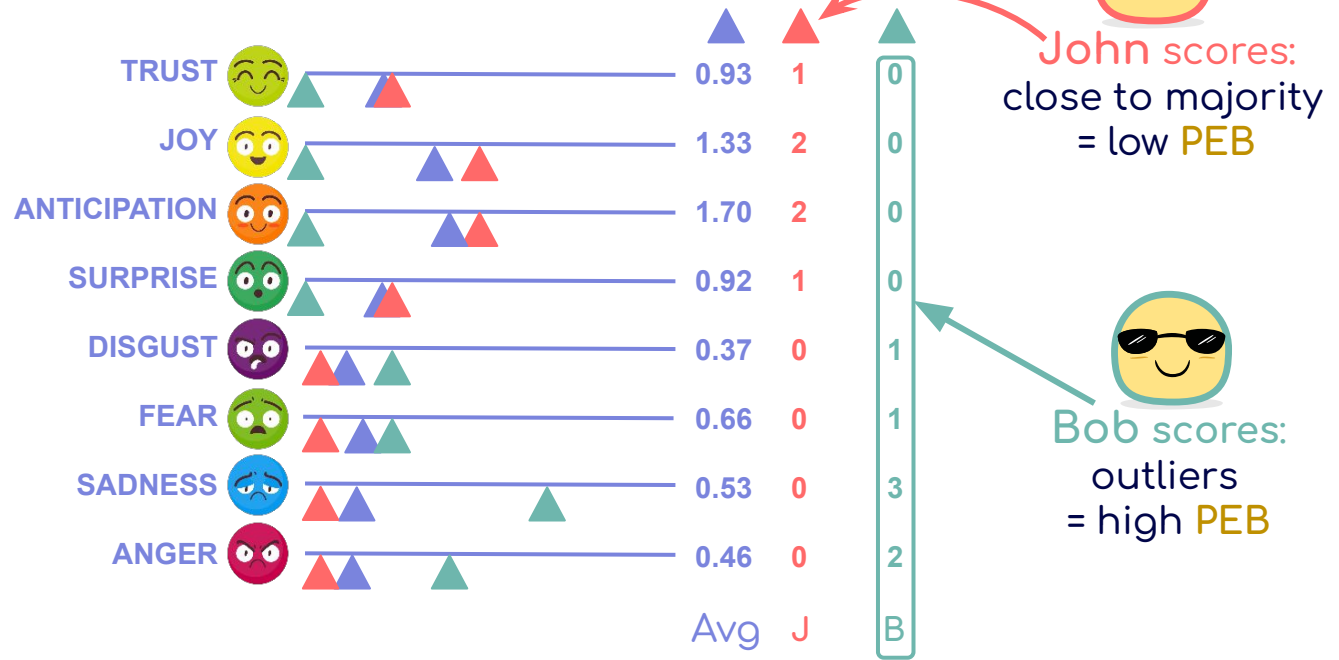
Different answers

“She closed an unsuccessful chapter in her life and decided to start all over again.”



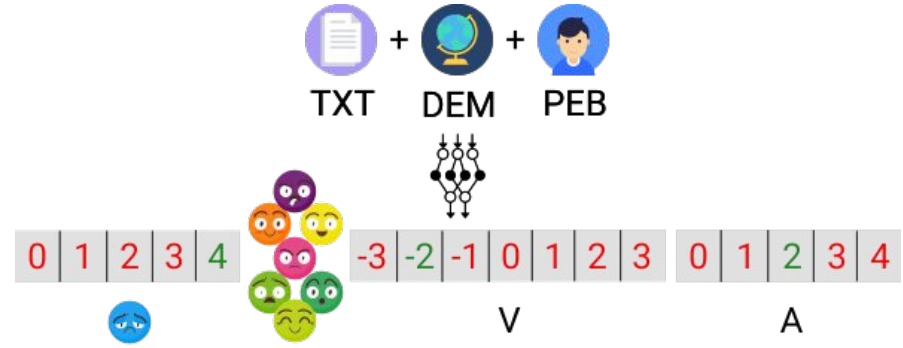
Different answers

“She closed an unsuccessful chapter in her life and decided to start all over again.”

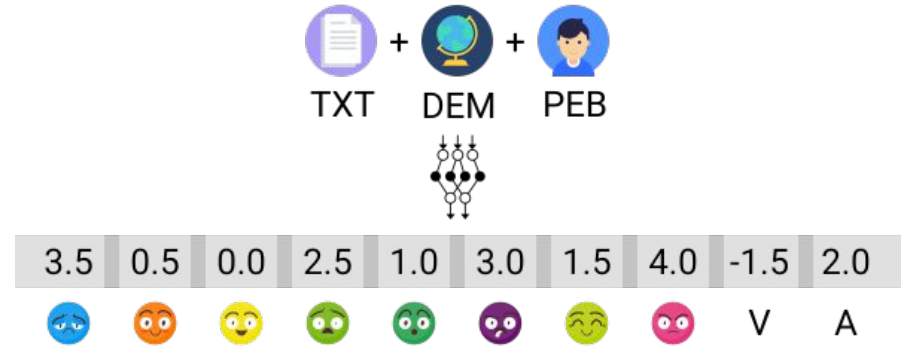


EMOTIONAL EXPERIMENTS

(1) Multi-task classification

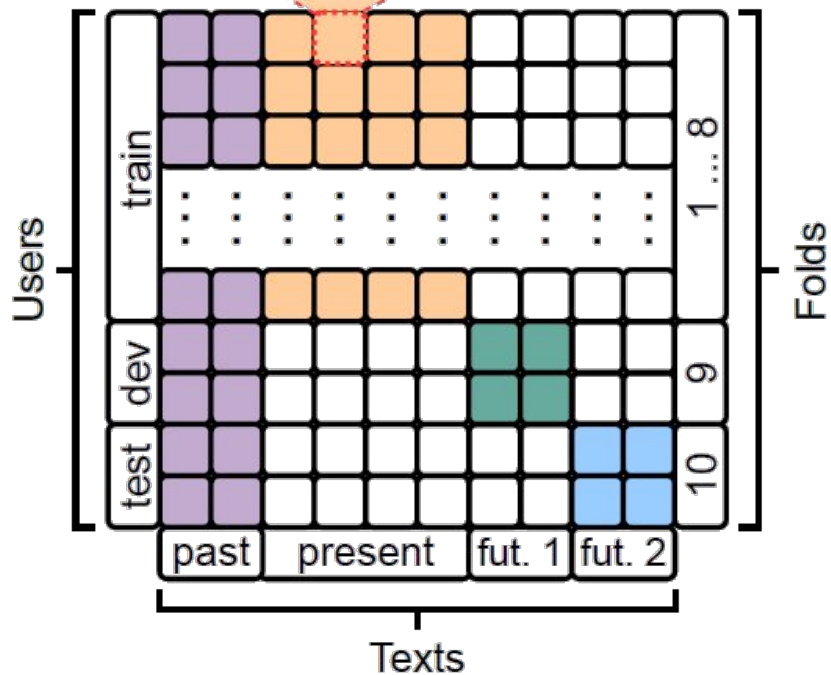


(2) Multivariate regression



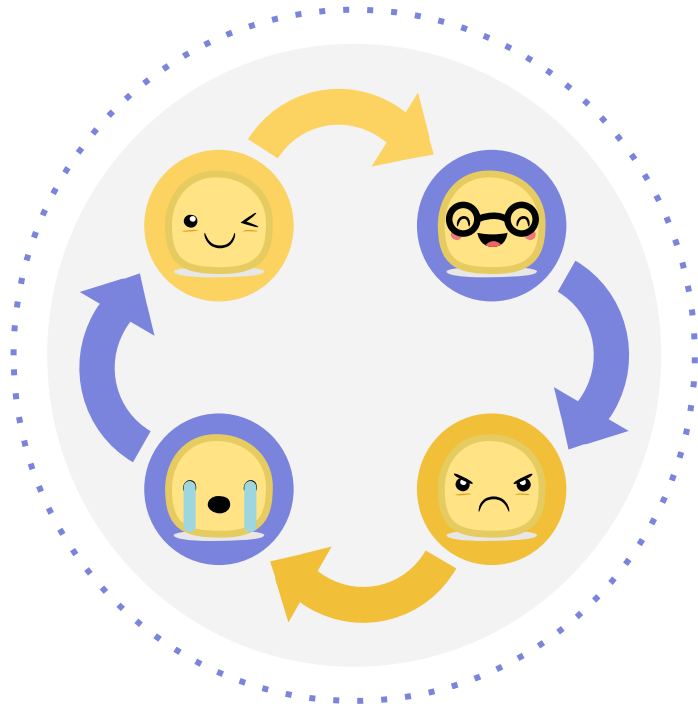
EMOTIONAL DATA SPLIT

Similar to offensive data but with 10 folds



PEB: Z-score

$$PEB(u, c) = \frac{\sum_{d \in D_u^{past}} \frac{v_{c,d,u} - \mu_{c,d}}{\sigma_{c,d}}}{|D_u^{past}|}$$



6a

**RESEARCH ON
EMOTIONS:
METHODS**

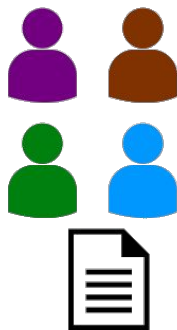
GENERALIZED vs. PERSONALIZED NLP



Generalized reasoning

Generalized rating


[2, 4, 3, 2, 3, 1, 2, 2, 1, 2]



