# Fake News: Fundamental Theories, Detection Strategies and Challenges

Xinyi Zhou, Reza Zafarani, Kai Shu, Huan Liu.

Syracuse University

ASU Arizona State University

# Meet our Team
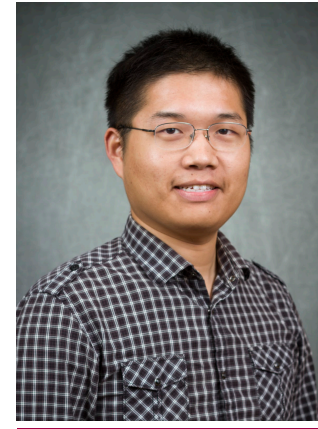
### Xinyi Zhou

Syracuse University

Ph.D. Student
Data Lab,
EECS Department

### Reza Zafarani

Syracuse University

Assistant Professor
Data Lab,
EECS Department

### Kai Shu

Arizona State University

Ph.D. Student
Computer science
and Engineering

### Huan Liu

Arizona State University

Professor
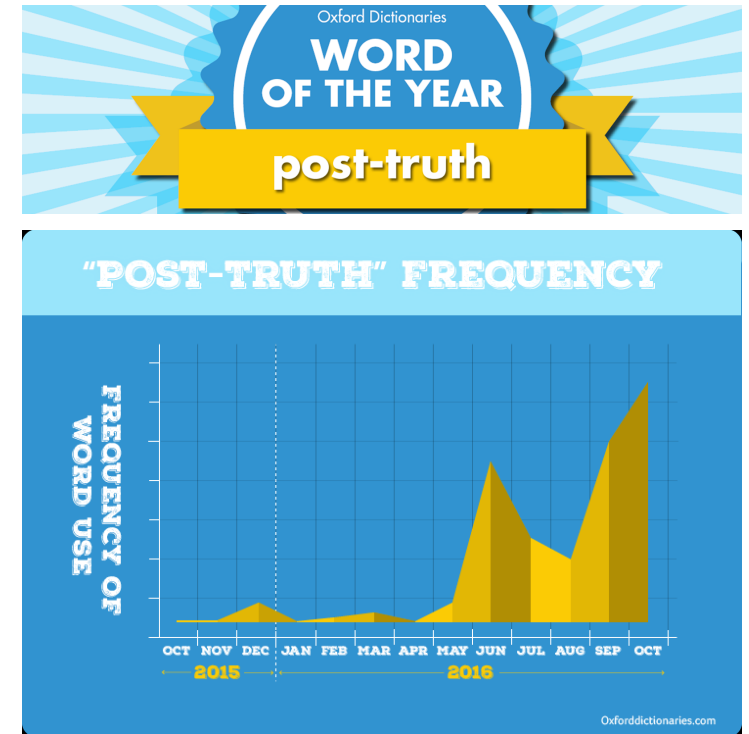Computer Science
and Engineering

# Introduction

- Research Background
- What is Fake News?
- Related Concepts
- Fundamental Theories

# Research Background

*Why Study Fake News?*

Fake news is now viewed as one of the greatest threats to **democracy**, **justice**, **public trust**, **freedom of expression**, **journalism** and **economy**.

- **Political** Aspects: May have had an impact on
  - "Brexit" referendum
  - 2016 U.S. presidential election
    - # Shares, reactions, and comments on Facebook.[1]
    - 8,711,000 for top 20 frequently-discussed **FAKE** election stories.
    - 7,367,000 for top 20 frequently-discussed **TRUE** election stories.

- Oxford Dictionaries international word of the year 2016:
  - **Post-Truth**: "Relating to or denoting circumstances in which objective facts are less influential in shaping public opinion than appeals to emotion and personal belief."
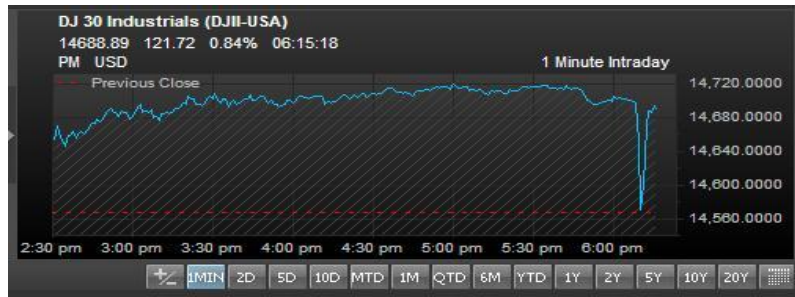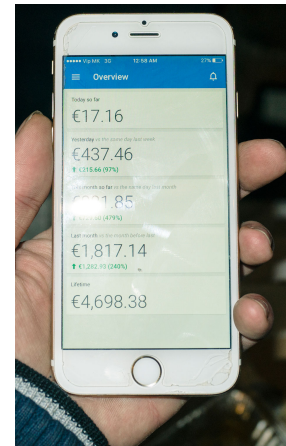


[1]C. Silverman. This analysis shows how viral fake election news stories outperformed real news on Facebook. BuzzFeed News, 2016.

X. Zhou, R. Zafarani, K. Shu, H. Liu

4

# Research Background

*Why Study Fake News?*

- **Economic** Aspects:
  - "Barack Obama was injured in an explosion" wiped out <u>$130 billion</u> in stock value.[1]
  - Dozens of "well-known" teenagers in Veles, Macedonia[2]
    - Penny-per-click advertising
    - During U.S. 2016 presidential Elections
    - Earning at least $60,000 in six months
    - Far outstripping their parents' income
    - Average annual wage in town: $4,800



[1] K. Rapoza. Can 'fake news' impact the stock market? 2017.
[2] S. Subramanian, Inside the Macedonnian Fake News Complex https://www.wired.com/2017/02/veles-macedonia-fake-news/

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Research Background

*Why Study Fake News?*

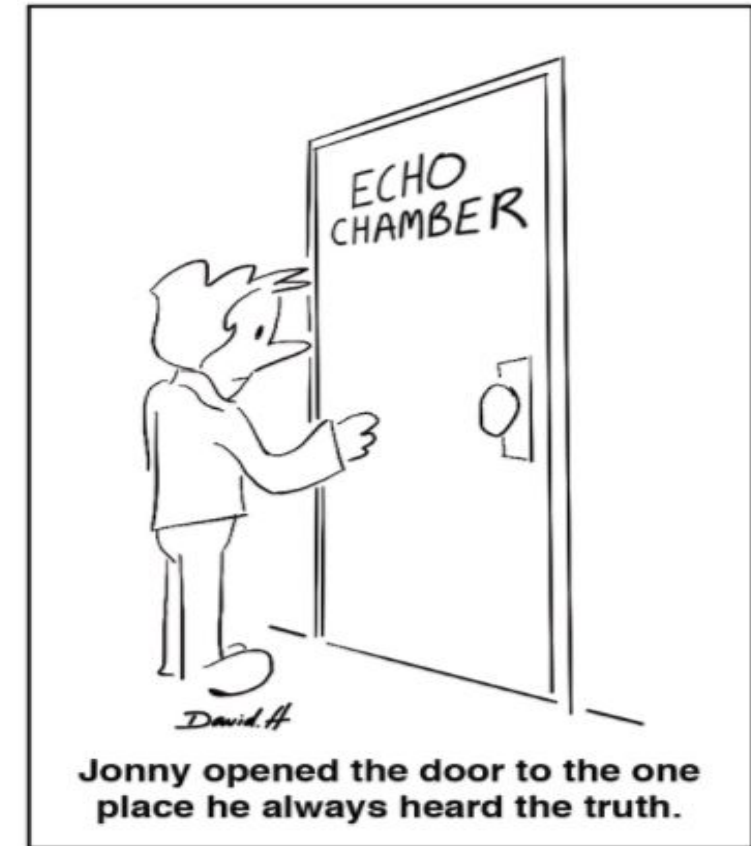- **Social/Psychological** Aspects:
  - Humans have been proven to be irrational/vulnerable when differentiating between truth/false news
    - Typical accuracy in the range of 55-58%
  - For fake news, it is relatively easier to obtain public trust
    - **Validity Effect:** individuals tend to trust fake news after repeated exposures
    - **Confirmation Bias:** individuals tend to believe fake news when it confirms their pre-existing knowledge
    - **Peer Pressure/Bandwagon Effect**



X. Zhou, R. Zafarani, K. Shu, H. Liu

# Research Background

*Why is Fake News attracting more public attention recently?*

- Fake news can now be created and published faster and cheaper

- The rise of Social Media and its popularity also plays an important role
  - As of Aug. 2017, 67% of Americans *get* their news from social media.[3]

- Social media accelerates fake news *dissemination*.
  - It breaks the physical distance barrier among individuals.
  - It provides rich platforms to share, forward, vote, and review to encourage users to participate and discuss online news.

- Social media accelerates fake news *evolution*.
  - **Echo chamber effect**: biased information can be amplified and reinforced within the social media.[4]
  - **Echo Chamber:** a situation in which beliefs are amplified or reinforced by communication and repetition inside a closed system



Jonny opened the door to the one place he always heard the truth.

---

[3]http://www.journalism.org/2017/09/07/news-use-across-social-media-platforms-2017/
[4]K. Jamieson and J. Cappella. Echo Chamber: Rush Limbaugh and the Conservative Media Establishment. Oxford University Press, 2008.

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News & Related Concepts

*Definition* of fake news

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

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

A more broad definition:

- *Fake news is false news*



FAKE NEWS DAILY
Misleading headline
Alternative facts



Pope Francis Shocks World, Endorses Donald Trump for President, Releases Statement

TOPICS: Pope Francis Endorses Donald Trump



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

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

X. Zhou, R. Zafarani, K. Shu, H. Liu

| | Authenticity | Intention | News? |
|---|---|---|---|
| **Fake news** | False | Bad | Yes |
| **False news** | False | Unknown | Yes |
| **Satire news** | Unknown | Not bad | Yes |
| **Disinformation** | False | Bad | Unknown |
| **Misinformation** | False | Unknown | Unknown |
| **Rumor** | Unknown | Unknown | Unknown |

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

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fundamental Theories

*Why is it necessary to study Fundamental Theories?*

**Fundamental human cognition and behavior theories** developed <u>across various discipline</u> such as psychology, philosophy, social science, and economics provide invaluable insights for fake news studies.

1. Provide opportunities for **qualitative and quantitative studies** of <u>big fake news data</u>;

2. Support to build **well-justified and explainable models** for fake news detection and intervention; and

...elop dat...                    ...d tru...                    ...ch

*[Udo] Undeutsch hypothesis:*
A **statement** based on a factual experience differs in **content and quality** from that of fantasy.

<u>Verification:</u>
Is a **fake news** article differs in **content and quality** from the truth?

<u>Utilizing:</u>
How to **detect fake news** based on its **content style and quality**?

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Style-Based Fundamental Theories

*Studying fake news from a style perspective, i..e, how it's written*

| | Term | Phenomenon |
|---|---|---|
| **Style-based** | *Undeutsch hypothesis* | A statement based on a factual experience differs in **content and quality** from that of fantasy |
| | *Reality monitoring* | Actual events are characterized by higher levels of **sensory-perceptual** information. |
| | *Four-factor theory* | Lies are expressed differently in terms of arousal, behavior control, **emotion**, and thinking from truth. |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Propagation-based Fundamental Theories

*Studying fake news based on how it spreads*

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

**"Fake news is incorrect but hard to correct"** [5]

It is difficult to correct users' perceptions after fake news has gained their trust.

⬇

**Fake News Early Detection!**

**Providing a solid foundation for epidemic models**

---

[5]A. Roets, et al. 'Fake news': Incorrect, but hard to correct. The role of cognitive ability on the impact of false information on social impressions. Intelligence, 2017.

X. Zhou, R. Zafarani, K. Shu, H. Liu
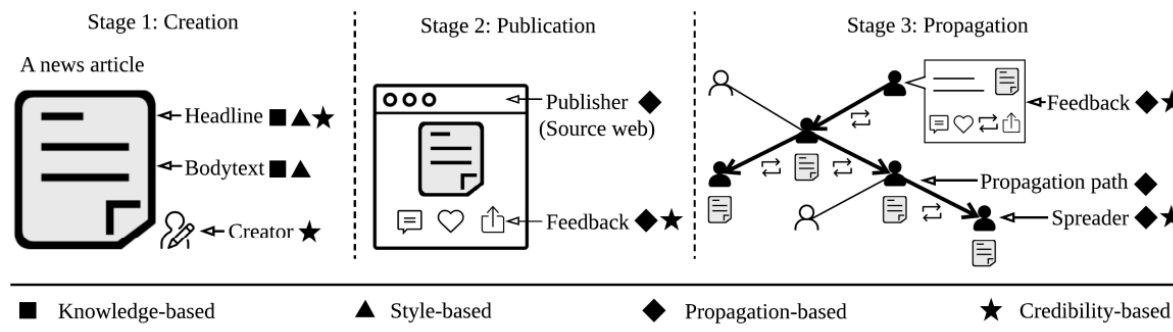
# User-based Fundamental Theories

*Studying fake news from a perspective of users:*
*How users engage with fake news and the role users play (or can play) in fake news creation, propagation, or intervention*

| | | Term | Phenomenon |
|---|---|---|---|
| **User-based** (User's Engagement and Role) | **Social influence** | *Attentional bias* | **Exposure frequency –** individuals tend to believe information is correct after repeated exposures. |
| | | *Validity effect* | |
| | | *Echo chamber effect* | |
| | | *Bandwagon effect* | **Peer pressure –** individuals do something primarily because others are doing it and to conform to be liked and accepted by others. |
| | | *Normative influence theory* | |
| | | *Social identity theory* | |
| | | *Availability cascade* | |
| | **Self-influence** | *Confirmation bias* | **Preexisting knowledge –** individuals tend to trust information that confirms their preexisting beliefs or hypotheses, which they perceive to surpass that of others. |
| | | *Illusion of asymmetric insight* | |
| | | *Naïve realism* | |
| | | *Overconfidence effect* | |
| | **Benefit Influence** | *Prospect theory* | **Loss and gains preference –** people make decisions based on the value of losses and gains rather than the outcome, and they tend to overestimate the likelihood of gains happening rather than losses. |
| | | *Valence effect* | |
| | | *Contrast effect* | |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Detection

- Knowledge-based Fake News Detection

- Style-based Fake News Detection

- Propagation-based Fake News Detection

- Credibility-based Fake News Detection

- Fake News Datasets & Tools

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Detection

- Knowledge-based Fake News Detection
- Style-based Fake News Detection
- Propagation-based Fake News Detection
- Credibility-based Fake News Detection
- Fake News Datasets & Tools



X. Zhou, R. Zafarani, K. Shu, H. Liu

# Knowledge-based Fake News Detection
*Overview*

Knowledge-based fake news detection aims to assess **news authenticity** by comparing the **knowledge** extracted from to-be-verified **news content** with known facts (i.e., true knowledge).

It is also known as **fact-checking**.

- *Manual Fact-checking* – providing ground truth.
- *Automatic Fact-checking* – a better choice for scalability.

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Knowledge-based Fake News Detection

*Overview*

Knowledge-based fake news detection aims to assess **news authenticity** by comparing the **knowledge** extracted from to-be-verified **news content** with known facts (i.e., true knowledge).

It is also known as **fact-checking**.

- *Manual Fact-checking* – providing ground truth.
- *Automatic Fact-checking* – a better choice for scalability.

