



Carnegie Mellon University  
School of Computer Science

# Interpreting Social Media

**Elijah Mayfield**

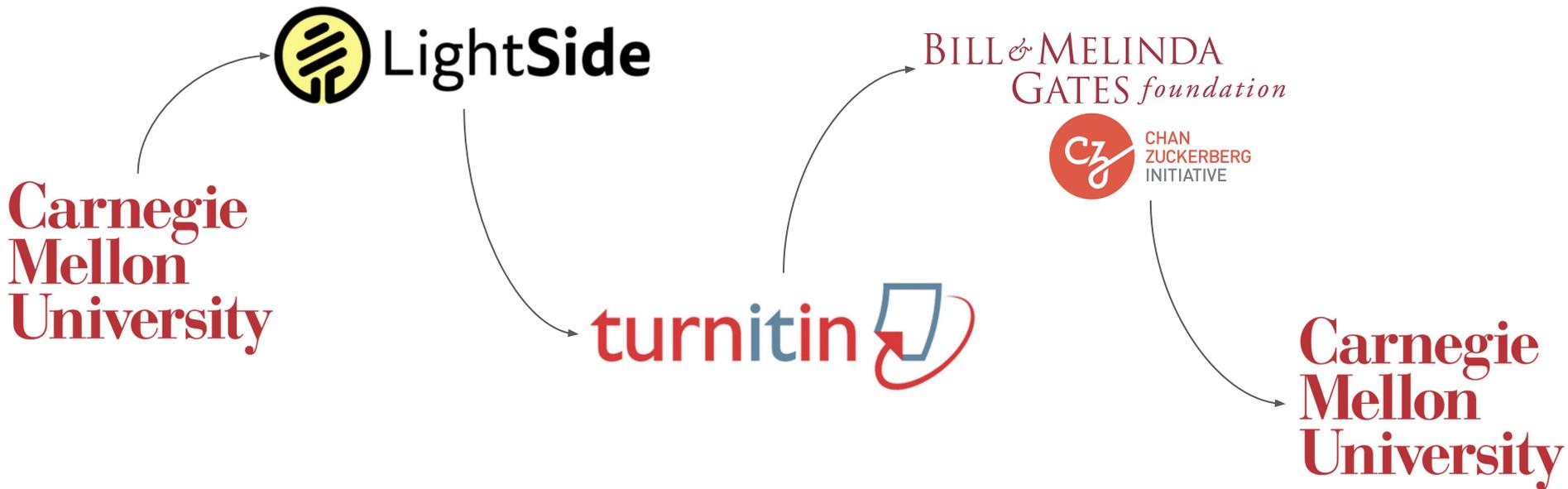
School of Computer Science  
Carnegie Mellon University  
elijah@cmu.edu

*(many slides borrowed with permission from **Diyi Yang, CMU → Google AI → GaTech**)*

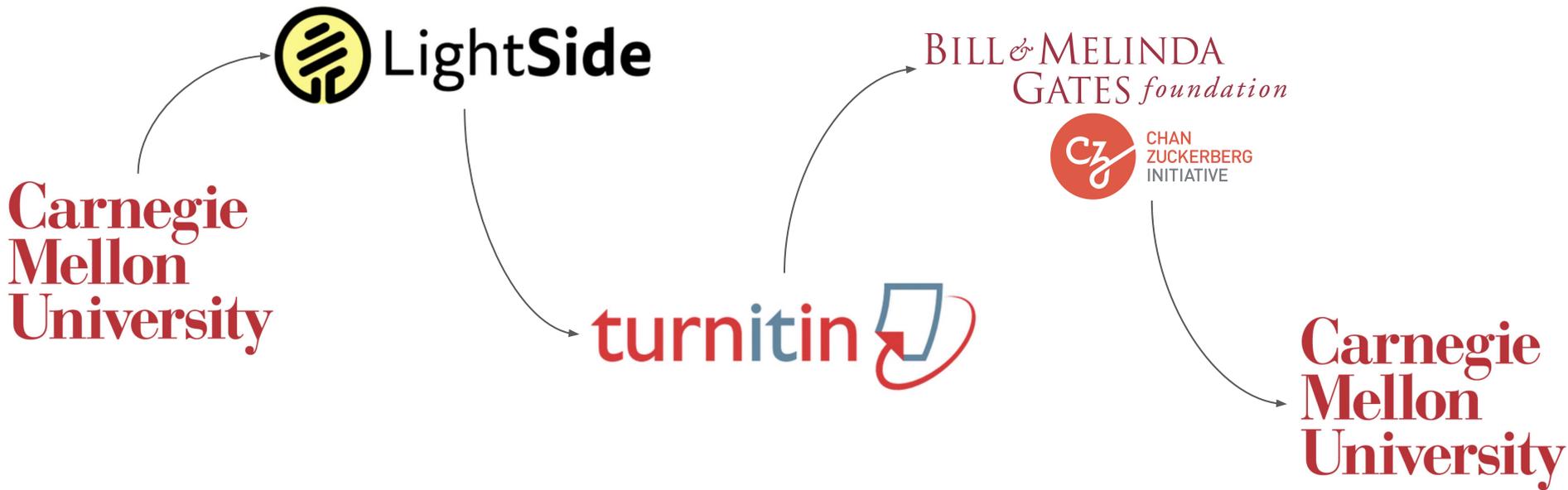
# Lecture Goals

1. Understand what it looks like to apply NLP on real-world data
  - What's different about online data compared to cleaner problems like newswire text?
  - What questions are you going to have to answer as part of working with online data?
  
2. What does a research project on social media data look like?
  - How are the projects designed and what are their goals?
  - What kind of findings we do come up with using NLP today?

# About Me



# About Me



Language Technologies Institute  
Ph.D. Student



Project Olympus / Swartz Center  
Entrepreneur-in-Residence

# Lecture Goals

1. Understand what it looks like to apply NLP on real-world data
  - What's different about online data compared to cleaner problems like newswire text?
  - What questions are you going to have to answer as part of working with online data?
  
2. What does a research project on social media data look like?
  - How are the projects designed and what are their goals?
  - What kind of findings we do come up with using NLP today?

Social Media generates

**BIG UNSTRUCTURED NATURAL LANGUAGE DATA**



Social Media generates

## **BIG UNSTRUCTURED NATURAL LANGUAGE DATA**

### **Volume**

2 billion  
monthly active  
FB users

### **Velocity**

2 Wikipedia  
revisions per  
sec

### **Variety**

tweets, articles,  
discussions,  
news

# What's different about online data?

- NLP researchers love benchmark corpora and standardized tasks
  - Preprocessing takes forever
  - Easy to measure improvement compared to prior approaches
  - Collection, transcription, annotation is unbelievably expensive.

*(computer vision believes all of these things even more than NLP does)*

# What's different about online data?

- NLP researchers love benchmark corpora and standardized tasks

***“Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29. Mr. Vinken is chairman of Elsevier N.V., the Dutch publishing group. [...]”***

# What's so different about online data?

- NLP researchers love benchmark corpora
  - (computer vision researchers love them even more)
- But for most applied work, you are going to be taking in unknown / weird text



# What's so different about online data?

- NLP researchers love benchmark corpora
  - (computer vision researchers love them even more)
- But for most applied work, you are going to be taking in unknown / weird text



HyperHampster **Mozambique Here!** 🦊 38 points · 6 hours ago

Man, whenever I get to top 2 the other squad is a full on Shroud, Dizzy, Viss team while i'm trying to carry tweedle dee and tweedle dum.

Reply Give Award Share Report Save



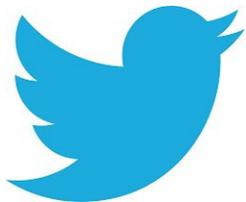
senzung 9 points · 4 hours ago

It's usually the 3rd or 4th squad, last squad are mostly blue armour campers.

Reply Give Award Share Report Save

# Formality online (and elsewhere) is a continuum

- Language varies based on who you're talking to and what you're doing.
- People are really good at “reading the room” and switching styles!
- NLP mostly does not have this ability on the fly yet, needs to be trained.



