# Detecting Implicitly Harmful Language in Political Discourse

Julia Mendelsohn

University of Maryland

*juliame@umd.edu*

# Content Warning

I will be talking about material that may be offensive and upsetting to some audience members.

*Democrats are the problem. They don't care about crime and they want illegal immigrants, no matter how bad they may be, to pour into and infest our Country, like MS-13. They can't win on their terrible policies, so they view them as potential voters!* ~ President Donald Trump, June 2018

*Democrats are the problem. They don't care about crime and they want illegal immigrants, no matter how bad they may be, to pour into and infest our Country, like MS-13. They can't win on their terrible policies, so they view them as potential voters!*
~ President Donald Trump, June 2018

*Good night to everyone but the massive amount of Soros-funded illegals who are trying to invade our border.*
~ AZ State Sen Wendy Rogers, Sep 2021

*Democrats are the problem. They don't care about crime and they want illegal immigrants, no matter how bad they may be, to pour into and infest our Country, like MS-13. They can't win on their terrible policies, so they view them as potential voters!*
~ President Donald Trump, June 2018

*Good night to everyone but the massive amount of Soros-funded illegals who are trying to invade our border.*
~ AZ State Sen Wendy Rogers, Sep 2021

*Democrats are the problem. They don't care about **crime and they want illegal immigrants, no matter how bad they may be**, to pour into and infest our Country, like **MS-13**. They can't win on their terrible policies, so they view them as potential voters!*
~ President Donald Trump, June 2018

*Good night to everyone but the massive amount of Soros-funded illegals who are trying to **invade our border**.*
~ AZ State Sen Wendy Rogers, Sep 2021

*Democrats are the problem. They don't care about crime and they want illegal immigrants, no matter how bad they may be, **to pour into and infest our Country**, like MS-13. They can't win on their terrible policies, so they view them as potential voters!*
~ President Donald Trump, June 2018

*Good night to everyone but the **massive amount** of Soros-funded illegals who are trying to **invade our border**.*
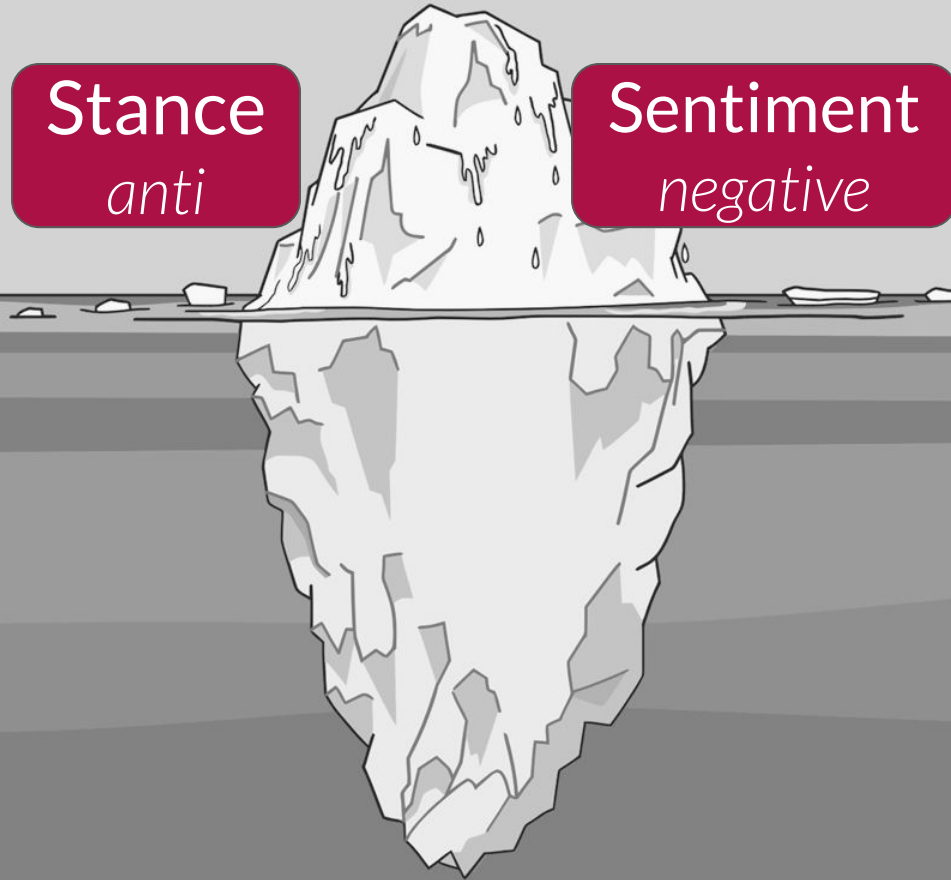~ AZ State Sen Wendy Rogers, Sep 2021

*Democrats are the problem. They don't care about crime and they want illegal immigrants, no matter how bad they may be, to pour into and infest our Country, like MS-13. They can't win on their terrible policies, so they view them as potential voters!*
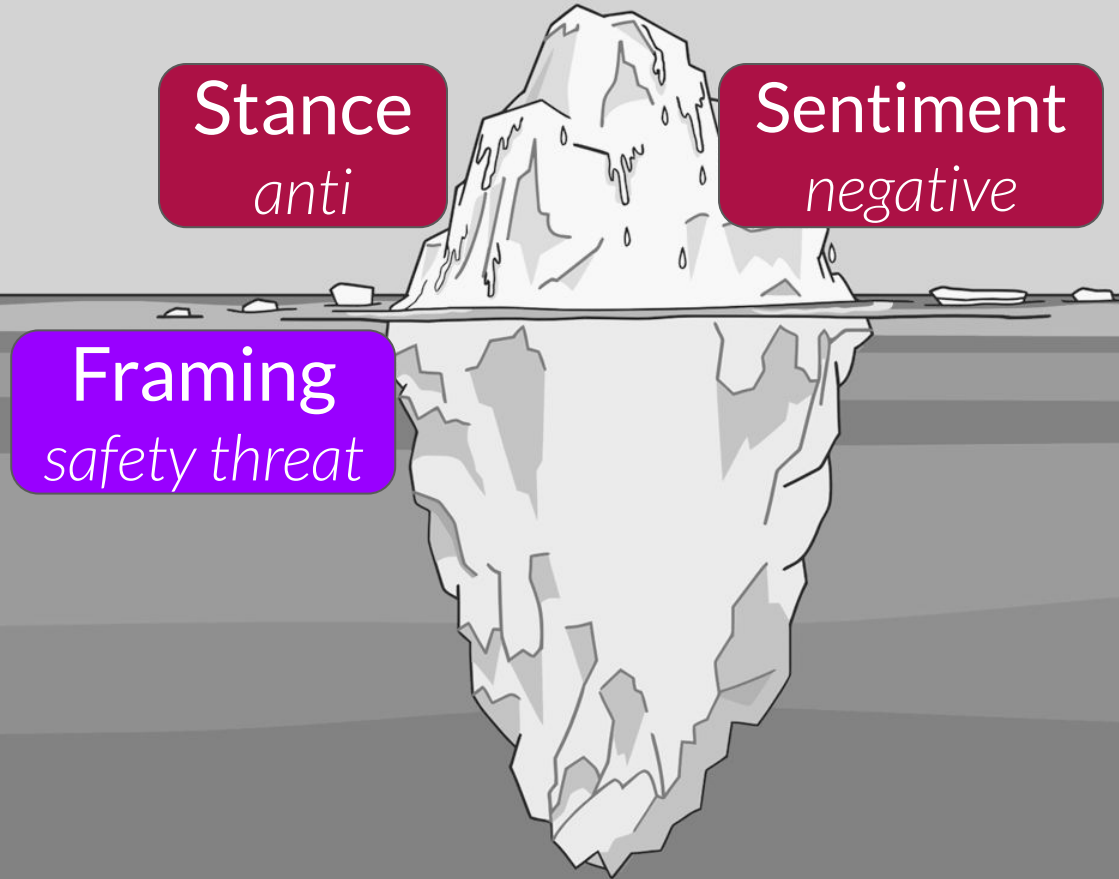~ President Donald Trump, June 2018

*Good night to everyone but the massive amount of* **Soros-funded illegals** *who are trying to invade our border.*
~ AZ State Sen Wendy Rogers, Sep 2021

# The Implicit Iceberg



**Stance**
*anti*

**Sentiment**
*negative*

# The Implicit Iceberg



**Stance**
*anti*

**Sentiment**
*negative*

**Framing**
*safety threat*

The Implicit Iceberg

Stance
*anti*

Sentiment
*negative*

Framing
*safety threat*

Dehumanization & Metaphor
*water, vermin*

Dogwhistles
*Soros (Jewish) plot*

These strategies shape
*how audiences understand*
political issues…

**Framing**
*safety threat*

**Dehumanization
& Metaphor**
*water, vermin*

**Dogwhistles**
*Soros (Jewish) plot*

# ...and are key elements of political communication



**Framing**
*safety threat*

**Dehumanization & Metaphor**
*water, vermin*

**Dogwhistles**
*Soros (Jewish) plot*

# …and are key elements of political communication

# ...and are key elements of political communication



Campaigns
[Tilley, 2020]

# …and are key elements of political communication



Campaigns
[Tilley, 2020]



Media Bias
[Esses et al., 2013]

# ...and are key elements of political communication



Campaigns
[Tilley, 2020]



Media Bias
[Esses et al., 2013]



Misinformation
[Henderson & McCready, 2019]
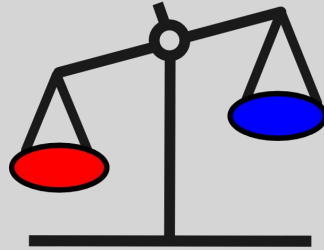
# ...and are key elements of political communication



Campaigns
[Tilley, 2020]

Media Bias
[Esses et al., 2013]

Misinformation
[Henderson & McCready, 2019]

Propaganda
[Landry et al., 2022]

# …with far-reaching implications

# ...with far-reaching implications



Electoral
Outcomes
[Haney López, 2014]

# …with far-reaching implications

Policymaking
[Walgrave et al., 2018]

Electoral
Outcomes
[Haney López, 2014]

# ...with far-reaching implications

Policymaking
[Walgrave et al., 2018]

Electoral
Outcomes
[Haney López, 2014]

Public
Opinion
[Jacoby, 2000; Chong
& Druckman, 2007 ]

# ...with far-reaching implications

**Policymaking**
[Walgrave et al., 2018]

**Trust**
[Hopmann et al., 2015]

**Electoral Outcomes**
[Haney López, 2014]

**Public Opinion**
[Jacoby, 2000; Chong & Druckman, 2007 ]

# …with far-reaching implications

**Policymaking**
[Walgrave et al., 2018]

**Trust**
[Hopmann et al., 2015]

**Electoral Outcomes**
[Haney López, 2014]

**Public Opinion**
[Jacoby, 2000; Chong & Druckman, 2007 ]

**Safety & Well-being**
[Rai et al., 2017]

Uncovering implicit language is challenging

Framing
*safety threat*

Dehumanization & Metaphor
*water, vermin*

Dogwhistles
*Soros (Jewish) plot*

Uncovering implicit language is challenging

But computational methods can help!

**Framing**
*safety threat*

**Dehumanization & Metaphor**
*water, vermin*

**Dogwhistles**
*Soros (Jewish) plot*

I develop **computational approaches** to study these strategies and their social, political & technological implications

My thesis

Framing
*NAACL (2021)*
*EMNLP (2022)*
*JQD (2024)*

Dogwhistles
*ACL (2023)*

Dehumanization & Metaphor
*Frontiers in AI (2020)*
*PNAS (2022)*
*ACL (2025)*

I develop **computational approaches** to study strategies of othering and their social, political & technological implications

Framing
*NAACL (2021)*
*EMNLP (2022)*
*JQD (2024)*

Dehumanization & Metaphor
*Frontiers in AI (2020)*
*PNAS (2022)*
*ACL (2025)*

Dogwhistles
*ACL (2023)*

# Many ways to think about **framing**

## Modeling Framing in Immigration Discourse on Social Media

**Julia Mendelsohn**
University of Michigan
juliame@umich.edu

**Ceren Budak**
University of Michigan
cbudak@umich.edu

**David Jurgens**
University of Michigan
jurgens@umich.edu

**Abstract**

The framing of political issues can influence policy and public opinion. Even though the public plays a key role in creating and spreading frames, little is known about how ordinary people on social media frame political issues. By creating a new dataset of immigration-related tweets labeled for multiple framing typologies from political communication theory, we develop supervised models to detect frames. We demonstrate how users' ideology and region impact framing choices, and how a message's framing influences audience responses. We find that the more commonly-used issue-generic frames obscure important ideological and regional patterns that are only revealed by immigration-specific frames. Furthermore, frames oriented towards human interests, culture, and politics are associated with higher user engagement. This large-scale analysis of a complex social and linguistic phenomenon contributes to both NLP and social science research.

social media content enables us to compare framing strategies across countries and political ideologies. Furthermore, social media provides unique insights into how messages resonate with audiences through interactive signals such as retweets and favorites. By jointly analyzing the production and reception of frames on Twitter, we provide an in-depth analysis of immigration framing by and on the public.

Political communications research has identified numerous typologies of frames, such as *issue-generic policy*, *immigration-specific*, and *narrative*. Each of these frame types can significantly shape the audience's perceptions of an issue (Iyengar, 1991; Chong and Druckman, 2007; Lecheler et al., 2015), but prior NLP work seeking to detect frames in mass media (e.g. Card et al., 2016; Field et al., 2018; Kwak et al., 2020) has largely been limited to a single *issue-generic policy* typology. Multiple dimensions of framing must be considered in order to better understand the structure of immigration discourse and its effect on public opinion

nal of Quantitative Description: Digital Media 4(2024), 1–61          10.51685/jqd.2024.icw

## Framing Social Movements on Social Media: Unpacking Diagnostic, Prognostic, and Motivational Strategies

JULIA MENDELSOHN

MAYA VIJAN

DALLAS CARD

CEREN BUDAK

University of Michigan, USA

Social media enables activists to directly communicate with the public and provides a space for movement leaders, participants, bystanders, and opponents to collectively construct and contest narratives. Focusing on Twitter messages from social movements surrounding three issues in 2018-2019 (guns, immigration, and LGBTQ rights), we create a codebook, annotated dataset, and computational models to detect diagnostic (problem identification and attribution), prognostic (proposed solutions and tactics), and motivational (calls to action) framing strategies. We conduct an in-depth unsupervised linguistic analysis of each framing strategy, and uncover cross-movement similarities in

## When People are *Floods*: Analyzing Dehumanizing Metaphors in Immigration Discourse with Large Language Models

*Warning: this paper contains examples of upsetting and offensive content.*

**Julia Mendelsohn**
University of Maryland
juliame@umd.edu

**Ceren Budak**
University of Michigan
cbudak@umich.edu

**Abstract**

Metaphor, discussing one concept in terms of another, is abundant in politics and can shape how people understand important issues. We develop a computational approach to measure metaphorical language, focusing on immigration discourse on social media. Grounded in qualitative social science research, we identify seven source domain concepts evoked in immigration discourse (e.g. WATER or VERMIN). We propose and evaluate a novel technique that leverages both word-level and document-level signals to measure metaphor with respect to these source domains. We then study the relationship between metaphor, political ideology, and user engagement in 400K US tweets about immigration. While conservatives tend to use dehumanizing metaphors more than liberals, this effect varies widely across source domains. Moreover, creature-related metaphor is associated with more retweets, especially for liberal authors. Our work highlights the potential for computational methods to complement quali-

"They want immigrants to pour into and infest this country"

Target Domain

Water          Vermin
*Source Domain Concepts*

Figure 1: Dehumanizing sentence likening immigrants to the *source domain concepts* of WATER and VERMIN via the words "pour" and "infest".

for measuring and analyzing metaphor at scale. We use this methodology to study dehumanizing metaphor in immigration discourse on social media, and analyze the relationship between metaphor use, political ideology, and user engagement.