# Manual Fact-checking

*Classification and comparison*

|  | Expert-based manual fact-checking | Crowd-sourced manual fact-checking |
|---|---|---|
| Fact-checker(s) | One or several domain-expert(s) | A large population of regular individuals |
| Easy to manage? | Yes | No |
| Credibility | High | Comparatively low |
| Scalability | Poor | Comparatively high |
| Current resources (e.g., websites) | Rich | Comparatively poor |

E.g., political bias and conflicting annotations of fact-checkers

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Expert-based Manual Fact-checking

*Current resources*

| | Topics Covered | Content Analyzed | Assessment Labels |
|---|---|---|---|
| **PolitiFact** | American politics | Statements | True; Mostly true; Half true; Mostly false; False; Pants on fire |
| **The Washington Post Fact Checker** | American politics | Statements and claims | One pinocchio; Two pinocchio; Three pinocchio; Four pinocchio; The Geppetto checkmark; An upside-down Pinocchio; Verdict pending |
| **FactCheck** | American politics | TV ads, debates, speeches, interviews and news | True; No evidence; False |
| **Snopes** | Politics and other social and topical issues | News articles and videos | True; Mostly true; Mixture; Mostly false; False; Unproven; Outdated; Miscaptioned; Correct attribution; Misattributed; Scam; Legend |
| **TruthOrFiction** | Politics, religion, nature, aviation, food, medical, etc. | Email rumors | Truth; Fiction; etc. |
| **FullFact** | Economy, health, education, crime, immigration, law | Articles | Ambiguity (no clear labels) |
| **HoaxSlayer** | Ambiguity | Articles and messages | Hoaxes, scams, malware, bogus warning, fake news, misleading, true, humour, spams, etc. |

Multilabel classification

Binary classification

across domains

Multi-modal

**Donald Trump's file**

**Republican from New York**

Donald Trump was elected the 45th president of the United States on Nov. 8, 2016. He has been a real estate developer, entrepreneur and host of the NBC reality show, "The Apprentice." Trump's statements were awarded PolitiFact's 2015 Lie of the Year. Born and raised in New York City, Trump is married to Melania Trump, a former model from Slovenia. Trump has five children and eight grandchildren. Three of his children, Donald Jr., Ivanka, and Eric, serve as executive vice presidents of the Trump Organization.

**The PolitiFact scorecard**

| | | |
|---|---|---|
| True | | 25 (5%) |
| Mostly True | | 61 (11%) |
| Half True | | 83 (15%) |
| Mostly False | | 118 (22%) |
| False | | 173 (32%) |
| Pants on Fire | | 78 (14%) |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Expert-based Manual Fact-checking

*Current resources*

*Reporters Lab – Duke University*

X. Zhou, R. Zafarani, K. Shu, H. Liu

https://reporterslab.org/fact-checking/

# Crowd-sourced Manual Fact-checking

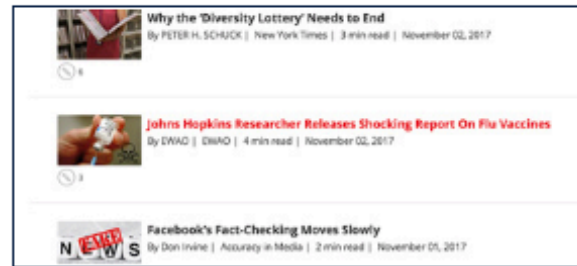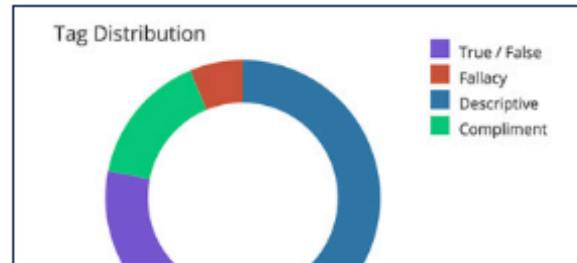*Current resources*

**1** Take an online article that you want to comment on, copy and paste the link into Fiskkit. This allows you to input the article into our system for you to comment on.

**OR** Click on an article you find interesting.

| TRUE/FALSE | FALLACY |
|---|---|
| True | Overly General |
| False | Cherry Picking |
| Matter of Opinion | Straw Man |

| DESCRIPTIVE | COMPLIMENTARY |
|---|---|
| Unsupported | Insightful |
| Overly Simplistic | Well Researched |
| Biased Wording | Funny |

**2** Rate any sentence inside the article by clicking on a sentence & choosing tags that best describe it. Add comments to support your arguments.

**3** See how the article has been rated by other people through our insights page. Share the article so that your friends can come comment too.

X. Zhou, R. Zafarani, K. Shu, H. Liu

http://www.fiskkit.com/

21

# Crowd-sourced Manual Fact-checking

*Current resources*

A. Zhang, et al. A structured response to misinformation: Defining and annotating credibility indicators in news articles. WWW'18 Companion

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Knowledge-based Fake News Detection
*Overview*

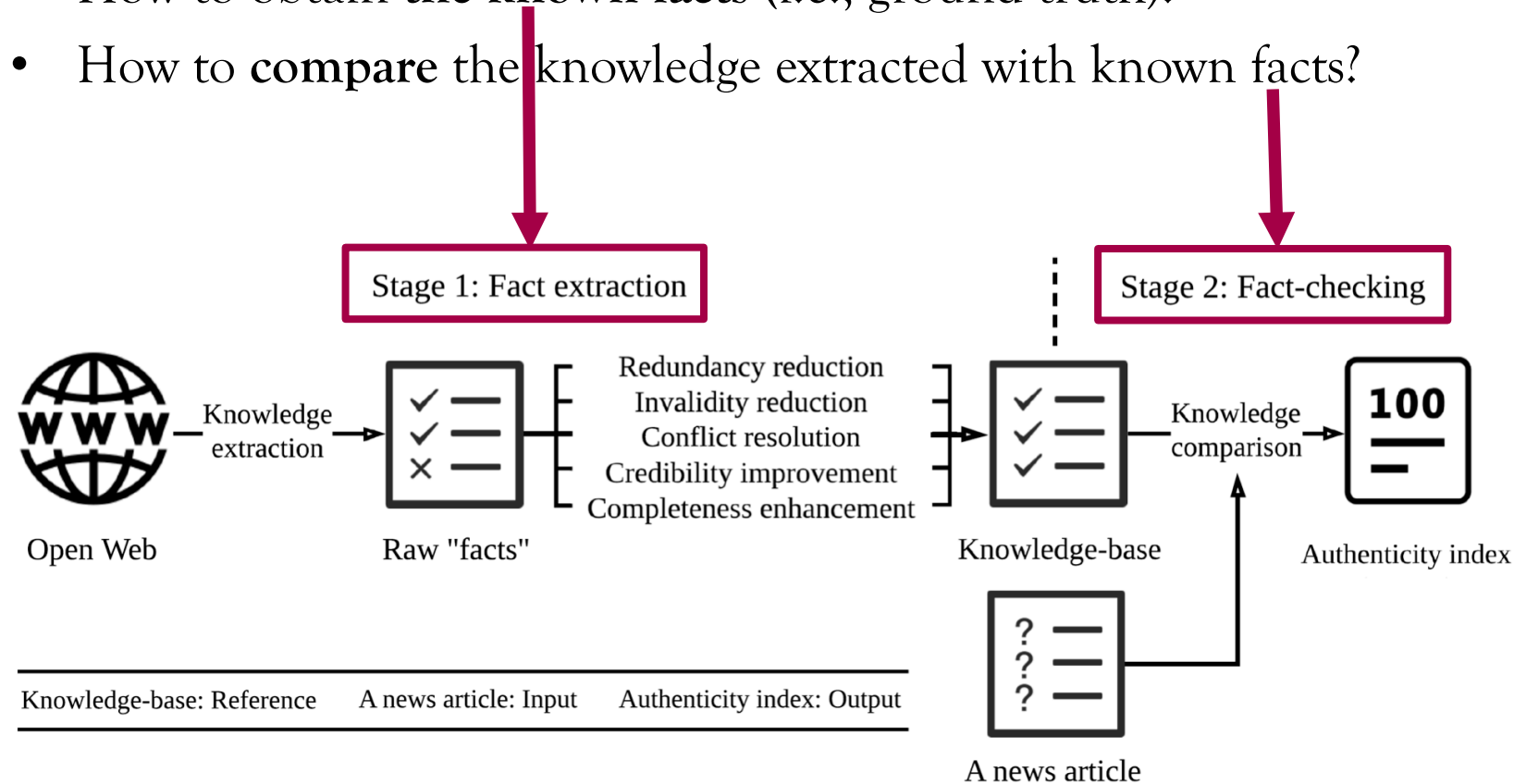Knowledge-based fake news detection aims to assess **news authenticity** by comparing the **knowledge** extracted from to-be-verified **news content** with known facts (i.e., true knowledge).

It is also known as **fact-checking**.

- *Manual Fact-checking* – providing ground truth.

- *Automatic Fact-checking* – a better choice for scalability.

X. Zhou, R. Zafarani, K. Shu, H. Liu

It aims to assess news authenticity by comparing the knowledge extracted from to-be-verified news content with known facts (i.e., true knowledge).

- How to represent "**knowledge**"?
- How to obtain **the known facts** (i.e., ground truth)?
- How to **compare** the knowledge extracted with known facts?
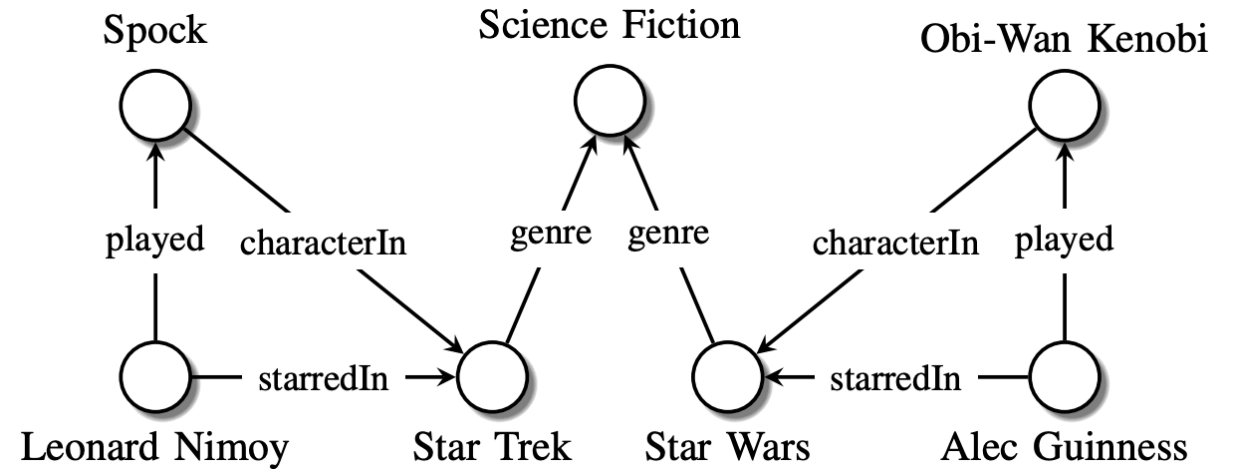
X. Zhou, R. Zafarani, K. Shu, H. Liu

# Knowledge Representation

Knowledge is represented as **a set of (Subject, Predicate, Object) (SPO) triples** extracted from the given information. For example,

*"Leonard Nimoy was an actor who played the character Spock in the science-fiction movie Star Trek"*

| subject | predicate | object |
|---|---|---|
| (LeonardNimoy, | profession, | Actor) |
| (LeonardNimoy, | starredIn, | StarTrek) |
| (LeonardNimoy, | played, | Spock) |
| (Spock, | characterIn, | StarTrek) |
| (StarTrek, | genre, | ScienceFiction) |

X. Zhou, R. Zafarani, K. Shu, H. Liu        The illustration is from: M. Nickel, et al. A Review of Relational Machine Learning for Knowledge Graphs, 2016
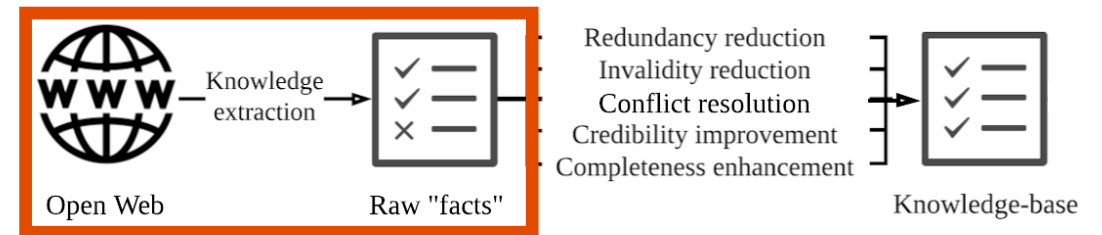
# Stage 1. Fact Extraction

*Constructing knowledge graph to obtain the known facts*

<u>Types</u> of Web content that contain relational information and can be utilized for knowledge extraction by different extractors: **text, tabular data, structured pages** and **human annotations.**[6]

<u>Source(s):</u>

- Single-source knowledge extraction
  - Rely on one comparatively reliable source (e.g., Wiki)
  - Efficient ⬆, Knowledge completeness ⬇
- Open-source knowledge extraction
  - Fuse knowledge from distinct knowledge
  - Efficient ⬇, Knowledge completeness ⬆



---

[6]X. Dong, et al.. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. KDD'14

T1: **Entity Resolution (deduplication/record linkage)** to reduce redundancy

- To identify mentions that refer to the same real-world entity, e.g., *(DonaldJohnTrump, profession, President)* & *(DonaldTrump, profession, President)* should be a redundant pair.

- Current techniques are often distance- or dependence-based.

- Often expensive (requires pairwise distance) computation

- Blocking/Indexing can be used to deal with complexity

T2: **Time Recording** to remove outdated knowledge

- E.g., *(Britain, joinIn, EuropeanUnion)* has been outdated.

- Use Compound Value Type (CVT): facts having beginning and end dates

- Timeliness studies are limited

T3: **Knowledge Fusion** to handle conflicts (often in open-source knowledge extraction)

- E.g., *(DonaldTrump, bornIn, NewYorkCity)* & *(DonaldTrump, bornIn, LosAngeles)* are a conflicting pair.

- Fix by having support values for facts (e.g., website credibility), or using ensemble methods

- Often correlated to (T4).

T4: **Credibility Evaluation** to improve the credibility of knowledge

- E.g., The knowledge extracted from The Onion[7].

- Often focus on analyzing the source website(s).

[7]A https://www.theonion.com/

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Stage 1. Fact Extraction

*Constructing knowledge graph to obtain the known facts*



Open Web — Knowledge extraction → Raw "facts" → Redundancy reduction / Invalidity reduction / Conflict resolution / Credibility improvement / Completeness enhancement → Knowledge-base

T5: *Knowledge Inference/Link Prediction* to infer new facts based on known ones

- Knowledge extracted from online resources, particularly, using a single source, are far from complete.