Personalized reasoning

Personalized rating

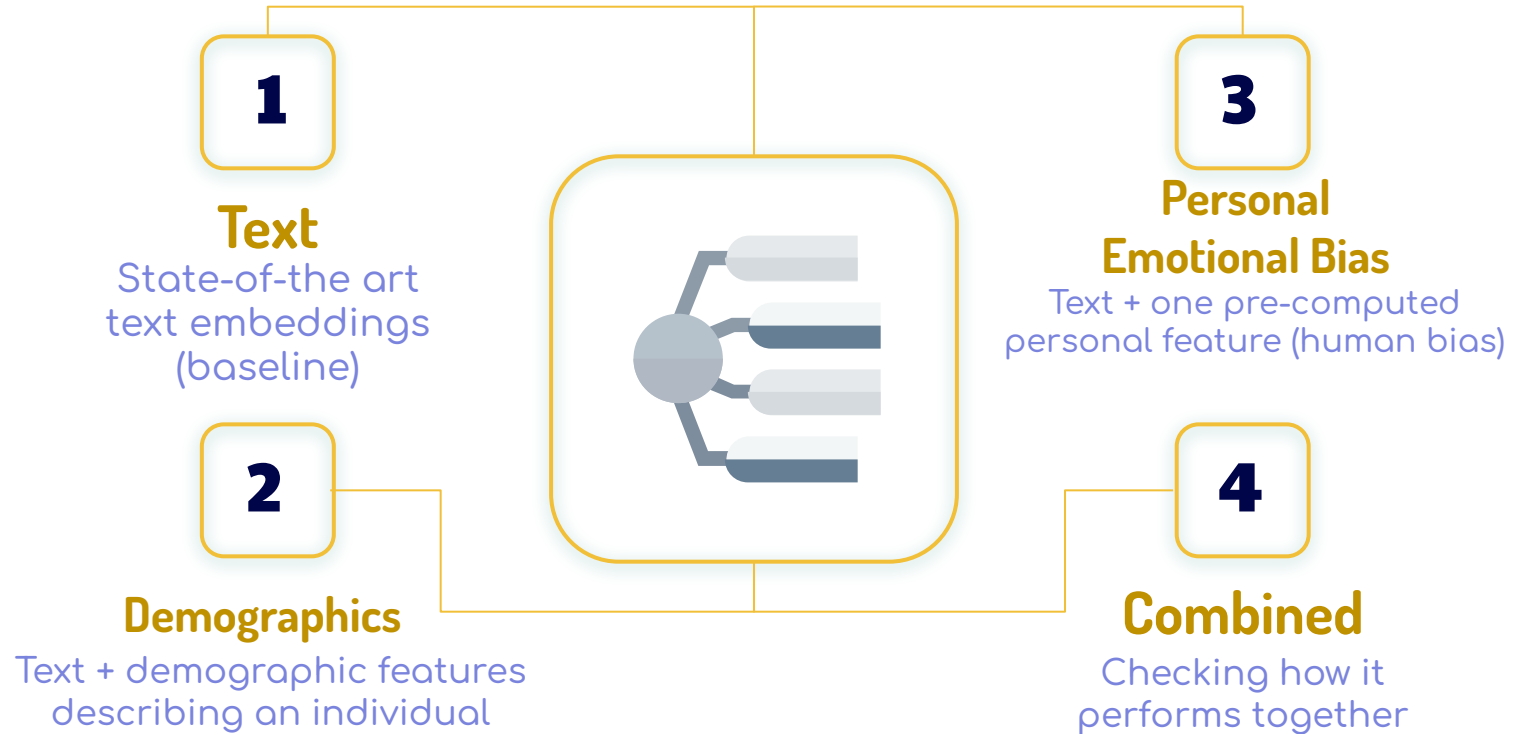

[0, 4, 2, 3, 4, 2, 3, 1, 0, 1]



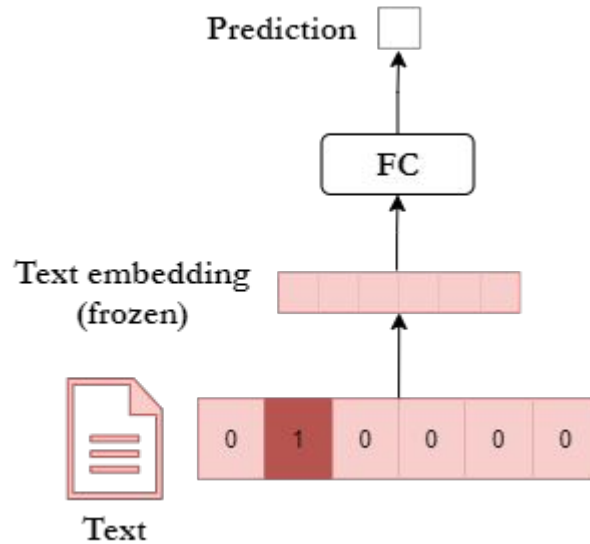

[4, 4, 4, 1, 2, 0, 1, 3, 2, 3]



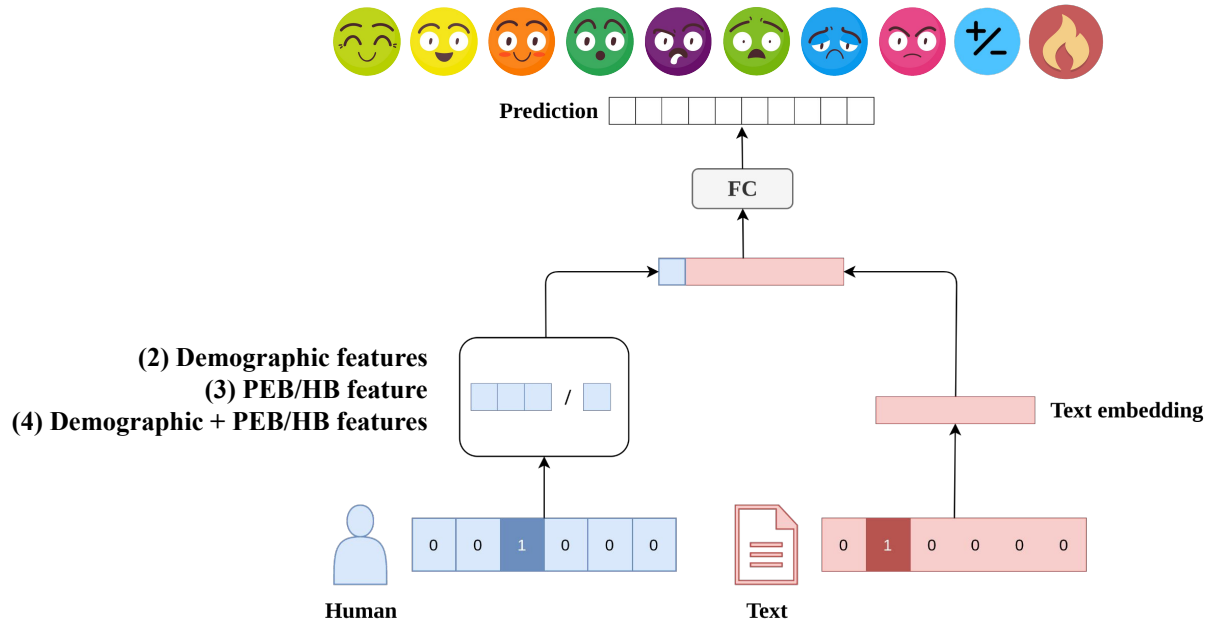
FOUR METHODS

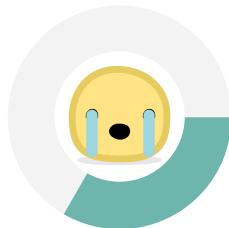
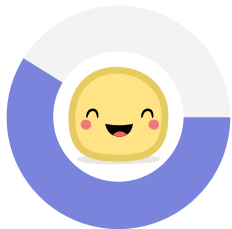


(1) TEXT ONLY: BASELINE



(2) DEMOGRAPHICS & (3) PERSONAL EMOTIONAL BIAS (PEB/HB) (4) ALL: demogr. + PEB feature





