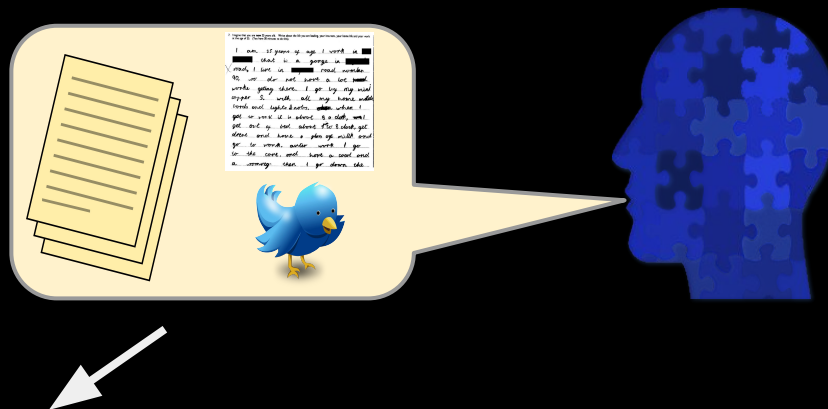


# Human-Centered Natural Language Processing

CSE 354

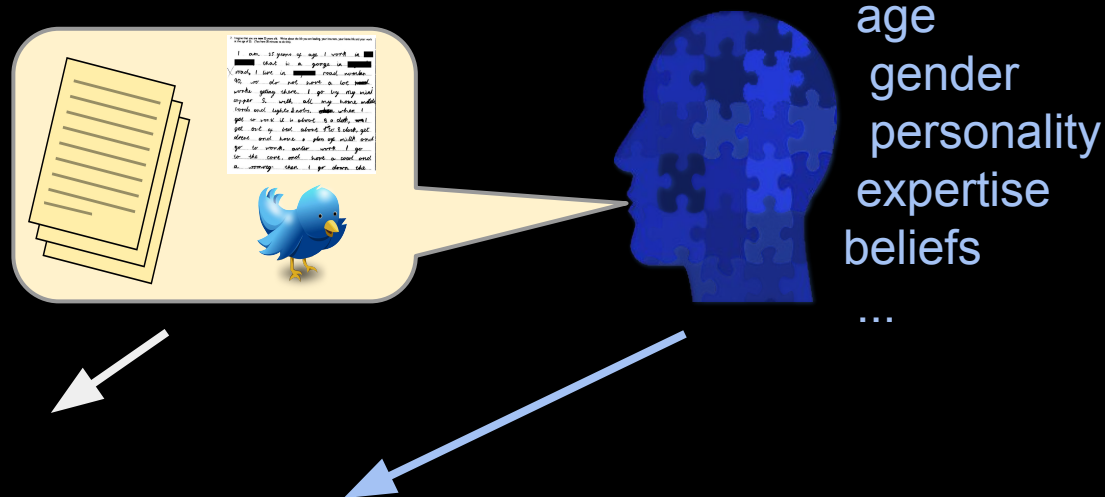
# The “Task” of human-centered NLP



Most NLP Tasks. E.g.

- POS Tagging
  - Document Classification
  - Sentiment Analysis
  - Stance Detection
  - Mental Health Risk Assessment
  - ...
- (language modeling, QA, ...)

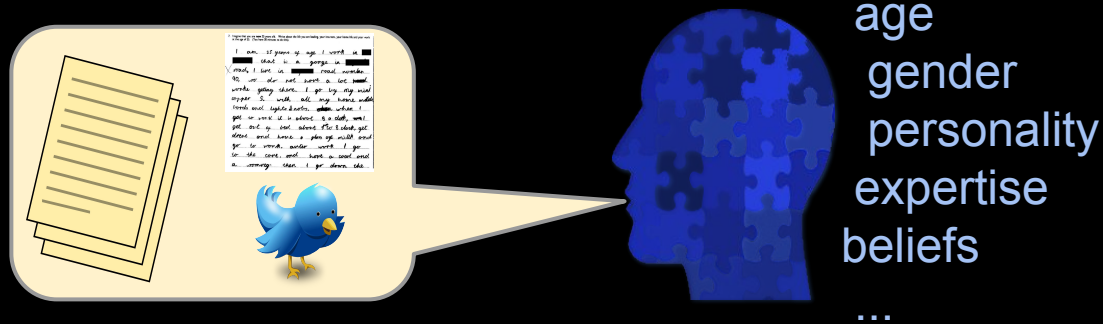
# The “Task” of human-centered NLP



## Most NLP Tasks. E.g.

- POS Tagging
  - Document Classification
  - Sentiment Analysis
  - Stance Detection
  - Mental Health Risk Assessment
  - ...
- (language modeling, QA, ...)

# The “Task” of human-centered NLP



## Most NLP Tasks. E.g.

- POS Tagging
  - Document Classification
  - Sentiment Analysis
  - Stance Detection
  - Mental Health Risk Assessment
  - ...
- (language modeling, QA, ...)

## How to include extra-linguistics?

- Additive Inclusion
- Adaptive Extralinguistics
  - Adapting Embeddings
  - Adapting Models
- Correcting for bias



**Natural  
Language  
Processing**



**Psychological  
& Health  
Sciences**





Natural  
Language  
Processing

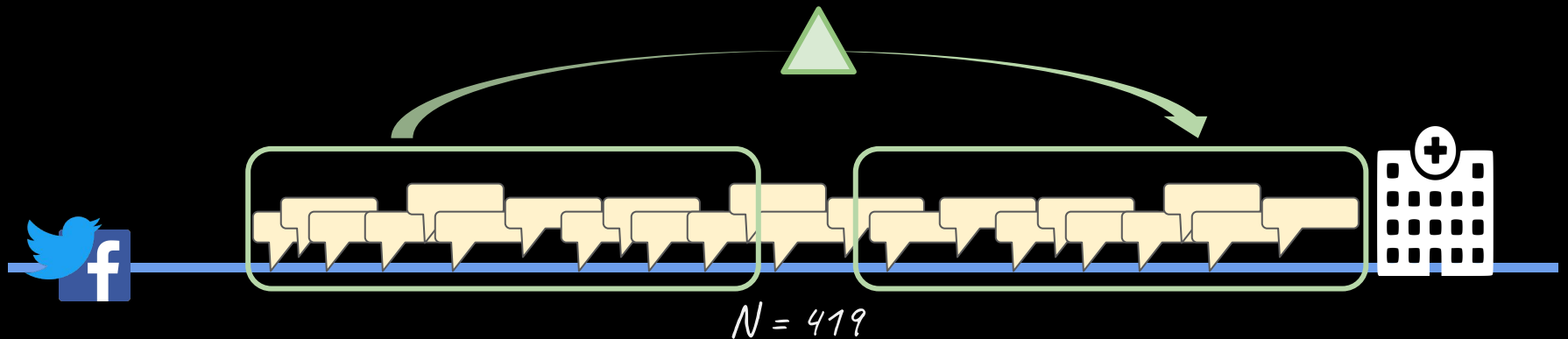
Psychological  
& Health  
Sciences





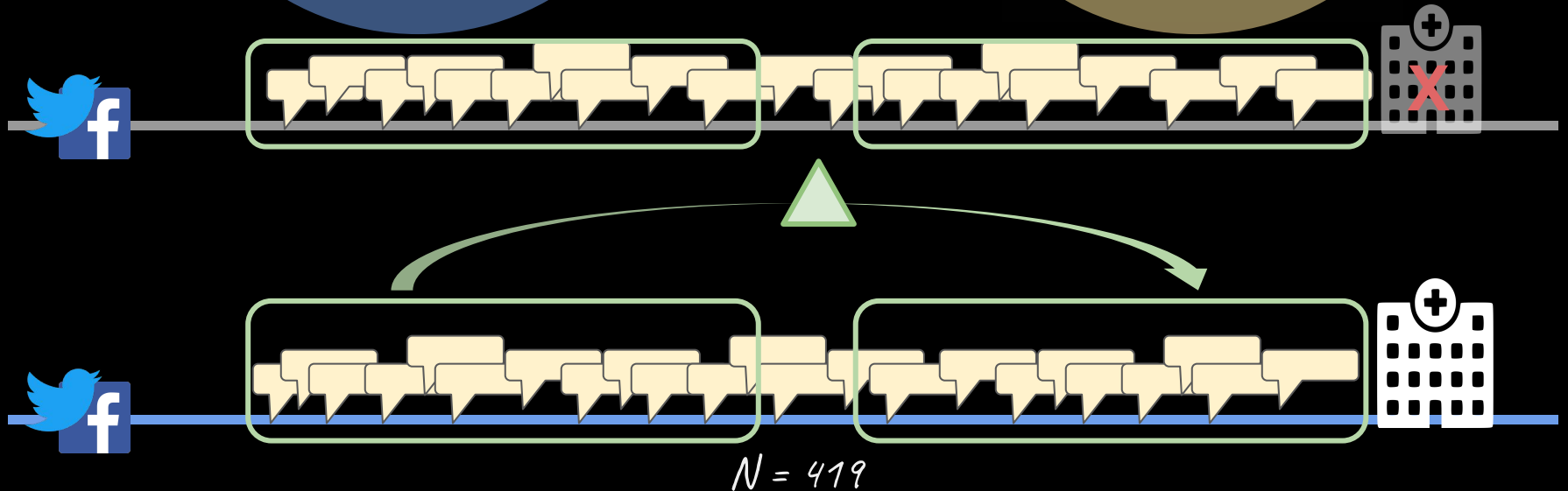
Natural  
Language  
Processing

Psychological  
& Health  
Sciences



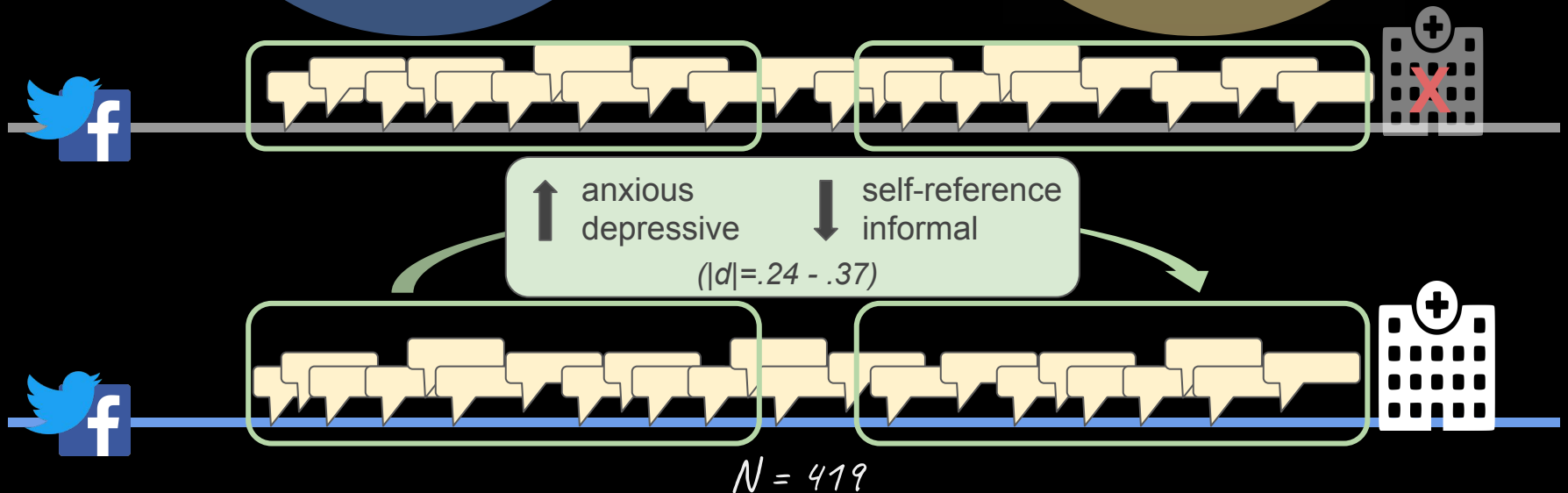
# Natural Language Processing

# Psychological & Health Sciences




# Natural Language Processing

# Psychological & Health Sciences





**Natural  
Language  
Processing**



**Psychological  
& Health  
Sciences**

# Problem

Natural language is written by

# Problem

Natural language is written by **people**.

# Problem

Natural language is written by **people**.

That's sick



# Problem

Natural language is written by **people**.

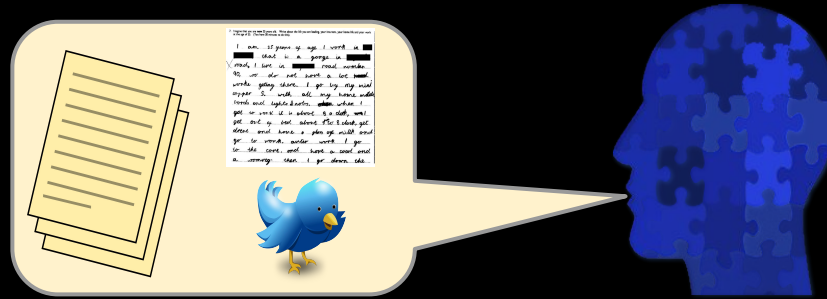


That's sick



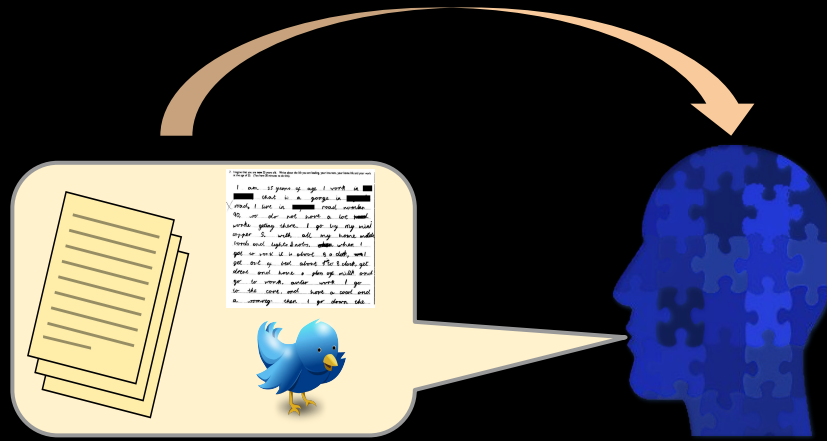


*Natural language is generated by people.*



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, ...

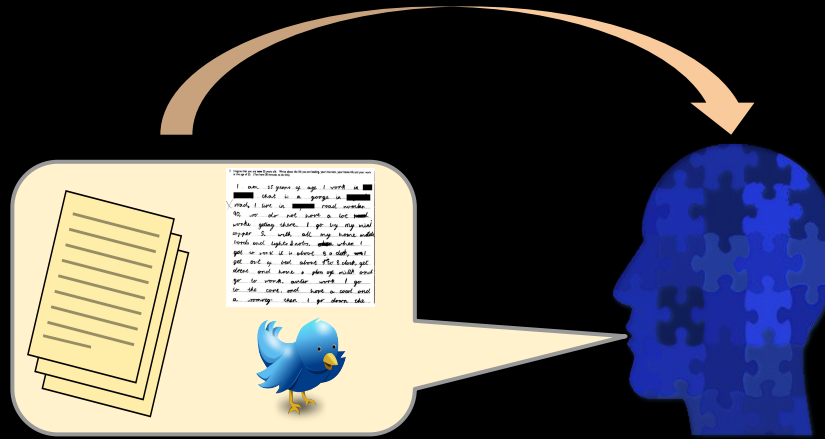
*Natural language is generated by people.*



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, ..., and our language reflects these differences.



Natural language is generated by people.



Shannon,  
1948

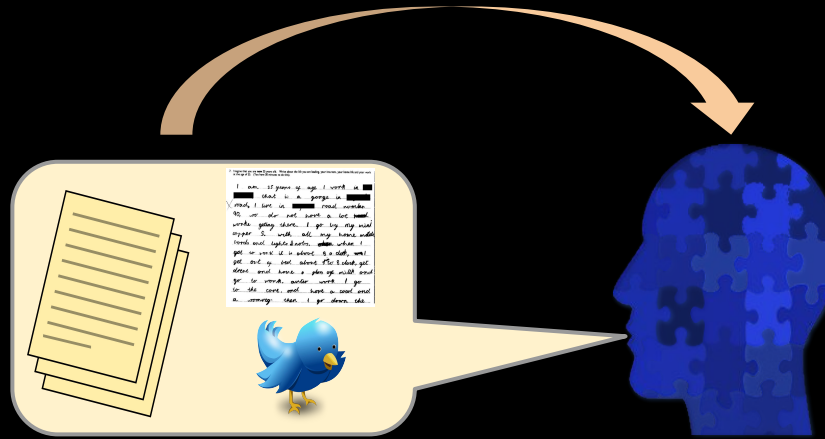
Mosteller &  
Wallace 1963

Clark &  
Schober, 1992

Mairesse, Walker,  
et al., 2007

Hovy & Soegaard,  
2015

# Natural language is generated by people.



*“The common misconception is that language has got to do with words and what they mean. It does not. It has to do with people and what they mean.”*

Shannon,  
1948

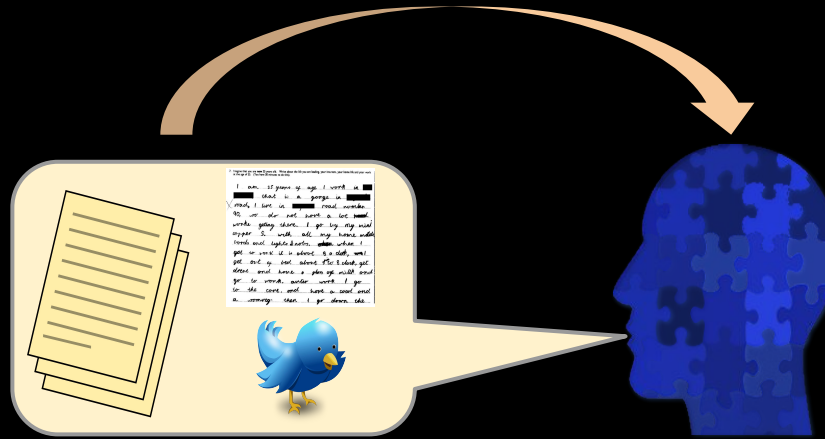
Mosteller &  
Wallace 1963

Clark &  
Schober, 1992

Mairesse, Walker,  
et al., 2007

Hovy & Soogaard,  
2015

Natural language is generated by people.



(MOST LINGUISTS)

Shannon,  
1948

Mosteller &  
Wallace 1963

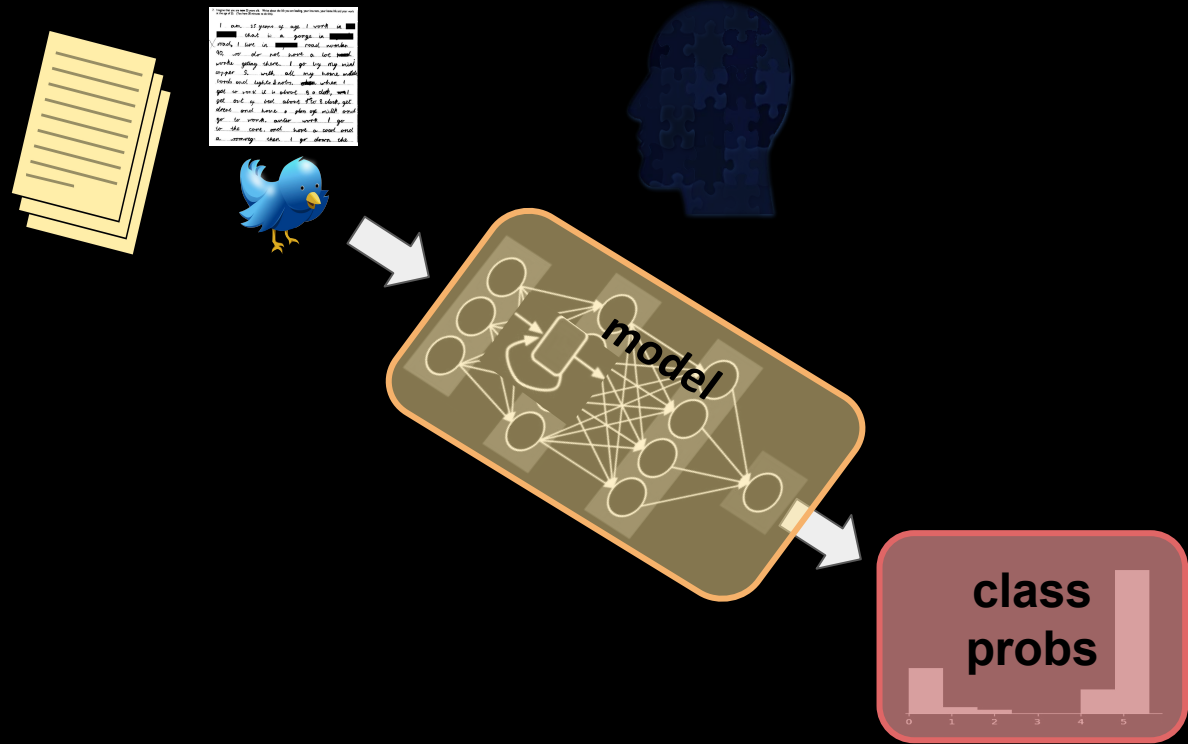
Clark &  
Schober, 1992

Mairesse, Walker,  
et al., 2007

Hovy & Soogaard,  
2015

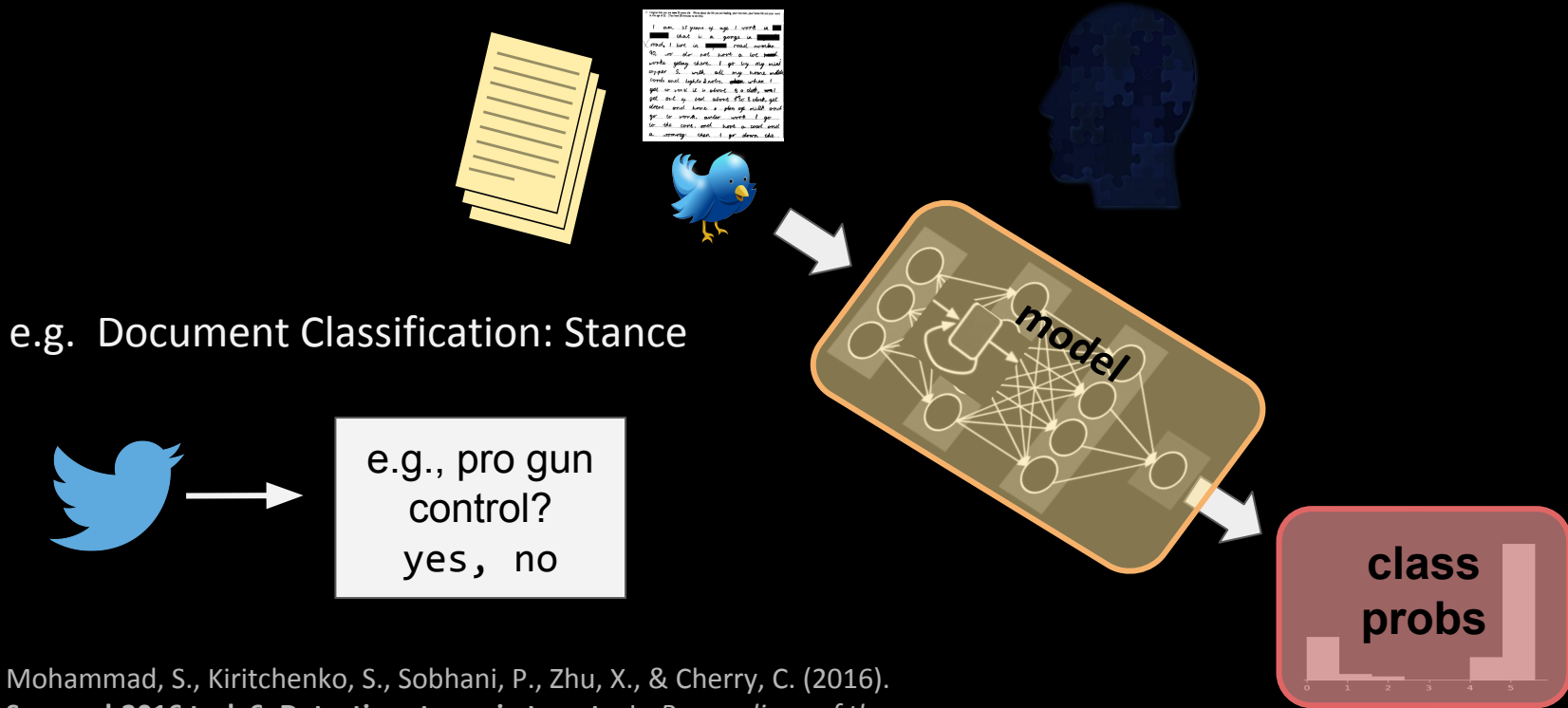
Natural language is generated by people.

