

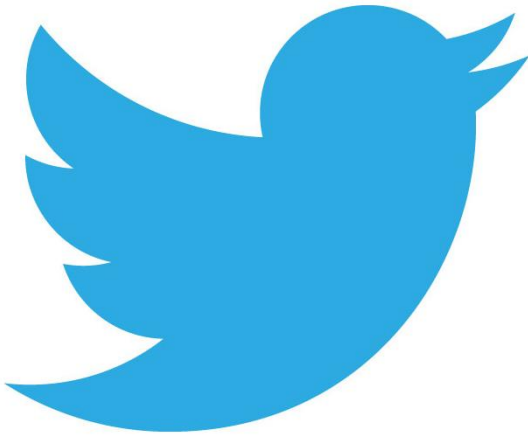
# **Personality, Gender and Age in the Language of Social Media**

Crowdsourcing and Human Computation  
12/4/2013

H. Andrew Schwartz

Core Collaborators:

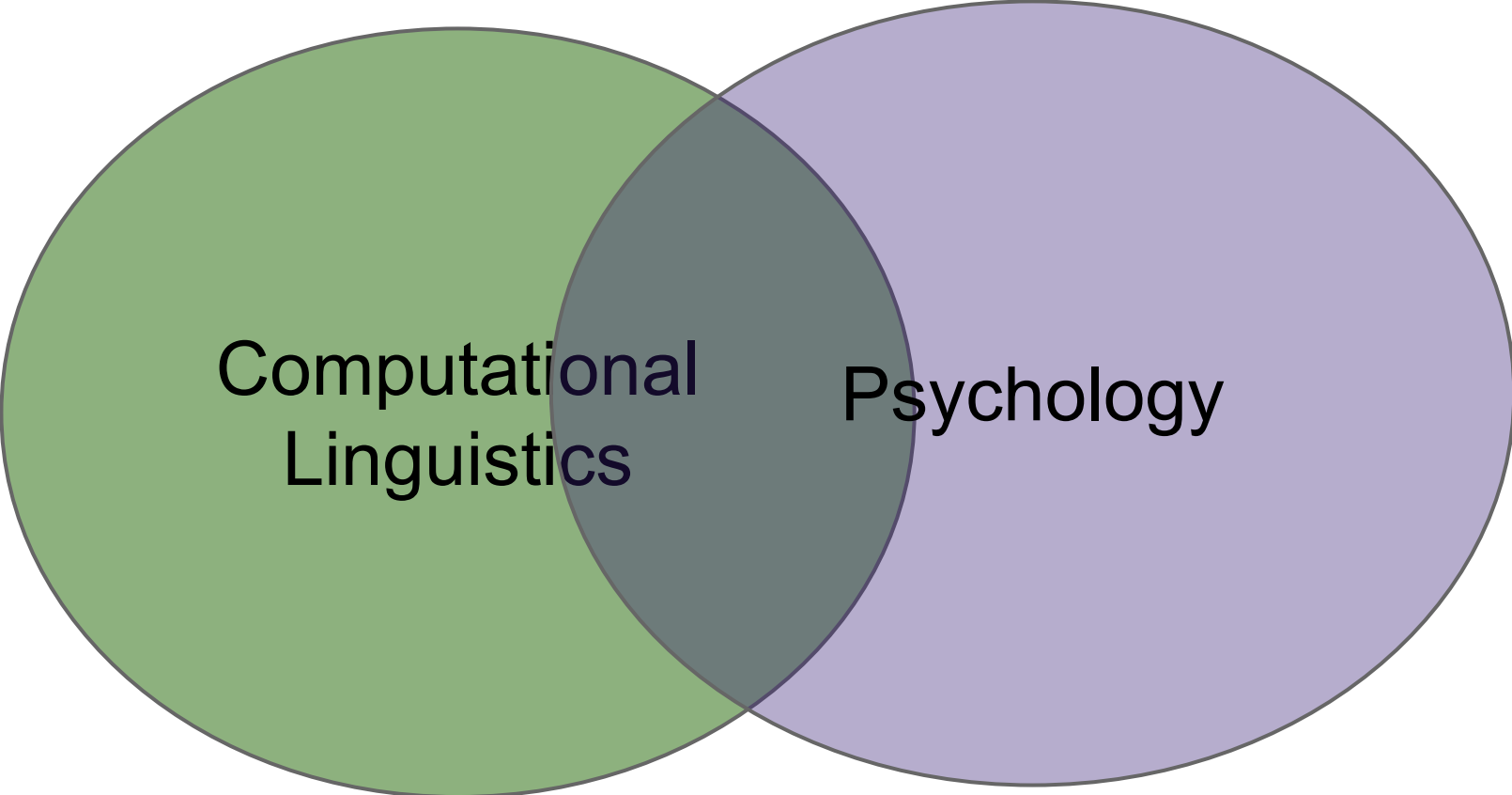
Johannes Eichstaedt, Margaret Kern, Gregory Park, Martin Seligman Lyle Ungar