6b

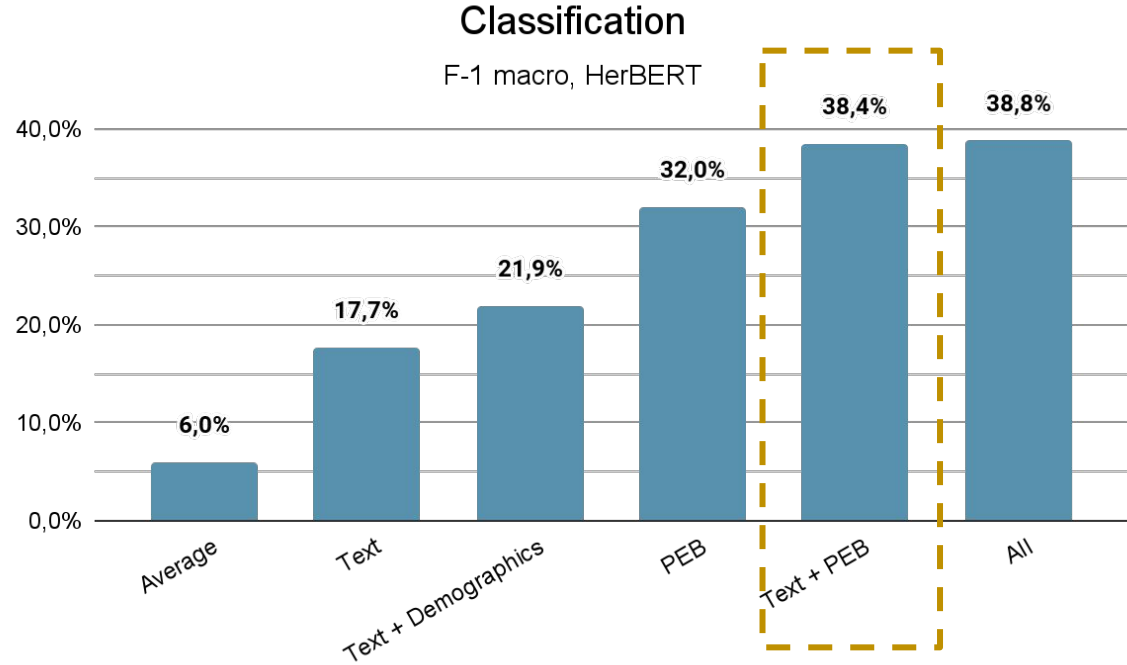
RESEARCH ON EMOTIONS: RESULTS

CLASSIFICATION: all emotions aggregated

Other language models:

- XLM-RoBERTa
- fastText + LSTM
- Polish RoBERTa

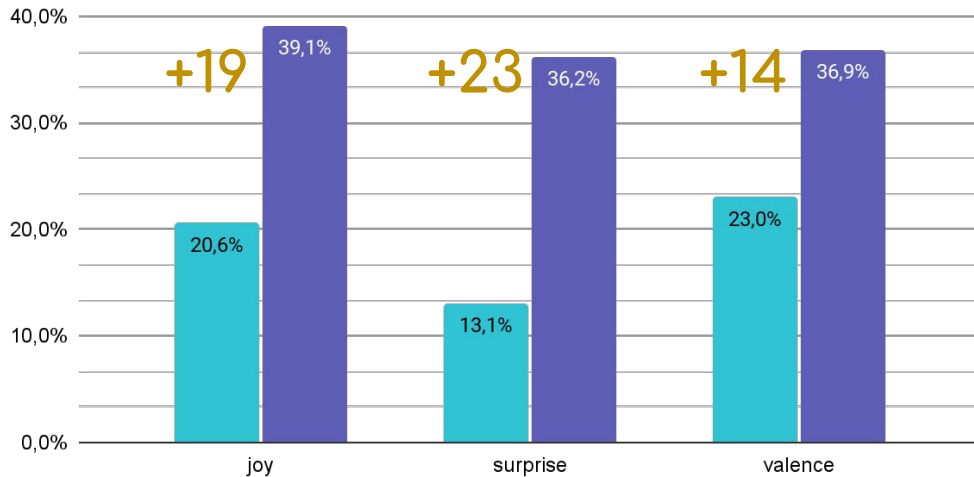
Worse by <1.5 p.p.



CLASSIFICATION: three emotional dimensions

Classification

F-1 macro, HerBERT (PL SOTA)



(1) Text only

Model based only on text embeddings



(3) Text and PEB

Model prepared on text embeddings and Personal Emotional Bias

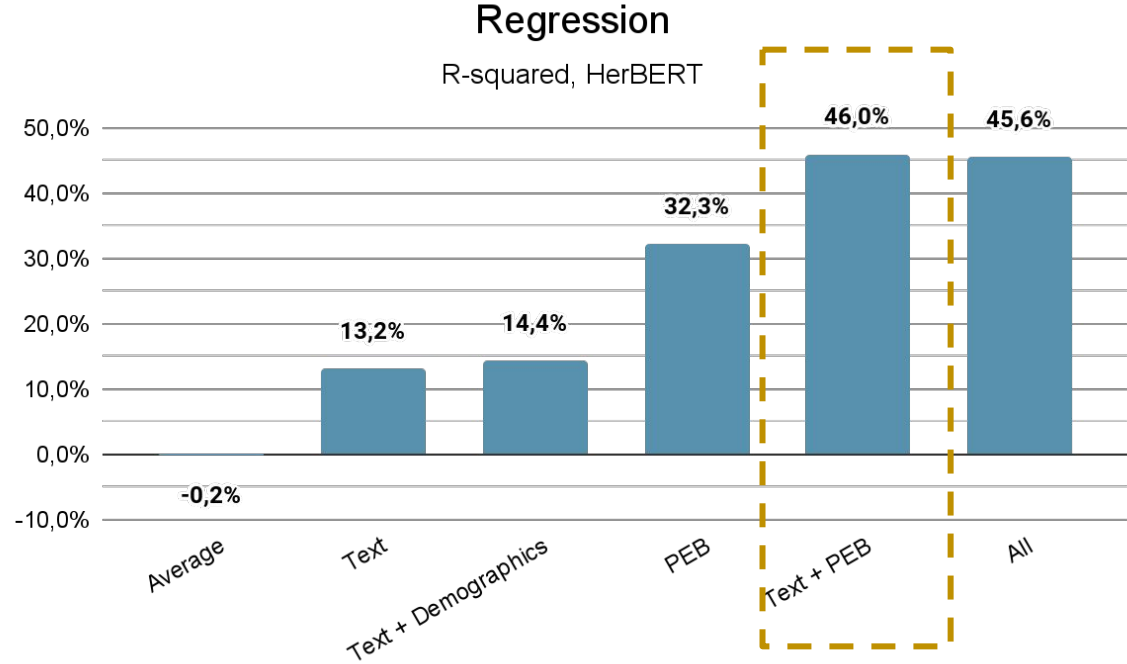


REGRESSION: all emotions aggregated

Other language models:

- XLM-RoBERTa
- fastText + LSTM
- Polish RoBERTa

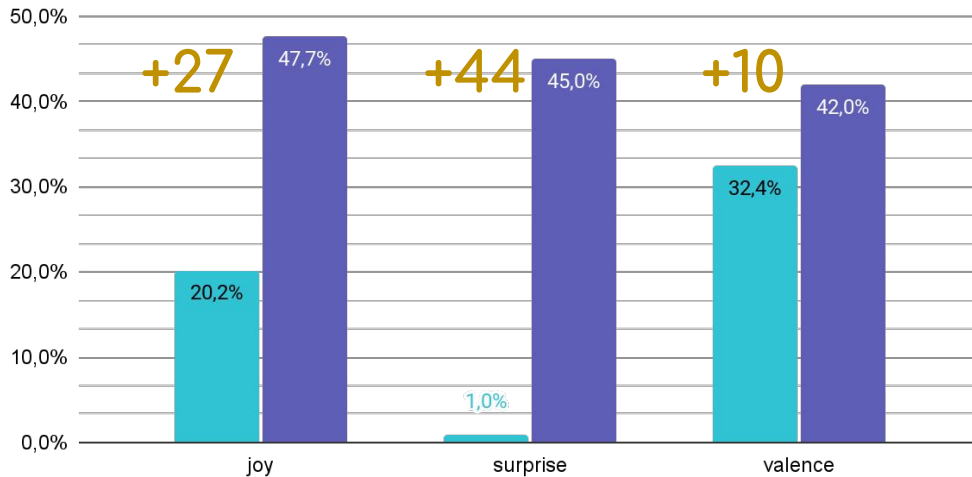
Worse by 3 p.p.



REGRESSION: three emotions

Regression

R-squared, HerBERT (PL SOTA)



(1) Text only

Model based only on text embeddings



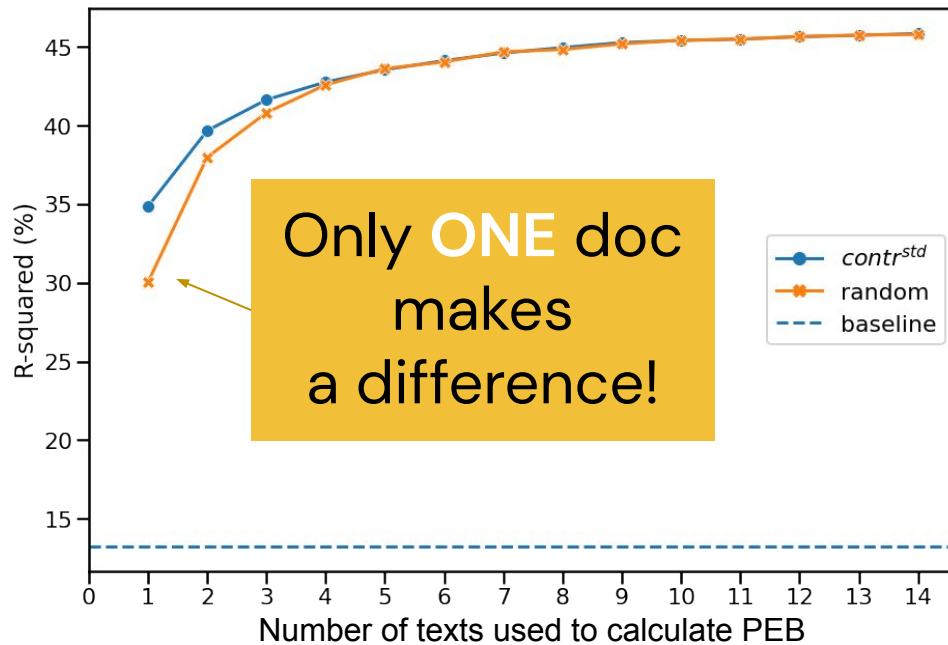
(3) Text and PEB

Model prepared on text embeddings and Personal Emotional Bias



How many texts are needed for PEB?

