

Fake news detection:

Limited Ground Truth,

Limited Text,

No Understanding of Spreading Intent

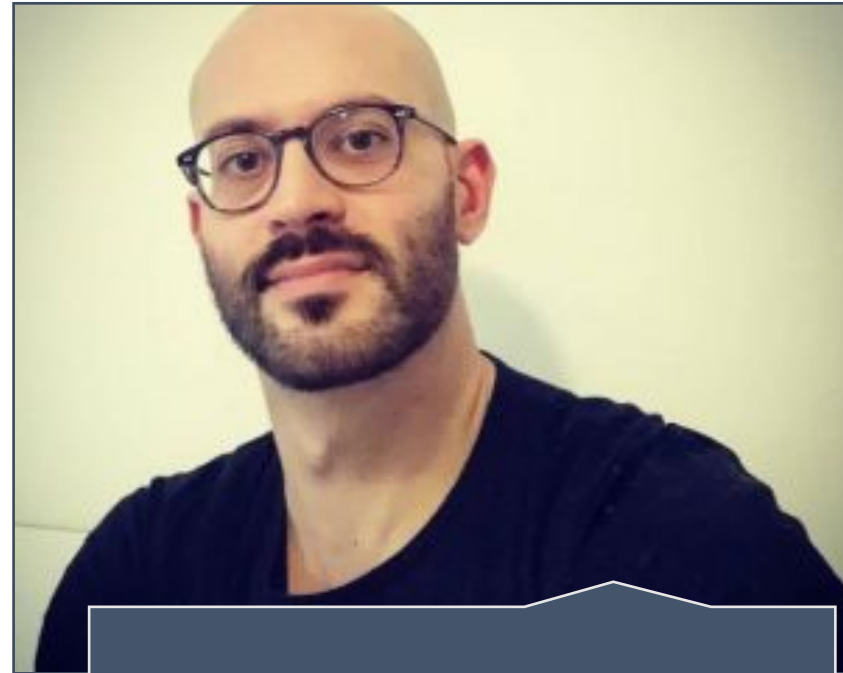
Reza Zafarani

ROMCIR 2022: the 2nd Workshop on Reducing Online
Misinformation Through Credible Information
Retrieval - April 10, 2022 | Stavanger, Norway

Thanks to the organizers



Marinella Petrocchi



Marco Viviani



What Is Fake News?

Fake News & Related Concepts

Definition of fake news

Fake news is **intentionally false** news published by a **news** outlet.

- **Intention** : Bad
- **Authenticity** : False
- **News** or not? News

A broader definition:

- *Fake news is false news*



Denzel Washington Backs Trump In The Most Epic Way Possible

While the rest of liberal Hollywood is still trying to demonize Donald Trump, Denzel Washington is speaking out in favor of the president-elect. "We need more and..."

AMERICANNEWS.COM

BREAKING: Obama And Hillary Now Promising Amnesty To Any Illegal That Votes Democrat

Posted by Alex Cooper | Nov 8, 2016 | Breaking News



All Breaking News Being Given Amnesty For Clinton Votes

Concept	Authenticity	Intention	News?
Deceptive news	Non-factual	Mislead	Yes
False news	Non-factual	Undefined	Yes
Satire news	Non-unified ²	Entertain	Yes
Disinformation	Non-factual	Mislead	Undefined
Misinformation	Non-factual	Undefined	Undefined
Cherry-picking	Commonly factual	Mislead	Undefined
Clickbait	Undefined	Mislead	Undefined
Rumor	Undefined	Undefined	Undefined

For example, **Disinformation** is **false information** [news or non-news] with a **bad intention** aiming to mislead the public.



Fake News & Related Concepts

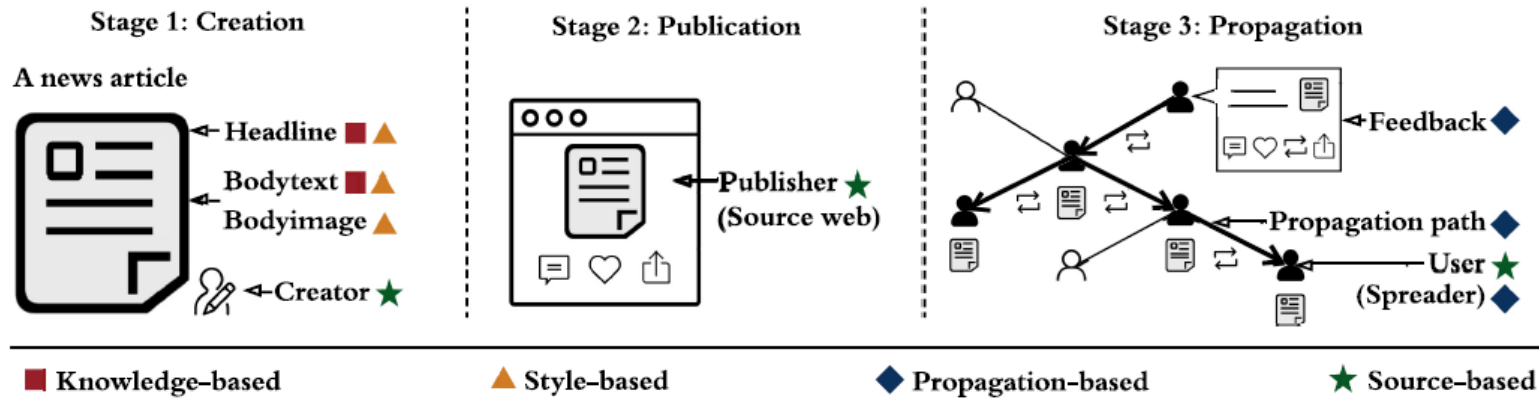
Distinguishing fake news from other related concepts

Open Problems:

- How similar are writing styles or propagation patterns?
- Can we use the same detection strategies?
- Can we distinguish between them? e.g., fake news from satire news

Fake News Detection

- **Knowledge-based** Fake News Detection
- **Style-based** Fake News Detection
- **Propagation-based** Fake News Detection
- **Source-based** Fake News Detection



Challenges and Highlights

- I. Limited Ground Truth**
- II. Limited Text**
- III. Unknown Intent of Fake News Spreaders**

I. Limited ground truth:

- *you can collect data:* ReCOVery dataset

Table 1: Data Statistics

	Reliable	Unreliable	Total
News articles	1,364	665	2,029
w/ images	1,354	663	2,017
w/ social information	1,219	528	1,747
Tweets	114,402	26,418	140,820
Users	78,659	17,323	93,761

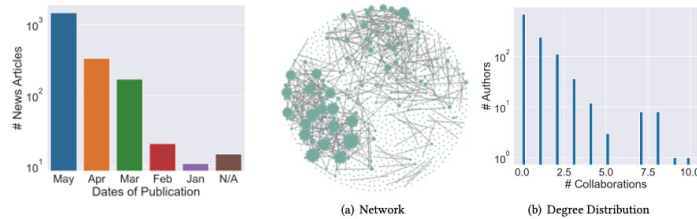


Figure 5: Publication Date

Figure 7: Author Collaborations

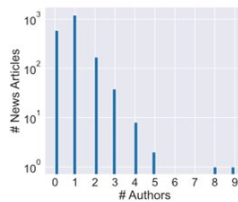


Figure 6: Author Count

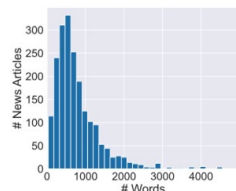


Figure 8: Word Count



Figure 9: Word Cloud

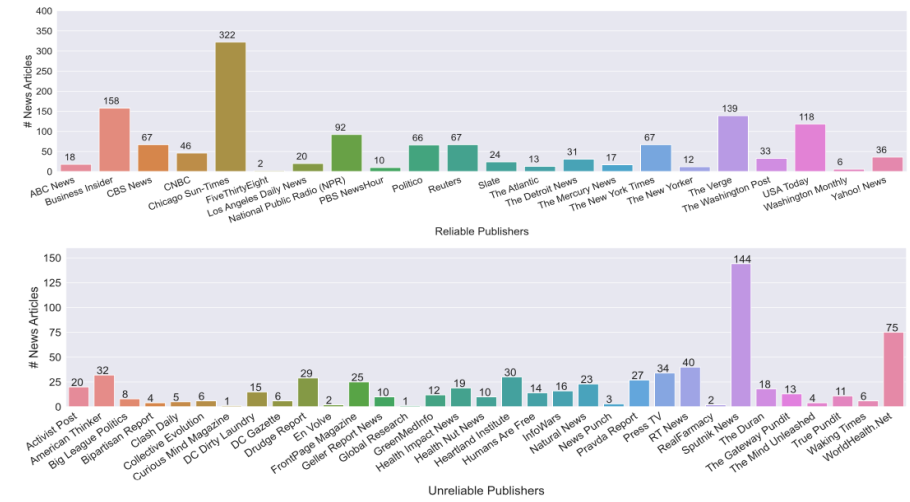
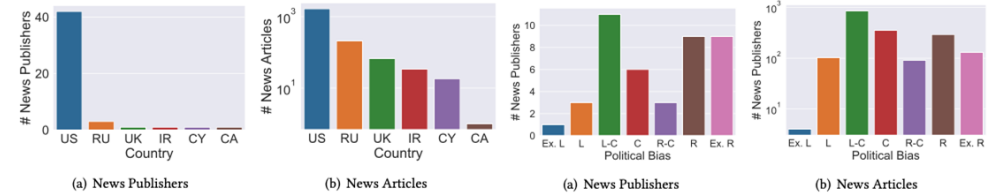


