



# Detecting and Understanding Information Disorder

---

Lin Ai

Senior Applied Scientist, Microsoft  
PhD work with Prof. Hirschberg, Columbia University



# What is information disorder?

## Content types







- Misinformation • Disinformation • Propaganda
- Manipulative persuasion

## How it looks

- False/unsupported claims
- Context stripping & selective framing
- Fabricated/edited media
- etc. ...



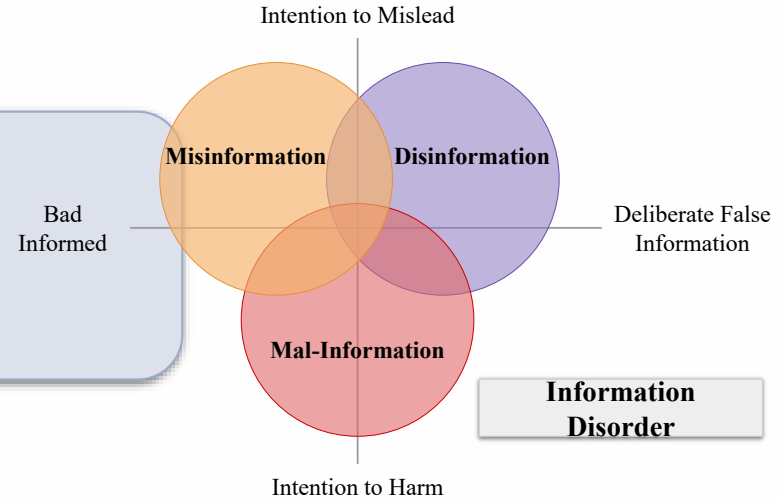
## 7 COMMON FORMS OF INFORMATION DISORDER

-  SATIRE OR PARODY
-  FALSE CONNECTION
-  MISLEADING CONTENT
-  FALSE CONTEXT
-  IMPOSTER CONTENT
-  MANIPULATED CONTENT
-  FABRICATED CONTENT

# What is information disorder?

## Intent dimension

- Intent to Mislead
- Intent to harm



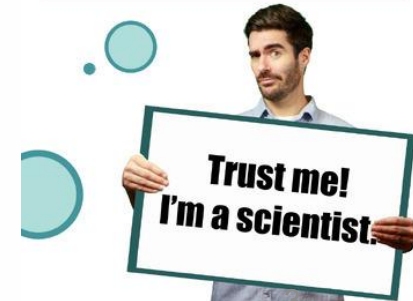
## How it works

- Emotional appeals, authority mimicry, etc
- Agenda-driven narratives
- Target specific tactics

### The Fallacy Understander



### Appeal to Authority



## Who uses it—and who it targets?

### Actors

- State-affiliated media & campaigns
- Influencers & politicians
- Issue networks & for-profit operations
- Emerging: LLM agents

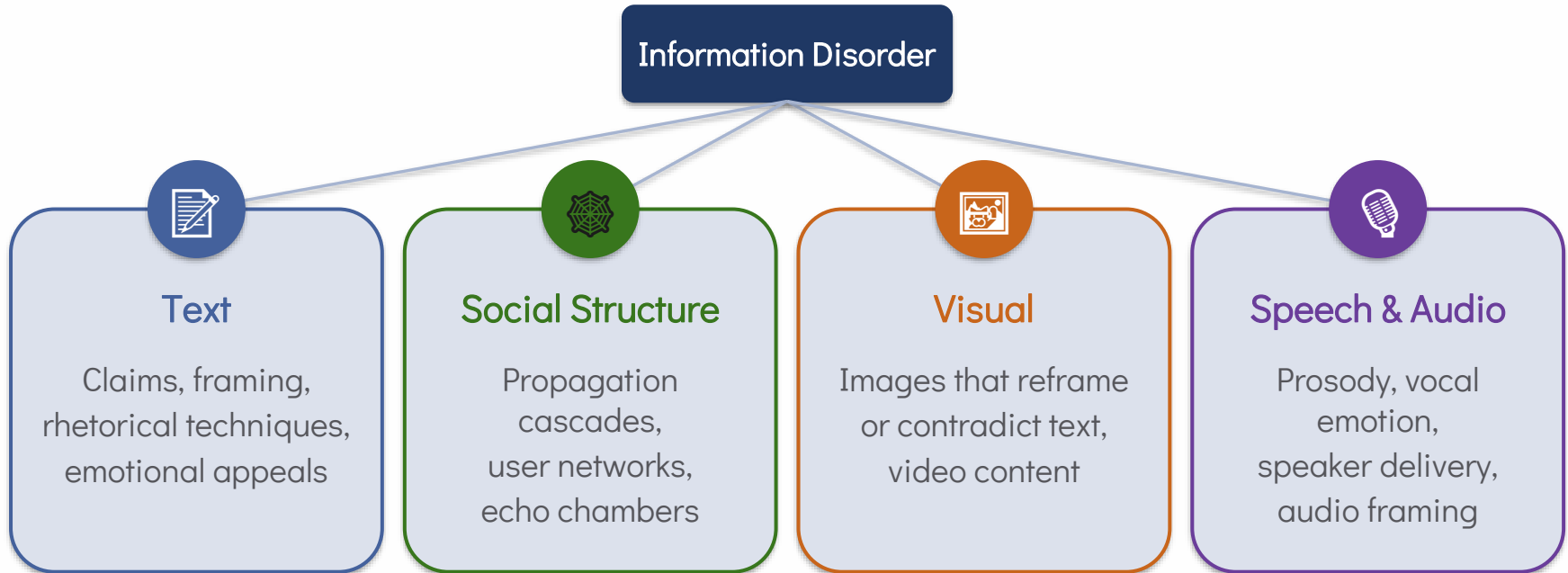
### Audiences

- General public segments
- High-susceptibility cohorts
- Issue-aligned communities
- Decision makers

### Intents

- Mislead or Deceive
- Propaganda
- Radicalize
- etc. ...

# Information disorder leaves traces across modalities

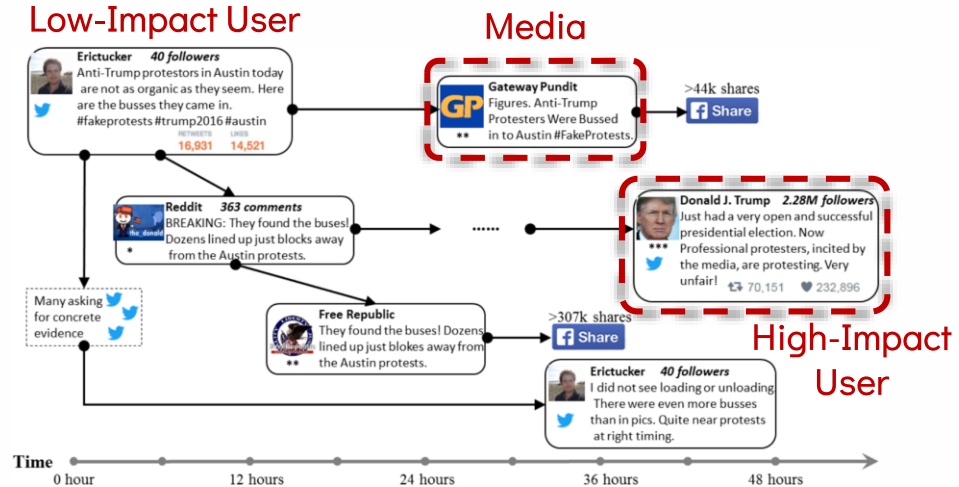
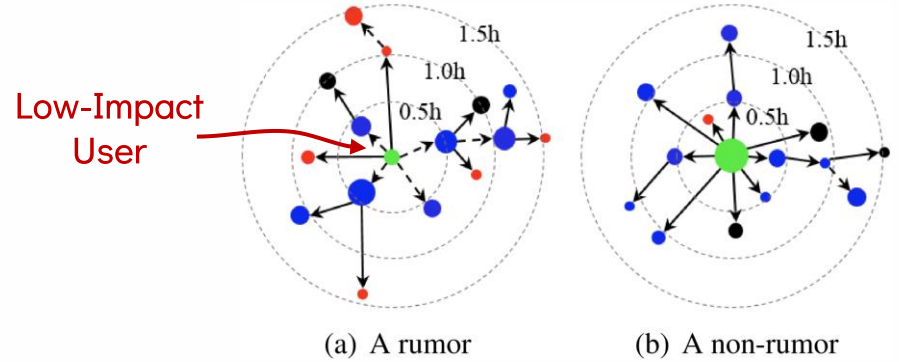