Medium



# Group Exercise: Spot the Difference

Welcome to the chat room!

dared1ev11: contrpick 🤪🤪

Searix: k

flexin\_on\_grandma: Get a haircut

Nodusman: someone won

HeyYoPancake: @dogdog have you tid ressurect priest? it's already good/fun but you could prolly make the current lists even better 😊

PickleVL: @dogdog How about imp warlock? XD

mdosrobbot: just play emeriss you coward

simbasir: Lul

Best\_gnar\_demaglio: @zeofar Dog won

zedkumo: @PhilFryer I want to see it

MaxThaGr8: nostam had a spicy meme with darkest hour

★ TholasLeon  
Subscribed with Twitch Prime

Deqnkata: zeofar we hit him in the face 😊

xbanng: Counter

MTGHELLION: are there even any good paladins

Morvalus: Hunting party hallazi and mountain giants

xbanng: Same player

 **Lin-Manuel Miranda** @Lin\_Manuel · 21m  
Grandma Tala just stunting on tourists ❤️❤️❤️

 **Diane Duane** @dduane  
Cripes, it's right out of MOANA. I wonder if .@Lin\_Manuel has seen this?  
[twitter.com/terrinakamura/...](https://twitter.com/terrinakamura/)

12 121 1.4K

 **Lin-Manuel Miranda** @Lin\_Manuel · 7h  
Long drive today, and I'm driving, see y'all tonight.

108 83 3.9K

 **Lin-Manuel Miranda** @Lin\_Manuel · 8h  
¡Tenemos un regalo para ti!  
¿Para mi?  
¡Claro que sí!  
Una copia del libro ¡Buen Día, Buenas Noches! De hecho, tenemos 10. Desde hoy, el 22 de abril, 2019, hasta el 29, ve al [sweeps.penguinrandomhouse.com/enter/buen-dia...](https://sweeps.penguinrandomhouse.com/enter/buen-dia...) Se escogerán basado en tus entradas. @VintageEspañol  
Buena suerte! 🍀🍀🍀🍀🍀

Translate Tweet

 **¡Buen día, buen concurso!** | [penguinrandomhouse.c...](https://penguinrandomhouse.com)  
Lin-Manuel Miranda y Vintage Español te quieren regalar una copia firmada de ¡Buen día, buenas noches! ¿Qué mejor forma de comenzar y terminar el día que con un...  
[sweeps.penguinrandomhouse.com](https://sweeps.penguinrandomhouse.com)

37 94 1.4K

Parking enforcement officers return to the car after the posted time for parking has passed, and if the chalk marks are still there—a sign that the vehicle has not moved—the officer issues a citation. Alison Taylor, a frequent recipient of parking tickets, sued the City and its parking enforcement officer Tabitha Hoskins, alleging that chalking violated her Fourth Amendment right to be free from unreasonable search. The City moved to dismiss the action. The district court granted the City's motion, finding that, while chalking may have constituted a search under the Fourth Amendment, the search was reasonable. Because we chalk this practice up to a regulatory exercise, rather than a community-caretaking function, we **REVERSE**.

## I. BACKGROUND

Between 2014 and 2017, Tabitha Hoskins chalked Taylor's tires on fifteen separate occasions and issued her citations in kind. Each citation included the date and time the chalk was placed on her vehicle's tires. The cost of a citation starts at \$15 and increases from there.

On April 5, 2017, Taylor filed this 42 U.S.C. § 1983 action against the City, alleging defendants violated her Fourth Amendment right against unreasonable searches by placing chalk marks on her tires without her consent or a valid search warrant. Taylor also sued Hoskins in her individual capacity. The defendants filed a motion to dismiss pursuant to Fed. R. Civ. P. 12(b)(6), asserting that chalking was not a search within the meaning of the Fourth Amendment, or alternatively, if it was a search, it was reasonable under the community caretaker exception.<sup>1</sup> Hoskins also asserted a qualified immunity defense.

# Group Exercise: Spot the Difference

What differences are easy to spot?

- [answers go here]
- [and here]
- [and here]

What differences are less obvious?

- [answers go here]
- [and here]

# Existing NLP for Social Media is... not good yet?

## ➤ **Machine Translation**

- Works for EN-FR in parliamentary documents
- Not so great for translating posts from Urdu Facebook

## ➤ **Part-of-Speech Tagging**

- Very nearly perfect for Wall Street Journal newstext
- Still plenty of work to do for Black Twitter

## ➤ **Sentiment Classification**

- Works for thumbs-up/down movie reviews
- Pretty bad at complex emotions, short chats, topical humor

# Lecture Goals

1. Understand what it looks like to apply NLP on real-world data
  - What's different about online data compared to cleaner problems like newswire text?
  - **What questions are you going to have to answer as part of working with online data?**
  
2. What does a research project on social media data look like?
  - How are the projects designed and what are their goals?
  - What kind of findings we do come up with using NLP today?

# What are **common tasks** in social media?

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

# Each task is composed of a **pipeline** of subtasks

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

Overlapping geographic locations, events

Identifying shared habits, mutual interests

Moods and mental health (e.g., depression)

Demographic attributes (gender, race, language)

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

# Each task is composed of a **pipeline** of subtasks

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

Factoid Extraction / Stance Classification

Formality / Politeness / Discourse Analysis

Source Reputation Ranking

Virality / Graph analytics

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

# Each task is composed of a **pipeline** of subtasks

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

Linguistic accommodation

Behaviors tied to retention

Homogeneity of population

Social roles / leadership

# Why do universities work on social media?

- It's incredibly convenient.
  - Data collection is **expensive**! Crawled/open data is free, relatively fast.
  - IRB approval for human subjects research is **slow**; public social media data (Twitter, Wikipedia, IMDB) is typically exempt or expedited.
- It acts as a “model organism.”
  - Looks more like real language in use than WSJ.
  - Fairly rapid transition to industry interventions.
  - Multilingual by nature in some cases.



# Why do **companies** fund the work?

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

Some tasks improve a site's **engagement** - companies get a direct, measurable outcome.

# Why do **companies** fund the work?

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis ←
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction) ←
- Learning to Rank (Information Retrieval / Search Engines) ←
- New member dynamics (Longitudinal/Survival analysis) ←

Some tasks are about **profiling** your user demographics and their intent.

Knowing who your users are, and what they want, lets you make your site more relevant.

# Why do **companies** fund the work?

## ➤ **Unsupervised Tasks**

- Trending Topic Clustering / Detection
- Friend / Article Recommendation

## ➤ **Classification Tasks**

- Sentiment Analysis
- “Fake News” Identification
- Hateful Content / Cyberbullying Detection

Some tasks are about preserving **reputation** - if your site is toxic and unmanaged, your community of users will abandon you for alternatives.

## ➤ **Structured Tasks**

- Text generation (Article Summarization)
- Knowledge base population (Information Extraction)
- Learning to Rank (Information Retrieval / Search Engines)
- New member dynamics (Longitudinal/Survival analysis)

# What's not guaranteed?

## ➤ User perceived value

**Snapchat base to decline for first time as users defect to Instagram**

## ➤ Legal accountability

*Facebook Engages in Housing Discrimination With Its Ad Practices, U.S. Says*

## ➤ Answers from the class

- [go here]
- [and here]
- [and here]

## ➤ University motives

- Convenient
- Authentic
- Generalizable

## ➤ Industry motives

- Engagement
- Profiles
- Reputation

# Summary of Part 1

- There are **enormous open opportunities** for NLP developers and scientists.
  - Difficult new domains for NLP models to improve.
  - Interesting, entwined pipelines of tasks that all need to work together.
  - Support from both academia and industry.
- But **blind spots** in task definition and data selection carry significant **risks**:
  - Data selection early in the field limited which language ‘worked’ with NLP tools; the lack of accessibility lasted decades (to today!)
  - Some tasks can put marginalized populations directly in harm’s way.

# Actionable Steps

- Identify **what population is represented** in your data.
  - *Who are your users? How do they self-identify?*
- Design and develop from a place of **deep expertise about that population**.
  - *Easiest, best way to do this: Make sure they're on your team!*
- **Make your goals explicit** about your NLP tools early and often.
  - *Why are we doing this? What metric will go up or down if we do/don't?*

# Break

Questions?

