

# Empath: Understanding Topic Signals in Large-scale Text

Ethan Fast, Binbin Chen, Michael Bernstein

STANFORD HCI GROUP

**“Language is rich in subtle signals.”**

rich → *wealth*

subtle → *cleverness*

language, signals → *communication*

emotional contagion  
(kramer et al., 2014)

linguistic correlates of deception  
(ott et al., 2011)

conversational signs of betrayal  
(niculae et al., 2015)

# LIWC: linguistic inquiry and word count

(Pennebaker et. al, 2001)

*anger* = {scream, war, mad, ...}

but what about other categories like violence or social media?

e.g., “paypal” not in *money* category

# Empath

generate categories from seed words

*twitter, facebook* → {tweet, instagram, selfie, comment...}

broad set of 200 built-in categories:

technology = {ipad, android, ...}

violence = {bleed, punch, ...}

government = {embassy, democrat, ...}

strength = {tough, forceful, ...}

**how Empath works**

how researchers can use it

how we evaluated it

# analysis via lexicon

“The CHI attendees **scream** in **rage** at the poor quality of the talk.”

2 anger words

2 (anger) / 13 (total words) = 0.15

normalized anger count

**open-ended category generation**

built-in categories

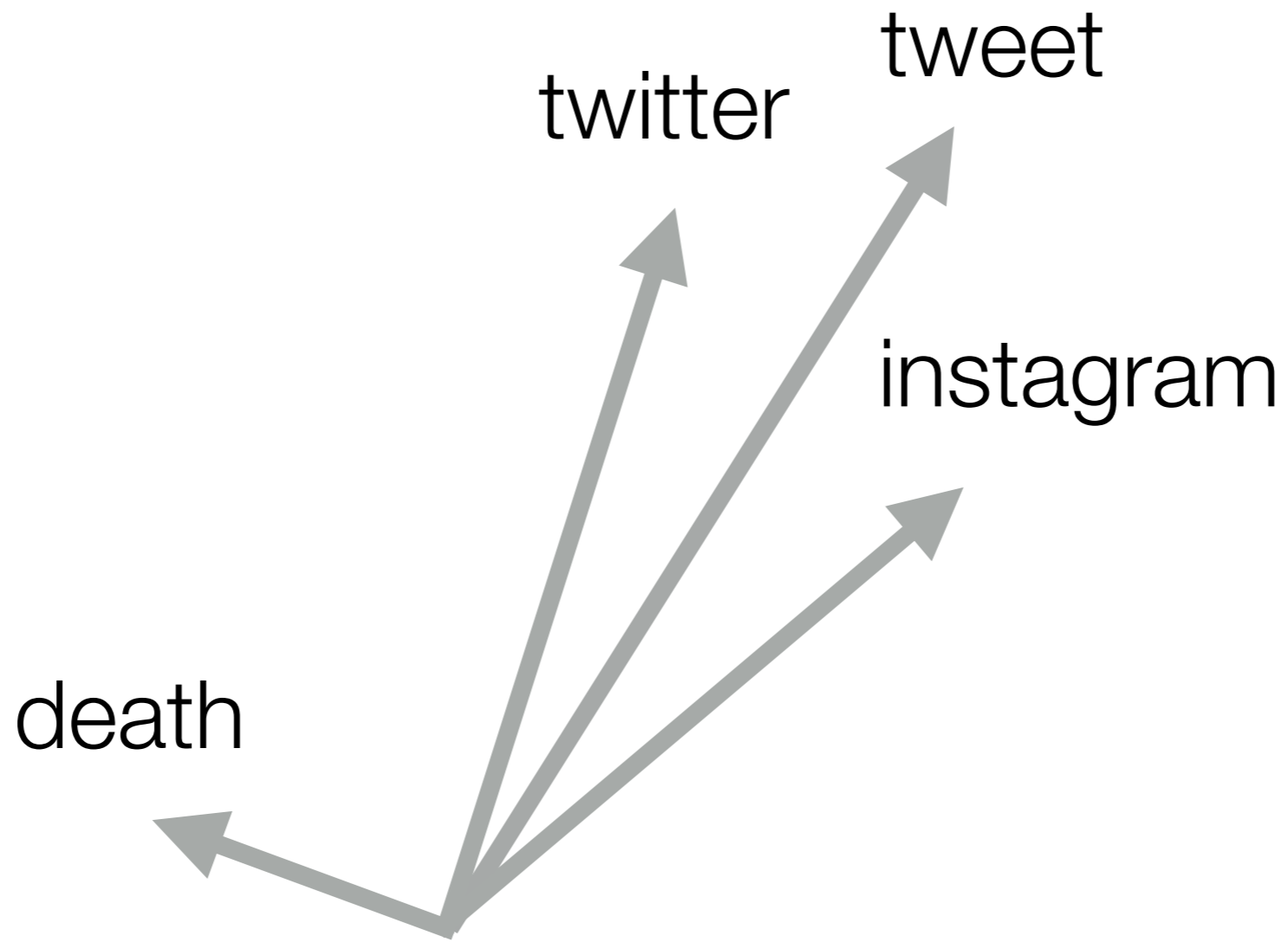


# facebook + twitter

twitter, facebook, instagram, tumblr, social\_media,  
twitter\_account, news\_feed, snapchat, tweet, tweets,  
newsfeed, tweeted, youtube, twitter\_page, mentions,  
facebook\_page, vine, dm, fan\_account, face\_book, timeline,  
notifications, fb, tweet, skype, tweeting, kik, app, notification,  
direct\_message, retweeted

continuous skip-gram neural  
embedding  
(mikolov, et al., 2013)

# words in a vector space



# crowd validation of categories using mechanical turk

social media

retweet



tweet



instagram



email

retweet

tweet

instagram

~~email~~

# models

wattpad (amateur fiction)

reddit (2008-2015)

new york times (1987-2007)

open-ended category generation

**built-in categories**

# seeding Empath's built-in categories

ConceptNet (liu and singh, 2004):  
{shirt, hat} are elements of {clothing}

clothing = empath.generate("shirt", "hat")  