Yet, our models:



# Natural language is generated by people.

## Yet, our models:

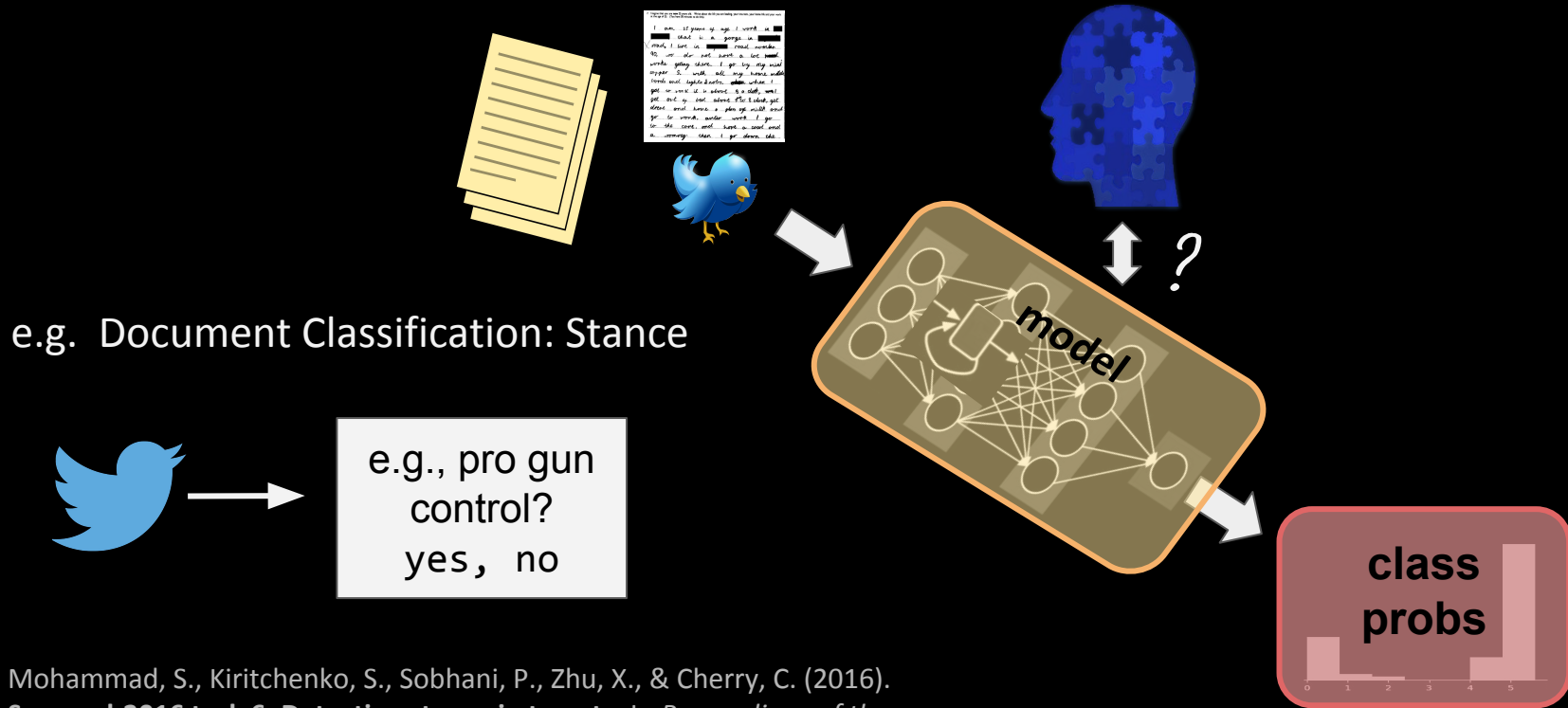


Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016). Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation*.



# Natural language is generated by people.

## Yet, our models:

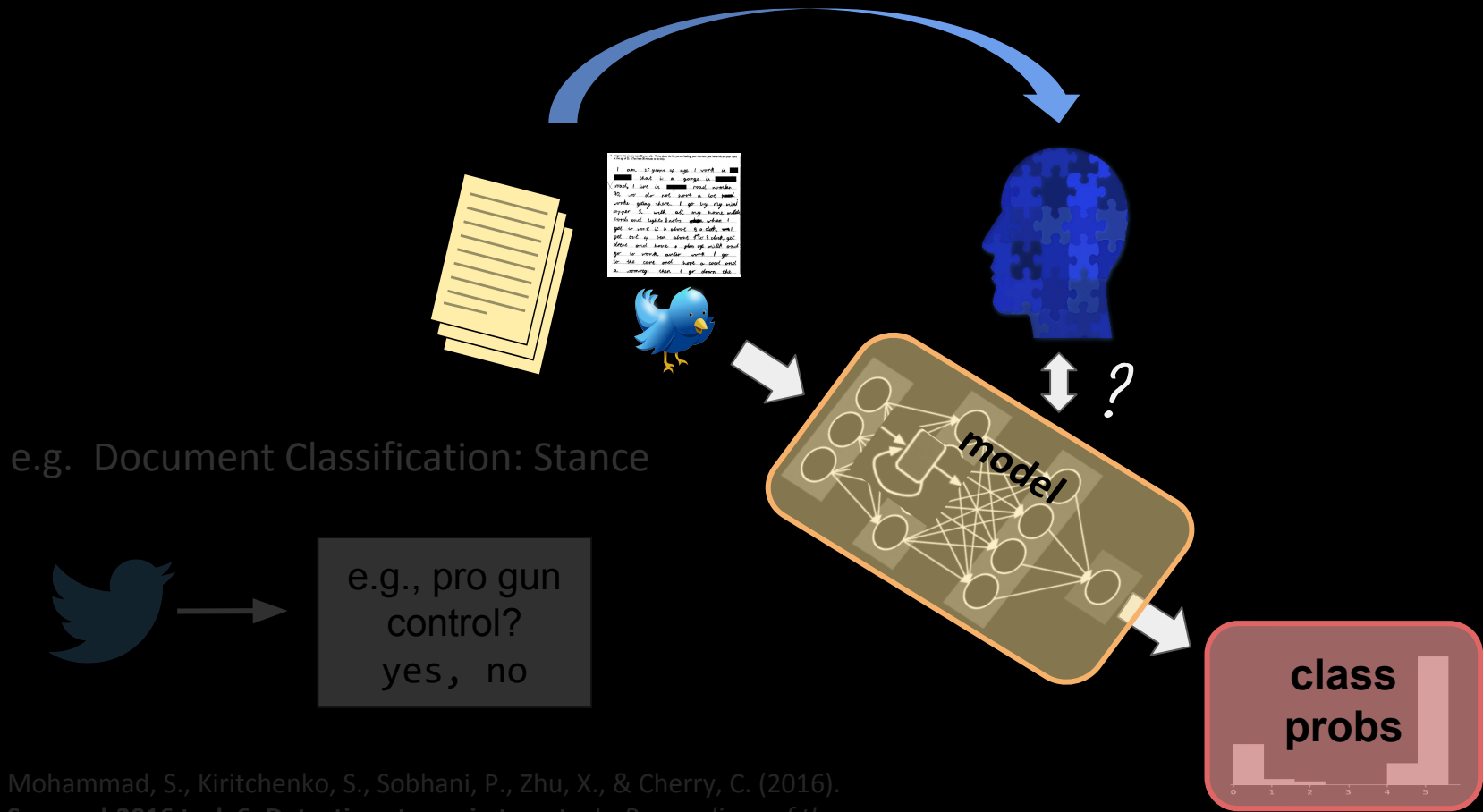


e.g. Document Classification: Stance

e.g., pro gun control?  
yes, no

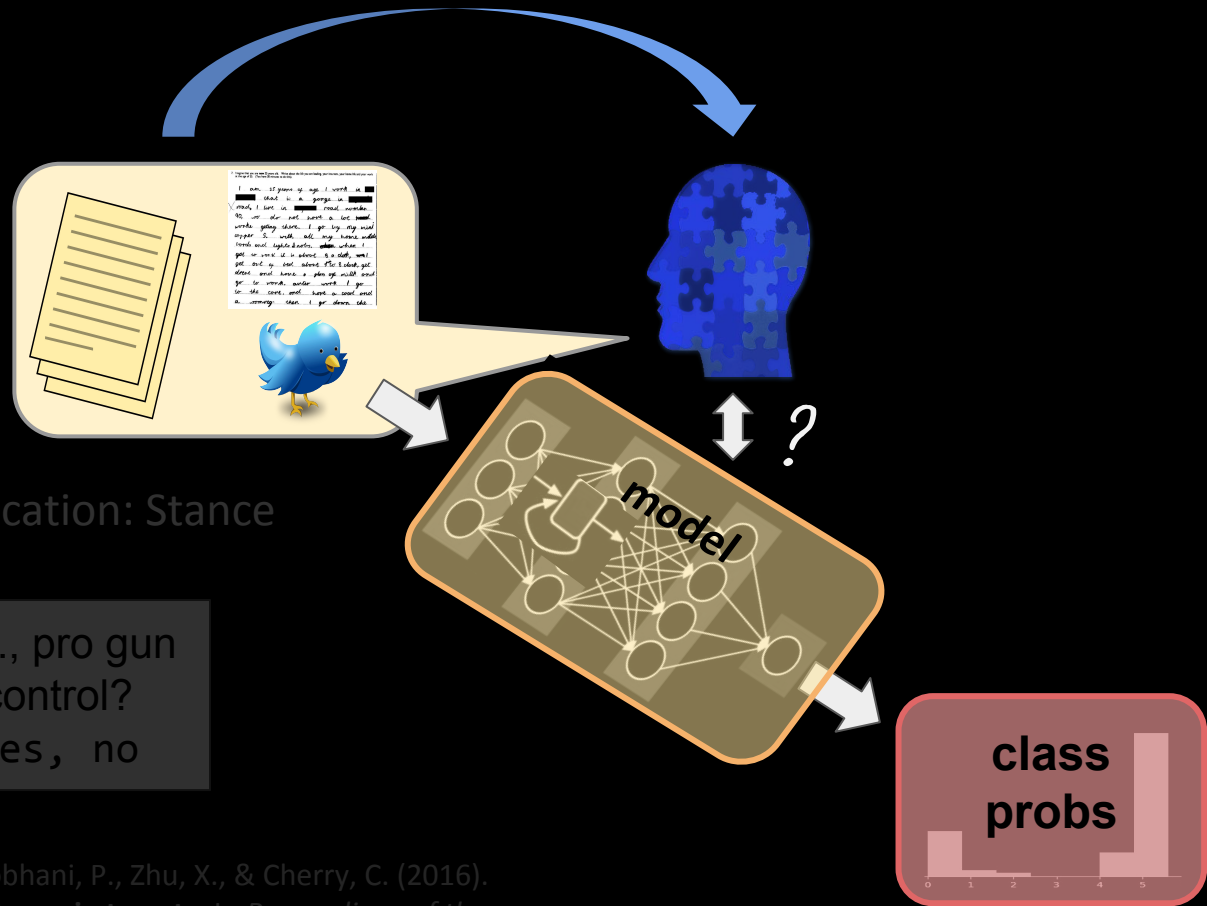
Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016). Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation*.

# Natural language is generated by people.



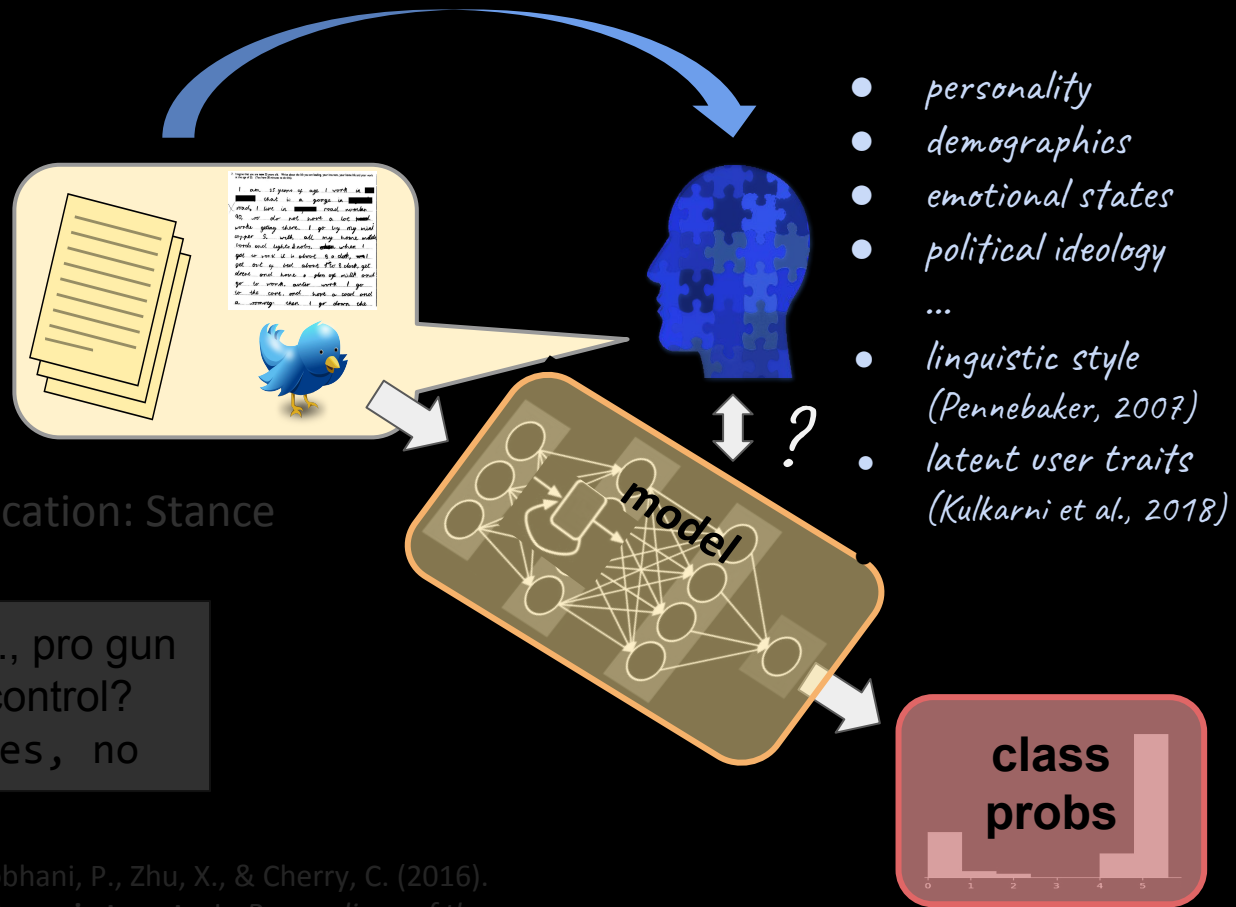
Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016). Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation*.

# Natural language is generated by people.



Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016). Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation*.

# Natural language is generated by people.



e.g. Document Classification: Stance



e.g., pro gun control?  
yes, no

*Natural language is generated by people.*

**What this means for NLP:**

- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



*Natural language is generated by people.*

**What this means for NLP:**

- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



*Natural language is generated by people.*

**What this means for NLP:**

- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



# Approaches to Human Factor Inclusion

1. Adaptive: Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)



# Approaches to Human Factor Inclusion

1. **Adaptive:** Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. **Additive:** Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)

# Approaches to Human Factor Inclusion

1. **Adaptive:** Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. **Additive:** Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. **Bias Correction:** Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)

# Approaches to Human Factor Inclusion

1. What are human “factors”?  
(e.g. age and distinguishing PTSD from Depression)
2. Additive: Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. Bias Correction: Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)

# Adaptation Approach: Domain Adaptation

Features for: source

|

$$\Phi^s(\mathbf{x}) = \langle \mathbf{x}, \mathbf{x}, \mathbf{0} \rangle,$$

target

|

$$\Phi^t(\mathbf{x}) = \langle \mathbf{x}, \mathbf{0}, \mathbf{x} \rangle$$

## Frustratingly Easy Domain Adaptation

**Hal Daumé III**

School of Computing

University of Utah

Salt Lake City, Utah 84112

me@hal3.name

### Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case

supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data from

# Adaptation Approach: Domain Adaptation

Features for: source

$$\Phi^s(\mathbf{x}) = \langle \mathbf{x}, \mathbf{x}, \mathbf{0} \rangle,$$

target

$$\Phi^t(\mathbf{x}) = \langle \mathbf{x}, \mathbf{0}, \mathbf{x} \rangle$$

```
newX = []
for all x in source_x:
    newX.append(x + x + [0]*len(x))
for all x in target_x:
    newX.append(x + [0]*len(x), x)
```

## Frustratingly Easy Domain Adaptation

Hal Daumé III

School of Computing

University of Utah

Salt Lake City, Utah 84112

me@hal3.name

### Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case

supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data from

# Adaptation Approach: Domain Adaptation

Features for: source

$$\Phi^s(\mathbf{x}) = \langle \mathbf{x}, \mathbf{x}, \mathbf{0} \rangle,$$

target

$$\Phi^t(\mathbf{x}) = \langle \mathbf{x}, \mathbf{0}, \mathbf{x} \rangle$$

```
newX = []
for all x in source_x:
    newX.append(x + x + [0]*len(x))
for all x in target_x:
    newX.append(x + [0]*len(x), x)
```

```
newY = source_y + target_y
```

```
model = model.train(newX,newY)
```

## Frustratingly Easy Domain Adaptation

Hal Daumé III

School of Computing

University of Utah

Salt Lake City, Utah 84112

me@hal3.name

### Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case

supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data from

# Human Factors

--- Any attribute, represented as a continuous or discrete variable, of the humans generating the natural language.

E.g.

- Gender
- Age
- Personality
- Ethnicity
- Socio-economic status

# Adaptation Approach: Factor Adaptation

## Human Centered NLP with User-Factor Adaptation

Veronica E. Lynn, Youngseo Son, Vivek Kulkarni  
Niranjan Balasubramanian and H. Andrew Schwartz  
Stony Brook University  
Stony Brook, NY  
{velynn, yson, vvkulkarni, niranjan, has}@cs.stonybrook.edu

### Abstract

We pose the general task of *user-factor adaptation* — adapting supervised learning models to real-valued user factors inferred from a background of their lan-

and Costa Jr., 1989; Ruscio and Ruscio, 2000; Widiger and Samuel, 2005).

Here, we ask how one can adapt NLP models to real-valued human *factors* — continuous valued attributes that capture fine-grained differences be-

## Residualized Factor Adaptation for Community Social Media Prediction Tasks

Mohammadzaman Zamani,<sup>1</sup> H. Andrew Schwartz,<sup>1</sup> Veronica E. Lynn,<sup>1</sup>  
Salvatore Giorgi,<sup>2</sup> and Niranjan Balasubramanian<sup>1</sup>  
<sup>1</sup>Computer Science Department, Stony Brook University  
<sup>2</sup>Department of Psychology, University of Pennsylvania  
mzamani@cs.stonybrook.edu

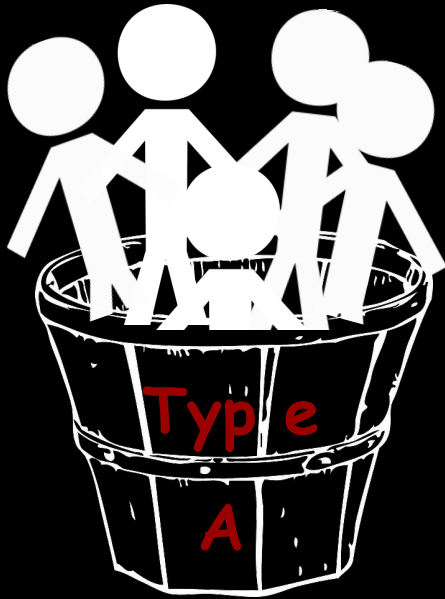
### Abstract

Predictive models over social media language promise in capturing community

linked to socio-demographic factors (age, gender, race, education, income levels) with many social scientific studies supporting their predictive value (Golder et al., 2002) and build the fun-



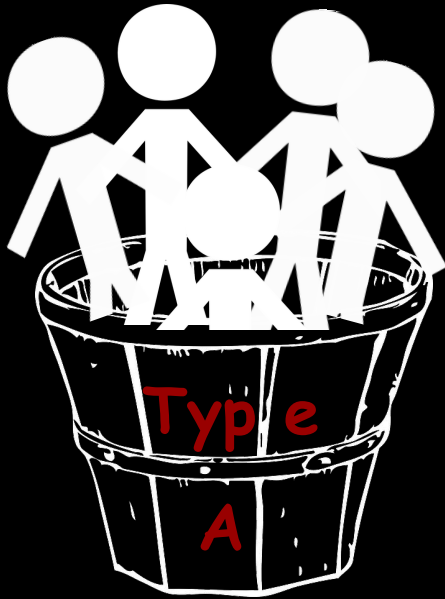
# Adaptation



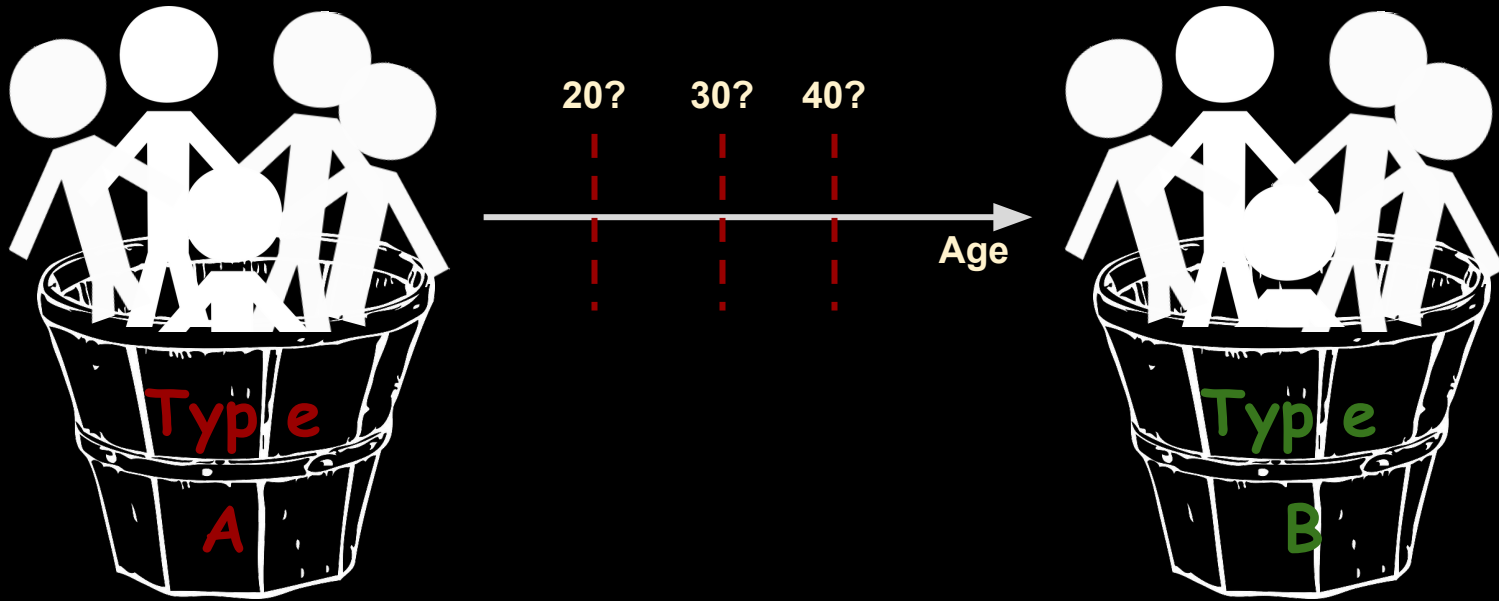
typically requires putting people into discrete bins