Part 2 (to come):

- Example project: Social Role Modeling on the Cancer Support Network

# Lecture Goals

1. Understand what it looks like to apply NLP on real-world data
  - What's different about online data compared to cleaner problems like newswire text?
  - What questions are you going to have to answer as part of working with online data?
  
2. **What does a research project on social media data look like?**
  - How are the projects designed and what are their goals?
  - What kind of findings we do come up with using NLP today?

# Modeling **Social Roles** in Online Cancer Support Groups - *Cancer Survivor Network*

---

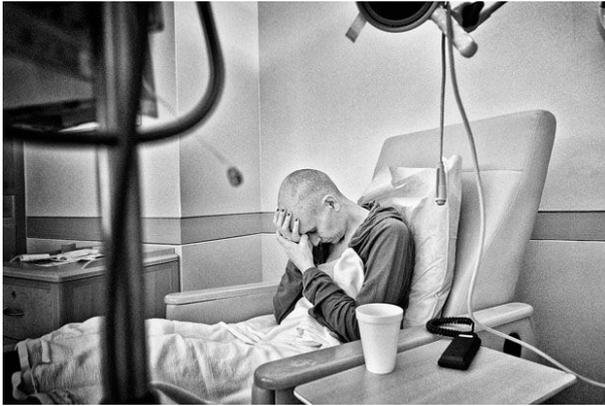
Diyi Yang, Robert Kraut, Tenbrock Smith, Elijah Mayfield, Dan Jurafsky. “Seekers, Providers, Welcomers, and Storytellers: Modeling Social Roles in Online Health Communities”. Proceedings of SIGCHI 2019.

## Problem Statement

**28%** of Internet users have used online support group for medical information (Fox 2009)

**How can NLP support patients and families?**

# Lots of Online Health Support Communities





I was diagnosed with Invasive Ductal Carcinoma grade 2. I'm told I will need chemo. I don't understand. Any words of that will help me wrap my head around this nightmare?

... Since you are a triple positive they can put you on hormones and the chance of recurrence is low. Listen to your chemo nurse ...



Sorry to hear..God bless you ..stay strong



This conversation has been paraphrased.

# Project History

- Long-studied research area (they all are)
- **Previous work:**
  - What do users of support groups **do**?
    - What kind of information do they share?
    - Which strategies reduce stress, promote self-efficacy?
  - Which users decide to stay?
    - What is the “lifecycle” of a user?
    - What events happen during those lifecycles, online or off?
- **New question:** what roles do users play?



I love your attitude. It gives me faith that you can have cancer, live a full life and have children. You give me hope and faith.

Emotional Support



Has your doctor tested your tumor for the Oncotype score? I believe they now do it on hormone negative tumors

Informational Support

# 9 Conversational Acts in Health Support Groups



Emotional Support

1. Seeking emotional support
2. Providing emotional support
3. Providing empathy
4. Providing appreciation
5. Providing encouragement



Informational Support

6. Seeking informational support
7. Providing informational support



Self-disclosure

8. Disclosing oneself positively
9. Disclosing oneself negatively

# Dataset: Text-based Cancer Support Groups

13-year data since 2005

66K users, 140K threads and 1.3M replies



## Cancer Survivors Network

**CSN Login** Username  Password   [Forgot username or password?](#)

**CSN**

- Discussion Boards
- Announcements
- Member Resource library

**CSN Home**

**Discussion boards**

- [Log in](#) to post new content in the forum.

# Dataset Construction

How much **self-disclosure and social support** does this message contain?

- Likert Scale: 1 (None) to 7 (a great deal)
- 1000 messages
- High reliability ( $r=0.92$ )

## Text to Features

| Feature Type   | Sample Feature Explanation |
|--|----------------------------|
| Generic Linguistic Inquiry and Word Count (Pennebaker, 1997) | I, my, we, our             |

“I love **your** attitude. It gives **me** faith that **you** can have cancer, live a full life and have children. **You** give **me** hope and faith. **You** are the greatest. ”

## Text to Features

| Feature Type   | Sample Feature Explanation |
|--|----------------------------|
| Generic Linguistic Inquiry and Word Count (Pennebaker, 1997) | I, my, we, our             |
| Topic Modeling (Wang et al., 2015)                           | Diagnose, treatment        |

“I was *diagnosed* with stage 2 *triple negative* with no *lymph* node involvement. I had the Red Devil first then 23 *radiations*.”

# Text to Features

| Feature Type   | Sample Feature Explanation      |
|--|---------------------------------|
| Generic Linguistic Inquiry and Word Count (Pennebaker, 1997) | I, my, we, our                  |
| Topic Modeling (Wang et al., 2015)                           | Diagnose, treatment             |
| Named Entity Recognition                                     | Person, organization, location  |
| Medicine/symptom via Freebase                                | Medicine, symptom names         |
| Word Embedding (medical domain)                              | Distributional semantic meaning |

# Predicting Conversational Acts in Messages

| 9 Conversational Acts              | Correlation (human, prediction) |
|------------------------------------|---------------------------------|
| Seeking informational support      | 0.729                           |
| Providing informational support    | 0.793                           |
| Seeking emotional support          | 0.637                           |
| Providing emotional support        | 0.748                           |
| Providing empathy                  | 0.723                           |
| Providing appreciation             | 0.669                           |
| Providing encouragement            | 0.641                           |
| Self-disclosing oneself positively | 0.719                           |
| Self-disclosing oneself negatively | 0.712                           |

# Modeling Social Roles on CSN



1. What roles do people occupy?
2. How do roles influence members' participation?

# Modeling Social Roles on CSN



## 1. What roles do people occupy?

### Methods:

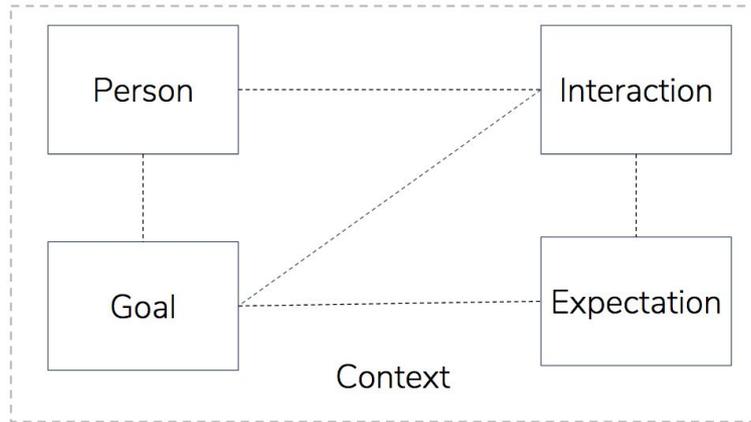
- Gaussian Mixture Model that identifies functional roles
- Interviews with active users, moderators, and clinicians for validation

## 2. How do roles influence members' participation?

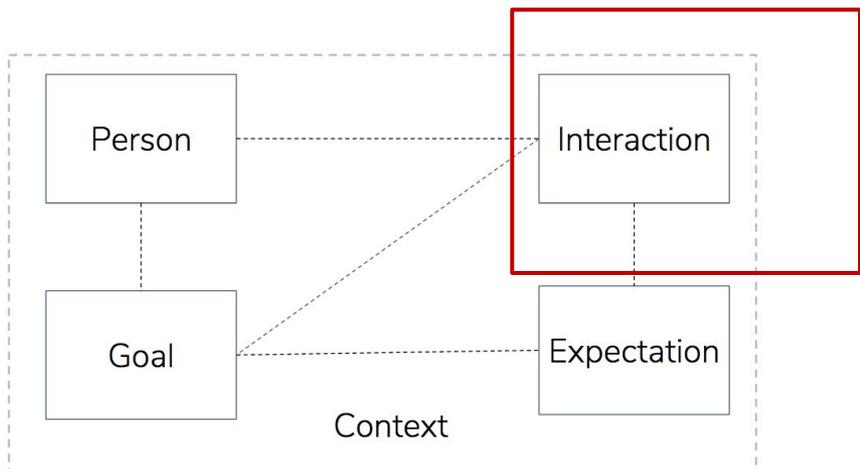
# Modeling Social Roles via Mixture Model