Figure 4: Distribution of News Publishers



(a) News Publishers

(b) News Articles

(a) News Publishers

(b) News Articles

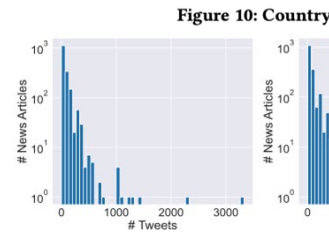


Figure 12: Spreading Frequency

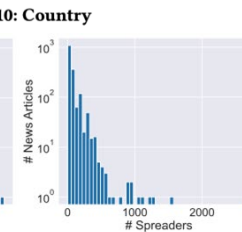


Figure 13: News Spreaders

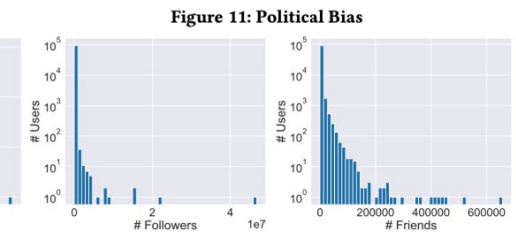


Figure 14: Follower Distribution

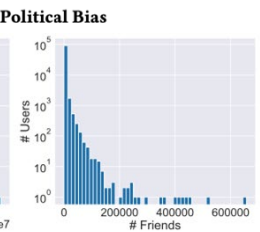


Figure 15: Friend Distribution

X. Zhou, A. Mulay, E. Ferrara, R. Zafarani

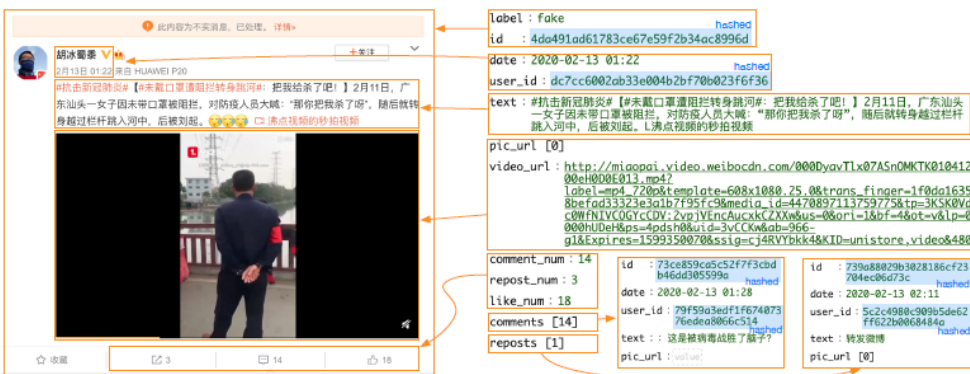
ReCOVery: A Multimodal Repository for COVID-19 News Credibility Research

I. and more data....

- CHECKED (Chinese COVID-19 Fake News Dataset) Dataset

Table 2 Statistics of CHECKED Data

	Real	Fake	All
# Microblogs	1,776	344	2,120
with images	1,153	53	1,206
with video	568	106	674
with reposts	1,167	229	1,396
with comments	1,167	292	1,459
# Reposts of microblogs	15,049	37,443	52,126
# Comments of microblogs	678,249	15,399	691,004
# Likes of microblogs	56,530,505	445,116	56,975,621
# Weibo users	690,755	51,674	737,347



(a) Fake Microblog

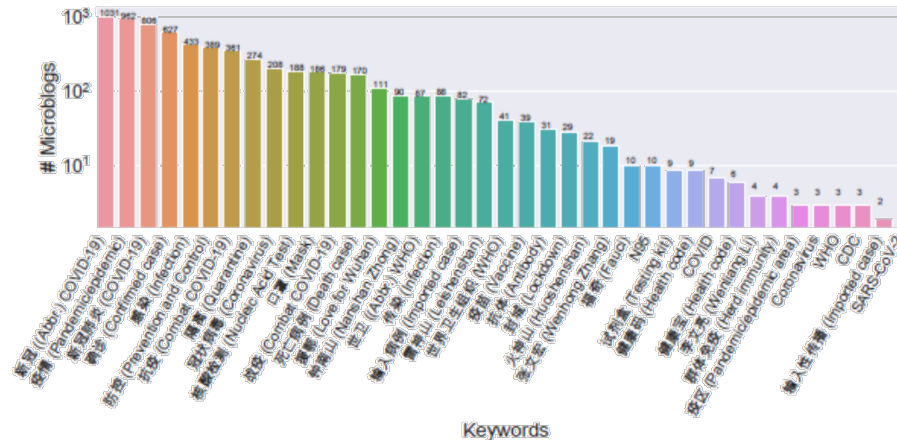


Fig. 2 Distribution of Selected Keywords in Collected Microblogs



Fig. 3 Word Cloud

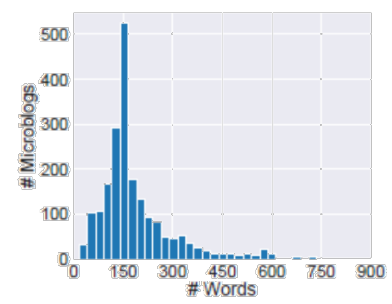


Fig. 4 Dist. of Words

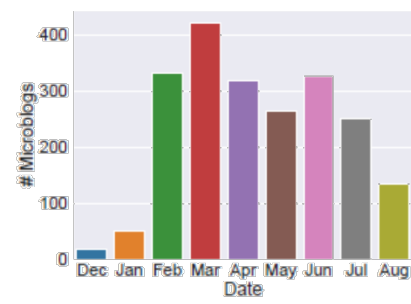


Fig. 5 Dist. of Dates Posted

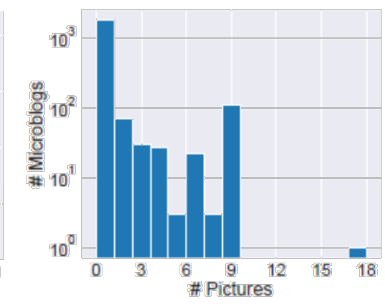


Fig. 6 Dist. of Images

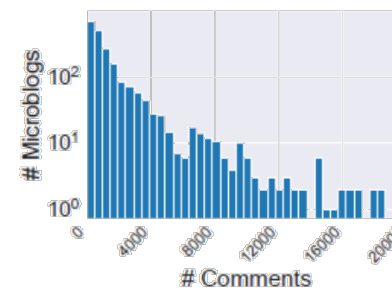


Fig. 7 Dist. of Comments

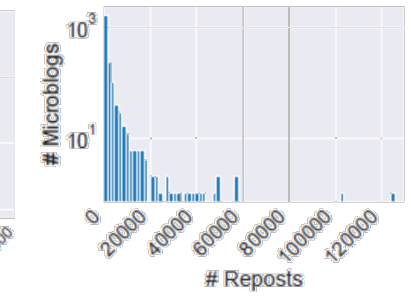


Fig. 8 Dist. of Reposts

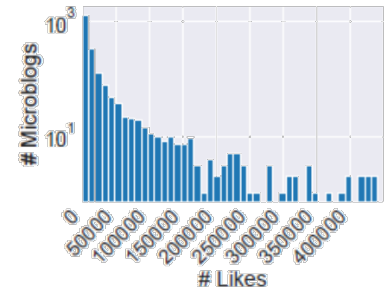


Fig. 9 Dist. of Likes

I. Or you can design methods that require limited data: **Fake News Early Detection**

Why is Fake News *Early Detection* important?