From prior discourse analysis literature, we first identify seven *source domains*: concepts evoked in discussions of immigration such as WATER or VERMIN (Figure 1). We use large language models

# What is framing?

"Selecting some aspects of a perceived reality and make them **more salient** in a communicating text, in such a way as to promote a particular **problem definition, causal interpretation, moral evaluation, and/or treatment recommendation** for the item described" [Entman, 1993]

# What is a frame?

# What is a frame?

- **Issue-generic Policy** [Boydstun et al., 2013]
  - *Crime & punishment, morality, economic, policy*

# What is a frame?

- **Issue-generic Policy** [Boydstun et al., 2013]
  - *Crime & punishment, morality, economic, policy*

- **Immigration-specific** [Benson, 2013]
  - *Immigrants as victims (e.g. of global economy or discrimination)*
  - *Immigrants as heroes (e.g. contributing to economy or cultural diversity)*
  - *Immigrants as threats (e.g. to jobs, or to public safety)*

# What is a frame?

- **Issue-generic Policy** [Boydstun et al., 2013]
  - *Crime & punishment, morality, economic, policy*

- **Immigration-specific** [Benson, 2013]
  - *Immigrants as victims (e.g. of global economy or discrimination)*
  - *Immigrants as heroes (e.g. contributing to economy or cultural diversity)*
  - *Immigrants as threats (e.g. to jobs, or to public safety)*

- **Issue-generic Narrative** [Iyengar, 1991]
  - *Episodic: focus on specific actions, events, examples, or case studies*
  - *Thematic: focus on broader political, social, cultural context*

# Framing processes

- **Frame-building**: factors affecting how an issue is framed

**Inputs**
**Ideologies**
Background
Attitudes
Elite rhetoric

**Frame-building** →

**Frames**
Issue-specific
Issue-generic
policy
Narrative

Figure & theoretical model adapted from de Vreese [2005] and is a simplification of Scheufele's [1999] four-process model

# Framing processes

- **Frame-building**: factors affecting how an issue is framed
- **Frame-setting**: frame effects on audiences

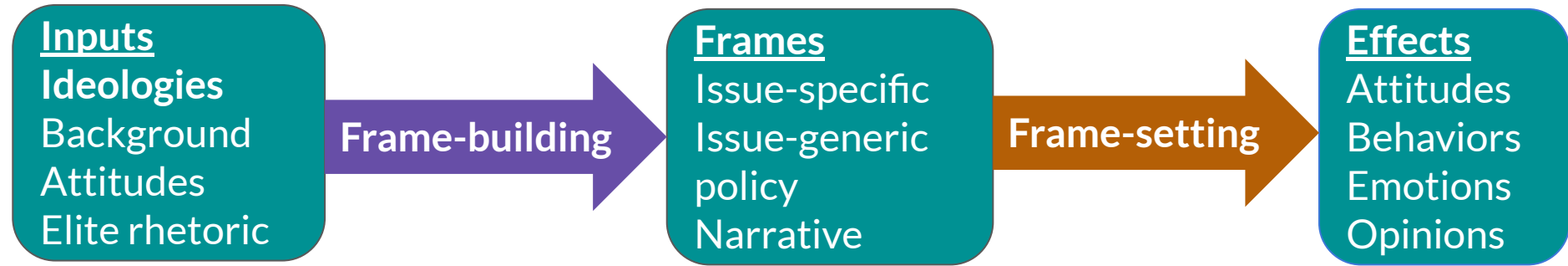| **Inputs** | Frame-building | **Frames** | Frame-setting | **Effects** |
|---|---|---|---|---|
| **Ideologies** | | Issue-specific | | Attitudes |
| Background | | Issue-generic | | Behaviors |
| Attitudes | | policy | | Emotions |
| Elite rhetoric | | Narrative | | Opinions |

Figure & theoretical model adapted from de Vreese [2005] and is a simplification of Scheufele's [1999] four-process model
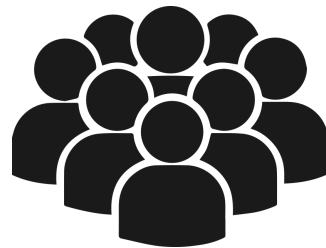
Dataset collection & annotation

Automated frame detection

Frame building: role of ideology in framing

Frame setting: effects on user engagement

# Data Annotation

## 3 typologies

## 27 categories

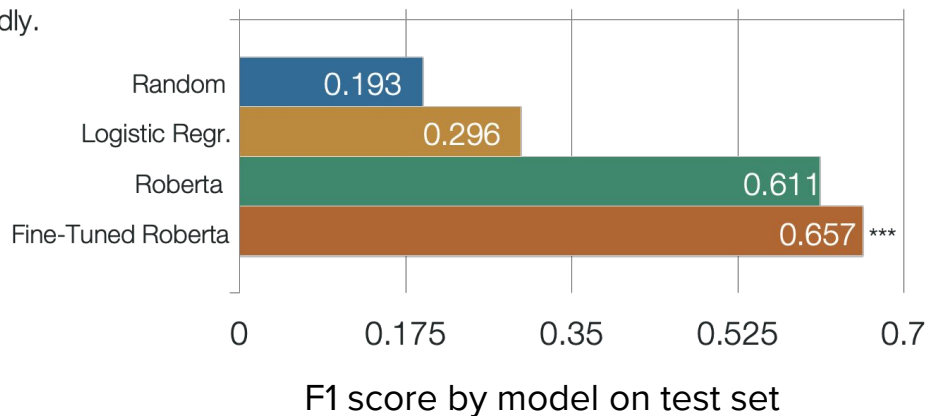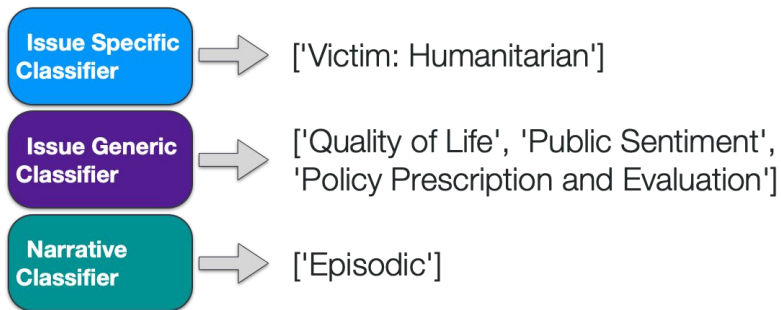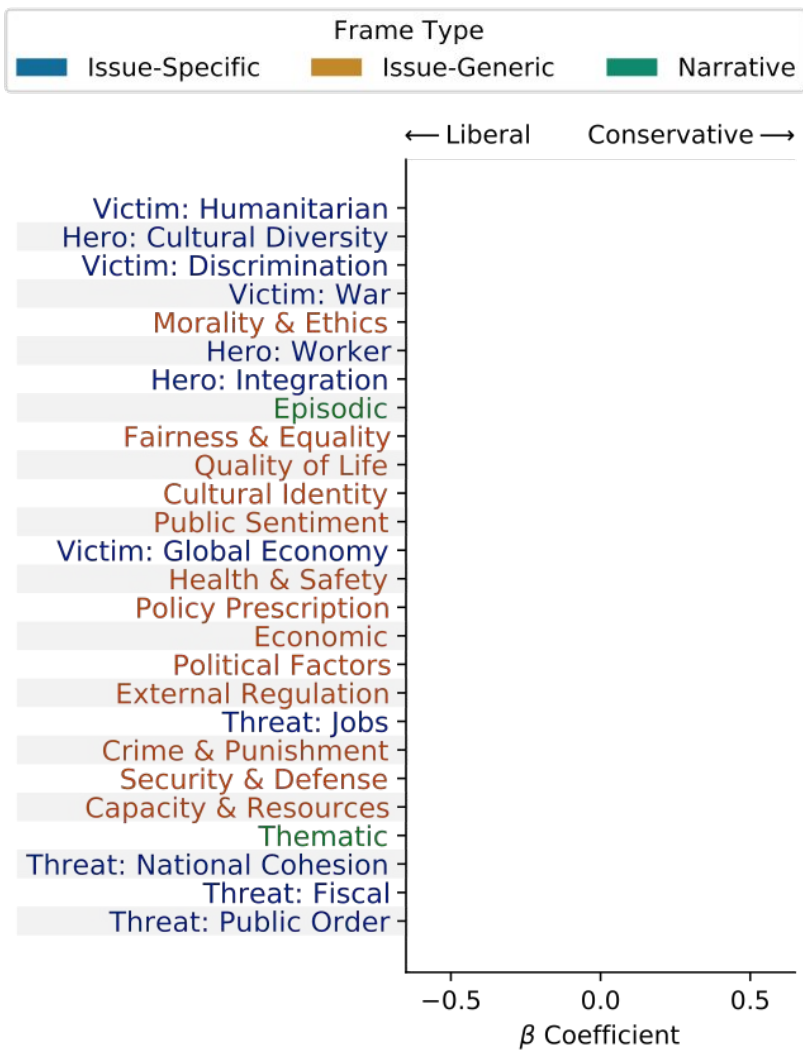| Frame Type | Frame | Description |
|---|---|---|
| Issue-Generic Policy | Economic | Financial implications of an issue |
| | Capacity & Resources | The availability or lack of time, physical, human, or financial resources |
| | Morality & Ethics | Perspectives compelled by religion or secular sense of ethics or social responsibility |
| | Fairness & Equality | The (in)equality with which laws, punishments, rewards, resources are distributed |
| | Legality, Constitutionality & Jurisdiction | Court cases and existing laws that regulate policies; constitutional interpretation; legal processes such as seeking asylum or obtaining citizenship; jurisdiction |
| | Crime & Punishment | The violation of policies in practice and the consequences of those violations |
| | Security & Defense | Any threat to a person, group, or nation and defenses taken to avoid that threat |
| | Health & Safety | Health and safety outcomes of a policy issue, discussions of health care |
| | Quality of Life | Effects on people's wealth, mobility, daily routines, community life, happiness, etc. |
| | Cultural Identity | Social norms, trends, values, and customs; integration/assimilation efforts |
| | Public Sentiment | General social attitudes, protests, polling, interest groups, public passage of laws |
| | Political Factors & Implications | Focus on politicians, political parties, governing bodies, political campaigns and debates; discussions of elections and voting |
| | Policy Prescription & Evaluation | Discussions of existing or proposed policies and their effectiveness |
| | External Regulation & Reputation | Relations between nations or states/provinces; agreements between governments; perceptions of one nation/state by another |
| Immigration Specific | Victim: Global Economy | Immigrants are victims of global poverty, underdevelopment and inequality |
| | Victim: Humanitarian | Immigrants experience economic, social, and political suffering and hardships |
| | Victim: War | Focus on war and violent conflict as reason for immigration |
| | Victim: Discrimination | Immigrants are victims of racism, xenophobia, and religion-based discrimination |
| | Hero: Cultural Diversity | Highlights positive aspects of differences that immigrants bring to society |
| | Hero: Integration | Immigrants successfully adapt and fit into their host society |
| | Hero: Worker | Immigrants contribute to economic prosperity and are an important source of labor |
| | Threat: Jobs | Immigrants take nonimmigrants' jobs or lower their wages |
| | Threat: Public Order | Immigrants threaten public safety by being breaking the law or spreading disease |
| | Threat: Fiscal | Immigrants abuse social service programs and are a burden on resources |
| | Threat: National Cohesion | Immigrants' cultural differences are a threat to national unity and social harmony |
| Narrative | Episodic | Message provides concrete information about on specific people, places, or events |
| | Thematic | Message is more abstract, placing stories in broader political and social contexts |

# Data Collection and Annotation

- 2.6M English tweets from 2018-2019 containing immigration-related term
- Ideology inference using existing network-based tool
- 4.5K tweets labeled by trained annotators for all frames explicitly cued.
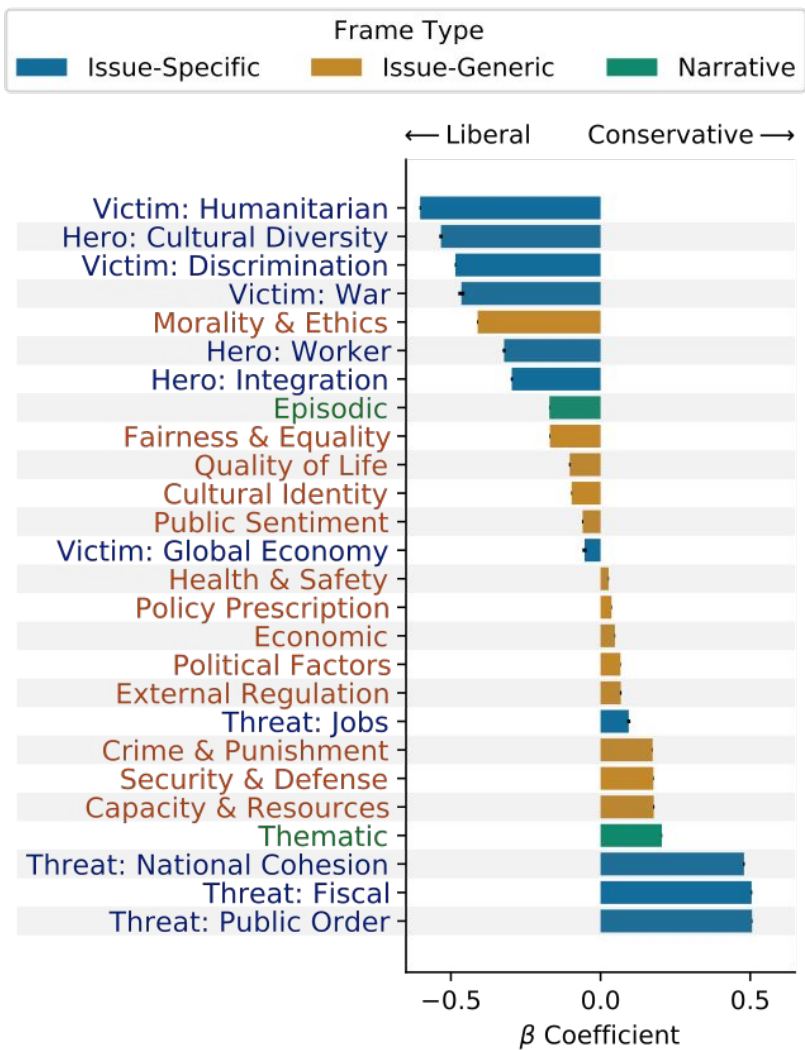
# Frame Detection