- Behavioral representation (features,  $\mathbf{X}$ )
- Observed user session structure
- The number of implicit roles  $\mathbf{K}$  (will

# Behavioral Representation X



# Behavioral Representation: Interaction

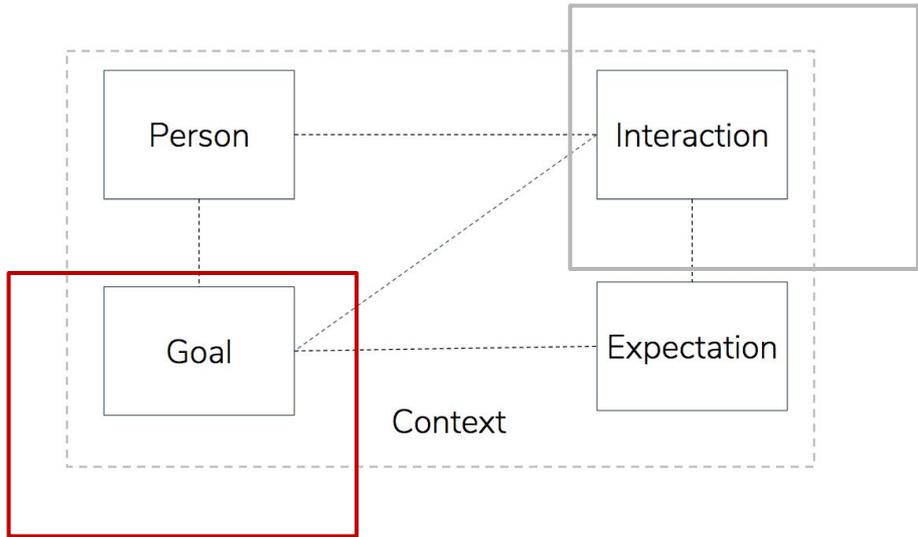


Network-based measures

Linguistic-based measures

- Emotional aspects: “anger”, “sadness”
- Social concerns: “friend”, “family”
- Self-focus: “I”, “you”, “he/she”
- Topics modeling

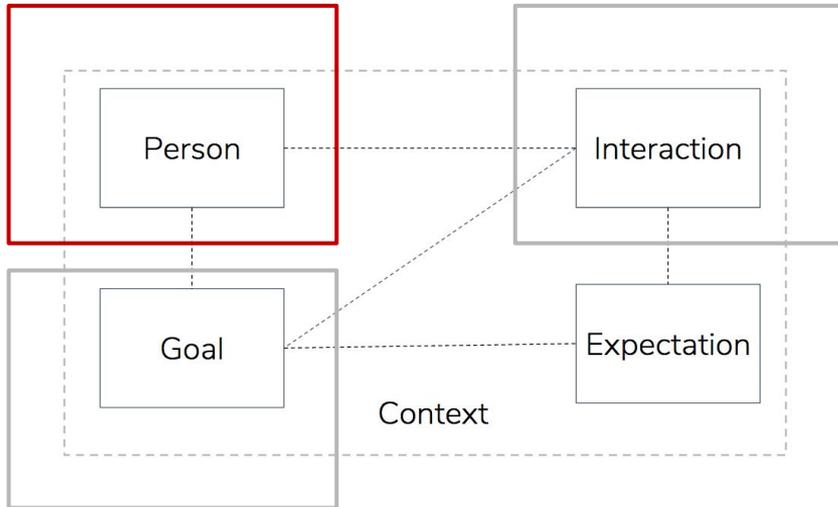
# Behavioral Representation: Goal



(Dindia +, 2002; Cohen and Syme, 1985)

- Seeking emo support
- Providing emo support
- Providing empathy
- Providing appreciation
- Providing encouragement
- Seeking info support
- Providing info support
- Disclosing oneself positively
- Disclosing oneself negatively

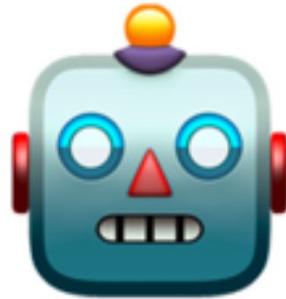
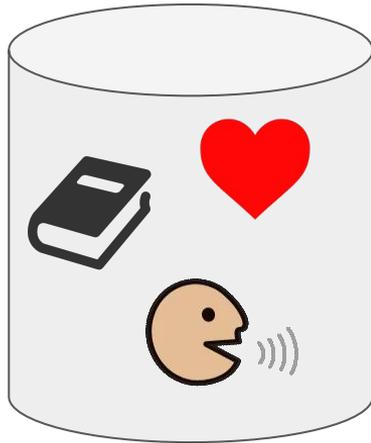
# Behavioral Representation: **Person**



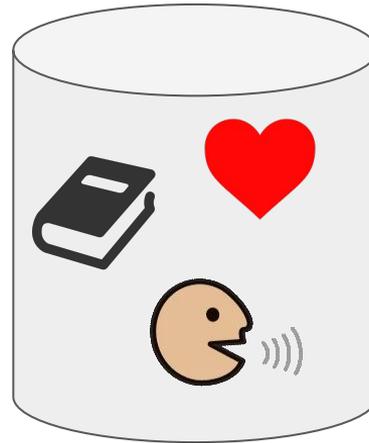
Infer users' attributes, including  
**gender, cancer status, and  
cancer type**  
based on their conversations

# Privacy-Preserving Modeling of DMs

Public Data



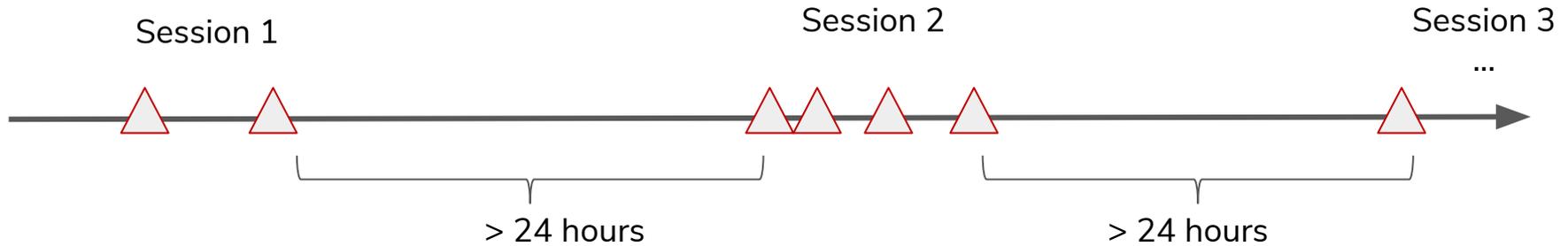
Private Data



- No human ever reads private data
- Labels are probably (?) still accurate
- Allows modeling to include more kinds of users

# The Length of User Representation

**Session:** a time interval where the time gap between any two adjacent actions in this session is less than  $t$  (e.g., 24 hours)



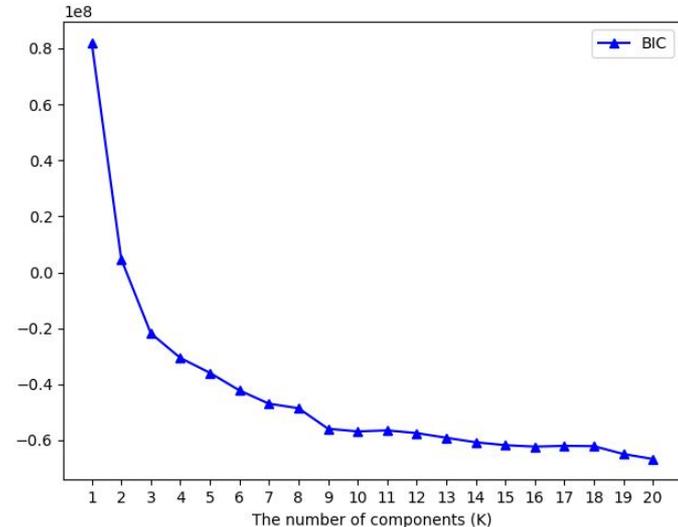
# The Number of Implicit Roles/Components

Quantitatively:

- Vary #components from 1 to 20
- Use BIC score to select models

Qualitatively:

- Validate with 6 moderators to assess the derived roles



# Derived Roles in Cancer Support Groups

|   |   |
|---|---|
|  Emotional Support Provider (33.3%)    |  Private Support Provider (5.3%) |
|  Newcomer Welcomer (15.9%)             |  All-round Expert (2.5%)         |
|  Informational Support Provider (13.3%) |  Newcomer Member (2.4%)          |
|  Story Sharer (10.2%)                  |  Knowledge Promoter (2.2%)       |
|  Informational Support Seeker (8.9%)   |  Private Networker (0.8%)        |
|  Private Communicator (5.3%)           |   |

# Qualitative Evaluation of Derived Roles

Work with 6 moderators on CSN to assess the derived roles



*“ It seems very comprehensive and there are so many different examples, so I feel like it is covered very well with your different roles and labels. ”*

The identified roles were mostly **comprehensive**

# Qualitative Evaluation of Derived Roles

Work with 6 moderators on CSN to assess the derived roles



*“The one that I think did not emerge is the policeman, these people complain to moderators when some people are doing things wrong or tell other people that they are violating norms.”*

Model failed to capture the “defenders”

# Modeling Social Roles on CSN



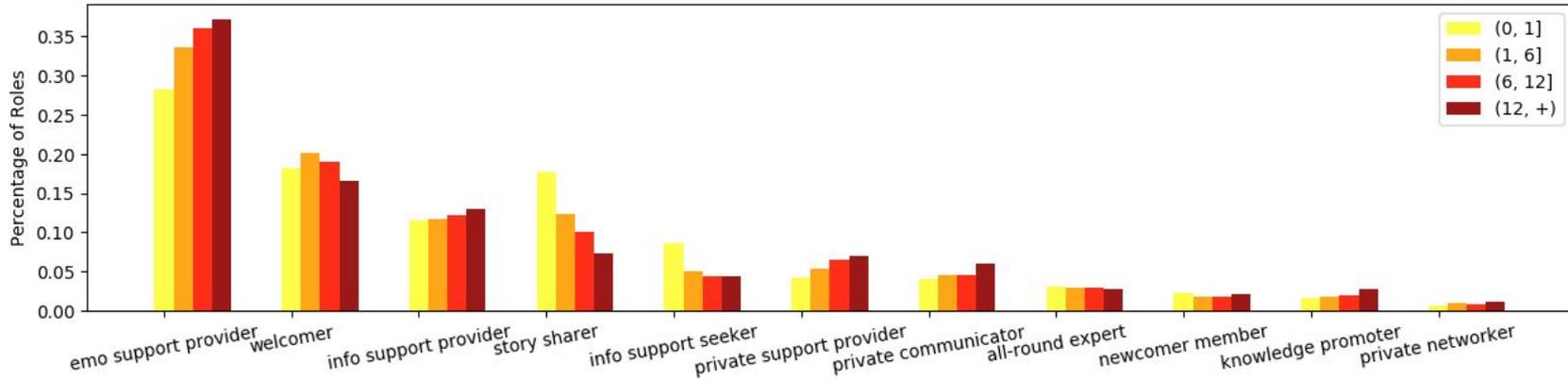
1. What roles do people occupy?

2. How do roles influence members' participation?

## Methods:

- Session-to-session transition matrix analysis
- More interviews with active users

# Dynamics of Role Occupations over Members' Tenure



(0, 1]: Users' first month; (1, 6]: from their second month to six months  
(6, 12]: from six months to a year; (12, +]: after one year

From roles seeking sources to ones offering help

# Top 8 Most Frequent Role Transition Patterns

---

Private communicator → private communicator (41.3% conditional probability)

---

Informational support provider → emotional support provider (36.2%)

---

Emotional support provider → emotional support provider (33.6%)

---

Welcomer → emotional support provider(33.5%)

---

Newcomer member → emotional support provider (33.0%)

---

Informational support seeker → emotional support provider(32.6%)

---

Private networker → private communicator (31.5%)

---

Story sharer → emotional support provider (31.2%)

---

\* Model role transition as a Markov process

# From Roles Seeking Sources to Ones Offering Help

12 interviews of users on Cancer Survivor Network



*I'm now looking for people who are seeking for advice to offer.*

*This message has been paraphrased.*

# From Roles Seeking Sources to Ones Offering Help

12 interviews of users on Cancer Survivor Network



*I initially stayed because information was important, but over time, I found talking with people who had similar experiences is more helpful*

This message has been paraphrased.

# Summary

- Years of research has given us **expectations and categories for behaviors**.
- **Latent behavioral roles** were discovered from our mixture modeling method.
  - *Those roles were **comprehensive and interpretable** by users in interviews.*
- Watching those roles change over time lets us **predict user retention**.
  - *In interviews, those automated discoveries **matched user intuition**.*

# Modeling **Impact** in Online Group Decision-Making - *Wikipedia, The Free Encyclopedia*

---

Elijah Mayfield and Alan W Black. "Stance Classification, Outcome Prediction, and Impact Assessment: NLP Tasks for Studying Group Decision-Making." Proceedings of NLP+CSS Workshop at NAACL 2019.

## **Problem Statement**

Many online communities are full of gatekeeping behaviors, and are difficult to enter and participate in.

**How can NLP open up contribution opportunities for newcomers?**

# Modeling Influence on Wikipedia



WIKIPEDIA  
The Free Encyclopedia

1. What behaviors/moves “work” in editor debates?
2. Who uses those behaviors?

# Modeling Influence on Wikipedia



WIKIPEDIA  
The Free Encyclopedia

## 1. What behaviors/moves “work” in editor debates?

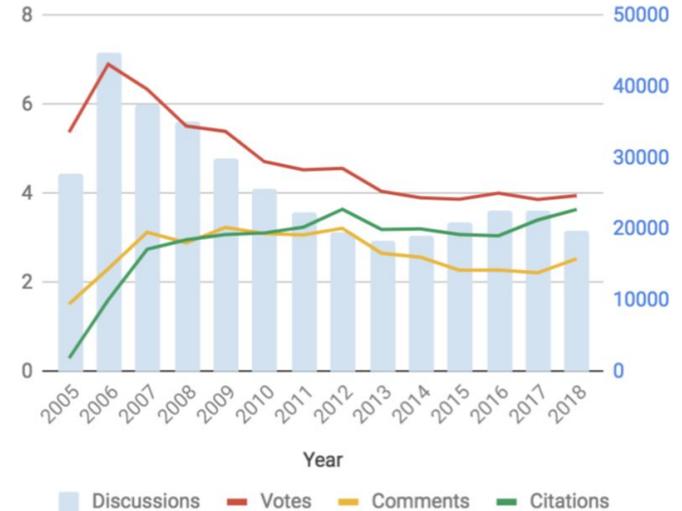
### Methods:

- Information extraction (policies, user tenure)
- Text classification (stance prediction, outcome prediction)
- Longitudinal measurement (macro / micro)

## 2. Who uses those behaviors?

# The Need from Wikipedia's Perspective

- Articles for Deletion - high traffic, dozens of debates per day
  - Articles can be nominated by anyone, with open debate for 7 days
  - Final decisions made by administrators based on discussion
  - High volume but with decline over time since 2007 (like the rest of the site)



# The Need from Wikipedia's Perspective

Wikipedia:Do not disrupt Wikipedia to illustrate a point

Wikipedia:Gaming the system

Wikipedia:Please do not bite the newcomers

Wikipedia:Disruptive editing



Wikipedia:Assume good faith



- Fairly **hostile** environment (from discussion with Wikimedia):
  - Intricate net of policies and guidelines
  - Unwritten or arcane rules about participating
  - Incentives not always aligned with optimal group discussion

# Characteristics of the text

The result was **delete**. - [Daniel Bryant](#) 08:47, 17 March 2007 (UTC)

**Comedian Hypnotist The Incredible BORIS** [ [edit](#) ]

[Comedian Hypnotist The Incredible BORIS](#) ([edit](#) | [talk](#) | [history](#) | [links](#) | [watch](#) | [logs](#) | [views](#)) – ([View log](#))