{shirt, hat, hoodie, jumper, sweater, t-shirt, jacket, ...}

business, fabric, banking, play, party, furniture, power,  
childish, home, hiking, joy, vehicle, fun, timidity,  
dominant\_personality, eating, musical, legend, prison, cold,  
school, night, breaking, lust, masculine, ridicule,  
positive\_emotion, kink, monster, cleaning, journalism, rural,  
fear, kill, driving, traveling, white\_collar\_job, phone, restaurant,  
emotional, optimism, disappointment, smell, beach,  
appearance, cheerfulness, youth, war, science, achievement,  
superhero, envy, shame, occupation, body, sadness,  
aggression, tourism, ancient, negative\_emotion, office, anger,  
trust, meeting, fire, attractive, suffering, listen, neglect, music,  
sailing, sports, clothing, exasperation, reading, warmth,  
children, affection, law, urban, strength, movement, college,  
contentment, communication, farming, anonymity, medieval,  
deception, work, health, money, economics, heroic,  
domestic\_work, injury, medical\_emergency, dispute, poor,  
anticipation, cooking, nervousness, ugliness, wedding, leader,  
weakness, programming, valuable, wealthy, shape\_and\_size



# violence

bully, rape, impact, dislocated, bruise, harshly, kick, agony, stabbing, dead, torment, hit, beat, injure, aggravate, fight, wince, fatal, wound, scarring, bash, inflict, sting, hurt, minor, beating, injury, shatters, senseless, bleeding, kill, scared, afraid, mean, trauma, abusing, slap, feel, bleed, cut, mad, suffering, toughen, bad, violence, threaten, resuscitate, severe, bruising, scratch, strangle, punch, harm, abuse, bloody, hurting, punching, wounded, painful, violent, stab, angry, tough, damage, death, damaged, injures, wreck, punish, struggle

# hipster

iconic, stylish, fashionable, eccentric, outfit, sophisticated,  
punk, indie, wannabe, trendy, snazzy, fashioned, geek,  
themed, stereotype, geeky, label, looking, hippie, grunge,  
design, artsy, costume, urban, preppy, wear, funky, stylishly,  
brand, chic, hipster, hairstyle, converse, retro, sneaker,  
alternative, hairdo, clothing, styled, flashy, attire, nerdy,  
fashion, vintage, 1950s, wardrobe