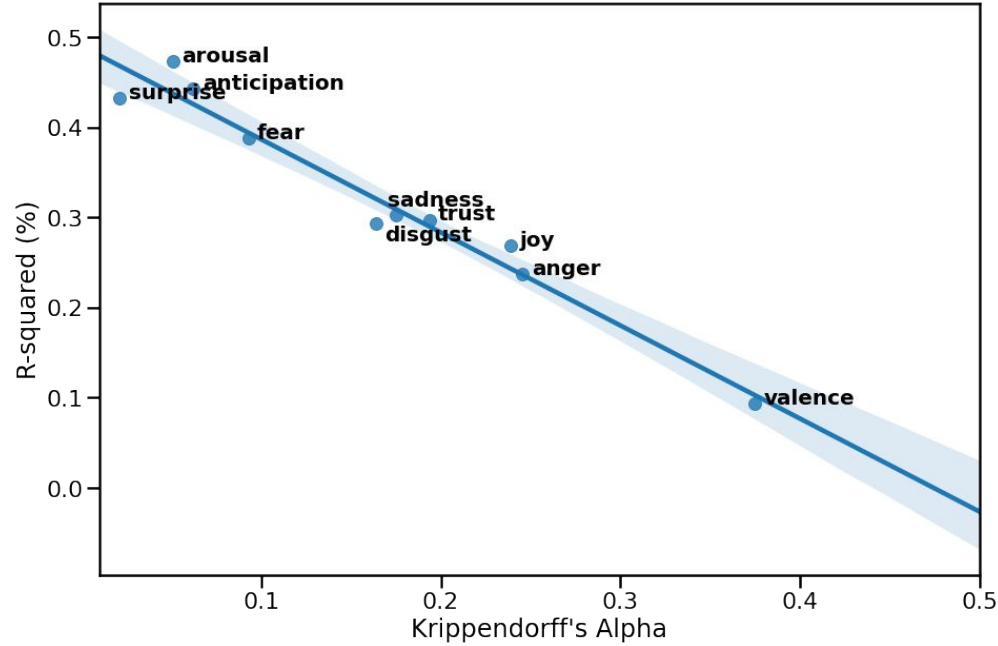
- (1) TXT - baseline
- (3) TXT+PEB:
 - random texts for PEB
 - most controversial texts for PEB



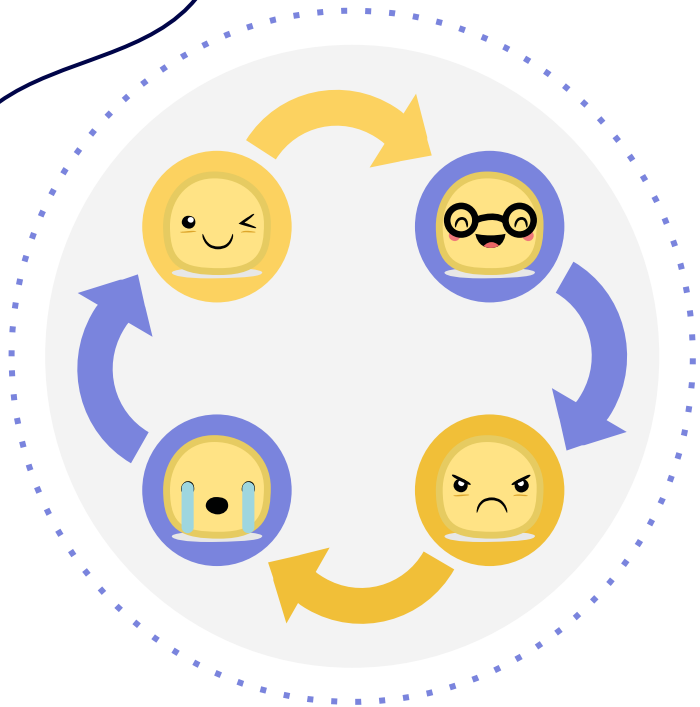
All emotions, HerBERT

AGREEMENT LEVEL (controversy) vs. performance

PEB-only model
Performance



Controversy in the collection



7

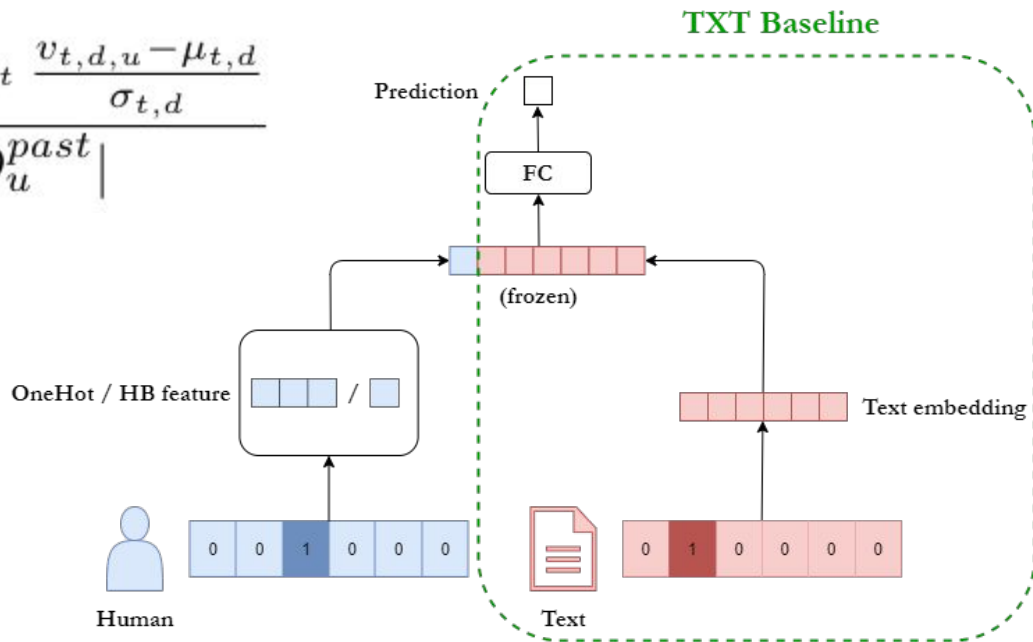
RESEARCH ON MULTIPLE TASKS AND MODELS

Wiki Detox: Attack,
Aggression, Toxicity
+ Emotions
ICDM2021: [Koc21b]

MODELS:

Baseline (TXT) & OneHot ID & HuBi-Formula

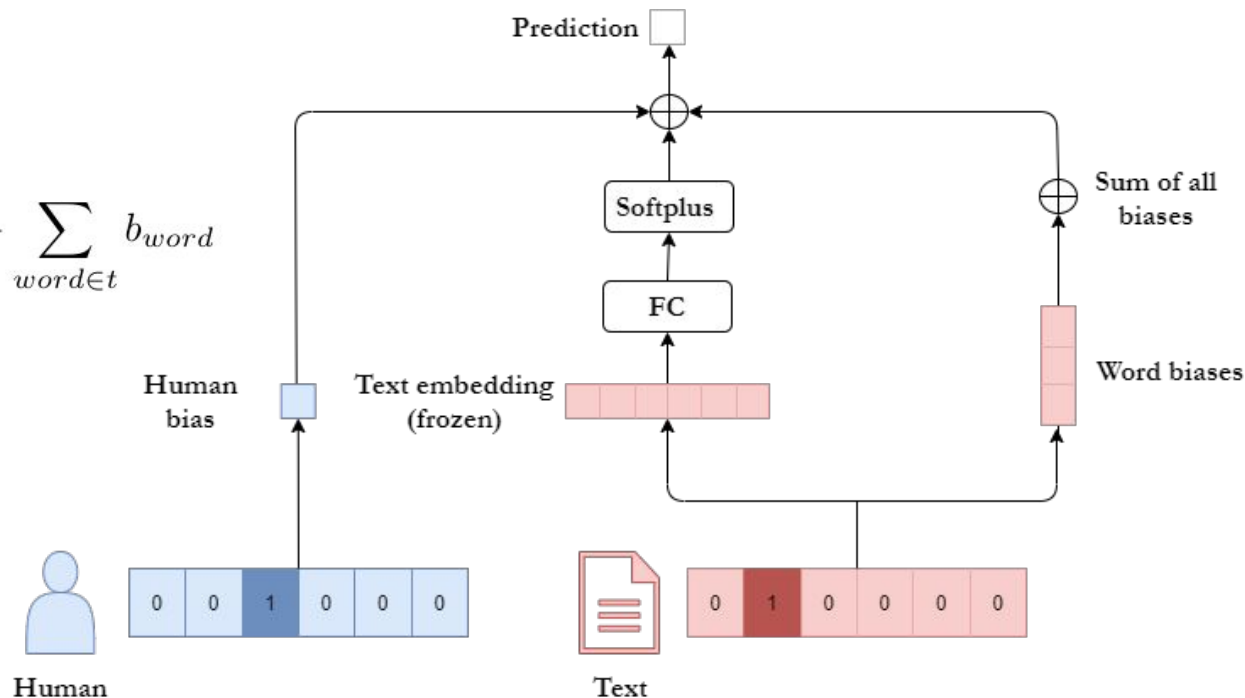
$$HB(u, t) = \frac{\sum_{d \in D_u^{past}} \frac{v_{t,d,u} - \mu_{t,d}}{\sigma_{t,d}}}{|D_u^{past}|}$$