*Today's talk draws on all four – including speech prosody and vocal emotion in the audience-perception study*

# Background

- Early work on information disorder detection rely primarily on linguistic features
- Propagation-based approaches make use of tree-structured propagation patterns of microblog posts and learn contextual representations

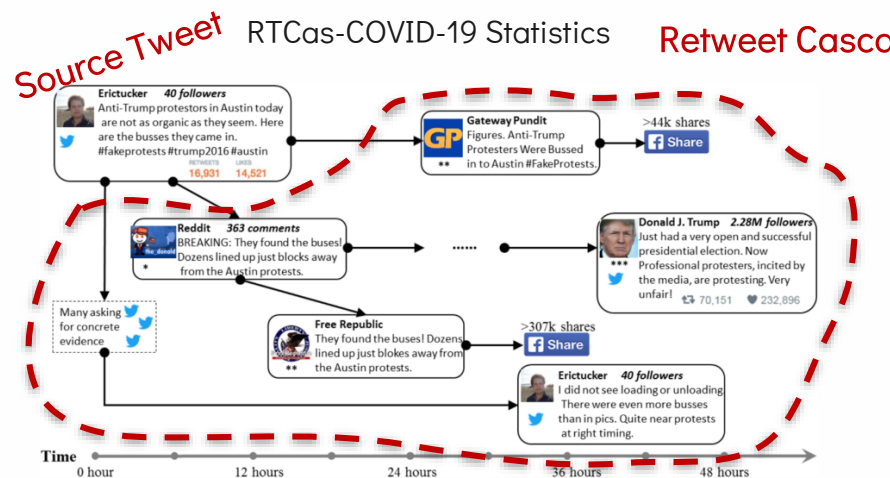
➔ We tackle the challenge by incorporating linguistic, social context, and tweet propagation patterns



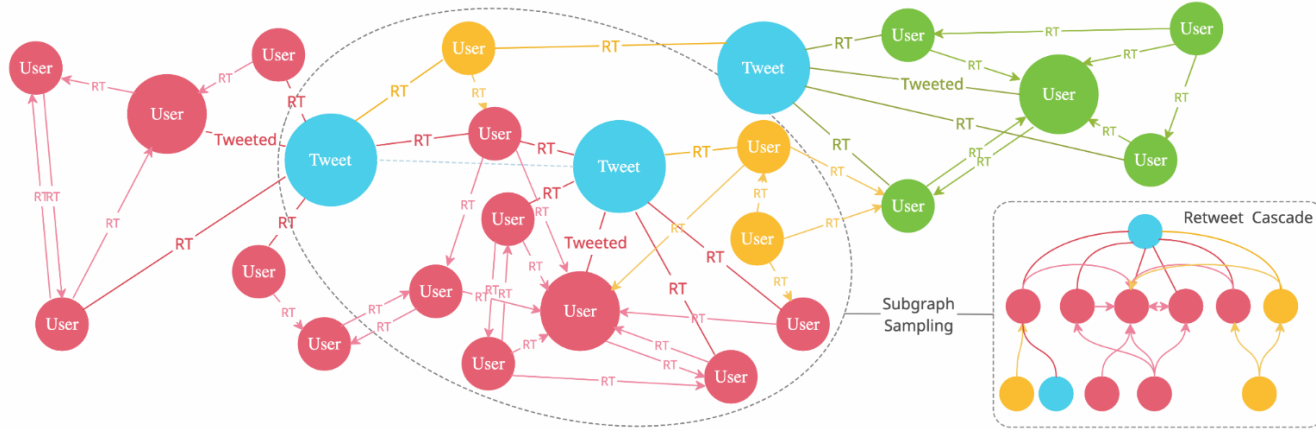
## RTCas-COVID-19 Corpus

- We collect and clean a COVID-19 corpus, **RTCas-COVID-19**, using two publicly available datasets
- Cleaned corpus: **35M tweets** (10M source tweets and retweet cascades)
- Retweet Cascade**: Tree-like structure representing how tweets spread through retweets and replies
- We weak-label a subset of tweets using URLs source credibility
  - Weak-labeled corpus: **2M tweets**

	Total	Source	Retweets
Full Corpus	35M	10M	25M
	Total	Trust	Untrust
Weak-Labeled	2M	1.64M	360K
Human-Annotated	380	215	165

