# DATA ANALYTICS ARE POWERFUL
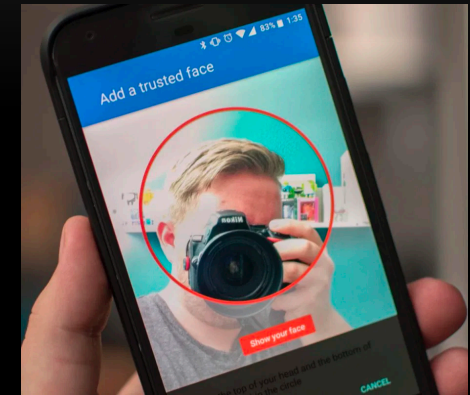
# HANDLE WITH CARE

Jeremy D. Wendt

Joint work with Philip Kegelmeyer, Ali Pinar, Tim Shead,

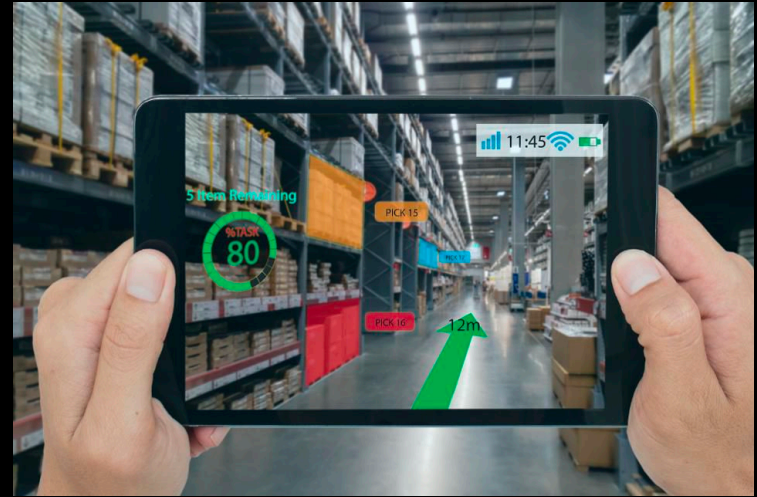Gary Saavedra, Cosmin Safta, Joe Bertino, and many others
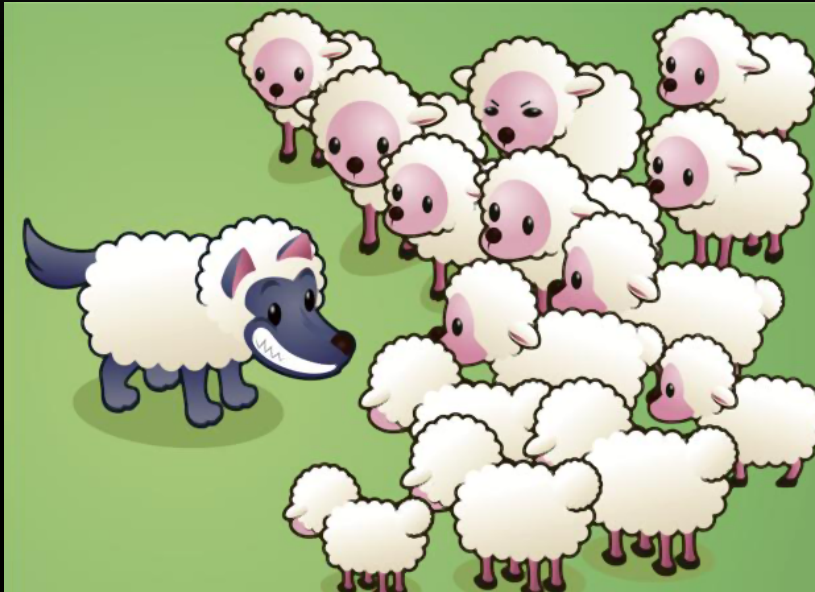
# NITROGLYCERIN



Ascanio Sobrero

# DATA ANALYTICS' DAILY EFFECTS - 2019

# DATA ANALYTICS DAILY EFFECTS – SOON?

# SUBVERT





e.g., Kegelmeyer, Shead, Crussell, Rodhouse, Robinson, Johnson, Zage, Davis, Wendt, Doak, Cayton, Colbaugh, Glass, Jones, and Shelburg, 2015.

# EVADE



e.g., Athalye, Engstrom, Ilyas, and Kwok, 2018.