MODELS:

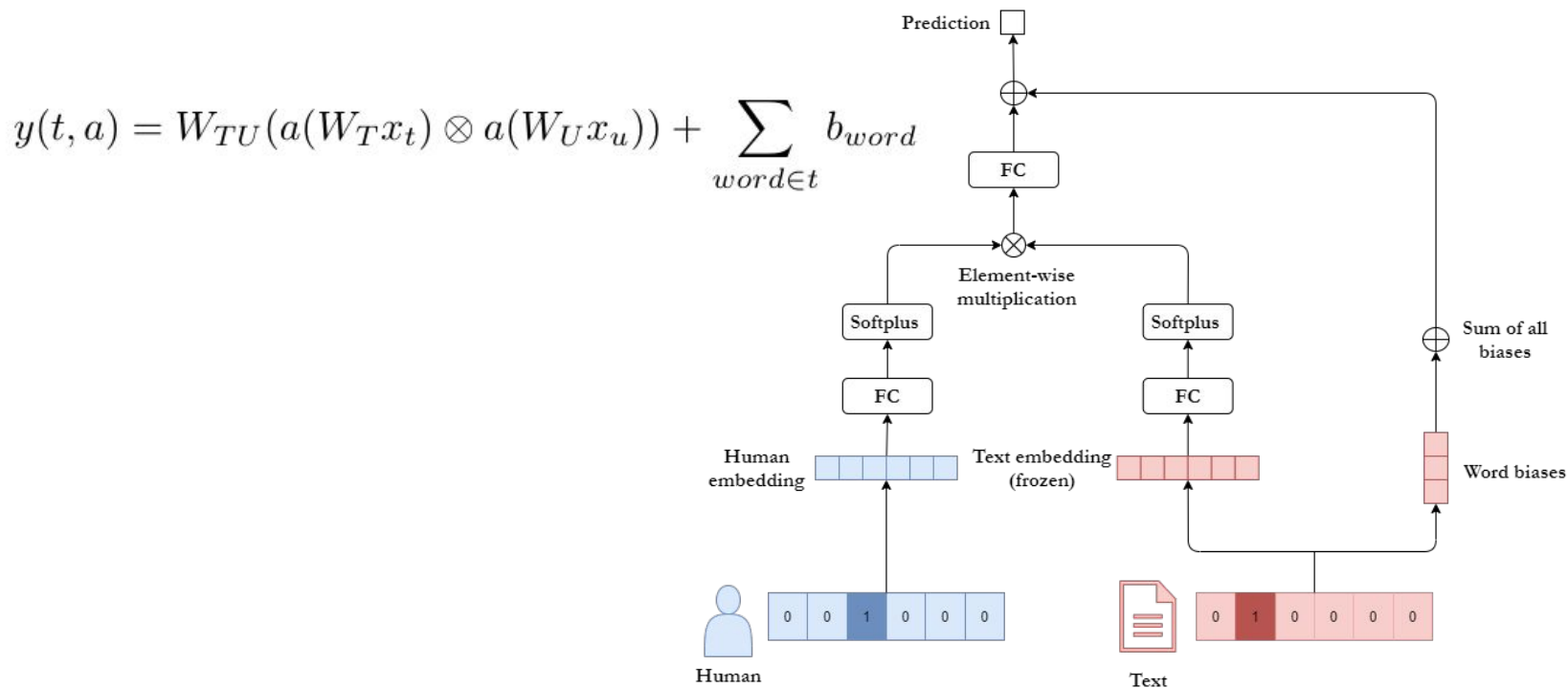
HuBi-Simple: learned human bias

$$y(t, u) = a(W_T x_t) + b_u + \sum_{word \in t} b_{word}$$



MODELS:

HuBi-Medium: learned human embedding

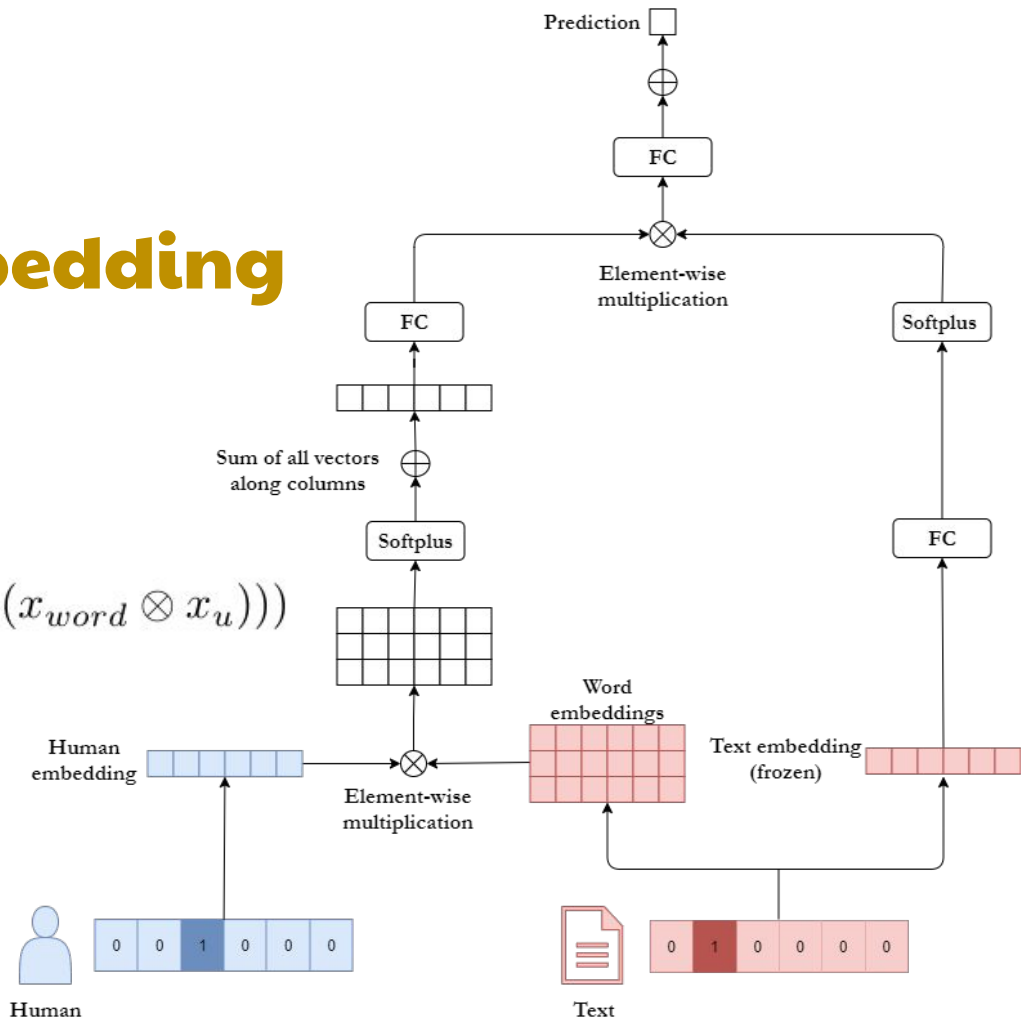


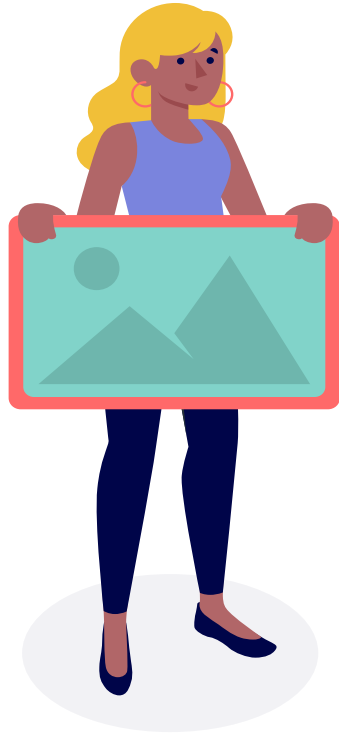
MODELS:

HuBi-Complex:

human-word embedding

$$y(t, a) = W(a(W_T x_t) \otimes W_{WU}(\sum_{word \in t} a(x_{word} \otimes x_u)))$$