- > 150 million active monthly users
- > 350 million messages a day



- > 1 billion active monthly users
- > 4 billion messages a day



Computational  
Linguistics

Psychology

**Computational Linguistics traditionally is focused on modeling and understanding language.**

- **Results: accurate predictions**

**Computational linguistics is traditionally focused on modeling and understanding language.**

- **Results: accurate predictions**

**Our Goal: Understand the people behind the Language.**

- **Results: Novel Insights**

**Computational linguistics is traditionally focused on modeling and understanding language.**

- **Results: accurate predictions**

**Our Goal: Understand the people behind the Language.**

- **Results: Novel Insights**

**Passive crowdsourcing the human condition?**

**Psychological language analyses typically limited to  
apriori language lexica**

LWC Category	Gender		Age		Extraversion		Agreeableness		Conscientious.		Neuroticism		Openness	
	[34] <i>d</i>	our $\beta$	[30] $\beta$	our $\beta$	[27] $\rho$	our $\beta$	[27] $\rho$	our $\beta$	[27] $\rho$	our $\beta$	[27] $\rho$	our $\beta$	[27] $\rho$	our $\beta$
Total function words	-	-0.04	-	0.16	-	-0.04	-	0.02	-	0.02	-	0.03	-	0.09
Total pronouns	0.36	0.07	-	-0.02	ns	ns	0.11	ns	ns	-0.03	ns	0.04	-0.21	0.07
Personal pronouns	-	0.14	-	-0.08	-	ns	-	ns	-	-0.04	-	0.04	-	0.05
1st pers singular	0.17	0.13	-0.14	-0.22	ns	ns	ns	-0.03	ns	-0.06	0.12	0.05	-0.16	0.05
1st pers plural	ns	ns	-0.13	0.21	0.11	0.03	0.18	0.05	ns	0.05	ns	-0.04	-0.1	ns
2nd person	-0.06	0.05	-	0.04	0.16	ns	ns	0.02	ns	ns	-0.15	ns	-0.12	0.02
3rd pers singular	-	0.09	-	0.15	-	ns	-	ns	-	ns	-	0.02	-	ns
3rd pers plural	-	-0.05	-	0.26	-	-0.06	-	-0.04	-	ns	-	0.02	-	0.03
3rd pers overall	0.2	-	-	-	ns	-	ns	-	ns	-	ns	-	ns	-
Impersonal pronouns	-	-0.09	-	0.11	-	-0.05	-	ns	-	ns	-	0.02	-	0.08
Articles	-0.24	-0.24	-	0.28	ns	-0.05	ns	ns	0.09	0.02	-0.11	-0.02	0.2	0.13
Common verbs	-	0.04	-	0.02	-	-0.03	-	ns	-	ns	-	0.04	-	0.03
Auxiliary verbs	-	0.02	-	0.08	-	-0.06	-	ns	-	ns	-	0.05	-	0.07
Past tense	0.12	-0.03	-0.16	ns	ns	-0.04	0.1	0.02	ns	-0.02	ns	ns	-0.16	ns
Present tense	0.18	0.08	0.04	ns	ns	ns	ns	ns	ns	ns	ns	0.04	-0.16	0.03
Future tense	ns	-0.07	0.14	0.09	ns	-0.05	ns	ns	ns	ns	ns	0.03	ns	0.05
Adverbs	-	0.05	-	-0.07	-	-0.04	-	ns	-	ns	-	0.05	-	0.04
Prepositions	-0.17	-0.13	-	0.27	ns	-0.04	ns	0.03	ns	0.06	ns	ns	0.17	0.06
Conjunctions	-	0.03	-	0.12	-	-0.02	-	0.02	-	0.02	-	0.02	-	0.06
Negations	0.11	ns	-	-0.12	ns	-0.06	ns	-0.05	-0.17	-0.03	0.11	0.07	-0.13	0.02
Quantifiers	-	-0.09	-	0.24	-	-0.02	-	0.03	-	0.05	-	ns	-	0.05
Numbers	-0.15	-0.13	-	0.05	-0.12	-0.06	0.11	0.02	ns	0.02	ns	ns	-0.08	0.06
Swear words	-0.22	-0.21	-	-0.17	ns	ns	-0.21	-0.15	-0.14	-0.09	0.11	0.06	ns	ns
Social processes	-	0.08	-0.13	0.21	0.15	0.04	0.13	0.02	ns	ns	ns	ns	-0.14	ns
Family	0.12	0.22	-	0.28	0.09	0.03	0.19	0.03	ns	0.03	ns	ns	-0.17	-0.12
Friends	0.09	0.08	-	0.26	0.15	0.05	0.11	0.04	ns	0.02	-0.08	ns	ns	-0.04
Humans	ns	0.04	-	0.06	0.13	0.06	ns	-0.05	-0.12	ns	ns	ns	-0.09	ns
Affective processes	0.11	0.11	-	-0.05	0.09	0.07	ns	0.02	ns	ns	ns	ns	-0.12	-0.04
Positive emotion	ns	0.21	0.12	0.14	0.1	0.13	0.18	0.13	ns	0.1	ns	-0.08	-0.15	-0.07
Negative emotion	0.1	-0.12	-0.05	-0.31	ns	-0.07	-0.15	-0.17	-0.18	-0.13	0.16	0.15	ns	0.03
Anxiety	0.16	0.08	-	-0.13	ns	-0.04	ns	-0.02	ns	-0.02	0.17	0.06	ns	0.07
Anger	ns	-0.22	-	-0.25	ns	-0.05	-0.23	-0.19	-0.19	-0.12	0.13	0.11	ns	0.02
Sadness	0.1	0.08	-	-0.15	ns	-0.04	ns	-0.02	-0.11	-0.04	0.1	0.09	ns	ns
Cognitive processes	0.07	-0.03	0.07	0.1	ns	-0.05	ns	0.02	-0.11	ns	0.13	0.04	-0.09	0.1
Insight	0.09	-0.05	0.11	0.04	ns	-0.09	ns	ns	ns	-0.02	ns	0.05	ns	0.13
Causation	ns	-0.05	ns	-0.01	-0.09	-0.06	-0.11	-0.02	-0.12	ns	0.11	0.02	ns	0.08
Discrepancy	0.07	ns	-	0.02	ns	-0.05	ns	-0.02	-0.13	-0.03	0.13	0.07	-0.12	0.02
Tentative	ns	-0.12	-	0.07	-0.11	-0.08	ns	ns	-0.1	-0.03	0.12	0.06	ns	0.07
Certainty	0.14	ns	-	0.09	0.1	ns	ns	0.03	-0.1	0.04	0.13	ns	ns	0.06
Inhibition	-	0.03	-	0.09	-0.13	ns	ns	ns	ns	0.04	0.09	ns	ns	ns
Inclusive	ns	0.04	-	0.23	0.09	0.04	0.18	0.05	ns	0.05	ns	-0.02	0.11	0.06
Exclusive	ns	-0.05	ns	ns	ns	-0.07	ns	ns	-0.16	-0.03	0.1	0.05	ns	0.05
Perceptual Processes	0.12	ns	-	-0.06	0.09	-0.04	ns	ns	-0.1	-0.07	ns	0.03	-0.11	0.1
See	ns	ns	-	ns	ns	-0.02	0.09	ns	ns	-0.04	ns	ns	ns	0.04
Hear	0.1	-0.07	-	-0.1	0.12	-0.04	ns	ns	-0.12	-0.06	ns	0.02	-0.08	0.08
Feel	0.17	0.04	-	-0.07	ns	-0.02	0.1	ns	ns	-0.04	0.1	0.03	ns	0.05
Biological processes	ns	0.05	-	-0.06	0.14	0.04	0.09	-0.06	ns	-0.06	ns	0.05	-0.09	0.02
Body	-	-0.02	-	-0.14	0.1	ns	0.09	-0.09	ns	-0.09	ns	0.06	-0.04	0.04
Health	-	0.05	-	0.07	-	ns	-	ns	-	ns	-	0.06	-	ns
Sexual	ns	0.05	-	-0.14	0.17	0.1	0.08	-0.04	ns	-0.04	ns	ns	ns	ns
Ingestion	-	0.02	-	0.12	-	ns	-	-0.03	-	-0.03	-	ns	-	0.03
Relativity	-	-0.06	-	0.16	-	ns	-	0.05	-	0.08	-	-0.03	-	-0.03
Motion	0.07	ns	-	0.12	-	0.02	-	0.05	-	0.07	-	-0.04	-	-0.04
Space	ns	-0.18	-	0.21	ns	ns	0.16	ns	ns	0.02	-0.09	ns	-0.11	0.07
Time	ns	0.02	-0.19	0.08	ns	ns	0.12	0.06	0.09	0.09	ns	-0.03	-0.22	-0.07
Work	-0.12	-0.08	-	-0.02	-0.08	-0.05	ns	0.03	ns	0.1	ns	-0.03	ns	-0.02
Achievement	-	-0.17	-	0.16	-0.09	ns	ns	0.05	0.14	0.11	ns	-0.06	ns	-0.02
Leisure	ns	-0.08	-	0.03	0.08	0.06	0.15	0.04	ns	0.03	ns	-0.07	-0.17	ns
Home	0.15	0.19	-	0.18	ns	ns	0.19	0.03	ns	0.04	ns	-0.02	-0.2	-0.06
Money	-0.1	-0.12	-	0.24	ns	ns	-0.11	-0.04	ns	0.03	ns	ns	ns	0.03
Religion	-	-0.03	-	0.21	0.11	ns	ns	0.06	ns	0.04	ns	-0.04	ns	ns
Death	-	-0.18	-	-0.1	ns	-0.08	-0.13	-0.09	-0.12	-0.08	ns	0.08	0.15	0.09
Assent	-	0.07	-	-0.22	ns	0.05	ns	0.04	-0.09	ns	ns	-0.04	-0.11	-0.05
Nonfluencies	-	-0.03	-	0.02	-	ns	-	ns	-	ns	-	0.03	-	ns
Filters	-	-0.02	-	-0.24	-	ns	-	-0.04	-	-0.08	-	0.03	-	0.04
participants (N)	9,130	74,859	3,087	74,859	576	72,709	576	72,772	576	72,781	576	71,968	576	72,809

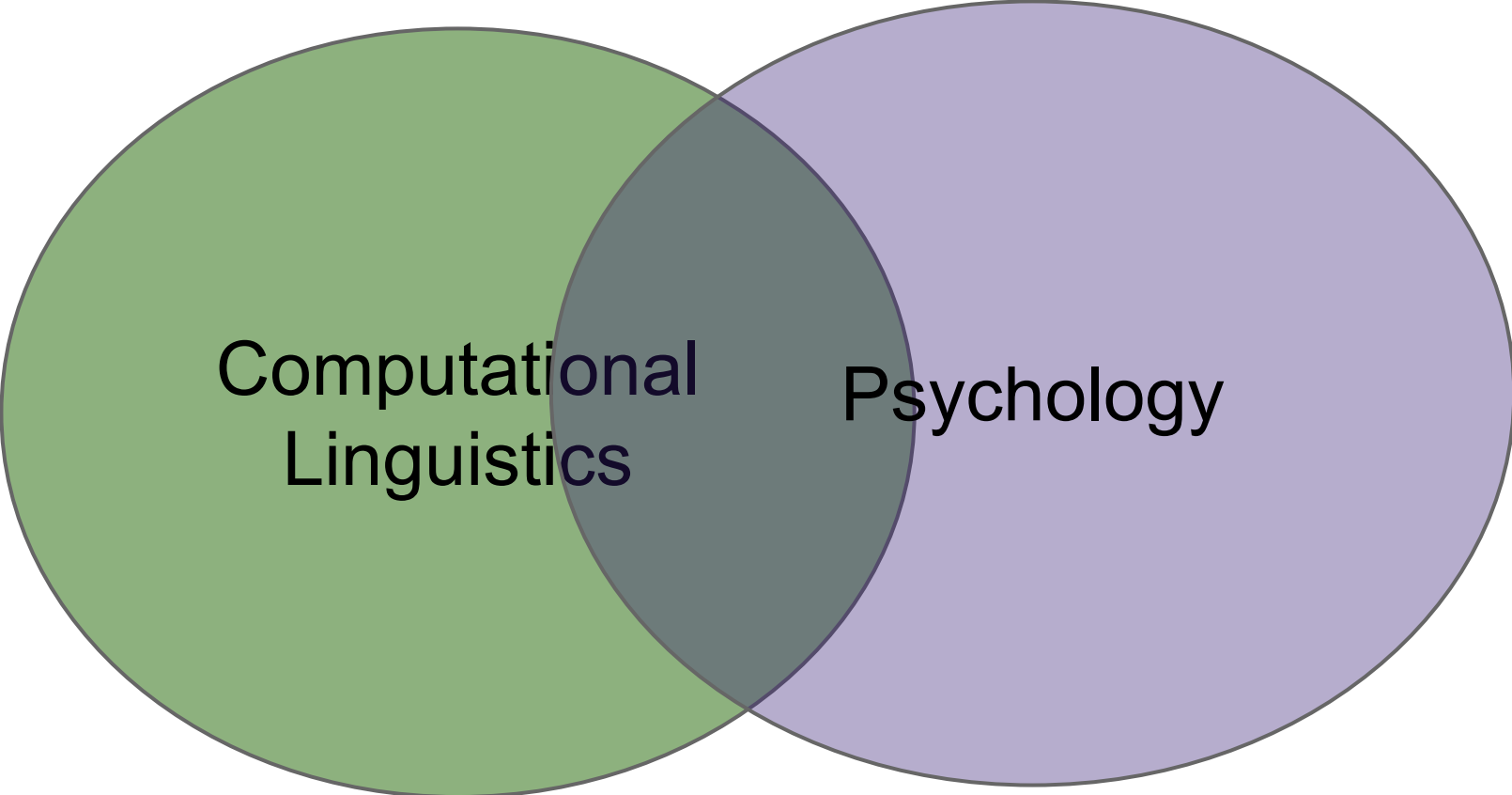


## **Psychological language analyses typically limited to apriori language lexica**

- **limited to pre-chosen hypotheses**
- **don't always measure what is expected**