- **Strong Delete**: Ludicrous Vanity. Delete! Delete! Delete!--[Jack Cox](#) 01:44, 17 March 2007 (UTC)
- **Strong Delete** - complete vanity page, no content. --[Haemo](#) 02:05, 17 March 2007 (UTC)
- **Strong Delete** Per Vanity page, there is no content. [Daniel5127](#) | [Talk](#) 02:39, 17 March 2007 (UTC)
- **Strong Delete**. Vanity. [Interlingua](#) [talk](#) [email](#) 02:44, 17 March 2007 (UTC)
- **Speedy Delete** assertions of notability are ludicrous. A7 this thing. --[NMChico24](#) 02:53, 17 March 2007 (UTC)
- **Comment**: Via Google I found a couple of very small blurbs in local free newspapers about his gigs (basically saying where he would be) but nothing at all that would confirm he was on all those TV shows that are claimed in the WP article and on his site. [LastChanceToBe](#) 03:31, 17 March 2007 (UTC)
- **Delete'** -- Vanity page. [Xdenizen](#) 03:45, 17 March 2007 (UTC)
- **Strong delete** vanity hoax ⇒ [SWATJester](#) [On Belay!](#) 04:14, 17 March 2007 (UTC)
- **Speedy delete** - [vanispamcruftisement](#). So tagged. [MER-C](#) 05:36, 17 March 2007 (UTC)

➤ Some contributions aren't really that helpful.

# Characteristics of the text

- Others produce strong controversy

Wikipedia:Articles for deletion/Justin Bieber on Twitter

---

# Characteristics of the text

➤ Others produce strong controversy

## Wikipedia:Articles for deletion/Justin Bieber on Twitter

The image displays a collage of 18 screenshots from the Wikipedia discussion page 'Wikipedia:Articles for deletion/Justin Bieber on Twitter'. Each screenshot captures a different user's contribution to the discussion, often featuring a list of numbered points or detailed arguments. The comments are scattered across the page, illustrating the 'strong controversy' mentioned in the text. The screenshots show various perspectives on the article's content, format, and notability, with users often providing specific examples or references to support their views. The text is dense and covers a wide range of topics related to the article's deletion, including concerns about redundancy, formatting, and the article's adherence to Wikipedia's guidelines.

# The most clear-headed ones rely on policy

**Edayilakkad** [ edit ]

Quality issues, possibly beyond any fixing. See discussion at talk:, and at

[Wikipedia:Village\\_pump\\_\(proposals\)#I.27ve\\_had\\_enough.\\_.22Approved\\_articles.22\\_clearly\\_no\\_better\\_than\\_ones\\_that\\_skip\\_it](#) [Andy Dingley \(talk\)](#) 00:11, 7 July 2017 (UTC)

- **Keep.** Per [WP:GEOLAND](#) an inhabited island is presumed Notable, and in my opinion that is an almost automatic qualification for inclusion if meaningful information can be verified. While I haven't yet found [WP:GNG](#)'s usual expectation for significant coverage in any particular source, I have been finding a fair number of sources with various brief mentions. Note that source searching is difficult because there are several variations on the spelling, and because useful search results tend to be heavily buried under garbage search results. [Asee \(talk\)](#) 02:35, 7 July 2017 (UTC)

P.S. An inhabited island in the U.S. would almost certainly be kept, and if this is deleted I'm sure it would just get re-created in a few years as India comes more online with sources. [Asee \(talk\)](#) 03:09, 7 July 2017 (UTC)

P.P.S. Here's the article at Malayalam language Wikipedia: [ml:ഇടയിലക്കട്](#). There are a mix of blog-sources as well as usable sources in that version. [Asee \(talk\)](#) 04:11, 7 July 2017 (UTC)

It's not about whether an island is implicitly notable, it's about whether this article passes our standards to adequately demonstrate that. [WP:RS](#) and [WP:V](#) are strong policy. [WP:OTHERSTUFFEXISTS](#) is not. [Andy Dingley \(talk\)](#) 09:11, 7 July 2017 (UTC)

[Andy Dingley](#), I do not disagree with your concerns about quality. However the excessive 11 keeps here indicate that you've missed a significant detail. Your first sentence got it backwards, it is about whether the island is implicitly Notable. Notability isn't a property of the article, it's a property of the topic. An article that contains zero evidence of notability is a Keep, if sources exist and the topic itself satisfies Notability. In the most extreme case you keep the article and delete all the junk down to a single sentence stub. An irredeemably promotional article on a company might get hit with an unsympathetic [TNT](#), but we're going to salvage anything we can for a desirable article on an inhabited place. Documenting significant geography is about as close to objectively-desirable as it gets. [Asee \(talk\)](#) 23:53, 13 July 2017 (UTC)

# Question: which policies are successful?

Wikipedia: Wikipedia is not for things made up one day

---

Wikipedia: No one cares about your garage band

---

Wikipedia: Do not create hoaxes 

---

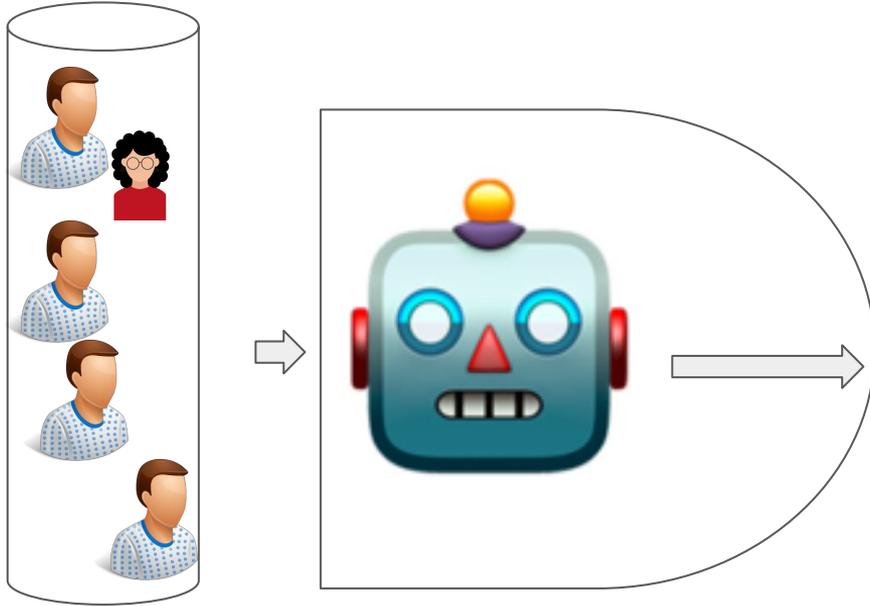


Wikipedia: Vanis spamcraftisement

---

- Policies that **everyone agrees on** win almost all the time.
- But is that really impact?

