



Modeling and Countering Misinformation in Adversarial Information Ecosystems

Arash Amini, Yigit Ege Bayiz, and Ufuk Topcu



AFOSR Center of Excellence in Assured Autonomy in Contested Environments

Defense for Special Operations and Low-Intensity Conflict, stated that adversarial use of *disinformation*, *misinformation*, and *propaganda* poses one of today's *greatest challenges* to the United States.

autonomy

Defense for Special Operations and Low-Intensity Conflict, stated that adversarial use of *disinformation*, *misinformation*, and *propaganda* poses one of today's *greatest challenges* to the United States.



autonomy

Defense for Special Operations and Low-Intensity Conflict, stated that adversarial use of *disinformation*, *misinformation*, and *propaganda* poses one of today's *greatest challenges* to the United States.



Defense for Special Operations and Low-Intensity Conflict, stated that adversarial use of *disinformation*, *misinformation*, and *propaganda* poses one of today's *greatest challenges* to the United States.



Defense for Special Operations and Low-Intensity Conflict, stated that adversarial use of *disinformation*, *misinformation*, and *propaganda* poses one of today's *greatest challenges* to the United States.







The OODA loop in theory appears as a consistent cycle of observe, orient, decide, and act. However, operational and adversarial inputs can distort this process in a number of ways.

How Disinformation Works





The OODA loop in theory appears as a consistent cycle of observe, orient, decide, and act. However, operational and adversarial inputs can distort this process in a number of ways.



How Disinformation Works





The OODA loop in theory appears as a consistent cycle of observe, orient, decide, and act. However, operational and adversarial inputs can distort this process in a number of ways.



How Disinformation Works

CENTER FOR autonomy



The OODA loop in theory appears as a consistent cycle of observe, orient, decide, and act. However, operational and adversarial inputs can distort this process in a number of ways.



















Misinformation







Misinformation



Malinformation













NEWS | Sept. 3, 2024

Russian Disinformation Campaign "DoppelGänger" Unmasked: A Web of Deception

















Initiation

Amini, Bayiz, and Topcu



One-Third Believe Or Are Unsure About Four Or More of Eight False Statements About COVID-19





CENTER FOR autonomy





Amini, Bayiz, and Topcu





CENTER FOR autonomy





Amini, Bayiz, and Topcu

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to Al.* Signal, 2024.

CENTER FOR

autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian



Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders

CENTER FOR

autonomy

¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to Al.* Signal, 2024.

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian



Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



CENTER FOR

autonomy

¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to Al.* Signal, 2024.

autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian



Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to AI.* Signal, 2024.

autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian



Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to Al.* Signal, 2024.

CENTER FOR autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian

0⁰0 000

Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to AI.* Signal, 2024.

CENTER FOR autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian

0⁰0 000

Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to AI.* Signal, 2024.

autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian

0⁰0 000

Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to AI.* Signal, 2024.

autonomy

- Totalitarian: centralized information management
- Democratic: distributed information management



Totalitarian

0⁰0 000

Democratic

Feasible Solutions

- Democratic
- Resilient
- Safe
- Applicable
- Considers all stakeholders



¹ Harari, Yuval Noah. *Nexus: A brief history of information networks from the Stone Age to Al.* Signal, 2024.

Modeling Information Environment

CENTER FOR autonomy

Role of Media Competition in Spreading of Misinformation			
Misinformation	Misinformation	Misinformation	
Initiation	Internalization	Reinforcement	

¹ Amini, Arash, et al. *How Media Competition Fuels the Spread of Misinformation*. Science Advances [Under review]

Amini, Bayiz, and Topcu
CENTER FOR autonomy

Role of Media Competition	on in Spreading of Misinformation		
Misinformation	Misinformation	Misinformation	
Initiation	Internalization	Reinforcement	

News Source Credibility

• What is the chance of a news source sharing misinformation?

autonomy

Role of Media Competition	in Spreading of Misinformation		
Misinformation	Misinformation	Misinformation	
Initiation	Internalization	Reinforcement	

News Source Credibility

• What is the chance of a news source sharing misinformation?