*“most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]”*

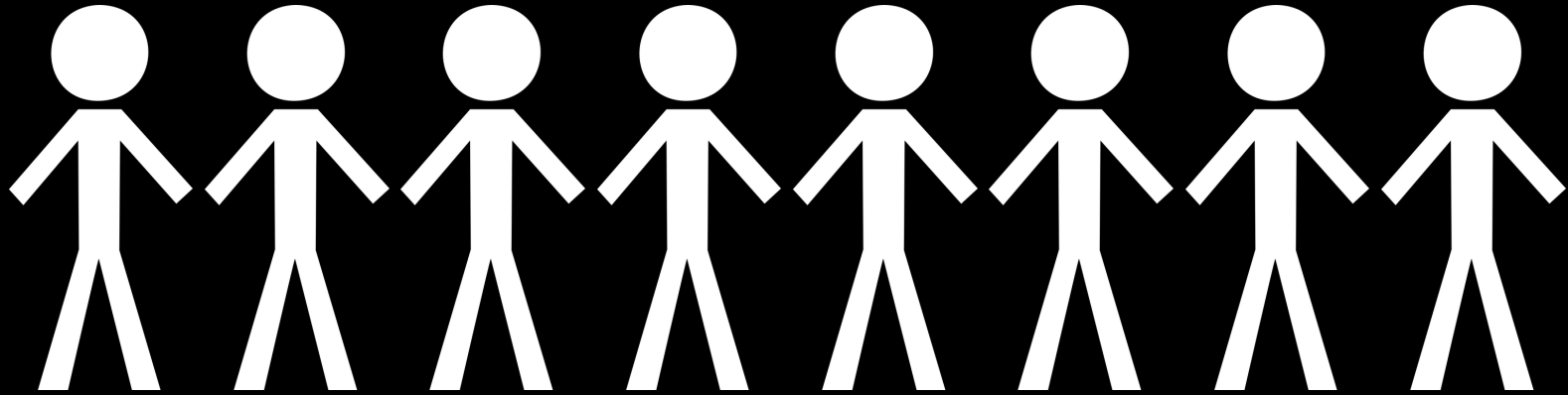
(Haslam et al., 2012)



*“most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]”*  
(Haslam et al., 2012)

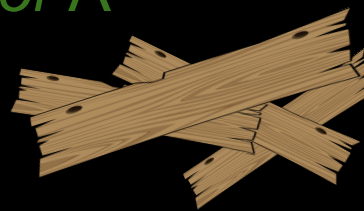
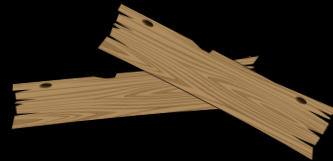


*“most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]”*  
(Haslam et al., 2012)

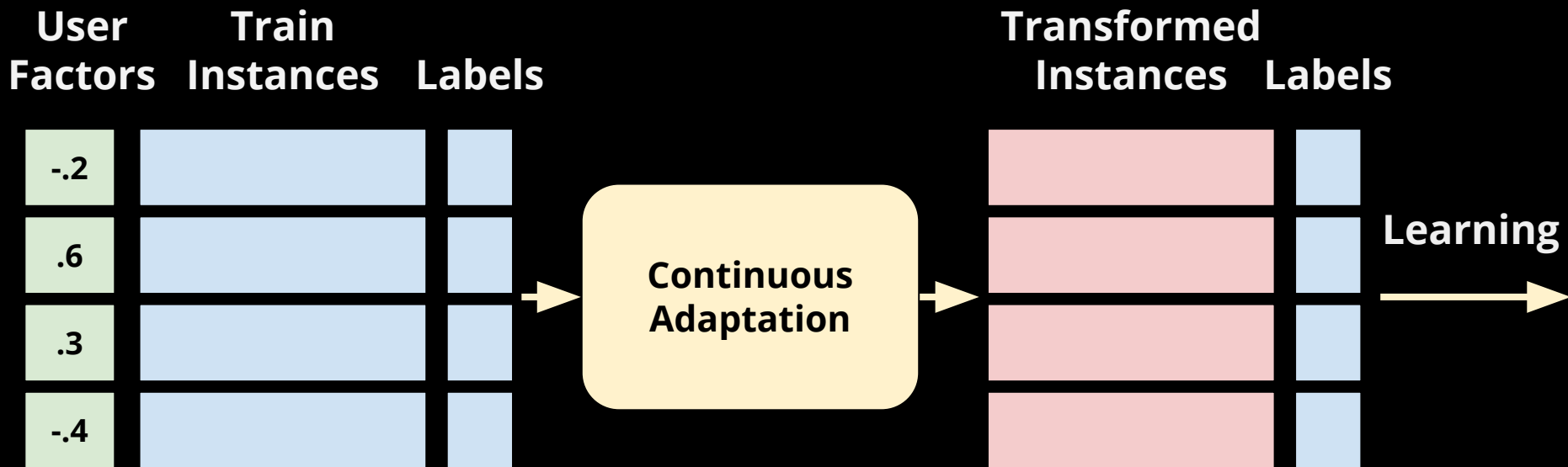


*Less Factor A*

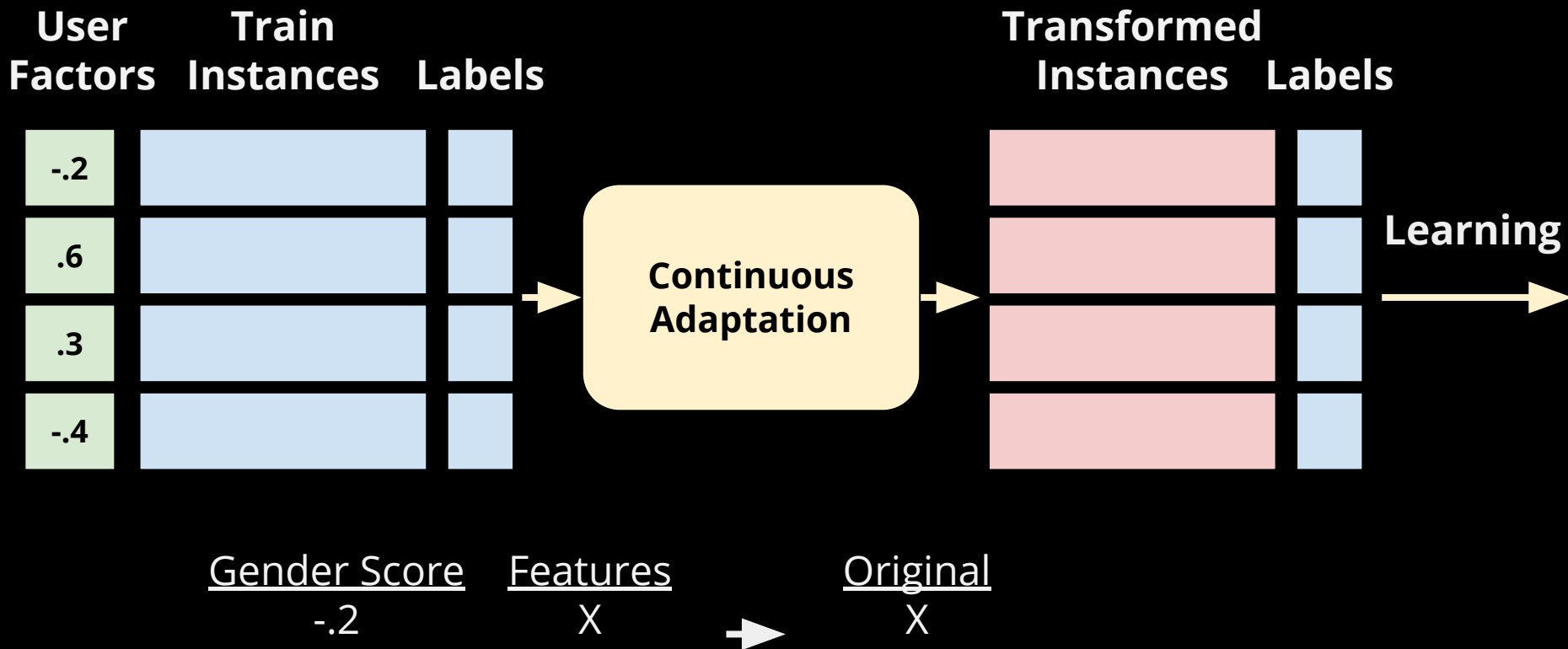
*More Factor A*



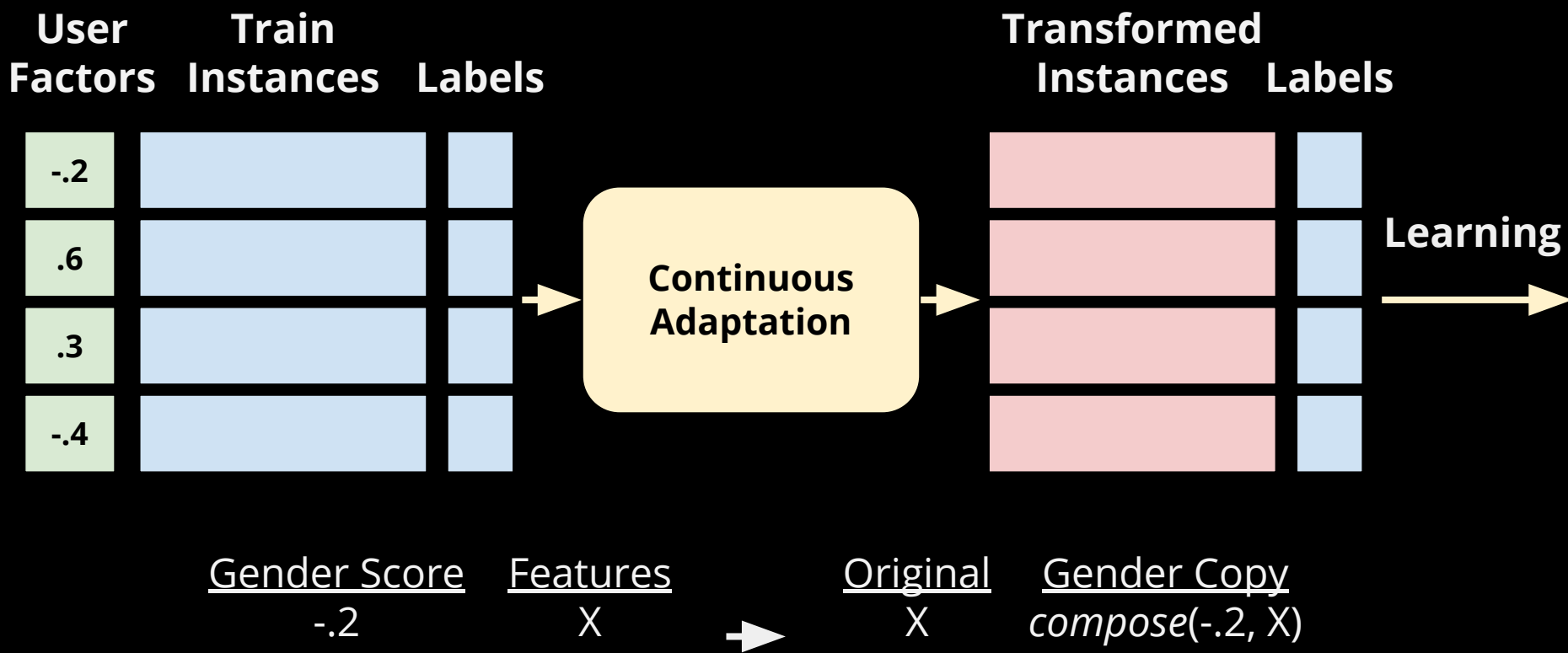
# Our Method: Continuous Adaptation



# Our Method: Continuous Adaptation



# Our Method: Continuous Adaptation



(Lynn et al., 2017)

# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function  $c$  combines  $d$  user factor scores  $f_{u,d}$  with original feature values  $\mathbf{x}$ :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \dots, c(f_{u,d}, \mathbf{x}) \rangle$$



# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function  $c$  combines  $d$  user factor scores  $f_{u,d}$  with original feature values  $\mathbf{x}$ :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \dots, c(f_{u,d}, \mathbf{x}) \rangle$$

User	Factor Classes	Augmented Instance $\Phi(\mathbf{x}, u)$
User 1	$F_1$	$\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0} \rangle$
User 2	$F_2$	$\langle \mathbf{x}, \mathbf{0}, \mathbf{x}, \mathbf{0}, \dots, \mathbf{0} \rangle$
User 3	$F_1, F_3$	$\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{x}, \dots, \mathbf{0} \rangle$
User 4	$F_k$	$\langle \mathbf{x}, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0}, \mathbf{x} \rangle$

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector  $\mathbf{x}$  under different factor class mappings. With  $k$  domains the augmented feature vector is of length  $n(k + 1)$ .

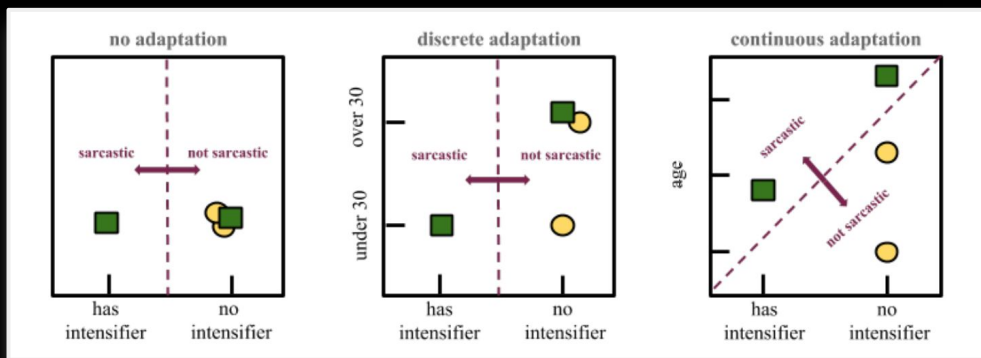
(Lynn et al., 2017)

# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function  $c$  combines  $d$  user factor scores  $f_{u,d}$  with original feature values  $\mathbf{x}$ :

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \dots, c(f_{u,d}, \mathbf{x}) \rangle$$



User	Factor Classes	Augmented Instance $\Phi(\mathbf{x}, u)$
User 1	$F_1$	$\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0} \rangle$
User 2	$F_2$	$\langle \mathbf{x}, \mathbf{0}, \mathbf{x}, \mathbf{0}, \dots, \mathbf{0} \rangle$
User 3	$F_1, F_3$	$\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{x}, \dots, \mathbf{0} \rangle$
User 4	$F_k$	$\langle \mathbf{x}, \mathbf{0}, \mathbf{0}, \dots, \mathbf{0}, \mathbf{x} \rangle$

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector  $\mathbf{x}$  under different factor class mappings. With  $k$  domains the augmented feature vector is of length  $n(k + 1)$ .

(Lynn et al., 2017)

# Main Results

Adaptation improves over unadapted baselines (Lynn et al., 2017)

Task	Metric	No Adaptation	Gender	Personality	Latent (User Embed)
Stance	F1	64.9	<b>65.1 (+0.2)</b>	<b>66.3 (+1.4)</b>	<b>67.9 (+3.0)</b>
Sarcasm	F1	73.9	<b>75.1 (+1.2)</b>	<b>75.6 (+1.7)</b>	<b>77.3 (+3.4)</b>
Sentiment	Acc.	60.6	<b>61.0 (+0.4)</b>	<b>61.2 (+0.6)</b>	<b>60.7 (+0.1)</b>
PP-Attach	Acc.	71.0	70.7 (-0.3)	70.2 (-0.8)	70.8 (-0.2)
POS	Acc.	91.7	<b>91.9 (+0.2)</b>	91.2 (-0.5)	90.9 (-0.8)

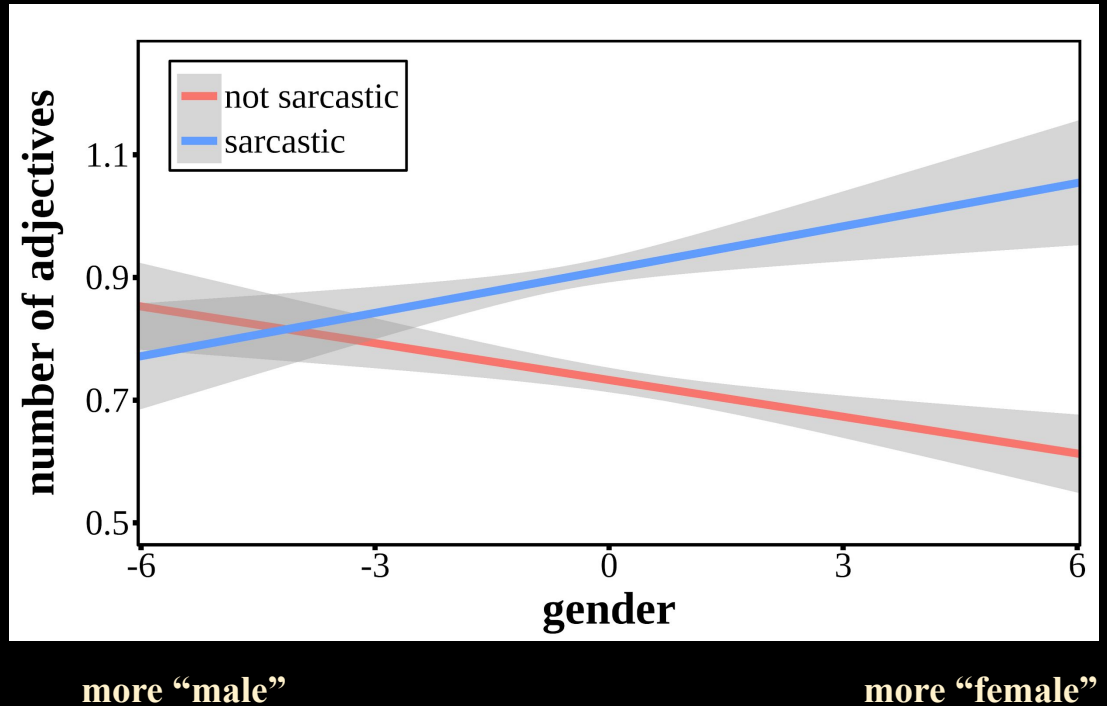
# Example: How Adaptation Helps

Women

more adjectives → sarcasm

Men

more adjectives → no sarcasm



# Problem

User factors are not always available.

# Solution: User Factor Inference

## past tweets

Niranjan @b\_niranjan · Sep 2

There must be a word for trending #hashtags that you know you will regret if you click. Is there?

Niranjan @b\_niranjan · Aug 31

Passwords spiral: Forget password for the acct you use twice a year. Ask for reset. Can't use previous. Create a new one to forget later.

Niranjan @b\_niranjan · Jul 31

Thrilled to hear @acl2017's diversity efforts as the first thing in the conference.



1



→ **inferred factors**

## Known

Age (Sap et al. 2014)

Gender (Sap et al. 2014)

Personality (Park et al. 2015)

## Latent

User Embeddings

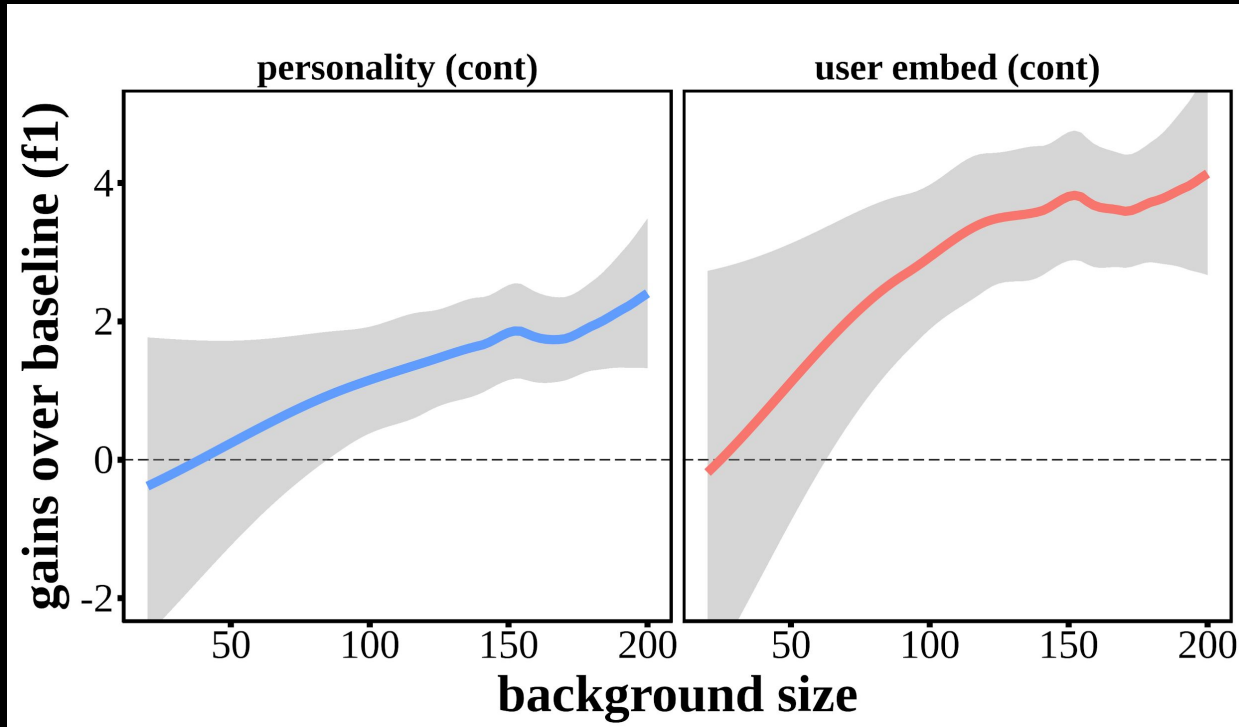
(Kulkarni et al. 2017)