# Classification task: Outcome Prediction

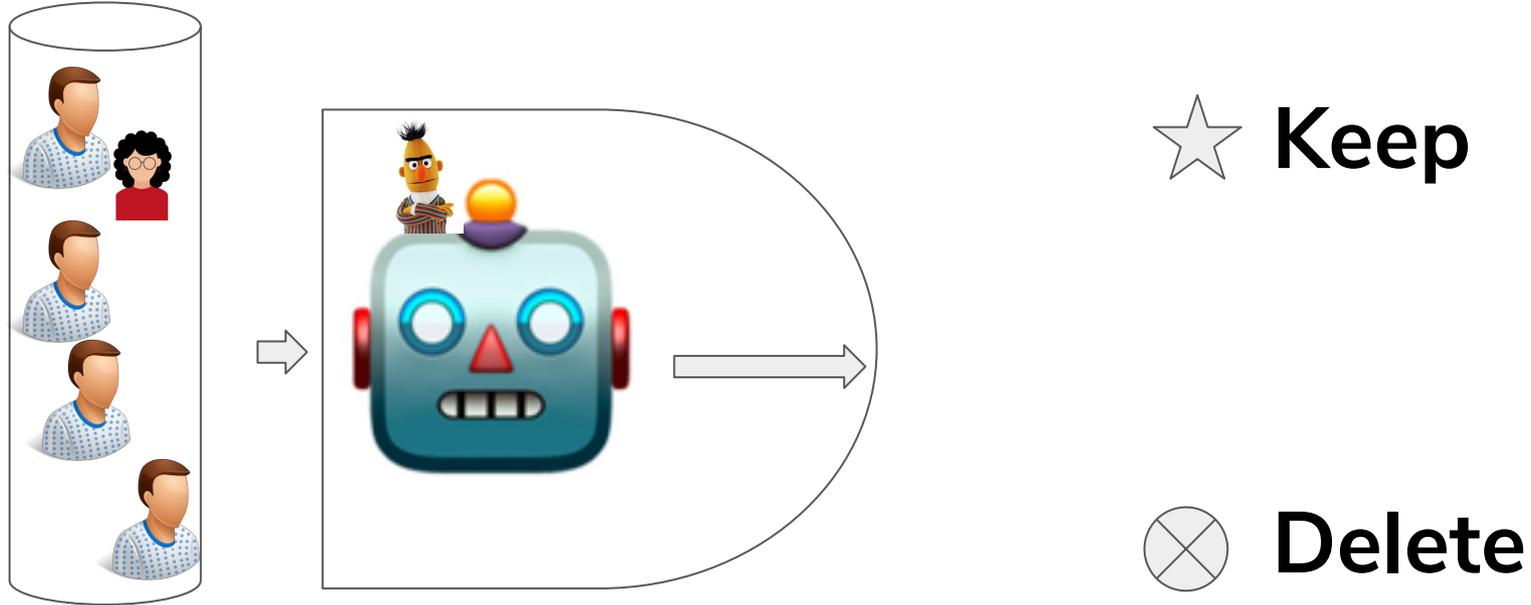


★ Keep

⊗ Delete

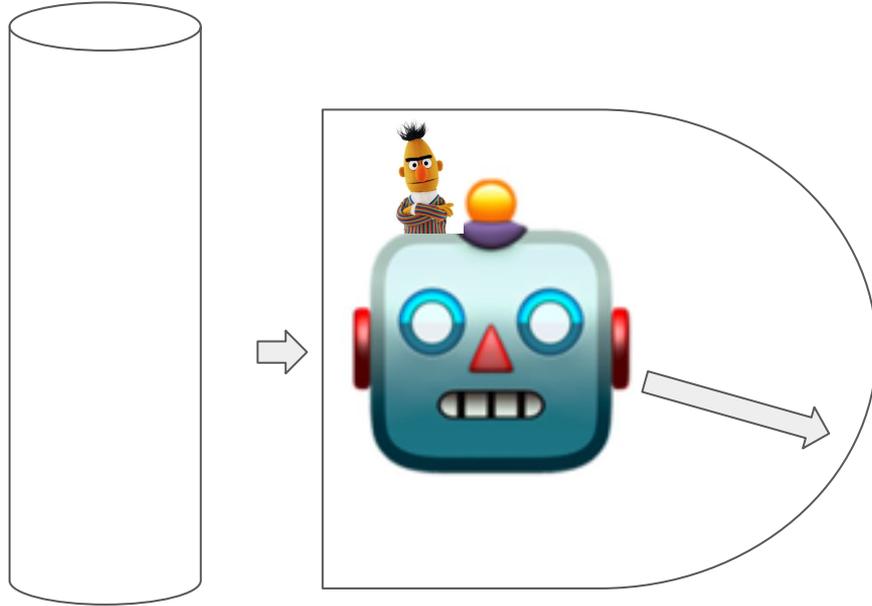
➤ Can we look at the debate and **predict** the final decision? (yes)

# Classification task: Outcome Prediction



- Can we look at the debate and **predict** the final decision? (yes)

# Classification task: Outcome Prediction

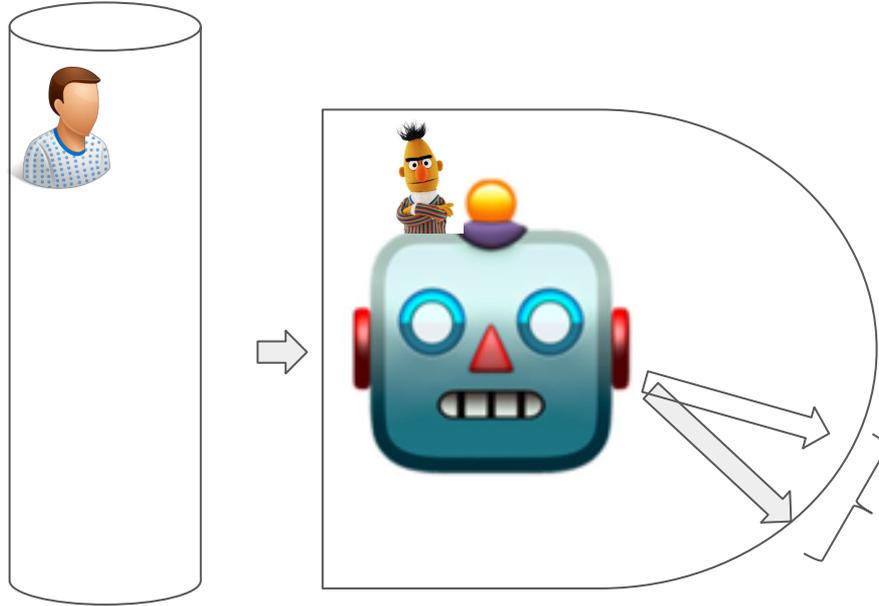


★ **Keep**

⊗ **Delete**

- Measure probabilities moment-by-moment to get impact?