**Relation machine learning**

**Latent Feature Models,** e.g., **RESCAL**

Assume the existence of knowledge-base triples is <u>conditionally independent</u> given <u>latent features</u> and parameters

**Graph Feature Models**, e.g., **PRA**

Assume the existence of triples is <u>conditionally independent</u> given observed <u>graph features</u> and parameters

**Markov Random Field (MRF) Models**

Assume the existing triples have local interactions

M. Nickel, et al. A Review of Relational Machine Learning for Knowledge Graphs, Proceedings of the IEEE, 2016

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Stage 1. Fact Extraction

*Constructing knowledge graph to obtain the known facts*



Open Web — Knowledge extraction → Raw "facts" → Redundancy reduction / Invalidity reduction / Conflict resolution / Credibility improvement / Completeness enhancement → Knowledge-base

# Stage 1. Fact Extraction

*Existing Knowledge Graphs*

| Name |
|------|
| *Knowledge Vault (KV)* |
| DeepDive [32] |
| NELL [8] |
| PROSPERA [30] |
| YAGO2 [19] |
| Freebase [4] |
| Knowledge Graph (KG) |

Table 1: Comparison of
Freebase and KG rely o
facts means with a prob

[a] Ce Zhang (U Wisconsin), private communication
[b] Bryan Kiesel (CMU), private communication
[c] Core facts, http://www.mpi-inf.mpg.de/yago-naga/yago/downloads.html
[d] This is the number of non-redundant base triples, excluding reverse predicates and "lazy" triples derived from flattening CVTs (complex value types).
[e] http://insidesearch.blogspot.com/2012/12/get-smarter-answers-from-knowledge_4.html

**Open issues:**

1. **Timeliness & Completeness of Knowledge Graphs**

2. **Domain-specific Knowledge Graphs for Fake News Detection**
   *Related tutorial:* X. Ren, et al., Scalable Construction and Querying of Massive Knowledge Bases, WWW tutorial, 2018.
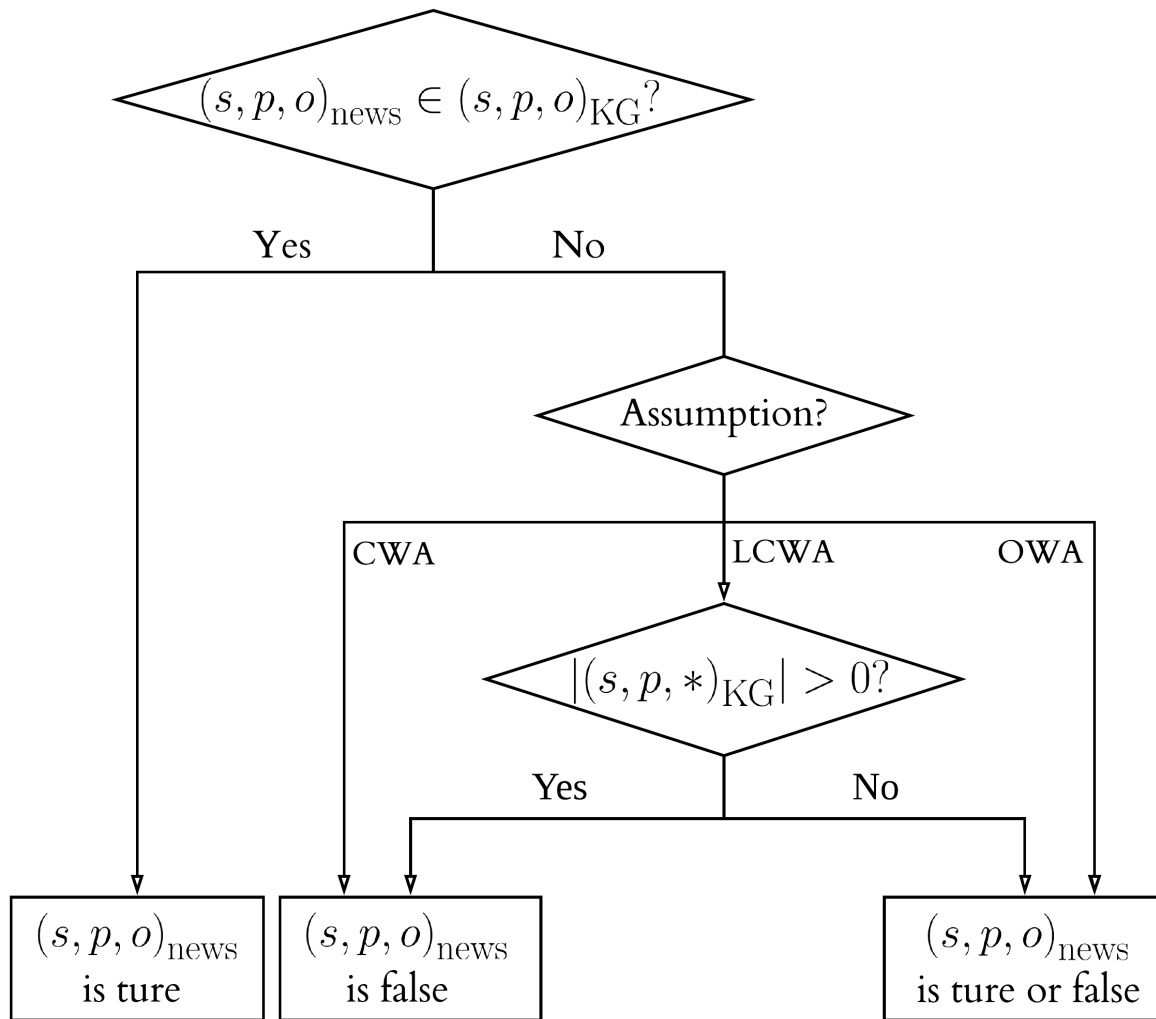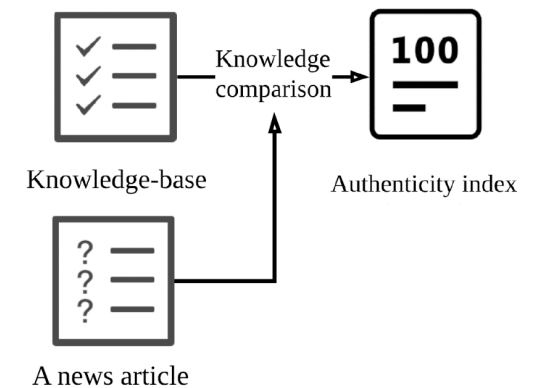
X. Zhou, R. Zafarani, K. Shu, H. Liu       Source: X. Dong, et al.. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. KDD'14

Stage 2. Fact-checking

*Comparing knowledge between news articles and knowledge graphs*

$(s, p, o)_{\text{news}} \in (s, p, o)_{\text{KG}}$?

Yes / No

Assumption?

CWA / LCWA / OWA

$|(s, p, *)_{\text{KG}}| > 0$?

Yes / No

$(s, p, o)_{\text{news}}$ is ture

$(s, p, o)_{\text{news}}$ is false

$(s, p, o)_{\text{news}}$ is ture or false

Knowledge Inference

KG: Knowledge Graph
CWA: Closed-World Assumption
LCWA: Local Closed-World Assumption
OWA: Open-World Assumption

X. Zhou, R. Zafarani, K. Shu, H. Liu

Knowledge comparison

Knowledge-base

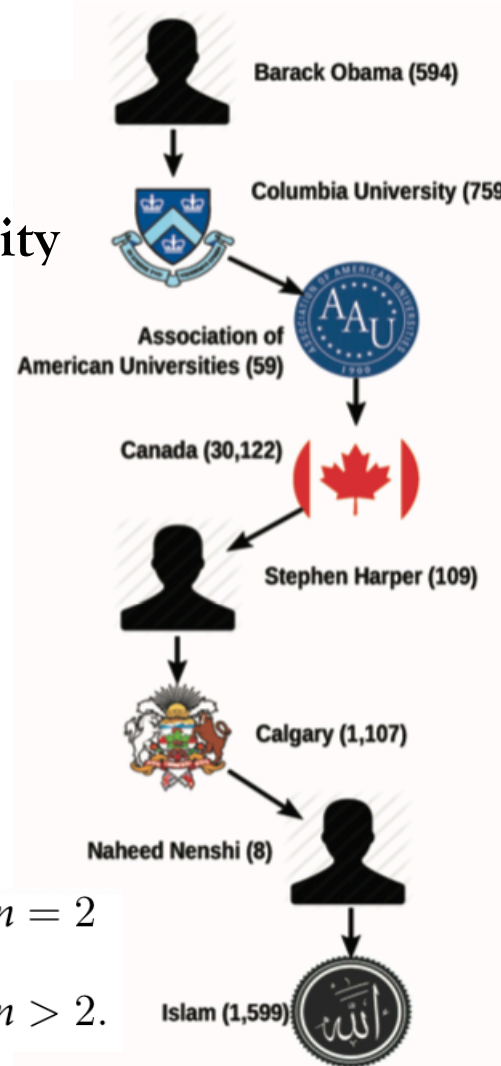Authenticity index

A news article

*Shortest path-based method:*

By finding the **shortest path** between concept nodes under properly defined **semantic proximity** metrics on knowledge graphs

$$\tau(e) = \max \mathcal{W}(P_{s,o}).$$

$$\mathcal{W}(P_{s,o}) = \mathcal{W}(v_1 \dots v_n) = \left[ 1 + \sum_{i=2}^{n-1} \log k(v_i) \right]^{-1}$$

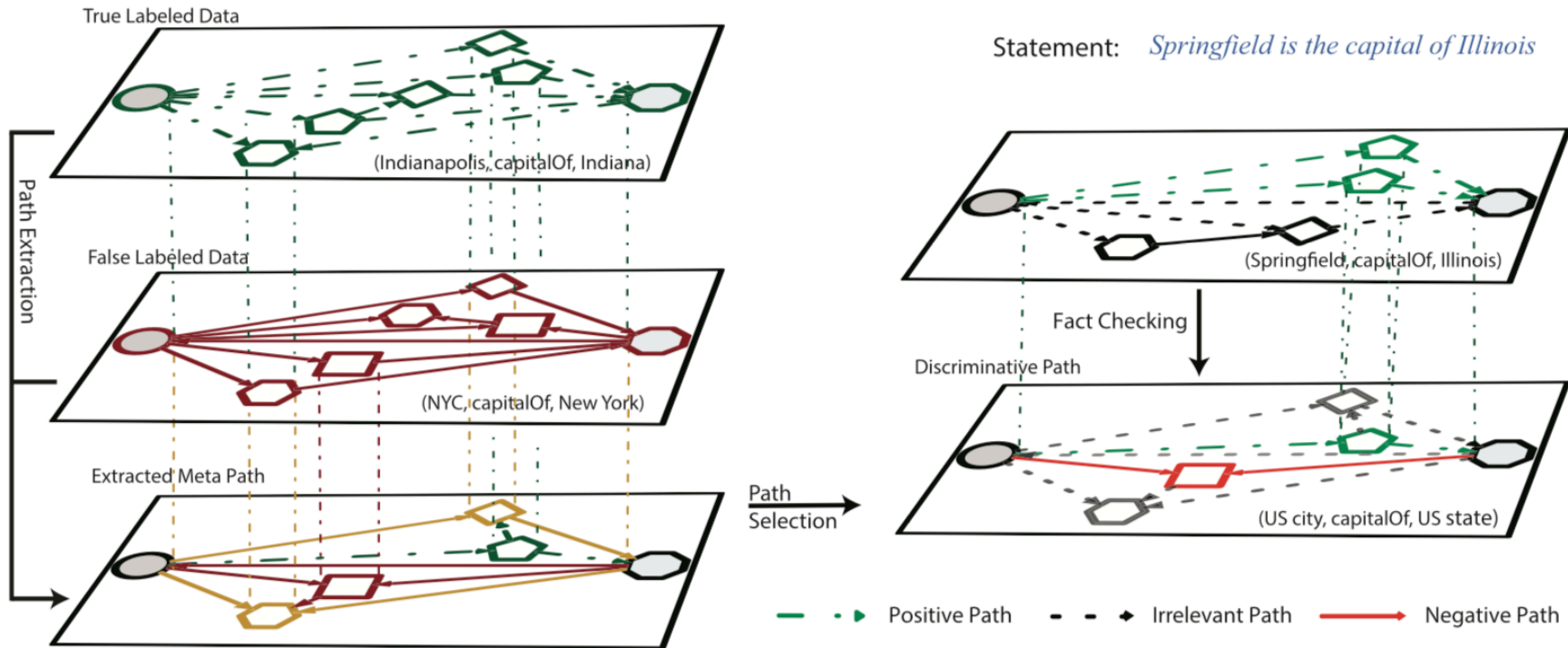An alternative formulation (widest bottleneck)

$$\mathcal{W}_u(P_{s,o}) = \mathcal{W}_u(v_1 \dots v_n) = \begin{cases} 1 & n = 2 \\ \left[ 1 + \max_{i=2}^{n-1} \left\{ \log k(v_i) \right\} \right]^{-1} & n > 2. \end{cases}$$



Barack Obama (594)

Columbia University (759)

Association of American Universities (59)

Canada (30,122)

Stephen Harper (109)

Calgary (1,107)

Naheed Nenshi (8)

Islam (1,599)

G. Ciampaglia, et al. Computational Fact Checking from Knowledge Networks, 2016

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Stage 2. Fact-checking

*Knowledge Inference for unknown SPO triples: Illustrated studies*

31

*Discriminative path-based method:*

*Knowledge Inference for unknown SPO triples: Illustrated studies*

B. Shi and T. Weninger, Discriminative predicate path mining for fact checking in knowledge graphs, 2015

# Knowledge Inference

*Comparison*

Knowledge inference can be conducted on both Stage I, when constructing knowledge graphs, and Stage II for fact-checking.

| Operation \ Stage | Knowledge Graph Construction | Fact-checking |
|---|---|---|
| **Entity/Node** | *Few* operations on entities | Generally requires *additional* operations on entities, e.g., entity matching |
| **Relationship/Edge** | Inference targets relationships between *each pair of* given entities | Inference only targets relationships among *partial* entities |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Detection

- Knowledge-based Fake News Detection
- Style-based Fake News Detection
- Propagation-based Fake News Detection
- Credibility-based Fake News Detection
- Fake News Datasets & Tools

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Style-based Fake News Detection
*Overview*

**Style-based Fake News Detection** is able to assess <u>news intention</u> by comparing the *writing style* extracted from to-be-verified *news content* with fake news style.

**Fake News Style** is a set of <u>machine learning features </u>that can well represent fake news and differentiate fake news from truth.

- *Textual* (*linguistic)* style features
- *Visual* style features

- *Manually* select features → Often within a *supervised* machine learning framework

- *Automatically* select features → Often within a *deep* machine learning framework

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Style-based Fake News Detection
*Overview*

**Style-based Fake News Detection** is able to assess <u>news intention</u> by comparing the *writing style* extracted from to-be-verified *news content* with fake news style.

**Fake News Style** is a set of <u>machine learning features</u> that can well represent fake news and differentiate fake news from truth.

- *Textual (linguistic)* style features

- *Visual* style features

"**More people watched President Trump's 2019 State of the Union address on television than watched Super Bowl Super Bowl LlII**"

X. Zhou, R. Zafarani, K. Shu, H. Liu

36

# Textual (Linguistic) Style of Fake News

*Structure-based* language features

| Level | Feature(s) | Technique(s) and Tool(s) | Reference(s) |
|---|---|---|---|
| Lexicon | Words | Bag of words | Perez-Rosas et al., 2017 |
| | | + n-gram to capture the word sequence | |
| | | + TF-IDF to unify the content length | |
| Syntax | Part-Of-Speech (POS) Tags | POS Taggers | Feng et al., 2012 Petrov and Klein, 2007 |
| | Context-Free Grammars (CFGs) | Probabilistic Context Free Grammars (PCFGs) Parsers | |
| Semantic | Psycholinguistic Words | Linguistic Inquiry and Word Count (LIWC) | Perez-Rosas et al., 2017 |
| Discourse | Rhetorical Relationships | Rhetorical Structure Theory (RST) Parser | Rubin and Lukoianova, 2015 Ji and Eisenstein, 2014 |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Textual (Linguistic) Style of Fake News

*Structure-based* language features

"The rat ate the cheese"

| Level | Feature(s) | Technique(s) and Tool(s) | Reference(s) |
|---|---|---|---|
| Lexicon | Words | Bag of words | Perez-Rosas et al., 2017 |
| | | + n-gram to capture the word sequence | |
| | | TF-IDF to unify the content length | |
| Syntax | Part-Of-Speech (POS) Tags | Taggers | Feng et al., 2012 Petrov and Klein, 2007 |
| | Context Free Grammars (PCFGs) Parsers | | |
| Semantic | | Linguistic Inquiry and Word Count (LIWC) | Perez-Rosas et al., 2017 |
| Discourse | Rhetorical Relationships | Rhetorical Structure Theory (RST) Parser | Rubin and Lukoianova, 2015 Ji and Eisenstein, 2014 |

"the": 2    "rat":  1
"ate":  1    "cheese": 1

P("ate"|"rat") = ?
P("cheese"|"ate") = ?