Word2Vec

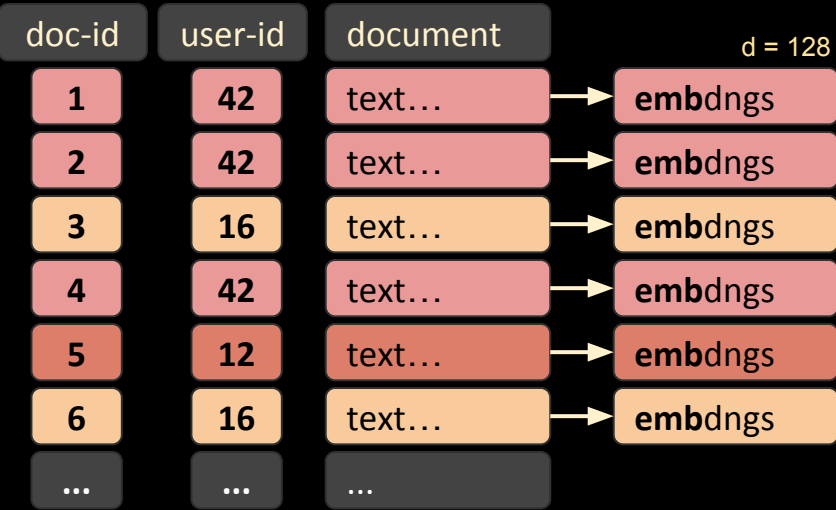
TF-IDF

# Background Size

Using more background tweets to infer factors produces larger gains



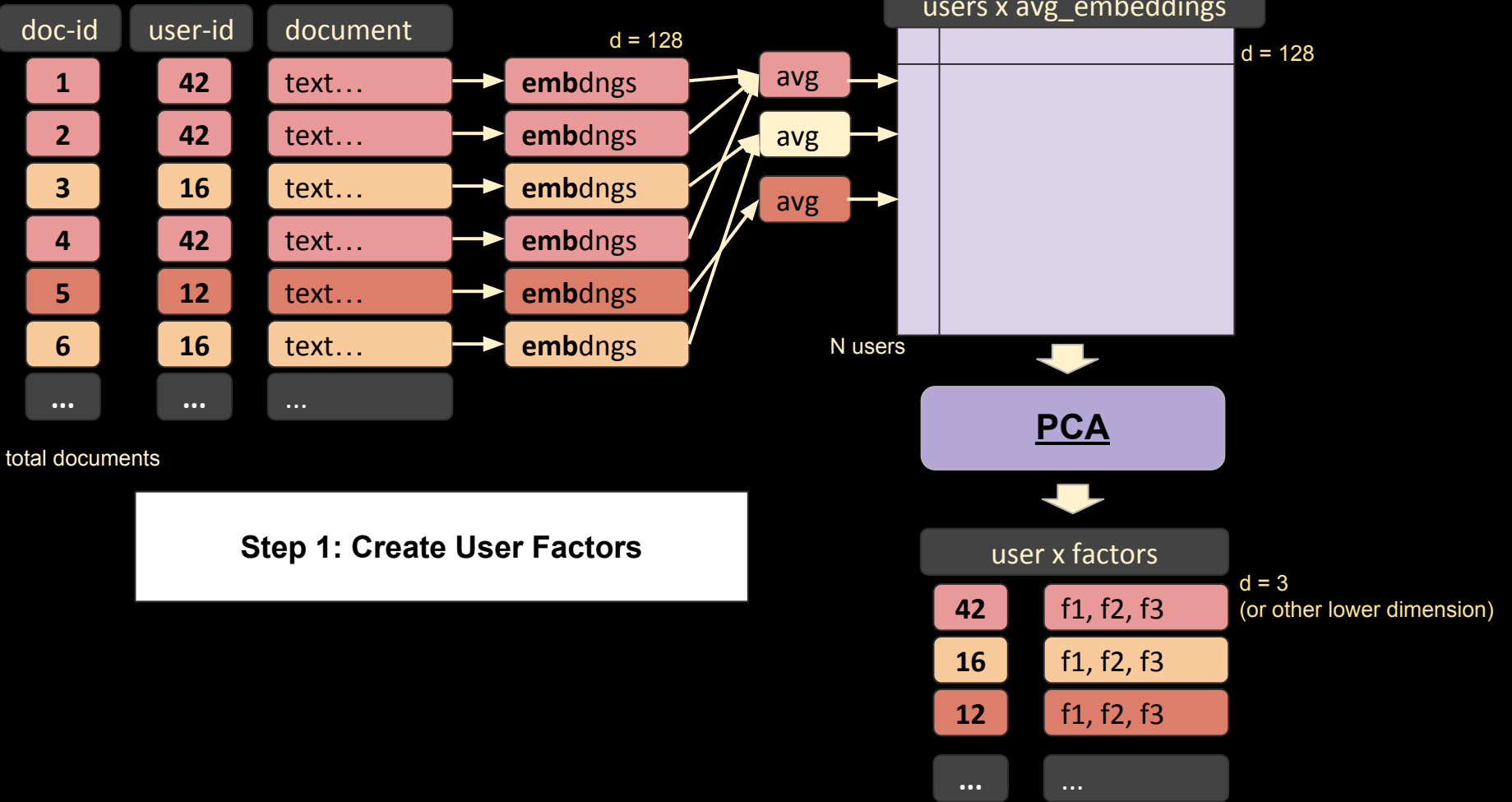
# Full User Factors Adaptation Pipeline: with latent factors from training



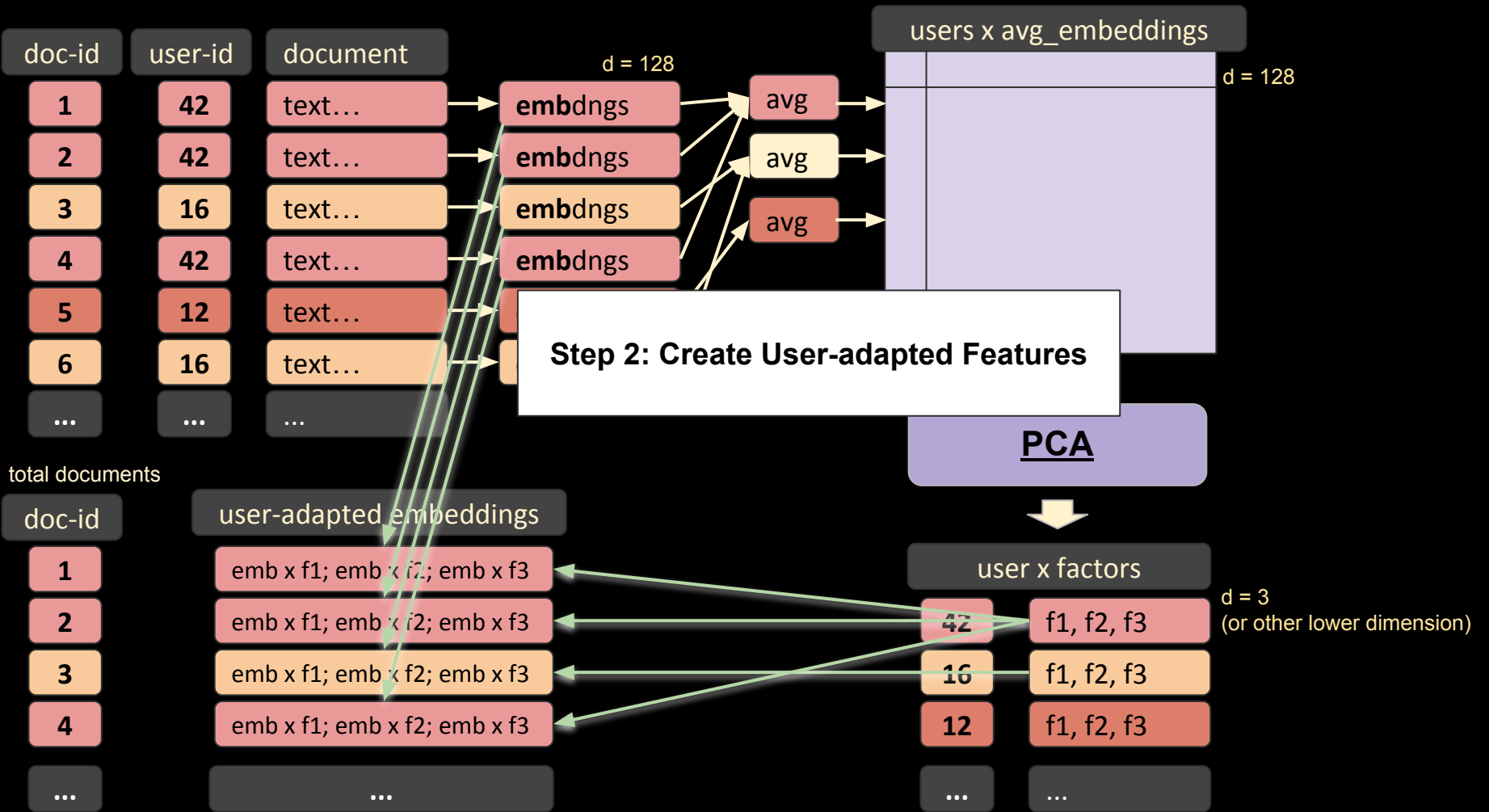
total documents



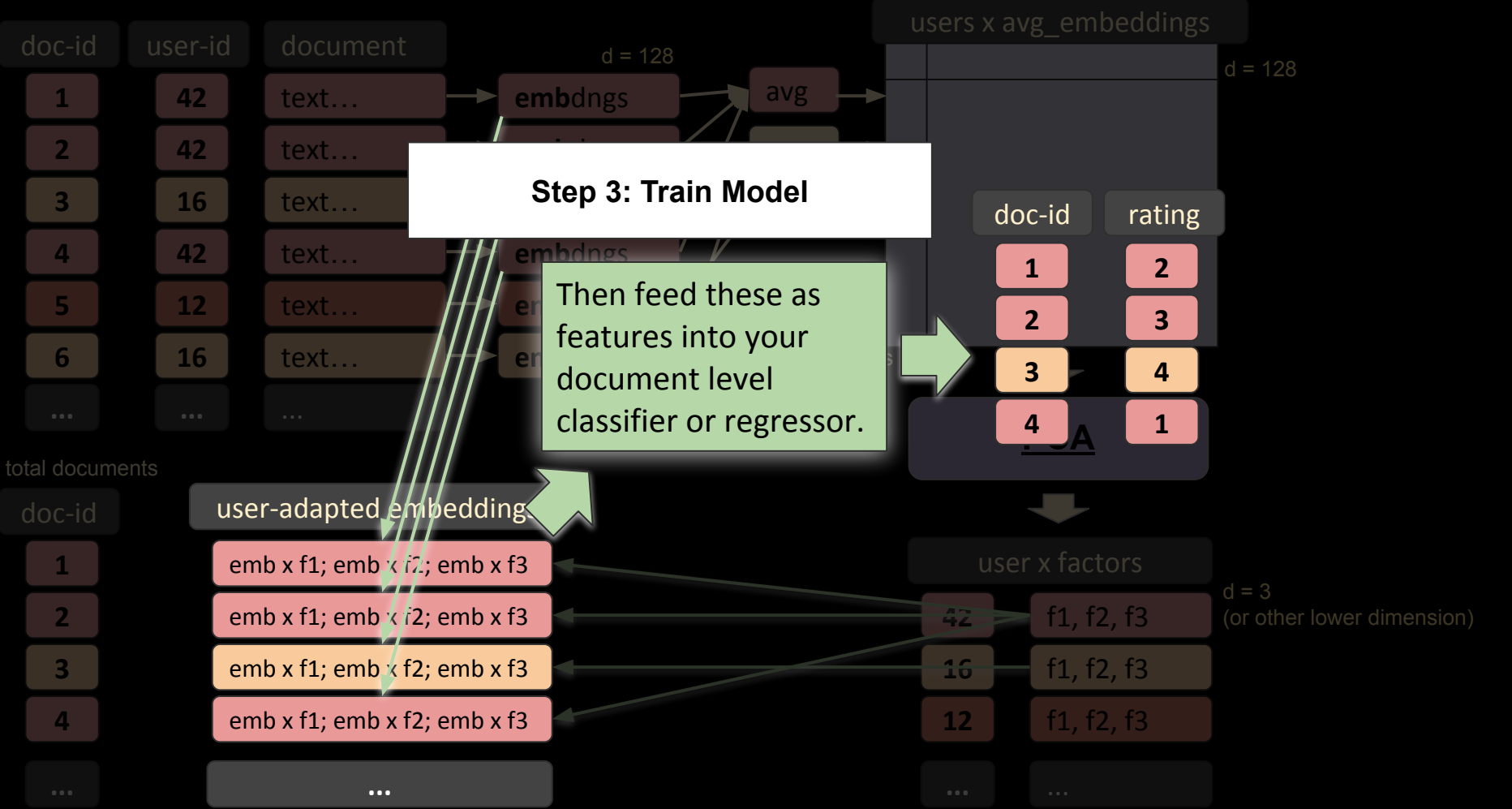
# Full User Factors Adaptation Pipeline: with latent factors from training



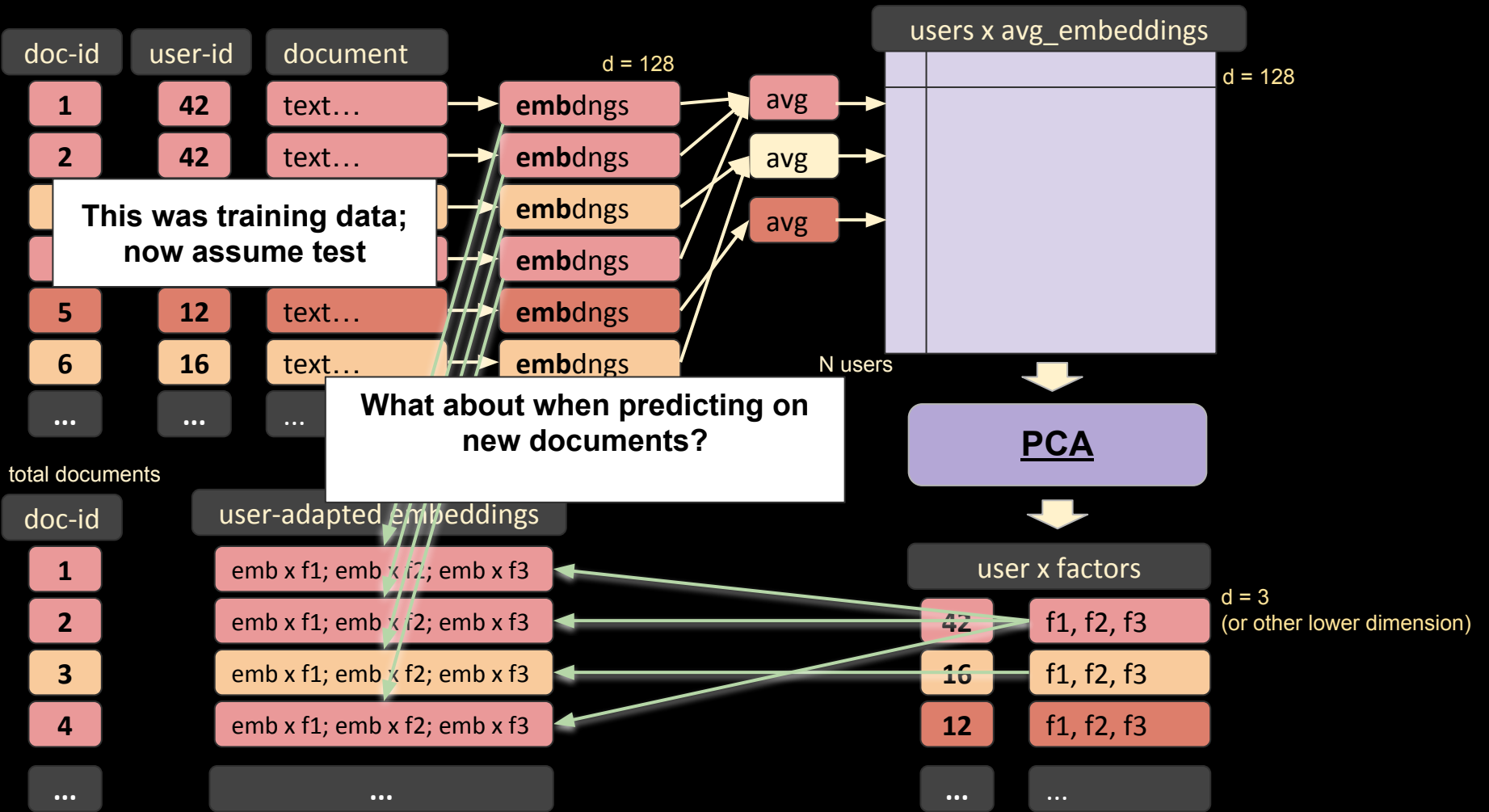
# Full User Factors Adaptation Pipeline: with latent factors from training



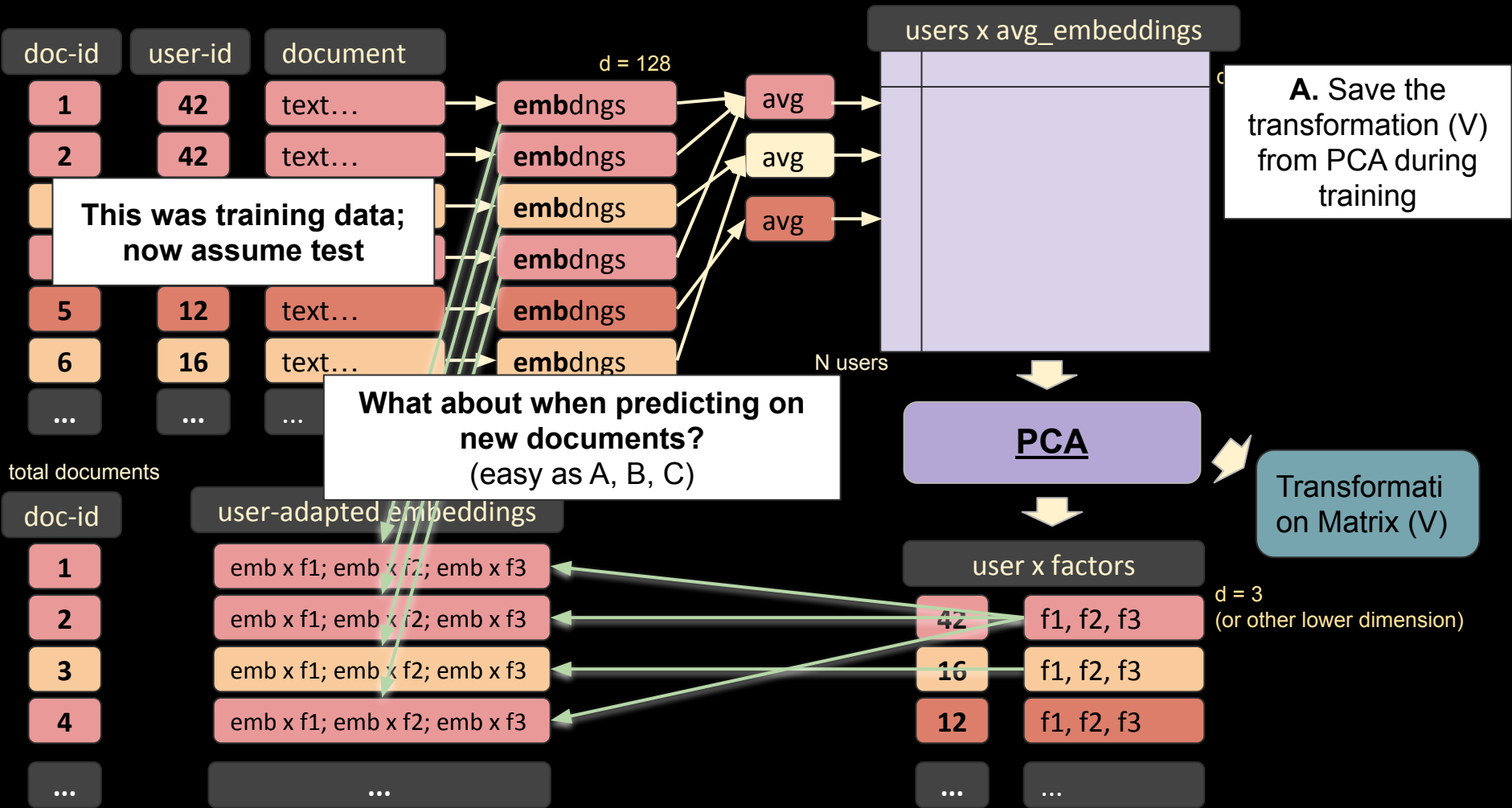
# Full User Factors Adaptation Pipeline: with latent factors from training



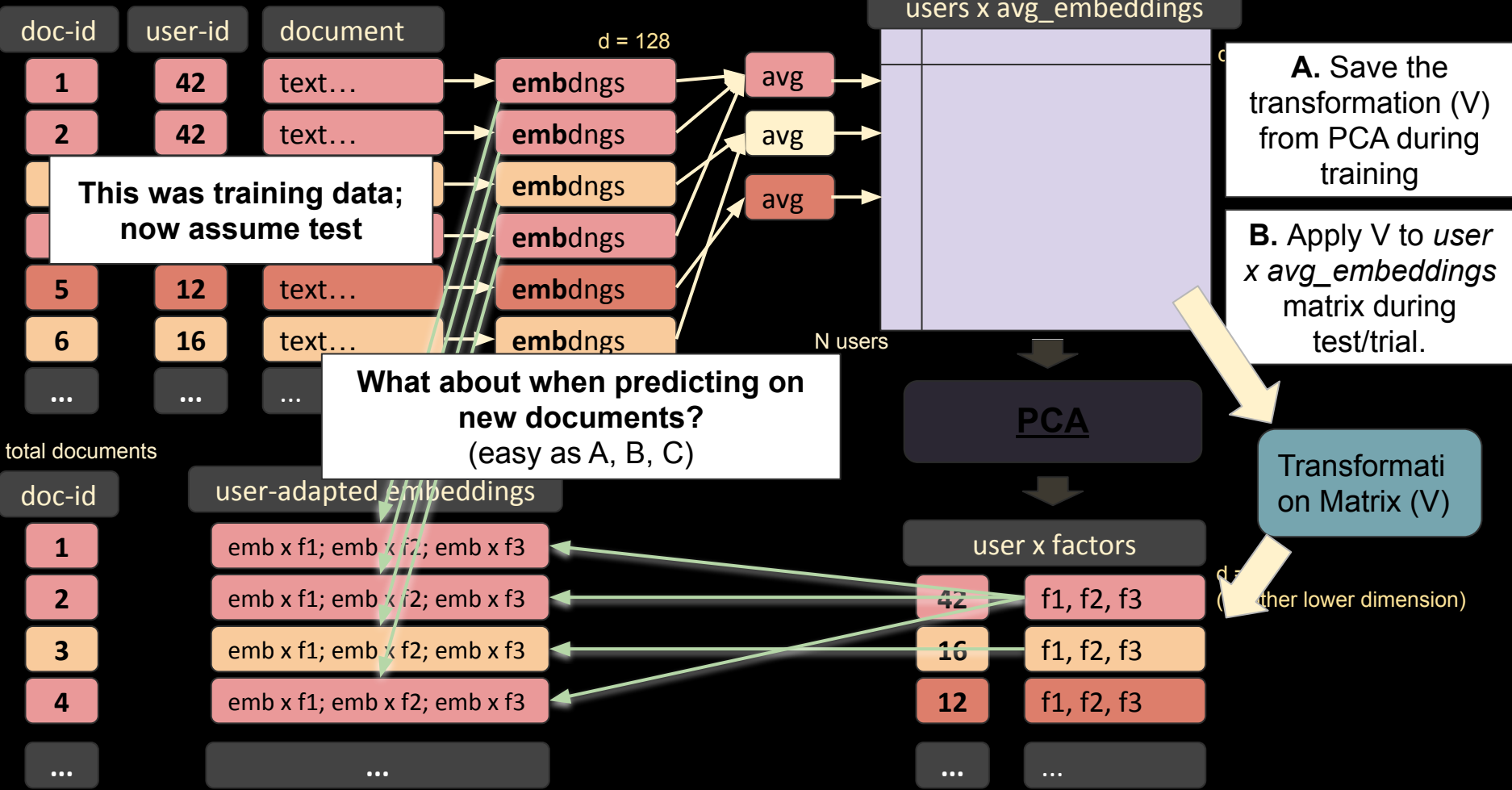
# Full User Factors Adaptation Pipeline: with latent factors from training