**“open vocabulary”**: let the data dictate the words, phrases, and linguistic features that matter.

- **transparent results**

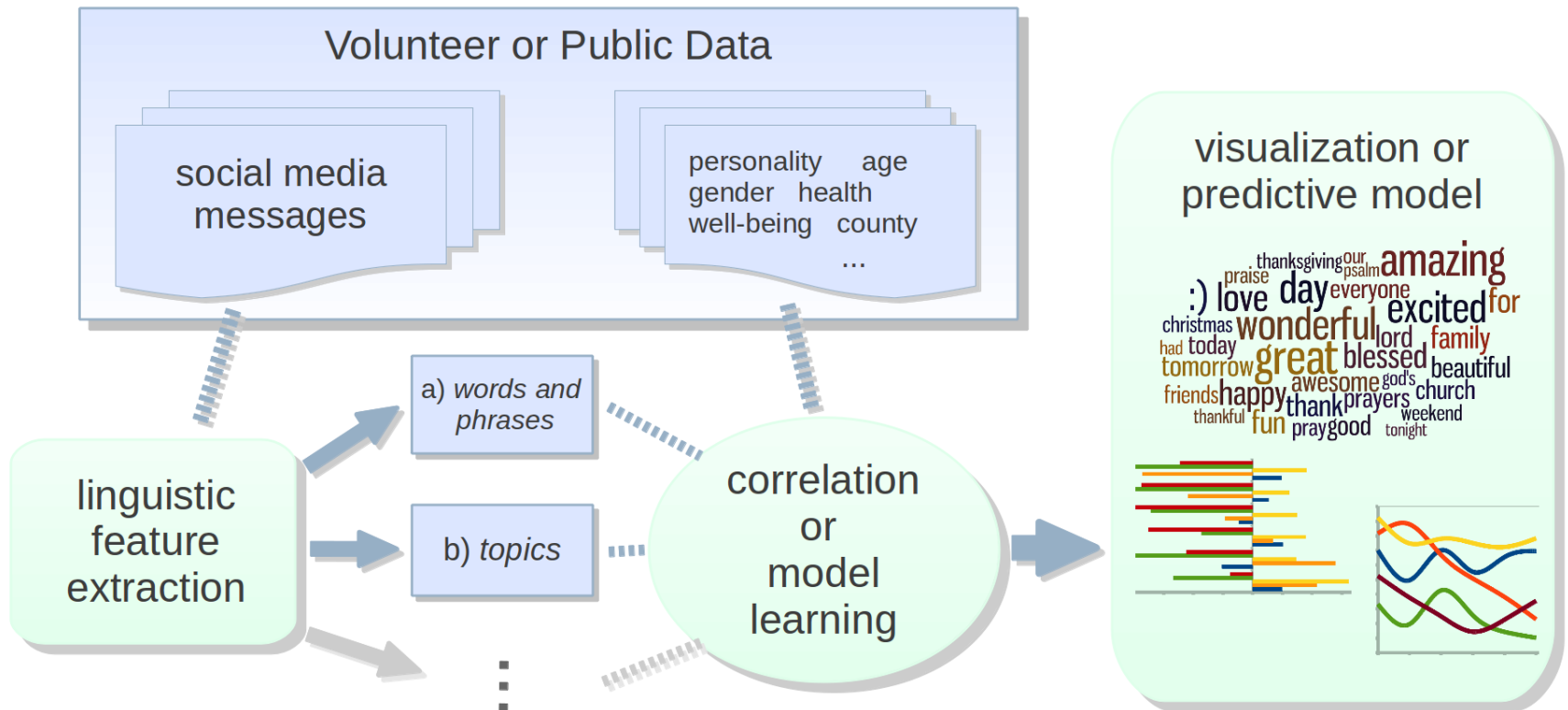


Computational  
Linguistics

Psychology

# Method

research mostly falls into this framework:



# Method

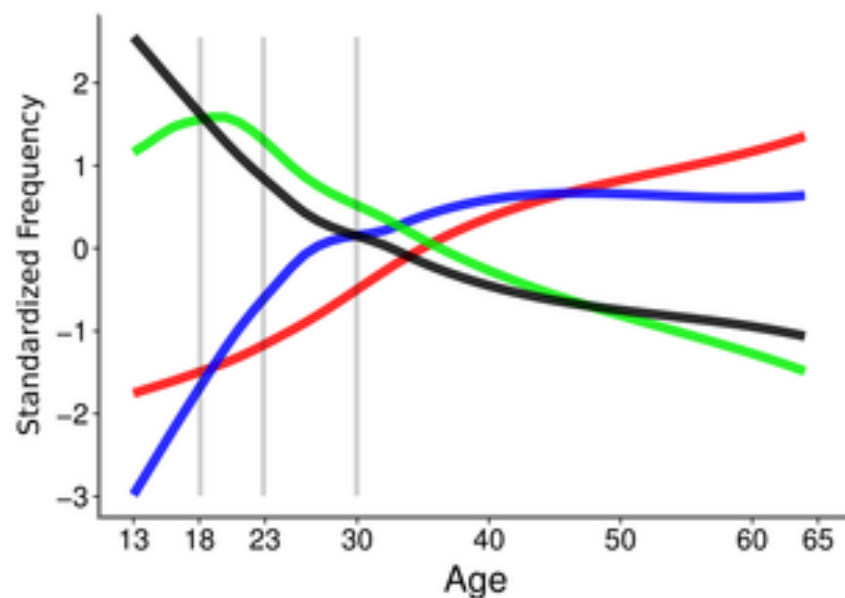
[http://prezi.com/rvpue6ucpci1/?utm\\_campaign=share&utm\\_medium=copy&rc=ex0share](http://prezi.com/rvpue6ucpci1/?utm_campaign=share&utm_medium=copy&rc=ex0share)



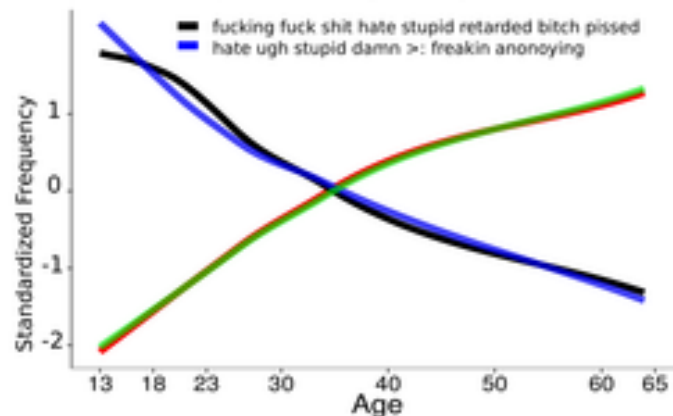


**A**

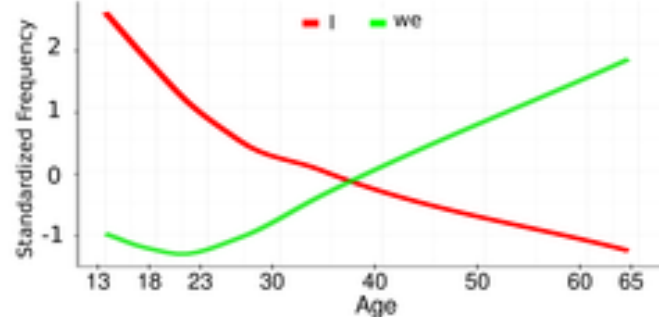
(30 to 65) ■ son daughter father mother proud oldest data youngest  
 (23 to 29) ■ job position company manager interview experience office assistant  
 (19 to 22) ■ classes semester class college schedule summer registered taking  
 (13 to 18) ■ haha lol :p :D :) hehe jk ;p