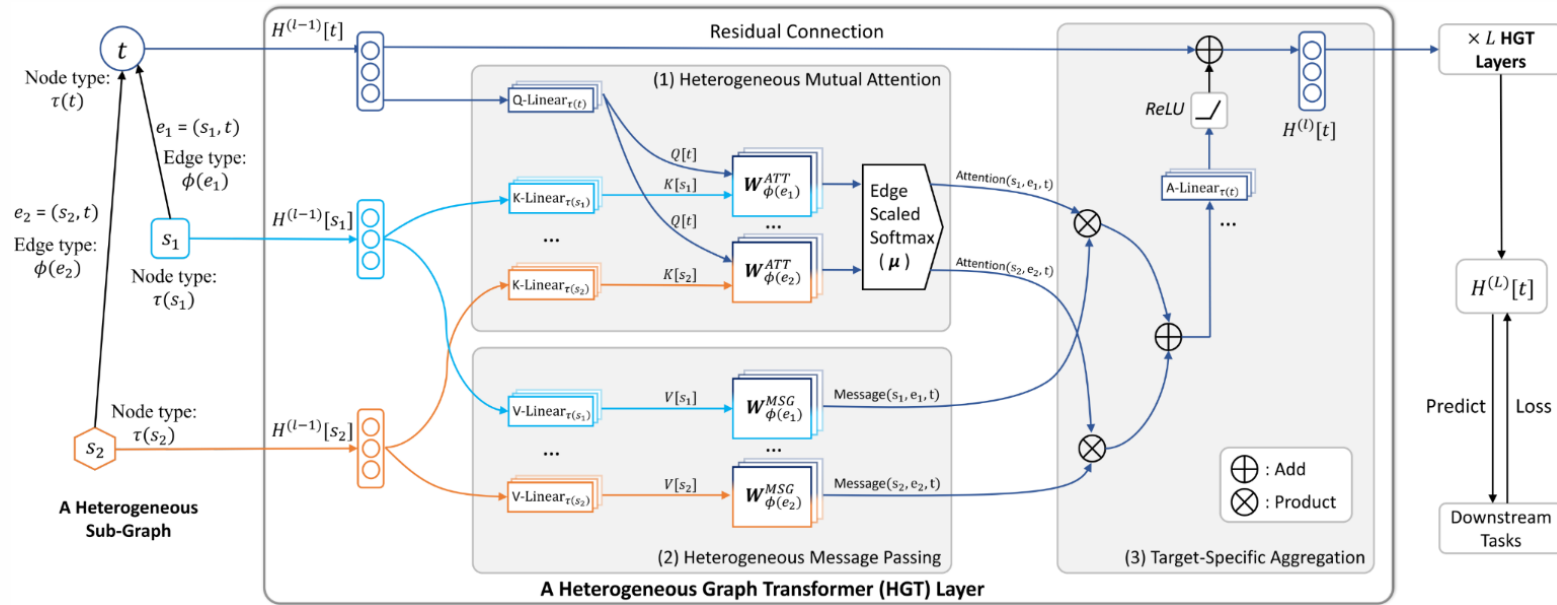


# Tweet-User Heterogeneous Graph



- **Echo chamber:** users tend to interact more with like-minded people
- Nodes are connected by "tweeting" and "retweeting"

# RTCS-HGT (Retweet Cascade Subgraph Sampling Heterogeneous Graph Transformer)



We apply an inductive heterogeneous graph transformer (HGT) for node representation learning

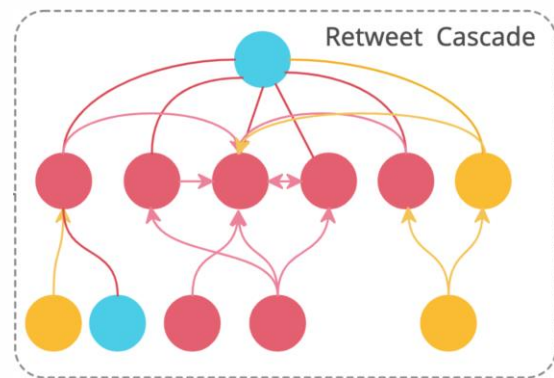
## RTCS-HGT (Retweet Cascade Subgraph Sampling Heterogeneous Graph Transformer)

- Supervised Tweet Classification Loss: A negative log-likelihood classification loss
- Unsupervised User Proximity Loss:

$$proxloss = \mathbb{E}[-\mathbb{E} \log \sigma(z_u^\top z_{v_p}) - Q \cdot \mathbb{E} \log \sigma(-z_u^\top z_{v_n})]$$

where  $P_u$  and  $N_u$  represent neighboring and non-neighboring nodes of user  $u$ , respectively, and  $z_v$  denotes the node representation

- To scale the model to large graphs and reduce training time, we introduce a tweet-centered **retweet cascade subgraph sampling (RTCS)** approach



# Experiments

- RTCS-HGT model outperforms all baselines on the weak-labeled RTCas-COVID-19 subset
- CT-BERT model performs comparably, it is limited to COVID-19 tweets and requires re-training for other topics
- HGATRD is a strong baseline but lacks scalability and cannot infer on unseen data

Model	Test Acc.	Macro F <sub>1</sub>	Trust F <sub>1</sub>	Untrust F <sub>1</sub>
RCNN-LSTM	0.844	0.844	0.846	0.842
RCNN-GRU	0.844	0.843	0.848	0.839
BERTweet	0.847	0.847	0.850	0.843
CT-BERT	0.893	0.893	0.894	0.893
HGATRD	0.894	0.894	0.894	0.895
HGT	0.908	0.908	0.910	0.906
RTCS-HGT	0.913	0.913	0.913	0.912
<b>RTCS-HGT</b> <i>(no proxloss)</i>	<b>0.918</b>	<b>0.918</b>	<b>0.918</b>	<b>0.918</b>

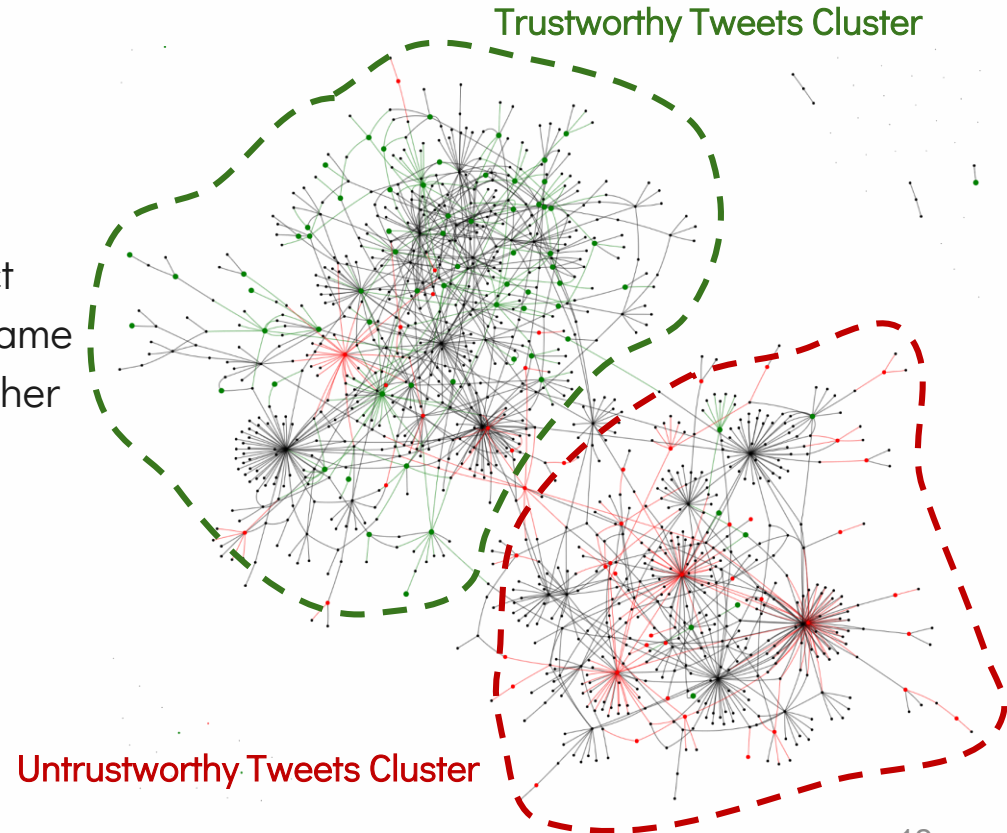
RTCS-HGT vs. Baselines on Weak-Labeled RTCas-COVID-19 Test Sets

## Analysis – Social Context Information

- Nodes naturally form 2 clusters
- Users within each cluster tend to interact more frequently with other users in the same cluster rather than with users from the other cluster



Echo Chamber Effect



## Multimodal Intent Detection in Social Media -- Background

- While previous framework detects untrustworthy information using linguistic and social context, **visual content** plays a key role in social media posts