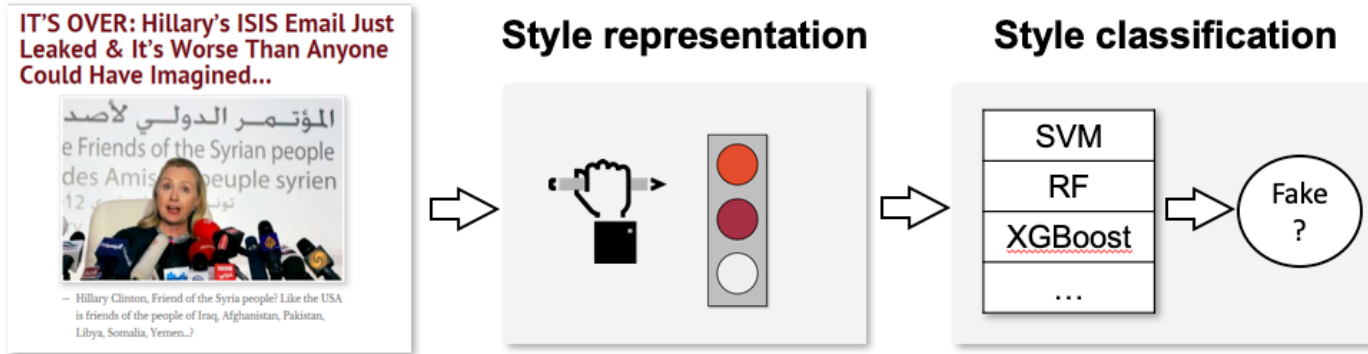
- The more fake news spreads, the more likely for people to trust it
- Once people have trusted the fake news, it can be difficult to correct users' perceptions

	Term	Phenomenon
Social influence	<i>Attentional bias</i>	Exposure frequency - individuals tend to believe information is correct after repeated exposures.
	<i>Validity effect</i>	
	<i>Echo chamber effect</i>	
	<i>Bandwagon effect</i>	Peer pressure - individuals do something primarily because others are doing it and to conform to be liked and accepted by others.
	<i>Normative influence theory</i>	
	<i>Social identity theory</i>	
<i>Availability cascade</i>		

Term	Phenomenon
<i>Backfire effect</i>	Given evidence against their beliefs, individuals can reject it even more strongly
<i>Conservatism bias</i>	The tendency to revise one's belief insufficiently when presented with new evidence.
<i>Semmelweis reflex</i>	Individuals tend to reject new evidence as it contradicts with established norms and beliefs.

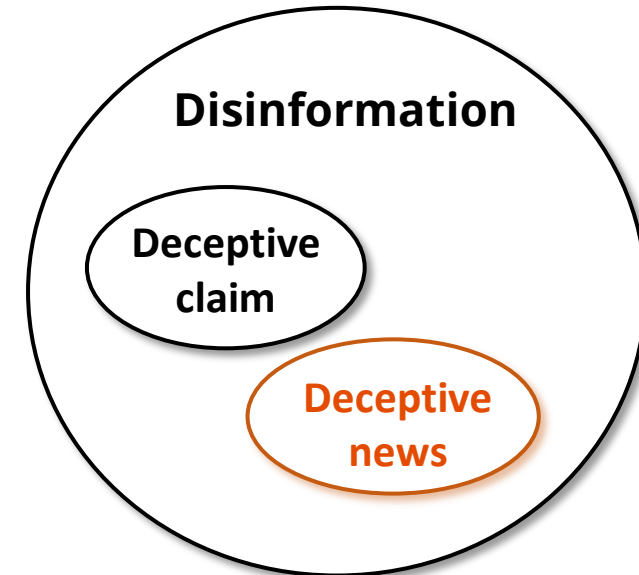
Fake News Early Detection: A Theory-driven Model

Xinyi Zhou, Atishay Jain, Vir V. Phoha, Reza Zafarani



- Interpretability
- Empirical relations

<i>Undeutsch hypothesis</i>	Deceptive statements differ in content style and quality from the truth.
<i>Reality monitoring</i>	Deceptive claims are characterized by higher levels of sensory-perceptual information.
<i>Four-factor theory</i>	Lies are expressed differently in emotion and cognitive process from the truth.
<i>Information Manipulation theory</i>	Extreme information quantity often exists in deception .

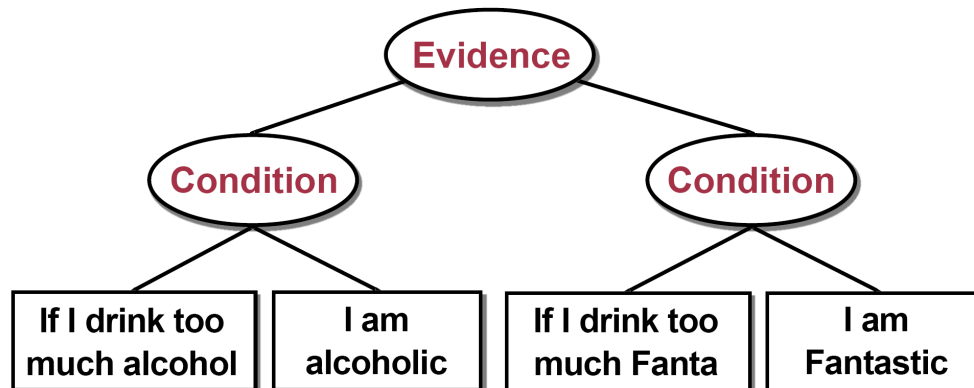
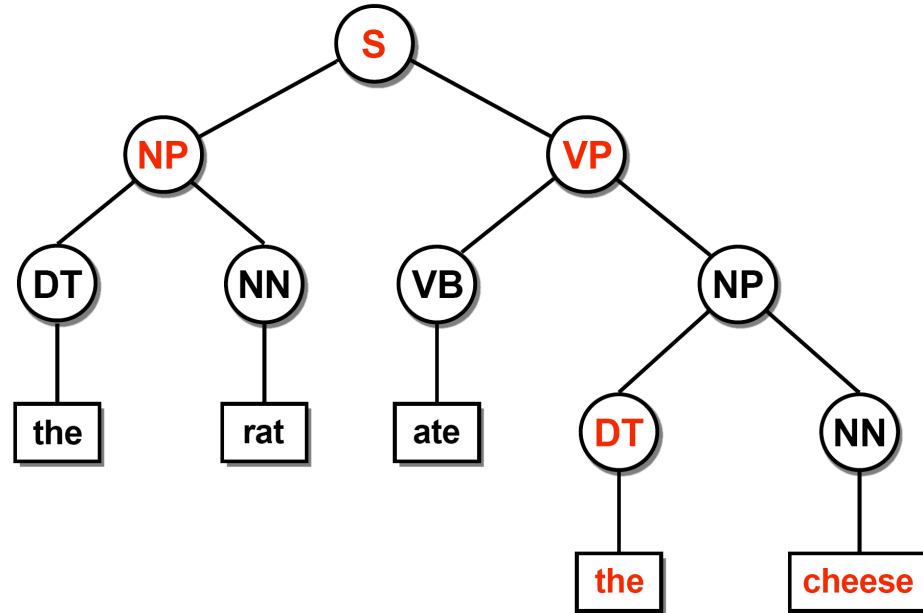


Fake News Early Detection: A Theory-driven Model

Xinyi Zhou, Atishay Jain, Vir V. Phoha, Reza Zafarani

I. Writing Style

Level	Feature(s)
Lexicon	BOWs
Syntax	POS Tags
	CFGs
Discourse	RRs



Lexicon	'rat'	1	x	x
	'cheese'	1	x	x
POS	noun	2	x	x
	verb	1	x	x
CFG	S → NP VP	1	x	x
	DT → 'the'	2	x	x
RR	Evidence	1	x	x
	Condition	2	x	x
		N ₁	N ₂	N ₃

Fake News Early Detection: A **Theory**-driven Model

Xinyi Zhou, Atishay Jain, Vir V. Phoha, Reza Zatarani

II. Content Quality

	Feature(s)	Example	Tool & Ref.
Informality	#/% Swear Words	"damn"	Linguistic Inquiry and Word Count (LIWC)
	#/% Netspeak	"btw"	
	#/% Assent	"OK"	
	#/% Nonfluencies	"umm"	
	#/% Fillers	"you know"	
	Overall #/% Informal Words	/	
Subjectivity	#/% Biased Lexicons	"attack"	[1]
	#/% Report Verbs	"announce"	[2]
	#/% Factive Verbs	"observe"	
Diversity	#/% Unique Words	/	/
	#/% Unique Content Words	"car"	LIWC
	#/% Unique Nouns	/	POS Taggers
	#/% Unique Verbs	/	
	#/% Unique Adjectives	/	
	#/% Unique Adverbs	/	



[1] Marta Recasens, et al. Linguistic Models for Analyzing and Detecting Biased Language. ACL, 2013.

[2] J Hooper. On Assertive Predicates in Syntax and Semantics, New York, 1975.