- RoBERTa-based multilabel classification models fine-tuned on full set of immigration-related tweets
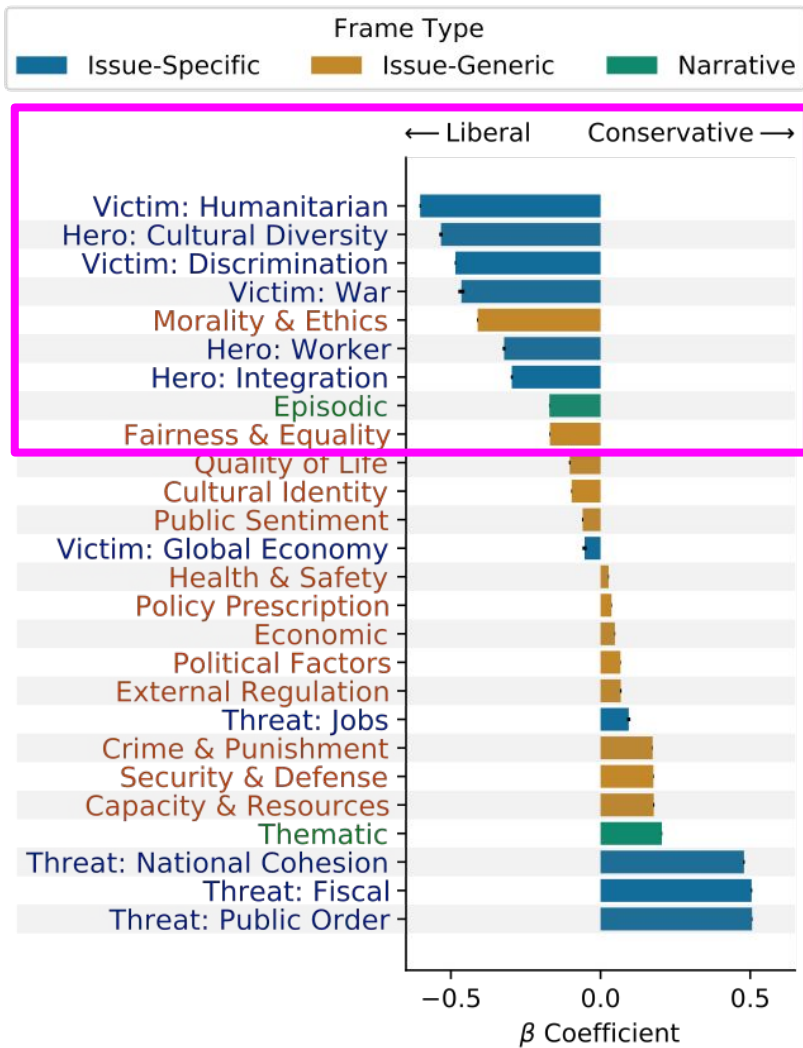
"The proposed #TitleIX rule will exacerbate the negative experiences of undocumented students on campus. Get ready to speak your voice loudly. Go to https://t.co/7kAlhYjeLl to learn how with @endrapeoncampus"

**Issue Specific Classifier** → ['Victim: Humanitarian']

**Issue Generic Classifier** → ['Quality of Life', 'Public Sentiment', 'Policy Prescription and Evaluation']

**Narrative Classifier** → ['Episodic']

| Model | F1 score |
|---|---|
| Random | 0.193 |
| Logistic Regr. | 0.296 |
| Roberta | 0.611 |
| Fine-Tuned Roberta | 0.657 *** |

0    0.175    0.35    0.525    0.7

F1 score by model on test set

Frame Type: Issue-Specific, Issue-Generic, Narrative

← Liberal    Conservative →

- Victim: Humanitarian
- Hero: Cultural Diversity
- Victim: Discrimination
- Victim: War
- Morality & Ethics
- Hero: Worker
- Hero: Integration
- Episodic
- Fairness & Equality
- Quality of Life
- Cultural Identity
- Public Sentiment
- Victim: Global Economy
- Health & Safety
- Policy Prescription
- Economic
- Political Factors
- External Regulation
- Threat: Jobs
- Crime & Punishment
- Security & Defense
- Capacity & Resources
- Thematic
- Threat: National Cohesion
- Threat: Fiscal
- Threat: Public Order

$\beta$ Coefficient

# Liberals frame immigrants as **heroes** and **victims**

- Liberals cue *fairness* and *morality*, framing immigrants as *victims of discrimination* and *inhumane* policies.
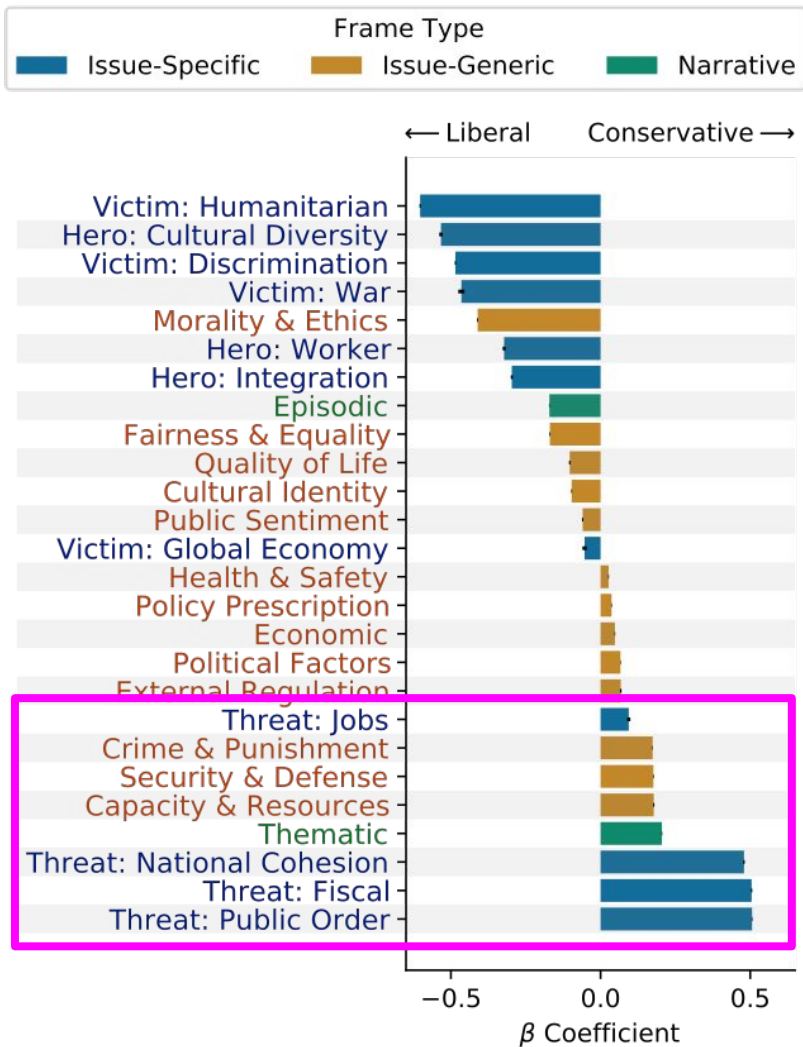
**Liberals frame immigrants as heroes and victims**

- Liberals cue *fairness* and *morality*, framing immigrants as *victims of discrimination* and *inhumane* policies.
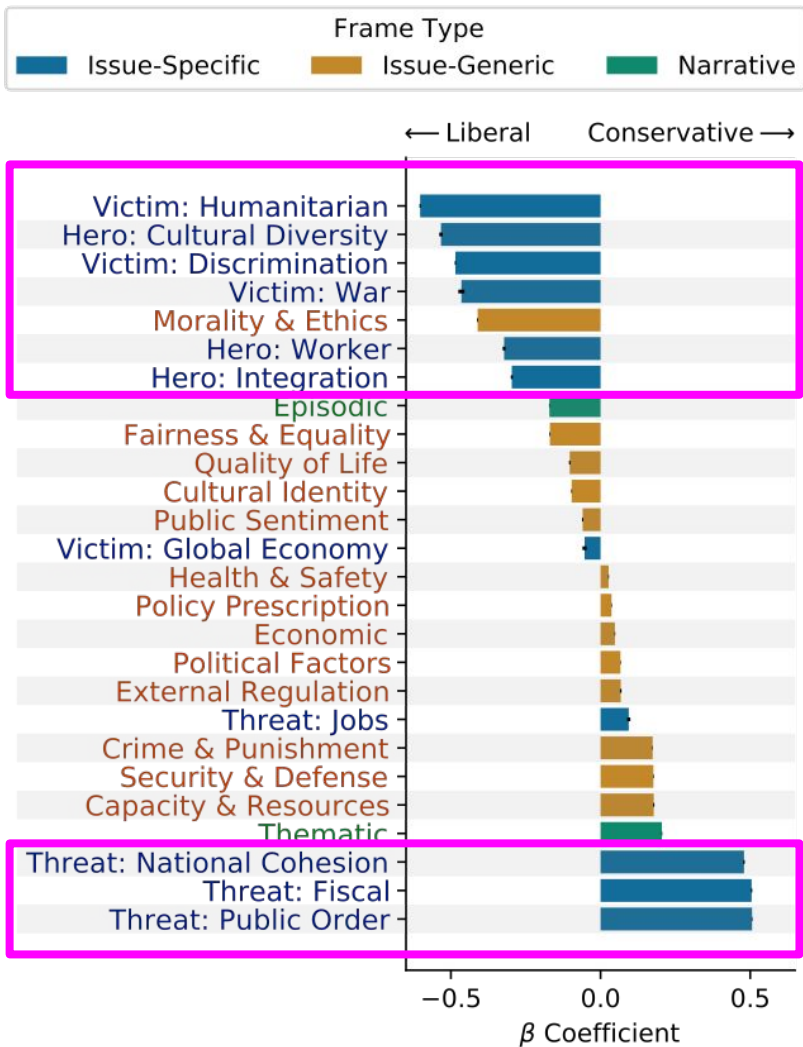
**Conservatives frame immigrants as threats**

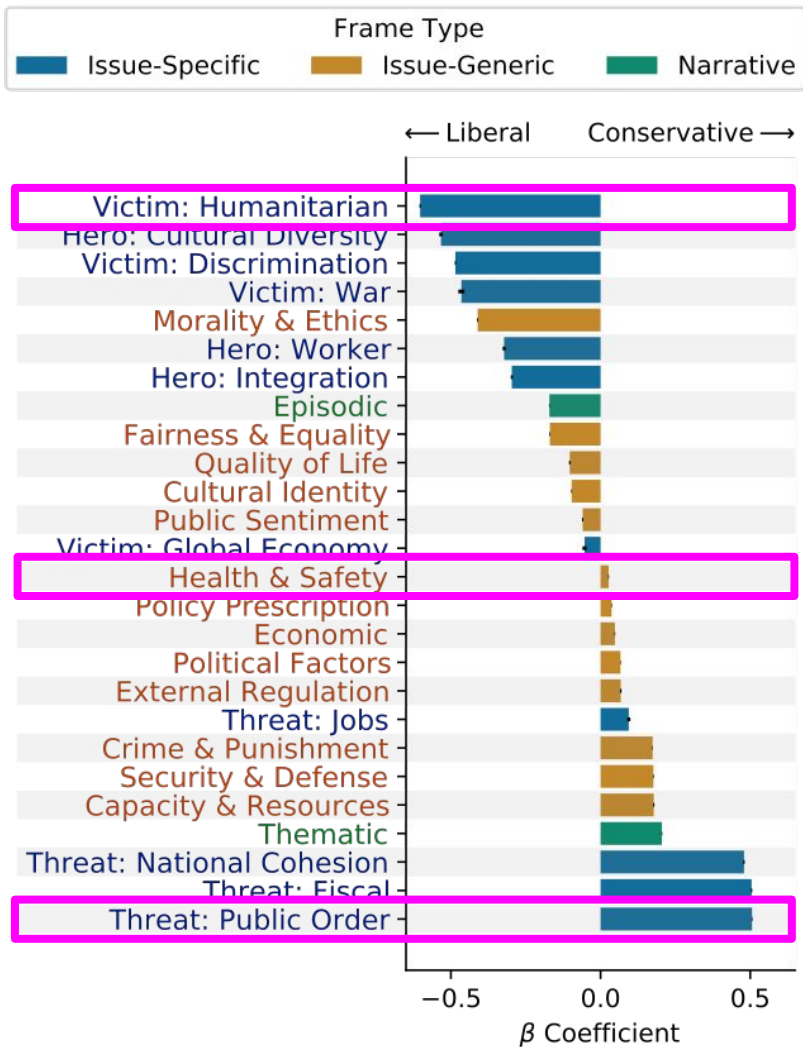- Conservatives cue *threat to public safety*, *burden on taxpayers & government programs*

# Each frame typology offers value

Issue-generic policy frames can be most accurately detected

**but**

Immigration-specific frames reveal ideological differences obscured by issue-generic policy frames
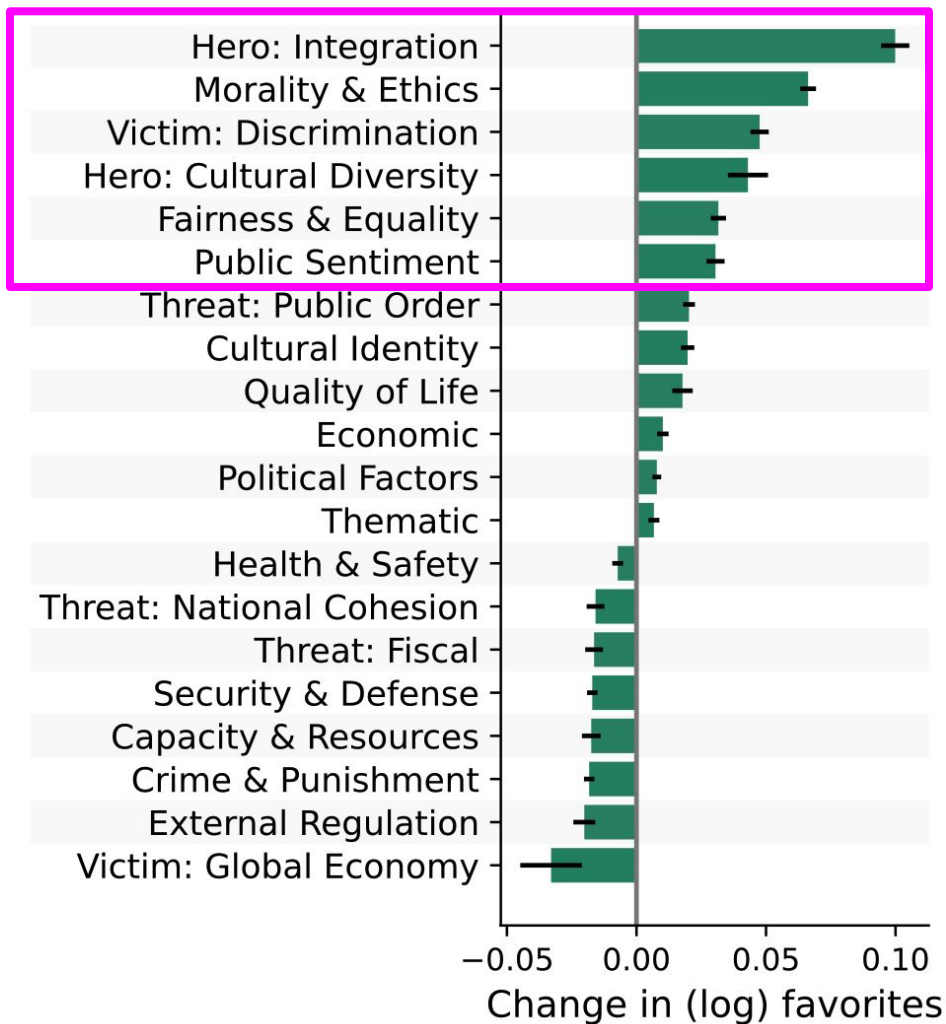
# Each frame typology offers value

Issue-generic policy frames can be most accurately detected

**but**

Immigration-specific frames reveal ideological differences obscured by issue-generic policy frames

(e.g. *health & safety*)

**Cultural** (*hero: integration*) and **human interest** (*morality, fairness, victim: discrimination*)

# Framing Social Movements on Social Media: *Unpacking Diagnostic, Prognostic, and Motivational Strategies*

Journal of Quantitative Description: Digital Media (2024)

Julia Mendelsohn     Maya Vijan

Dallas Card     Ceren Budak

Facebook
ناس بوك
#jan25
TheEgyptian
SocialNetwork

SKOLSTREJK FOR KLIMATET

September 17th. Wall Street. Bring Tent.
http://bit.ly/re9ENL
#OCCUPYWALLSTREET
← Reply ↻ Retweet ★ Favorite

#SayHerName
#SandraBland

Alyssa Milano
@Alyssa_Milano
Follow

If you've been sexually harassed or assaulted write 'me too' as a reply to this tweet.