Happy Earth 🌍 Day! #EarthDay #HappyEarthDay  
#ClimateChange #GlobalWarming



**BENIGN**

Happy Earth 🌍 Day! #EarthDay #HappyEarthDay  
#ClimateChange #GlobalWarming



**MALICIOUS**

# Multimodal Intent Detection in Social Media -- Corpus

- We collect a multi-modal social media corpus of approximately **13K posts**
  - Text + image posts from Twitter and Facebook
- Topics: COVID-19, Climate Change
- Each post is weak-labeled as malicious/benign with 5 fine-grained intent categories: *polarization*, *call-to-action*, *virality*, *sarcasm*, *humor*
  - Posts are weak-labeled using hashtags
  - Viral posts are those among the top 1,500 most retweeted

	Data →	Covid-19	Climate Change
<b>Malicious</b>	Polarizing	1938	1725
	Call-To-Action	2454	-
	Viral	1500	56
	<i>Total</i>	<i>5892</i>	<i>1781</i>
<b>Benign</b>	Sarcasm	2361	323
	Humor	1500	1497
	<i>Total</i>	<i>3861</i>	<i>1820</i>
<b>Total</b>		<b>9753</b>	<b>3601</b>

Detailed statistics of the 13K posts

# Multimodal Intent Detection in Social Media -- Corpus

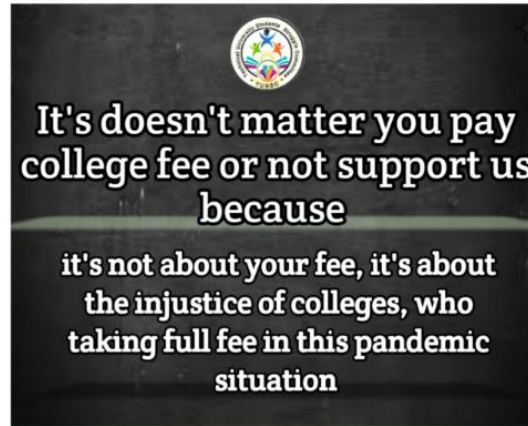
**⚠ Content Disclaimer:** *examples shown contain harmful/misleading content involving real public figures and do **NOT** reflect the views of the presenter or co-authors, nor any factual claims about the individuals depicted.*

Pro vaxxers are always demanding proof then dismiss it when you give it to them. They say vaccines have been “proven” safe but when asked to show said “proof” they give 🙏🙏🙏🙏🙏. Then they bow down and worship their lord and savior #BillGatesBioTerrorist #sad



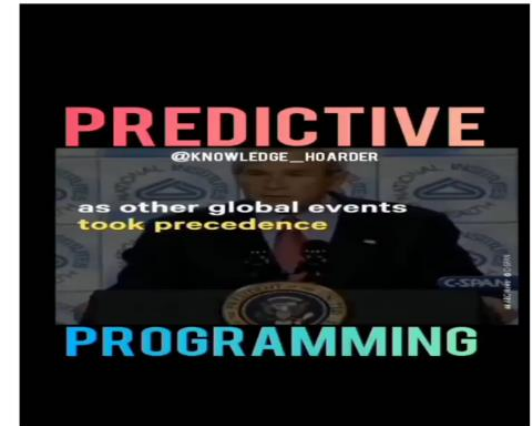
**Polarizing**

Support us on Twitter campaign  
Speakup against fee hike in colleges  
#tussc #ReduceCollegeFees  
#StudentProblems #Standup  
#COVID19 #supportus  
#ThursdayThoughts



**Call-to-Action**

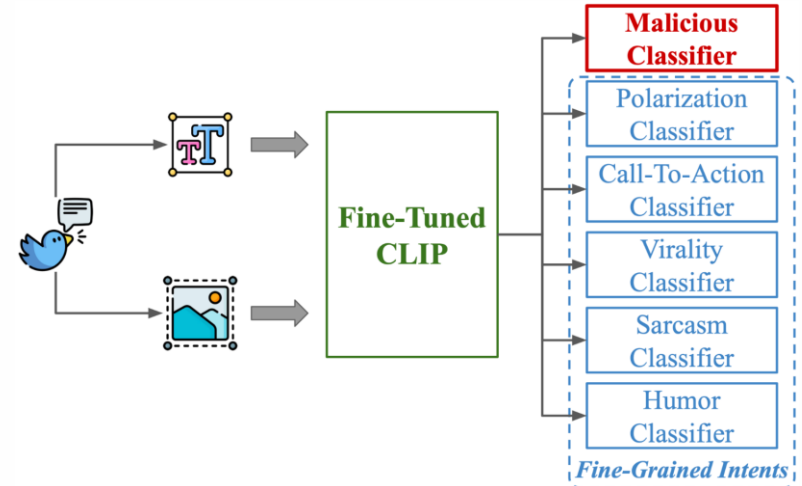
You do not understand what a global reset is.....it is a plan to kill you and your family and your generations to come..... It is a depopulation plan... It is a global slavery and surveillance plan.... It is a plan to create a one world government.... Covid was only a trigger..



**Viral**


## CLIP-MTL

- We fine-tune CLIP's pre-trained text and image encoders to extract linguistic and visual features
- We train the model using **multi-task learning (MTL)**, with shared CLIP weights and separated intent classifiers
- Each task benefits from the shared knowledge of the others



## Experiments and Analysis

- The model achieves a **test accuracy of 0.978** on the **maliciousness classification** task
- It attains an average test accuracy of 0.977 across the five intent classifiers

 Future improvements include addressing potential overfitting and enhancing generalizability

<b>Task ↓</b>	<b>Accuracy</b>	<b>Neg Macro F1</b>	<b>Pos Macro F1</b>
Polarization	0.970	0.981	0.932
Call-To-Action	0.979	0.988	0.923
Virality	0.985	0.992	0.914
Sarcasm	0.983	0.990	0.934
Humor	0.966	0.979	0.905
<b>Maliciousness</b>	<b>0.978</b>	<b>0.979</b>	<b>0.977</b>

Performance on malicious classification and each specific intent classification

## Warning from Heterogeneous Signals

- Combining textual cues with propagation lifts early detection and reveals community patterns missed by text only
- Visual features aid to uncover malicious narrative intent
- High-recall first line of defense at scale