Fake News Early Detection: A **Theory**-driven Model

Xinyi Zhou, Atishay Jain, Vir V. Phoha, Reza Zafarani

III. Perceptual Process

#/% See	LIWC
#/% Hear	
#/% Feel	
Overall #/% Perceptual Processes	

IV. Sentiment

#/% Positive Words	LIWC
#/% Negative Words	
#/% Anxiety Words	
#/% Anger Words	
#/% Sadness Words	
Overall #/% Emotional Words	NLTK
Avg. Sentiment Score of Words	

V. Cognitive Process

#/% Insight	"think"	LIWC
#/% Causation	"because"	
#/% Discrepancy	"should"	
#/% Tentative	"perhaps"	
#/% Certainty	"always"	
#/% Differentiation	"but"	
Overall #/% Cognitive Processes		

VI. Quantity

Characters
Words
Sentences
Paragraphs
Avg. # Characters Per Word
Avg. # Words Per Sentence
Avg. # Sentences Per Paragraph

Fake News Early Detection: A **Theory**-driven Model

Xinyi Zhou, Atishay Jain, Vir V. Phoha, Reza Zafarani

Within/Across-level Performance

	Language Level	Feature Group	PolitiFact				BuzzFeed			
			XGBoost		RF		XGBoost		RF	
			Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
Within Levels	Lexicon	BOW	.856	.858	.837	.836	.823	.823	.815	.815
	Shallow Syntax	POS	.755	.755	.776	.776	.745	.745	.732	.732
	Deep Syntax	CFG	.877	.877	.836	.836	.778	.778	.845	.845
	Semantic	DIA+CBA	.745	.748	.737	.737	.722	.750	.789	.789
	Discourse	RR	.621	.621	.633	.633	.658	.658	.665	.665
Across Two Levels	Lexicon+Syntax	BOW+POS+CFG	.858	.860	.822	.822	.845	.845	.871	.871
	Lexicon+Semantic	BOW+DIA+CBA	.847	.820	.839	.839	.844	.847	.844	.844
	Lexicon+Discourse	BOW+RR	.877	.877	.880	.880	.872	.873	.841	.841
	Syntax+Semantic	POS+CFG+DIA+CBA	.879	.880	.827	.827	.817	.823	.844	.844
	Syntax+Discourse	POS+CFG+RR	.858	.858	.813	.813	.817	.823	.844	.844
	Semantic+Discourse	DIA+CBA+RR	.855	.857	.864	.864	.844	.841	.847	.847
Across Three Levels	All-Lexicon	All-BOW	.870	.870	.871	.871	.851	.844	.856	.856
	All-Syntax	All-POS-CFG	.834	.834	.822	.822	.844	.844	.822	.822
	All-Semantic	All-DIA-CBA	.868	.868	.852	.852	.848	.847	.866	.866
	All-Discourse	All-RR	.892	.892	.887	.887	.879	.879	.868	.868
	Overall		.865	.865	.845	.845	.855	.856	.854	.854

Within-level

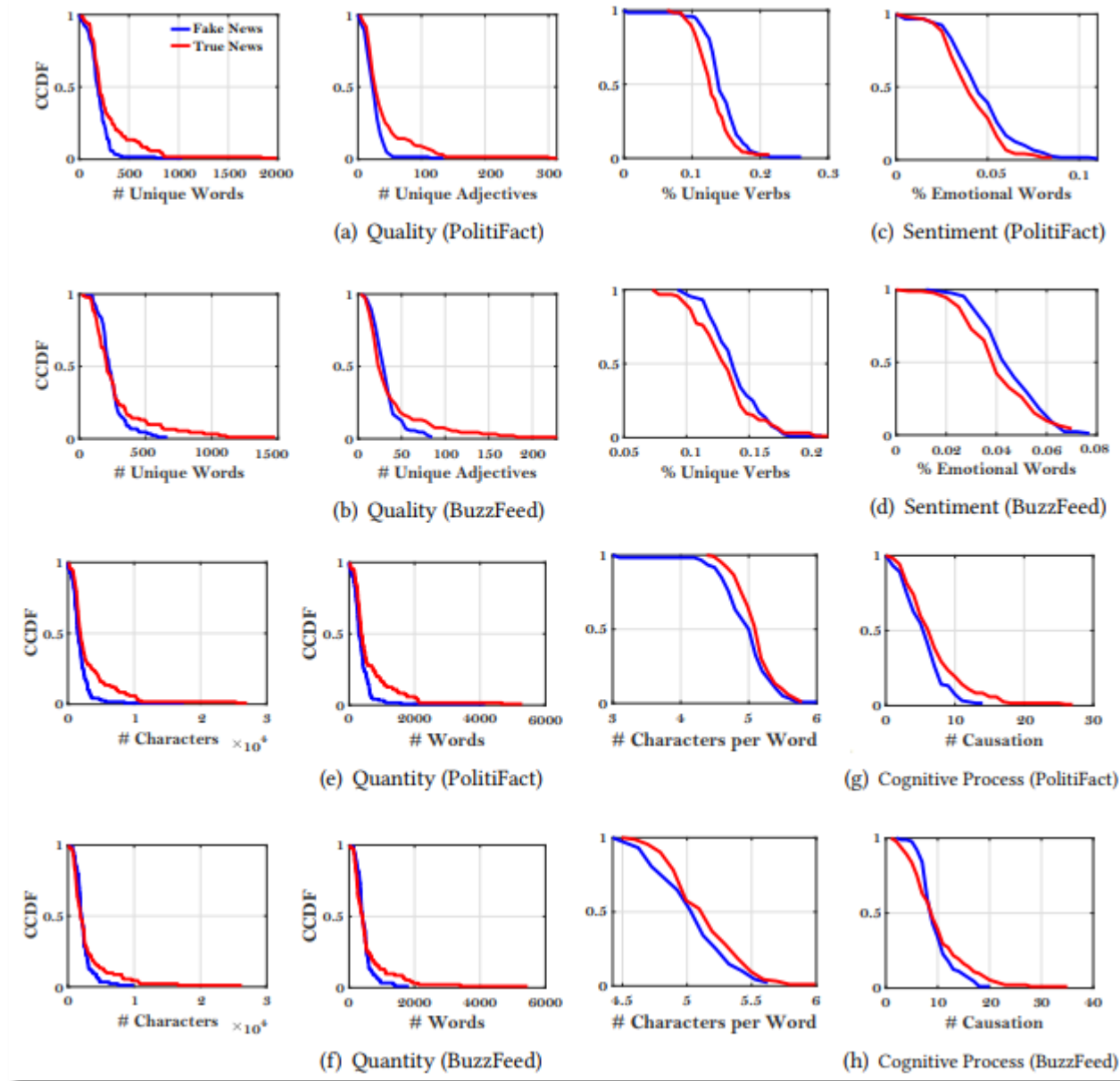
- 1. Lexicon / Deep Syntax**
(80%~90%)
- 2. Semantic / Shallow Syntax**
(70%~80%)
- 3. Discourse**
(60%~70%)

Across-level > Within-level
(exclude RRs)

Fake News & Deception

Supportive Theory	Deception	Fake News
<i>Undeutsch hypothesis</i>	Differs in content style and quality from truth	😊 Consistent
<i>Reality monitoring</i>	Has a higher levels of sensory-perceptual information than truth	😬 Similar levels to the truth
<i>Four-factor theory</i>	Differs in cognitive process from the truth	😊 Carries less cognitive information than truth
<i>Information Manipulation theory</i>	Often refers to extreme information quantity	🤖 More words in headlines while less in body-text.

$p\text{-value} < 0.1$



II. Limited Text

- Pursue Multi-Modal Fake News Detection

- Few existing studies have explored **the relationship (similarity) between news text and images** to help detect fake news.

Washington State Legislature votes to change its name because George Washington owned Slaves

The legislature of Washington State has met in special session and overwhelming voted to change the name of the State. Since George Washington owned Slaves, it is improper for this State to be named after him. Due to the great support provided to the cause of eliminating the history of slavery in the United States by George Soros, the Legislature of Washington has chosen the new name of Soros State. The change in name will take effect on November 1st, 2017 once the Governor of Washington signs the bill



Motivation

Why is such similarity worth exploring?