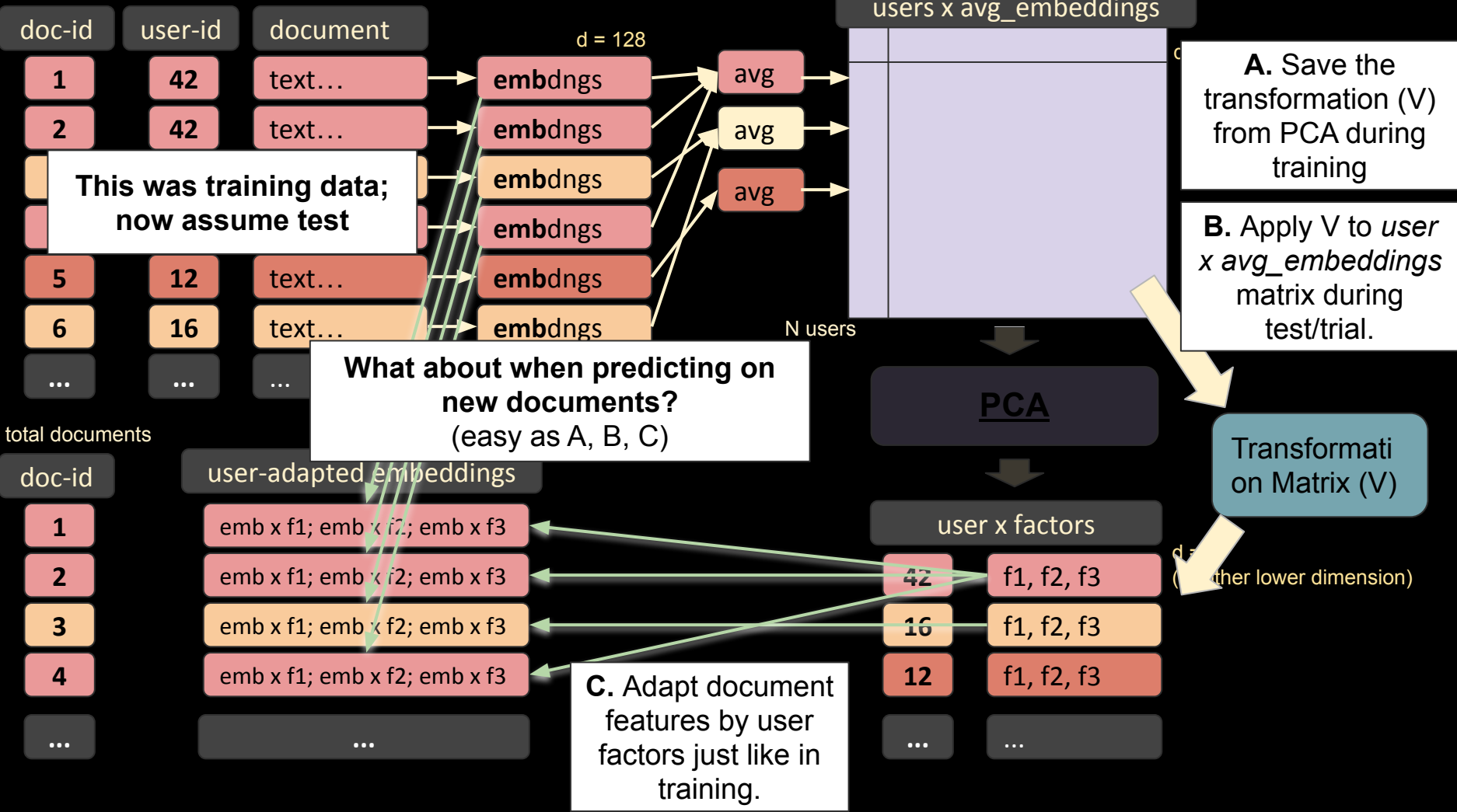
# Full User Factors Adaptation Pipeline: with latent factors from training



# Full User Factors Adaptation Pipeline: with latent factors from training



# Full User Factors Adaptation Pipeline: with latent factors from training



# Approaches to Human Factor Inclusion

1. **Adaptive:** Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. **Additive:** Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. **Bias Correction:** Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)



# Approaches to Human Factor Inclusion

1. **Adaptive:** Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. **Additive:** Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. **Bias Correction:** Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)

# Ethics in NLP to Human Factor Inclusion

1. Adaptive: Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. Additive: Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. Bias Correction: Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)

# Ethics in NLP

Bias

Privacy

Ethical Research

# Ethics in NLP

## Types of bias in NLP tasks:

- Outcome Disparity: Predicted distribution given A, are dissimilar from ideal distribution given A
  - Selection bias
  - Label bias
  - Over-amplification
- Error Disparity: Predicts less accurate for authors of given demographics.
- Semantic Bias: Representations of meaning store demographic associations.

# Two Examples

## The WSJ Effect

model  
accuracy

Jørgensen/Hovy/Sogaard, 2015  
Hovy & Sogaard, 2015

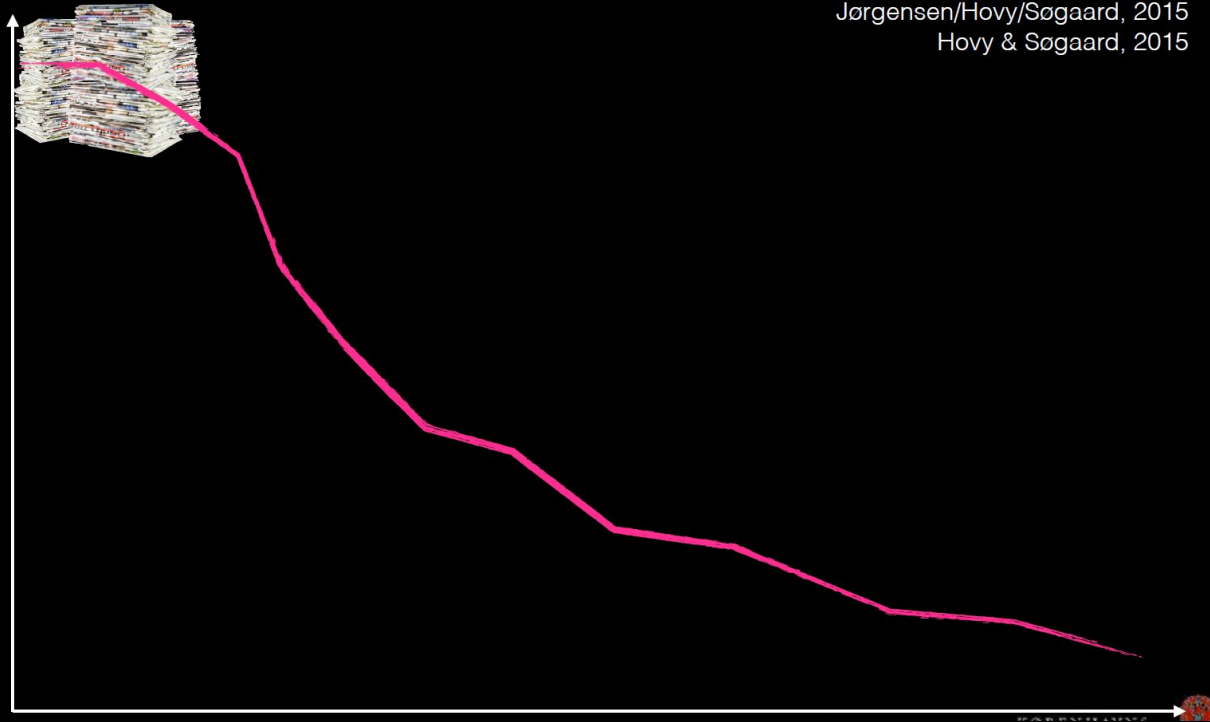


distance from "standard" WSJ author demographics

# Two Examples

## The WSJ Effect

model  
accuracy



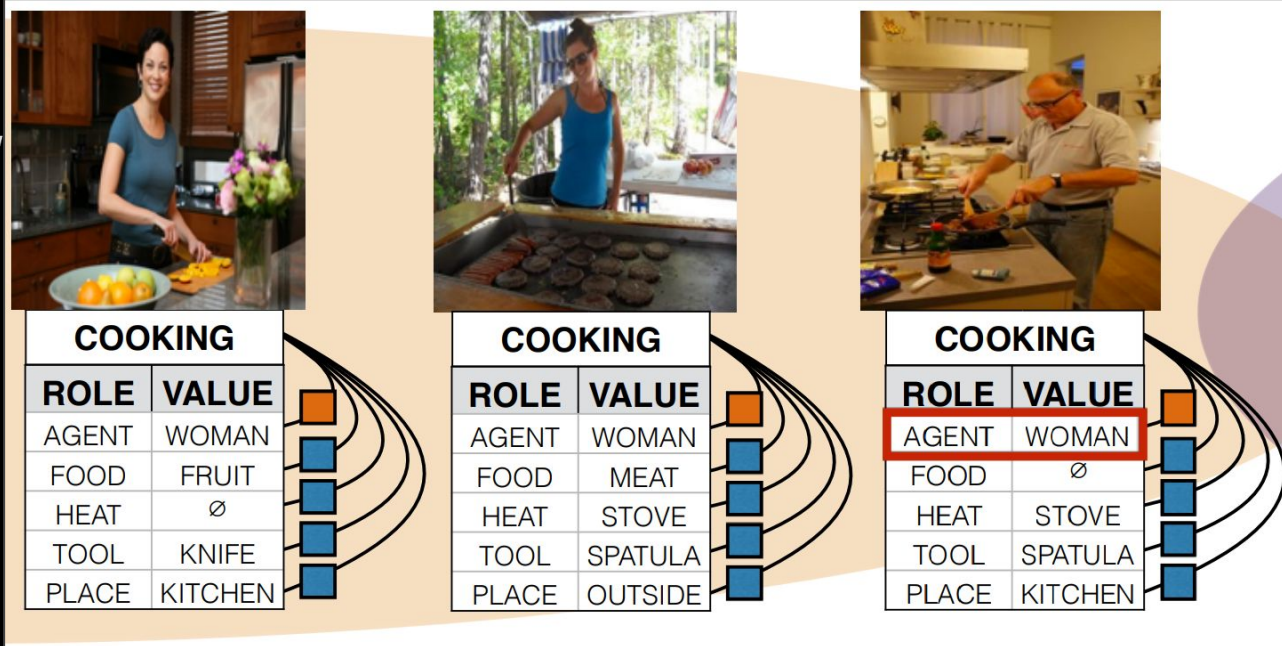
Jørgensen/Hovy/Sogaard, 2015  
Hovy & Sogaard, 2015

distance from "standard" WSJ author demographics

# Two Examples

The W

model accuracy



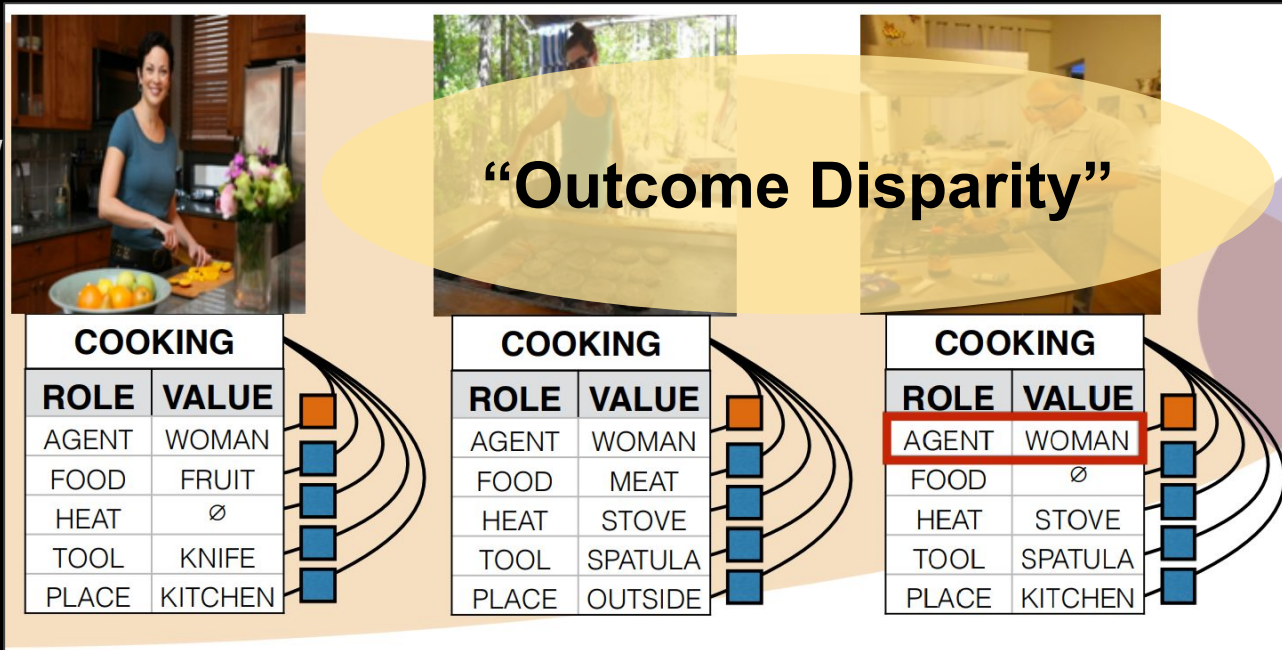
Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics

# Two Examples

## The W

model accuracy



**“Error Disparity”**

Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics



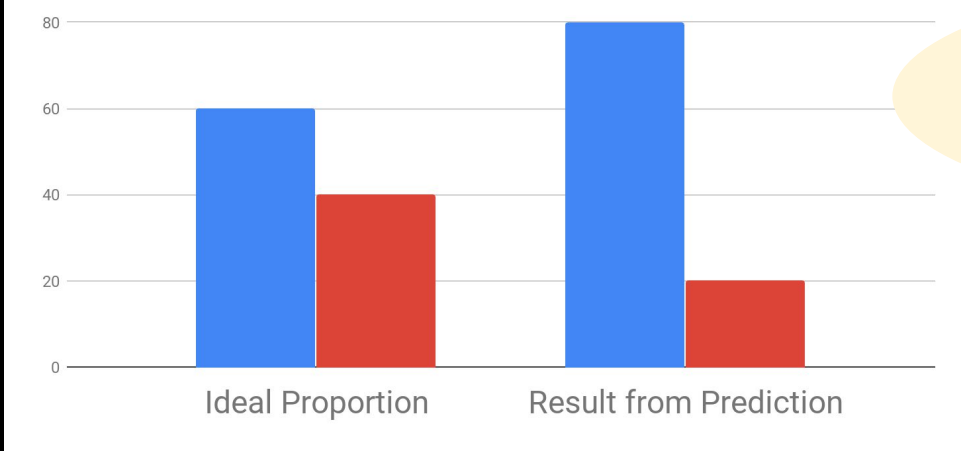
*Our data and models are (human) biased.*

**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

**“Error Disparity”**

*Our data and models are (human) biased.*

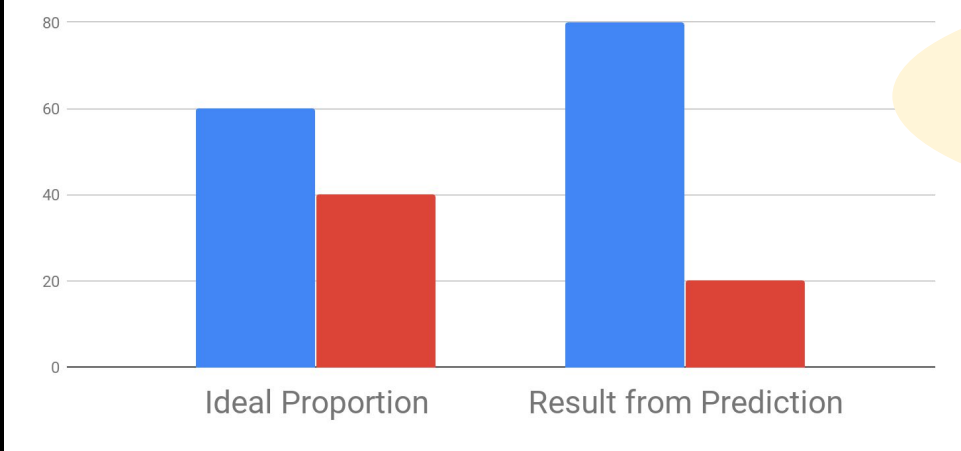


**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

**“Error Disparity”**

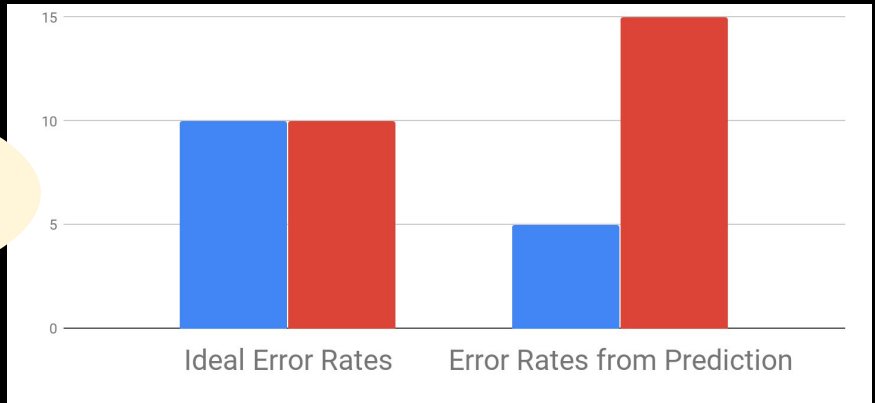
*Our data and models are (human) biased.*



**“Outcome Disparity”**

Person-level  
■ attribute = 1  
■ attribute = 2

**“Error Disparity”**



### outcome disparity

The distribution of outcomes, given attribute  $A$ , is dissimilar than the *ideal distribution*:

$$Q(\hat{Y}_i|A) \neq P(Y_i|A)$$

### Target Population

features

$X_{target}$

(Application Side)

predict

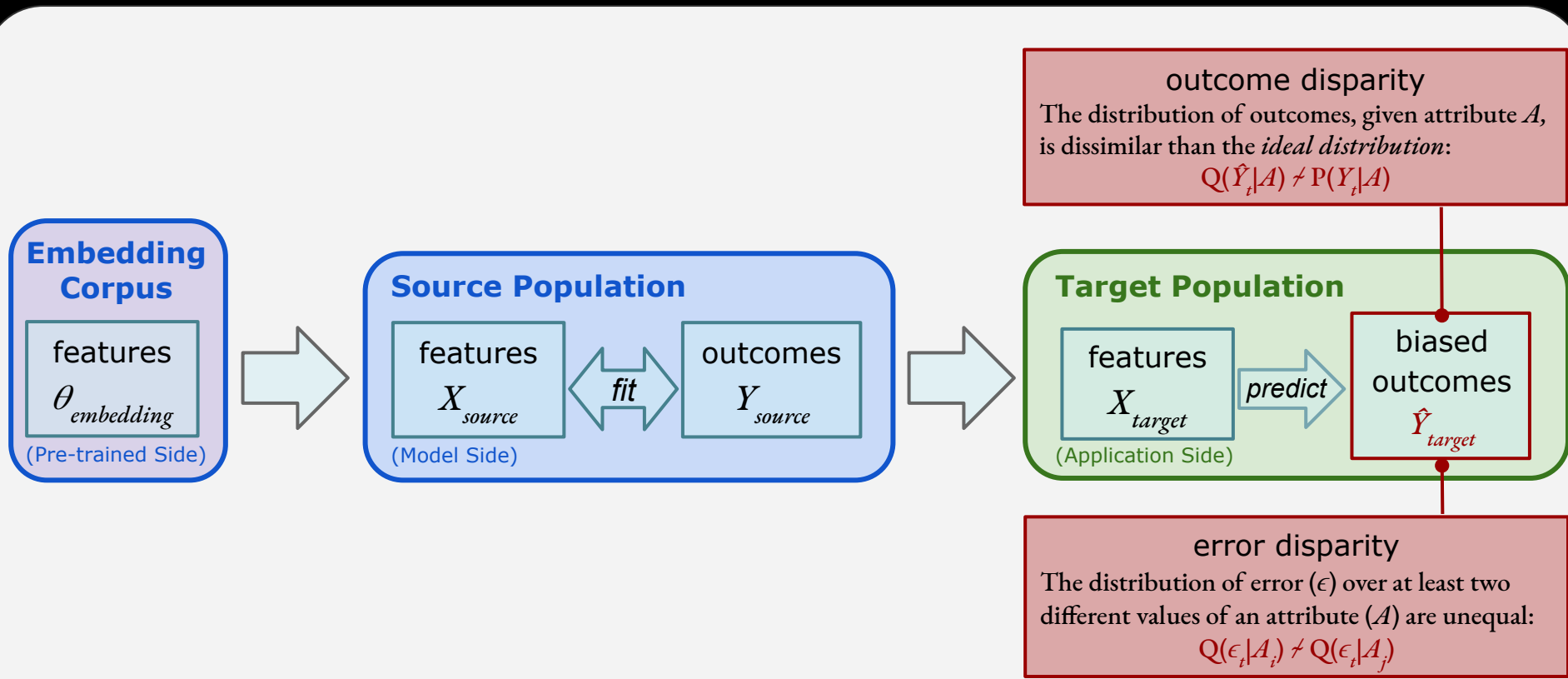
biased  
outcomes

$\hat{Y}_{target}$

### error disparity

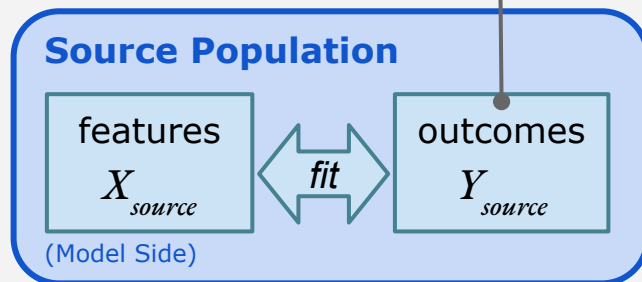
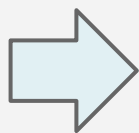
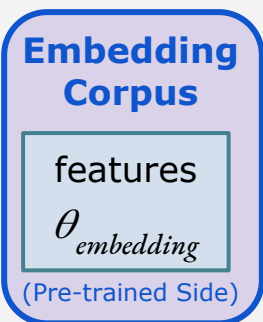
The distribution of error ( $\epsilon$ ) over at least two different values of an attribute ( $A$ ) are unequal:

$$Q(\epsilon_i|A_i) \neq Q(\epsilon_i|A_j)$$



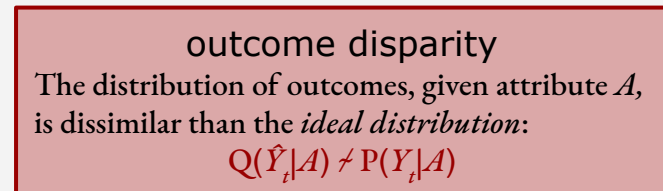
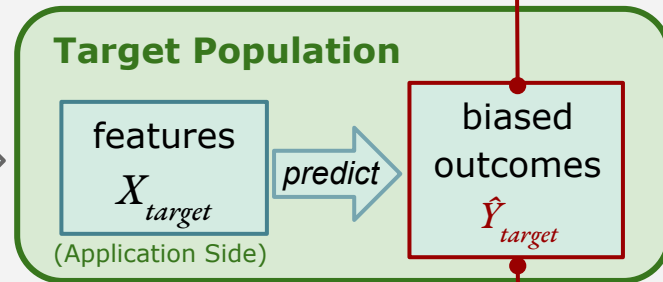
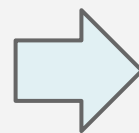
 **potential origin**

 **consequence**



**label bias**

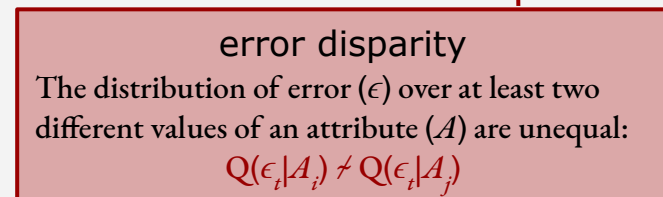
Biased annotations,  
interaction, or latent bias  
from past classifications.