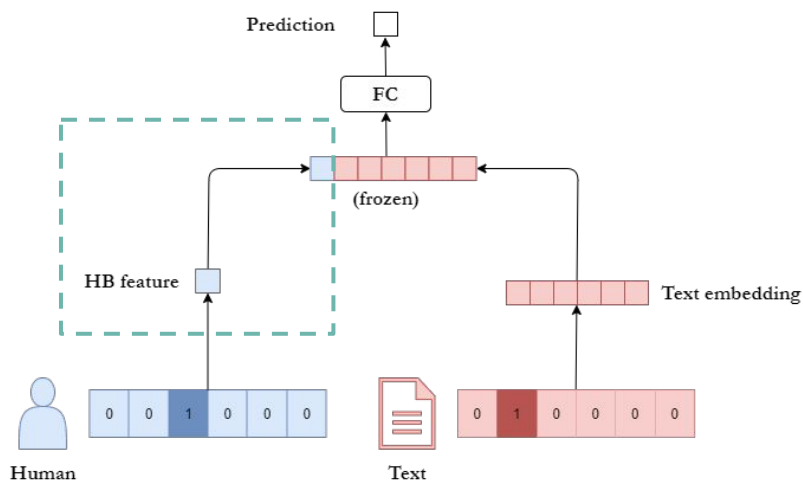
7a

MULTIPLE TASKS: RESULTS

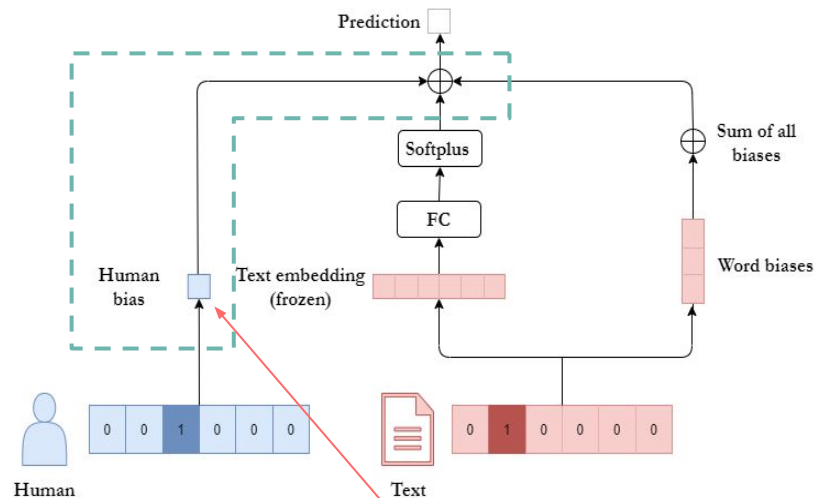
Wiki Detox + Emotions

FORMULA vs. LEARNED BIAS

HB feature vs. HuBi-Simple (learned bias)



VS.



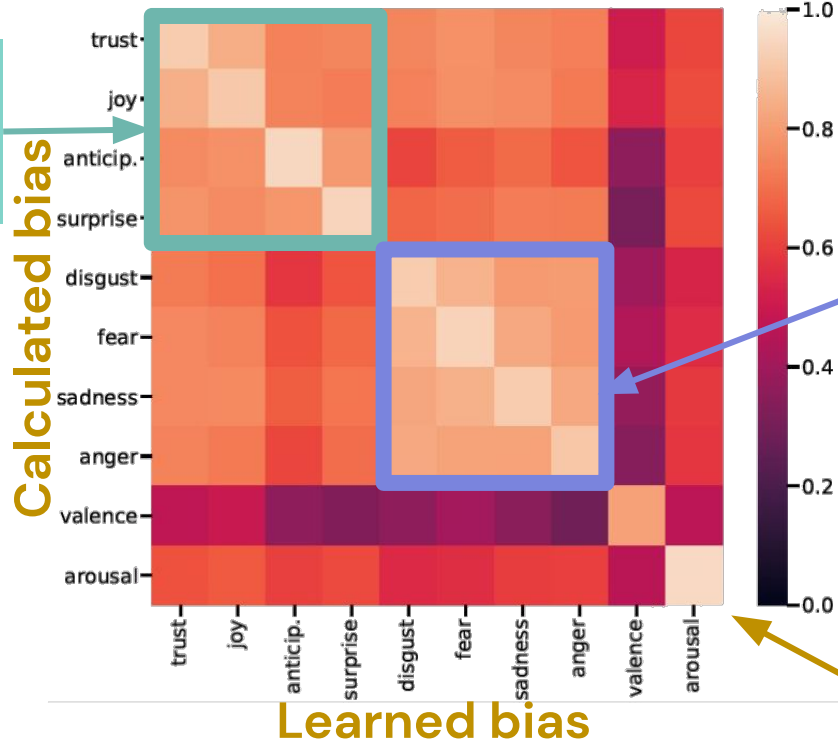
HB calculated feature (**formula**)

HuBi-Simple: **learned** human bias

FORMULA vs. LEARNED BIAS

Correlation between biases

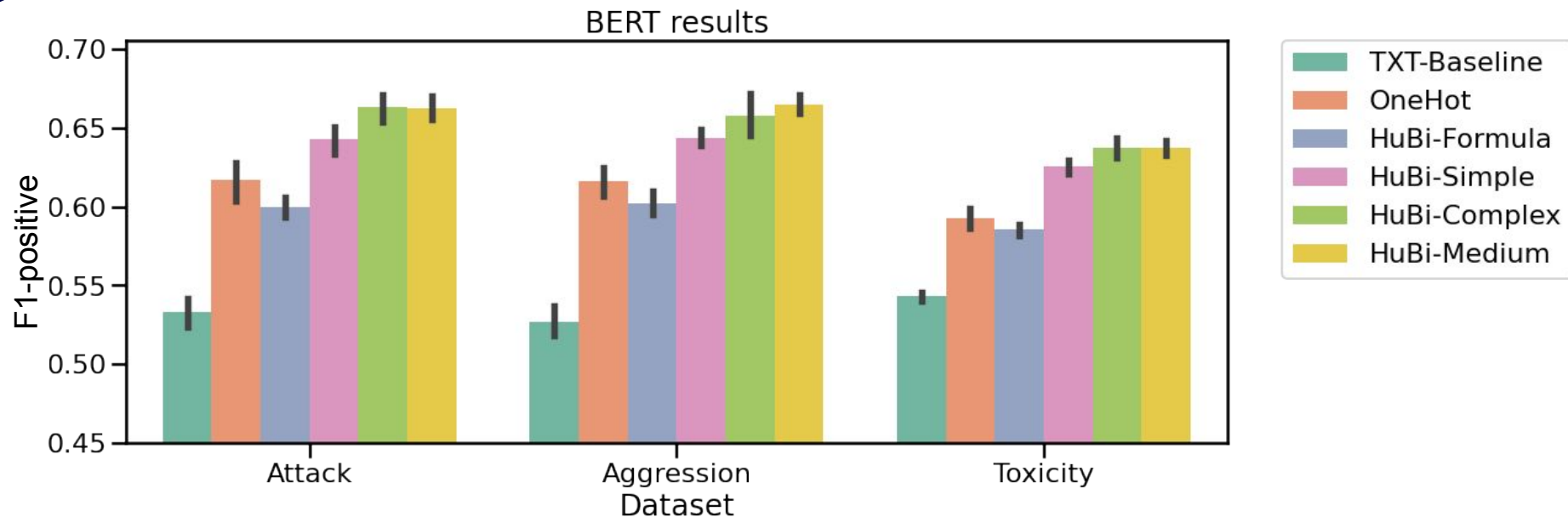
Positive emotions are highly correlated 73% and more



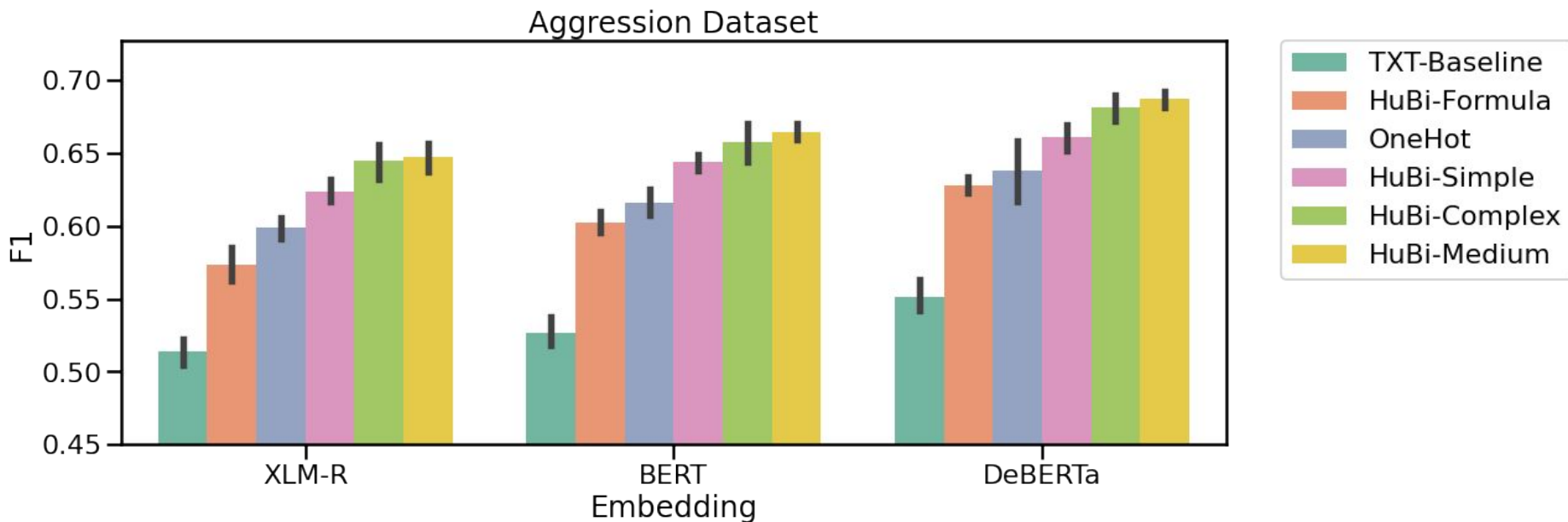
Negative emotions are highly correlated 80% and more