- Fake news writers **actively** use attractive but **irrelevant** textual and visual information to form a false story
 - To attract the public attention
- Sometimes it is **passive** behavior
 - Cannot find related and non-manipulated images to support false claims

Angelina Jolie & Jared Leto Dating After

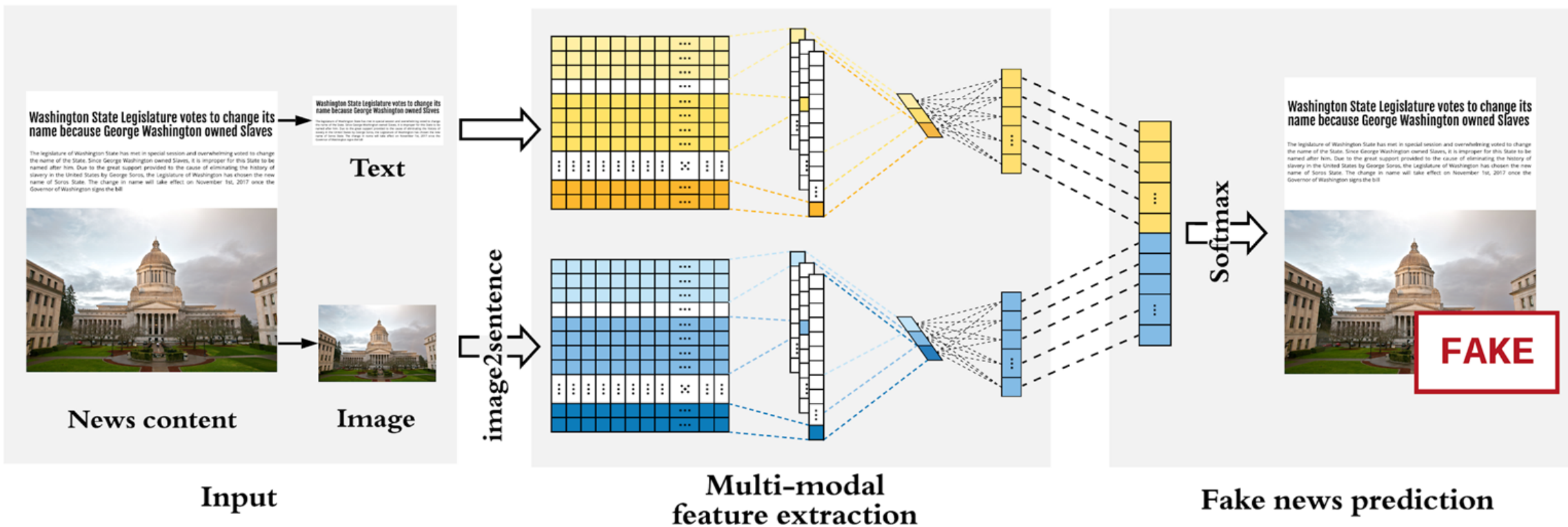
Brad Pitt Divorce — Report



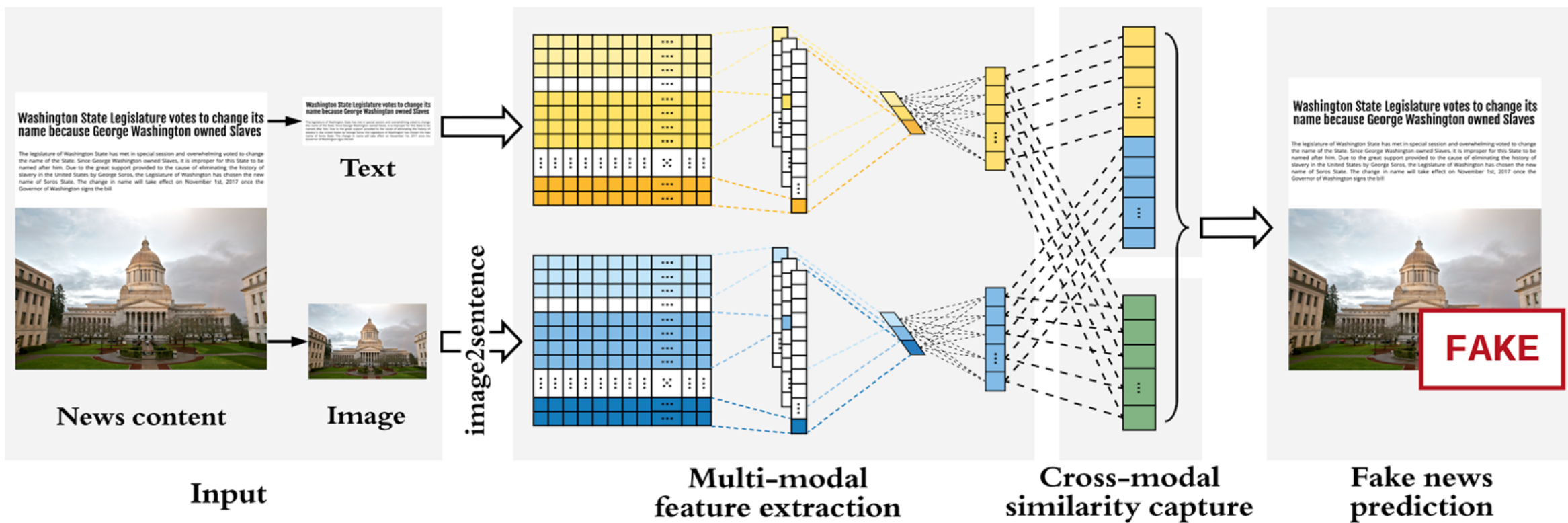
Chrissy Teigen and John Legend Have First Date Night Since Welcoming Son Miles: Pic!



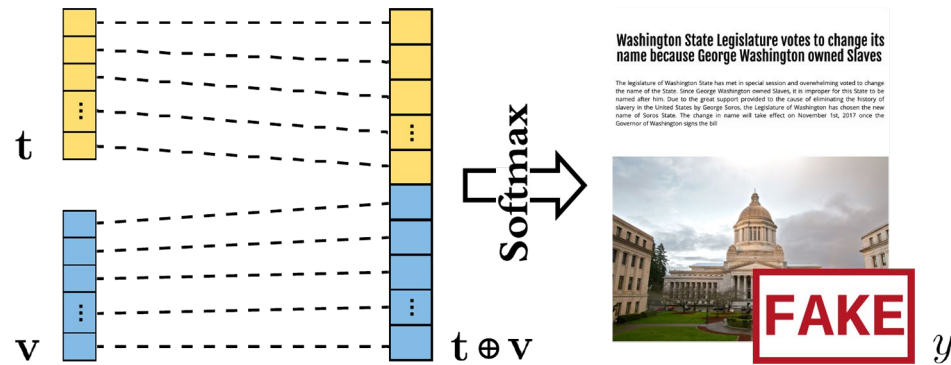
SAFE: Predicting Process



SAFE: Learning Process



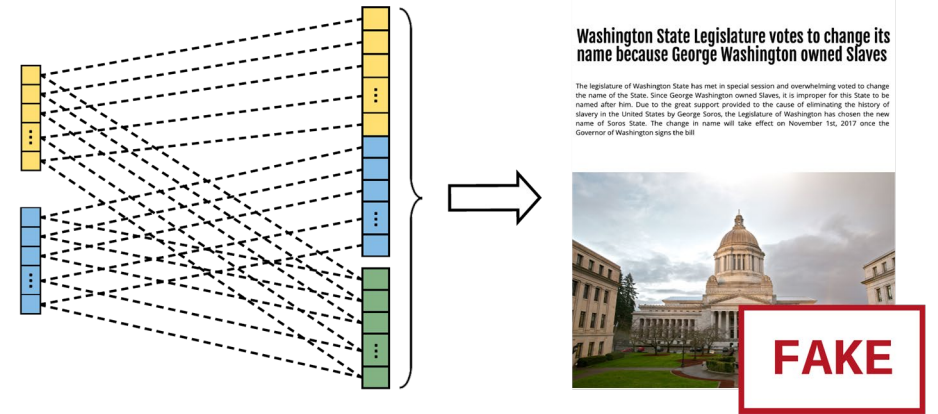
SAFE: Feature Representativeness/Joint Learning



$$\mathcal{M}_p(\mathbf{t}, \mathbf{v}) = \mathbf{1} \cdot \text{softmax}(\mathbf{W}_p(\mathbf{t} \oplus \mathbf{v}) + \mathbf{b}_p),$$

$$\mathcal{L}_p(\theta_t, \theta_v, \theta_p) = -\mathbb{E}_{(a,y) \sim (A,Y)} (y \log \mathcal{M}_p(\mathbf{t}, \mathbf{v}) + (1 - y) \log(1 - \mathcal{M}_p(\mathbf{t}, \mathbf{v}))),$$

$$(\hat{\theta}_t, \hat{\theta}_v, \hat{\theta}_p) = \arg \min_{\theta_t, \theta_v, \theta_p} \mathcal{L}_p(\theta_t, \theta_v, \theta_p).$$



$$\mathcal{L}(\theta_t, \theta_v, \theta_p) = \alpha \mathcal{L}_p(\theta_t, \theta_v, \theta_p) + \beta \mathcal{L}_s(\theta_t, \theta_v),$$

$$(\hat{\theta}_t, \hat{\theta}_v, \hat{\theta}_p) = \arg \min_{\theta_t, \theta_v, \theta_p} \mathcal{L}(\theta_t, \theta_v, \theta_p).$$