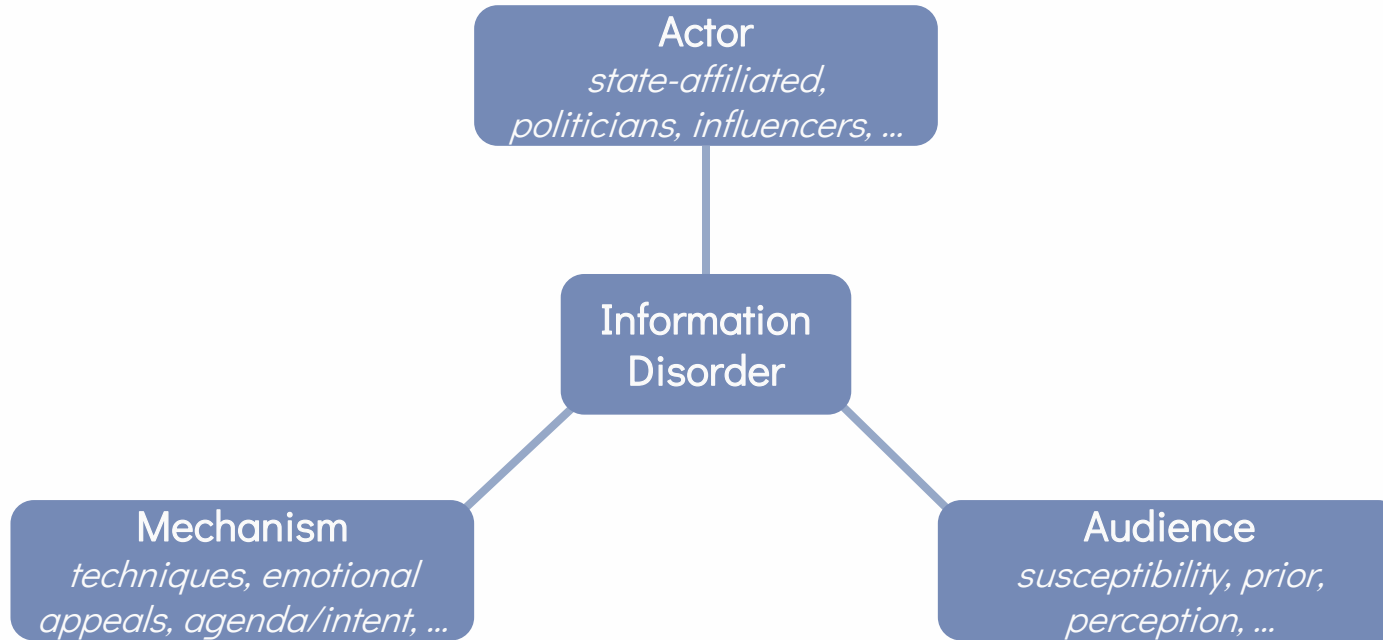
Findings / Takeaways



- Detection ≠ persuasion model – we still don't know why the content works
- We don't identify who is most susceptible (audience traits, perceptions)
- Labels hint at intent, but we need actor-level narratives and technique/appeal structure

Gaps → Next

## Fine-Grained Understand



# Manipulative Narrative of State Actors

Actor  
*state-affiliated,  
politicians, influencers, ...*

- We dive deeper into **why** and **how** specific content is spread to convey actors' messages
- Researchers have also focused on analyzing **propaganda** campaigns, notably from Russia

➔ Our work aims to dive deeper into this, focusing on narratives from both Russia and Ukraine



Kyiv, our splendid, peaceful city, survived another night under attacks by Russian ground forces, missiles. One of them has hit a residential apartment in Kyiv. I demand the world: fully isolate Russia, expel ambassadors, oil embargo, ruin its economy. Stop Russian war criminals! #russianwarcriminals #stoprussia



This tweet is a **call-to-action**, but **what specific actions** are being urged, and who is the intended **target audience**?

**⚠ Research Disclaimer:** *this work analyzes state-affiliated messaging from both Russian and Ukrainian government accounts as a computational case study. It does **NOT** take a political position on the conflict, endorse any actor's narrative, or equate the actions of the parties involved.*

# Fine-Grained Intent

- **Call To Action (CTA)**
  - *Called Subjects*: the intended audience
  - *Called Actions*: specific actions
- **Discredit Entity (DE)**
  - *Discredited Subjects*: the entities being undermined
  - *Discrediting Phrases*: specific language used to discredit these entities

Kyiv, our splendid, peaceful city, survived another night under attacks by Russian ground forces, missiles. One of them has hit a residential apartment in Kyiv. I demand **the world**<sup>1</sup>: fully isolate Russia, expel ambassadors, oil embargo, ruin its economy.<sup>2</sup> Stop Russian war criminals!<sup>3</sup>

1. Called Subject
2. Called Actions, using dictatorship, black-and-white fallacy, thought-terminating cliché
3. Called Actions, using slogans, flag-waving

Example Tweet Annotated as **CTA**

Dear Black people, Don't let Starr convince you to defend the **US/NATO**<sup>1</sup> aggressions that caused this conflict. US/NATO **turned Libya into a hell scape**<sup>2</sup> for Black people and they're **at fault in Ukraine**<sup>3</sup> too. There are neo-Nazis in Ukraine government. Fact.  
 @terrelljstarr #NotoNATO

1. Discredited Subjects
2. Discrediting Phrases, using causal oversimplification, obfuscation
3. Called Actions, using causal oversimplification, obfuscation

Example Tweet Annotated as **DE**

# TweetIntent@Crisis



## Collect

- 67 RU / 12 UA gov-affiliated Twitter accounts



## Seed human labels

- 5,000 pre-filtered by GPT-4 → 3,691 valid
- CTA: 93 (≈1.9%) DE: 411 (≈8.8%)

**Scarce Positives**



## Model Selection

- (a) Intent (b) Interrogative span localization
- Fine-tuned GPT-3.5-Turbo > GPT-4 / GPT-4+ICL



## Machine annotation

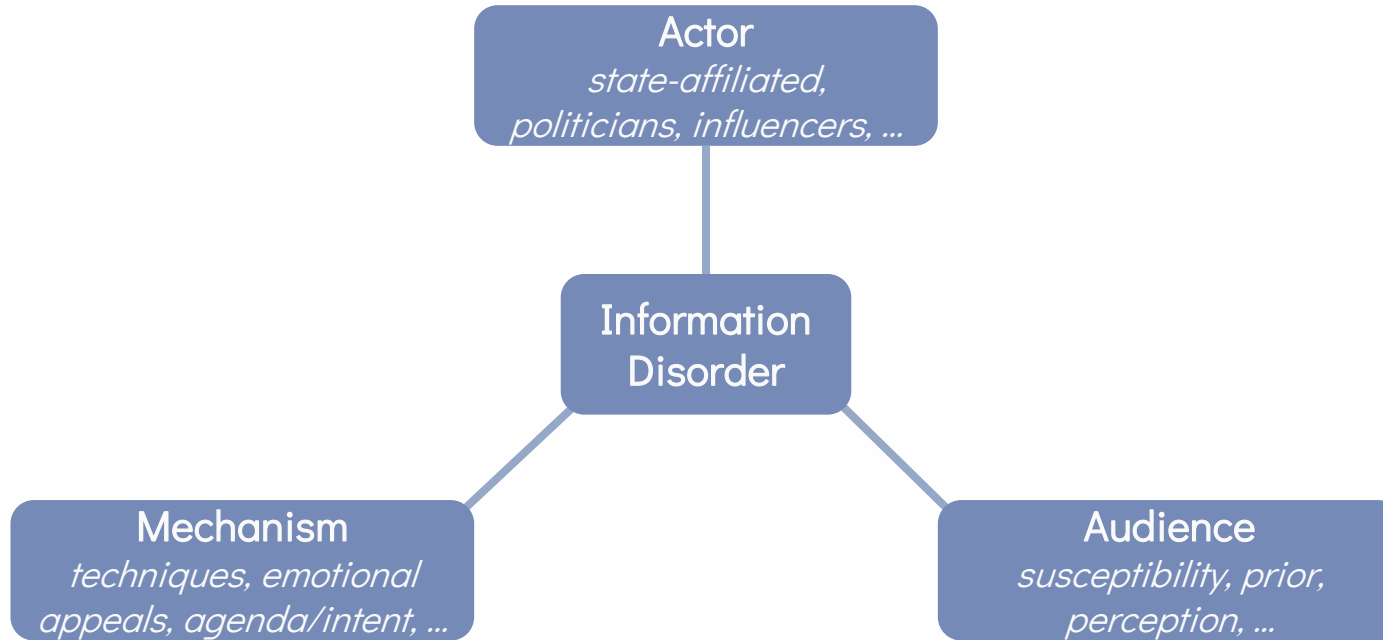
- 14K tweets auto-labeled with FT GPT-3.5-Turbo

Full dataset #Tweets	17,854	
	Human Annotated	Machine Annotated
#Tweets	3,691	14,163
#CTA Tweets	93	307
#CTA Text Spans	196	537
#DE Tweets	411	767
#DE Text Spans	1,292	2,447

Dataset Statistics




## Fine-Grained Understand



## Background

- Moving forward, we propose a **generalizable framework for analyzing propaganda and the specific intent** driving each propaganda attempt

- 
- Telling non-expert readers only what techniques are used is far from satisfying
  - More research is needed to understand the intent behind these techniques

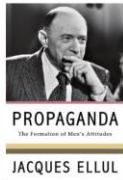
... A government spokesman said individuals whose presence “is not conducive to the public good” could be excluded by the home secretary. He added: “We condemn all those whose behaviours and views run counter to our shared values and will not stand for extremism in any form.” ...