**outcome disparity**

The distribution of outcomes, given attribute  $A$ ,  
is dissimilar than the *ideal distribution*:

$$Q(\hat{Y}_i|A) \neq P(Y_i|A)$$



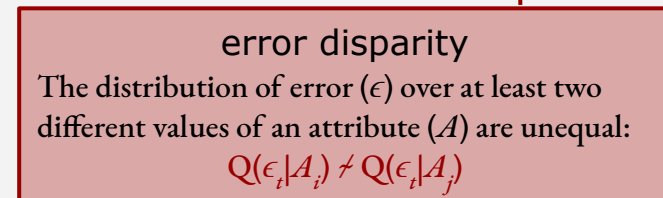
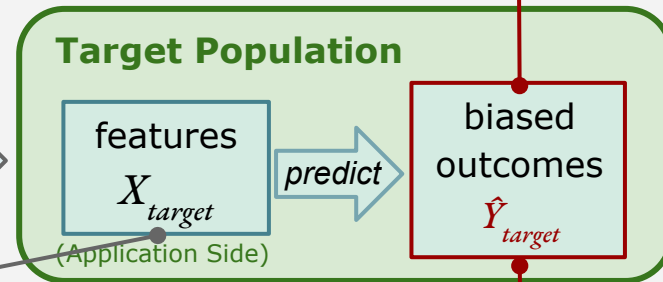
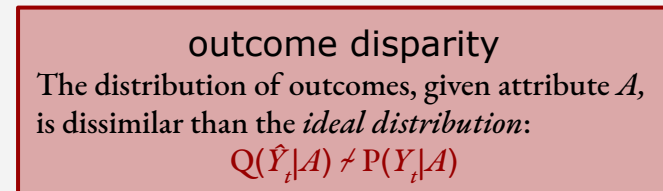
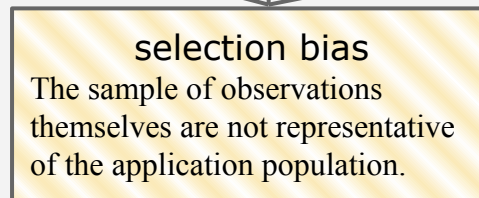
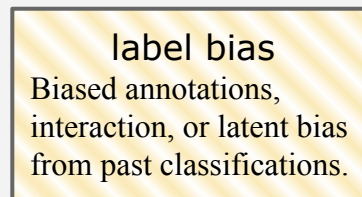
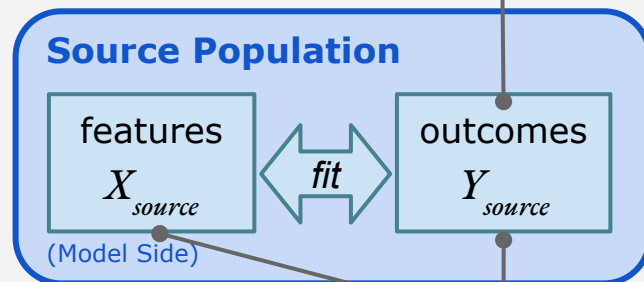
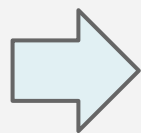
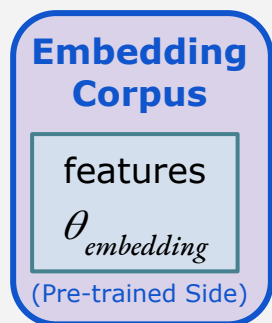
**error disparity**

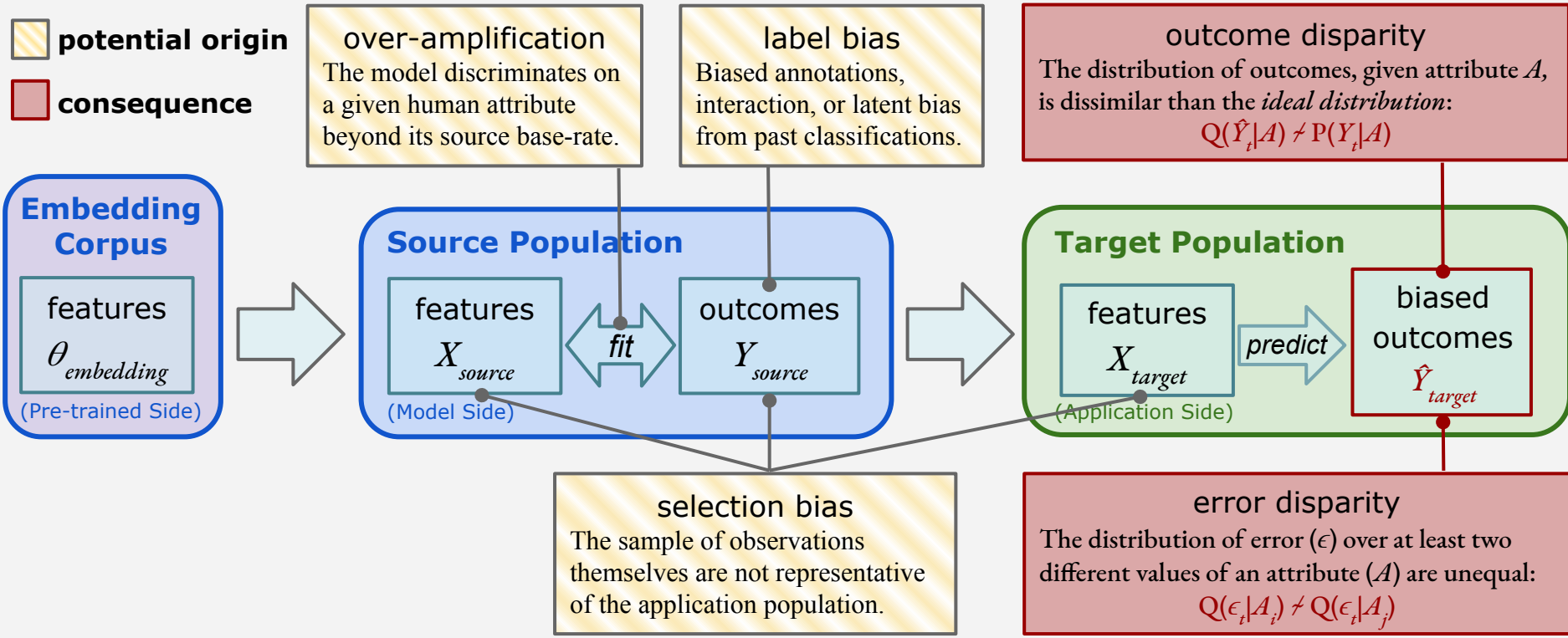
The distribution of error ( $\epsilon$ ) over at least two  
different values of an attribute ( $A$ ) are unequal:

$$Q(\epsilon_i|A_i) \neq Q(\epsilon_i|A_j)$$

 **potential origin**

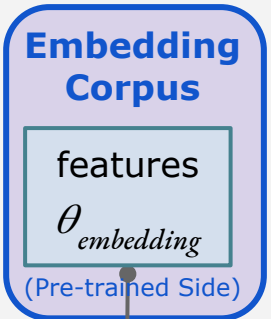
 **consequence**





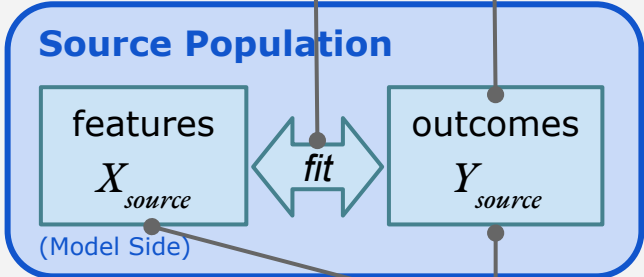


▨ **potential origin**  
■ **consequence**

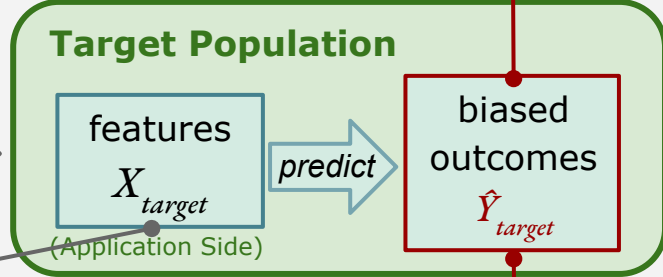


**over-amplification**  
 The model discriminates on a given human attribute beyond its source base-rate.

**label bias**  
 Biased annotations, interaction, or latent bias from past classifications.



**outcome disparity**  
 The distribution of outcomes, given attribute  $A$ , is dissimilar than the *ideal distribution*:  
 $Q(\hat{Y}_i|A) \neq P(Y_i|A)$



**semantic bias**  
 Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

**selection bias**  
 The sample of observations themselves are not representative of the application population.

**error disparity**  
 The distribution of error ( $\epsilon$ ) over at least two different values of an attribute ( $A$ ) are unequal:  
 $Q(\epsilon_i|A_i) \neq Q(\epsilon_i|A_j)$

E.g. Coreference resolution:  
connecting entities to references (i.e. pronouns).

*“The doctor told Mary that she had run some blood tests.”*

### semantic bias

Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

### selection bias

The sample of observations themselves are not representative of the application population.

### error disparity

The distribution of error ( $\epsilon$ ) over at least two different values of an attribute ( $A$ ) are unequal:

$$Q(\epsilon_t | A_i) \neq Q(\epsilon_t | A_j)$$

# Ethics in NLP

## Types of bias in NLP tasks:

- Outcome Disparity: Predicted distribution given A, are dissimilar from ideal distribution given A
  - Selection bias
  - Label bias
  - Over-amplification
- Error Disparity: Predicts less accurate for authors of given demographics.
- Semantic Bias: Representations of meaning store demographic associations.

# Ethics in NLP

## Types of bias in NLP tasks:

E.g. Coreference resolution: connecting entities to references (i.e. pronouns).

*“The doctor told Mary that she had run some blood tests.”*

- Semantic Bias: Representations of meaning store demographic associations.

# Ethics in NLP

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion



# Ethics in NLP

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion
- Mitigation strategies:
  - Informed consent -- let participants know
  - Do not share / secure storage
  - *Federated learning* -- separate and obfuscate to the point of preserving privacy
  - Transparency in information targeting  
“You are being shown this ad because ...”



# Ethics in NLP

Human Subjects Research

Observational versus Interventional

# Ethics in NLP

## Human Subjects Research

### Observational versus Interventional

(The Belmont Report, 1979)

- (i) Distinction of research from practice.
- (ii) Risk-Benefit criteria
- (iii) Appropriate selection of human subjects for participation in research
- (iv) Informed consent in various research settings.







*Natural language is generated by people.*

**What this means for NLP:**

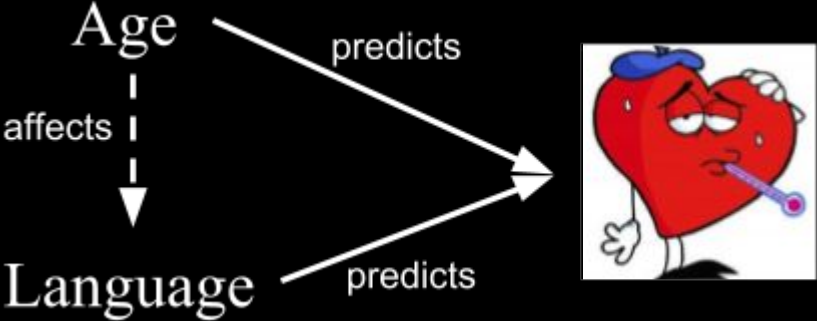
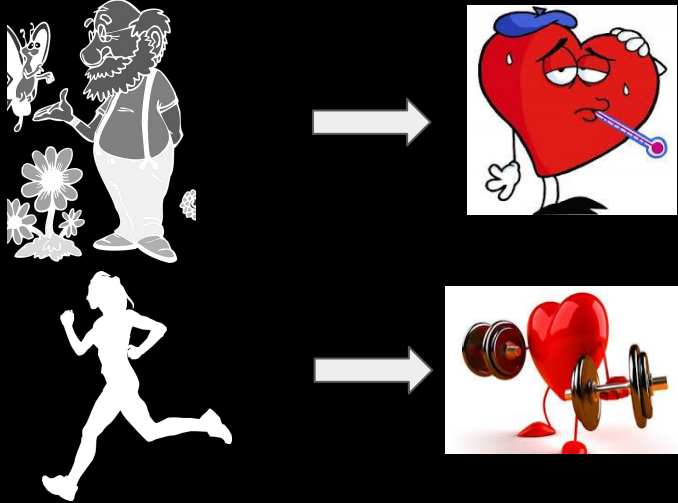
- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



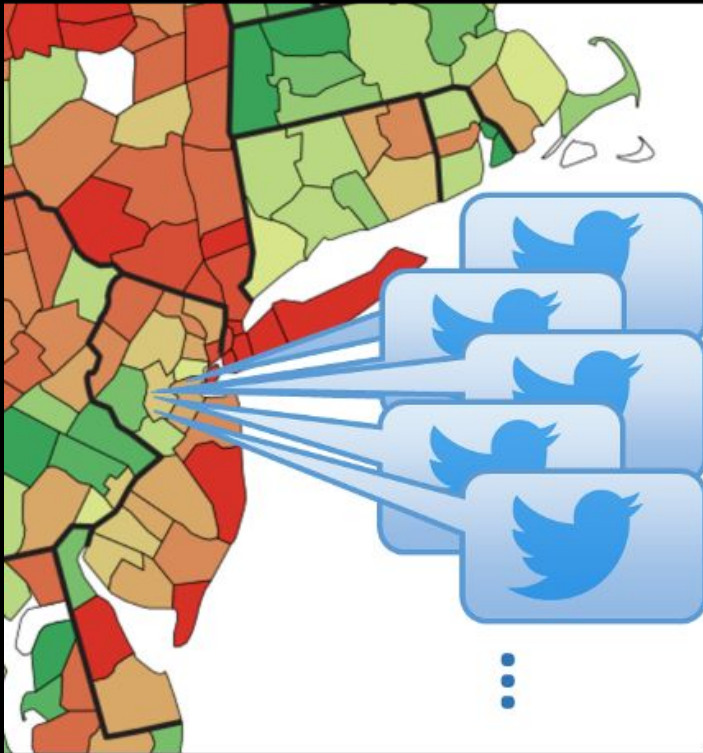
# Approaches to Human Factor Inclusion

1. **Adaptive:** Allow meaning of language to change depending on human context. (also called “compositional”)  
(e.g. “sick” said from a young individual versus old individual)
2. **Additive:** Include direct effect of human factor on outcome.  
(e.g. age and distinguishing PTSD from Depression)
3. **Bias Correction:** Optimize so as not to pick up on unwanted relationships.  
(e.g. image captioner label pictures of men in kitchen as women)

# Example 1: Individual Heart Disease



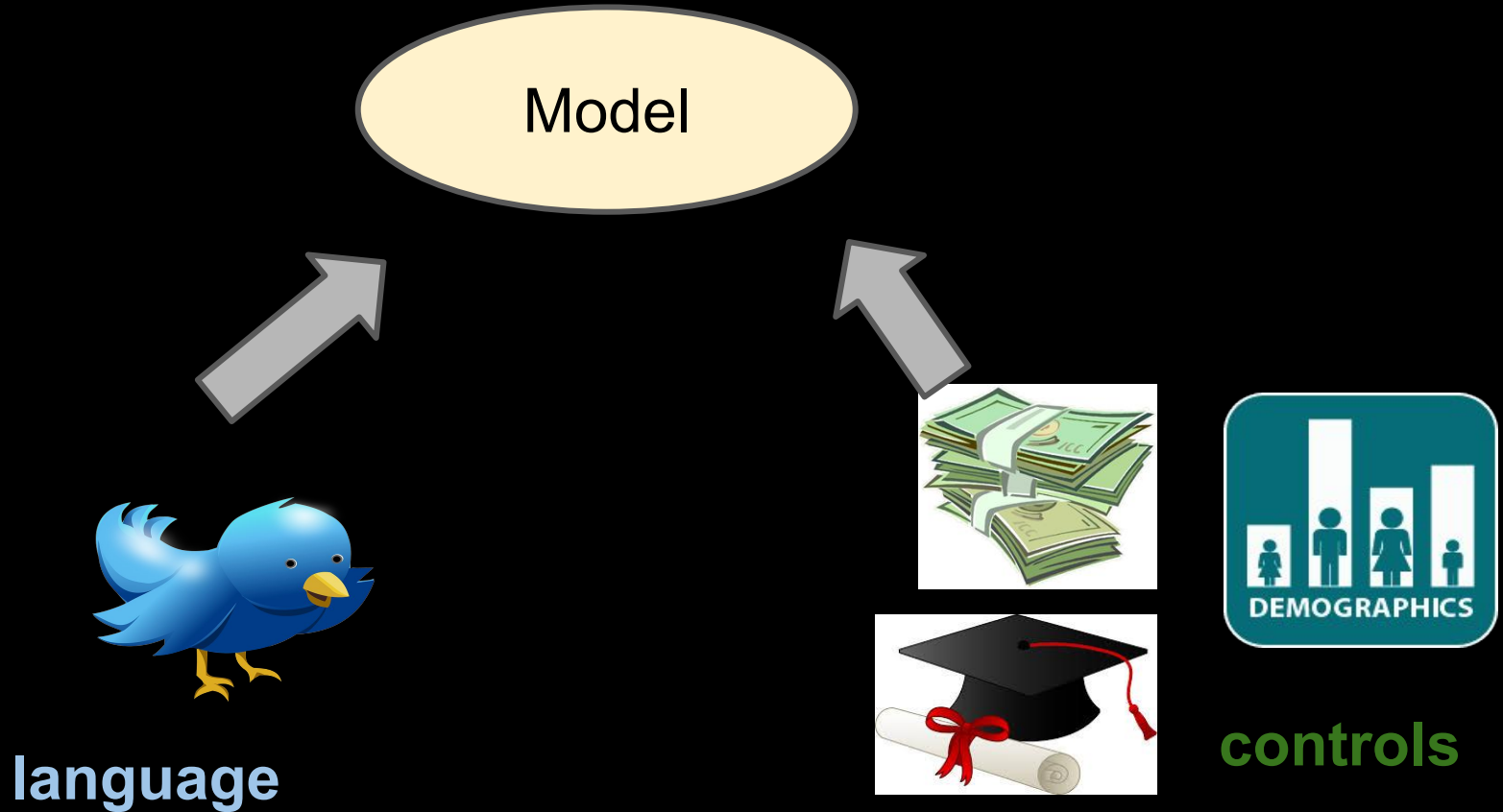
# Example 2: Twitter Language + Socioeconomics



0.0852
0.8794
0.1415
0.1996
0.4561
0.3556
0.7532
0.2703
0.6872
0.2623
0.3795
0.6451
0.2032
0.4075
0.5010
0.4783
0.9845
0.6314



# Additive (Residualized Control)



# Additive (Residualized Control)

## Challenges:

High-dimensional,  
sparse, and noisy.



language

few and  
well estimated



controls



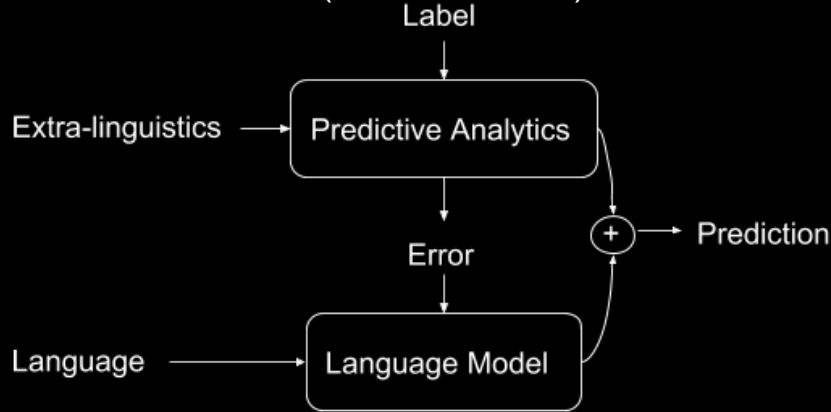
# Additive (Residualized Control)

Effectively use both low dimensional control features and high-dimensional, noisy language features:

1. **Train a control model** using the control values
2. **Calculate the residual error** and consider it as the new label
3. **Train a language model over the new labels**

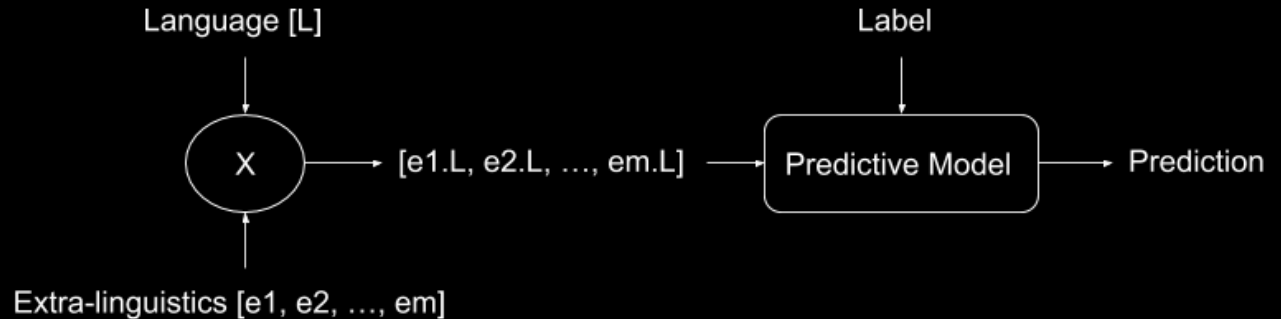
# Additive (Residualized Control)

Residualize control (additive model):



(Zamani et al., EACL 2017)

Adaptive model:



# Additive (Residualized Control)

Effectively use both low dimensional control features and high-dimensional, noisy language features:

1. **Train a control model** using the control values
2. **Calculate the residual error** and consider it as the new label
3. **Train a language model over the new labels**

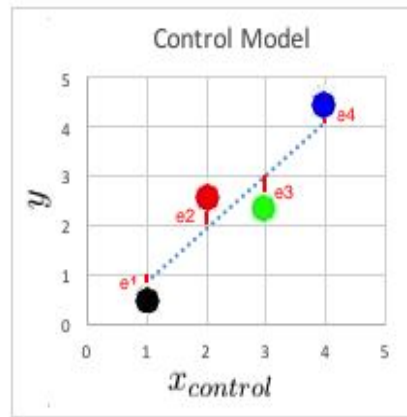
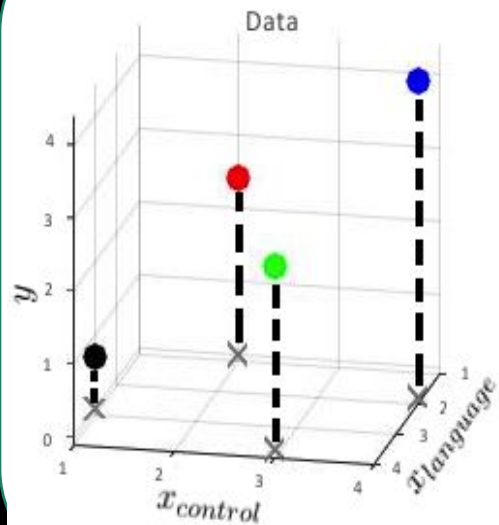
Model:

$$Y = \alpha x_1 + \beta x_2 + \gamma$$

Both learn the same linear model above, but

- Different learning algorithms per variable type.
- Different penalization methods