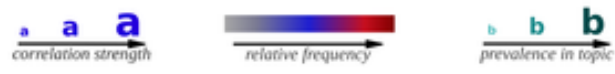
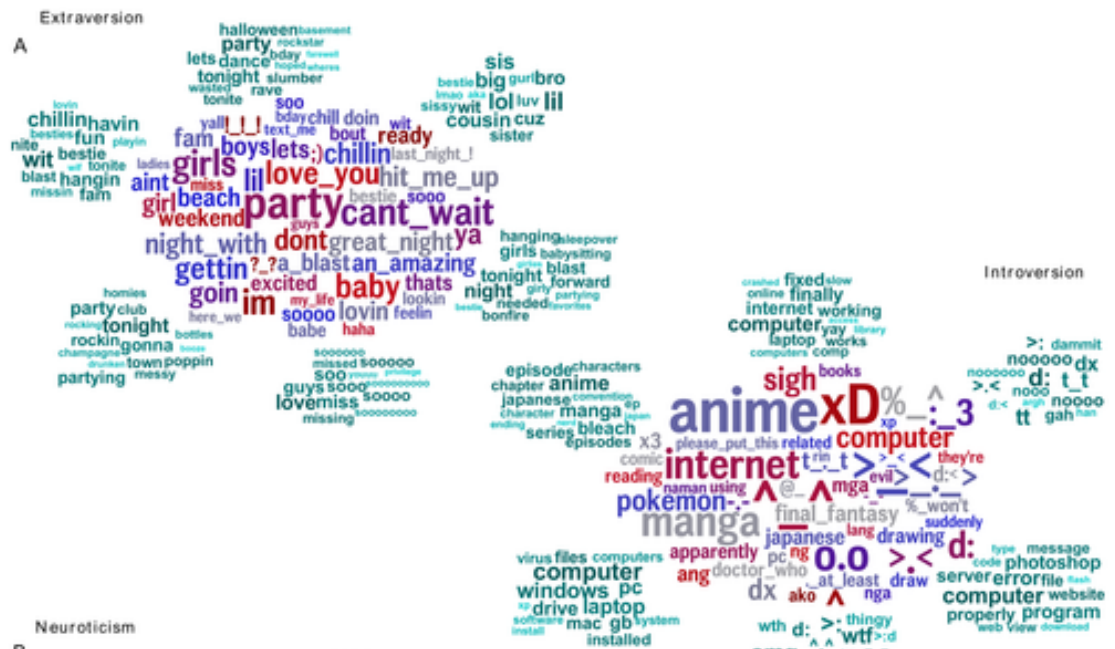
**B**

■ family friends wonderful blessed amazing thankful loving  
■ kind person word loving act caring honest touch  
■ fucking fuck shit hate stupid retarded bitch pissed  
■ hate ugh stupid damn >: freakin annoying

**C**

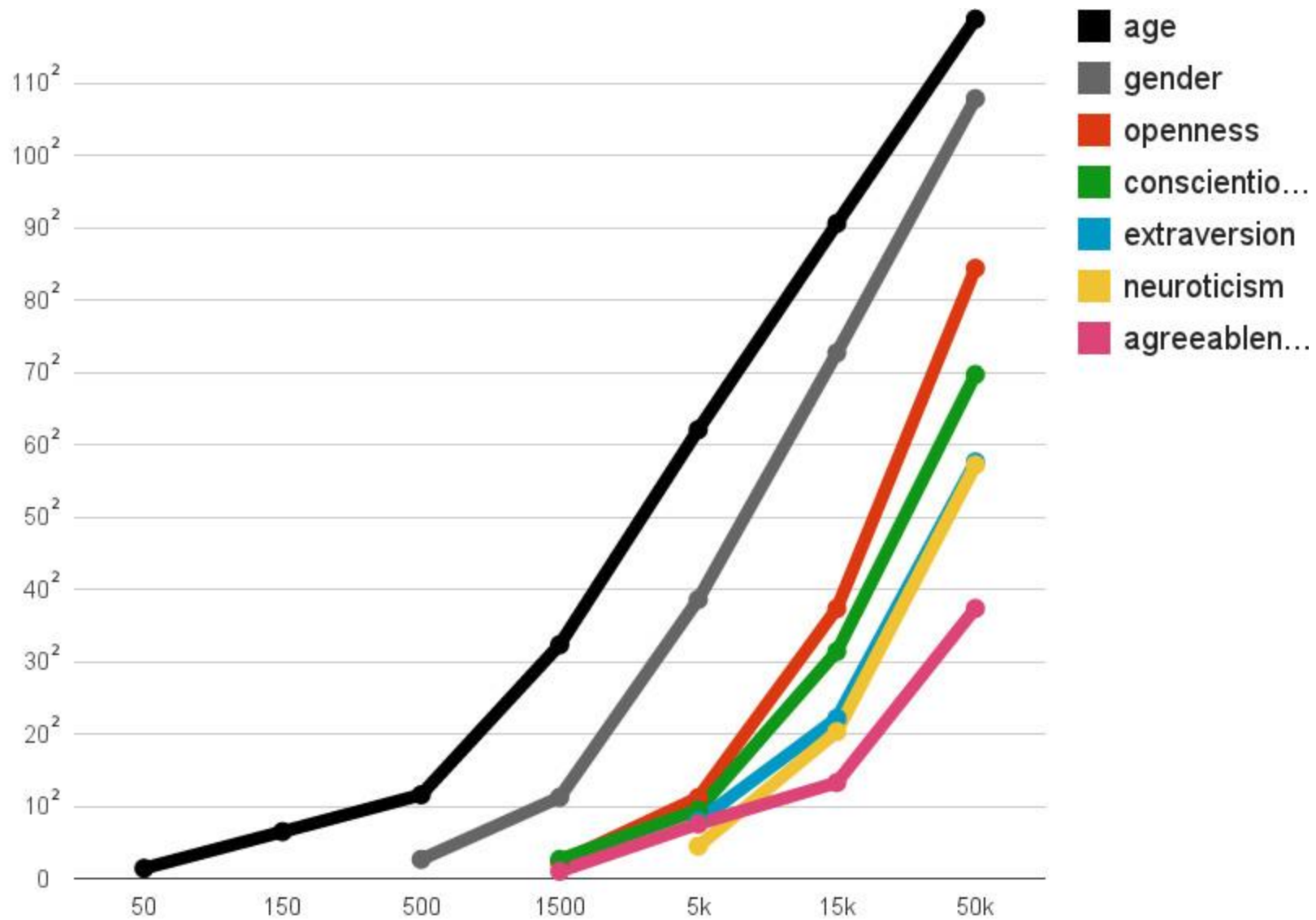
■ I ■ we







significantly correlated features



sample size

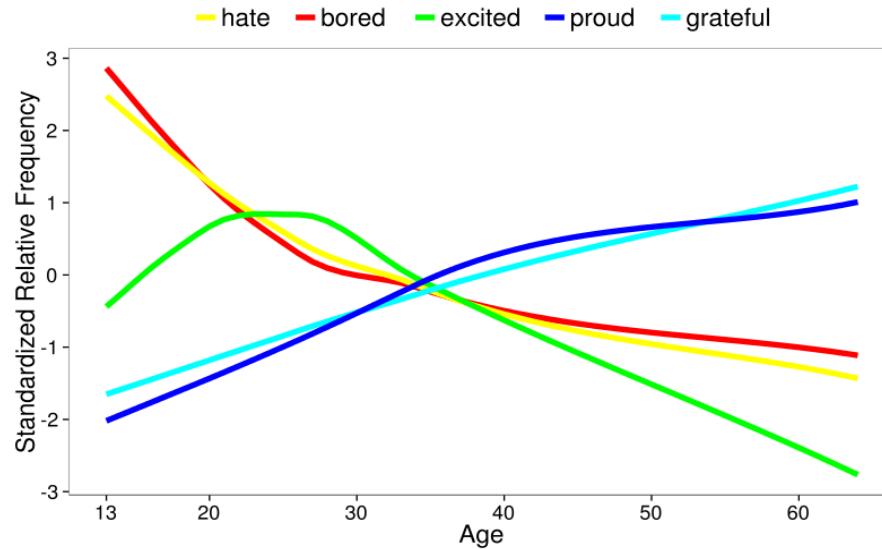
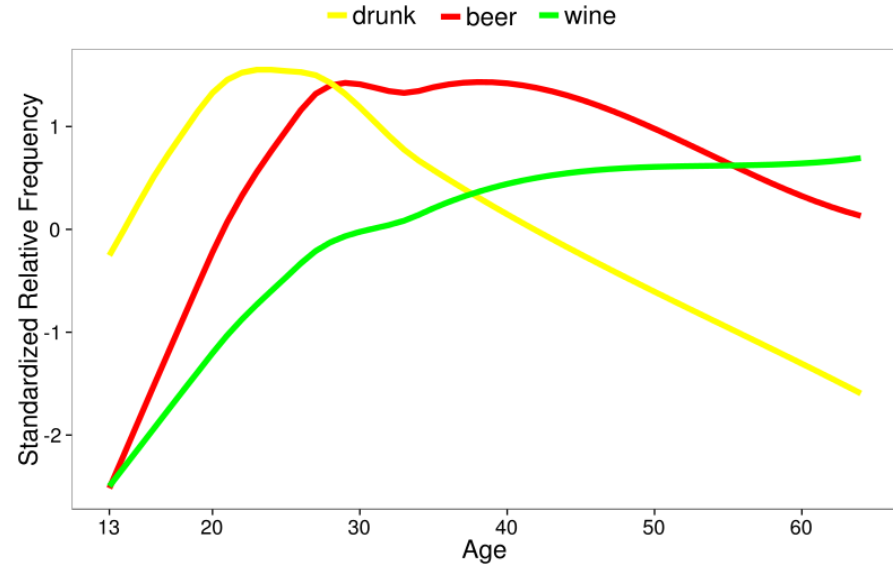
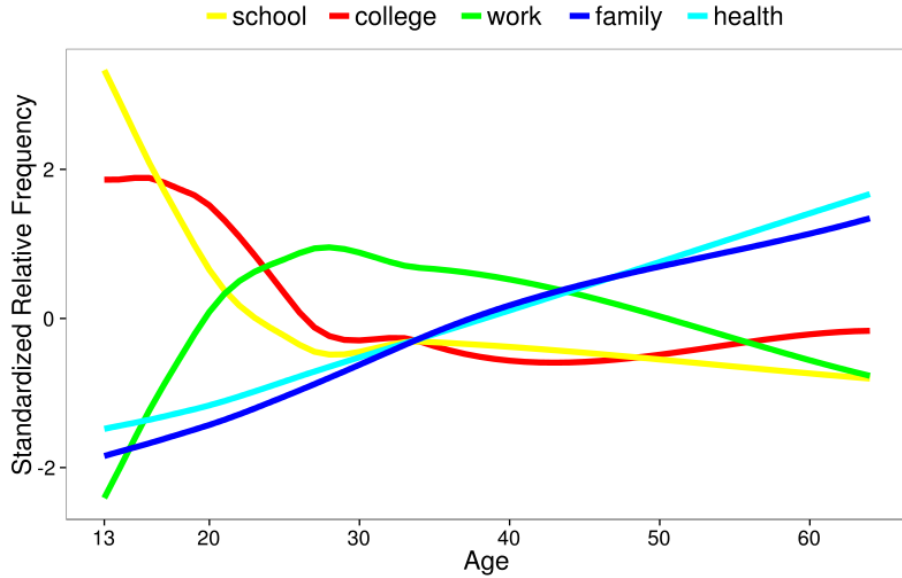