# Textual (Linguistic) Style of Fake News

*Structure-based* language features

"The rat ate the cheese"

| Level | Feature(s) | Technique(s) and To... |
|---|---|---|
| Lexicon | Words | Bag of words |
| | | + n-gram to capture t... |
| | | + TF-IDF to unify the... |
| Syntax | Part-Of-Speech (POS) Tags | POS Taggers |
| | Context-Free Grammars (CFGs) | Probabilistic Context... |
| Semantic | Psycholinguistic Words | Linguistic Inquiry an... |



NN: 2 ("rat", "cheese")
DT: 1 ("the")
VB: 1 ("ate")

S ➔ NP VP    NP ➔ DT NN
VP ➔ VB NP    DT ➔ the
NN ➔ rat        VB ➔ ate
NN ➔ cheese

X. Zhou, R. Zafarani, K. Shu, H. Liu

39

# Textual (Linguistic) S

*Structure-based* language features

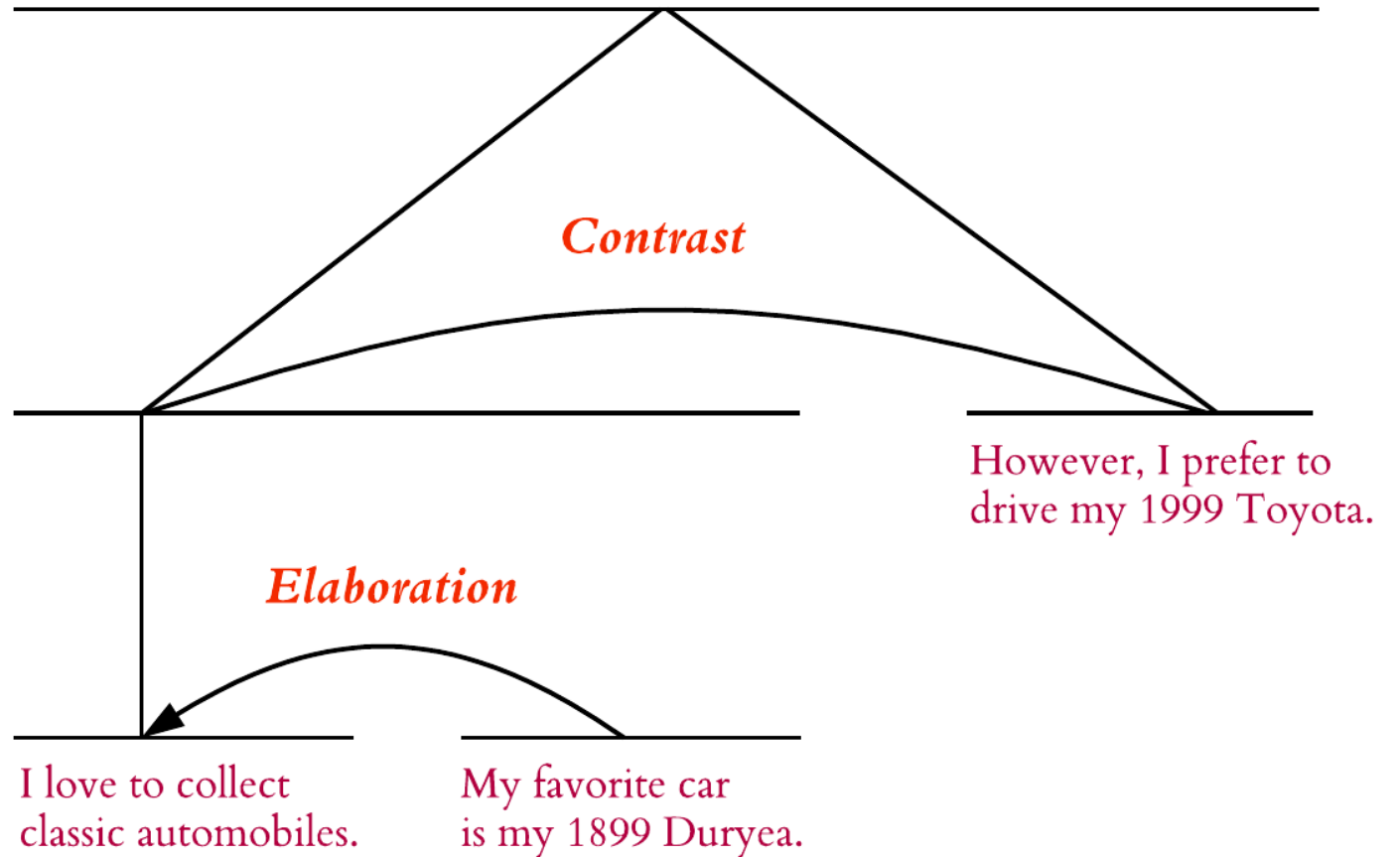| Level | Feature(s) |
|-------|-----------|
| Lexicon | Words |
| Syntax | Part-Of-Speech (POS) Tags |
| | Context-Free Grammars (CFGs) |
| Semantic | Psycholinguistic Words |
| Discourse | Rhetorical Relationships |

X. Zhou, R. Zafarani, K. Shu, H. Liu

| Category | Abbrev | Examples |
|----------|--------|----------|
| Word count | WC | - |
| **Summary Language Variables** | | |
| Analytical thinking | Analytic | - |
| Clout | Clout | - |
| Authentic | Authentic | - |
| Emotional tone | Tone | - |
| Words/sentence | WPS | - |
| Words > 6 letters | Sixltr | - |
| Dictionary words | Dic | - |
| **Linguistic Dimensions** | | |
| Total function words | funct | it, to, no, very |
| Total pronouns | pronoun | I, them, itself |
| Personal pronouns | ppron | I, them, her |
| 1st pers singular | i | I, me, mine |
| 1st pers plural | we | we, us, our |
| 2nd person | you | you, your, thou |
| 3rd pers singular | shehe | she, her, him |
| 3rd pers plural | they | they, their, they'd |
| Impersonal pronouns | ipron | it, it's, those |
| Articles | article | a, an, the |
| Prepositions | prep | to, with, above |
| Auxiliary verbs | auxverb | am, will, have |
| Common Adverbs | adverb | very, really |
| Conjunctions | conj | and, but, whereas |
| Negations | negate | no, not, never |
| **Other Grammar** | | |
| Common verbs | verb | eat, come, carry |
| Common adjectives | adj | free, happy, long |
| Comparisons | compare | greater, best, after |
| Interrogatives | interrog | how, when, what |
| Numbers | number | second, thousand |
| Quantifiers | quant | few, many, much |
| **Psychological Processes** | | |
| Affective processes | affect | happy, cried |
| Positive emotion | posemo | love, nice, sweet |
| Negative emotion | negemo | hurt, ugly, nasty |
| Anxiety | anx | worried, fearful |
| Anger | anger | hate, kill, annoyed |
| Sadness | sad | crying, grief, sad |
| Social processes | social | mate, talk, they |
| Family | family | daughter, dad, aunt |

| Category | Abbrev | Examples |
|----------|--------|----------|
| Friends | friend | buddy, neighbor |
| Female references | female | girl, her, mom |
| Male references | male | boy, his, dad |
| Cognitive processes | cogproc | cause, know, ought |
| Insight | insight | think, know |
| Causation | cause | because, effect |
| Discrepancy | discrep | should, would |
| Tentative | tentat | maybe, perhaps |
| Certainty | certain | always, never |
| Differentiation | differ | hasn't, but, else |
| Perceptual processes | percept | look, heard, feeling |
| See | see | view, saw, seen |
| Hear | hear | listen, hearing |
| Feel | feel | feels, touch |
| Biological processes | bio | eat, blood, pain |
| Body | body | cheek, hands, spit |
| Health | health | clinic, flu, pill |
| Sexual | sexual | horny, love, incest |
| Ingestion | ingest | dish, eat, pizza |
| Drives | drives | |
| Affiliation | affiliation | ally, friend, social |
| Achievement | achieve | win, success, better |
| Power | power | superior, bully |
| Reward | reward | take, prize, benefit |
| Risk | risk | danger, doubt |
| Time orientations | TimeOrient | |
| Past focus | focuspast | ago, did, talked |
| Present focus | focuspresent | today, is, now |
| Future focus | focusfuture | may, will, soon |
| Relativity | relativ | area, bend, exit |
| Motion | motion | arrive, car, go |
| Space | space | down, in, thin |
| Time | time | end, until, season |
| Personal concerns | | |
| Work | work | job, majors, xerox |
| Leisure | leisure | cook, chat, movie |
| Home | home | kitchen, landlord |
| Money | money | audit, cash, owe |
| Religion | relig | altar, church |
| Death | death | bury, coffin, kill |
| Informal language | informal | |
| Swear words | swear | fuck, damn, shit |
| Netspeak | netspeak | btw, lol, thx |
| Assent | assent | agree, OK, yes |
| Nonfluencies | nonflu | er, hm, umm |
| Fillers | filler | Imean, youknow |

# Textual (Linguistic) Style of Fake News

*Structure-based* language features

| Level | Feature(s) |
|---|---|
| Lexicon | Words |
| Syntax | Part-Of-Speech (POS) Tags |
| | Context-Free Grammars (CFGs) |
| Semantic | Psycholinguistic Words |
| Discourse | Rhetorical Relationships |



*Contrast*

However, I prefer to drive my 1999 Toyota.

*Elaboration*

I love to collect classic automobiles.

My favorite car is my 1899 Duryea.

X. Zhou, R. Zafarani, K. Shu, H. Liu

39

# Textual (Linguistic) Style of Fake News

*Performance of structure-based language features*

| | Level(s) | Feature(s) | [Ott et al 2011] | [Feng et al 2012a] | [Shojaee et al. 2013] | [Mukherjee et al. 2013b] | [Li et al 2014] | [Pérez-Rosas and Mihalcea 2014] | [Pérez-Rosas et al. 2015] | [Pérez-Rosas and Mihalcea 2015] | [Li et al 2017b] | [Ott et al 2011] | [Shojaee et al. 2013] | [Li et al 2014] | [Pérez-Rosas et al. 2015] | [Abouelenien et al. 2017] | [Braud and Søgaard 2017] | [Pérez-Rosas et al. 2015] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Within Levels | Lexicon | UG | .884 | .729 | | .663 | .668 | .691 | .609 | .695 | .825 | .884 | | .645 | .763 | .585 | .717 | .678 |
| | | BG | .896 | .708 | | .661 | | | | | .804 | .889 | | | | | .696 | |
| | | UG+BG | | .738 | | | | | | | .637 | | | | | | | |
| | | Others | | | .810 | | | | | | | | .700 | | | | | |
| | Syntax | POS | .730 | | | .564 | .638 | | | .695 | | | | .690 | | .513 | .717 | |
| | | CFG | | .742 | | | | | | .654 | | | | | | .513 | | |
| | | Others | .768 | | .760 | | | | | | .525 | | .690 | | .627 | | | .534 |
| | Semantic | LIWC | | | | | .633 | .691 | .602 | .534 | | | | | .695 | .500 | .504 | .661 |
| | Discourse | RR | | | | | | | | | | | | | | | .553 | |
| Across Levels | Lexicon + Syntax | UG+POS | | .733 | | | | | | | .831 | | | | | | | |
| | | UG+CFG | | .769 | | | | | | | | | | | | | | |
| | | BG+POS | | | | .664 | | | | | .808 | | | | | | | |
| | | BG+CFG | | | | .659 | | | | | | | | | | | | |
| | | UG+BG+POS | | | | | | | | | | | | | | | .760 | |
| | | Others+Others | | | .840 | | | | | | | | .740 | | | | | |
| | Lexicon + Semantic | UG+LIWC | | | | | | | | | .622 | | | | | .594 | | |
| | | BG+LIWC | .898 | | | .661 | | | | | | | | | | | | |
| | Lexicon + Syntax + Semantic | UG+POS+ LIWC | | | | | | | | | .653 | | | | .636 | | | .576 |

UG: Unigram    BG: Bigram    POS: Part-of-Speech tags    CFG: Context-Free Grammar (particularly refers to lexicalized production rules)
LIWC: Linguistic Inquiry and Word Count    RR: Rhetorical Relations

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Textual (Linguistic) Style of Fake News

*Attribute-based* *language features*

- Most related studies belong to the general area of **Deception Detection**.

- Deception is **disinformation**, including **fake statements, fake reviews, <u>fake news</u>**, etc.

- Attributes are generally inspired from **forensic psychological theories**, e.g.,

| Term | Phenomenon |
|---|---|
| *Undeutsch hypothesis* | A statement based on a factual experience differs in content and **quality** from that of fantasy |
| *Reality monitoring* | <u>Actual events</u> are characterized by higher levels of **sensory-perceptual** information. |
| *Four-factor theory* | <u>Lies</u> are expressed differently in terms of arousal, behavior control, **emotion**, and thinking from truth. |

# Textual (Linguistic) Style of Fake News

*Attribute-based language features*

| | Attribute Type | Feature |
|---|---|---|
| 1 | **Quantity** | Character count |
| | | Word count |
| | | Noun count |
| | | Verb count |
| | | Number of noun phrases |
| | | Sentence count |
| | | Paragraph count |
| | | Number of modifiers (e.g., adjectives and adverbs) |
| 2 | **Complexity** | Average number of clauses per sentence |
| | | Average number of words per sentence |
| | | Average number of characters per word |
| | | Average number of punctuations per sentence |
| 3 | **Uncertainty** | Percentage of modal verbs "Can"; "May"; "Shall" |
| | | Percentage of centainty terms "Always"; "Never" |
| | | Percentage of generalizing terms "Generally"; "All"; "Many" |
| | | Percentage of tentative terms "Possibly"; "Probably" |
| | | Percentage of numbers and quantifiers |
| | | Number of question marks |
| 4 | **Subjectivity** | Percentage of subjective verbs "Feel"; "Indicate"; "Believe" |
| | | Percentage of report verbs "Suggest"; "Speculate" |
| | | Percentage of factive verbs "Accept"; "Note"; "Confirm" |
| | | Percentage of imperative commands "Give"; "Do" |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Textual (Linguistic) Style of Fake News

*Attribute-based language features*

| Attribute Type | Feature |
|---|---|
| **Non-immediacy** (5) | Percentage of passive voice |
| | Percentage of rhetorical questions |
| | Self reference: 1st person singular pronouns |
| | Group reference: 1st person plural pronouns |
| | Other reference: 2nd and 3rd person pronouns |
| | Number of quotations |
| **Sentiment** (6) | Percentage of positive words |
| | Percentage of negative words |
| | Number of exclamation marks |
| | Activation: the dynamics of emotional state |
| **Diversity** (7) | Lexical diversity: unique words or terms (%) |
| | Content word diversity: unique content words (%) |
| | Redundancy: unique function words (%) |
| **Informality** (8) | Typographical error ratio: misspelled words (%) |
| **Specificity** (9) | Temporal ratio |
| | Spatial ratio |
| | Sensory ratio |
| | Causation terms |
| | Exclusive terms |
| **Readablity** (10) | (e.g., Flesch-Kincaid and Gunning-Fog index) |

"Car"; "Red"
"Are"; "An"

$0.4[(\#words/\#sentences)+(\#long\_words/\#words)$

The general construct of immediacy and nonimmediacy refers to (non-)verbal behaviors that create **a psychological sense of closeness or distance**.

X. Zhou, R. Zafarani, K. Shu, H. Liu

45

# Textual (Linguistic) Style of Fake News

*Performance of attribute-based language features*

| Attribute Type | [Newman et al. 2003] | [Fuller et al. 2009] | [Matsumoto and Hwang 2015] | [Derrick et al. 2013] | [Zhou et al. 2004b] | [Hancock et al. 2007] | [Anderson and Simester 2014] | [Braun and Van Swol 2016] | [Bond and Lee 2005] | [Zhou and Zenebe 2008] | [Ali and Levine 2008] | [Humpherys et al. 2011] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Quantity** | | + | + | − | + | + | + | − | | + | + | + |
| **Complexity** | | | | | − | | | | | | | + |
| **Uncertainty** | | − | | | + | + | | + | | | − | − |
| **Non-immediacy** | + | + | + | | + | + | + | + | + | + | | + |
| **Sentiment** | − | + | − | | | − | | + | − | | + | + |
| **Diversity** | | − | | − | − | | − | | | − | − | − |
| **Informality** | | | | | + | | | | | + | | |
| **Specificity** | − | − | + | | − | | | | | − | | − |

+: The attribute is positively related to the existence of deception;
−: The attribute is negatively related to the existence of deception.

- Quantity ↑
- Non-immediacy ↑
- Informality ↑
- Diversity ↓
- Specificity ↓

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Style-based Fake News Detection

*Overview*

**Style-based Fake News Detection** is able to assess <u>news intention</u> by comparing the *writing style* extracted from to-be-verified *news content* with fake news style.

**Fake News Style** is a set of <u>machine learning features</u> that can well represent fake news and differentiate fake news from truth.