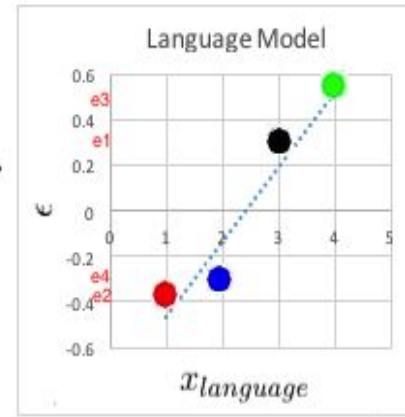
# Residualized Control Model



$$y = \alpha \times x_{control} + \beta + \epsilon$$

$$\alpha = 1, \beta = 0$$

error  
→



$$\epsilon \simeq \gamma \times x_{language} + \lambda$$

$$\gamma = 0.3, \lambda = -0.7$$

$$y \simeq \alpha \times x_{control} + \gamma \times x_{language} + \beta + \lambda$$

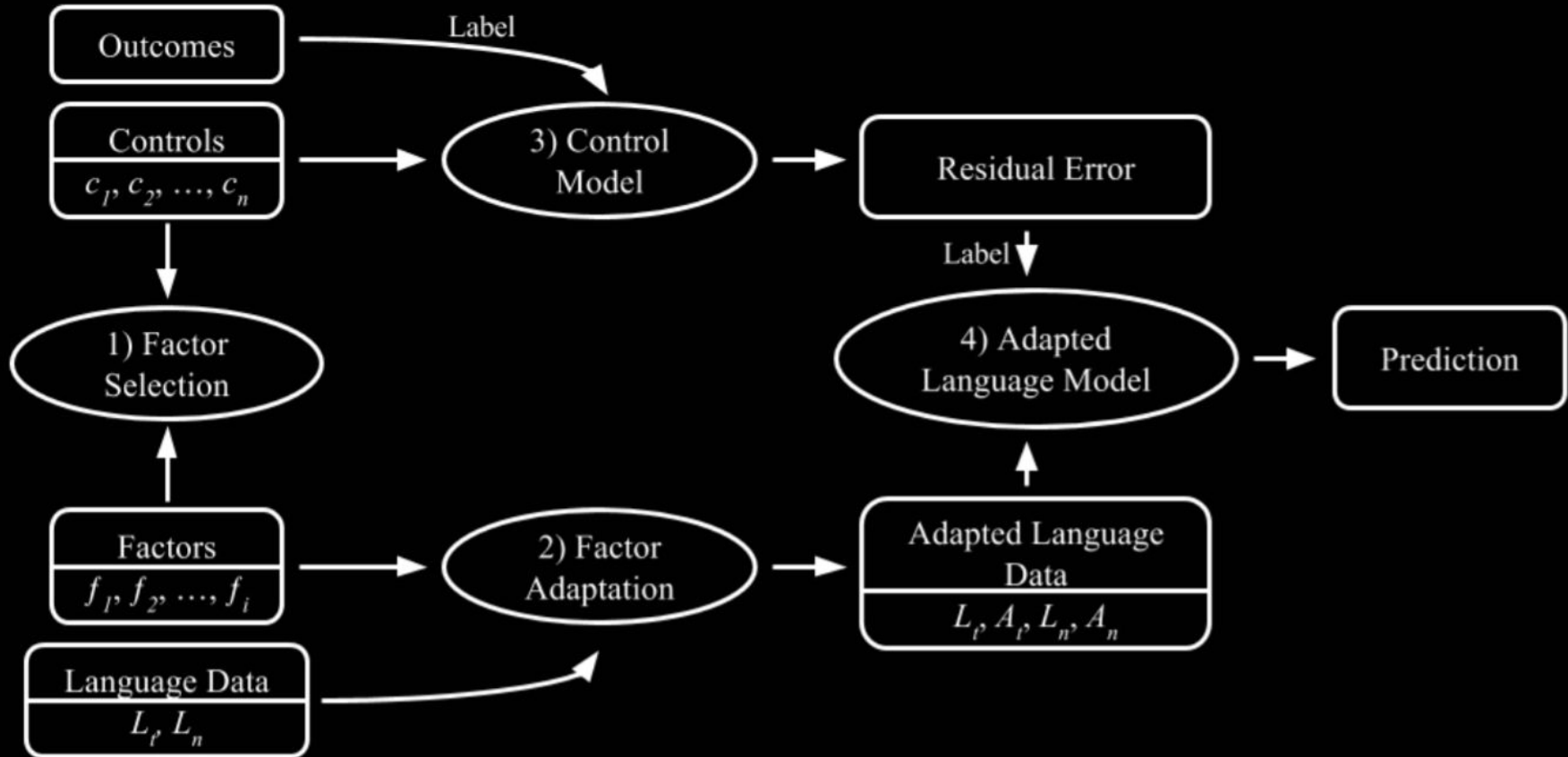
$$\Rightarrow \beta + \lambda = -0.7$$

# Combining Adaptive and Additive

Two Goals:

1. **Adaptive:** adapt to given human attributes  
(*user factor adaptation*;  
Lynn, Balasubramanian, Son, Kulkarni & Schwartz,  
*EMNLP 2017*)
2. **Additive:** predict beyond given attributes  
(*residualized control*; Zamani & Schwartz, *EACL 2017*)

# Solution: Residualized Factor Adaptation



# Results: County Health Predictions

	<b>Lang.</b>		
		<i>Controls Only</i>	<i>Added- Controls</i>
Heart Dis	0.585	0.514	0.608
Suicide	0.414	0.307	0.431
Poor Health	0.602	0.609	0.641
Life Satis.	0.209	0.329	0.335
Avg.	0.453	0.440	0.503



# Results: County Health Predictions

	<b>Lang.</b>	<b>All Factors</b>		
		<i>Controls Only</i>	<i>Added- Controls</i>	<i>Res- Control</i>
Heart Dis	0.585	0.514	0.608	0.628
Suicide	0.414	0.307	0.431	0.460
Poor Health	0.602	0.609	0.641	0.661
Life Satis.	0.209	0.329	0.335	0.372
Avg.	0.453	0.440	0.503	0.530

# Results: County Health Predictions

	<b>Lang.</b>	<b>All Factors</b>			
		<i>Controls Only</i>	<i>Added- Controls</i>	<i>Res- Control</i>	<i>FA</i>
Heart Dis	0.585	0.514	0.608	0.628	0.635
Suicide	0.414	0.307	0.431	0.460	0.494
Poor Health	0.602	0.609	0.641	0.661	0.674
Life Satis.	0.209	0.329	0.335	0.372	0.352
Avg.	0.453	0.440	0.503	0.530	0.539

# Results: County Health Predictions

	<b>Lang.</b>	<b>All Factors</b>				
		<i>Controls Only</i>	<i>Added-Controls</i>	<i>Res-Control</i>	<i>FA</i>	<i>RFA</i>
Heart Dis	0.585	0.514	0.608	0.628	0.635	<b>0.655</b>
Suicide	0.414	0.307	0.431	0.460	0.494	<b>0.510</b>
Poor Health	0.602	0.609	0.641	0.661	0.674	0.682
Life Satis.	0.209	0.329	0.335	0.372	0.352	<b>0.396</b>
Avg.	0.453	0.440	0.503	0.530	0.539	0.560

# Results: County Health Predictions

	<b>Lang.</b>	<b>All Factors</b>				
		<i>Controls Only</i>	<i>Added-Controls</i>	<i>Res-Control</i>	<i>FA</i>	<i>RFA</i>
Heart Dis	0.585	0.514	0.608	0.628	0.635	<b>0.655</b>
Suicide	0.414	0.307	0.431	0.460	0.494	<b>0.510</b>
Poor Health	0.602	0.609	0.641	0.661	0.674	0.682
Life Satis.	0.209	0.329	0.335	0.372	0.352	<b>0.396</b>
Avg.	0.453	0.440	0.503	0.530	0.539	0.560

variance explained ( $R^2$ )

*Natural language is generated by people.*

**What this means for NLP:**

- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



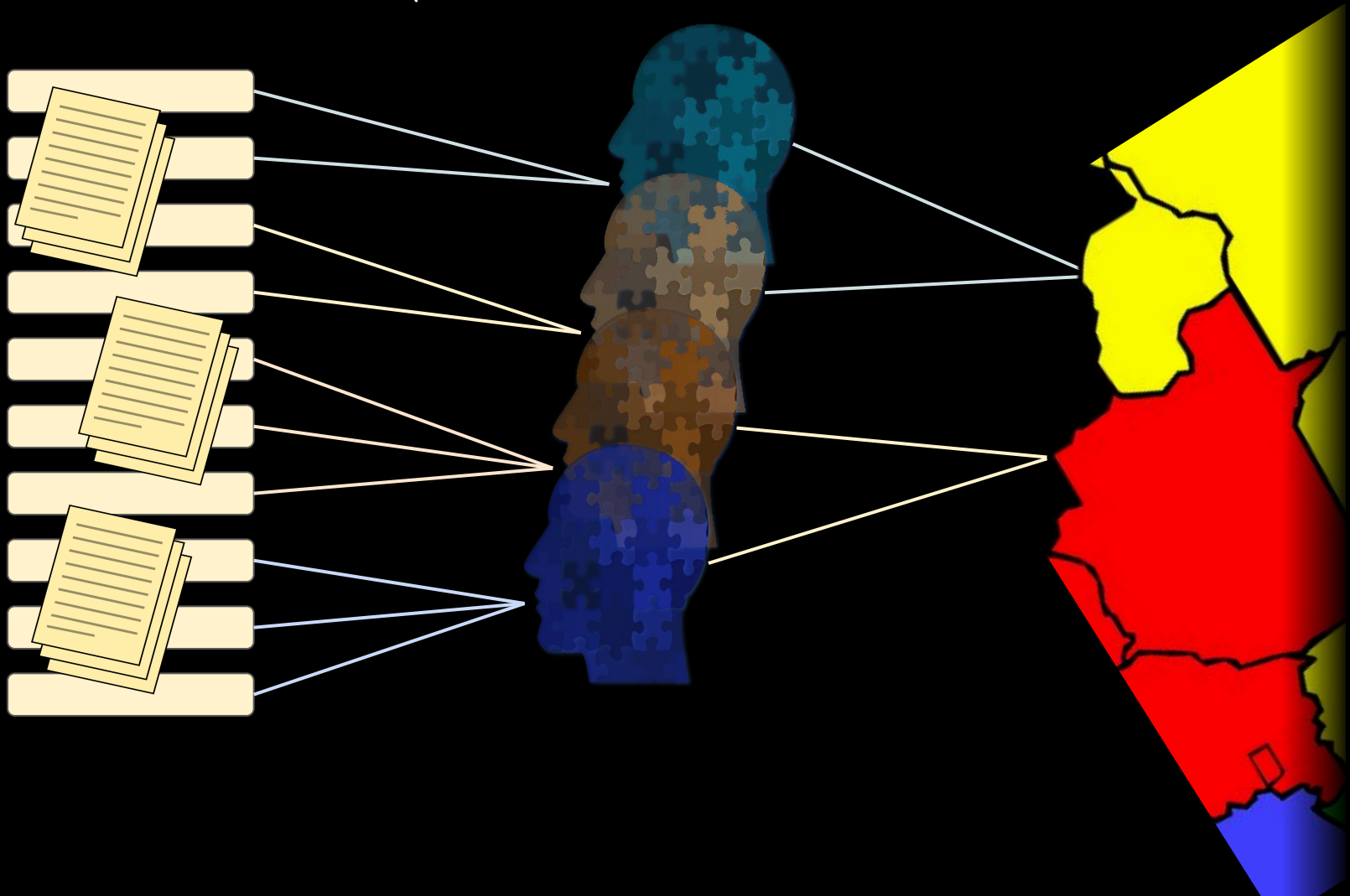
*Natural language is generated by people.*

**What this means for NLP:**

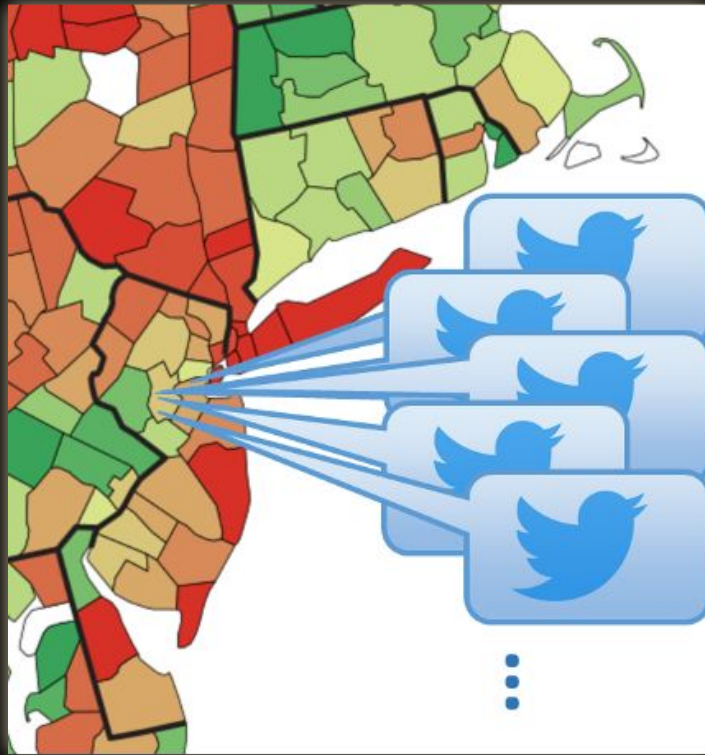
- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



*Data are inherently multi-level.*



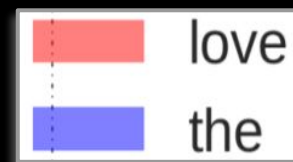
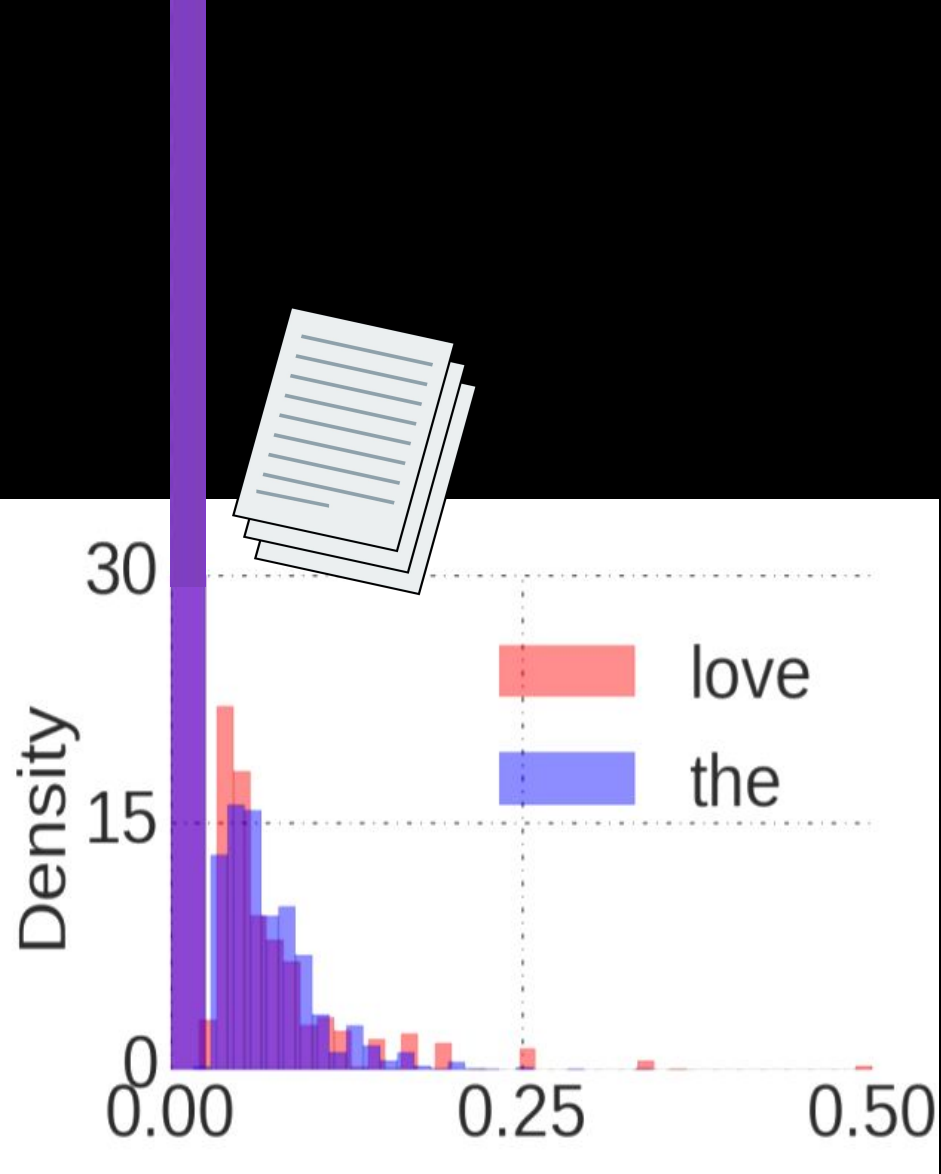
1,639,750 tweets from 5,226 users in 420 counties



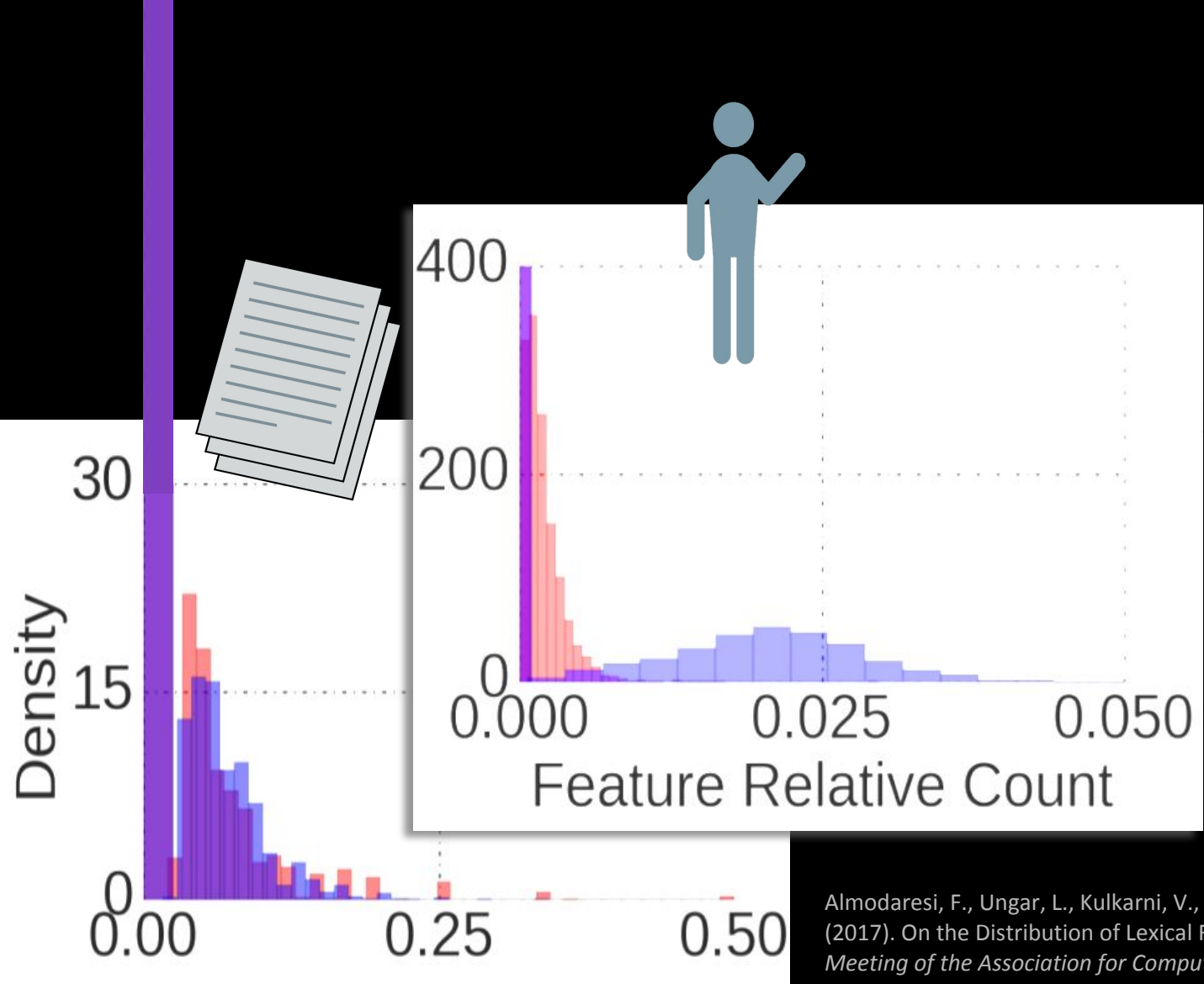
0.0852
0.8794
0.1415
0.1996
0.4561
0.3556
0.7532
0.2703
0.6872
0.2623
0.3795
0.6451
0.2032
0.4075
0.5010
0.4783
0.9845
0.6314



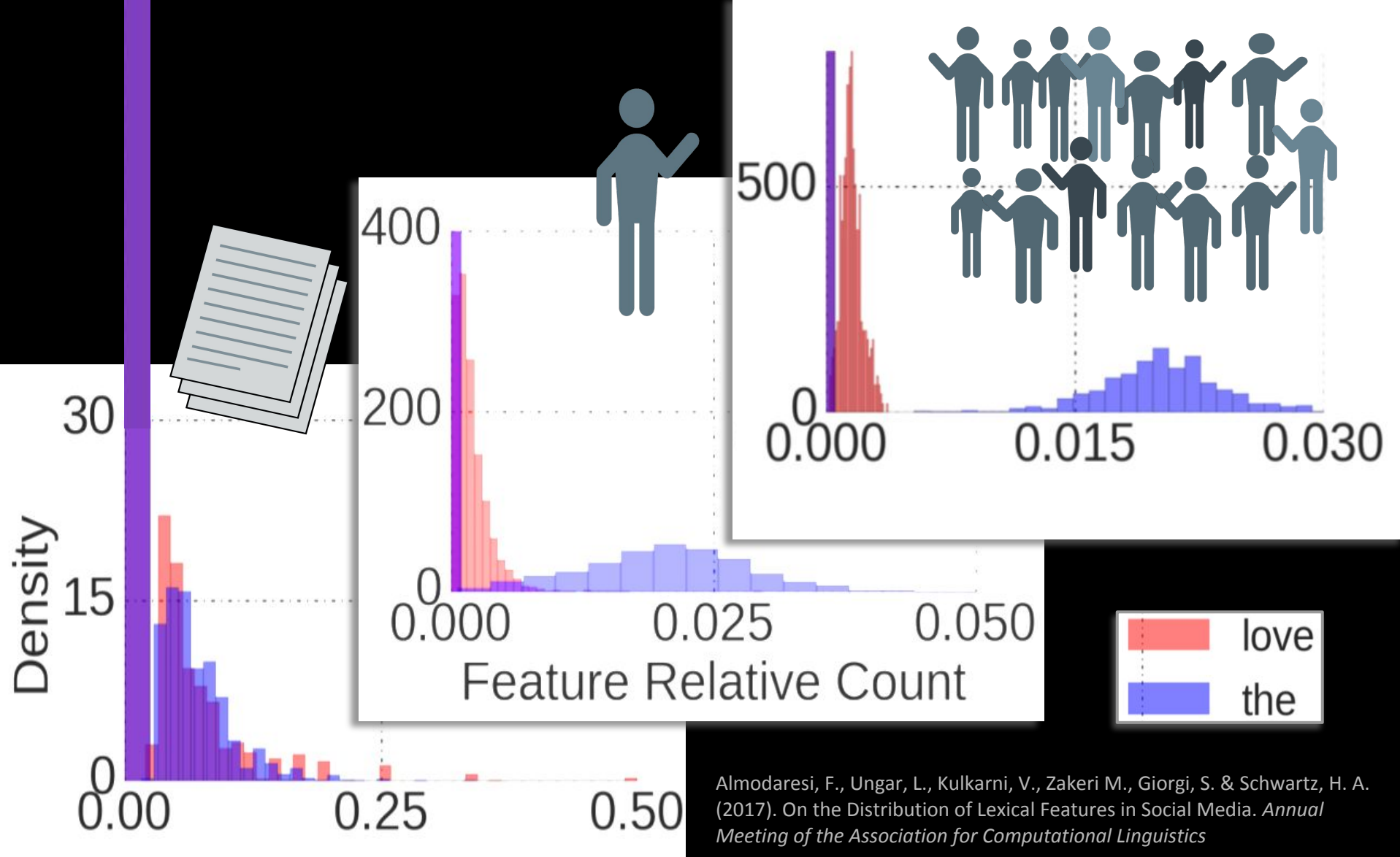




Almodaresi, F., Ungar, L., Kulkarni, V., Zakeri M., Giorgi, S. & Schwartz, H. A. (2017). On the Distribution of Lexical Features in Social Media. *Annual Meeting of the Association for Computational Linguistics*