# technology

robot, handheld, install, online, console, desktop, radar, keyboard, download, microchip, processor, database, inventor, simulator, cable, website, battery, scanning, hack, grid, transmitter, screen, spacecraft, data, interactive, computer, mobile, digital, network, prototype, technology, virtual, innovative, automate, mainframe, optical, technological, scientific, programming, scientist, outdated, module, communication, hacking, solar, scanner, binary, nexus, camcorder, connector, server, malfunction, machinery, compute, browser, advanced, technical, laptop, tablet, manufacture, engineering, web, interface, glitch, multiplayer, laboratory, experimental, research, wireless

how Empath works

**how researchers can use it**

how we evaluated it

**what kinds of words  
accompany our lies?**

(ott et al, 2011)

run Empath's built-in categories  
across the data  
(when comparing make corrections,  
e.g., bonferroni)

# trends in deceptive language

tormented (2.5 odds)

“it was **torture** hearing the sounds of the elevator which would never stop”

joyous (2.5 odds)

“I got a **great** deal and I am so **happy** I stayed here”

# trends in truthful language

ocean (1.6 odds)

“it seemed like a nice enough place  
with reasonably close **beach** access”

vehicles (2.5 odds)

“they took forever to Valet our **car**”



testing new hypotheses outside the  
scope of traditional lexicons

# **circular + big + small**

small, large, circular, huge, massive, gigantic, giant, center, big, circular, tiny, rectangular, enormous, centre, rectangle, wooden, size, marble, compact, oak, oval, shaped, structure, columns, triangle, square, very\_center, miniature, bordered, white\_stone, towers, decoration, exterior, granite, ginormous, white, shiny, brass, antique, shape, bronze, left\_side, adorned, plush, middle, ornate, smaller, squares, pillars, interior, square, sized, decorated, spanned, largest, near, flooring, lining, skeleton, larger, above, carpeted, branching, smallest, decorative, circumference, sized, ...

spatial language

(1.2 times more likely for truthful reviews,  $p < 0.001$ )

how Empath works

how researchers can use it

**how we evaluated it**

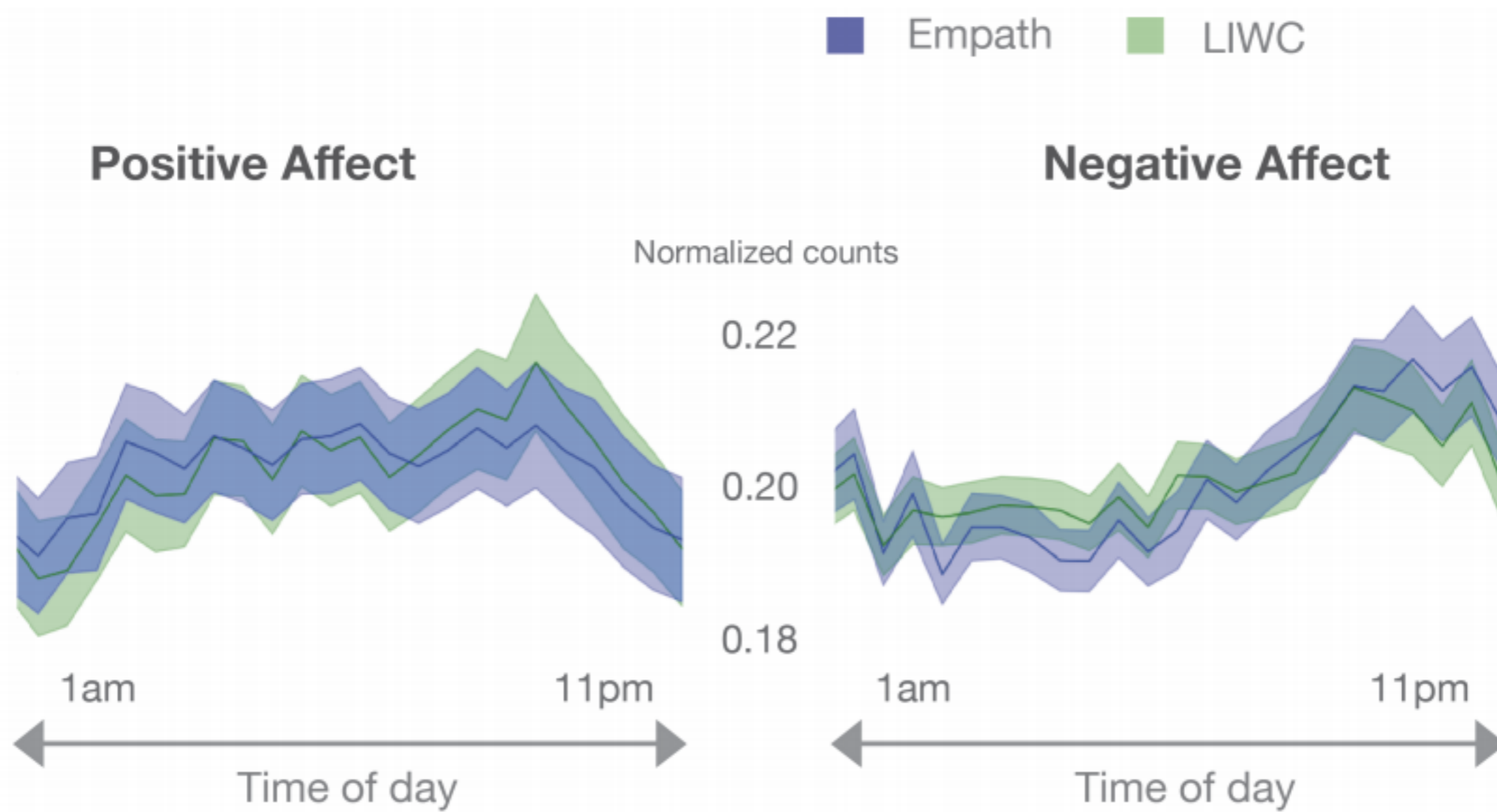
agreement with LIWC?