- *Textual (linguistic)* style features
- *Visual* style features

X. Zhou, R. Zafarani, K. Shu, H. Liu

A fully connected layer with softmax

Fake News Detector

Word Embedding

Text Feature

$R_T$

Reddit, has, found, a, much, clearer, photo…

Text-CNN

Multimodal Feature

pred-fc

Concatenation

$R_F$

VGG-19

vis-fc

$R_V$

Visual Feature

reversal

adv-fc1

adv-fc2

Event Discriminator

Multimodal Feature Extractor

Two fully connected layers with activation functions

W. Yaqing, et al., EANN: Event Adversarial Neural Networks for Multi-Modal Fake News Detection. *KDD'18*

# Visual Style of Fake News

*An illustration: EANN*

EANN:
multi-modal;
**adversarial network inspired**;
fake news early detection

Fake News Early Detection:
extract a set of **generalizable** and **discriminable** features to represent news content and detect fake news

# Knowledge- & Style-based Fake News Detection

*Summary*

> How to involve *social context information* of fake news, e.g., its propagation patterns on social networks?

| | Knowledge-based fake news detection | Style-based fake news detection |
|---|---|---|
| Information utilized | News content | News content |
| Modality involved | Single: only text | Single or multi: text, visual, etc. |
| Objective(s) evaluated | News authenticity | News authenticity and intention |
| Framework for solving the problem | Link prediction | Machine learning |
| Related topic | Fact-checking | Deception detection |
| Open issues | Timeliness and completeness of knowledge graphs | Cross-domain, language, topic fake news studies |

# Fake News Detection

- Knowledge-based Fake News Detection
- Style-based Fake News Detection
- Propagation-based Fake News Detection
- Credibility-based Fake News Detection
- Fake News Datasets & Tools

50

# Propagation-based Fake News Detection

*Overview*

**Propagation-based Fake News Detection** utilizes <u>social context information</u> to explore the relationships among entities in news propagation.

- Entities, e.g., spreaders (users) of news, publishers of news, posts of users

- Relationships among the same or different entities

Basis of propagation-based fake news detection approaches

- **News cascades (propagation trees)** – a *direct* way to present news propagation

- **Self-defined graphs (networks)** – an *indirect* way to present news propagation

X. Zhou, R. Zafarani, K. Shu, H. Liu

51

# Propagation-based Fake News Detection

*Overview*

**Propagation-based Fake News Detection** utilizes <u>social context information</u> to explore the relationships among entities in news propagation.

- Entities, e.g., spreaders (users) of news, publishers of news, posts of users

- Relationships among the same or different entities

Basis of propagation-based fake news detection approaches

- **News cascades (propagation trees)** – a *direct* way to present news propagation

- **Self-defined graphs (networks)** – an *indirect* way to present news propagation

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Cascade

*Definition*

**A news cascade:** One propagation path of a news article

**Root node:** The original post of user related to the news article

**Other node:** The re-post of the post of parent node

**Directed Edge:** Post → repost relationships

X. Zhou, R. Zafarani, K. Shu, H. Liu

**Fake news spreads deeper than the truth**

**Fake news spreads farther than the truth**

# News Cascade

*Illustrated studies –*

*A. Cascade-based pattern discovering*

- Depth: 3
- Breadth: 1,2,3,1
- Size: 7

S. Vosoughi, et al. The spread of true and false news online. Science, 2018

X. Zhou, R. Zafarani, K. Shu, H. Liu

**Fake news spreads more broadly than the truth**

**Fake news spreads faster than the truth**

# News Cascade

*Illustrated studies –*
*A. Cascade-based **pattern***
*discovering of fake news*

- Depth: 3
- Breadth: 1,2,3,1
- Size: 7

S. Vosoughi, et al. The spread of true
and false news online. Science, 2018

X. Zhou, R. Zafarani, K. Shu, H. Liu

55

# News Cascade

*Illustrated studies –*
*A. Cascade-based **pattern***
*discovering of fake news*

- Depth: 3
- Breadth: 1,2,3,1
- Size: 7

Hop

0    1    2    3

**Political** fake news spreads **deeper, farther, more broadly** and **faster** than fake news in other domains

S. Vosoughi, et al. The spread of true and false news online. Science, 2018

X. Zhou, R. Zafarani, K. Shu, H. Liu

56

# News Cascade

*Illustrated studies –*
*B. Fake news detection based on cascade **similarity***



- **Approval score**
- **Doubt score**
- **Sentiment score**

**Opinion leader**

**Normal user**

**Random walk graph kernel**

Challenges:
**Computational expense,** as similarity will be computed between pairwise cascades.

K. Wu, et al. False Rumors Detection on Sina Weibo by Propagation Structures, ICDE'15

X. Zhou, R. Zafarani, K. Shu, H. Liu

A **TF-IDF** vector as post representation

**GRU** to learn hidden state (pass important information) of post

Pooling

Softmax

News label (*true* or *fake*)

# News Cascade

*Illustrated studies –*

*C. Fake news detection based on cascade* **representation**

Challenges:

**Cascade depth sensitivity,** as the depth of cascade is equivalent to that of neural network.

J. Ma, et al. Rumor Detection on Twitter with Tree-structure Recursive Neural Networks, ACL'18

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Propagation-based Fake News Detection
*Overview*

- **Homogeneous Networks** contain a single type of nodes and edge.
- **Heterogeneous Networks** contain multiple types of nodes or edges.
- **Hierarchical Networks**, whose various nodes and edges form set-subset relationships.

Basis of propagation-based fake news detection approaches

- **News cascades (propagation trees)** – a *direct* way to present news propagation
- **Self-defined graphs (networks)** – an *indirect* way to present news propagation

X. Zhou, R. Zafarani, K. Shu, H. Liu

59

**Stance Network**

News article

Similarity of text, stance, topic, etc.

X. Zhou, R. Zafarani, K. Shu, H. Liu

**Stance Network**

User post

Similarity of text, stance, topic, etc.

# *Homogeneous Network*

*Illustrations of homogeneous networks*

X. Zhou, R. Zafarani, K. Shu, H. Liu

LDA

Jensen-Shannon Distance

Topic-Viewpoint

$+$ Topic-Viewpoint $-$

$+$ $-$ Topic-Viewpoint

$+$ $-$ $+$ $-$

Topic-Viewpoint $-$ Topic-Viewpoint

$$\arg\min_{\mathbf{c}} \underbrace{\mu||\mathbf{c} - \mathbf{c}_0||^2}_{\text{Fitting constraint}} \underbrace{+ (1 - \mu) \sum_{i,j=1}^{n} \mathbf{A}_{ij}\left(\frac{\mathbf{c}_i}{\sqrt{\mathbf{D}_{ii}}} - \frac{\mathbf{c}_j}{\sqrt{\mathbf{D}_{jj}}}\right)^2}_{\text{Smoothness constraint}}$$

Z. Jin, et al. News Verification by Exploiting Conflicting Social Viewpoints in Microblogs, AAAI'16

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Homogeneous Network

*Illustrations of related studies*

conflicting viewpoints mining
tweets

Assumption:
Posts with the same (contradicting) viewpoints rise (weaken) each other's credibility.

original credibility network | credibility network with conflicting relations

**Friendship Network**

**User/Spreader**

**Friend relationship**

# *Homogeneous Network*

*Illustrations of homogeneous networks*

X. Zhou and R. Zafarani, Fake News in Networks: Patterns, Representation and Detection.

X. Zhou, R. Zafarani, K. Shu, H. Liu

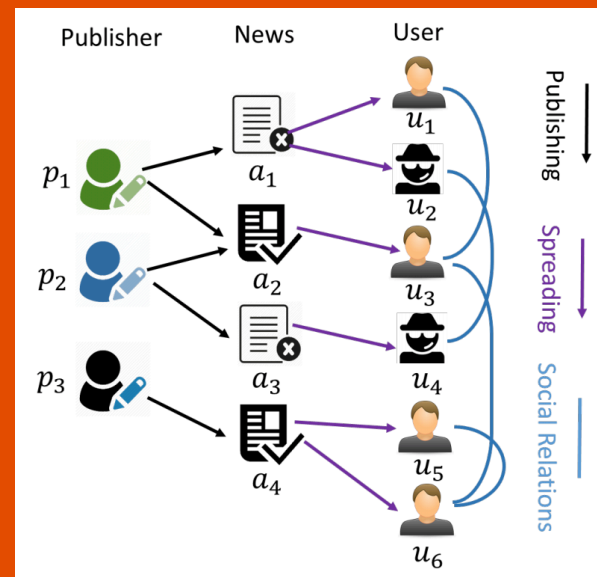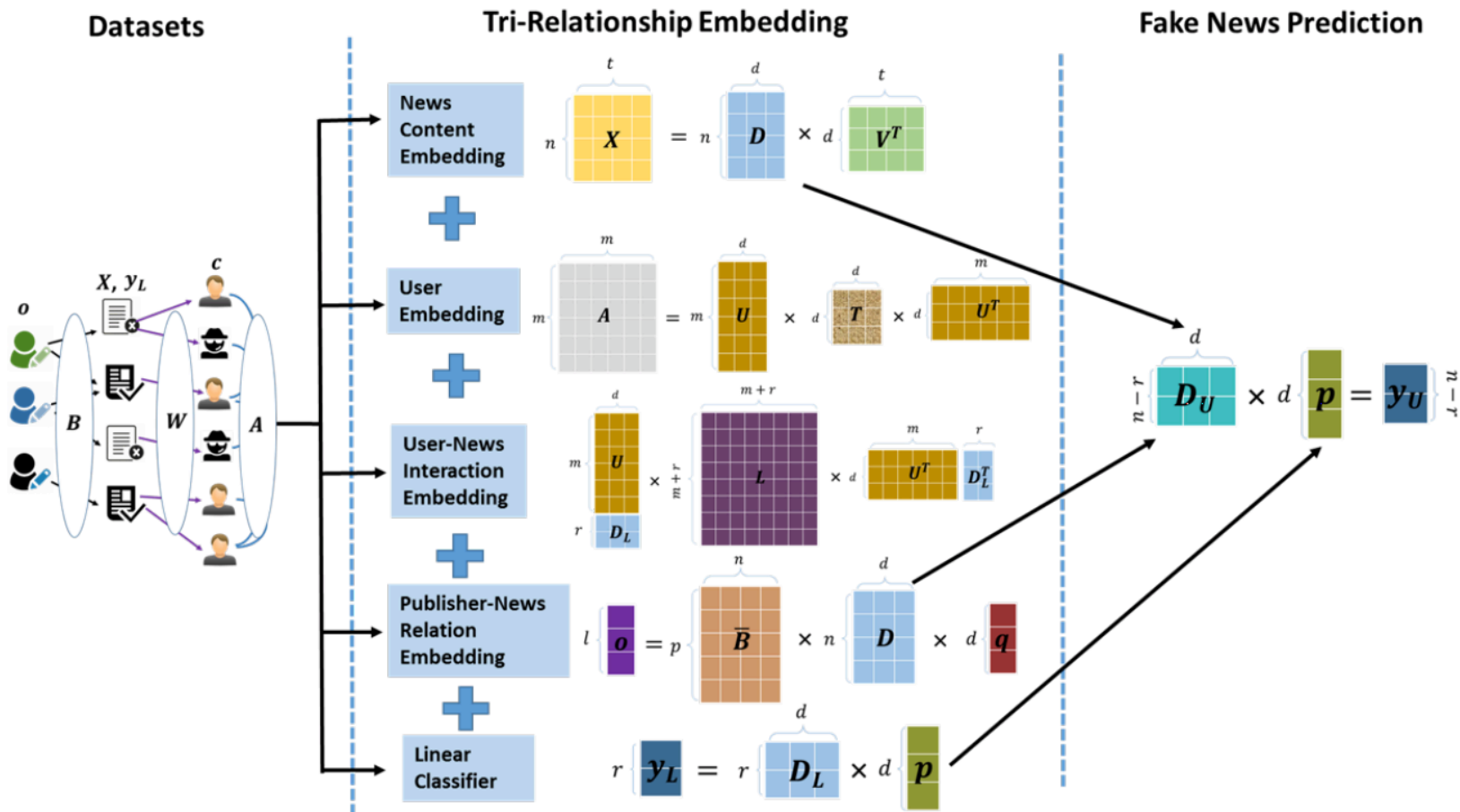# Heterogeneous Network

*Illustrations of homogeneous networks and related studies*

Assumption:
- Credible user → Credible tweets
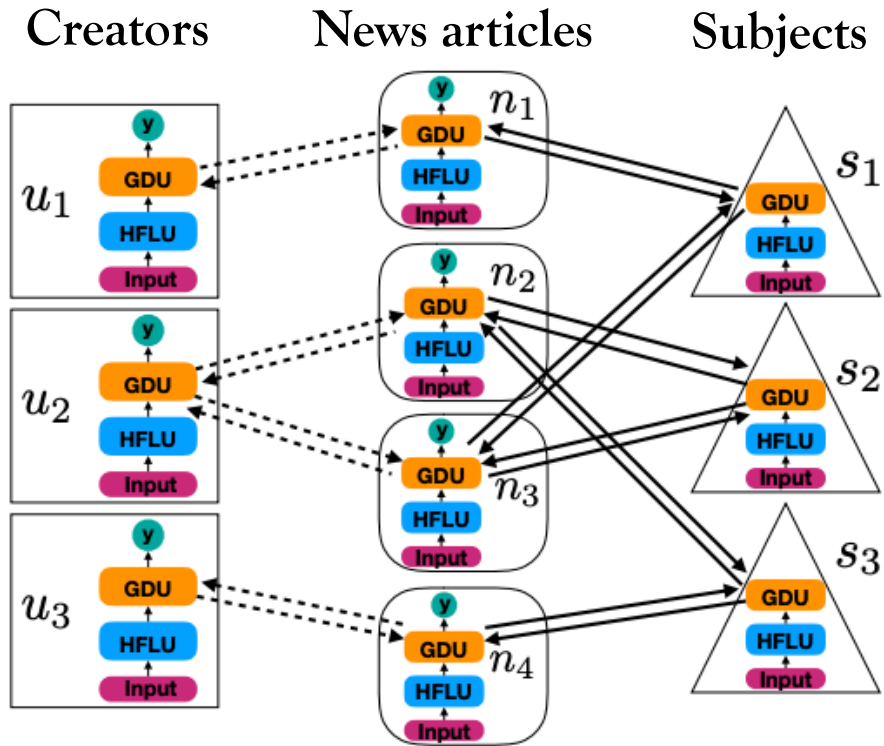- Average credibility of tweets: Credible events > Incredible events

M. Gupta, et al. Evaluating Event Credibility on Twitter, SDM'12

# Heterogeneous Network

*Illustrations of homogeneous networks and related studies*



X. Zhou, R. Zafarani, K. Shu, H. Liu

K. Shu, et al. Beyond News Contents: The Role of Social Context for Fake News Detection, WSDM'19.

# Heterogeneous Network

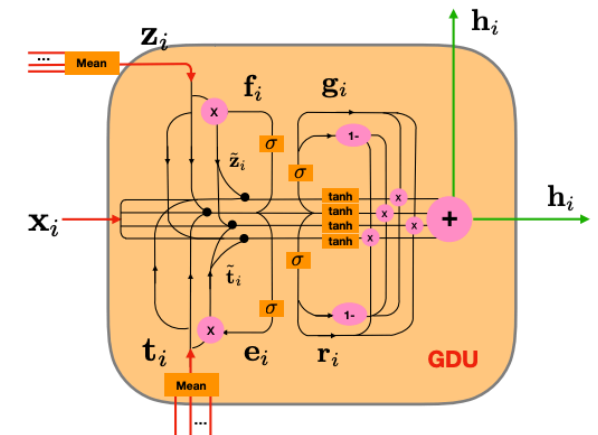*Illustrations of homogeneous networks and related studies*



Creators          News articles          Subjects
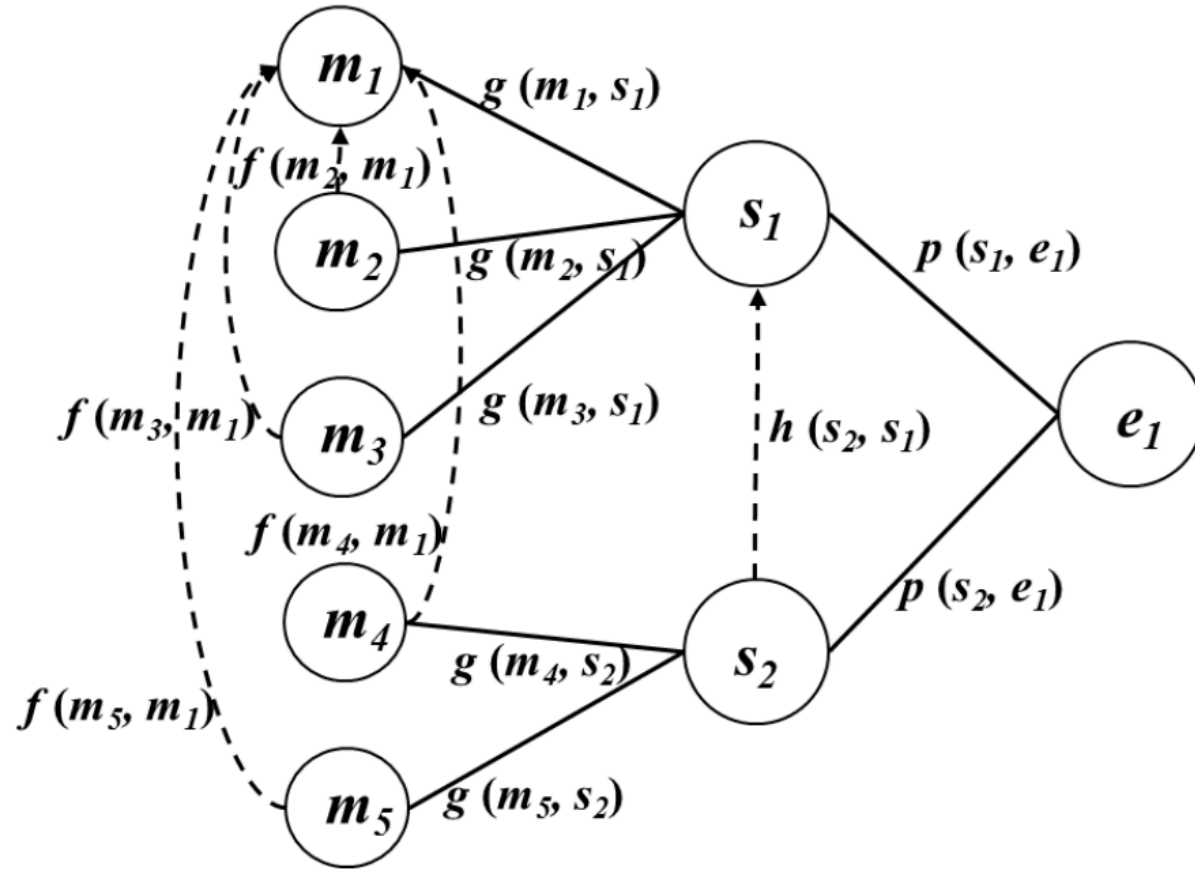
(c) Framework Architecture.