Hmm... What exactly is the ‘black-and-white fallacy,’ and why is it used?”

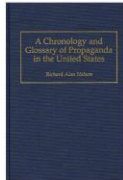


The highlighted text uses **black-and-white fallacy**.

# Established Research on Propaganda – Social Science



The expression of **opinions or actions carried out deliberately** by individuals or groups with a view to **influencing the opinions or actions** of other individuals or groups for **predetermined ends**.



A **systematic form of purposeful persuasion** that attempts to **influence the emotions, attitudes, opinions, and actions** of specified target audiences for **ideological, political or commercial purposes**.



A **deliberate, systematic attempt** to **shape perceptions, manipulate cognitions, and direct behaviors** to achieve a response that furthers the **desired intent**.

## Propaganda Techniques

The tools or instruments of persuasion or propaganda



## Arousal Appeals

Influences on the readers

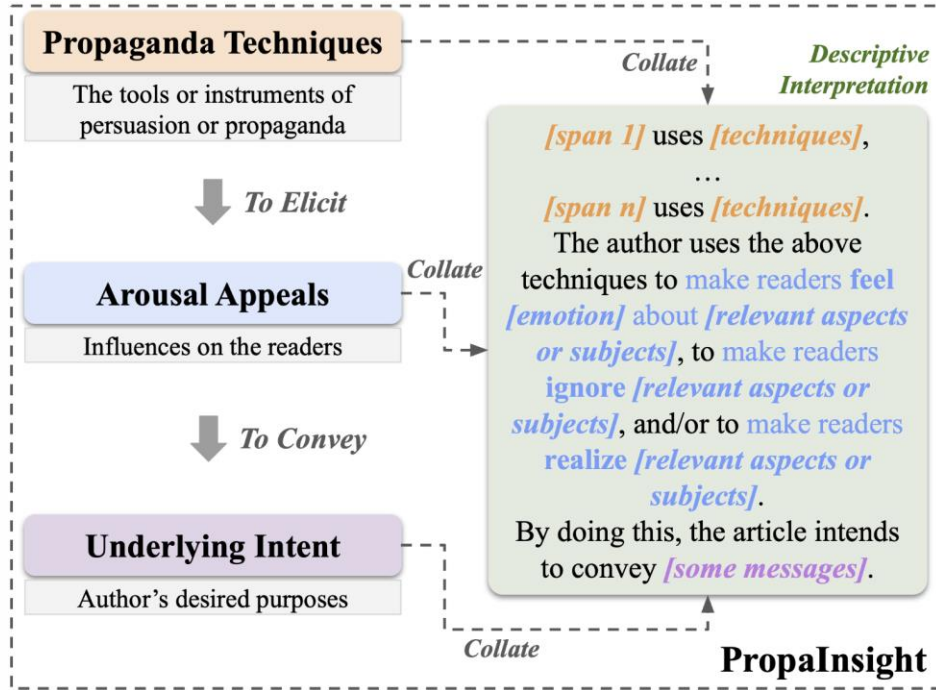


## Underlying Intent

Author's desired purposes

- Drawing from science research on propaganda, we identify three key elements of each propaganda attempt:
  - Propaganda Techniques (instrument)
  - Arousal Appeals (persuasion strategy)
  - Underlying Intent (purpose)
- We incorporate the elements into a *conceptual framework of propaganda analysis*

# Propalnsight: Framework of Propaganda Analysis



Step 1: Identify and classify techniques



Step 2: Infer evoked arousal appeals

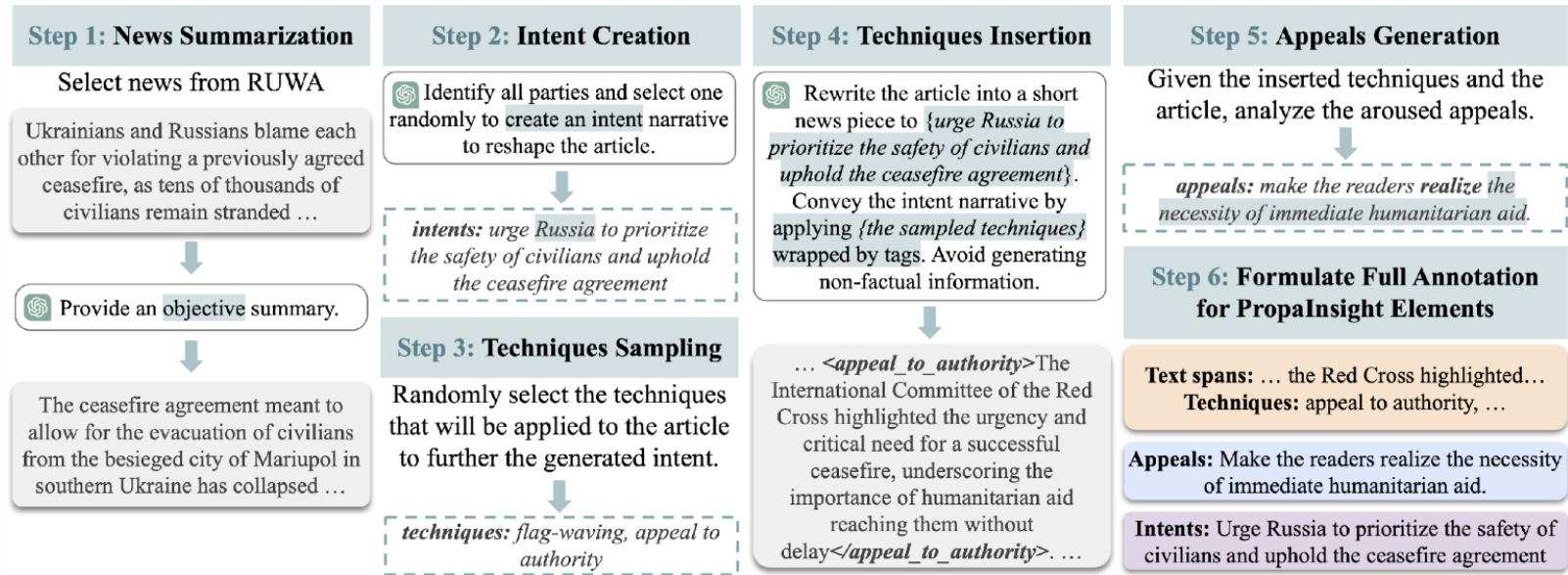


Step 3: Deduce underlying intent



Step 4: Consolidate elements into structured natural language paragraph

# PropaGaze: A Novel Dataset for Propaganda Analysis



- We construct *RUWA-Gaze* and *Politifact-Gaze* using a controlled data generation pipeline
- Rewrite articles into objective summaries, and then insert sampled propaganda techniques to shape the article's intent