Biases are very highly correlated 90% and more (diagonal)

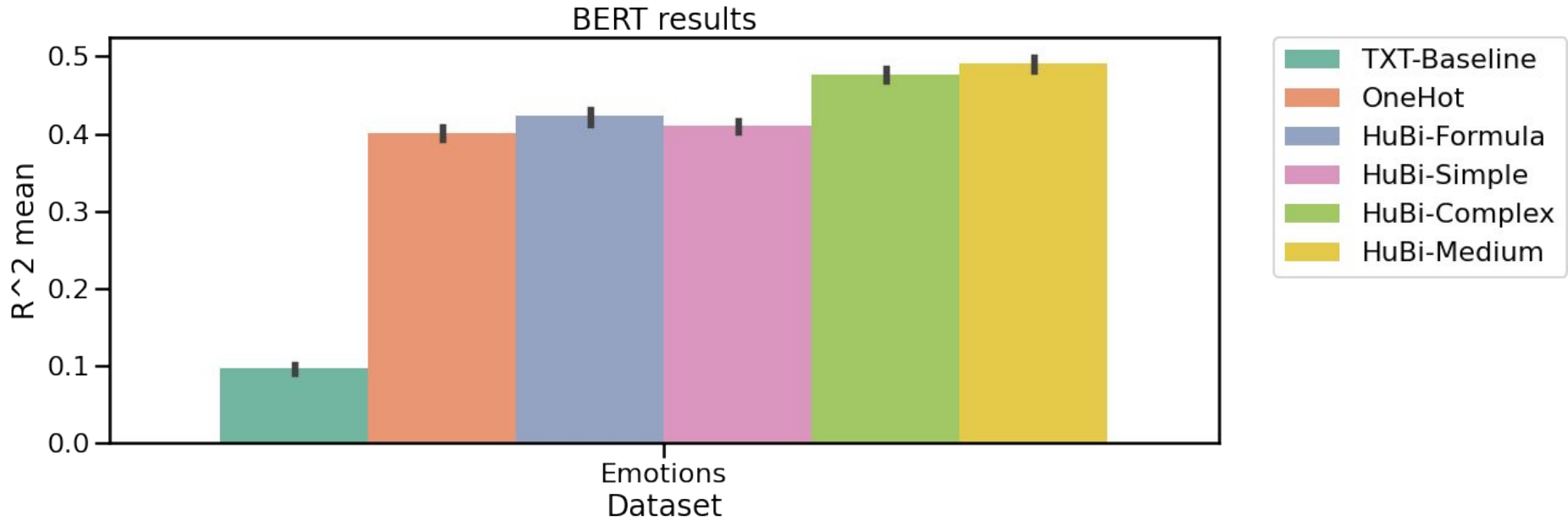
WIKI: results on three datasets



WIKI: Results on Aggression Data



EMOTIONS: Results

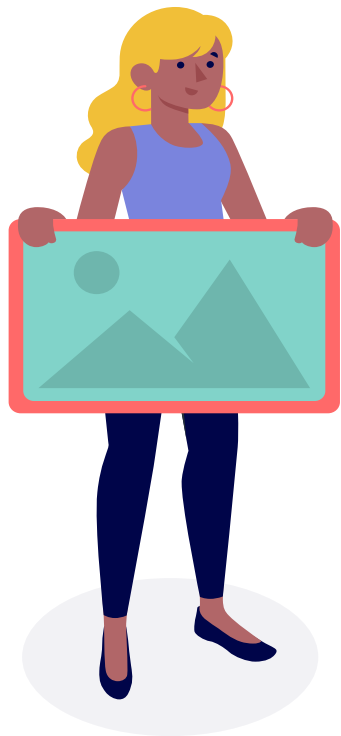


Multivariate regression

EMOTIONS: Results



Multivariate regression

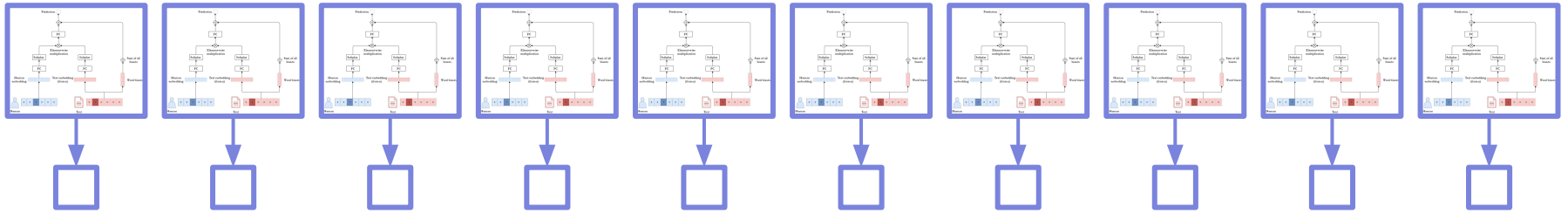


7b

MULTI-TASK vs. SINGLE TASK studies on emotions

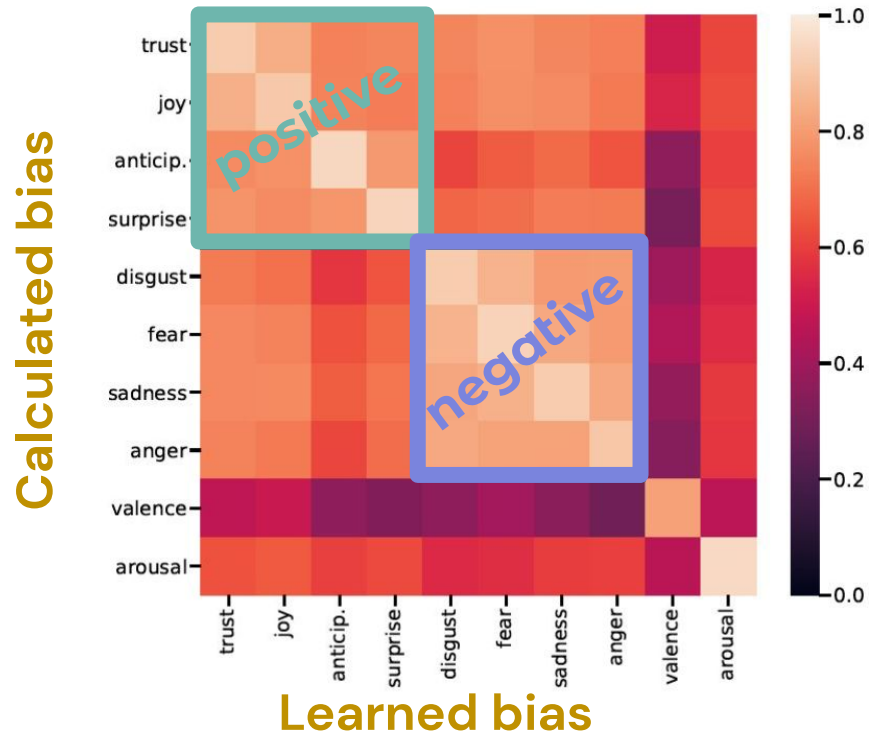
EmotionAware at PerCom
[Mit22]