User Susceptibility

• Chance of a user getting influence under exposure to information from a low credibility source

¹ Amini, Arash, et al. How Media Competition Fuels the Spread of Misinformation. Science Advances [Under review]

autonomy

Misinformation Misinformation Misinformation	
Initiation Internalization Reinforcement	

News Source Credibility

• What is the chance of a news source sharing misinformation?

User Susceptibility

• Chance of a user getting influence under exposure to information from a low credibility source

Sharing misinformation: gain attention loose integrity

¹ Amini, Arash, et al. How Media Competition Fuels the Spread of Misinformation. Science Advances [Under review]

CENTER FOR autonomy

Role of Media Competit	ion in Spreading of Misinformation		
Misinformation	Misinformation	Misinformation	
Initiation	Internalization	Reinforcement	

News Source Credibility

• What is the chance of a news source sharing misinformation?

User Susceptibility

• Chance of a user getting influence under exposure to information from a low credibility source

Sharing misinformation: gain attention loose integrity



High-Credibility News Source(S) Disseminates Misinformation



Low-Credibility News Source(S) Disseminates Misinformation



High-Credibility News Source(S) Disseminates Factual Information



Low-Credibility News Source(S) Disseminates Factual Information

CENTER FOR autonomy

Role of Media Competition in	Spreading of Misinformation		
Misinformation	Misinformation	Misinformation	
Initiation	Internalization	Reinforcement	

News Source Credibility

• What is the chance of a news source sharing misinformation?

User Susceptibility

Chance of a user getting influence under exposure to information from a low credibility source

Sharing misinformation: gain attention loose integrity

How users beliefs change under exposure to misinformation?

$$x_{t+1}^{i} = \frac{a}{A_{t}^{i}} \sum_{j=1}^{N} \phi(|x_{t}^{i} - x_{t}^{j}|)(x_{t}^{j} - x_{t}^{i}) + \frac{b}{B_{t}^{i}} \sum_{m=1}^{M} \psi(|x_{t}^{i} - y^{m}|, c_{t}^{m}, a_{t}^{m}, s^{i})(y^{m} - x_{t}^{i}) + \sigma w_{t}^{i},$$

Social influence

Media influence



High-Credibility News Source(S) Disseminates Misinformation



Low-Credibility News Source(S) Disseminates Misinformation



High-Credibility News Source(S) Disseminates Factual Information



Low-Credibility News Source(S) Disseminates Factual Information

¹ Amini, Arash, et al. How Media Competition Fuels the Spread of Misinformation. Science Advances [Under review]

Amini, Bayiz, and Topcu

Ĵ





Source Credibility



¹ Amini, Arash, et al. How Media Competition Fuels the Spread of Misinformation. Science Advances [Under review]

CENTER FOR autonomy

В

Credibility(c)

0

D

Misinformation Exposure (Γ)

0-

-1

-1

0

Opinion(x)

0

Opinion(x)

Population

40

· 30

· 20

- 10

n

400

t = 400

0

1 -1



CENTER FOR autonomy



CENTER FOR **autonomy**



Gain Influence Loose Credibility

CENTER FOR autonomy



¹ Amini, Arash, et al. How Media Competition Fuels the Spread of Misinformation. Science Advances [Under review]

Gain Influence Loose Credibility

CENTER FOR autonomy



Arm Race

If a player **increases misinformation dissemination** below equilibrium, the **optimal response** of the other player leads to greater misinformation sharing.



Arm Race

If a player increases misinformation dissemination below equilibrium, the optimal response of the other player leads to greater misinformation sharing.



12

10

6 Credibility gain (ζ) ò

Opinion(x)

Misinformation gain (ŋ) 3 –

2 ·

0.

С

Misinformation gain (ŋ)

3

0 -

0

0

Arm Race

If a player **increases misinformation dissemination** below equilibrium, the **optimal response** of the other player leads to *greater misinformation sharing*.