# Off-the-Shelf LLMs Performance

- Zero-shot LLMs struggle with propaganda analysis
- Zero-shot LLMs struggle to pinpoint techniques and to infer arousal appeals
- Models perform better at inferring intent
- Few-shot prompting improves LLM performance in analyzing propaganda elements

Data Setting ↓	Dataset → Model ↓	RUWA-Gaze		Politifact-Gaze		PTC-Gaze	
		Span	Techniques	Span	Techniques	Span	Techniques
		Avg. IoU	Macro F1	Avg. IoU	Macro F1	Avg. IoU	Macro F1
<b>No Training Data</b>	GPT-4-Turbo <sub>0s</sub>	0.073	0.097	0.152	0.226	0.124	0.068
	GPT-4-Turbo <sub>1s</sub>	0.132	0.145	0.183	0.269	0.165	0.171
<b>Data-Sparse Training</b>	MGNN	0.089	0.139	0.160	0.159	0.140	<b>0.206</b>
	Llama-7B-Chat <sub>ft</sub>	0.230	0.210	0.253	0.281	<b>0.179</b>	<u>0.191</u>
<b>Data-Rich Training</b>	MGNN	<b>0.545</b>	<u>0.591</u>	<b>0.449</b>	<b>0.461</b>	-	-
	Llama-7B-Chat <sub>ft</sub>	<u>0.506</u>	<b>0.607</b>	<u>0.409</u>	<u>0.453</u>	-	-

Model Performance on *Propaganda Technique Identification* Subtask

Dataset →	Model ↓	RUWA-Gaze		Politifact-Gaze		PTC-Gaze	
		Appeals	Intents	Appeals	Intents	Appeals	Intents
		BertScore	BertScore	BertScore	BertScore	BertScore	BertScore
	GPT-4-Turbo <sub>0s</sub>	0.282	0.849	0.298	0.863	0.228	<u>0.869</u>
	GPT-4-Turbo <sub>1s</sub>	<u>0.324</u>	<b>0.879</b>	0.345	<b>0.875</b>	<b>0.331</b>	<b>0.881</b>
	Llama-7B-Chat <sub>ft</sub> ( <i>Data-Sparse</i> )	0.313	0.851	0.342	0.860	<u>0.249</u>	0.843
	Llama-7B-Chat <sub>ft</sub> ( <i>Data-Rich</i> )	<b>0.612</b>	<u>0.861</u>	<b>0.495</b>	<u>0.864</u>	-	-

Model Performance on *Appeal Analysis* and *Intent Analysis* Subtasks

# PropaGaze Effectiveness

- Models: **Llama-7B-Chat** and Multi-Granularity Neural Networks (**MGNN**)
- PropaGaze** substantially improves the overall propaganda analysis performance

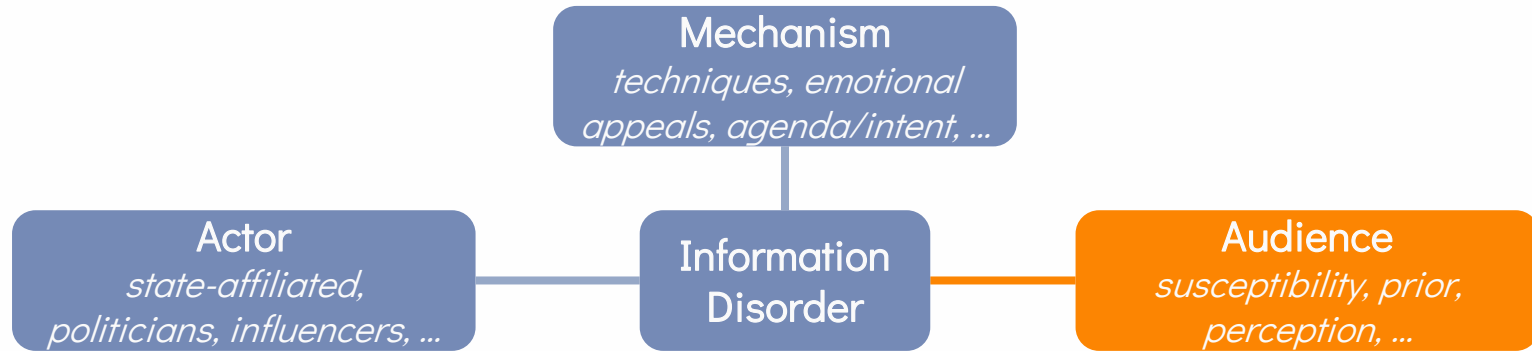
Data Setting ↓	Dataset →	<i>RUWA-Gaze</i>		<i>Politifact-Gaze</i>		<i>PTC-Gaze</i>	
	Model ↓	Span	Techniques	Span	Techniques	Span	Techniques
		Avg. IoU	Macro F1	Avg. IoU	Macro F1	Avg. IoU	Macro F1
<i>No Training Data</i>	GPT-4-Turbo <sub>0s</sub>	0.073	0.097	0.152	0.226	0.124	0.068
	GPT-4-Turbo <sub>1s</sub>	0.132	0.145	0.183	0.269	0.165	0.171
<i>Data-Sparse Training</i>	MGNN	0.089	0.139	0.160	0.159	0.140	<b>0.206</b>
	Llama-7B-Chat <sub>ft</sub>	0.230	0.210	0.253	0.281	<b>0.179</b>	<u>0.191</u>
<i>Data-Rich Training</i>	MGNN	<b>0.545</b>	<u>0.591</u>	<b>0.449</b>	<b>0.461</b>	-	-
	Llama-7B-Chat <sub>ft</sub>	<u>0.506</u>	<b>0.607</b>	<u>0.409</u>	<u>0.453</u>	-	-

Model Performance on *Propaganda Technique Identification* Subtask

Dataset →	<i>RUWA-Gaze</i>		<i>Politifact-Gaze</i>		<i>PTC-Gaze</i>	
Model ↓	Appeals	Intents	Appeals	Intents	Appeals	Intents
	BertScore	BertScore	BertScore	BertScore	BertScore	BertScore
GPT-4-Turbo <sub>0s</sub>	0.282	0.849	0.298	0.863	0.228	<u>0.869</u>
GPT-4-Turbo <sub>1s</sub>	<u>0.324</u>	<b>0.879</b>	0.345	<b>0.875</b>	<b>0.331</b>	<b>0.881</b>
Llama-7B-Chat <sub>ft</sub> ( <i>Data-Sparse</i> )	0.313	0.851	0.342	0.860	<u>0.249</u>	0.843
Llama-7B-Chat <sub>ft</sub> ( <i>Data-Rich</i> )	<b>0.612</b>	<u>0.861</u>	<b>0.495</b>	<u>0.864</u>	-	-

Model Performance on *Appeal Analysis* and *Intent Analysis* Subtasks

## Multimodal Cues in Radicalization: An Audience-Centric Study



1. **RQ1:** What viewer traits, such as personality traits and media consumption, are associated with their video preferences?
2. **RQ2:** What video characteristics, such as speaker traits, video quality, and arousing emotions, are correlated with viewers' perception?
3. **RQ3:** Which modality features affect viewers' perception the most?