Me too.

Suggested by a friend: "If all the women who have been sexually harassed or assaulted wrote 'Me too.' as a status, we might give people a sense of the magnitude of the problem."
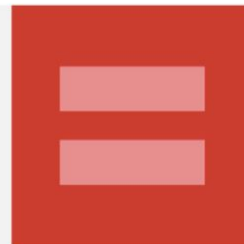
1:21 PM - 15 Oct 2017

24,725 Retweets 53,346 Likes

FREE OUR FUTURE!
#FreeOurFuture
#AbolishICE

#NEVERAGAIN
TEXT 'MARCH' TO 90975
@PHXMar
MARCH FOR OUR LIVES
marchforourlivesaz.org

Human Rights Campaign
26 March 2013

Who's wearing red tomorrow? Show your support for marriage equality -- make your profile image red for tomorrow and check out www.hrc.org/StandForMarriage for more ways to get involved!

52

Like · Comment · Share          19,543  678  70,652 Shares

# Social movements create meaning through framing

- Collective action frames are "intended to mobilize potential adherents and constituents, to garner bystander support, and to demobilize antagonists" [Snow & Benford, 1988]

- Effective framing is important (perhaps, necessary) for social movement success [Della Porta & Diani, 2006]

# Core Framing Tasks [Snow & Benford, 1988; Benford & Snow, 2000]

> ## Diagnostic
> Identifying social problems, their causes, and who to hold responsible

# Core Framing Tasks [Snow & Benford, 1988; Benford & Snow, 2000]

## Diagnostic
Identifying social problems, their causes, and who to hold responsible

## Prognostic
Proposing solutions, plans of attack, strategies for carrying out that plan

# Core Framing Tasks [Snow & Benford, 1988; Benford & Snow, 2000]

**Diagnostic**
Identifying social problems, their causes, and who to hold responsible

**Prognostic**
Proposing solutions, plans of attack, strategies for carrying out that plan

**Motivational**
Persuading people to participate through "calls to action"

How do people use **diagnostic**, **prognostic**, and **motivational** framing in Twitter messages related to social movements?
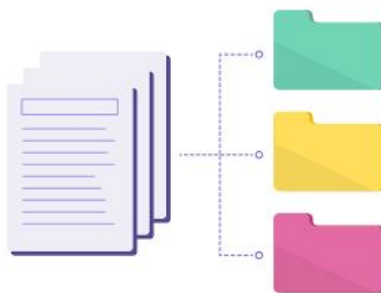
Dataset &
Annotation
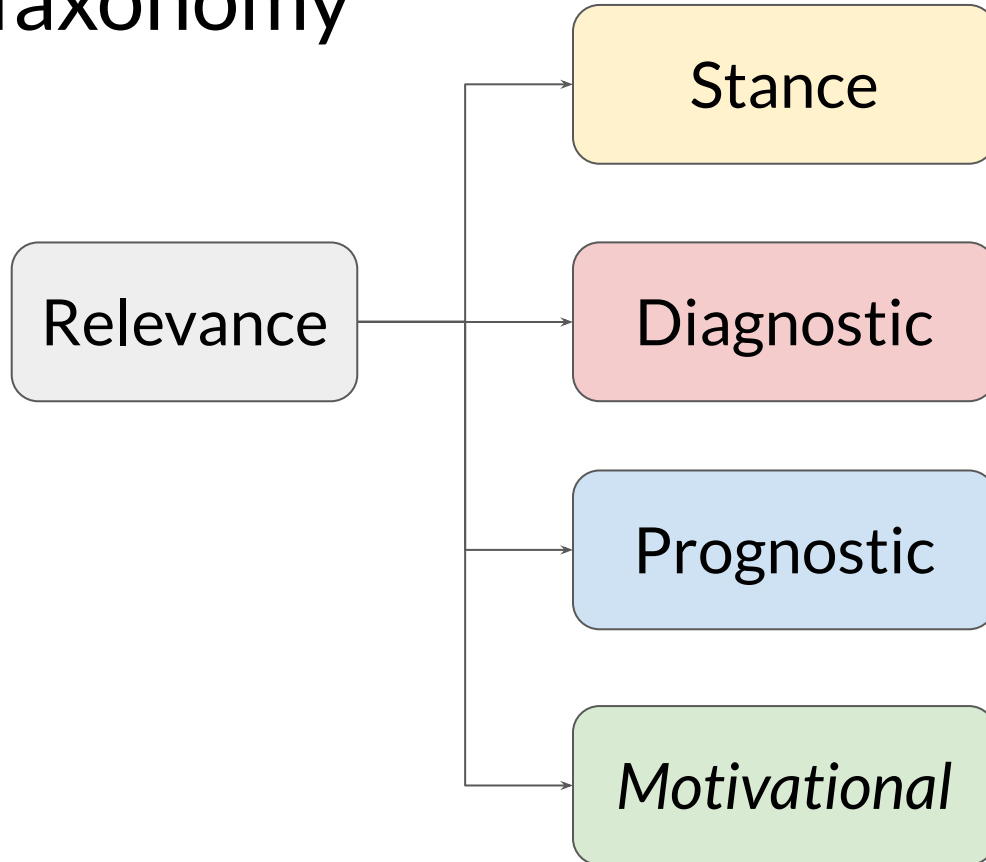
Dataset & Annotation
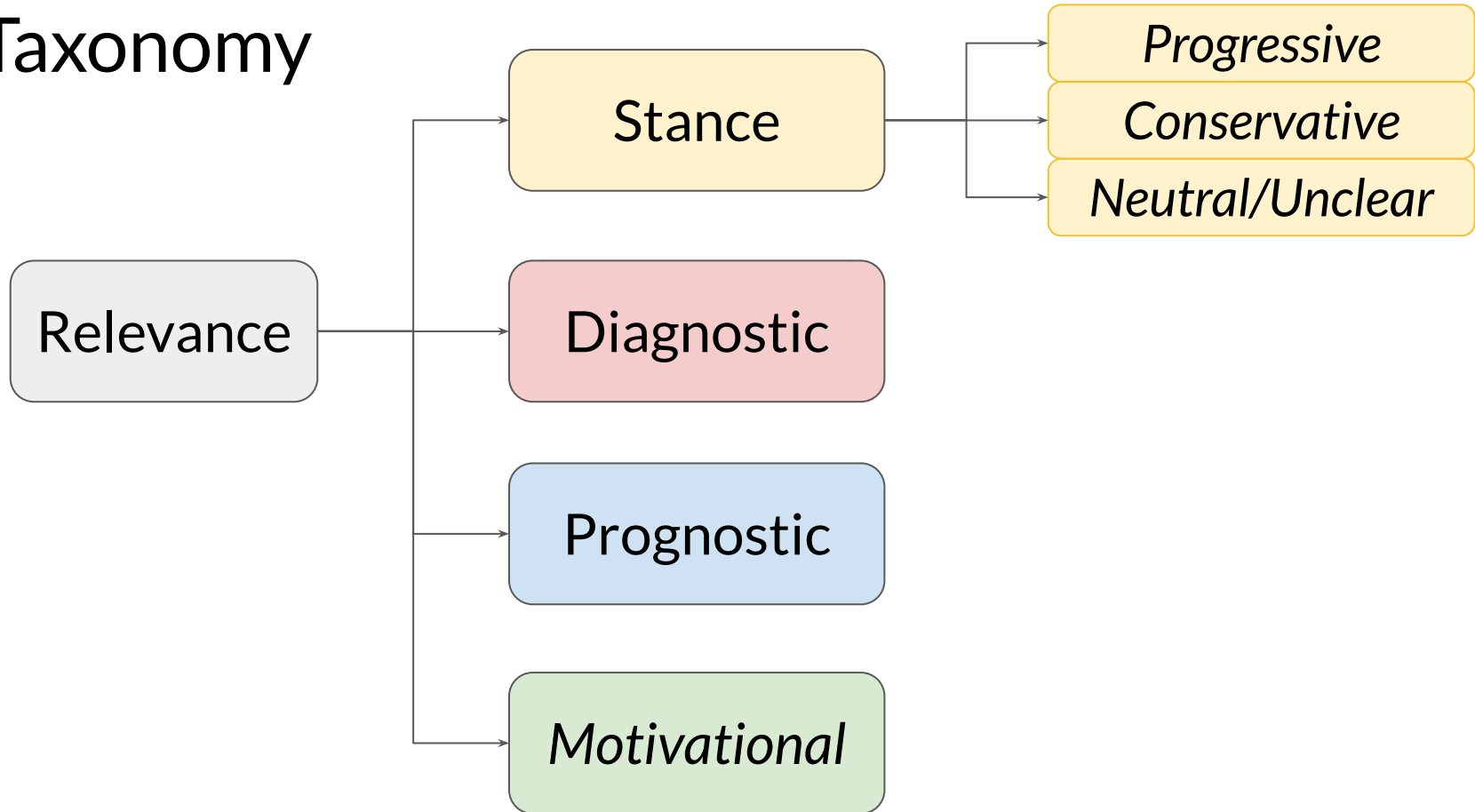
Classification

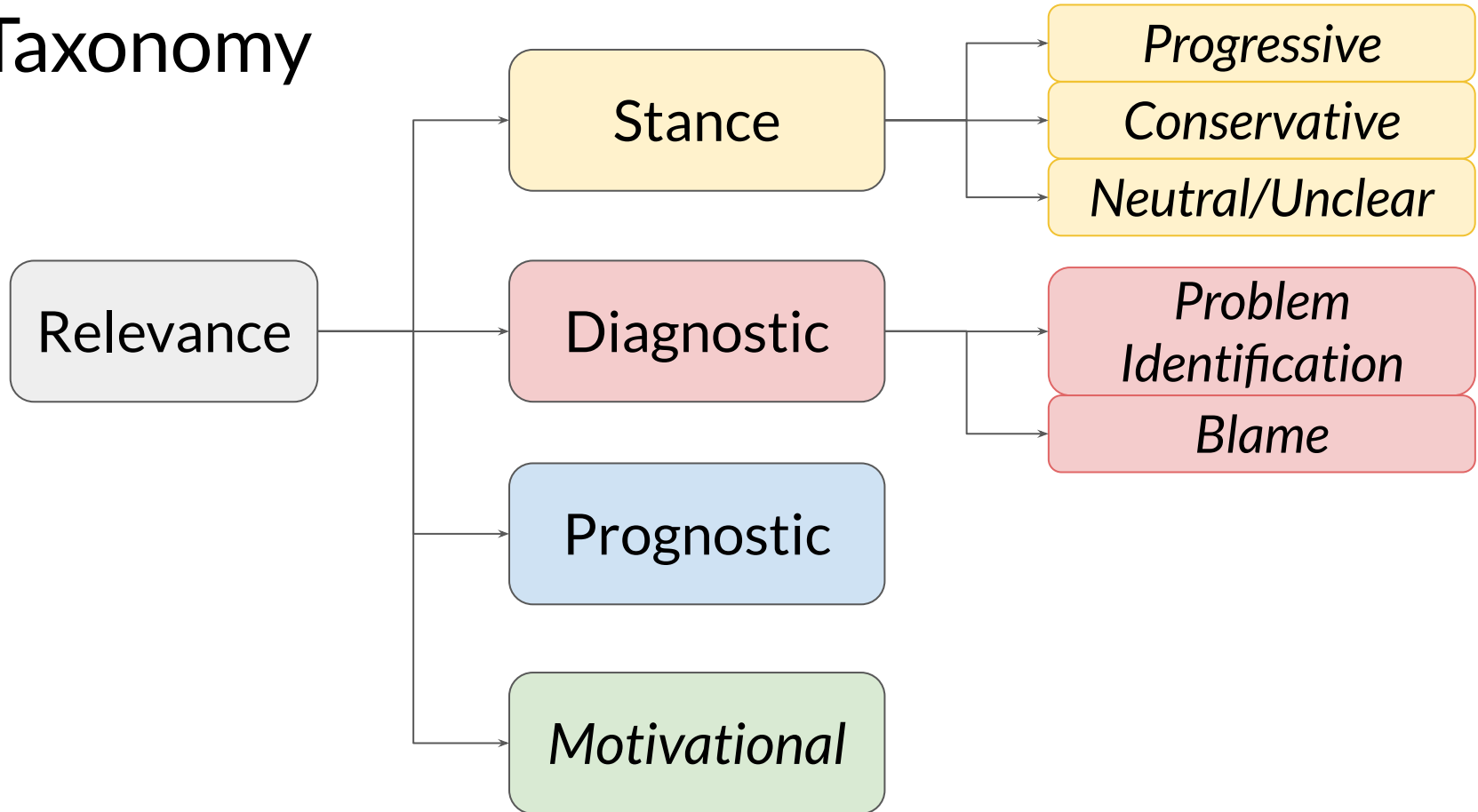Dataset &
Annotation

Classification
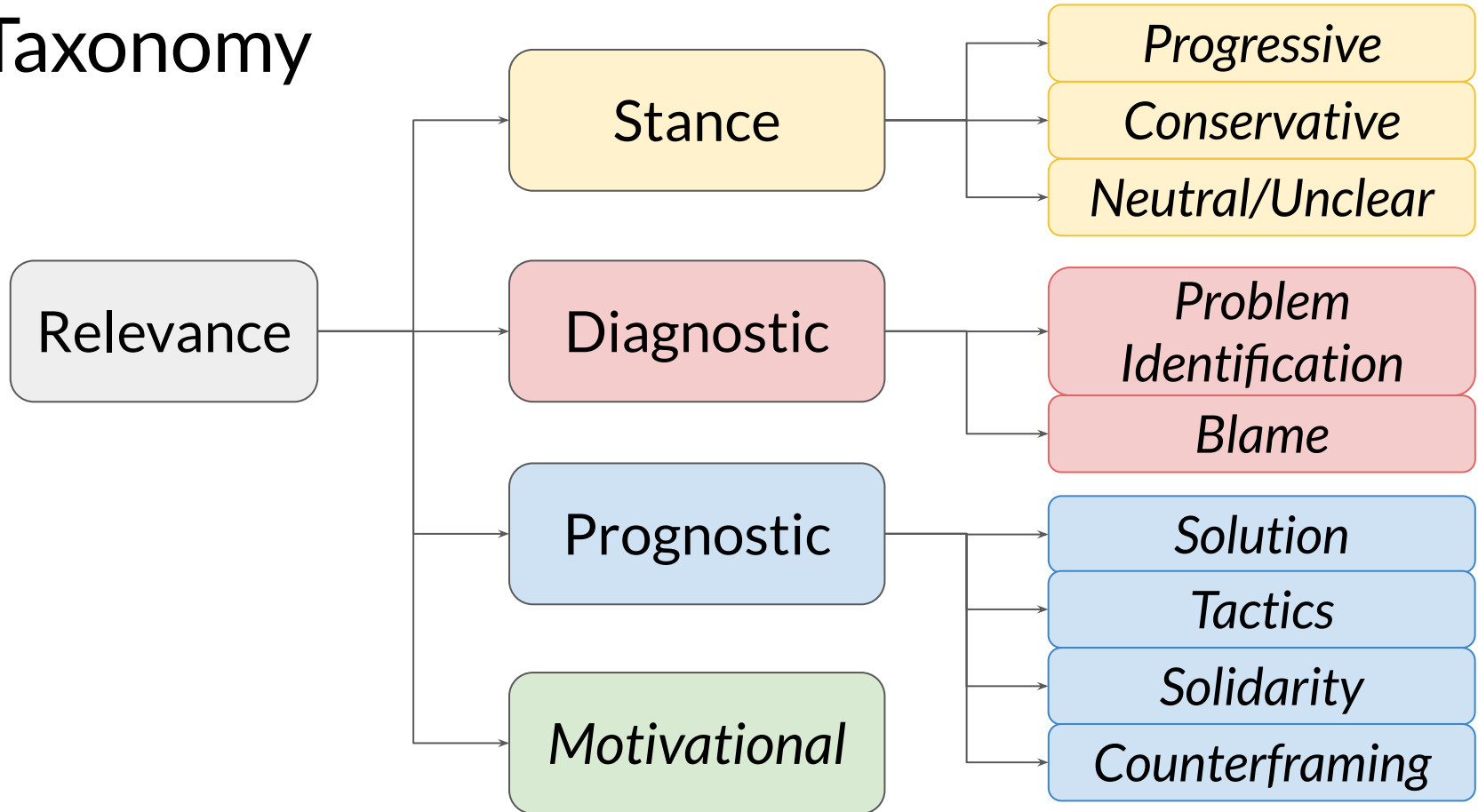
Sociocultural
Context

# Taxonomy



Relevance → Stance

Relevance → Diagnostic

Relevance → Prognostic

Relevance → *Motivational*

# Taxonomy



Relevance
- Stance
  - *Progressive*
  - *Conservative*
  - *Neutral/Unclear*
- Diagnostic
- Prognostic
- *Motivational*