**Supporting Table 2. Prediction results when selecting features via differential language analysis.**

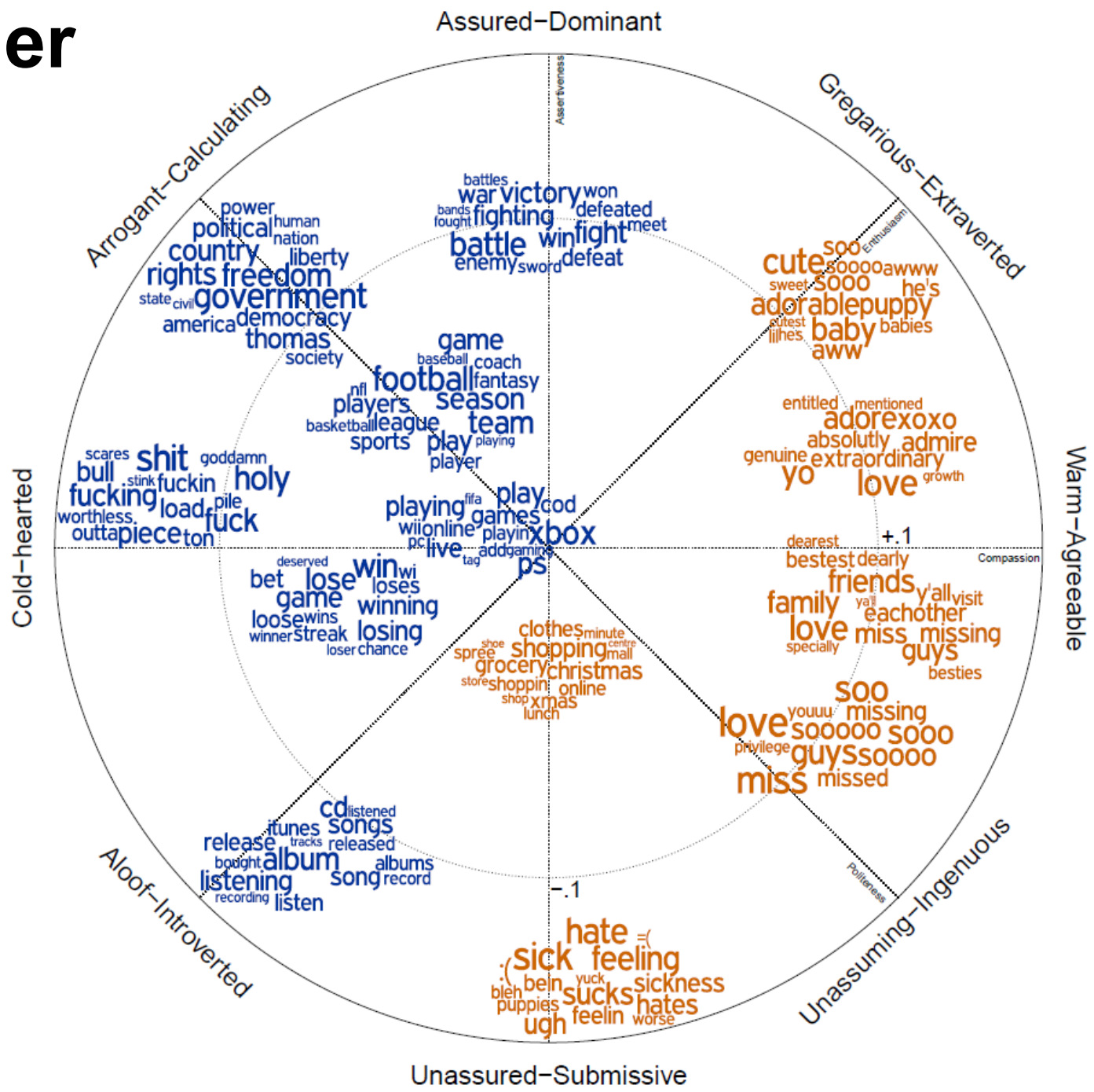
features	Gender <i>accuracy</i>	Age <i>R</i>	Extraversion <i>R</i>	Agreeableness <i>R</i>	Conscientious. <i>R</i>	Neuroticism <i>R</i>	Openness <i>R</i>
<i>LIWC</i>	77.7%	.65	.25	.25	.29	.22	.28
<i>Topics</i>	<b>88.2%</b>	<b>.79</b>	<b>.34</b>	<b>.28</b>	<b>.34</b>	<b>.28</b>	<b>.39</b>
<i>WordPhrases</i>	<b>91.8%</b>	<b>.81</b>	<b>.37</b>	<b>.27</b>	<b>.34</b>	<b>.28</b>	<b>.40</b>
<i>WordPhrases + Topics</i>	<b>92.0%</b>	<b>.82</b>	<b>.38</b>	<b>.29</b>	<b>.35</b>	<b>.30</b>	<b>.41</b>
<i>Topics + LIWC</i>	<b>89.2%</b>	<b>.80</b>	<b>.35</b>	<b>.28</b>	<b>.34</b>	<b>.28</b>	<b>.40</b>
<i>WordPhrases + LIWC</i>	<b>91.8%</b>	<b>.81</b>	<b>.38</b>	<b>.28</b>	<b>.34</b>	<b>.29</b>	<b>.40</b>
<i>WordPhrases + Topics + LIWC</i>	<b>92.0%</b>	<b>.82</b>	<b>.38</b>	<b>.30</b>	<b>.35</b>	<b>.30</b>	<b>.41</b>

*accuracy*: percent predicted correctly (for discrete binary outcomes). *R*: Square-root of the coefficient of determination (for sequential / continuous outcomes). *LIWC*: *A priori* word-categories from Linguistic Inquiry and Word Count. *Topics*: Automatically created *LDA* topic clusters. *WordPhrases*: words and phrases (n-grams of size 1 to 3 passing a collocation filter). Bold indicates significant ( $p < .01$ ) improvement over the baseline set of features (use of *LIWC* alone). Differential language analysis was run over the training set, and only those features significant at Bonferonni-corrected  $p < 0.001$  were included during training and testing. No controls were used so as to be consistent with the evaluation in the main paper, and so one could consider this a univariate feature selection. On average results are just below those of not using *differential language analysis* to select features but there is no significant difference.