Almodaresi, F., Ungar, L., Kulkarni, V., Zakeri M., Giorgi, S. & Schwartz, H. A. (2017). On the Distribution of Lexical Features in Social Media. *Annual Meeting of the Association for Computational Linguistics*



Almodaresi, F., Ungar, L., Kulkarni, V., Zakeri M., Giorgi, S. & Schwartz, H. A. (2017). On the Distribution of Lexical Features in Social Media. *Annual Meeting of the Association for Computational Linguistics*

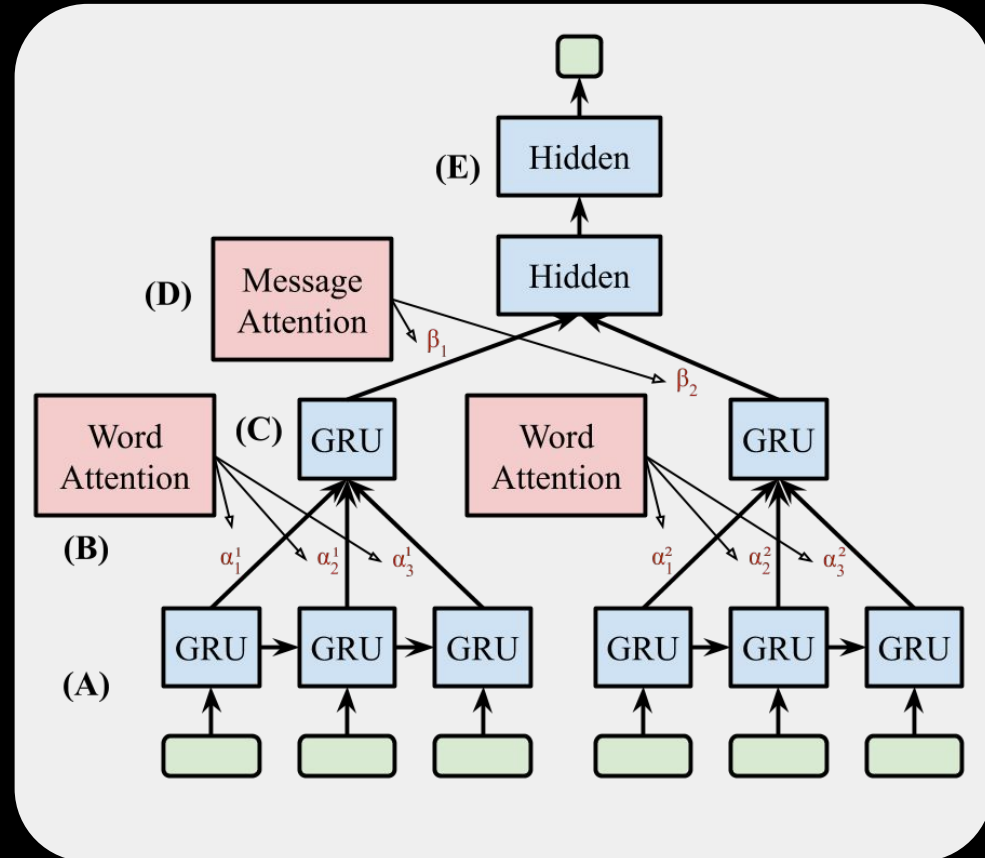
*Data are inherently multi-level.*

Distribution	Message			User			County		
	1-gram	topic	Lex.	1-gram	topic	Lex.	1-gram	topic	Lex.
Power Law	<b>.71</b>	.10	.00	.04	.00	.00	.07	.00	.00
Log-Normal	.25	<b>.89</b>	<b>1.00</b>	<b>.96</b>	<b>.97</b>	<b>.64</b>	<b>.92</b>	<b>.86</b>	.44
Normal	.04	.01	.00	.00	.03	.36	.01	.14	<b>.56</b>

Proportion best fit by the given distribution.

Data are inherently multi-level.

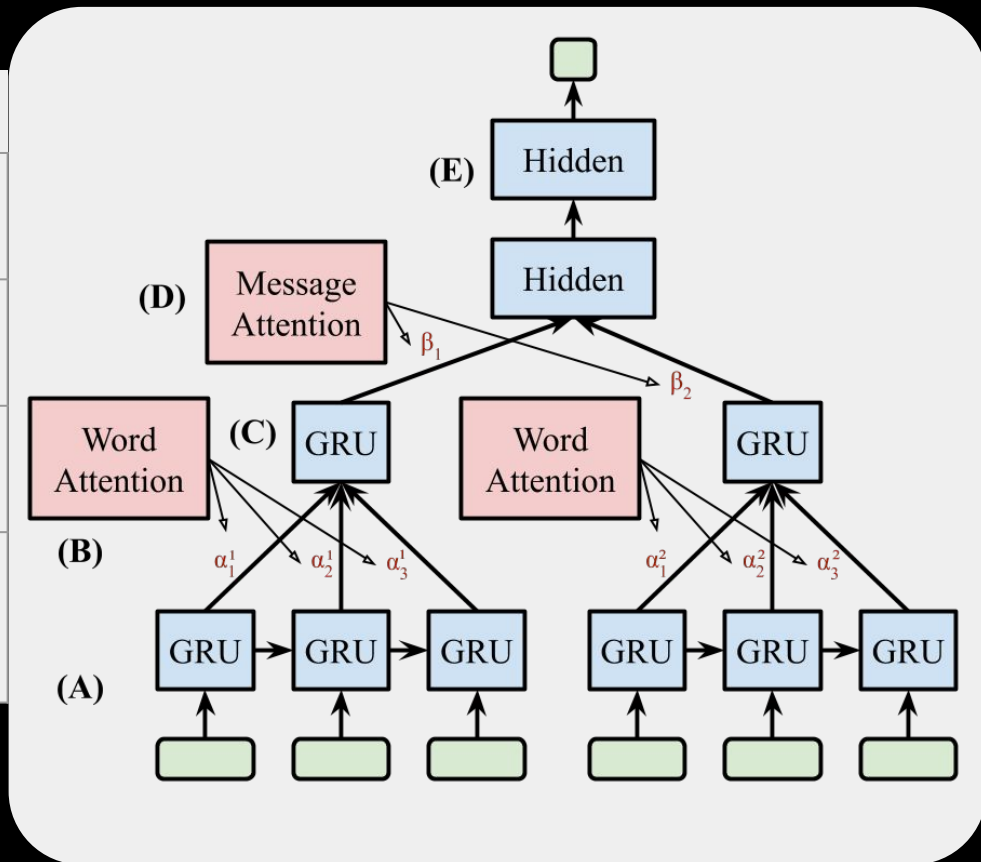
## Multi-level Attention and Sequence Model



Data are inherently multi-level.

	<i>d</i>	<i>O</i>	<i>C</i>	<i>E</i>	<i>A</i>	<i>N</i>
<b>BERT + DAN</b>	768	.60	.51	.54	.51	.52
<b>(Park et al., 2015)</b>	5106	.63	.52	.56	.54	.53
<b>Multi-level Attention</b>	200	.63	.52	.55	.51	.54
<b>Multi-level Attention + Ridge</b>	5306	<b>.66</b>	<b>.54</b>	<b>.58</b>	<b>.56</b>	<b>.56</b>

bold:  $p < .01$

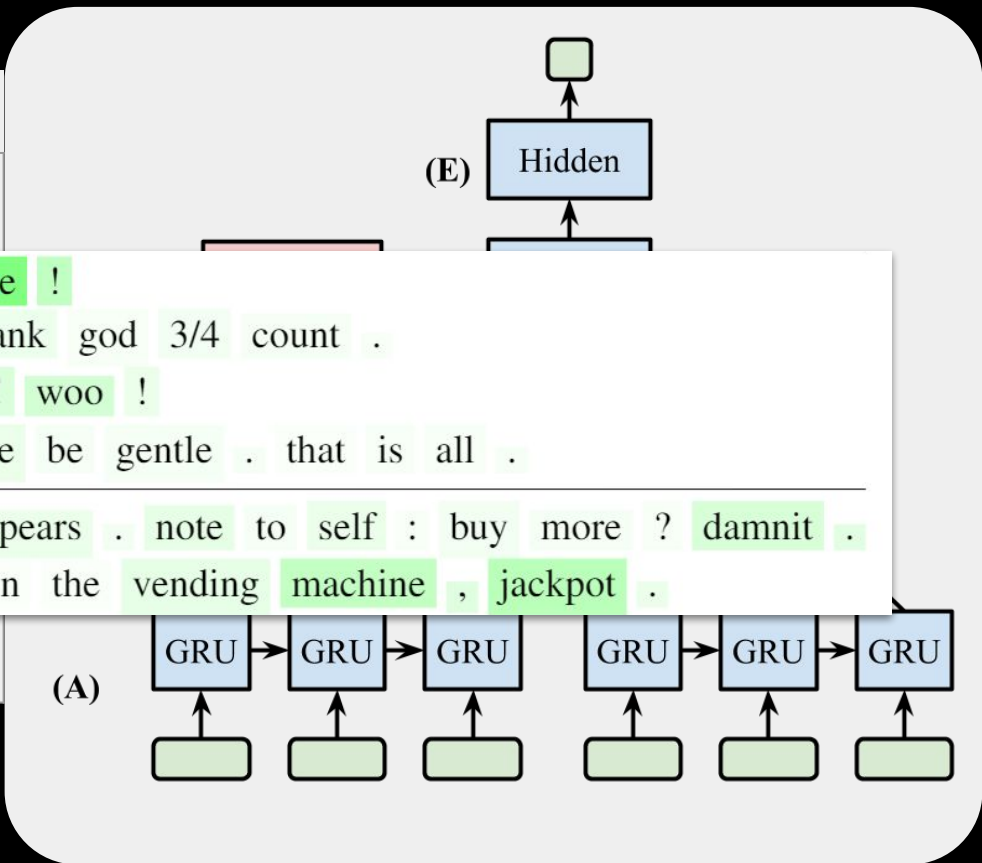


Lynn, V., Balasubramanian, N., & Schwartz, H. A. (2020). Hierarchical Modeling for User Personality Prediction: The Role of Message-Level Attention. *Annual Meeting of the Association for Computational Linguistics*.

# Data are inherently multi-level.

	<i>d</i>	<i>O</i>	<i>C</i>	<i>E</i>	<i>A</i>	<i>N</i>
<b>BERT + DAN</b>	768	.60	.51	.54	.51	.52
<b>al.:</b>						
<i>High</i>						
<i>CON</i>						
<b>Mult</b>						
<b>Att</b>						
<b>Mult</b>						
<b>Attention</b>						
<b>+ Ridge</b>	5306	.66	.54	.56	.56	.56

bold:  $p < .01$



# Differential Language Analysis

Input:

Linguistic features

Human or community attribute

Output:

Features distinguishing attribute

Goal: Data-driven insights about an attribute





# Differential Language Analysis

Input:

Linguistic features

Human or community attribute

Output:

Features distinguishing attribute

Goal: Data-driven insights about an attribute



# Differential Language Analysis

Methods of Correlation Analysis:

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

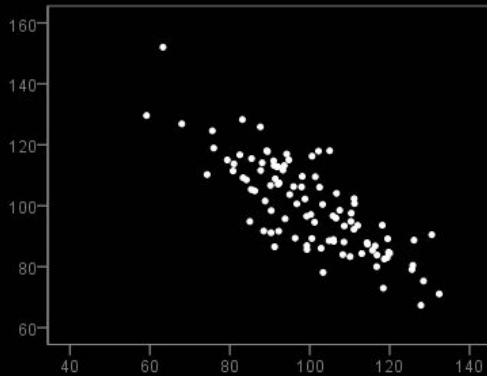
$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

# Differential Language Analysis

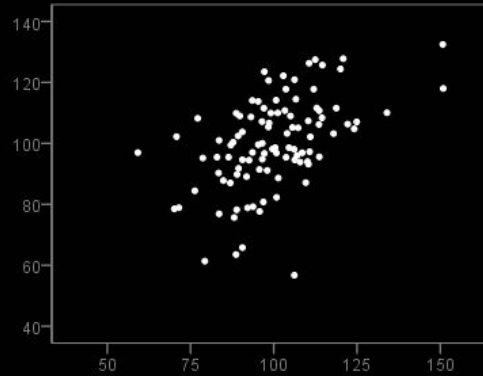
Methods of Correlation Analysis:

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

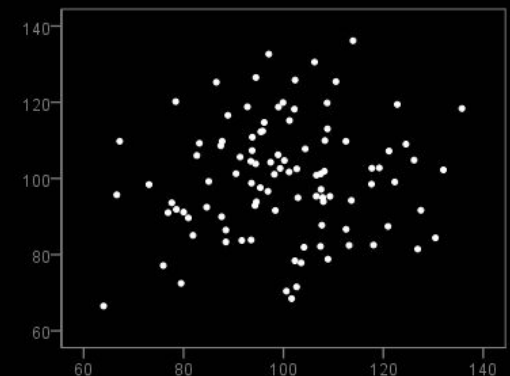
$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$



**r = -0.8**



**r = 0.5** © 2017 www.sj



**r = 0.1**

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- **Standardized** Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Adjust all variables to have “mean center” and “unit variance”:

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- **Standardized** Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Adjust all variables to have “mean center” and “unit variance”:

$$z = \frac{x - \mu}{\sigma}$$

$\mu$  = Mean

$\sigma$  = Standard Deviation



# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

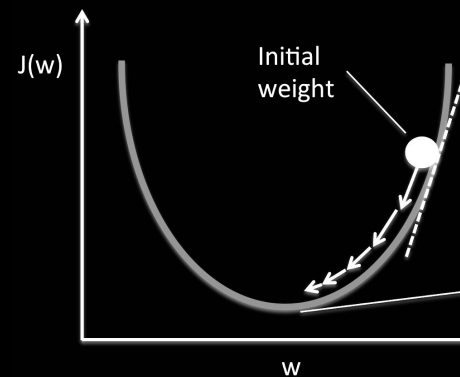
- Standardized Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Option 1: Gradient Descent:

$$J = \sum (y - \hat{y})^2 \text{ -- "Sum of Squares" Error}$$



# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression

Fit the model:

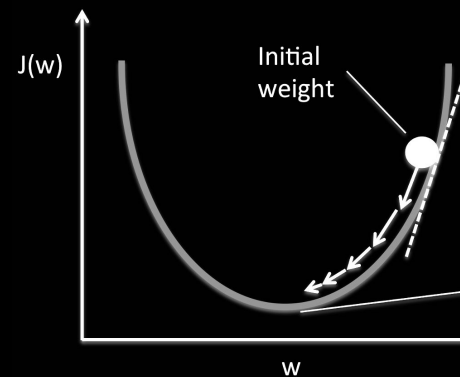
$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Option 1: Gradient Descent:

$$J = \sum (y - \hat{y})^2 \text{ -- "Sum of Squares" Error}$$

Option 2: Matrix model:

$$Y = X\beta + \epsilon$$



# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Option 1: Gradient Descent:

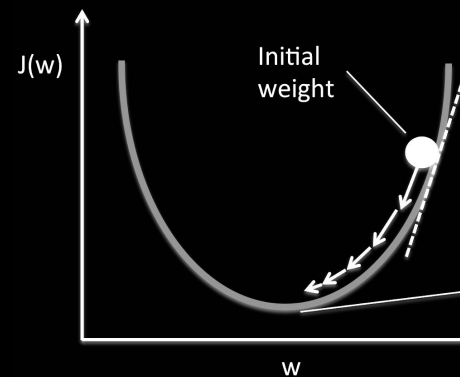
$$J = \sum (y - \hat{y})^2 \text{ -- "Sum of Squares" Error}$$

Option 2: Matrix model:

$$Y = X\beta + \epsilon$$

Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation  
Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression

Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im} + \epsilon_i$$

Option 1: Gradient Descent:

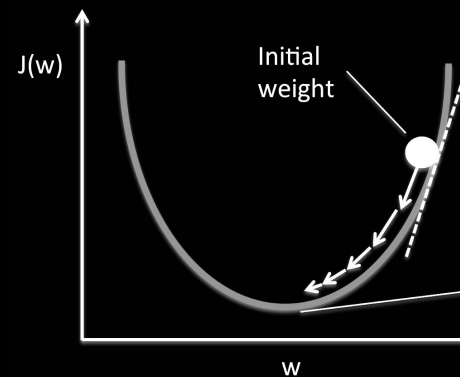
$$J = \sum (y - \hat{y})^2 \text{ -- "Sum of Squares" Error}$$

Option 2: Matrix model:

$$Y = X\beta + \epsilon$$

Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio

$$\frac{\frac{\text{countA}(\text{"horrible"})}{NA}}{1 - \frac{\text{countA}(\text{"horrible"})}{NA}}$$

---

$$\frac{\frac{\text{countB}(\text{"horrible"})}{NB}}{1 - \frac{\text{countB}(\text{"horrible"})}{NB}}$$

# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio

$$\frac{\frac{\frac{\text{countA}(\text{"horrible"})}{NA}}{1 - \frac{\text{countA}(\text{"horrible"})}{NA}}}{\frac{\frac{\text{countB}(\text{"horrible"})}{NB}}{1 - \frac{\text{countB}(\text{"horrible"})}{NB}}} \propto \log \left( \frac{\frac{\text{countA}(\text{"horrible"})}{NA}}{1 - \frac{\text{countA}(\text{"horrible"})}{NA}} \right) - \log \left( \frac{\frac{\text{countB}(\text{"horrible"})}{NB}}{1 - \frac{\text{countB}(\text{"horrible"})}{NB}} \right)$$
$$= \log \left( \frac{\text{countA}(\text{"horrible"})}{NA - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{NB - \text{countB}(\text{"horrible"})} \right)$$

# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio using Informative Dirichlet Prior  $\log \left( \frac{\text{countA}(\text{"horrible"})}{NA - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{NB - \text{countB}(\text{"horrible"})} \right)$

$$\hat{\delta}_w^{(i-j)} = \log \left( \frac{y_w^i + \alpha_w}{n^i + \alpha_0 - (y_w^i + \alpha_w)} \right) - \log \left( \frac{y_w^j + \alpha_w}{n^j + \alpha_0 - (y_w^j + \alpha_w)} \right)$$

# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio using Informative Dirichlet Prior  $\log \left( \frac{\text{countA}(\text{"horrible"})}{NA - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{NB - \text{countB}(\text{"horrible"})} \right)$

$$\hat{\delta}_w^{(i-j)} = \log \left( \frac{y_w^i + \alpha_w}{n^i + \alpha_0 - (y_w^i + \alpha_w)} \right) - \log \left( \frac{y_w^j + \alpha_w}{n^j + \alpha_0 - (y_w^j + \alpha_w)} \right)$$



# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio using Informative Dirichlet Prior

$$\log \left( \frac{\text{countA}(\text{"horrible"})}{NA - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{NB - \text{countB}(\text{"horrible"})} \right)$$

$$\hat{\delta}_w^{(i-j)} = \log \left( \frac{y_w^i + \alpha_w}{n^i + \alpha_0 + (y_w^i + \alpha_w)} \right) - \log \left( \frac{y_w^j + \alpha_w}{n^j + \alpha_0 + (y_w^j + \alpha_w)} \right)$$

Bayesian term for “smoothing”: accounts for uncertainty as a function of less events (i.e. words observed less) by integrating “prior” beliefs mathematically.

# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio using Informative Dirichlet Prior

$$\log \left( \frac{\text{countA}(\text{"horrible"})}{NA - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{NB - \text{countB}(\text{"horrible"})} \right)$$

$$\hat{\delta}_w^{(i-j)} = \log \left( \frac{y_w^i + \alpha_w}{n^i + \alpha_0 + (y_w^i + \alpha_w)} \right) - \log \left( \frac{y_w^j + \alpha_w}{n^j + \alpha_0 + (y_w^j + \alpha_w)} \right)$$

Bayesian term for “smoothing”: accounts for uncertainty as a function of less events (i.e. words observed less) by integrating “prior” beliefs mathematically.

“Informative”: the prior is based on past evidence. Here, the total frequency of the word.

# Differential Language Analysis

Methods of “Correlation” Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio using Informative Dirichlet Prior  $\log \left( \frac{\text{countA}(\text{"horrible"})}{\text{NA} - \text{countA}(\text{"horrible"})} \right) - \log \left( \frac{\text{countB}(\text{"horrible"})}{\text{NB} - \text{countB}(\text{"horrible"})} \right)$

$$\hat{\delta}_w^{(i-j)} = \log \left( \frac{y_w^i + \alpha_w}{n^i + \alpha_0 - (y_w^i + \alpha_w)} \right) - \log \left( \frac{y_w^j + \alpha_w}{n^j + \alpha_0 - (y_w^j + \alpha_w)} \right)$$

( $n^i$  is the size of corpus  $i$ ,  $n^j$  is the size of corpus  $j$ ,  $y_w^i$  is the count of word  $w$  in corpus  $i$ ,  $y_w^j$  is the count of word  $w$  in corpus  $j$ ,  $\alpha_0$  is the size of the background corpus, and  $\alpha_w$  is the count of word  $w$  in the background corpus.)

$$\sigma^2 \left( \hat{\delta}_w^{(i-j)} \right) \approx \frac{1}{y_w^i + \alpha_w} + \frac{1}{y_w^j + \alpha_w}$$

- Final statistic for a word: z-score of its log-odds-ratio:

$$\frac{\hat{\delta}_w^{(i-j)}}{\sqrt{\sigma^2 \left( \hat{\delta}_w^{(i-j)} \right)}}$$

(Monroe et al., 2010; Jurafsky, 2017)

*Natural language is generated by people.*

**What this means for NLP:**

- 1. Our data are inherently multi-level.*
- 2. Often, there are "already-available" human attributes.*
- 3. Our data and models are (human) biased.*



*Natural language is generated by people.*

## What this means for NLP:

### **Practical implication**

- 1) More accurate models
- 2) Increased fairness in applications

*Considering the people behind the language not only offers opportunities for improved accuracies but it could be fundamental to NLP's role in our increasingly digital world.*