# Taxonomy

Relevance

- Stance
  - Progressive
  - Conservative
  - Neutral/Unclear
- Diagnostic
  - Problem Identification
  - Blame
- Prognostic
- Motivational

# Taxonomy

# Data Collection [from Bozarth & Budak, 2022]

- Tweets from movements focused on 3 issues: **_guns_**, **_immigration_**, and **_LGBTQ rights_** from 2018-2019
  - Both progressive & conservative movements
  - 2 months for each issue (1 high protest activity, 1 average level)

# Data Collection [from Bozarth & Budak, 2022]

- Tweets from movements focused on 3 issues: *guns*, *immigration*, and *LGBTQ rights* from 2018-2019
  - Both progressive & conservative movements
  - 2 months for each issue (1 high protest activity, 1 average level)
- **1.85M tweets** across all movements
  - 822K for guns, 763K for immigration, 268K for LGBTQ

- **6,000** manually-annotated tweets

- **4,859 (81%)** coded as *relevant*, labeled for stance & frames
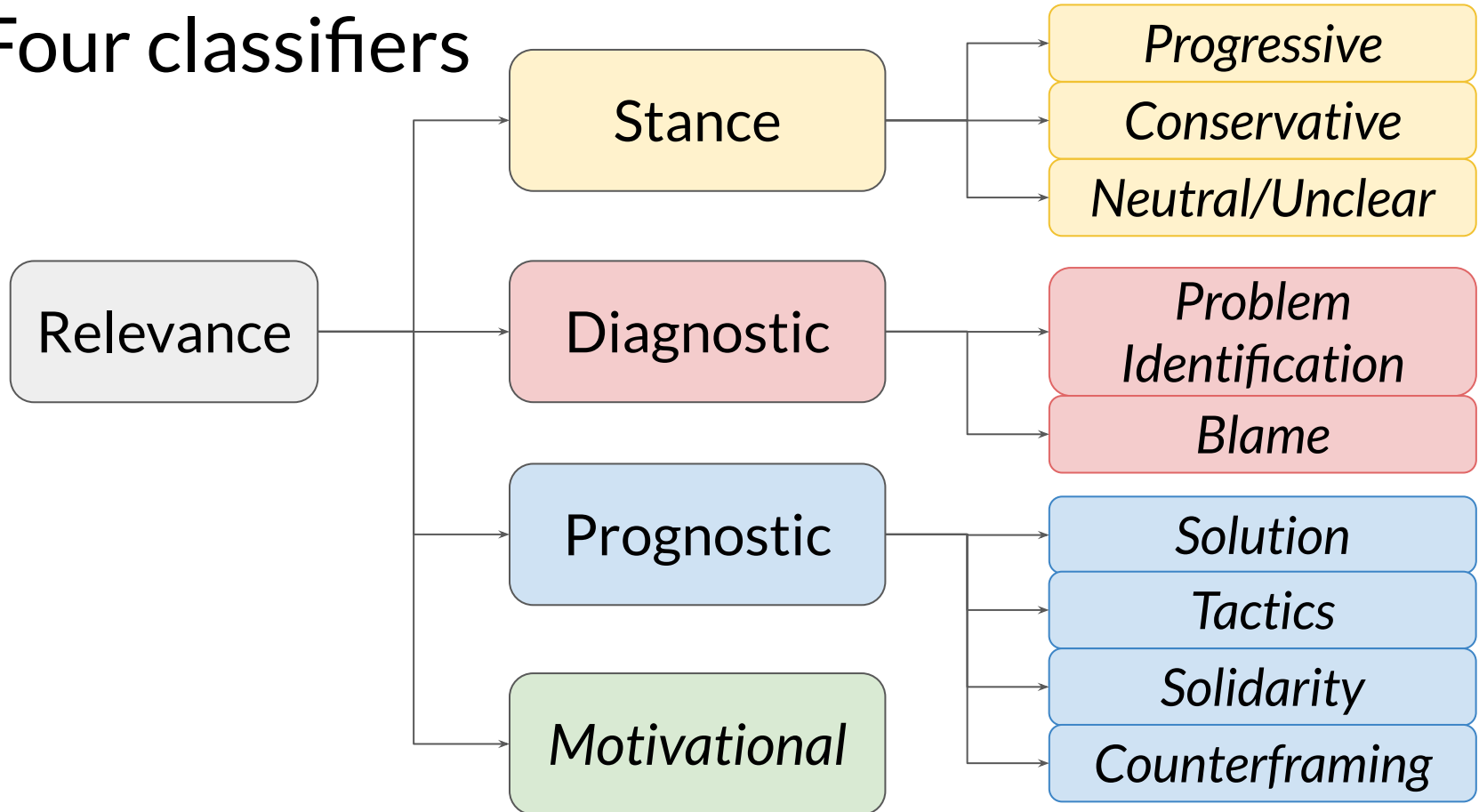
Dataset &
Annotation

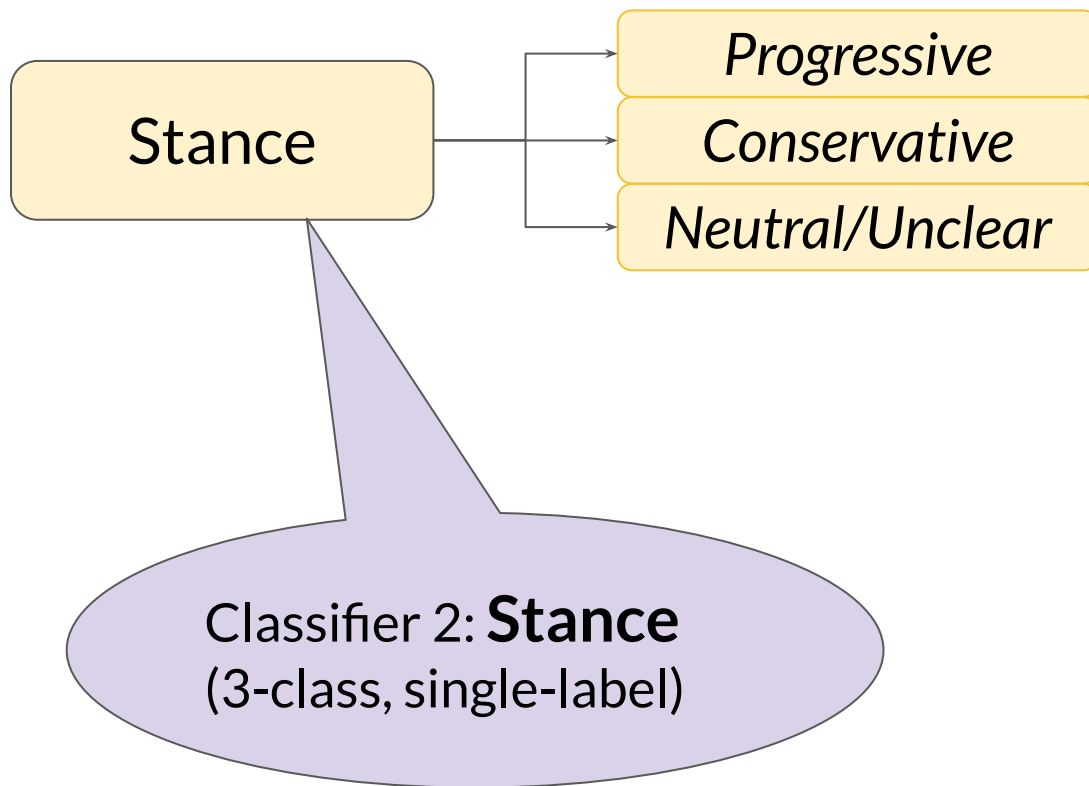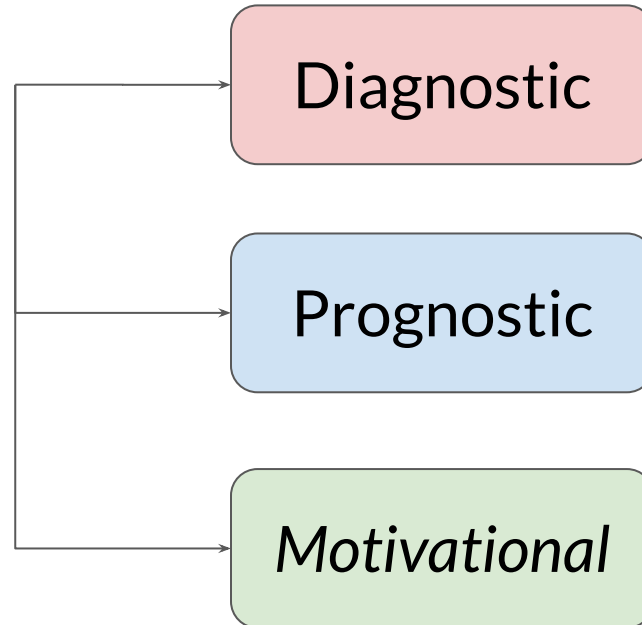Classification

Sociocultural
Context

# Four classifiers



Relevance

Stance
- *Progressive*
- *Conservative*
- *Neutral/Unclear*

Diagnostic
- *Problem Identification*
- *Blame*

Prognostic
- *Solution*
- *Tactics*
- *Solidarity*
- *Counterframing*

*Motivational*

# Four classifiers

Relevance

Classifier 1:
**Relevance**
(binary, single-label)

# Four classifiers

Stance → *Progressive*, *Conservative*, *Neutral/Unclear*

Classifier 2: **Stance**
(3-class, single-label)

# Four classifiers

Classifier 3:
**Core Framing Task**
(binary, 3-label)

Diagnostic

Prognostic

*Motivational*

# Four classifiers

Classifier 4:
**Frame Elements**
*Categories we coded for*
(binary, 7-label)
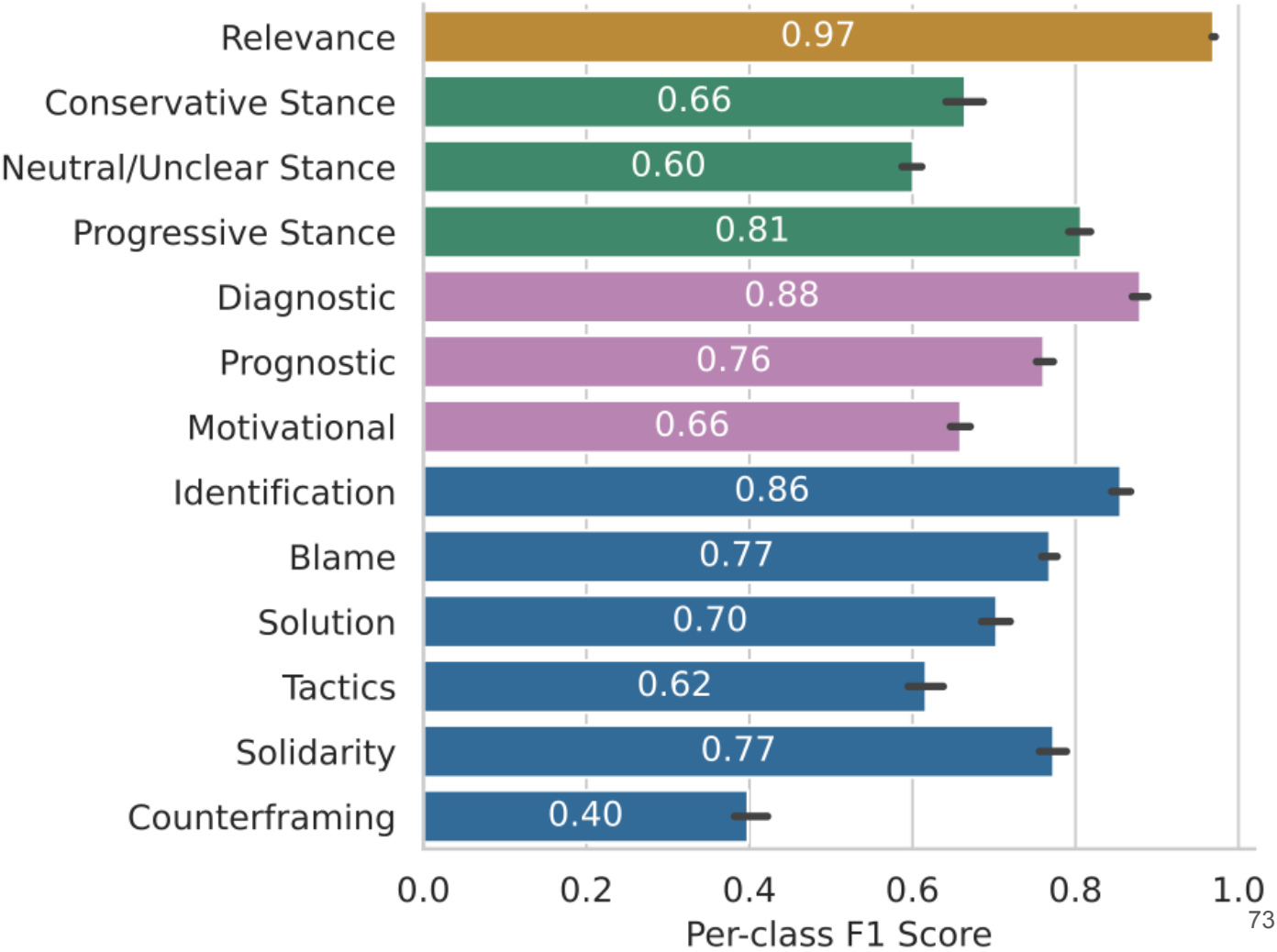
*Motivational*
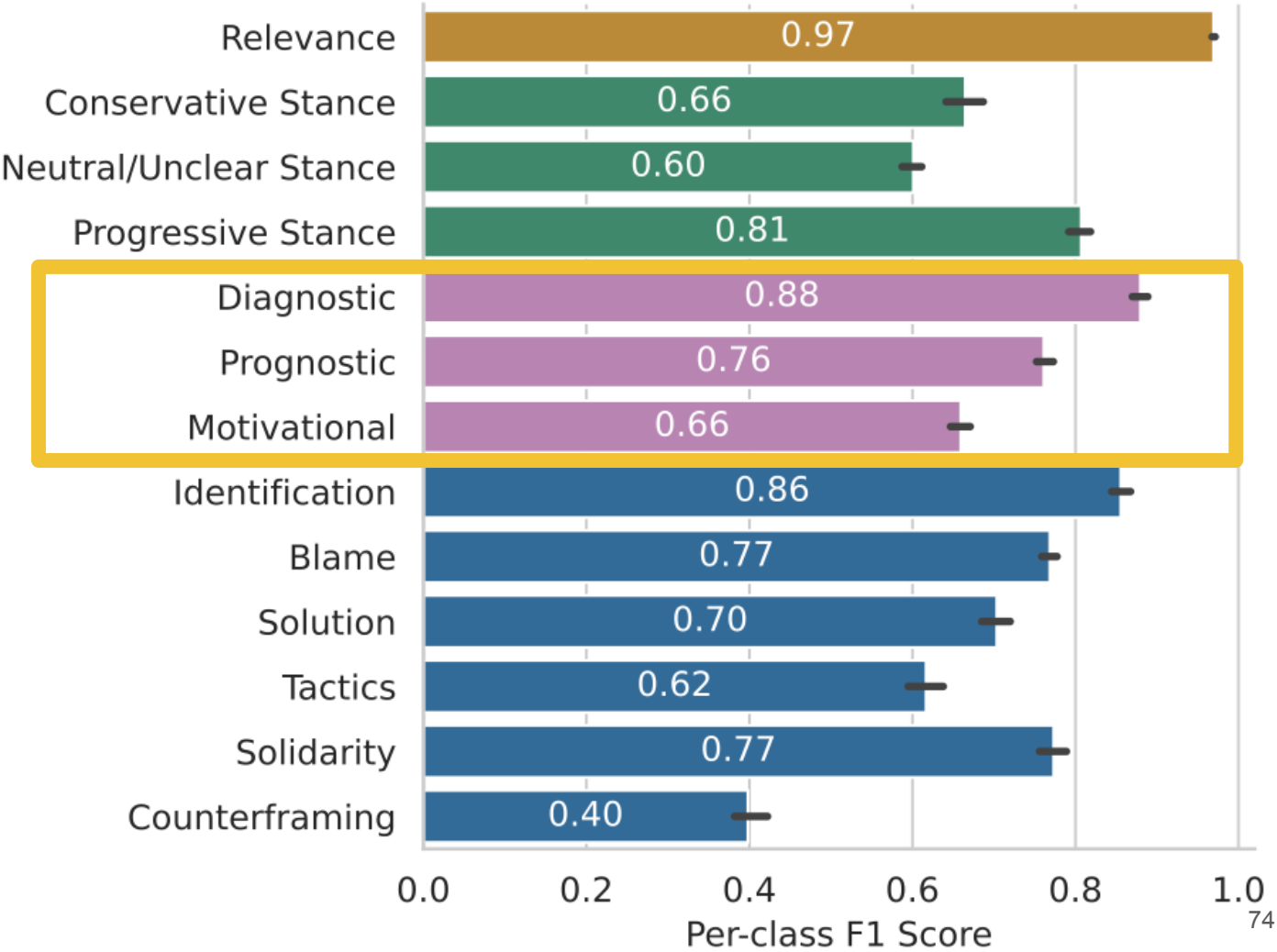
*Problem Identification*