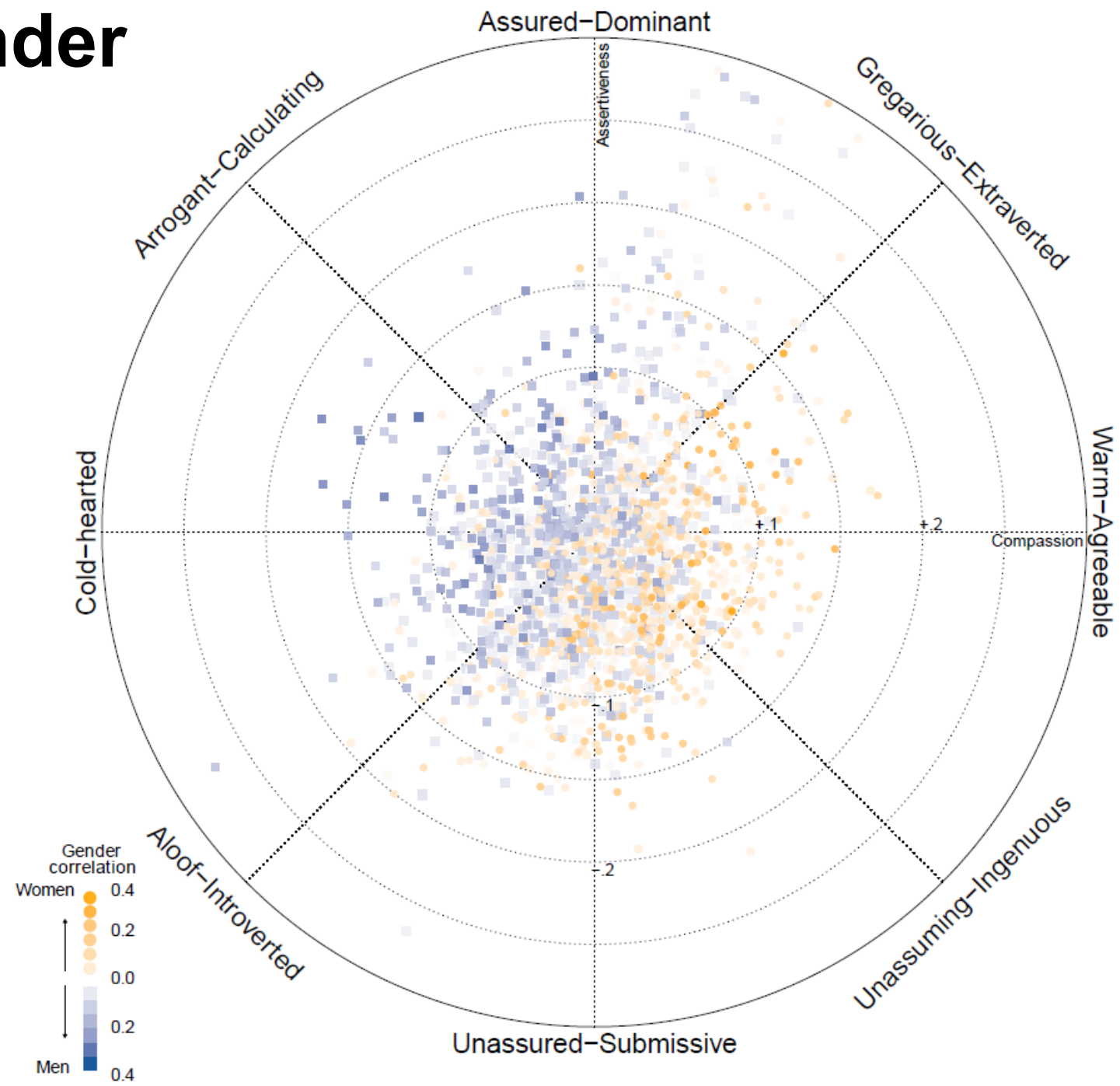
# Development



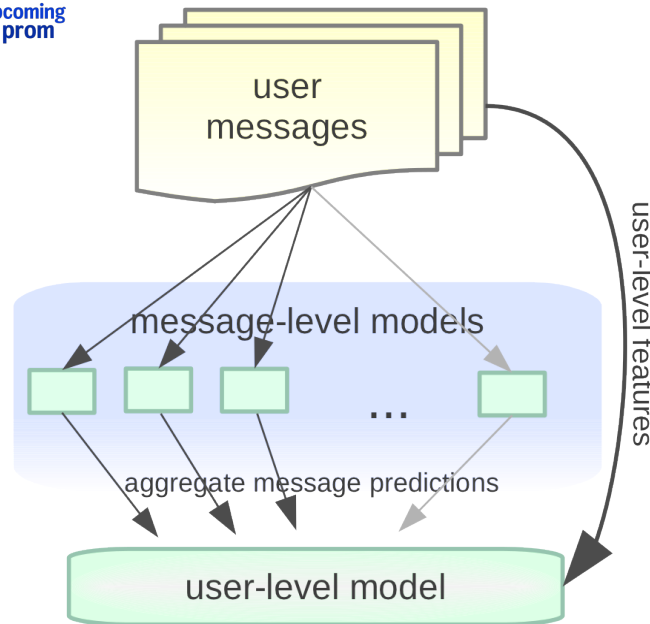
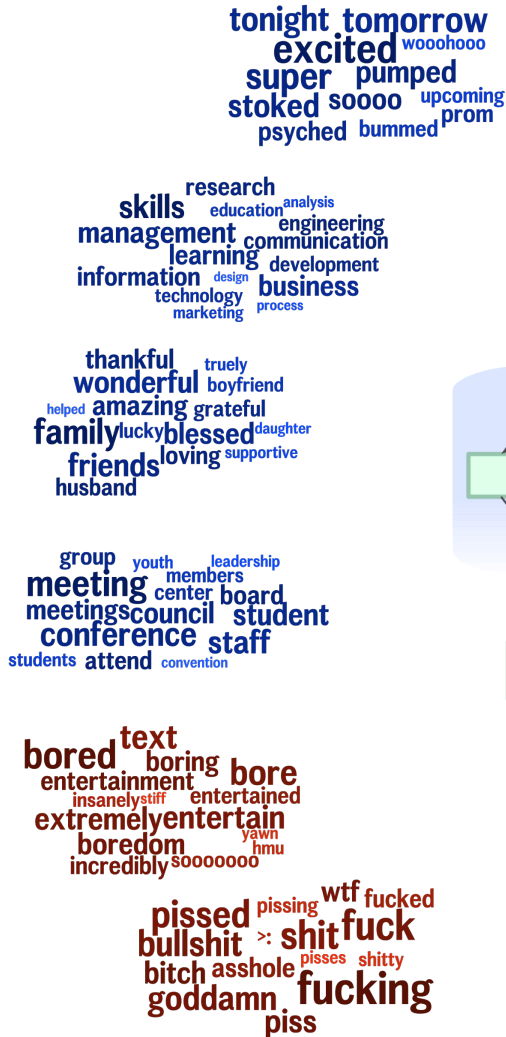
# Gender



# Gender



# Individual Well-Being: message to user-level



baselines	$r$
(mean)	.000
<i>lexica: GNH</i>	.210
<i>lexica: Hedonometer</i>	.108

message and user-level

# Optimism/Pessimism - CAVE

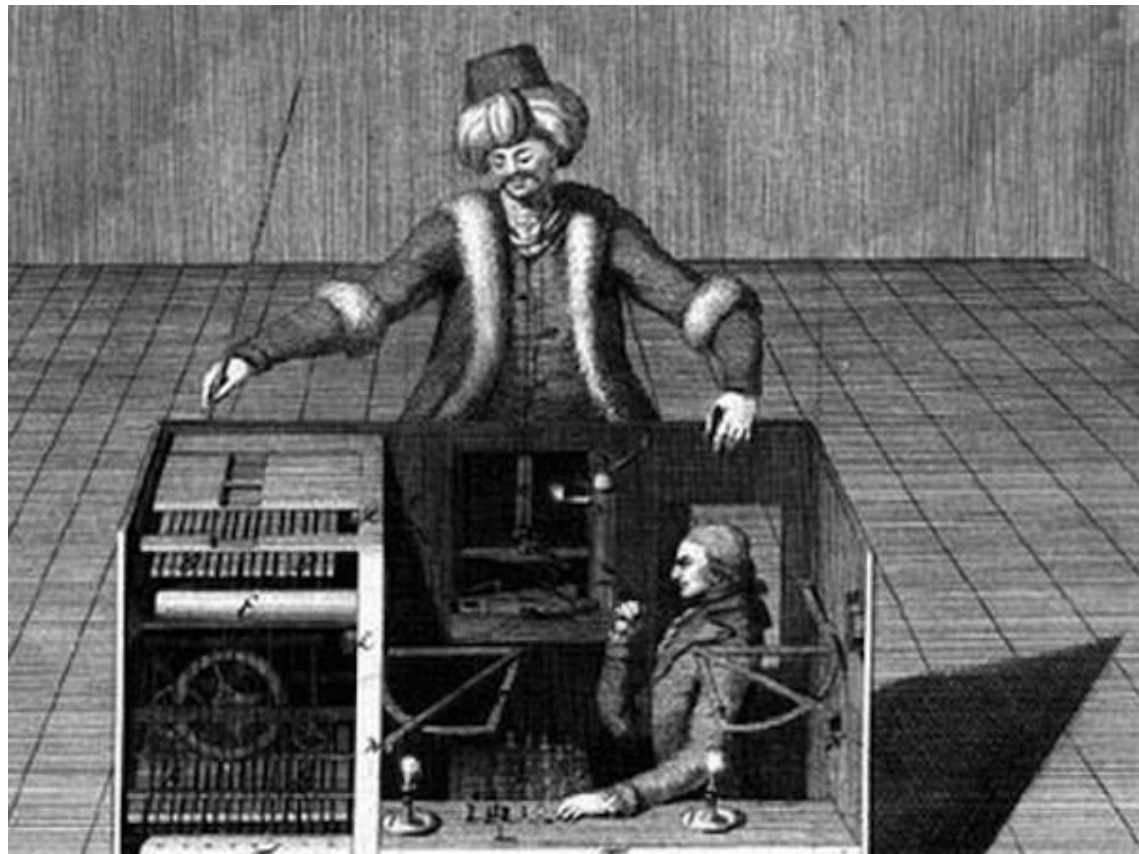
- Uses explanatory style
- Three factors:
  - Internality: whether the cause of the event implicates something about the speaker or is due to situational characteristics
  - Stability: whether the cause of the event persists across time
  - Globality: whether the cause of the event persists across situation (often covaries with stability)
- Positive events: high scores indicate optimism
- Negative events: low scores indicate optimism



