

We're living in a post truth era → disinfo has consistently stayed a head of research

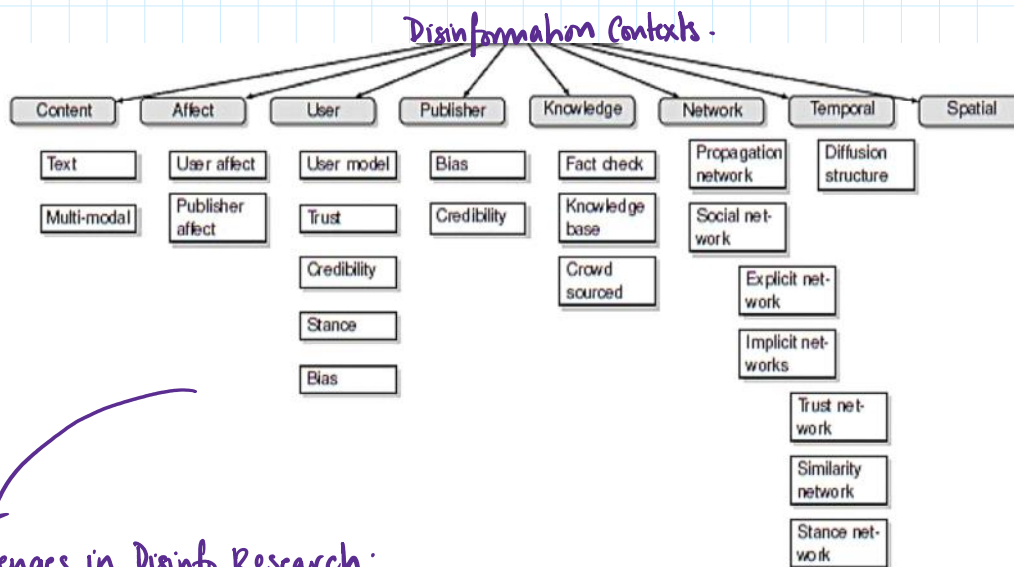
fake news was just in 2017.

Now we look at it w/ diff perspectives → fakeness
→ newsworthiness
→ shareworthiness

Memes, satire, humor ALSO influence p well.

wrong/incorrect
↑
MISINFO — not to mislead
↑
apart/away
↑
DISINFO — To mislead
↑
reversal
↑
MALINFO — to harm
↑
evil

| Term | Information type | Example | Intent | Authenticity |
|--|--------------------------------|------------|--------------------------------|------------------------------|
| Misinformation | False | Rumour | Not to mislead | False |
| Disinformation | False | Fake news | To mislead | False |
| Malinformation | Genuine | True news | To harm | True |
| Coordinated inauthentic behavior (CIB) | Can be mix of genuine and fake | Propaganda | To influence, harm and mislead | Can be mix of True and False |



Challenges in Disinfo Research :

① Datasets → all context support?
→ explainability?

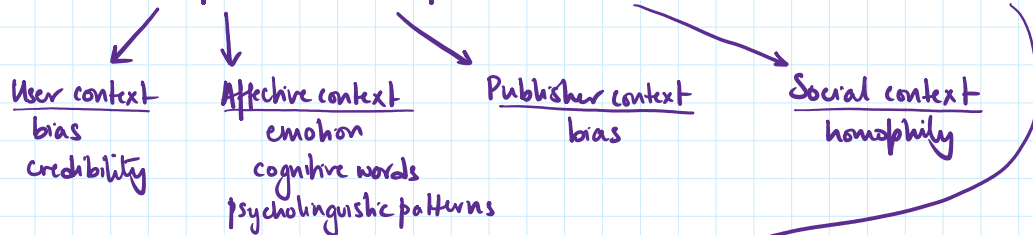
② Existing interpretability techniques

- feature analysis
- post-hoc
- attenⁿ + cross attⁿ
- Sub models
- anomaly detection
- Zero shot learning
- fine grained reasoning

③ Early detection w/ context → content
→ propagation
→ user
→ knowledge
→ affect.

④ Weak social supervision — depends on labels derived from other models trained on limited data

Weak social supervision depends on labels derived from other models trained on limited data



Okay but how do you learn w/o labelled data?
Unsupervised? Hell no. Too many contexts.

- How do you work w/ multiple contexts?
- ① Multiview learning
 - co training
 - view aggregation
 - view discordance
 - view ensemble
 - cross view
 - ② Multitask learning
 - ③ Multitask multiview
- Zero/few shot — works w/ typologically similar content.
 - Transfer learning — catastrophic forgetting in domain adaptation for disinfo?
 - SSL — Deep SSL + graph based SSL.

LLMs + disinfo

- ① Weakness:
 - can't reason
 - factual errors, bias, + discrimination
 - can't capture nuanced understanding
- ② Fake news detection:
 - LLM makes stylistically consistent text. (models which depend on diff styles fail)
 - but diff feature distributions
 - ↳ also: selection + integration of rationale — LLM cannot.
- ③ Rumor detection:
 - not dependent on world knowledge
 - work w/ emotion, rebuttal, support } LLM can.
- ④ Propaganda detection:
 - works w/ english. Not other langs
- ⑤ Dataset generation:
 - adversarial prompting!
- ⑥ Use zero + few shot learning!
 - LLM for weak supervision + credibility signals.