(a) Hybrid Feature Learning Unit (HFLU).

(b) Gated Diffusive Unit (GDU).

J. Zhang, et al. Fake News Detection with Deep
Diffusive Network Model, arXiv: 1805.08751, 2018

X. Zhou, R. Zafarani, K. Shu, H. Liu

66

Message Layer    Sub-event Layer    Event Layer

# Hierarchical Network

*Illustrations of hierarchical networks and related studies*



Z. Jin, et al. News Credibility Evaluation on Microblog with a Hierarchical Propagation Model, ICDM'14

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Detection

- Knowledge-based Fake News Detection
- Style-based Fake News Detection
- Propagation-based Fake News Detection
- Credibility-based Fake News Detection
- Fake News Datasets & Tools

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Credibility-based Fake News Detection

*Overview*

**Credibility-based Fake News Detection** also involve <u>social context information</u>

- <u>Credibility</u> of entities, e.g., **news headlines**, **comments** and **spreaders**
- Relationships among the <u>credibility</u> of the same or different entities

Overlaps with propagation-based fake news detection

Clickbait detection

Review spam(mer) detection

Bot detection;

X. Zhou, R. Zafarani, K. Shu, H. Liu

69

This is your brain on clickbait

intrigued   excited   disappointed   angry   depressed

approximately 3 seconds

FORTUNE.COM

# News Headline Credibility

*Clickbait*

**Clickbait** is <u>headlines</u> whose main purpose is to <u>attract the attention of visitors and encourage them to click on a link to a particular web page</u>.

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Headline Credibility

*Clickbait & Fake News*

When news articles meet clickbait:

- Attract eyeballs but are rarely newsworthy

- Increase click rate and **further gain the public trust**

| Term | Phenomenon |
|---|---|
| *Attentional bias* | **Exposure frequency** - individuals tend to believe information is correct after repeated exposures. |
| *Validity effect* | |
| *Echo chamber effect* | |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Headline Credibility

*By detecting clickbait*

**Feature engineering** within a supervised machine learning framework[8]

- N-gram and POS tags → Structure-based style features

- Informality, readability and immediacy → Attribute-based style features

- **Similarity between news headline and body-text**

News with clickbait < News without clickbait

**Deep clickbait detection**

---

[8]P. Biyani, et al., "8 Amazing Secrets for Getting More Clicks": Detecting Clickbaits in News Streams Using Article Informality . AAAI'16

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Comment Credibility

*Review Spam Detection*

- **Content-based / Style-based** models

- **Behavior-based** models

- **Graph-based** models

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Comment Credibility

*Review Spam Detection*

- **Content-based / Style-based**
- **Behavior-based** models
- **Graph-based** models

| Category | Features |
|---|---|
| **Burstiness** | Measuring the sudden promotion or descent of average rating, number of reviews, etc. for a product. This category of features emphasize on the *collective* behavior among reviewers |
| **Activity** | Measuring the total or maximum number of reviews a reviewer writes for a single product or products in a fixed time interval. This category of features emphasize on the *individual* behavior of reviewers |
| **Timeliness** | Measuring how early a product has received the review(s), or one reviewer has posted the reviews for products |
| **Similarity** | Measuring the (near) duplicate reviews written by a single reviewer or for a product, or measuring the rating deviation of one reviewer from the others for a product |
| **Extremity** | Measuring the ratio or number of extreme positive or negative reviews of a product, or for a reviewer among products |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Comment Credibility

*Review Spam Detection*

- **Content-based / Style-based** models
- **Behavior-based** models
- **Graph-based** models



**Probabilistic Graphical Models**

**Web ranking algorithm**

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Spreader Credibility

*User Classification*

User credibility score: low → high

### Malicious users

- **Intentionally** engage in fake news activities

### Susceptible users

- **Unintentionally** engage in fake news activities

### Insusceptible users

- Immune to fake news

X. Zhou, R. Zafarani, K. Shu, H. Liu

# News Spreader Credibility

*Why normal users can unintentionally engage in spreading fake news?*

| | Term | Phenomenon |
|---|---|---|
| **Social influence** | *Attentional bias* | **Exposure frequency -** individuals tend to believe information is correct after repeated exposures. |
| | *Validity effect* | |
| | *Echo chamber effect* | |
| | *Bandwagon effect* | **Peer pressure -** individuals do something primarily because others are doing it and to conform to be liked and accepted by others. |
| | *Normative influence theory* | |
| | *Social identity theory* | |
| | *Availability cascade* | |
| **Self-influence** | *Confirmation bias* | **Preexisting knowledge -** individuals tend to trust information that confirms their preexisting beliefs or hypotheses, which they perceive to surpass that of others. |
| | *Illusion of asymmetric insight* | |
| | *Naïve realism* | |
| | *Overconfidence effect* | |

**Social Influence** ➔ How widely the news article has been spread?

**Self-influence** ➔ What preexisting knowledge a user has?

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Beyond News Contents:
# The Role of Social Context for Fake News Detection
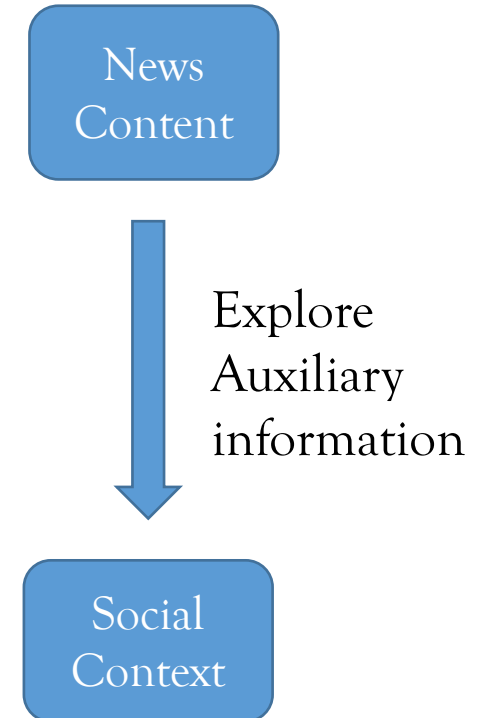
Kai Shu, Suhang Wang and Huan Liu

WSDM 2019

# Fake News Detection on Social Media - Challenges

- ## News Content
  - Fake news pieces are intentionally written to mislead users
  - Diverse in terms of topics, styles, and media platforms

- ## Social Context
  - Social engagements are massive, incomplete, unstructured, and noisy
  - Effective methods are sought to differentiate credible users, extract useful post features, and exploit network interactions

News Content

Explore Auxiliary information

Social Context

X. Zhou, R. Zafarani, K. Shu, H. Liu

79

# Fake News Detection – Multi-Source

- A typical news dissemination system on social media
  - Entities: publisher p, news a, and social media users u
  - Relations: **publishing**, spreading, social relations

> ➤ *Publishing* Publisher with partisan bias are more likely to post fake news

e.g., $p_1 \rightarrow a_1$   $p_2 \rightarrow a_3$

$p_3 \rightarrow a_4$



Publisher   News   Social Engagements

$p_1$   $a_1$   $u_1$
$p_2$   $a_2$   $u_2$   $u_3$
$p_3$   $a_3$   $u_4$
        $a_4$   $u_5$   $u_6$

Publishing   Spreading   Social Relations

> ➤ *spreading*

Low credibility users on social media are likely to share fake news, **e.g.,** $a_1 \rightarrow u_2$ $a_3 \rightarrow u_2$

> ➤ *social*

Users form relationship with like-minded people

**e.g.,** $u_2 \leftrightarrow u_4$ $u_3 \leftrightarrow u_1$

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Tri-Relationship Embedding (TriFN)

- News content embedding
  - Content modeling
  - Publisher news relation embedding
- Social Context embedding
  - Basic user feature representation
  - User news engagement modeling
- We jointly combine news content embedding and social context embedding for fake news detection

$$\min_{\mathbf{D}, \mathbf{V} \geq 0} \| \mathbf{X} - \mathbf{D}\mathbf{V}^T \|_F^2 + \lambda(\|\mathbf{D}\|_F^2 + \|\mathbf{V}\|_F^2)$$

$$\min \| \bar{\mathbf{B}}\mathbf{D}\mathbf{Q} - \mathbf{o} \|_2^2 + \lambda\|\mathbf{Q}\|_2^2$$

$$\min_{\mathbf{U}, \mathbf{T} \geq 0} \|\mathbf{Y} \odot (\mathbf{A} - \mathbf{U}\mathbf{T}\mathbf{U}^T)\|_F^2 + \lambda(\|\mathbf{U}\|_F^2 + \|\mathbf{T}\|_F^2)$$

$$\min \underbrace{\sum_{i=1}^m \sum_{j=1}^r \mathbf{W}_{ij}\mathbf{c}_i(1 - \frac{1+\mathbf{y}_{Lj}}{2})\|\mathbf{U}_i - \mathbf{D}_{L_j}\|_2^2}_{\text{True news}}$$
$$+ \underbrace{\sum_{i=1}^m \sum_{j=1}^r \mathbf{W}_{ij}(1 - \mathbf{c}_i)(\frac{1+\mathbf{y}_{Lj}}{2})\|\mathbf{U}_i - \mathbf{D}_{L_j}\|_2^2}_{\text{Fake news}}$$

X. Zhou, R. Zafarani, K. Shu, H. Liu

81

# Evaluation Setting

- Datasets: FakeNewsNet with information for news conten social context and ground truth labels from fact-checking websites
- Compared baselines:
  - RST: rhetorical relations among the words in the text
  - LIWC: lexicons falling into psycholinguistic categories
  - Castillo: features from user profiles, social networks
  - RST+Castillo
  - LIWC+Castillo

News Content + Social Context

**Table 1: The statistics of FakeNewsNet dataset**

| Platform | BuzzFeed | PolitiFact |
|---|---|---|
| # Users | 15,257 | 23,865 |
| # Engagements | 25,240 | 37,259 |
| # Social Links | 634,750 | 574,744 |
| # Candidate news | 182 | 240 |
| # True news | 91 | 120 |
| # Fake news | 91 | 120 |
| # Publisher | 9 | 91 |

News Content

Social Context

X. Zhou, R. Zafarani, K. Shu, H. Liu

82

# Evaluation Results - Detection Performance

- Social context based features are more effective than news content based features
- TriFN performs the best than other methods using both news content and social context information

Table 2: Performance comparison for fake news detection

| Datasets | Metric | RST | LIWC | Castillo | RST+Castillo | LIWC+Castillo | TriFN |
|----------|--------|-----|------|----------|--------------|---------------|-------|
| BuzzFeed | Accuracy | $0.610 \pm 0.023$ | $0.655 \pm 0.075$ | $0.747 \pm 0.061$ | $0.758 \pm 0.030$ | $0.791 \pm 0.036$ | $\mathbf{0.864 \pm 0.026}$ |
|          | Precision | $0.602 \pm 0.066$ | $0.683 \pm 0.065$ | $0.735 \pm 0.080$ | $0.795 \pm 0.060$ | $0.825 \pm 0.061$ | $\mathbf{0.849 \pm 0.040}$ |
|          | Recall | $0.561 \pm 0.057$ | $0.628 \pm 0.021$ | $0.783 \pm 0.048$ | $0.784 \pm 0.074$ | $0.834 \pm 0.094$ | $\mathbf{0.893 \pm 0.013}$ |
|          | F1 | $0.555 \pm 0.057$ | $0.623 \pm 0.066$ | $0.756 \pm 0.051$ | $0.789 \pm 0.056$ | $0.802 \pm 0.023$ | $\mathbf{0.870 \pm 0.019}$ |
| PolitiFact | Accuracy | $0.571 \pm 0.039$ | $0.637 \pm 0.021$ | $0.779 \pm 0.025$ | $0.812 \pm 0.026$ | $0.821 \pm 0.052$ | $\mathbf{0.878 \pm 0.020}$ |
|          | Precision | $0.595 \pm 0.032$ | $0.621 \pm 0.025$ | $0.777 \pm 0.051$ | $0.823 \pm 0.040$ | $0.856 \pm 0.071$ | $\mathbf{0.867 \pm 0.034}$ |
|          | Recall | $0.533 \pm 0.031$ | $0.667 \pm 0.091$ | $0.791 \pm 0.026$ | $0.792 \pm 0.026$ | $0.767 \pm 0.120$ | $\mathbf{0.893 \pm 0.023}$ |
|          | F1 | $0.544 \pm 0.042$ | $0.615 \pm 0.044$ | $0.783 \pm 0.015$ | $0.793 \pm 0.032$ | $0.813 \pm 0.070$ | $\mathbf{0.880 \pm 0.017}$ |

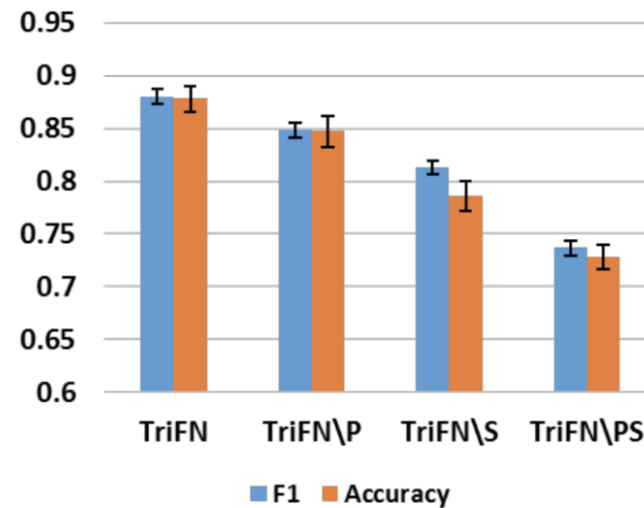News Content          Social Context          News Content + Social Context

X. Zhou, R. Zafarani, K. Shu, H. Liu

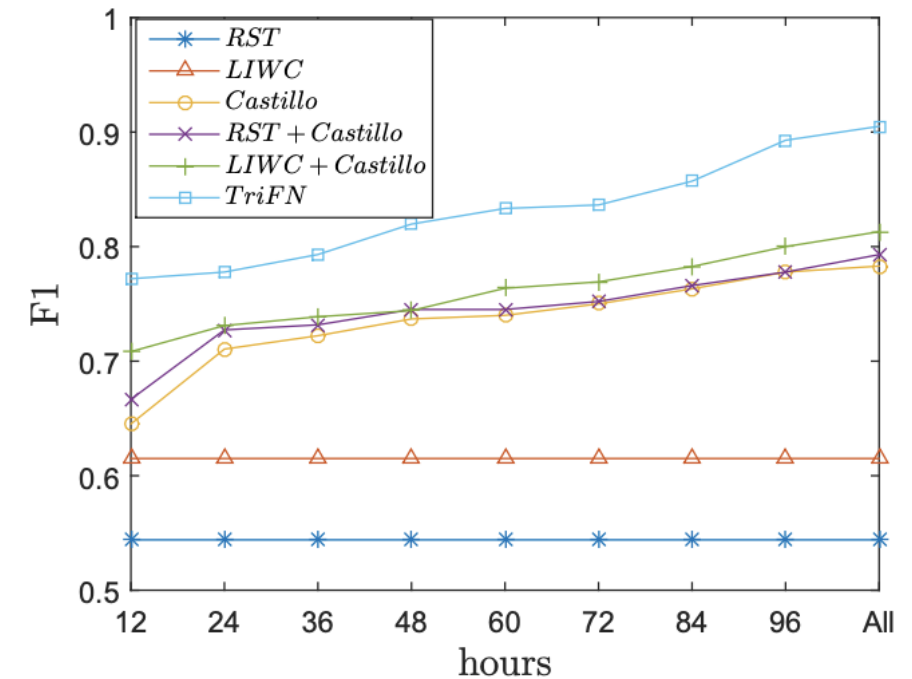# Evaluation Results - Component Analysis and Early Detection

- Both publisher-news and news-user relations can contribute to the performance improvement of TriFN

- TriFN consistently achieves best performances in the early stage of news dissemination



(a) BuzzFeed

(b) PolitiFact

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Summary

- Social context information brings additional signals to fake news detection
- It is important to capture the relations among publishers, news pieces, and users to detect fake news
- The proposed TriFN framework is effective to model tri-relationships through heterogeneous network embedding

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Unsupervised Fake News Detection:
# A Generative Approach

Shuo Yang, Kai Shu, Suhang Wang, Renjie Gu, Fan Wu, and Huan Liu

AAAI 2019

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Unsupervised Fake News Detection

- Existing methods are mainly supervised, which require extensive amount of time and labor to build a reliably annotated dataset.
- We aim to build an unsupervised fake news detection method by modeling user opinions and user credibility



Janie Johnson ✓ @jjauthor · 4 Nov 2016
Not shocking! Vote Babies!