# Mechanical Turk

**Problem:** Manually created lists of words are not always used as expected.

**Solution:** Annotate thousands of samples of real Twitter and Facebook posts for expressions of well-being.



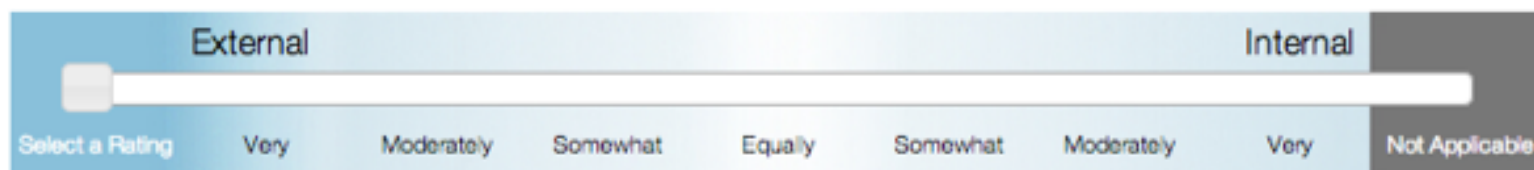
RT @XSTROLOGY: #Leo can not stand you bad-mouthing their lover.&lt;~true its goin be a fight if u talkin sideways bout my boo lol

Does the above message include a causal explanation?

Yes  No

To what extent does the causal explanation of this message indicate **that the cause of the event is internal (e.g. cause of the event placed on speaker)** or **that the cause of the event is external (e.g. cause of the event placed on environment)**:

*If you cannot assess a degree of internality or externality, rate as 'Not Applicable'.*



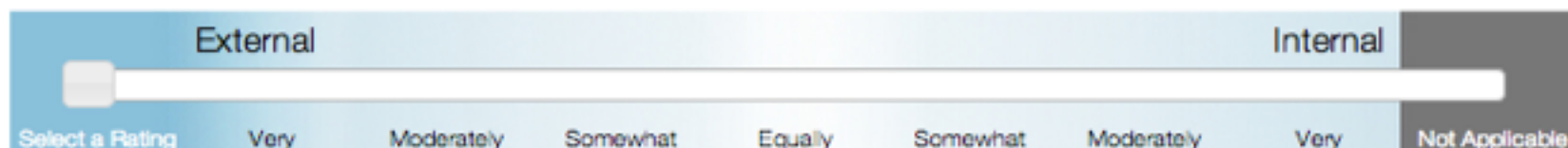
RT @XSTROLOGY: #Leo can not stand you bad-mouthing their lover.&lt;~true its goin be a fight if u talkin sideways bout my boo lol

Does the above message include a causal explanation?

Yes  No

To what extent does this message indicate **internality** or **externality**:

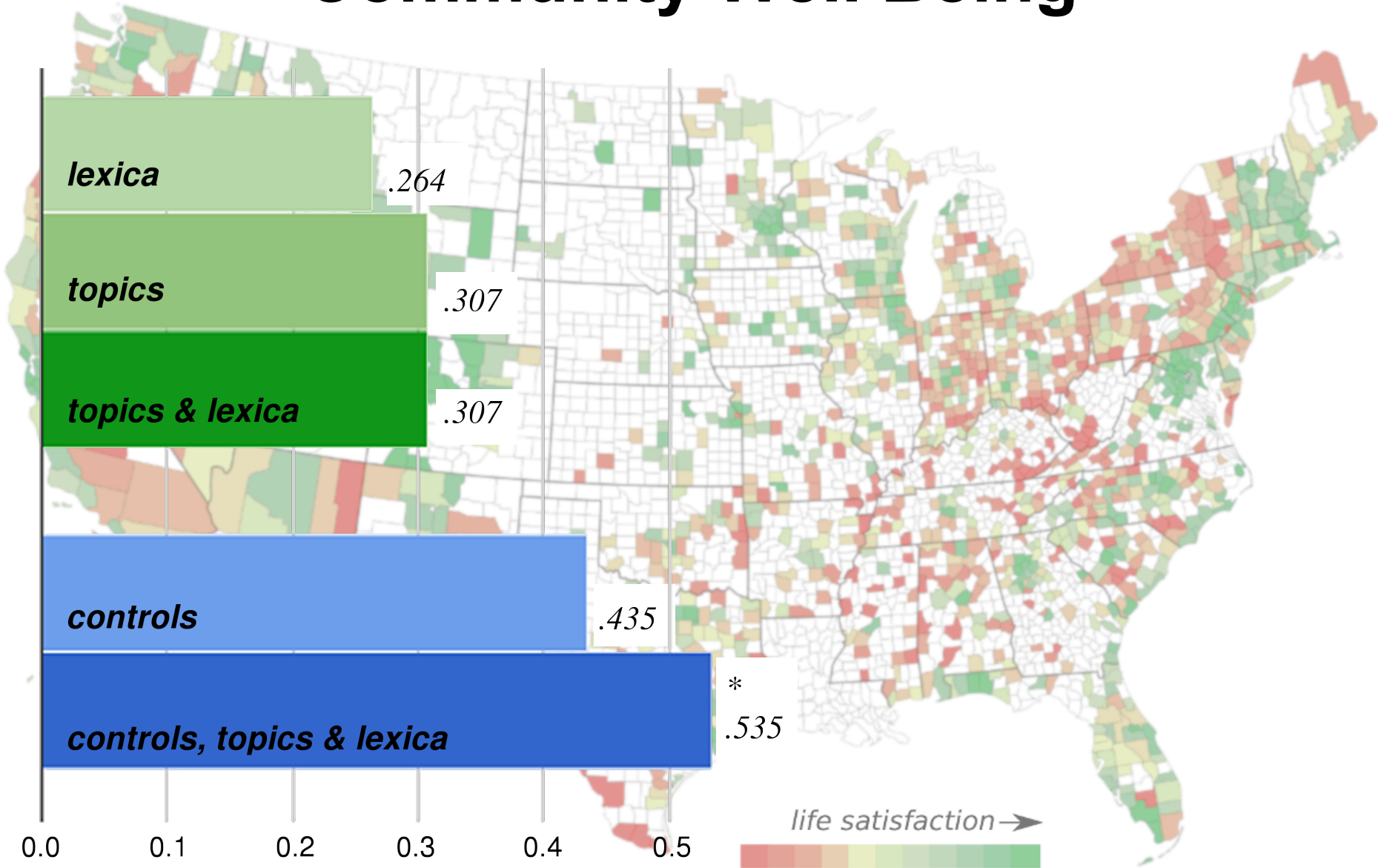
*If you cannot assess a degree of internality or externality, rate as 'Not Applicable'.*



# End of talk

(remaining slides were not presented but might look fun!)