• **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.
- Arms Race Dynamics: Reducing credibility can trigger a misinformation escalation, while credible biased sources encourage factual reporting.



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.
- Arms Race Dynamics: Reducing credibility can trigger a misinformation escalation, while credible biased sources encourage factual reporting.

Conclusion



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.
- Arms Race Dynamics: Reducing credibility can trigger a misinformation escalation, while credible biased sources encourage factual reporting.

Conclusion

• Misinformation is a *dynamic* problem.



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.
- Arms Race Dynamics: Reducing credibility can trigger a misinformation escalation, while credible biased sources encourage factual reporting.

Conclusion

- Misinformation is a *dynamic* problem.
- Is this outcome the *fault* of the media or the information ecosystem?



- **Debunking Challenges**: While fact-checking reduces misinformation, its impact depends on source credibility and audience alignment.
- **Source Accountability**: Increasing credibility gain by holding sources accountable reduces misinformation and polarization.
- **Community Resilience**: Lowering susceptibility through education and media literacy fosters a less polarized, more informed society.
- Arms Race Dynamics: Reducing credibility can trigger a misinformation escalation, while credible biased sources encourage factual reporting.

Conclusion

- Misinformation is a *dynamic* problem.
- Is this outcome the *fault* of the media or the information ecosystem?
- The current information enviromen is explotaible

autonomy

Points of intervention

Misinformation initiation

Misinformation internalization

Misinformation reinforcement

CENTER FOR autonomy

Points of intervention

Misinformation initiation

Misinformation internalization

Misinformation reinforcement

Option 1: Debunking

Identify misinformed users.

Correct internalized misinformation by debunking their perceived beliefs.

autonomy

Points of intervention



CENTER FOR autonomy

Points of intervention



autonomy

Points of intervention



¹ van der Linden, S. *Misinformation: susceptibility, spread, and interventions to immunize the public.* Nature Medicine, 2022

Countering Misinformation: Pre-bunking

autonomy

Points of intervention







Overusing pre-bunking negatively affects user experience

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025





Overusing pre-bunking negatively affects user experience



Multiple misinformation sources

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025







Overusing pre-bunking negatively affects user experience

Goal 1: Ensure pre-bunks are delivered before misinformation.

Eg: $C_1, P_1, C_2, C_1, M_1, C_2, P_2, M_2, C_4, \dots$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025





Content Sequence: $Q = q_1, q_2, q_3, ...$ Misinformation Misc. Content $q_i \in \mathbf{M} \cup \mathbf{P} \cup \mathbf{C}$ Pre-bunk Example: $Q = C_1, C_2, C_1, M_1, C_2, M_2, C_4, ...$ Multiple misinformation sources

Overusing pre-bunking negatively affects user experience

Goal 1: Ensure pre-bunks are delivered before misinformation.

Goal 2: Minimize pre-bunk concentration in feed.

Eg: $C_1, P_1, C_2, C_1, M_1, C_2, P_2, M_2, C_4, \dots$

$$\min_{a} \max_{i} c_{i+1} = \beta c_i + \mathbf{1}_{\{q_i \in \mathbf{P}\}}$$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025

autonomy

Each M_i propagates based on an **SI** epidemic model.

At each time we can decide on which pre-bunks to deliver based on policy $\pi(I_t^1, \ldots, I_t^k)$.



¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025
CENTER FOR autonomy

Each M_i propagates based on an **SI** epidemic model.

At each time we can decide on which pre-bunks to deliver based on policy $\pi(I_t^1, \ldots, I_t^k)$.



$$\begin{array}{ll} \min\max_{\pi\in\Pi} & c(\pi,t) \\ \text{s.t.} & A_1, \dots, A_T \sim \pi \left(I_t^1 \dots I_t^k \right), \\ & \forall k \in \mathbb{N}, \quad \mathbb{P}\left(P_k \in \bigcup_{t=0}^{X_k-1} A_t \right) = 1. \end{array}$$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025

Each M_i propagates based on an **SI** epidemic model.