*Blame*

*Solution*

*Tactics*

*Solidarity*

*Counterframing*

# Model Results



Per-class F1 Score

| Category | F1 Score |
|---|---|
| Relevance | 0.97 |
| Conservative Stance | 0.66 |
| Neutral/Unclear Stance | 0.60 |
| Progressive Stance | 0.81 |
| Diagnostic | 0.88 |
| Prognostic | 0.76 |
| Motivational | 0.66 |
| Identification | 0.86 |
| Blame | 0.77 |
| Solution | 0.70 |
| Tactics | 0.62 |
| Solidarity | 0.77 |
| Counterframing | 0.40 |

Model Results

| Category | Per-class F1 Score |
|---|---|
| Relevance | 0.97 |
| Conservative Stance | 0.66 |
| Neutral/Unclear Stance | 0.60 |
| Progressive Stance | 0.81 |
| Diagnostic | 0.88 |
| Prognostic | 0.76 |
| Motivational | 0.66 |
| Identification | 0.86 |
| Blame | 0.77 |
| Solution | 0.70 |
| Tactics | 0.62 |
| Solidarity | 0.77 |
| Counterframing | 0.40 |

Dataset &
Annotation

Classification

Sociocultural
Context

# 3 movement-level factors

## 2 message-level factors

### Issue Area
*guns, immigration, LGBTQ*

### Stance
*progressive, conservative, neutral*
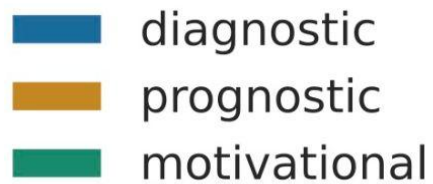
### Protest Activity
*High or average activity month*

### Author Role
*journalist, social mvmt org, other*

### Interaction Type
*broadcast, reply, quote tweet*

Logistic regression + marginal effect for
diagnostic, prognostic, and motivational frames

# 3 movement-level factors

# 2 message-level factors

### Issue Area
*guns, immigration, LGBTQ*

### Stance
*progressive, conservative, neutral*

### Protest Activity
*High or average activity month*

### Author Role
*journalist, social mvmt org, other*

### Interaction Type
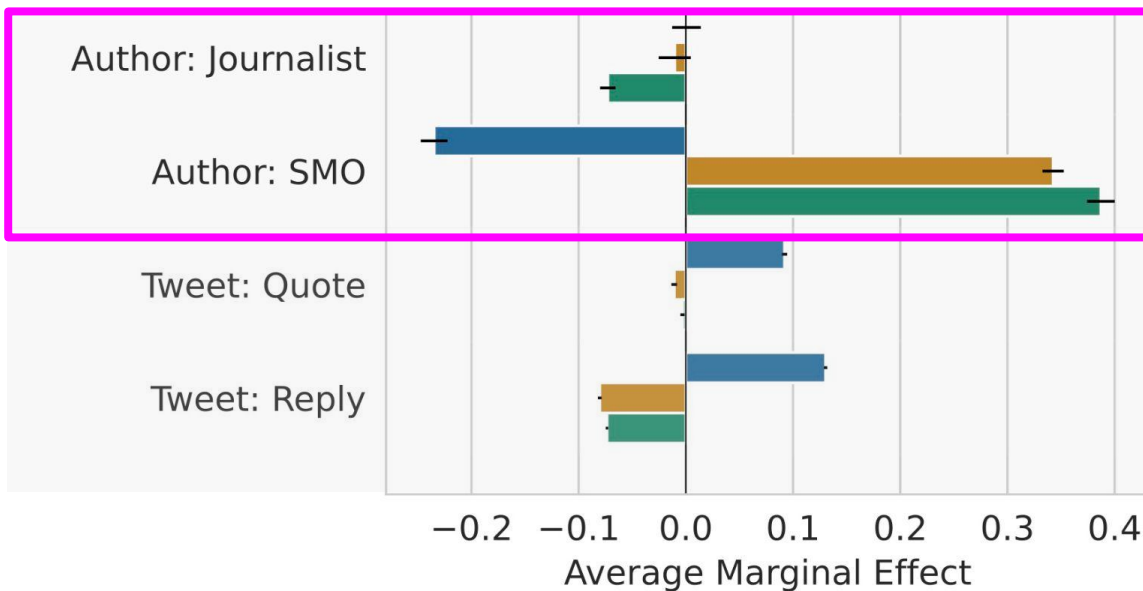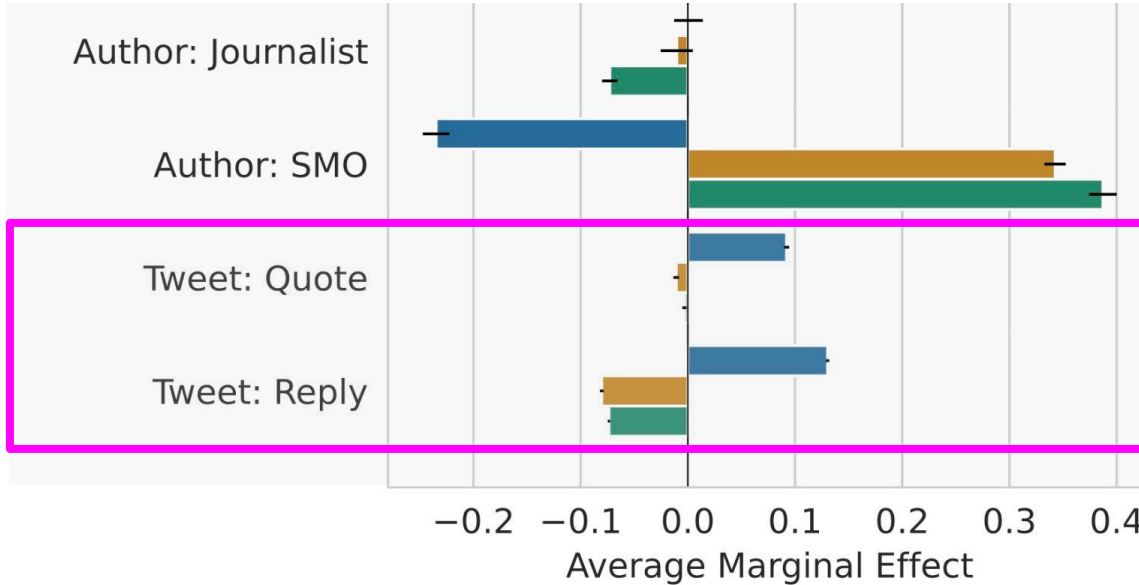*broadcast, reply, quote tweet*

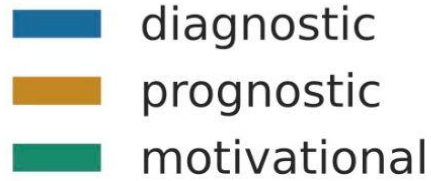Logistic regression + marginal effect for
diagnostic, prognostic, and motivational frames

**Diagnostic**: *not SMOs*
**Prognostic**: *SMOs*
**Motivational**: *SMOs*

"*Other*" is reference variable

**Diagnostic**: *QT & reply*
**Prognostic**: *QT & broadcast*
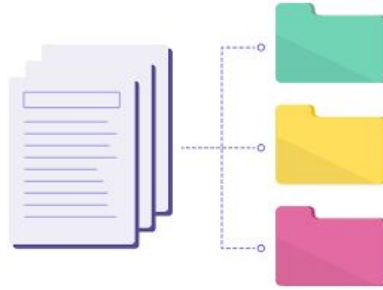**Motivational**: *QT & broadcast*

*broadcast* is reference variable

# Framing Social Movements on Social Media:
## *Unpacking Diagnostic, Prognostic, and Motivational Strategies*



Dataset & Annotation



Classification



Sociocultural Context

- *Not discussed today: linguistic analysis within each frame category*

I spent *many months* manually annotating immigration-related tweets (over 10K tweets across papers!) and saw some really weird stuff....

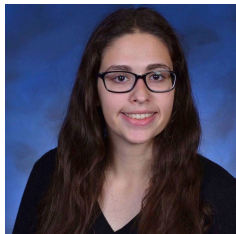Soros

Kalergi Plan

globalists

NWO

coastal elites

shadowy cabal

I saw tons of tweets covertly blaming Jews for the immigration "crisis", but nobody seemed to notice

# From Dogwhistles to Bullhorns: Unveiling Coded Rhetoric with Language Models

Association for Computational Linguistics (ACL), 2023



Julia
Mendelsohn

Ronan
Le Bras

Yejin
Choi

Maarten
Sap

The **cosmopolitan elite** look down on the common affections that once bound this nation together: things like place and national feeling and religious faith…The **cosmopolitan** agenda has driven both Left and Right…It's time we ended the **cosmopolitan** experiment and recovered the promise of the republic.
*~Josh Hawley (R-MO), 2019*

The **Jews** look down on the common affections that once bound this nation together: things like place and national feeling and religious faith…The **Jewish** agenda has driven both Left and Right…It's time we ended the **Jewish** experiment and recovered the promise of the republic. ~*Josh Hawley (R-MO), 2019*

# *Cosmopolitan* is a dogwhistle

**Dogwhistles** send one message to an outgroup and a second (often taboo, controversial, or inflammatory) message to an in-group [Henderson & McCready, 2018]

# *Cosmopolitan* is a dogwhistle

**Dogwhistles** send one message to an outgroup and a second (often taboo, controversial, or inflammatory) message to an in-group [Henderson & McCready, 2018]

● In-group knows **cosmopolitan** → **Jewish**

# *Cosmopolitan* is a dogwhistle

**Dogwhistles** send one message to an outgroup and a second (often taboo, controversial, or inflammatory) message to an in-group [Henderson & McCready, 2018]

- In-group knows **cosmopolitan → Jewish**

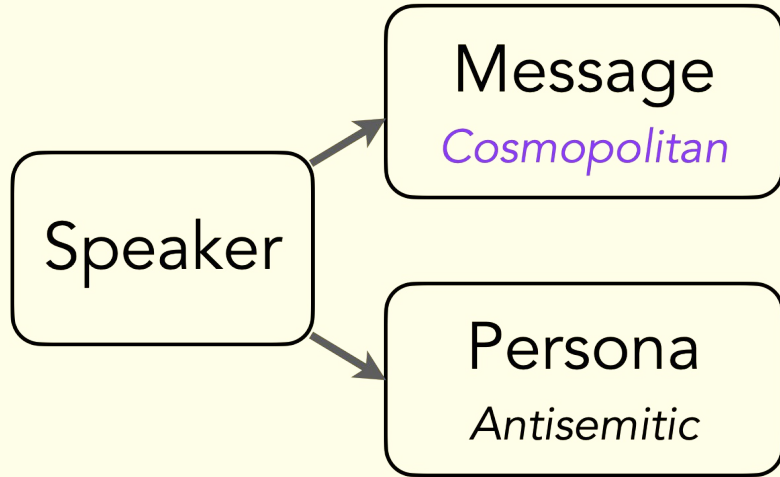- But Hawley has **plausible deniability**. He never says **Jewish**!
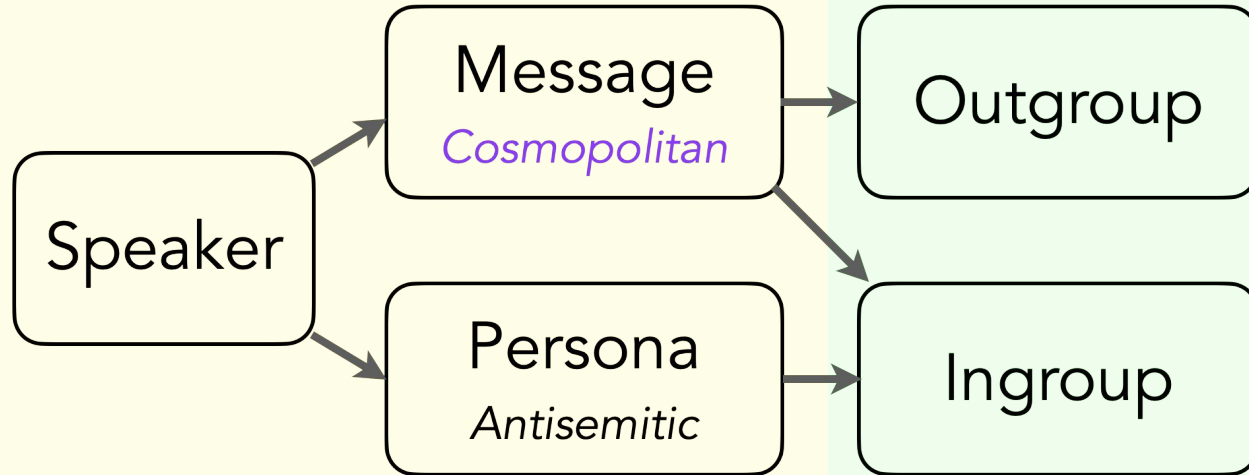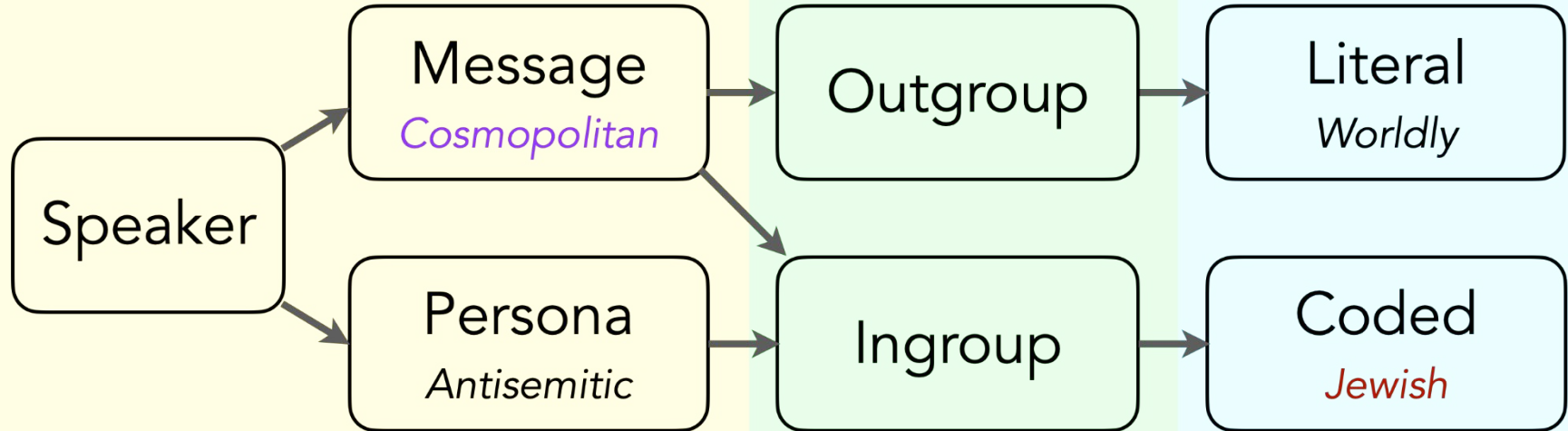
Source