Single task

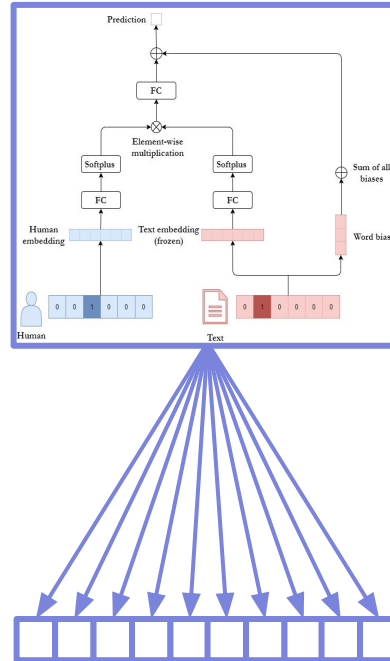


Each emotional category learned separately

Sub-multi-task: motivation - correlation



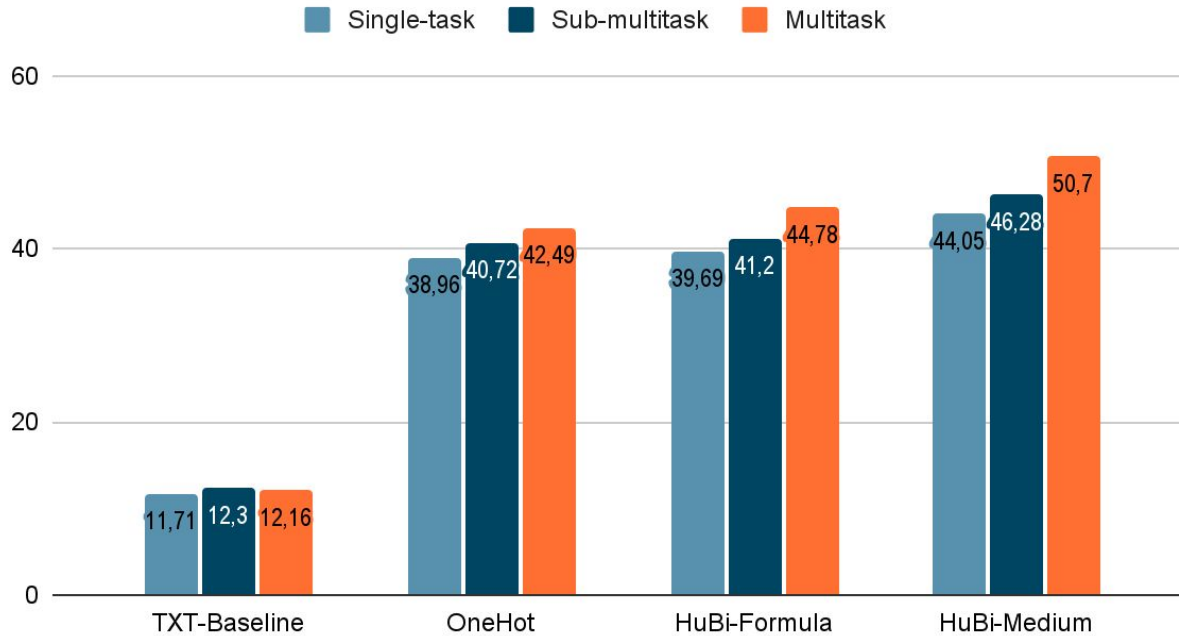
Multi-task



All emotional category learned together

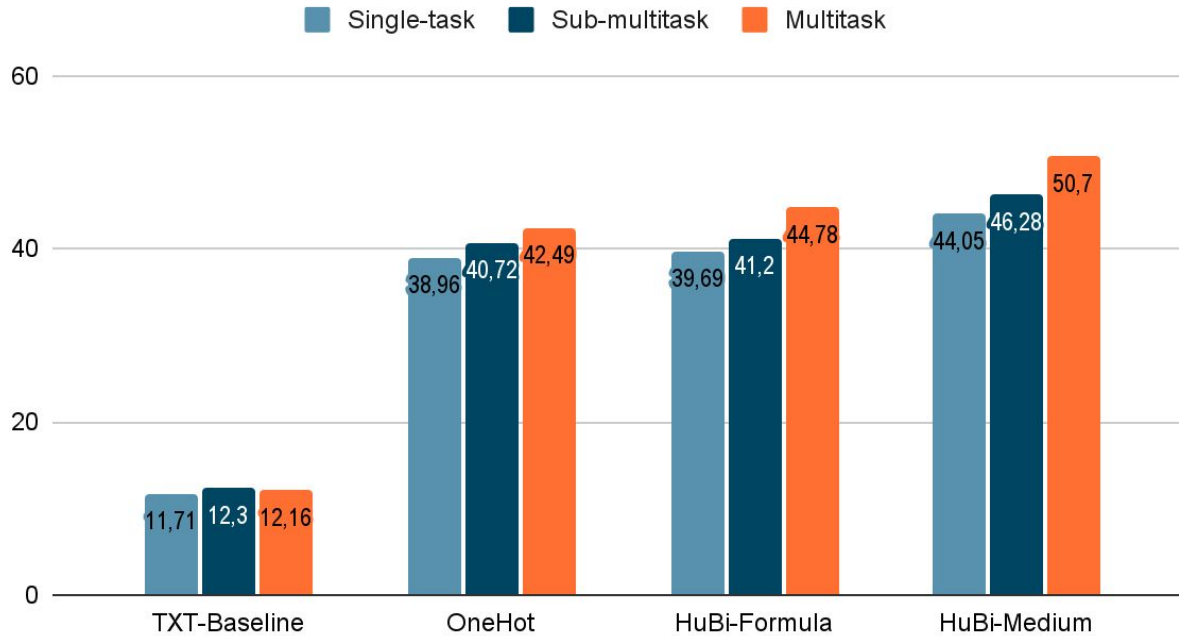
Multi-task: results

Regression results in R-squared for emotion prediction



Multi-task: results

Regression results in R-squared for emotion prediction



+15%

Multi-task improvement
over single task



8

CONCLUSIONS



CONCLUSIONS #1



PNLP vs. GNLN

Personalized methods **ALWAYS** perform better than the generalized ones



Diversity

Conformity, Controversy and **Human Bias** deliver vital information about the user



PNLP vs. language

Each PNLN method gains **much more than** language models



Few docs is enough

Even four docs provide user information that improves reasoning (5-6 docs for emotional texts)



CONCLUSIONS #2



Validation

Train/dev/test split should be based on **users** instead of texts



Application

Our PNLN methods can be applied to **any subjective** task



Demographics

Demographic data only slightly improves reasoning



Data

Human-centered annotations are crucial for personalised NLP

TEAM



Przemysław Kazienko



Jan Kocoń



Kamil Kanclerz



Julita Bielaniewicz



Marcin Gruza



Piotr Miłkowski

BIBLIOGRAPHY

- [Koc21a] Kocoń J., Figas A., Gruza M., Puchalska D., Kajdanowicz T., Kazienko P.: *Offensive, aggressive, and hate speech analysis: from data-centric to human-centred approach*. **Information Processing and Management**, 58(5) 2021, art. 102643.
- [Kan21] Kanclerz K., Figas A., Gruza M., Kajdanowicz T., Kocoń J., Puchalska D., Kazienko P.: *Controversy and Conformity: from Generalized to Personalized Aggressiveness Detection*. **ACL 2021**, 5915–5926.
- [Mił21] Miłkowski P., Gruza M., Kanclerz K., Kazienko P., Grimling D., Kocoń J.: *Personal Bias in Prediction of Emotions Elicited by Textual Opinions*. **ACL 2021**, Student Research Workshop, 248–259.
- [Koc21b] Kocoń J., Gruza M., Bielaniewicz J., Grimling D., Kanclerz K., Miłkowski P., Kazienko P.: *Learning Personal Human Biases and Representations for Subjective Tasks in Natural Language Processing*, IEEE **ICDM'21** 2021, IEEE, 1168–1173.
- [Mił22] Miłkowski P., Saganowski S., Gruza M., Kazienko P., Piasecki M., Kocoń J.: *Multitask Personalized Recognition of Emotions Evoked by Textual Content*. **EmotionAware'22 at PerCom'22**, IEEE, 2022, 347–352.

Other two papers in reviews, SRW at ACL 2022 (**active learning** and **humor recognition** for PNL), one in AfCAL.

OTHER RESEARCH

Sentiment analysis

1. Kocoń J., Baran J., Gruza M., Janz A., Kajstura M., Kazienko P., Korczyński W., Miłkowski P., Piasecki M., Szołomicka J.: *Neuro-symbolic Models for Sentiment Analysis*. **ICCS 2022**.
2. Miłkowski P., Gruza M., Kazienko P., Szołomicka J., Woźniak S., Kocoń J.: *MultiEmo: Language-agnostic Sentiment Analysis*. **ICCS 2022**.
3. Augustyniak Ł., Kajdanowicz T., Kazienko P.: *Comprehensive analysis of aspect term extraction methods using various text embeddings*. **Computer Speech & Language**, Vol. 69, 2021, 101217

OTHER RESEARCH

Emotion recognition from physiological signals

1. Saganowski S., et al.: *The cold start problem and per-group personalization in real-life emotion recognition with wearables*. **WristSense 2022 at PerCom 2022**, IEEE, 2022, 812–817. **BEST PAPER AWARD**.
2. Saganowski S., et al.: *Emognition dataset: emotion recognition with self-reports, facial expressions, and physiology using wearables*. **Scientific Data**, 9, 158 (2022).
3. Saganowski S., Perz B., Polak A., Kazienko P.: *Emotion Recognition for Everyday Life Using Physiological Signals from Wearables: A Systematic Literature Review*. **IEEE Transactions on Affective Computing**, in reviews, 2022.
4. Saganowski S., et al.: *A system for collecting emotionally annotated physiological signals in daily life using wearables*. **ACII 2021**.
5. Dzieżyc M., et al.: *How to catch them all? Enhanced data collection for emotion recognition in the field*. **PerCom 2021**, IEEE, 2021, 348–351.
6. Saganowski S., et al.: *Consumer Wearables and Affective Computing for Wellbeing Support*. **MobiQuitous 2020**, 482–487.
7. Saganowski S., et al.: *Emotion Recognition Using Wearables: A Systematic Literature Review – Work-in-progress*. **EmotionAware 2020 at PerCom 2020**, IEEE, 2020, 1–6.
8. Dzieżyc M., Gjoreski M., Kazienko P., Saganowski S., Gams M.: *Can we ditch feature engineering? End-to-End Deep Learning for Affect Recognition from Physiological Sensor Data*. **Sensors**, 2020, 20(22), 6535.

Take-home message

***Personalized NLP
is much better than
generalized for all
subjective tasks***





Thank you for your attention!

Q & A