effect of crowdsourcing on  
categories?

| <b>LIWC Category</b> | <b>Empath (r-value)</b> |
|----------------------|-------------------------|
| Positive Emotion     | 0.944                   |
| Negative Emotion     | 0.941                   |
| Sadness              | 0.890                   |
| Anger                | 0.889                   |
| Achievement          | 0.915                   |
| Religion             | 0.893                   |
| Work                 | 0.859                   |
| Home                 | 0.919                   |
| Money                | 0.902                   |
| Health               | 0.866                   |
| Sex                  | 0.928                   |
| Death                | 0.856                   |
| Average              | 0.900                   |

| <b>LIWC Category</b> | <b>Empath (r-value)</b> |
|----------------------|-------------------------|
| Positive Emotion     | 0.944                   |
| Negative Emotion     | 0.941                   |
| Sadness              | 0.890                   |
| Anger                | 0.889                   |
| Achievement          | 0.915                   |
| Religion             | 0.893                   |
| Work                 | 0.859                   |
| Home                 | 0.919                   |
| Money                | 0.902                   |
| Health               | 0.866                   |
| Sex                  | 0.928                   |
| Death                | 0.856                   |
| Average              | 0.900                   |

Baseline: LIWC agrees at  $r=0.899$  with General Inquirer and  $r=0.876$  with Emolex



strongly correlated positive ( $r=0.87$ ) and negative ( $r=0.90$ ) sentiment (golder and macy, 2011)



| <b>LIWC Category</b> | <b>Empath (r-value)</b> | <b>Em+Crowd (r-value)</b> |
|----------------------|-------------------------|---------------------------|
| Positive Emotion     | 0.944                   | 0.950                     |
| Negative Emotion     | 0.941                   | 0.936                     |
| Sadness              | 0.890                   | 0.907                     |
| Anger                | 0.889                   | 0.894                     |
| Achievement          | 0.915                   | 0.903                     |
| Religion             | 0.893                   | 0.908                     |
| Work                 | 0.859                   | 0.820                     |
| Home                 | 0.919                   | 0.941                     |
| Money                | 0.902                   | 0.878                     |
| Health               | 0.866                   | 0.898                     |
| Sex                  | 0.928                   | 0.935                     |
| Death                | 0.856                   | 0.901                     |
| Average              | 0.900                   | 0.906                     |

demo: CHI and CSCW abstracts

<http://hci.st/empath>

```
pip install empath
```

This work is supported by an NSF Fellowship and the Stanford Medical Scientist Training Program.