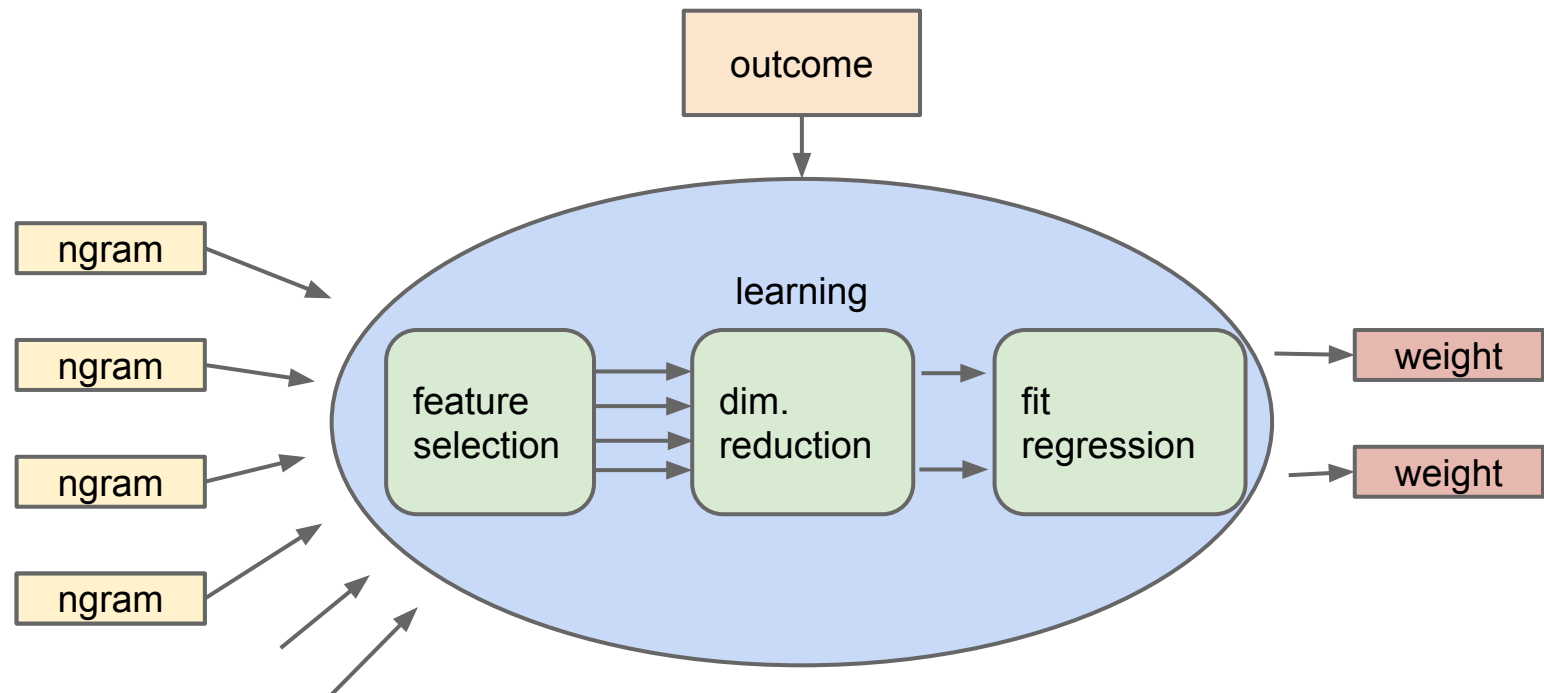
# Community Well Being



\*significant improvement over controls alone

# Generating Lexica from \*most\* Supervised n-gram Models

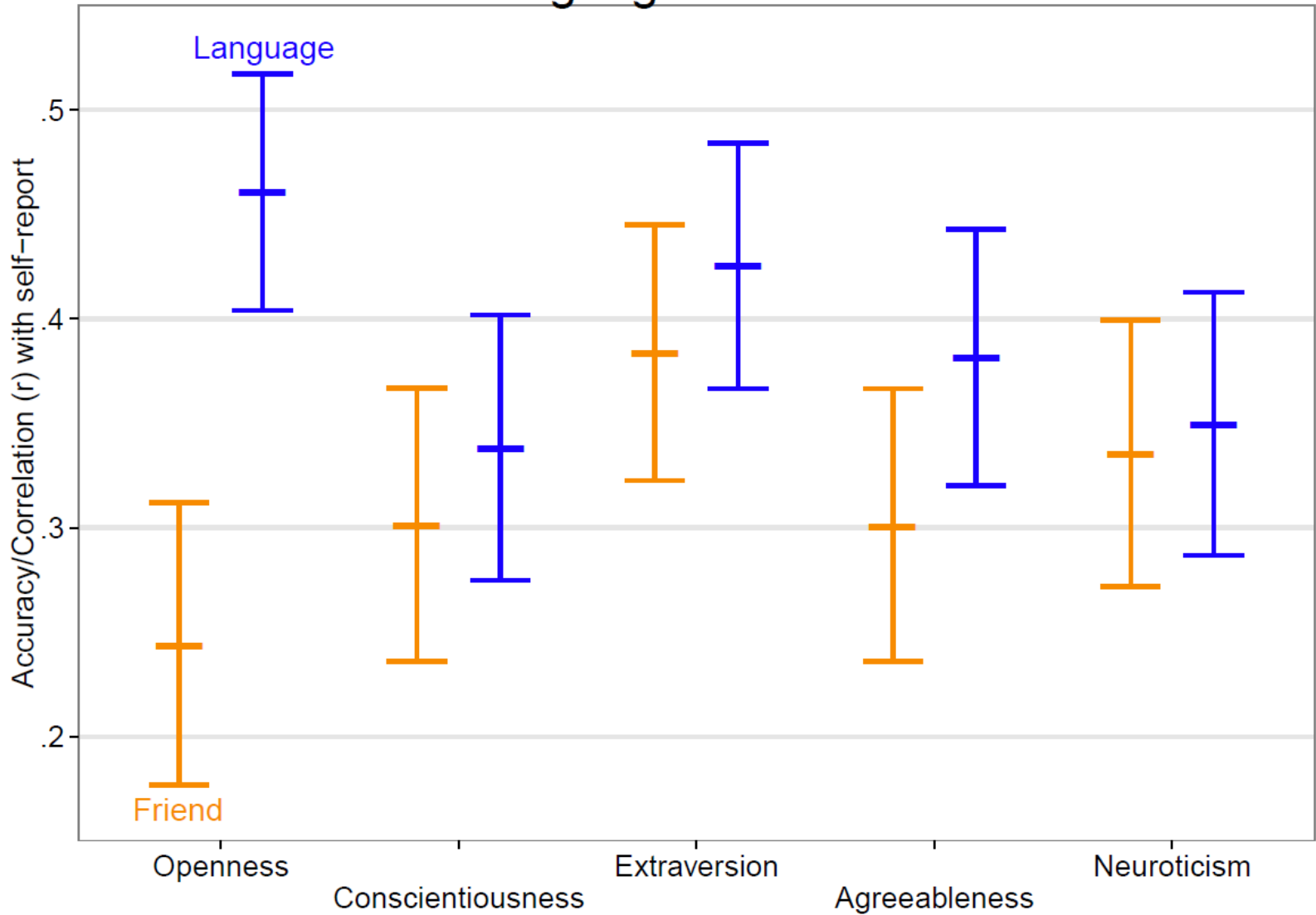
- Generalize multi-variate regression model into lexica.
- Works at multiple levels:
  - Hand annotated messages or users
    - OCEAN: User-level Cambridge data set
    - PERMA: MTurk-ed messages





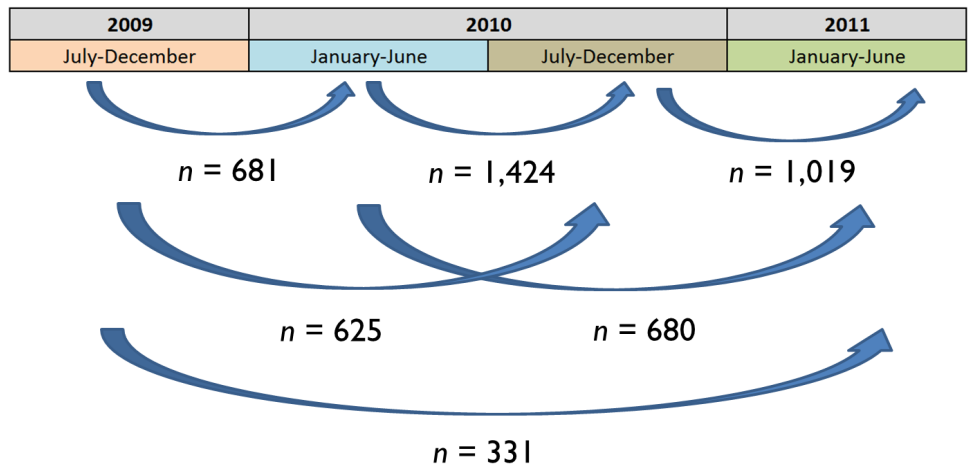
# Language-Based Psychometrics

## Predicting Personality Traits: Language vs. Friends

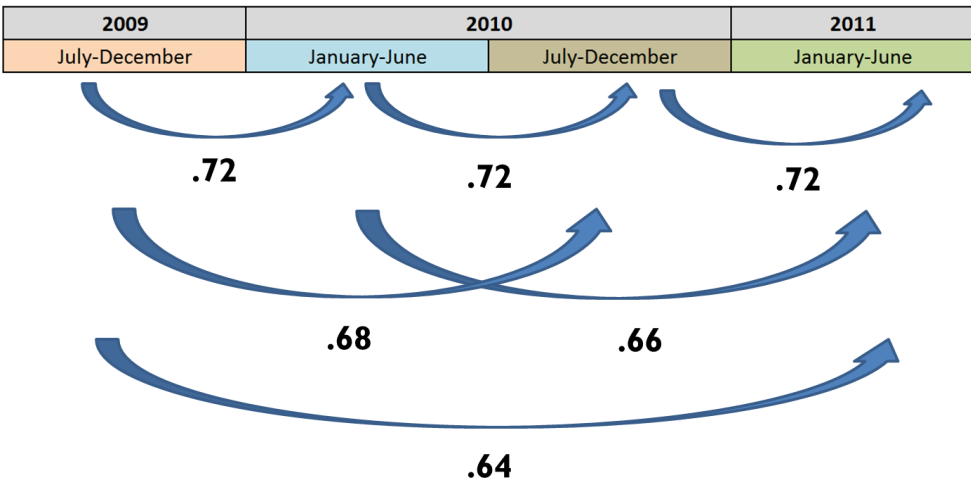


# Language-Based Psychometrics

## Sample sizes



## Extraversion



	Time 2	Time 3	Time 4
Time 1	.71	.68	.64
O Time 2		.74	.71
Time 3			.76
Time 1	.75	.74	.70
C Time 2		.76	.72
Time 3			.76
Time 1	.72	.68	.64
E Time 2		.72	.66
Time 3			.72
Time 1	.65	.61	.55
A Time 2		.64	.57
Time 3			.65
Time 1	.62	.57	.51
N Time 2		.62	.61
Time 3			.63