# Classification task: Outcome Prediction

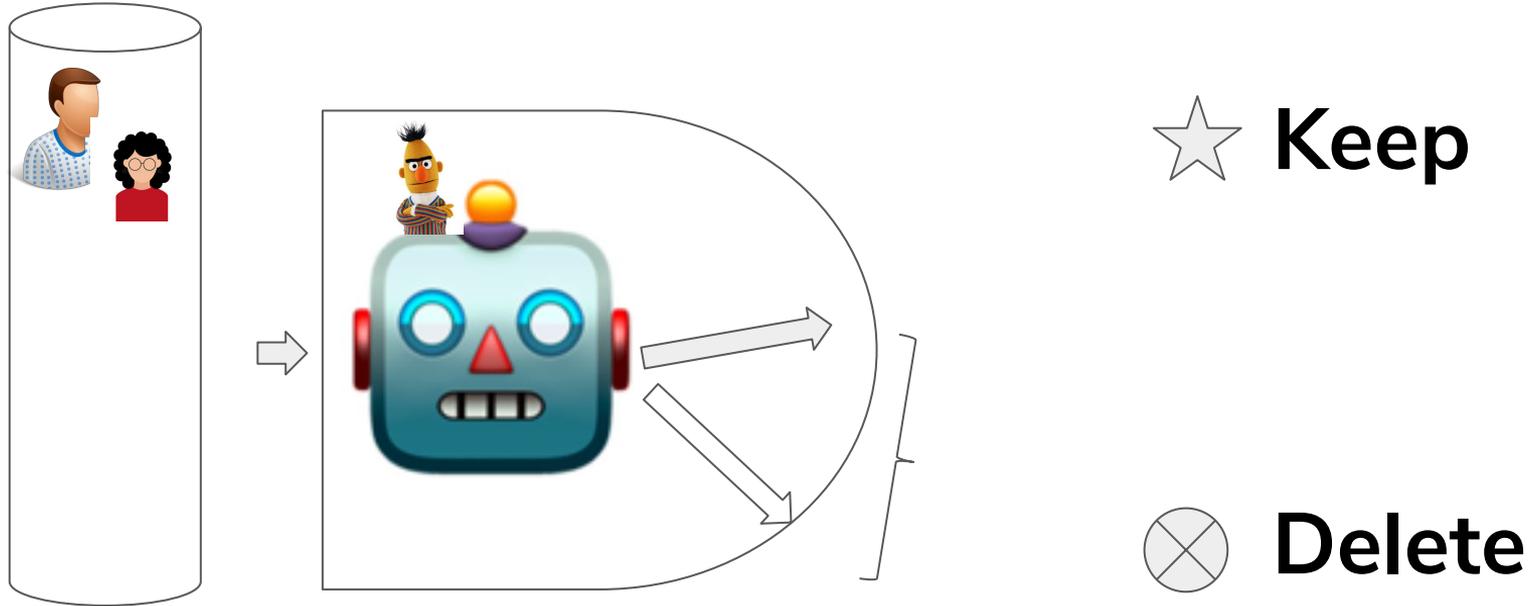


★ Keep

⊗ Delete

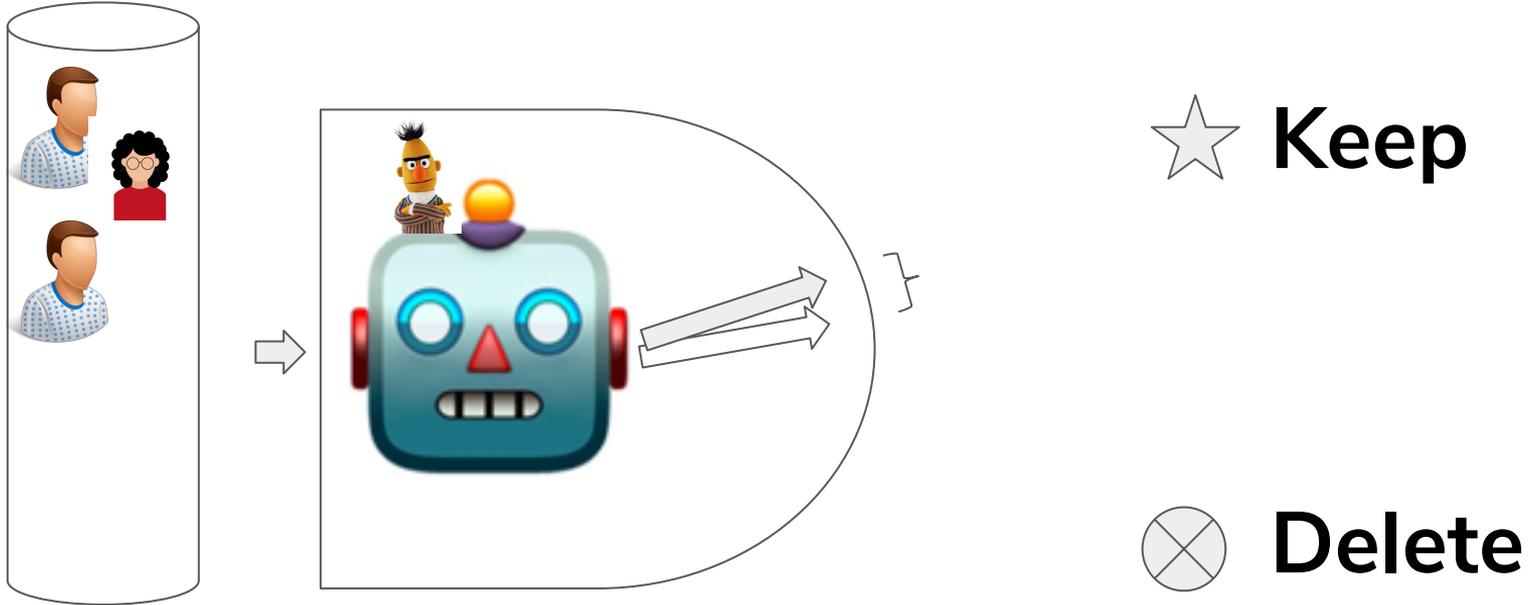
- Measure probabilities moment-by-moment to get impact?

# Classification task: Outcome Prediction



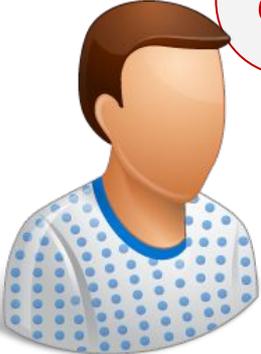
- Measure probabilities moment-by-moment to get impact?

# Classification task: Outcome Prediction



- Measure probabilities moment-by-moment to get impact?

# Question: which are successful and impactful?

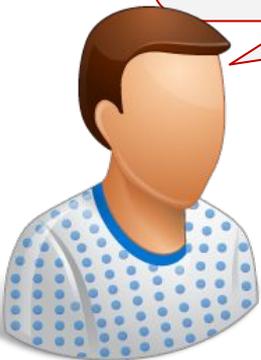


Keep. Per **WP:GEOLAND** an inhabited island is presumed Notable, and in my opinion that is an almost automatic qualification for inclusion

- Specific policies that domain experts can lean on for structural support.

# Question: which are successful and impactful?

- Specific policies that domain experts can lean on for structural support.

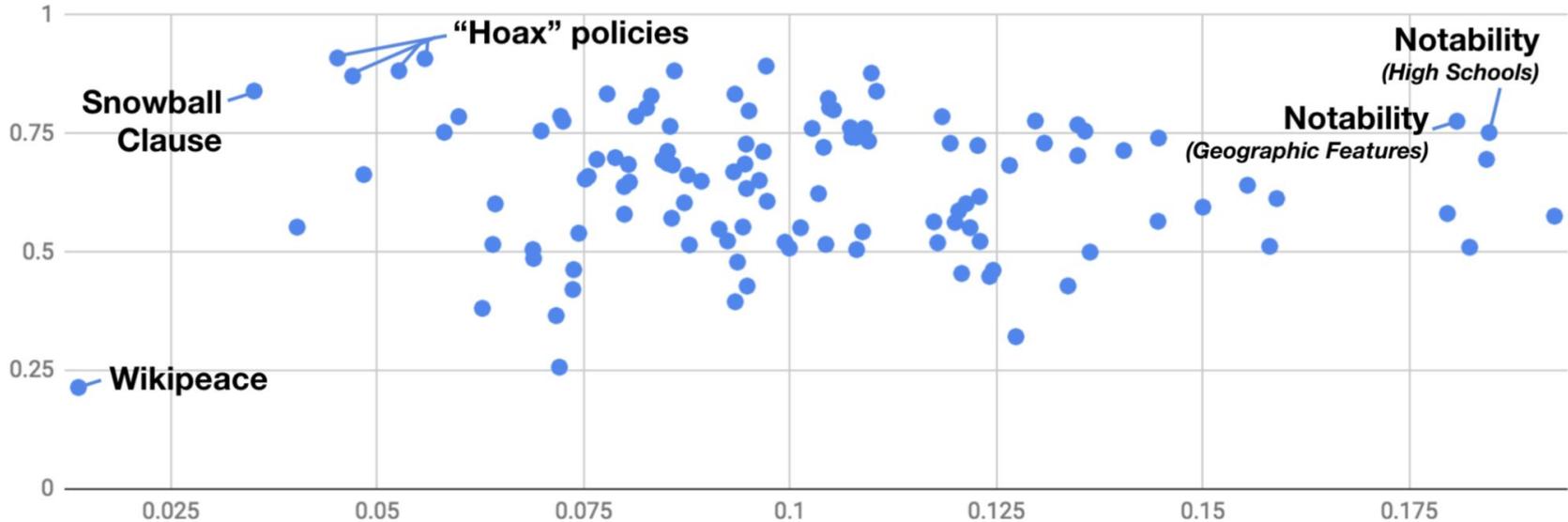


Delete. Non-notable badminton player. Lacks **WP:GNG** to justify an article



Keep Notable badminton player. meet **WP:NBADMINTON** #2 and 3.

# Modeling Influence on Wikipedia



- This gives us a meaningful, interpretable space of impact measurement for studying individual posts, strategies, users

# Modeling Influence on Wikipedia



WIKIPEDIA  
The Free Encyclopedia

1. What behaviors/moves “work” in editor debates?
- 2. Who uses those behaviors?**

# Fall 2019 Project

- Use API to extract user characteristics from public **self-disclosed profiles**.
- Align public profiles to **participation in debates**.
- **Measure correlations** between impactful behaviors and profile characteristics.
- Use quantitative outcomes to make **design recommendations** for Wikimedia.

# Followup / Contact

- I'm [elijah@cmu.edu](mailto:elijah@cmu.edu)
  
- Topics I know things about:
  - Online data: Wikipedia, Cancer Support Network
  - Educational data: student writing, discussion groups, tutoring systems
  - Fairness and equity topics in NLP
  - Entrepreneurship: Startups, Investing, Grantwriting (especially related to NLP/ML!)