Experiments: General Performance

Result on multiple modalities:

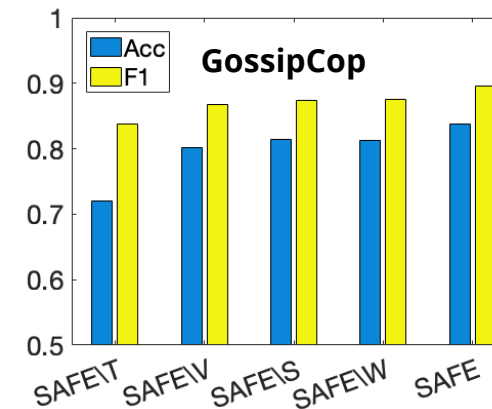
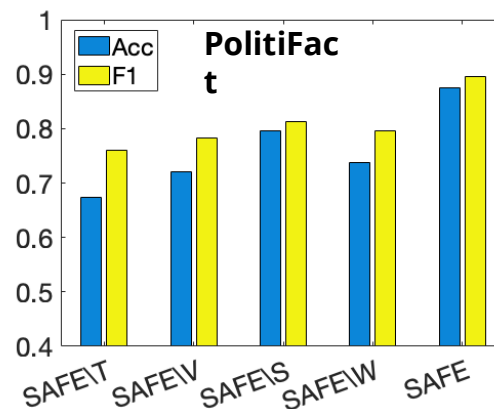
- Textual + Visual + Relational > Textual + Visual information
 - SAFE vs att-RNN, SAFE\S, SAFE\W
- Textual + Visual \approx Relational information
 - SAFE\S vs SAFE\W

Multi-modal > Single-modal methods

- Multi-modal > Single-modal information
 - SAFE, SAFE\S, SAFE\W, att-RNN vs LIWC, VGG-19, SAFET, SAFEW

Among single-modal methods

- Textual > Visual infor.
 - LIWC vs VGG-19
 - SAFEW vs SAFET



		LIWC [†]	VGG-19 [‡]	att-RNN [‡]	SAFE [‡]
Politi-Fact	Acc.	0.822	0.649	0.769	0.874
	F ₁	0.815	0.720	0.826	0.896
Gossip-Cop	Acc.	0.836	0.775	0.743	0.838
	F ₁	0.466	0.862	0.846	0.895

†: Text-based ‡: Image-based ‡: Multi-modal

Experiments: Case Studies

Examples of **true** news articles:

"Face the Nation" transcripts, August 26, 2012: Rubio, Priebus, Barbour, Blackburn



(a) $s = 0.966$

98 Degrees' 2017 Macy's Parade Performance Will Take You Right Back To The '90s



(b) $s = 0.975$

Chrissy Teigen and John Legend Have First Date Night Since Welcoming Son Miles: Pic!



(c) $s = 0.983$

Examples of **fake** news articles:

Washington State Legislature votes to change its name because George Washington owned Slaves



(a) $s = 0.024$

MORGUE EMPLOYEE CREMATED BY MISTAKE WHILE TAKING A NAP

Beaumont, Texas | An employee of the Jefferson County morgue died this morning after being accidentally cremated by one of his coworkers.



(b) $s = 0.044$

Angelina Jolie & Jared Leto Dating After Brad Pitt Divorce — Report



(c) $s = 0.001$

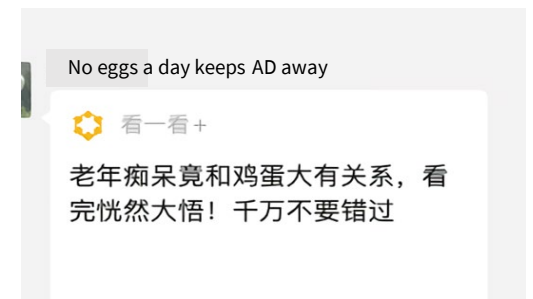
III. Assessing Intent of Fake News Spreaders

A frequently observed Phenomenon:

Individuals can spread fake news unintentionally without recognizing its falsehood

Our goal is to address some research questions:

1. **Why** does an individual unintentionally spread fake news?
2. How can we **model and assess** the intent of fake news spreaders?
3. Where can we obtain the ground-truth **data to evaluate** such models?
 - If no such data is available, how can one collect it from scratch?
4. How does modeling the intention of news spreaders help **fake news detection and mitigation**?



Xinyi Zhou, Kai Shu, Vir V. Phoha, Huan Liu, Reza Zafarani, "This is Fake! Shared it by Mistake": Assessing the Intent of Fake News Spreaders, TheWeb Conference 2022

Why? Psychological Interpretations for Unintentional Fake News Spreading

- **External Influence:** a user trusting/spreading a frequently-posted idea due to
 - *Peer pressure*, conforming to the behavior of others for being accepted by the community (*social identity theory* [1]).
 - *Social Exposure*, where more exposure increases one's perceived accuracy of fake news and leads to unintentional spreading (e.g., due to *validity effect* [2])
- **Internal Influence:** a user would trust and spread a fake story that matches his or her preexisting knowledge
 - Individuals tend to believe fake news articles that confirm their preexisting values and beliefs [3])

[1] Michael A Hogg. 2020. Social identity theory. Stanford University Press

[2] Gordon Pennycook, Tyrone D Cannon, and David G Rand. 2018. Prior exposure increases perceived accuracy of fake news. *Journal of experimental psychology: general* 147, 12 (2018), 1865.

[3] Sendhil Mullainathan and Andrei Shleifer. 2005. The market for news. *American Economic Review* 95, 4 (2005), 1031–1053.

Modeling Intention of Fake News Spreaders

- Fake news spreading is more unintentional if the posting behavior is affected more
 - **Externally** (by the similar behavior of other users) and/or
 - **Internally** (by the user's similar past behavior)
- **Rough Idea:** Constructing an **influence graph** of posts to capture pairwise influence among posts, where a (directed) edge between two posts indicates the (external or internal) influence flow from one post to the other.

Intention Modeling of Fake News Spreaders on Social Media

Consider a pair of posts p_i and p_j ...

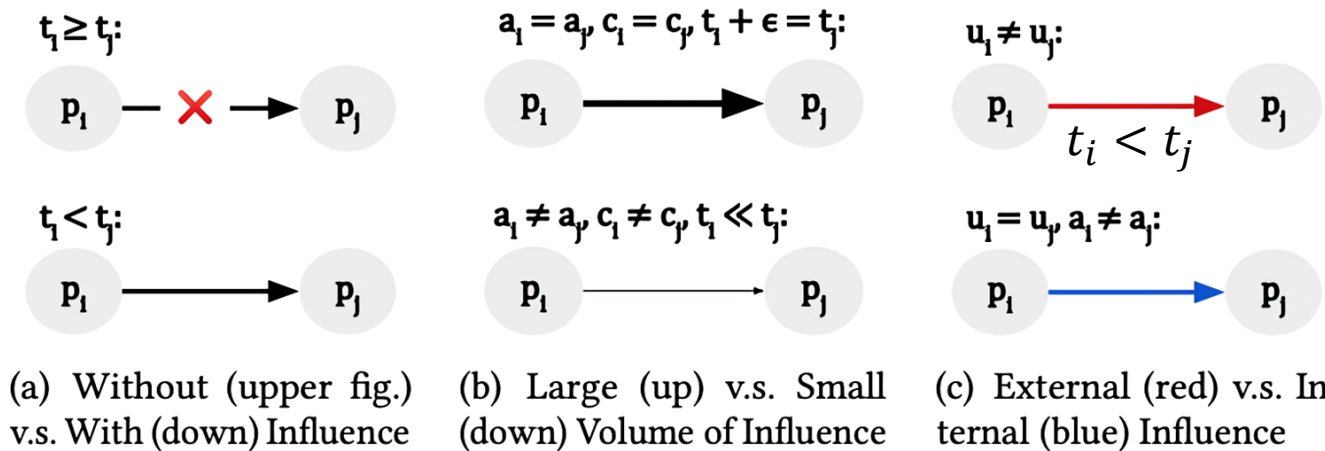


Figure 3: Pairwise Influence of Posts p_i and p_j : (a) decides if there is an edge from p_i to p_j in influence graph; (b) determines the edge weight; and (c) identifies the edge attribute.