**Pope Francis Shocks World, Endorses Donald Trump for President, Releases Statement** endingthefed.com/pope-francis-s...

💬 12    🔁 58    ♡ 46    ✉

Agreeing the authenticity of the news

iYamWhatIYam @MRIrene · 21 Oct 2016
FALSE: **Pope Francis Shocks World, Endorses Donald Trump for President**
Trumpbots getting desperate and creative. go.shr.lc/2cNK449

💬    🔁 4    ♡ 3    ✉

Doubting the authenticity of the news

X. Zhou, R. Zafarani, K. Shu, H. Liu

87

# Unsupervised Fake News Detection - challenges

- User social engagements are usually unstructured, large-scale, and noisy
- User opinions may be <span style="color:purple">conflicting</span> and <span style="color:purple">unreliable</span>, as the users usually have different degrees of credibility in identifying fake news

- The relationships among news, tweets, and users on social media form more complicated topologies
- Existing truth discovery methods mainly focus on "source-item" paths, and cannot be directly applied

X. Zhou, R. Zafarani, K. Shu, H. Liu

# The hierarchical user engagement structure

- We build a hierarchical user engagement structure for each news
  - $x_i$ is a random variable denoting the label of $news_i$
  - $y_{i,j}$ denotes the opinion with sentiment of verified user $j$ to $news_i$
  - $z_{i,j,k}$ is the opinion of unverified user $k$ to $news_i$
    - Like: opinion same with $y_{i,j}$
    - Reply: sentiment score of the reply
    - Retweet: opinion same with $y_{i,j}$

Verified User

Unverified User

X. Zhou, R. Zafarani, K. Shu, H. Liu

# The Proposed Probabilistic Model (UFD)

- For each news $i$, $x_i$ is generated from Bernoulli distribution

$$x_i \sim \text{Bernoulli}(\theta_i)$$

- For verified user $j$      $y_{i,j} \sim \text{Bernoulli}(\phi_j^{x_i})$

  - $\phi_j^1$ ($\phi_j^0$) the probability that the user $j$ thinks a news piece is real given the truth estimation of the news is true and fake

- For unverified $k$,    $z_{i,j,k} \sim \text{Bernoulli}(\psi_k^{x_i,y_{i,j}})$

  - the opinion is likely to be influenced by the news itself and the verified users' opinions

$$\psi_k^{0,0} := p(z_{i,j,k} = 1 | x_i = 0, y_{i,j} = 0)$$
$$\psi_k^{0,1} := p(z_{i,j,k} = 1 | x_i = 0, y_{i,j} = 1)$$
$$\psi_k^{1,0} := p(z_{i,j,k} = 1 | x_i = 1, y_{i,j} = 0)$$
$$\psi_k^{1,1} := p(z_{i,j,k} = 1 | x_i = 1, y_{i,j} = 1)$$



X. Zhou, R. Zafarani, K. Shu, H. Liu

90

# Evaluation Results - Detection Performance

- Majority voting achieves the worst performance since it equally aggregates the users' opinions without considering user's credibility degree
- The proposed framework UFD can achieve best performance comparing with other unsupervised truth discovery methods
- We can also discover the top-k creidible users, and these users are mostly expert journalists, professional news reporters

Table 2: Performance comparison on LIAR dataset

| Methods | Accuracy | True | | | Fake | | |
|---|---|---|---|---|---|---|---|
| | | Precision | Recall | F1-score | Precision | Recall | F1-score |
| Majority Voting | 0.586 | 0.624 | 0.628 | 0.626 | 0.539 | 0.534 | 0.537 |
| TruthFinder | 0.634 | 0.650 | 0.679 | 0.664 | 0.615 | 0.583 | 0.599 |
| LTM | 0.641 | 0.654 | 0.691 | 0.672 | 0.624 | 0.583 | 0.603 |
| CRH | 0.639 | 0.653 | 0.687 | 0.669 | 0.621 | 0.583 | 0.601 |
| **UFD** | **0.759** | **0.766** | **0.783** | **0.774** | **0.750** | **0.732** | **0.741** |

Table 4: Top accurate verified users on two datasets

| User | Accuracy | Sensitivity | Specificity |
|---|---|---|---|
| amy_hollyfield | 1.0 | 1.0 | 1.0 |
| politico | 0.909 | 0.833 | 1.0 |
| loujacobson | 0.84 | 0.842 | 0.833 |
| dcexaminer | 0.833 | 0.818 | 0.857 |
| FoxNews | 0.818 | 0.714 | 1.0 |

# Summary

- We study the novel problem of unsupervised fake news detection, a much desired scenario in the real world
- We propose a probabilistic model to consider the user opinions and user credibility in a hierarchical engagement structure
- We demonstrate the effectiveness of the proposed framework in real-world datasets
- **Future work**
  - Incorporating user profiles and news contents into unsupervised models
  - Building semi-supervised models with limited engagements information

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Deep Headline Generation for Clickbait Detection

Kai Shu, Suhang Wang, Thai Le, Dongwon Lee, and Huan Liu

ICDM 2018

# Clickbaits

- Clickbaits are catchy social media posts or sensational headlines that attempt to lure the readers to click



- Clickbaits can have negative societal impacts
  - clickbaits may contain sensational and inaccurate information to mislead readers and spread fake news
  - clickbaits may be used to perform clickjacking attacks by redirecting users to phishing websites

# Clickbait Detection

- Existing approaches mainly focus on extracting hand-crafted linguistic features (as traditionally done so) or building sophisticated predictive models such as deep neural networks
- However, these methods may face following limitations
  - Scale: datasets with labels are often limited
  - Distribution: imbalanced distribution of clickbaits and non-clickbaits

> We aim to generate synthetic headlines with specific styles and exploit the utility to improve clickbait detection

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Headline Generation from Documents

- Goal: Generate stylized headlines that also preserve document contents

Document → Clickbait headline

Non-Clickbait headline

- Stylized headlines can help augment training data for clickbait detection
- Content preserved headlines make it possible to suggest a non-clickbait headline to readers after we detect a clickbait

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Problem Definition

- Let $\{x_1, x_2, \ldots, x_m\}$ $\{h_1, h_2, \ldots, h_m\}$ and $\{y_1, y_2, \ldots, y_m\}$ denote the set of $m$ documents, and corresponding headlines and labels

- Giving $S = \{(x_i, h_i) | i = 1, \ldots, m\}$, learn a generator that can generate stylized headlines given a document and a style label, i.e., $o_i = f(x_i, y_i)$

- Challenges
  - How to generate realistic and readable headlines from original documents?
  - How to utilize generated headlines to augment training data for clickbait detection
  - How to generate new headlines that can preserve the content of documents and transfer the style of original headlines

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Stylized Headline Generation (SHG)

- We propose a deep learning model to generate both click-baits and non-clickbaits with style transfer
  - Generator Learning: a document autoencoder $A$, a headline generator $G$
  - Discriminator Learning: a transfer discriminator $D_T$, a style discriminator $D_S$, a pair discriminator $D_P$

# Generator Learning

- Document autoencoder $A$ extract document representation by minimizing the reconstruction error

$$\mathcal{L}_{rec}(\theta_e, \theta_d) = -\sum_{i=1}^{m} \log p(\hat{x}_i | x_i; \theta_d, \theta_e)$$

- Headline generator $G$
  - Generate stylized headline by minimizing the reconstruction error of original headline

$$\mathcal{L}_G(\theta_G) = \mathbb{E}_{(x,h) \in \mathcal{S}} [-\log p_G(h | \mathbf{y}^L, \mathbf{z}))]$$

- Generate a set of new headlines $O$ with the styles $y^U$ opposite to the original headlines

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Discriminator Learning

- Discriminators regularize the representation learning of document $z$, original headline $s_N$, and generated headline $\tilde{s_N}$

- Transfer discriminator $D_T$ : discriminate original data samples with generated data samples

<span style="color:orange">Original clickbaits and generated non-clickbaits</span>

$$\mathcal{L}_{D_T} = \boxed{\mathcal{L}_{D_T^{(1)}}(\theta_{D_T^{(1)}})} + \boxed{\mathcal{L}_{D_T^{(2)}}(\theta_{D_T^{(2)}})}$$

<span style="color:purple">Original non-clickbaits and generated clickbaits</span>

- Style discriminator $D_S$: assign a correct label of styles for both original headlines and generated headlines

<span style="color:orange">Original clickbaits and original non-clickbaits</span>

$$\mathcal{L}_{D_S}(\mathbf{W}, \mathbf{b}) = \boxed{\mathcal{L}_{D_S}^{(1)}} + \boxed{\mathcal{L}_{D_S}^{(2)}}$$

X. Zhou, R. Zafarani, K. Shu, H. Liu

<span style="color:purple">Generated clickbaits and generated non clickbaits</span>

100

# Discriminator Learning

- Pair discriminator $D_P$ ensures that the correspondences of documents and headlines are maintained

Proximity function $\quad p(h_i, x_j) = \dfrac{1}{1 + \exp(-\mathbf{s}^{(i)}\mathbf{Q}\mathbf{z}^{(j)})}$

Document representation

Headline representation

- Maximizing the proximity of (document, headline) pairs with negative sampling

$$\mathcal{L}_{D_P} = -\log \sigma(\mathbf{s}^{(i)}\mathbf{Q}\mathbf{z}^{(i)}) - \sum_{k=1}^{K} \mathbb{E}_{x_k \sim P_n(x)}\left[\log \sigma(-\mathbf{s}^{(i)}\mathbf{Q}\mathbf{z}^{(k)})\right]$$

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Experiments Setting

**TABLE I: The statistics and descriptions of the datasets**

| Dataset | Source | # Clickbaits | # Non-clickbaits |
|---|---|---|---|
| $P$ | **Professional Writers** | 5,000 | 16,933 |
| $M$ | **Social Media Users** | 4,883 | 16,150 |

- Datasets
  - Professional writers (P):

    Reporters or editors generate clickbaits for their news pieces
  - Social media users (M):

    Clickbaits to lure people to click their posts on social media.
- Baselines
  - SeqGAN [AAAI'17] : Text generation using GAN with reinforcement learning
  - SVAE [CONLL'16]: Sentence generation using Variational AutoEncoder (VAE)
  - CrossA [NIPS'17]: Generating sentences across different styles

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Experiments - Evaluation questions

- **Consistency**: are generated clickbaits/non-clickbaits consistent with the original datasets?
- **Readability**: are generated headlines readable or not?
- **Similarity**: are generated headlines semantically similar to original documents?

Data Quality

- **Differentiability**: are generated clickbaits/non-clickbaits differentiable?
- **Accuracy**: can generated clickbaits/non-clickbaits help improve the detection performance?

Data Utility

# Experimental Results - Data Quality

- **Similarity**: evaluate the semantic similarity of headlines and documents
  - Bilingual Evaluation Understudy (BLEU) score
  - Uni_sim: similarity of universal text embedding
- SHG achieves better performances to preserve document content than CrossA

TABLE V: **EQ3**: The Average BLEU (BLEU-4) Score Comparison of Generated Headlines. $\mathcal{H}$ indicates original headlines, and $\mathcal{O}$ represents the generated headlines.

| Data | Headlines | Methods | Clickbait | Non-Clickbait |
|------|-----------|---------|-----------|---------------|
| P | $\mathcal{H}$ | | 0.555 | 0.527 |
| | $\mathcal{O}$ | CrossA | 0.407 | 0.432 |
| | | SHG | **0.453** | **0.446** |
| M | $\mathcal{H}$ | | 0.541 | 0.534 |
| | $\mathcal{O}$ | CrossA | 0.432 | 0.437 |
| | | SHG | **0.451** | **0.442** |

TABLE VI: **EQ3**: The Average Uni_sim Value Comparison of Generated Headlines. $\mathcal{H}$ indicates original headlines, and $\mathcal{O}$ represents the generated headlines.

| Data | Headlines | Methods | Clickbait | Non-Clickbait |
|------|-----------|---------|-----------|---------------|
| P | $\mathcal{H}$ | | 0.63 | 0.81 |
| | $\mathcal{O}$ | CrossA | 0.20 | 0.22 |
| | | SHG | **0.37** | **0.40** |
| M | $\mathcal{H}$ | | 0.64 | 0.81 |
| | $\mathcal{O}$ | CrossA | 0.26 | 0.34 |
| | | SHG | **0.34** | **0.38** |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Experimental Results - Data Utility

- **Accuracy**: improvement comparison of original headlines on AUC
  - The headlines generated by SVAE, CrossA, and SHG can increase the performance of clickbait detection to some extent
  - SHG consistently outperforms SVAE and CrossA

| Data | Classifier | Org | SeqGAN | SVAE | CrossA | SHG |
|------|-----------|-----|--------|------|--------|-----|
| *P* | LogReg | 0.928 | 0.900 (↓ 3.02%) | 0.933 (↑ 0.54%) | 0.932 (↑ 0.64%) | **0.936 (↑ 0.86%)** |
| | DTree | 0.894 | 0.882 (↓ 1.34%) | 0.908 (↑ 1.57%) | 0.900 (↑ 0.67%) | **0.910 (↑ 1.79%)** |
| | RForest | 0.900 | 0.893 (↓ 0.78%) | 0.912 (↑ 1.33%) | 0.916 (↑ 1.78%) | **0.925 (↑ 2.78%)** |
| | XGBoost | 0.919 | 0.914 (↓ 0.54%) | 0.923 (↑ 0.43%) | 0.926 (↑ 0.76%) | **0.928 (↑ 0.98%)** |
| | AdaBoost | 0.917 | 0.896 (↓ 2.29%) | 0.921 (↑ 0.44%) | 0.921 (↑ 0.44%) | **0.931 (↑ 1.64%)** |
| | SVM | 0.904 | 0.898 (↓ 0.66%) | 0.917 (↑ 1.44%) | 0.920 (↑ 1.77%) | **0.923 (↑ 2.10%)** |
| | GradBoost | 0.921 | 0.914 (↓ 0.76%) | 0.924 (↑ 0.33%) | 0.926 (↑ 0.54%) | **0.928 (↑ 0.76%)** |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Summary