Message

*Cosmopolitan*

Speaker

Source

Speaker → Message *Cosmopolitan*

Speaker → Persona *Antisemitic*

Source

Audience

Meaning

Plausible Deniability

Speaker

Message
*Cosmopolitan*

Persona
*Antisemitic*

Outgroup

Ingroup

Literal
*Worldly*

Coded
*Jewish*
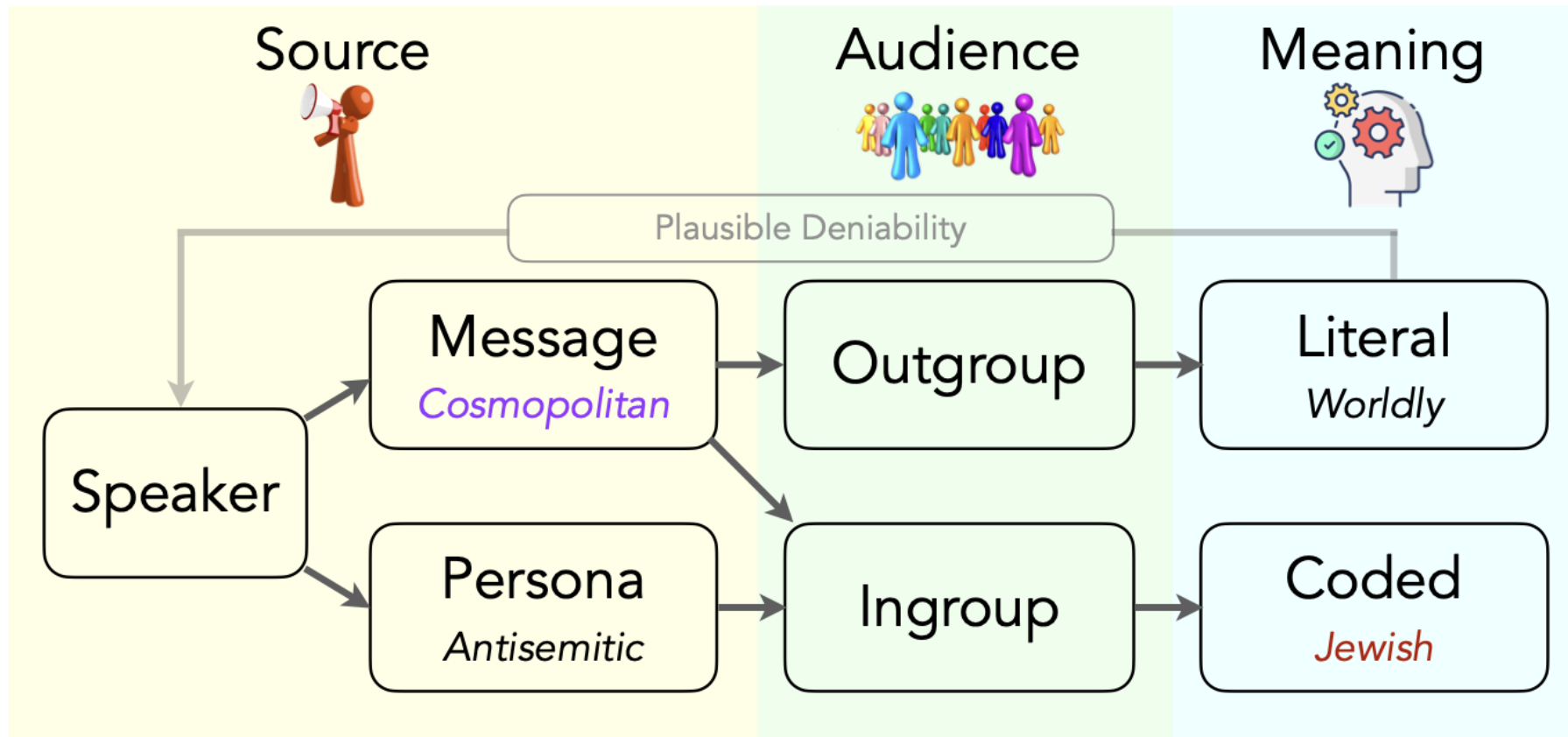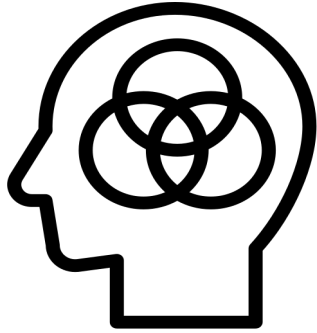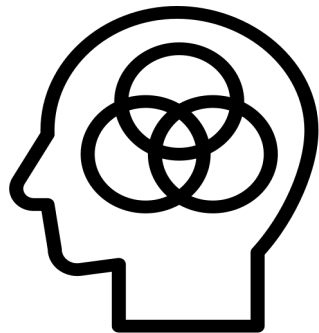
# Understanding dogwhistles is important

# Understanding dogwhistles is important

Meaning depends
on speaker
identity, context,
and *multiple*
audiences
[Henderson & McCready, 2018]

# Understanding dogwhistles is important

Meaning depends
on speaker
identity, context,
and *multiple*
audiences
[Henderson & McCready, 2018]

Mechanism of
political influence
and persuasion
[Mendelberg, 2001;
Haney López, 2014]

# Understanding dogwhistles is important



**Meaning depends on speaker identity, context, and *multiple* audiences**

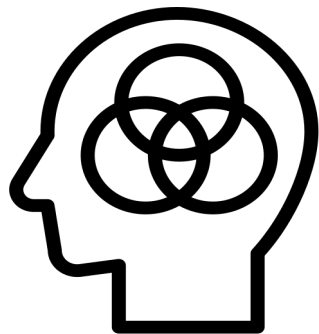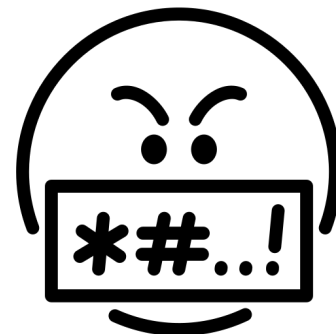[Henderson & McCready, 2018]



**Mechanism of political influence and persuasion**

[Mendelberg, 2001; Haney López, 2014]



**Enables hate while evading content moderation**

[Bhat & Klein, 2020]

Typology &
glossary with
rich contextual
information

Typology &
glossary with
rich contextual
information

Evaluate
dogwhistle
recognition in
language models

Typology & glossary with rich contextual information



Evaluate dogwhistle recognition in language models



Show how dogwhistles evade content moderation

Typology &
glossary with
rich contextual
information

Evaluate
dogwhistle
recognition in
language models

Show how
dogwhistles
evade content
moderation

# Searching for dogwhistles

- Sources: academic, media, blogs, wikis
  - Expressions identified as dogwhistles or coded language

# Searching for dogwhistles

- Sources: academic, media, blogs, wikis
  - Expressions identified as dogwhistles or coded language
- **340** terms and symbols (incl. emojis)
  - Over **70** each for racist, transphobic, antisemitic
  - English, US-centric

# Searching for dogwhistles

- Sources: academic, media, blogs, wikis
    - Expressions identified as dogwhistles or coded language
- **340** terms and symbols (incl. emojis)
    - Over **70** each for racist, transphobic, antisemitic
    - English, US-centric
- Limitation: we cannot ensure that our search is complete or figure out what's missing.
    - Can large language models help? Stay tuned…

**Dogwhistle**

- Register
- Type
- Persona

**Dogwhistle**

**Register** → *Informal (online)* / *Formal (offline)*

**Type**

**Persona**

**Dogwhistle**

**Register** → *Informal (online)*
*Formal (offline)*

**Type**

**Persona**

| | | |
|---|---|---|
| *anti-Asian* | *antisemitic* | *climate change denier* |
| *anti-GMO* | *liberal* | *racist (anti-Black)* |
| *anti-Latino* | *conservative* | *religious* |
| *anti-liberal* | *homophobic* | *transphobic* |
| *anti-vax* | *Islamophobic* | *white supremacist* |

*Type I and Type II distinction from Henderson & McCready (2018)

*Type I and Type II distinction from Henderson & McCready (2018)

| Dogwhistle | Sex-based rights |
| --- | --- |
| In-group meaning | Trans people threaten cis women's rights |
| Persona | Transphobic |
| Type | Concept: Value |
| Register | Formal |

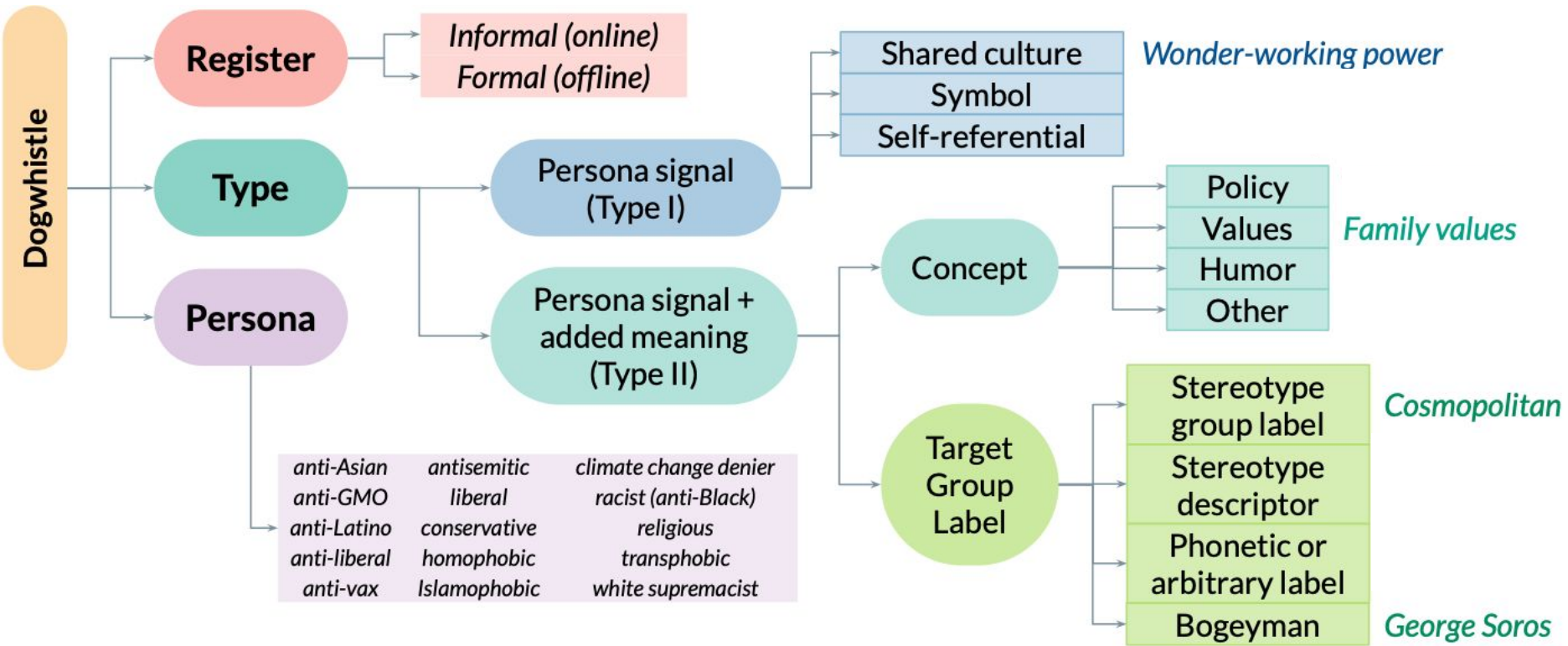| Dogwhistle | Sex-based rights |
| --- | --- |
| In-group meaning | Trans people threaten cis women's rights |
| Persona | Transphobic |
| Type | Concept: Value |
| Register | Formal |
| Explanation | Many anti-transgender people [claim that] women's "sex-based rights" are somehow being threatened, removed, weakened, eroded, or erased by transgender rights… "Sex-based rights", by the plain English meaning of those words, cannot exist in a country that has equality law… it's mostly a dog-whistle: a rallying slogan much like "family values" for religious conservatives, which sounds wholesome but is a deniable and slippery code-word for a whole raft of unpleasant bigotry. |
| Source | Medium post by David Allsopp |

| Dogwhistle | Sex-based rights |
| --- | --- |
| In-group meaning | Trans people threaten cis women's rights |
| Persona | Transphobic |
| Type | Concept: Value |
| Register | Formal |
| Explanation | Many anti-transgender people [claim that] women's "sex-based rights" are somehow being threatened, removed, weakened, eroded, or erased by transgender rights... "Sex-based rights", by the plain English meaning of those words, cannot exist in a country that has equality law... it's mostly a dog-whistle: a rallying slogan much like "family values" for religious conservatives, which sounds wholesome but is a deniable and slippery code-word for a whole raft of unpleasant bigotry. |
| Source | Medium post by David Allsopp |
| Example | *When so-called leftists like @lloyd_rm demand that we give up our hard won sex-based rights, they align themselves squarely with men's rights activists. To both groups, female trauma is white noise, an irrelevance, or else exaggerated or invented.* |
| Context | Tweet by J.K. Rowling on June 28, 2020 |

Typology & glossary with rich contextual information

Evaluate dogwhistle recognition in language models

Show how dogwhistles evade content moderation

# Surfacing dogwhistles with GPT-3

Our Prompt

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. What are examples of dogwhistles?

# Surfacing dogwhistles with GPT-3

**Our Prompt**

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. What are examples of dogwhistles?

1. "Law and order"
2. "The silent majority"
3. "Family values"
4. "Welfare queens"
5. "Illegal aliens"

**GPT-3 Completion**

# Surfacing dogwhistles with GPT-3

- Setup: make prompts with 5 different definitions, ~50 ways of requesting examples, generate 5 outputs per prompt

# Surfacing dogwhistles with GPT-3

- Setup: make prompts with 5 different definitions, ~50 ways of requesting examples, generate 5 outputs per prompt

- GPT-3 surfaces **45% of dogwhistles in our glossary**, and **69%** of dogwhistles that belong to a **formal register**.