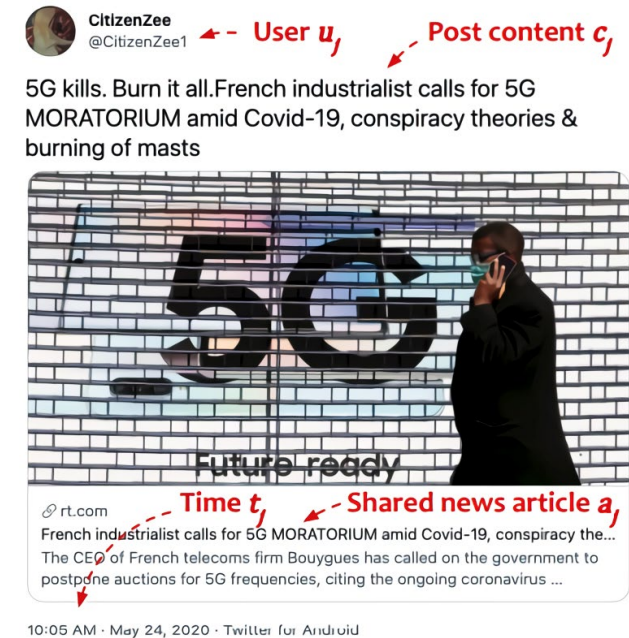


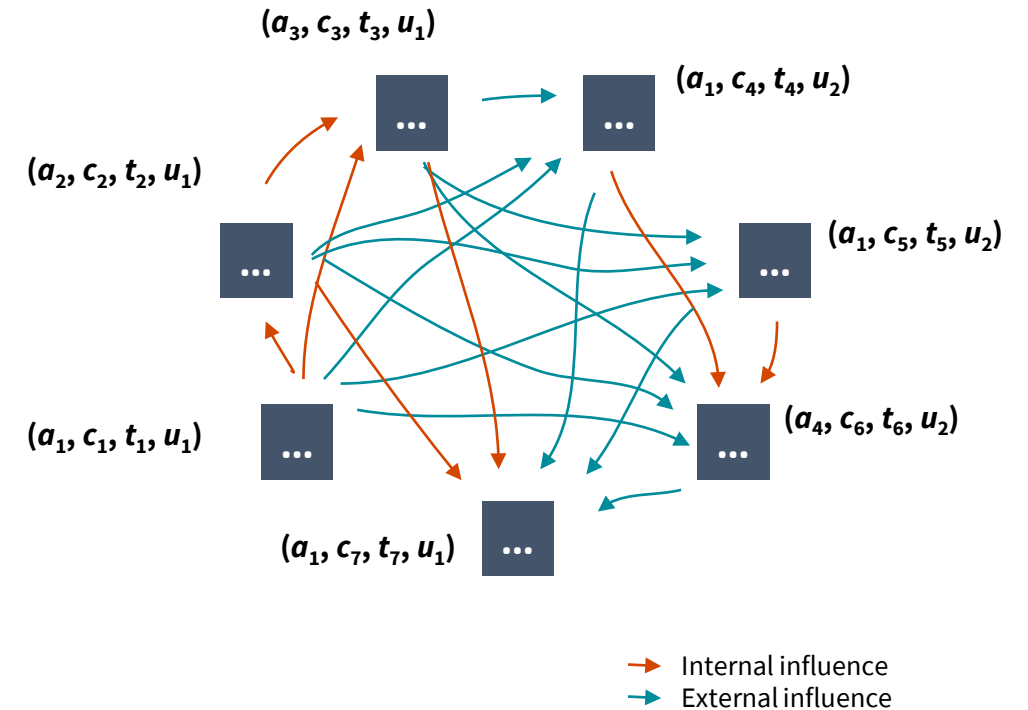
Figure 2: An Illustration of a Post $p_j = (a_j, c_j, t_j, u_j)$

Intention Modeling of Fake News Spreaders on Social Media

Influence Graph $G = (V, E, W)$

- $V = \{p_1, p_2, \dots, p_n\}$
- $(p_i, p_j) \in E \iff t_i < t_j \text{ and } a_i \neq a_j$
- $W_{ij} = S(a_i, a_j) \times S(p_i, p_j) \times T(t_i, t_j)$
 - $S(.,.)$: Similarity function (mostly by designing deep learning models for image/text)
 - $T(.)$: A self-defined monotonically decreasing decay function to capture users' forgetting

Using derived weights, we can compute the overall influence on each post (denoted as **affected degree**)



Method Evaluation: Data & Annotation

Evaluation data is required that contains the ground-truth label on

- News credibility, i.e., whether a news article is **fake news or the truth**; and
- User intention, i.e., whether a user spreads a fake news article **intentionally or unintentionally**.

Such datasets do not exist!

- **Our strategy:** Extend current datasets by annotating intention of fake news spreaders.

How?

Data for Method Evaluation

Request labels: news credibility + spreader intent

Manual annotation

- 2 well-trained annotators
- 300 posts randomly sampled
 - Intent: *intentional / unintentional*
 - Confidence: 0 / 0.5 / 1
 - Justification & Time
- Cohen's kappa: 0.61 (substantial)
- 119 posts: agree on intent with conf. ≥ 0.5
 - Small-scale, gold-standard, balanced



Time-consuming: 5 min per post
4-5 months in total if annotating
24/7

		MM- COVID	Re- COVery
# News	Fake	355	535
	True	448	1,231
# Tweets	Sharing Fake News	16,500	26,657
	Sharing True News	20,905	117,087

- Intentional spreaders: Bots + trolls + correctors
- Unintentional spreaders: Others

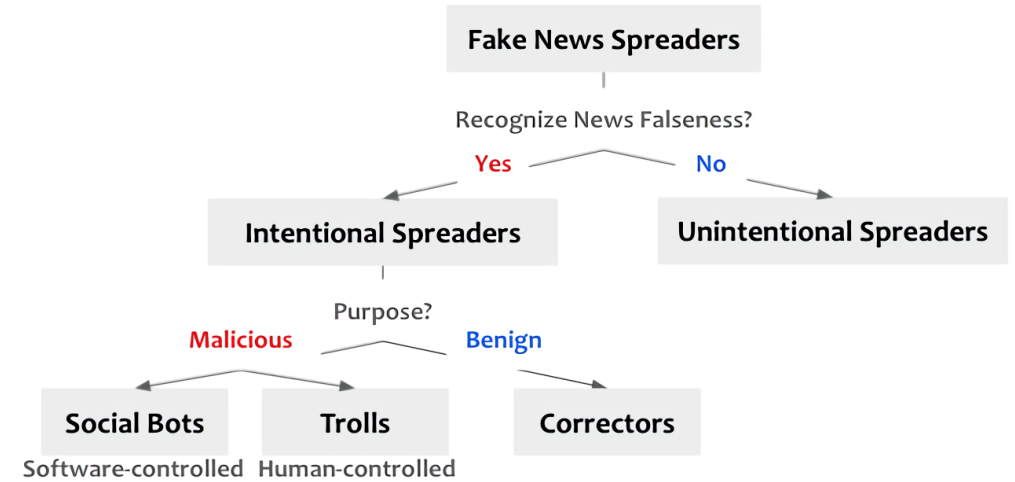
The task boils down to identifying bots, trolls, correctors and corresponding correction tweets....

Algorithm to simulate manual annotation


- Intentional: bots, trolls, correctors
 - Bots & trolls: often suspended, cannot be educated
 - Correctors: no need to be educated
- Unintentional: others

Table 1: Performance of Algorithmic Annotations on Intent of Fake News Spreaders

	AUC Score	Cohen's κ
MM-COVID + ReCOVery	0.8824	0.7482
MM-COVID	0.8857	0.7520
ReCOVery	0.8000	0.6484




Corrector (verifier) →



@CoreenaSuarez
@CoreenaSuarez2

Correction (verification) to a fake claim

Coronavirus before reaching the lungs remains in the throat for four days and at this time the person begins to cough and have throat pains. Drinking a lot of water, gargling with warm water mixed with salt or vinegar eliminates the virus = False



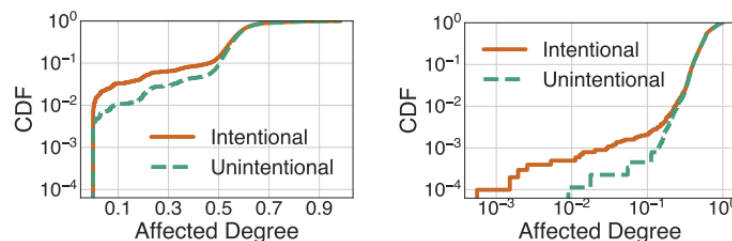
newsmeter.in
Fact Check: Can gargling with warm salt water prevent Coronavirus?