- We study the problem of generating clickbaits/nonclickbaits from original documents for clickbait detection
- We propose a novel deep generative model with adversarial learning

- **Future work**
  - Explore the generalization capacity of SHG on other styles such as positive-negative sentiment style and academic-news reporting style
  - Investigate the strategy of learning the disentangled representations of content and style

# Summary and Comparison for Fake News Detection

| | **Knowledge-based** fake news detection | **Style-based** fake news detection | **Propagation-based** fake news detection | **Credibility-based** fake news detection |
|---|---|---|---|---|
| Information Utilized | News content | | News content & Social context information | |
| Techniques | Graph models | Feature-based methods | Graph models & Feature-based methods | |
| Resources | Knowledge graphs | Fundamental theories | | |
| Related Topic(s) | Fact-checking | Deception detection | Rumor detection | Clickbait/bot/review spam detection |

# Fake News Detection

- Knowledge-based Fake News Detection
- Style-based Fake News Detection
- Propagation-based Fake News Detection
- Credibility-based Fake News Detection
- Fake News Datasets & Tools

X. Zhou, R. Zafarani, K. Shu, H. Liu

# FakeNewsNet: A Data Repository with News Content, Social Context and Dynamic Information for Studying Fake News on Social Media

Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, Huan Liu



kaggle  Search

Dataset

**FakeNewsNet**
Fake News, MisInformation, Data Mining

https://github.com/KaiDMML/FakeNewsNet

https://www.kaggle.com/mdepak/fakenewsnet

# How unique is FakeNewsNet?

- A comprehensive data repository that contains news contents, social context, and spatiotemporal information

Table 1: Comparison with existing fake news detection datasets

| Features / Dataset | News Content | | Social Context | | | | Spatiotemporal Information | |
|---|---|---|---|---|---|---|---|---|
| | Linguistic | Visual | User | Post | Response | Network | Spatial | Temporal |
| BuzzFeedNews | ✓ | | | | | | | |
| LIAR | ✓ | | | | | | | |
| BS Detector | ✓ | | | | | | | |
| CREDBANK | ✓ | | ✓ | ✓ | | | ✓ | ✓ |
| BuzzFace | ✓ | | | ✓ | ✓ | | | ✓ |
| FacebookHoax | ✓ | | ✓ | ✓ | ✓ | | | |
| **FakeNewsNet** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Data Integration

- **News Content:** we utilize fact-checking websites

  to obtain news

  contents for fake news and true news
- **Social Context:** collecting user

  engagements from Twitter using

  the headlines of news articles
- **Spatiotemporal Information:** spatial info

  and temporal data

  from meta data of Twitter



X. Zhou, R. Zafarani, K. Shu, H. Liu

111

# Data Analysis

- **User profiles:** users who share real news pieces tend to have longer register time than those who share the fake news on average

- **User engagements:** fake news pieces tend to have fewer replies and more retweets; real news pieces have more ratio of likes than fake news pieces do



X. Zhou, R. Zafarani, K. Shu, H. Liu

- **A case study of temporal engagements for fake news and real news**
  - For fake news, a sudden increase in the number of retweets and remain constant beyond a short time
  - For real news, the number of retweets increases steadily
  - Fake news pieces tend to receive fewer replies than real news



Fake News

Real News

X. Zhou, R. Zafarani, K. Shu, H. Liu

113

# Potential Applications for FakeNewsNet

- **Fake News Detection**
  - News content, social context based
  - Early fake news detection
- **Fake News Evolution**
  - Temporal, Topic, Network, evolution
- **Fake News Mitigation**
  - Provenances, persuaders, clarifiers
  - Influence minimization, mitigation campaign
- **Malicious Account Detection**
  - Detecting bots that spread fake news

# FakeNewsTracker: A Tool for Fake News Collection, Detection, and Visualization

## Kai Shu, Deepak Mahudeswaran, and Huan Liu

## SBP 2018

SBP Disinformation Challenge Winner

http://blogtrackers.fulton.asu.edu:3000

X. Zhou, R. Zafarani, K. Shu, H. Liu

# An end-to-end framework for fake news collection, detection, and visualization

- **Data Collection:** collecting fake and real news articles from fact-checking websites and related social engagements from social media

- **Fake News Detection:** finding fake news with advanced machine learning methods, such as deep neural networks

- **Fake News Visualization**: visualization on data attributes and model performance



X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Detection

- Detect fake news with fusion of news content and social context
    - **News representation:**
      Represent news content using autoencoders
    - **Social engagement representation:**
      Represent social engagements using RNNs
    - **Social Article Fusion:**
      Combine both news and social engagement features to detect fake news



X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Visualization



Trends on Twitter



Topics of Fake news vs Real News



Geolocation of Fake News vs Real News



Social Network on Users Spreading Fake/Real news

# Recent work at DMML on Fake News Detection

- Survey: Fake News Detection on Social Media: A Data Mining Perspective
- Data repository: FakeNewsNet, [Github], [Kaggle], [Paper]
- Software: FakeNewsTracker
- Book chapter: Studying Fake News via Network Analysis: Detection and Mitigation
- Other Publications: related publications are updated at:

  http://www.public.asu.edu/~skai2/

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Challenges and Highlights

- Fake News Early Detection
- Identify Check-worthy Content
- Cross-domain, -topic, -language Fake News Studies
- Deep Learning for Fake News Studies

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Early Detection

*Why is Fake News Early Detection is important?*

- The more fake news spreads, the more likely for people to trust it

- Once people have trusted the fake news, it is difficult to correct users' perceptions

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

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

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Fake News Early Detection

*How to achieve Fake News Early Detection?*

I. **Verification Efficiency**, e.g., compare knowledge in the framework that

- Knowledge graphs with timely ground truth
- To-be-verified news content is check-worthy – *Check-worthy content identification*

II. **Feature Compatibility**, e.g., to extract features that can capture

- The <u>generality</u> of deceptive content styles *across* domain, topic, and language[9]
- The <u>evolution</u> of deceptive content styles *within* domain, topic, and language

III. **Information Availability**, e.g., detect fake news with limited propagation information

[9]W. Yaqing, et al., EANN: Event Adversarial Neural Networks for Multi-Modal Fake News Detection. KDD'18

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Check-worthy Content Identification

*How to measure* *Check-worthy* *content?*

I.   **News-worthiness or Potential Influence on the Society,** e.g., if it is related to national affairs

II.  **Spammer Preference**, i.e., news historical likelihood of being fake
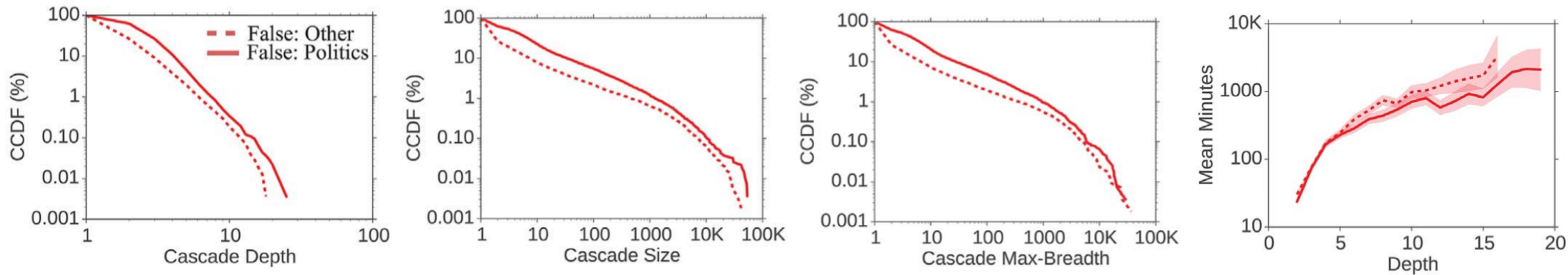


(a) (Expert-based) PolitiFact: the PolitiFact scorecard

*Related Studies:*
- N. Hassan, et al. Detecting Check-worthy Factual Claims in Presidential Debates, CIKM'15
- N. Hassan et al., Toward Automated Fact-Checking: Detecting Check-worthy Factual Claims by ClaimBuster, KDD'17

X. Zhou, R. Zafarani, K. Shu, H. Liu

123

# Cross-domain, -topic, -language

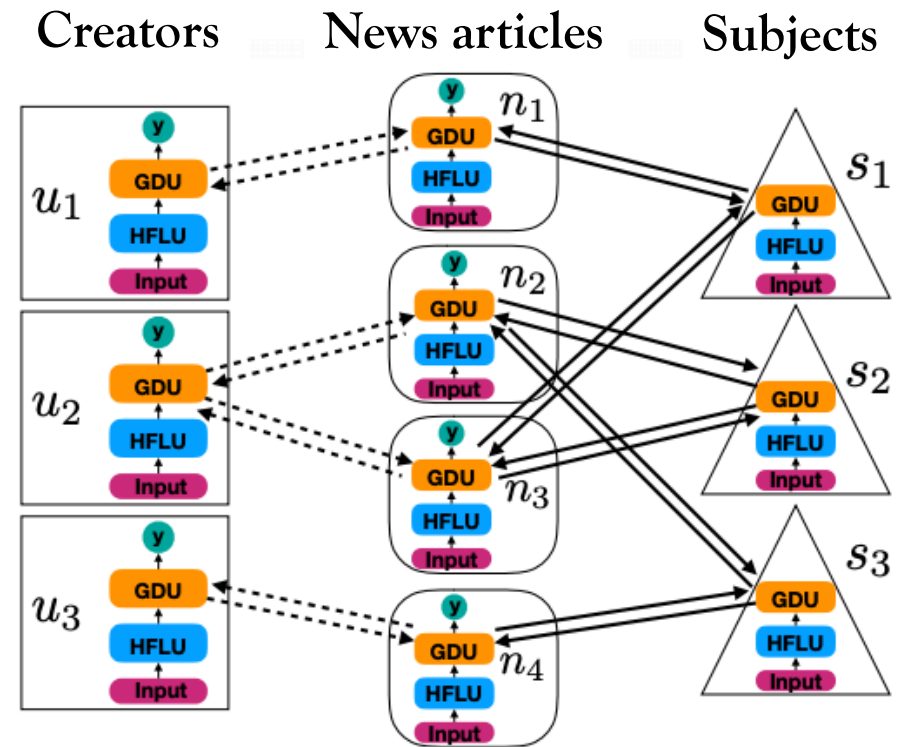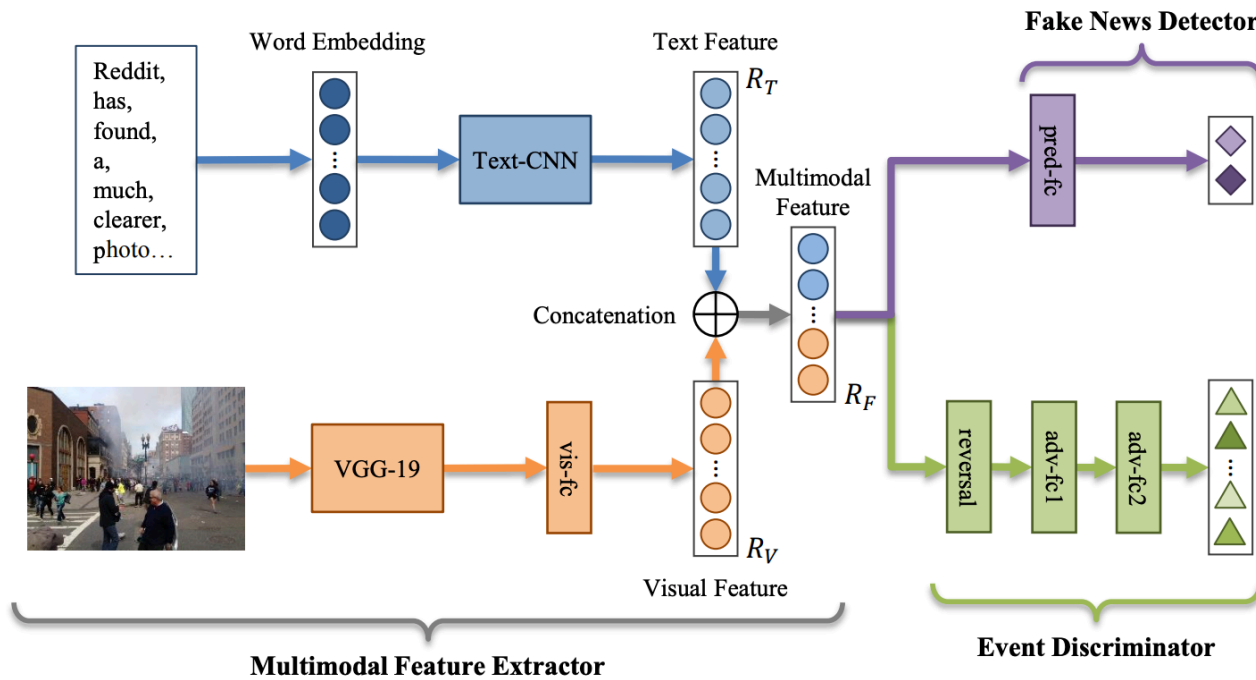*How to facilitate Cross-domain, -topic, -language Fake News Studies?*

I.   Develop **fake news datasets** containing cross-domain, -topic, -language data

II.  Explore **patterns** among fake news within different domains, topics and languages



III. Develop **techniques** enables cross-domain, -topic, -language fake news detection

Figures are from: S. Vosoughi, et al. The spread of true and false news online. Science, 2018

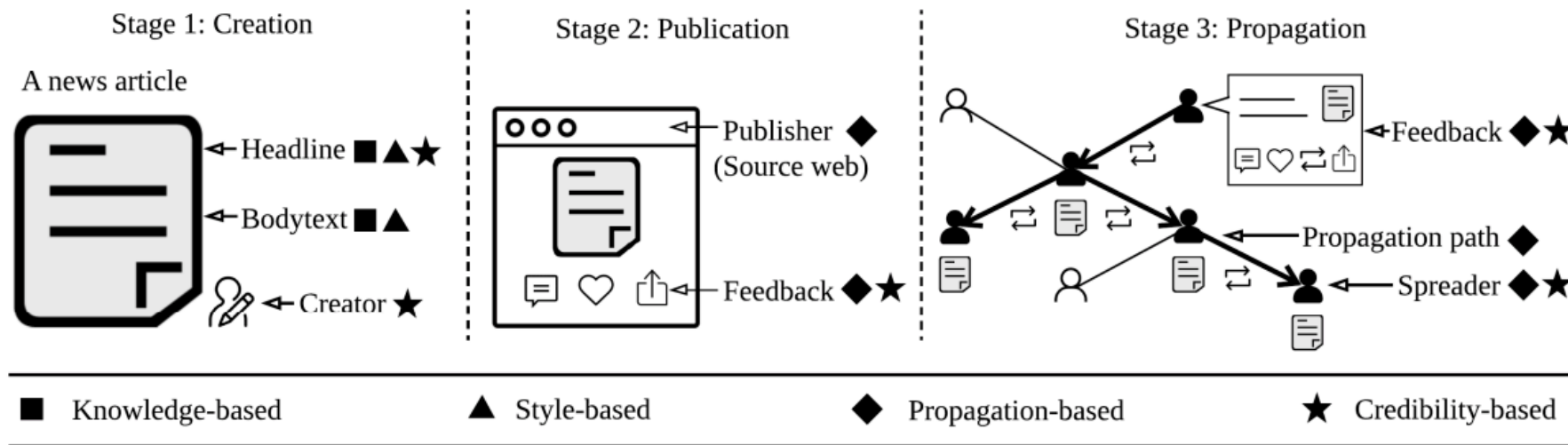# Deep Learning for Fake News Detection



W. Yaqing, et al., EANN: Event Adversarial Neural Networks for Multi-Modal Fake News Detection. *KDD'18*

J. Zhang, et al. Fake News Detection with Deep Diffusive Network Model, arXiv: 1805.08751, 2018

X. Zhou, R. Zafarani, K. Shu, H. Liu

# Summary

I.   **Fundamental Theories** encourage <u>interdisciplinary research</u> of fake news

II.  **Fake News Detection** from various perspectives



III. **Challenges and Highlights** for potential research opportunities for fake news studies

X. Zhou, R. Zafarani, K. Shu, H. Liu