At each time we can decide on which pre-bunks to deliver based on policy $\pi(I_t^1, \ldots, I_t^k)$.



¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



$$\forall k \in \mathbb{N}, \quad \mathbb{P}\left(P_k \in \bigcup_{t=0}^{X_k - 1} A_t\right) = 1.$$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



$$\forall k \in \mathbb{N}, \quad \mathbb{P}\left(P_k \in \bigcup_{t=0}^{X_k-1} A_t\right) = 1.$$

Policy 2: Deliver pre-bunk P_k whenever there are misinformed (I_t^k) nodes in the **neighborhood**.

$$A_{t} = \left\{ P_{k} \in P \mid I_{t}^{k} \cap N_{\mathsf{in}}(c) \neq \emptyset, P_{k} \notin \bigcup_{\tau=1}^{t-1} A_{\tau} \right\}$$



¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



$$\forall k \in \mathbb{N}, \quad \mathbb{P}\left(P_k \in \bigcup_{t=0}^{X_k - 1} A_t\right) = 1.$$

Policy 2: Deliver pre-bunk P_k whenever there are misinformed (I_t^k) nodes in the **neighborhood**.

$$A_{t} = \left\{ P_{k} \in P \mid I_{t}^{k} \cap N_{\mathsf{in}}(c) \neq \emptyset, P_{k} \notin \bigcup_{\tau=1}^{t-1} A_{\tau} \right\}$$

 Ogen

 Dser

 Neighborhood

Policy 1: Deliver pre-bunk P_k whenever $I_t^k \neq \emptyset$ in the **observed network.**

$$A_{t} = \left\{ P_{k} \in P \mid I_{t}^{k} \neq \emptyset, P_{k} \notin \bigcup_{\tau=1}^{t-1} A_{\tau} \right\}$$



¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



$$\forall k \in \mathbb{N}, \quad \mathbb{P}\left(P_k \in \bigcup_{t=0}^{X_k - 1} A_t\right) = 1.$$

Policy 2: Deliver pre-bunk P_k whenever there are misinformed (I_t^k) nodes in the **neighborhood**.

$$A_{t} = \left\{ P_{k} \in P \mid I_{t}^{k} \cap N_{\mathsf{in}}(c) \neq \emptyset, P_{k} \notin \bigcup_{\tau=1}^{t-1} A_{\tau} \right\}$$

OurUserNeighborhood

Policy 1: Deliver pre-bunk P_k whenever $I_t^k \neq \emptyset$ in the **observed network.**

$$A_{t} = \left\{ P_{k} \in P \mid I_{t}^{k} \neq \emptyset, P_{k} \notin \bigcup_{\tau=1}^{t-1} A_{\tau} \right\}$$



Both policies yield feasible but suboptimal solutions.

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



Proxy Problem

Optimally deliver $\lfloor \epsilon T \rfloor$ pre-bunks until time *T*.

$$\min_{t_1...t_{\lfloor \in T \rfloor}} \max_{n \in \{0... \lfloor \epsilon T \rfloor\}} \sum_{i=1}^n \beta^{t_n - t_i},$$
s.t. $0 \le t_{i-1} \le t_i \le T, \quad \forall i \in \{1... \lfloor \epsilon T \rfloor\}$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025

autonomy

Proxy Problem

Optimally deliver $\lfloor eT \rfloor$ pre-bunks until time *T*.

$$\min_{t_1...t_{\lfloor \in T \rfloor}} \max_{n \in \{0...\lfloor \epsilon T \rfloor\}} \sum_{i=1}^n \beta^{t_n - t_j},$$
s.t. $0 \le t_{i-1} \le t_i \le T, \quad \forall i \in \{1...\lfloor \epsilon T \rfloor\}$