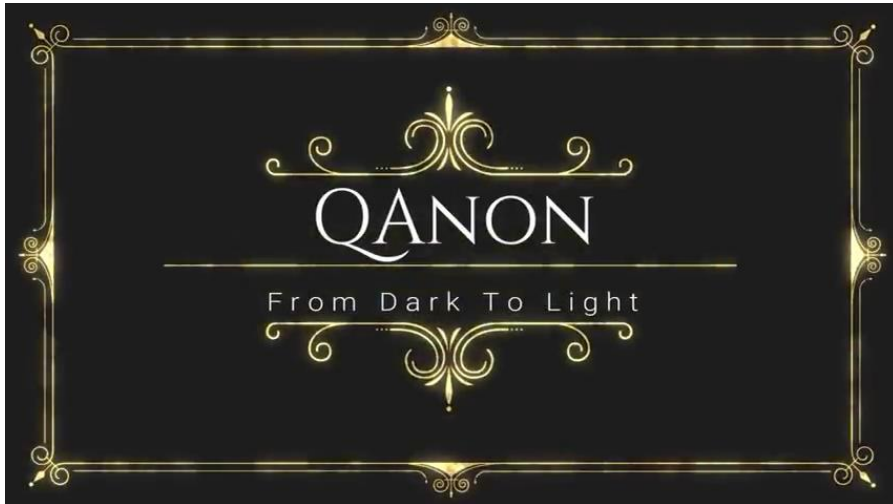
# Multimodal Cues in Radicalization: An Audience-Centric Study

**⚠ Content Disclaimer:** *examples shown contain harmful/misleading content and do **NOT** reflect the views of the presenter or co-authors.*

## Study at a glance



- 🎬 6 exemplar videos on QAnon
- 👤 46 participants (questionnaire)
- ☐ **Metrics:** Enjoyment, Content, Actions

- We focus on **QAnon**, a conspiracy-driven radicalization group
- We select 6 videos based on relevance, duration, diversity of styles, content quality, and popularity



# Multimodal Cues in Radicalization: An Audience-Centric Study

## Study at a glance

-  6 exemplar videos on QAnon
-  46 participants (questionnaire)
- Metrics:** Enjoyment, Content, Actions

1. Did you understand the video?  
 Yes  
 No

2. Do you think the video was professionally produced with good quality?  
 Yes  
 No

3. Who do you think the video was trying to appeal to?: \_\_\_\_\_

4. Was there any violence displayed in the video?  
 Yes  
 No

5. Was there any music in video?  
 Yes  
 No

6. Did any of the following objects appear in the video? Choose all that apply.  
 Guns  
 Swords  
 Other Weapons  
 Flags  
 Symbols of the Group  
 None of the Above

7. How likely do you think it is that the people in the video will become involved in the following actions?

	Not at All Likely	Not Much Likely	Undecided	Somewhat Likely	Very Much Likely
Protests	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Violence	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Illegal Acts	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>



## Enjoyment Score

- How much viewers enjoy watching each video

## Content Score

- Whether viewers think a video is persuasive, trustworthy, logical, and professionally created

## Actions Score

- posting a criticizing comment [score -2]
- disliking the video [score -1]
- liking the video [score 1]
- posting a supporting comment [score 2]
- considering joining the group [score 3]

# Multimodal Cues in Radicalization: An Audience-Centric Study



Enjoyment on Pro-QAnon Videos			Enjoyment on Anti-QAnon Videos		
Feature	Corr	p-value	Feature	Corr	p-value
Opinion_CNN	0.329	2.55E-02	Opinion_Antifa	0.368	1.19E-02
Opinion_WSJ	0.298	4.40E-02			

Content of Pro-QAnon Videos			Content of Anti-QAnon Videos		
Feature	Corr	p-value	Feature	Corr	p-value
Opinion_Fox	0.487	5.92E-04	Researved	0.339	2.13E-02
Opinion_NPR	-0.376	1.00E-02			
Opinion_AP	-0.33	2.53E-02			

Actions after Pro-QAnon Videos			Actions after Anti-QAnon Videos		
Feature	Corr	p-value	Feature	Corr	p-value
Opinion_OathKeepers	0.37	1.14E-02	Disorganized	0.318	3.12E-02
Opinion_Fox	0.358	1.45E-02	Sympathetic	-0.317	3.21E-02
Opinion_CNN	0.298	4.42E-02			

Significant viewer traits

## Viewer traits → Perception/Actions

- ↑ Media trust/outlet alignment → ↑Enjoyment/Actions
- ↑ Partisan priors → ↑ Actions toward aligned videos
- ↑ “Reserved” → ↑ Content agreement (anti-videos)

Takeaway: Susceptibility is not uniform across audiences

# Multimodal Cues in Radicalization: An Audience-Centric Study

- **Video characteristics:** arousal emotions, speaker characteristics of the videos, etc.
- **Textual features:** LIWC, Grievance Dictionary, VADER
- **Acoustic features** (*This is where the speech signal comes in 😊*)
  - Acoustic-prosodic features, such as pitch and intensity
  - SpeechBrain's emotion detection model
- **Visual features**
  - Pre-trained FER mode
  - Clarifai's weapon detector model

## Video cues → Perception/Actions

- ↑ Perceived trust/persuasion/valid points → ↑ Actions
- ↑ Perceived trust/persuasion/valid points → ↑ Actions
- ↑ Weapons / negative facial expressions → ↓ Perception

Takeaway: Presentation/style shifts perception

# Multimodal Cues in Radicalization: An Audience-Centric Study

## Study at a glance

- 📺 6 exemplar videos on QAnon
- 👤 46 participants (questionnaire)
- ☐ **Metrics:** Enjoyment, Content, Actions

## Viewer traits → Perception/Actions

- ↑ Media trust/outlet alignment → ↑ Enjoyment/Actions
- ↑ Partisan priors → ↑ Actions toward aligned videos
- ↑ “Reserved” → ↑ Content agreement (anti-videos)

**Takeaway:** Susceptibility is not uniform across audiences

## Video cues → Perception/Actions

- ↑ Perceived trust/persuasion/valid points → ↑ Actions
- ↑ Perceived trust/persuasion/valid points → ↑ Actions
- ↑ Weapons / negative facial expressions → ↓ Perception

**Takeaway:** Presentation/style shifts perception



## Why these matter

- **Targeted interventions:** tailor counterspeech
- **Evidence selection:** down-weight high-arousal/low-trust cues; up-weight credible clips
- **Persona-aware systems:** feed signals into steerable pipelines

## Conclusion: Toward a fuller picture of information disorder

### ACTOR

#### Who & why

- RTCS-HGT: untrustworthy tweet detection via text + propagation
- TweetIntent@Crisis: narrative-level CTA / DE intents in RU-UA conflict

### MECHANISM

#### How it persuades

- CLIP-MTL: visual cues uncover malicious intent
- PropaInsight: techniques + appeals + intent in propaganda

### AUDIENCE

#### Who it reaches

- QAnon study: viewer traits shape susceptibility
- Multimodal cues – including speech prosody & vocal emotion – shift perception

### What's next

Persona-aware, steerable interventions • Tighter integration of acoustic & visual signals for model reasoning • Real-time, scalable defense systems

Thank you – questions?