Method Evaluation

Our goal: unintentional fake news spreaders have significantly greater affected degrees than intentional ones

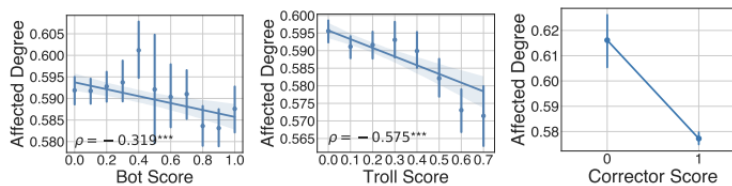
Results: unintentional fake news spreaders have greater affected degrees than intentional ones (bots, trolls, or correctors)

- Manual & algorithmic annotation
- *Statistically significant*
- Results are *stable* even when changing the hyperparameters in the annotation algorithm

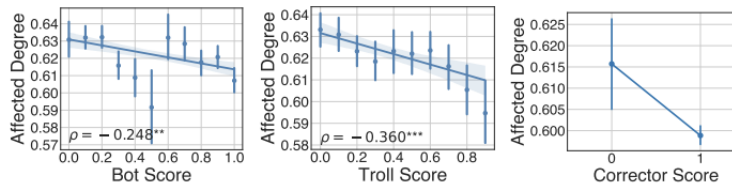


(a) MM-COVID ($p \ll 0.001$ with t-test) (b) ReCOVeRY ($p < 0.01$ with t-test)

Figure 4: Distribution of Affected Degree: Intentional Fake News Spreaders v.s. Unintentional Fake News Spreaders

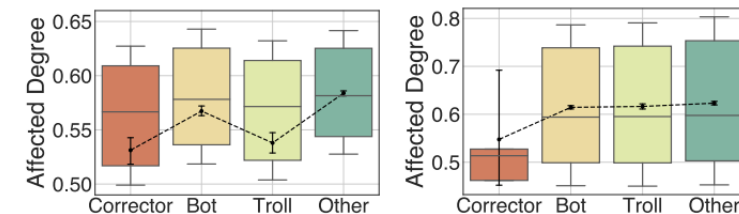


(a) MM-COVID ($p \ll 0.001$ using t-test for the right)



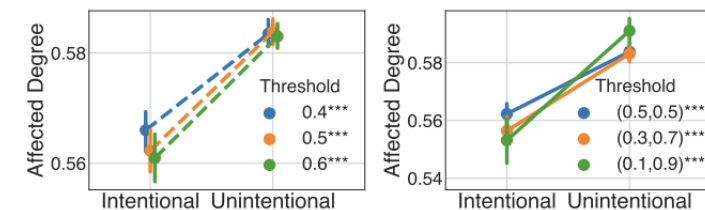
(b) ReCOVeRY ($p \ll 0.001$ using t-test for the right)

Figure 6: Relation between Affected Degree and (L) Bot Score, (M) Troll Score, and (R) Corrector Score. ρ : Spearman's Correlation Coefficient. *: $p < 0.001$; **: $p < 0.01$; and *: $p < 0.05$.**

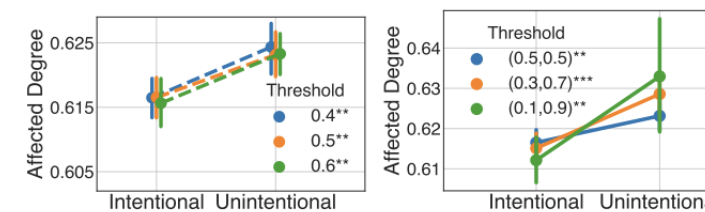


(a) MM-COVID ($p \ll 0.001$ by ANOVA) (b) ReCOVeRY ($p < 0.01$ by ANOVA)

Figure 5: Affected Degree of Bots, Trolls, Correctors, and Others (First Three: Intentional Fake News Spreaders; Others: Unintentional Fake News Spreaders)



(a) MM-COVID



(b) ReCOVeRY

Figure 7: Method Performance with Various Thresholds (*: $p < 0.001$; **: $p < 0.01$; and *: $p < 0.05$)**

I. Affected degree + traditional machine learning

Features: affected degree + content + propagation patterns (109 features)

Classifier: XGBoost

Table 3: Method Performance with Hand-crafted Features in Fake News Detection. Here, K : the first (earliest) K posts spreading the news available for news representation; Ranking: feature importance ranking of affected degree of posts in the prediction model.

	K	AUC Score	Ranking
MM-COVID	10	0.918 (± 0.009)	2
	20	0.912 (± 0.015)	2
	30	0.927 (± 0.021)	2
	40	0.923 (± 0.012)	2
	All	0.935 (± 0.005)	3
ReCOVery	10	0.891 (± 0.007)	5
	20	0.898 (± 0.007)	3
	30	0.903 (± 0.004)	3
	40	0.909 (± 0.014)	4
	All	0.925 (± 0.009)	5

Intent + Fake News Detection

II. Influence graph + deep learning

Features learned by HetGNN¹ & classified by XGBoost

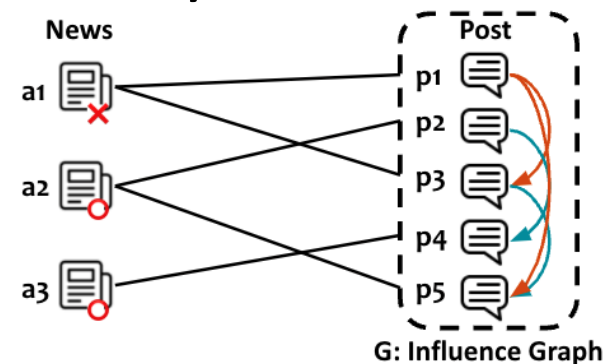


Table 2: Method Performance (Using AUC Scores) with Heterogeneous Graph Neural Networks (HetGNN) in Fake News Detection

	MM-COVID				ReCOVery			
% Labeled News	20%	40%	60%	80%	20%	40%	60%	80%
G_{RANDOM}	0.829	0.856	0.876	0.902	0.647	0.654	0.660	0.674
$G_{SUBGRAPH}$	0.817	0.861	0.890	0.915	0.820	0.845	0.869	0.908
G	0.869	0.864	0.902	0.905	0.825	0.863	0.883	0.881

Method's Prospects in Fake News Mitigation

Personalized intervention: Developing diverse strategies for fake news spreaders with various intentions to effectively and reasonably intervene with the spread of fake news on social media. For example,

- Removing and blocking bots and trolls, as intentional and malicious spreaders;
- Educating and correcting unintentional fake news spreaders.

Can there be a new recommendation algorithm that not only recommend interesting topics but also correction posts?

How effective are such algorithm in intervening with the spread of fake news?



- **Zhou, X.**, Shu, K., Phoha, V. V., Liu, H., & Zafarani, R. (2022). "This is Fake! Shared it by Mistake": Assessing the Intent of Fake News Spreaders. *arXiv preprint arXiv:2202.04752*.
- **Zhou, Xinyi**, and Reza Zafarani. "Fake news: A survey of research, detection methods, and opportunities." *arXiv preprint arXiv:1812.00315* (2018).
- **Zhou, Xinyi**, Jindi Wu, and Reza Zafarani. "SAFE: Similarity-aware multi-modal fake news detection." *arXiv preprint arXiv:2003.04981* (2020).
- **Zhou, Xinyi**, et al. "Fake News Early Detection: A Theory-driven Model." *arXiv* (2019): arXiv-1904.
- **Zhou, Xinyi**, and Reza Zafarani. "Network-based Fake News Detection: A Pattern-driven Approach." *arXiv preprint arXiv:1906.04210* (2019).
- **Zhou, Xinyi**, et al. "ReCOVery: A Multimodal Repository for COVID-19 News Credibility Research." *arXiv preprint arXiv:2006.05557* (2020).
- **Zhou, Xinyi**, et al. "Fake news: Fundamental theories, detection strategies and challenges." *Proceedings of the twelfth ACM international conference on web search and data mining*. 2019.
- Yang, Chen, et al. "CHECKED: Chinese COVID-19 Fake News Dataset." *arXiv preprint arXiv:2010.09029* (2020).

WEBSITES

<https://xinyizhou.xyz/papers/xzhou-kdd19-slides.pdf>



Xinyi Zhou (<https://xinyizhou.xyz/>)
zhouxinyi@data.syr.edu



Apurva Mulay



C. Mohan



Vir Phoha



Atishay Jain



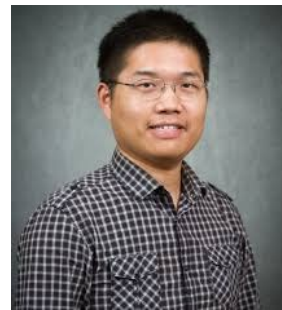
Jindi Wu



Chen Yang



Niraj Sitaula



Kai Shu



Emilio Ferrara



Huan Liu



Jennifer Grygiel