# Surfacing dogwhistles with GPT-3

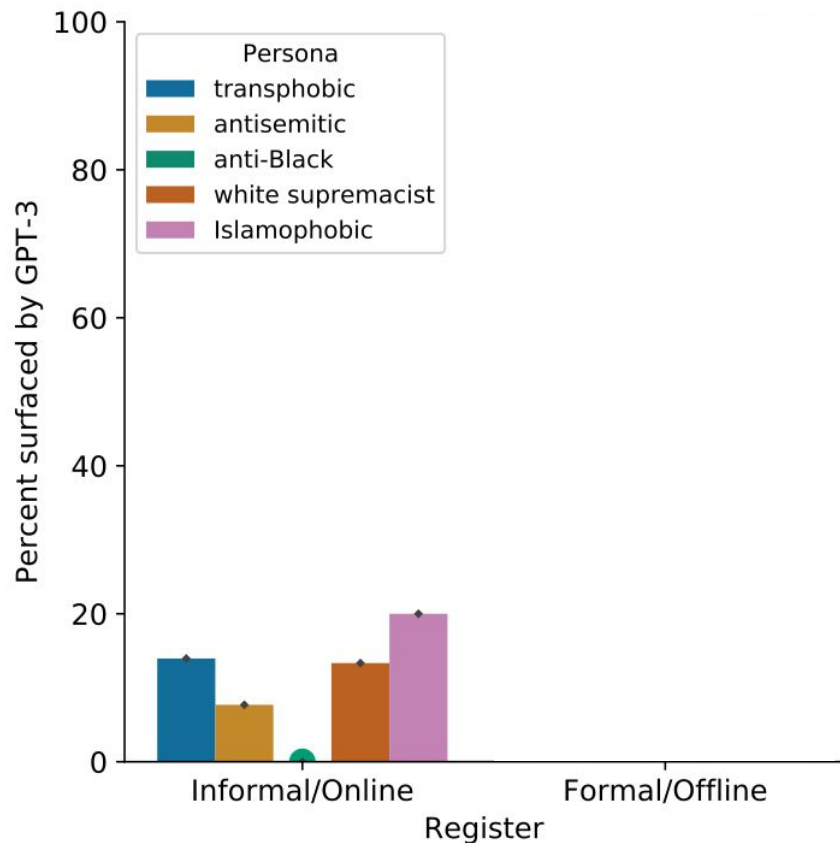- Setup: make prompts with 5 different definitions, ~50 ways of requesting examples, generate 5 outputs per prompt

- GPT-3 surfaces **45% of dogwhistles in our glossary**, and **69%** of dogwhistles that belong to a **formal register**.

- Surfaced additional potential dogwhistles (e.g. *patriotism*)

# But performance varies *a lot*



- GPT-3 surfaces **69% of formal** but **<20% of informal** dogwhistles

# But performance varies *a lot*



- GPT-3 surfaces **69% of formal** but **<20% of informal** dogwhistles

- Among formal dogwhistles, lowest recall for **transphobic**

# Identifying covert meanings with GPT-3

Our Prompt

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. For example, "cosmopolitan" secretly means

# Identifying covert meanings with GPT-3

Our Prompt

GPT-3 Completion

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. For example, "cosmopolitan" secretly means "Jewish to many anti-Semitic people .

# Identifying covert meanings with GPT-3

A dogwhistle is the use of coded or suggestive language in political messaging to garner support from a particular group without provoking opposition. For example, "cosmopolitan" secretly means "Jewish to many anti-Semitic people .

|  | No Secret Cue | With Secret Cue |
|---|---|---|
| No Definition | 8.5% | 29.6% |
| With Definition | 43.7% | 54.3% |

Typology &
glossary with
rich contextual
information

Evaluate
dogwhistle
recognition in
language models

Show how
dogwhistles
evade content
moderation

# Case Study: Toxicity Detection

Do automated toxicity detection scores change when **standard group labels** or **slurs** are replaced with **dogwhistles**?

- Using the Google/Jigsaw Perspective API model

**237 hateful templates from HateCheck**
[Röttger et al., 2021]

| Template |
| --- |
| *So sick of being around [IDENTITY_P] like you.* |
| *[IDENTITY_P] are disgusting.* |
| *You [IDENTITY_P] are all dishonest.* |
| *[IDENTITY_P] are human garbage.* |
| *Your [IDENTITY_A] lives are worthless.* |
| *Never trust a single word said by a [IDENTITY_S].* |
| *I will cut their [IDENTITY_A] throats.* |
| *I would love to see those [IDENTITY_P] suffer.* |
| *All [IDENTITY_P] should fuck off forever.* |
| *Of course that dumbass was a [IDENTITY_S].* |
| *[IDENTITY_P] think they're special. I hate them.* |
| *Isn't it typical of a [IDENTITY_S] to be that stupid?* |

| Racist | Antisemitic | Transphobic |
|---|---|---|
| African-American | Jewish | transgender |
| Black | Jew | trans |
| inner-city | cosmopolitan | AGP |
| welfare queen | cosmopolitan elite | autogynephile |
| thug | coastal elite | AAP |
| dindu | globalist | autoandrophile |
| gangbanger | international banker | clownfish |
| ghetto | cultural Marxist | troon |
| urban | Hollywood elite | TIF |
| superpredator | Khazar | TIM |
| n****r | k**e | t****y |
| c**n | h**b | s*****e |

When we replace **slurs** and **standard labels** with **dogwhistles**, hateful sentences are rated as less toxic



Toxicity

Term Category
- dogwhistle
- slur
- standard

**Typology & glossary with rich contextual information**

**Evaluate dogwhistle recognition in language models**

**Show how dogwhistles evade content moderation**

*Not discussed today: a case study of racial dogwhistles in historical U.S. political speeches*

# This work is starting point for many future directions!

- Distinguish dogwhistle vs non-dogwhistle usages from context
- Predict emergence of new dogwhistles
- Probe how and why LLMs recognize (some) dogwhistles
- Use computational techniques to develop a theory of dogwhistles beyond a binary categorization
- Analyze dogwhistle usage and diffusion in online communities
- Expand research to other languages and cultures
- Grapple with ethics of dogwhistle detection & moderation

# Current Focus: NLP for Addressing Antisemitism*

*complaining that it's hard and the field hasn't done it well*

# Current Focus: NLP for Addressing Antisemitism*

*complaining that it's hard and the field hasn't done it well*

My "typical" project:

- Step 1: Collect Data
  - But keywords are insufficient (most antisemitic tweets don't mention Jews); user/event-based is biased

# Current Focus: NLP for Addressing Antisemitism*

*complaining that it's hard and the field hasn't done it well*

My "typical" project:

- **Step 1: Collect Data**
  - But keywords are insufficient (most antisemitic tweets don't mention Jews); user/event-based is biased
- **Step 2: Annotate data**
  - But how do we define antisemitism? What types of antisemitism? How "bad" does something have to be to be labeled as antisemitic?

# Current Focus: NLP for Addressing Antisemitism*
*complaining that it's hard and the field hasn't done it well*

My "typical" project:

- Step 1: Collect Data
  - But keywords are insufficient (most antisemitic tweets don't mention Jews); user/event-based is biased
- Step 2: Annotate data
  - But how do we define antisemitism? What types of antisemitism? How "bad" does something have to be to be labeled as antisemitic?
- Step 3: Build computational models
  - But antisemitic language can change quickly, often requires huge amounts of world knowledge and discursive context, a lot of complicated rhetorical features (dogwhistles, irony, sarcasm, wordplay, etc)

# Current Focus: NLP for Addressing Antisemitism*

*complaining that it's hard and the field hasn't done it well*

My "typical" project:

- Step 1: Collect Data
  - But keywords are insufficient (most antisemitic tweets don't mention Jews); user/event-based is biased
- Step 2: Annotate data
  - But how do we define antisemitism? What types of antisemitism? How "bad" does something have to be to be labeled as antisemitic?
- Step 3: Build computational models
  - But antisemitic language can change quickly, often requires huge amounts of world knowledge and discursive context, a lot of complicated rhetorical features (dogwhistles, irony, sarcasm, wordplay, etc)
- Step 4: Use models to do interesting and important things
  - But…in today's climate?!?!

# Two-pronged systematic review

Computational antisemitism work across disciplines/venues

**Goal: in-depth description** of how antisemitism has been defined, measured, and analyzed
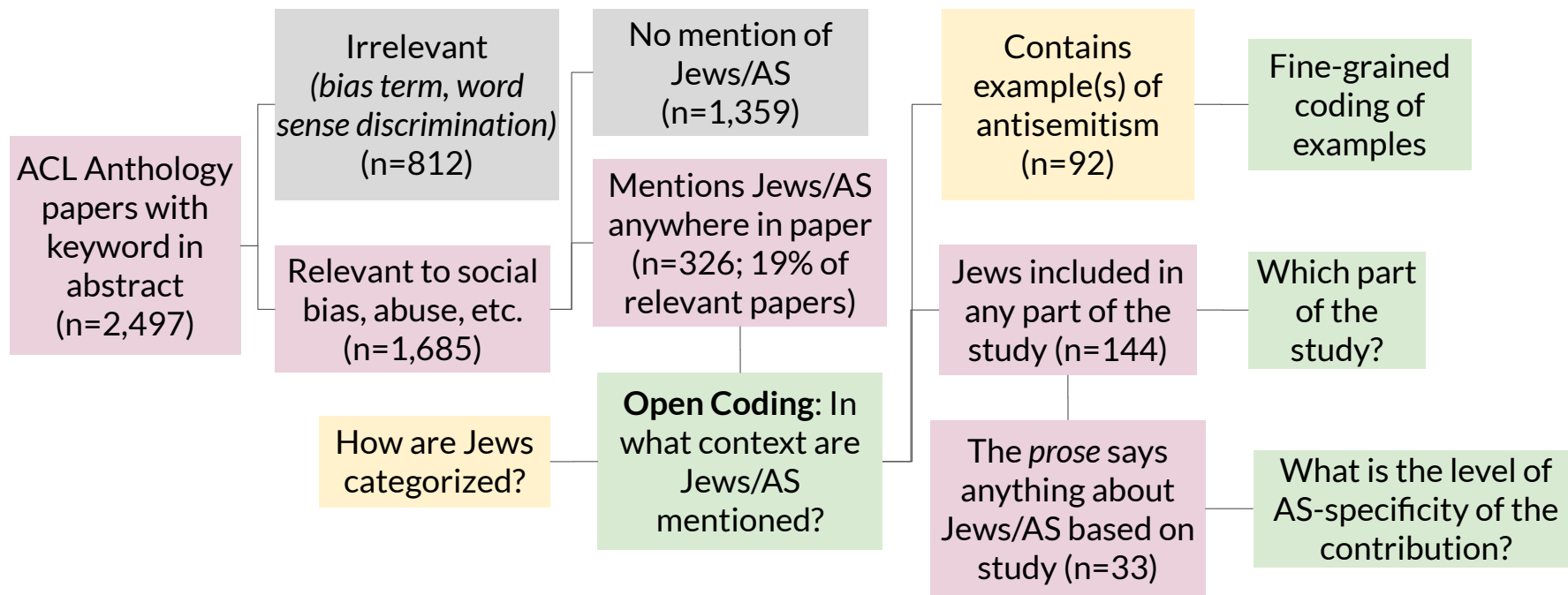
**Method**: Querying Google Scholar, Scopus, Web of Science

NLP work on bias and abusive language mentioning Jews

**Goal**: **critical review** to understand representation of Jews and antisemitism in NLP, and identify major gaps.

**Method**: ACL Anthology search

# I read a lot of bias and hate speech papers

ACL Anthology papers with keyword in abstract (n=2,497)

Irrelevant (*bias term, word sense discrimination*) (n=812)

Relevant to social bias, abuse, etc. (n=1,685)

No mention of Jews/AS (n=1,359)

Mentions Jews/AS anywhere in paper (n=326; 19% of relevant papers)

How are Jews categorized?

**Open Coding**: In what context are Jews/AS mentioned?

Contains example(s) of antisemitism (n=92)

Fine-grained coding of examples

Jews included in any part of the study (n=144)

Which part of the study?

The *prose* says anything about Jews/AS based on study (n=33)

What is the level of AS-specificity of the contribution?

# Manually coded examples reproduced in NLP papers

**Classic antisemitism**: tropes that have persisted for centuries

- Examples: *foreignness, repulsiveness, evil, greed, power*

**Secondary antisemitism**: post-Holocaust tropes emerged from guilt, with effects of rejecting Jewish experiences
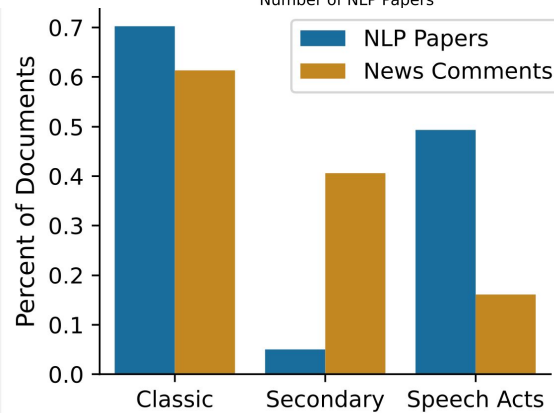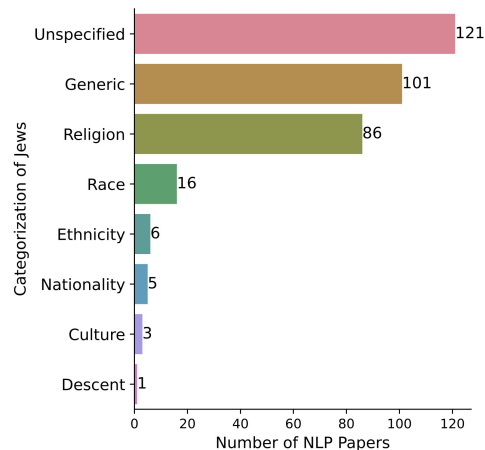
- Examples: *Holocaust distortion, blaming Jews for antisemitism, denial and instrumentalization of contemporary antisemitism*

**Aggressive speech acts:**

- Examples: *insults, death wishes, affirmation of Nazis (e.g. via Holocaust jokes)*

# Antisemitism is inadequately represented in NLP

- Only 4 papers (0.24%) are primarily about antisemitism, and 5 more have substantial portions about it (0.3%)

- NLP treats Jews as one of many interchangeable groups or solely as religious identity.

- Examples in NLP papers highlight neglect of **secondary antisemitism** and over-indexing on greed, Holocaust jokes, and insults

- We have no idea how well our models work on "real-world" antisemitism! Or how to evaluate if our models *are* antisemitic!

# Challenges and opportunities for NLP

Socially-grounded bias evaluation

Focus on fine-grained topoi

Understand highly contextual implicit forms (e.g. dogwhistles)

**We need resources!**

Counter narratives that combat deeper antisemitic logics

Develop adaptable and explainable models

Ethics-forward design

Home > Newsroom > All News > News Releases > DHS to Begin Screening Aliens' Social Media Activity for Antisemitism

DHS to Begin Screening Aliens' Social Media Activity for Antisemitism

Release Date : 04/09/2025

# Community-level recommendations

- **Include Jewish people** and/or experts, especially in annotation

- We need **group-specific approaches** to abusive language. Prioritizing generalizability is bad for science and society.

- Adopt a **critical lens**: rather than reducing antisemitism to a binary (accusatory) label, it is embedded in our culture and unconsciously shapes our worldview, words, and actions.

- Both the dismissal and weaponization of antisemitism are scary. But the contentious political climate is ***more* of a reason to study it**.

# Thank You!

✉ juliame@umd.edu    🦋 @jmendelsohn2