Theorem: Let $t_0 \dots t_{\lfloor \in T \rfloor} \in \mathbb{R}$ be an optimal solution to above problem , then there exists h such that $t_i - t_{i-1}$ is constant across all $i = h + 1, \dots, \lfloor \epsilon T \rfloor$. Specifically there exists a constant α such that¹,

$$t_h - t_{h-1} = \log_\beta(\alpha - 1) - \log_\beta(h),$$

$$t_i - t_{i-1} = \log_\beta(\alpha - 1) - \log_\beta(\alpha).$$

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025

autonomy

Proxy Problem

Optimally deliver $\lfloor eT \rfloor$ pre-bunks until time *T*.

$$\min_{t_1...t_{\lfloor \in T \rfloor}} \max_{n \in \{0...\lfloor \epsilon T \rfloor\}} \sum_{i=1}^n \beta^{t_n - t_j},$$
s.t. $0 \le t_{i-1} \le t_i \le T, \quad \forall i \in \{1...\lfloor \epsilon T \rfloor\}$

Theorem: Let $t_0...t_{\lfloor \in T \rfloor} \in \mathbb{R}$ be an optimal solution to above problem , then there exists h such that $t_i - t_{i-1}$ is constant across all $i = h + 1, ..., \lfloor \epsilon T \rfloor$. Specifically there exists a constant α such that¹,

$$t_h - t_{h-1} = \log_\beta(\alpha - 1) - \log_\beta(h),$$

$$t_i - t_{i-1} = \log_\beta(\alpha - 1) - \log_\beta(\alpha).$$

Optimal solution of the proxy problem evenly distributes pre-bunks.

¹ Y. E. Bayiz, U. Topcu, *Prebunking Design as a Defense Mechanism Against Misinformation Propagation on Social Networks,* Submitted to CDC 2025



Assumption: Content arrives already sorted in a **content feed** based on original value assignment $V(\sigma)$



¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Re-ranking for Credibility

autonomy

Assumption: Content arrives already sorted in a content feed based on original value assignment $V(\sigma)$



¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Re-ranking for Credibility

autonomy

Assumption: Content arrives already sorted in a content feed based on original value assignment $V(\sigma)$



Assume that user sees p^{th} post with probability

$$q_p = \lambda^p \longrightarrow \text{persistence}$$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Assumption: Content arrives already sorted in a **content feed** based on original value assignment $V(\sigma)$

Identity permutation optimizes platform objectives



Assume that user sees p^{th} post with probability

$$q_p = \lambda^p \longrightarrow \text{persistence}$$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

autonomy

),

Assumption: Content arrives already sorted in a **content feed** based on original value assignment $V(\sigma)$

Identity permutation optimizes platform objectives



$$V(\sigma) = \sum_{k=1}^{K} q_{p} v(\sigma(p))$$

p=1

 $e = \arg \max V(\sigma)$

Assume that user sees p^{th} post with probability

$$q_p = \lambda^p \longrightarrow \text{persistence}$$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Assumption: Content arrives already sorted in a **content feed** based on original value assignment $V(\sigma)$

Identity permutation optimizes platform objectives



¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Fake News
- Propaganda
- Scams
- Rumors

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Fake News
- Propaganda
- Scams
- Rumors

Instead of labeling **misinformation** we use a probabilistic **credibility** measure

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Fake News
- Propaganda
- Scams
- Rumors

Democratic View

Suppose a randomly selected user labels post p as $Z_p \sim \text{Bernoulli}(c)$

Define credibility as $E[Z_p] = c$

Instead of labeling **misinformation** we use a probabilistic **credibility** measure

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Fake News
- Propaganda
- Scams
- Rumors

Instead of labeling misinformation we use a — probabilistic credibility measure

Democratic View

Suppose a randomly selected user labels post p as $Z_p \sim \text{Bernoulli}(c)$

Define credibility as $E[Z_p] = c$

Evidence-based View

Parse top k news articles

Determine credibility based on the consensus

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025







- Propaganda
- Scams
- Rumors

Instead of labeling **misinformation** we use a — probabilistic **credibility** measure

Democratic View

Suppose a randomly selected user labels post p as $Z_p \sim \text{Bernoulli}(c)$

Define credibility as $E[Z_p] = c$

Evidence-based View

Parse top k news articles

Determine credibility based on the consensus

We use a hybrid approach that utilizes both views.

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





$$\hat{c}_{p}^{\text{human}} = \frac{\sum_{i=1}^{n} \left(z_{p,i} \cdot h_{p,i}^{+} + (1 - z_{p,i}) \cdot h_{p,i}^{-} \right)}{\sum_{i=1}^{n} \left(h_{p,i}^{+} + h_{p,i}^{-} \right)},$$

- Each post *p* comes with community notes
- Each community note *i* comes with a verdict $z_{p,i}$
- Eat community note also comes with ratings:
 - $h_{p,i}^+$: Number of 'helpful' ratings
 - $h_{p,i}^-$: Number of 'not helpful' ratings

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Each post *p* comes with community notes
- Each community note *i* comes with a verdict $z_{p,i}$
- Eat community note also comes with ratings:
 - $h_{p,i}^+$: Number of 'helpful' ratings
 - $h_{p,i}^-$: Number of 'not helpful' ratings



Total number of ratings across all posts

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025





- Each post *p* comes with community notes
- Each community note *i* comes with a verdict $z_{p,i}$
- Eat community note also comes with ratings:
 - $h_{p,i}^+$: Number of 'helpful' ratings
 - $h_{p,i}^-$: Number of 'not helpful' ratings



Total number of ratings across all posts

Assuming ratings and notes come from uniformly sampled users

 \hat{c}_p^{human} is an unbiased estimator of $c = E[Z_p]$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025



Generate artificial community notes using retrieval augmented generation (RAG)

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Enhancing Human Scores

autonomy

Generate artificial community notes using retrieval augmented generation (RAG)



¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Enhancing Human Scores

autonomy





¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

Enhancing Human Scores

autonomy



Generate artificial community notes using retrieval augmented generation (RAG)

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds*, Submitted to SIGIR 2025

CENTER FOR autonomy

Assumption: Content arrives already sorted in a content feed based on original value assignment $V(\sigma)$



$$\max_{\sigma \in S_K} \quad \alpha V(\sigma) + \sum_{p=1}^K q_p c_{\sigma(p)}$$

Problem: $V(\sigma)$ is not known

($V(\sigma)$ may not even exist)

Assume identity permutation optimizes platform objectives

$$e = \arg \max_{\sigma} V(\sigma)$$
 $V(\sigma) = \sum_{p=1}^{K} q_p v(\sigma(p)),$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

CENTER FOR autonomy

Assumption: Content arrives already sorted in a content feed based on original value assignment $V(\sigma)$



 $\max_{\sigma \in S_K} \quad \alpha V(\sigma) + \sum_{p=1}^K q_p c_{\sigma(p)}$

Problem: $V(\sigma)$ is not known

($V(\sigma)$ may not even exist)

Solution: Solve a surrogate problem

$$\min_{\sigma} D(\sigma) - \sum_{p=1}^{K} q_p c_{\sigma(p)}$$
$$D(\sigma) = \sum_{i=1}^{K} |\sigma(i) - i|$$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025



Theoretically "Optimal" solutions: $\max_{\sigma \in S_K} \alpha V(\sigma) + \sum_{p=1}^{\infty} q_p c_{\sigma(p)}$

¹ Y. E. Bayiz, U. Topcu, *Re-ranking with Semi-automated Misinformation Detection for Increasing Credibility in Social Media Feeds,* Submitted to SIGIR 2025

autonomy

• Challenges:

- Information ecosystem is not completely understood
- Information ecosystem is fragile against adversarial attacks
- Adversaries have initiatives
- Future Work:
 - How to design a robust social network?
 - Long term evolution of the network
 - How to integrate multiple information modalities?

Conclusion

Control and decision making enable new solutions for this challenge that goes well beyond fact-checking.