# REVEAL



e.g., Fredrikson, Jha, and Ristenpart, 2015, or Shorki, Stronati, Song, Shmatikov, 2017.
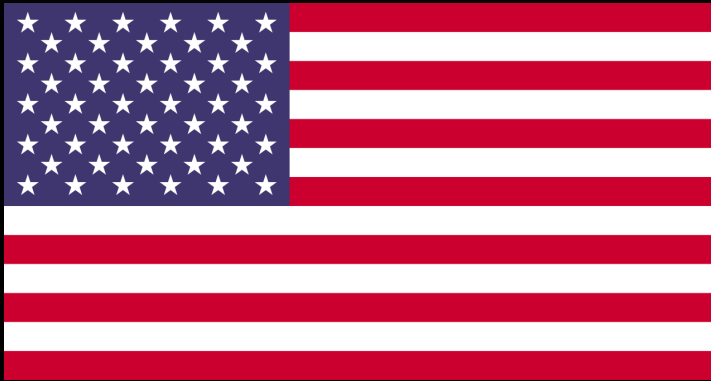
# APPLY



Jordan Peele uses AI, President Obama in fake news PSA

e.g., Isola, Zhu, Zhou, and Efros, 2016 or Kim, Garrido, Tewari, Xu, Thies, NieSSner, Perez, Richardt, Zollhofer, Theobalt, 2018

# WHO ARE ADVERSARIES?

- It gets complicated really quickly

# ADVERSARY'S ABILITIES



Credit: Bernard Goldbach

- In all analyses that follow, we provided the adversary significant abilities:

    - Full awareness of our chosen technique and implementation

    - Full understanding of the data to be analyzed

    - Some ability to alter the data or labels to affect our analytics
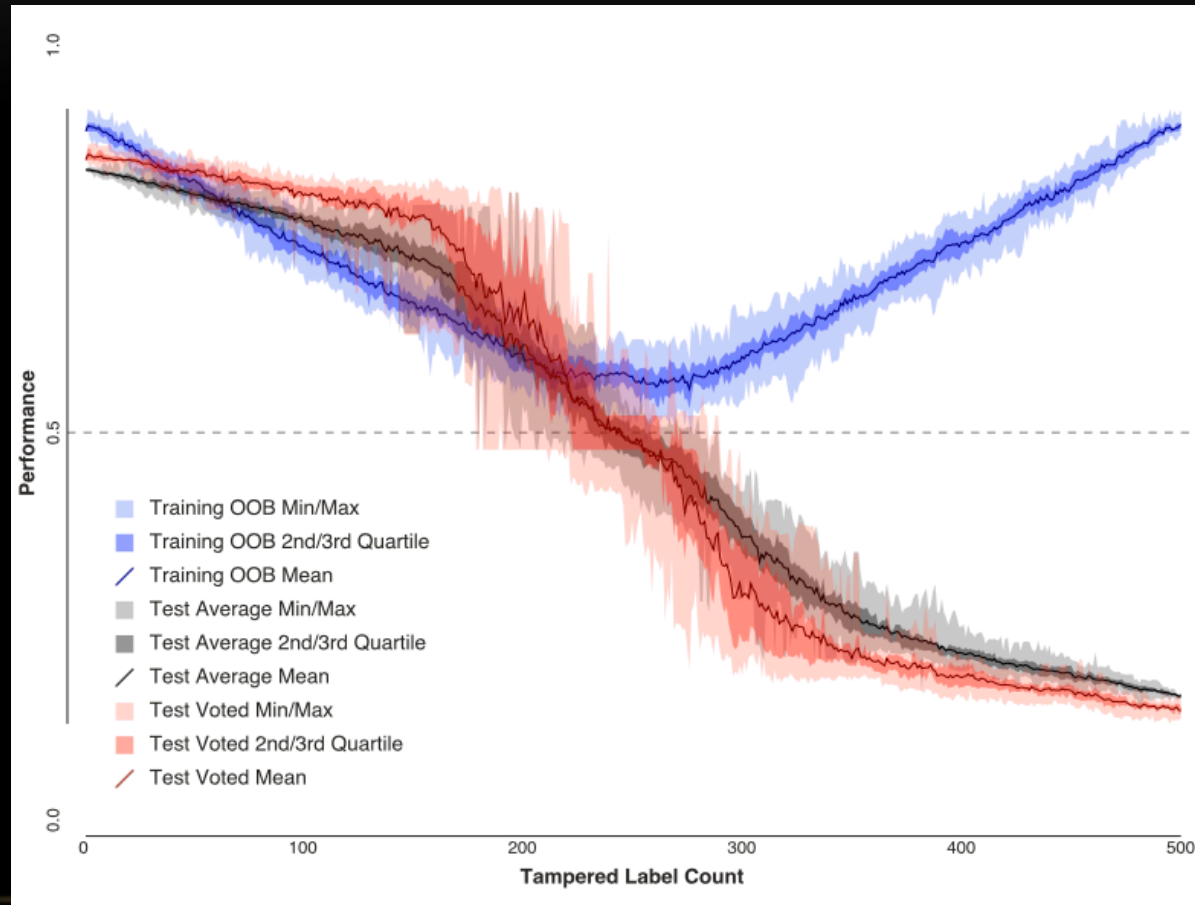
# SUPERVISED LEARNING



What happens if you tell the system this is a '1' or '8' ?
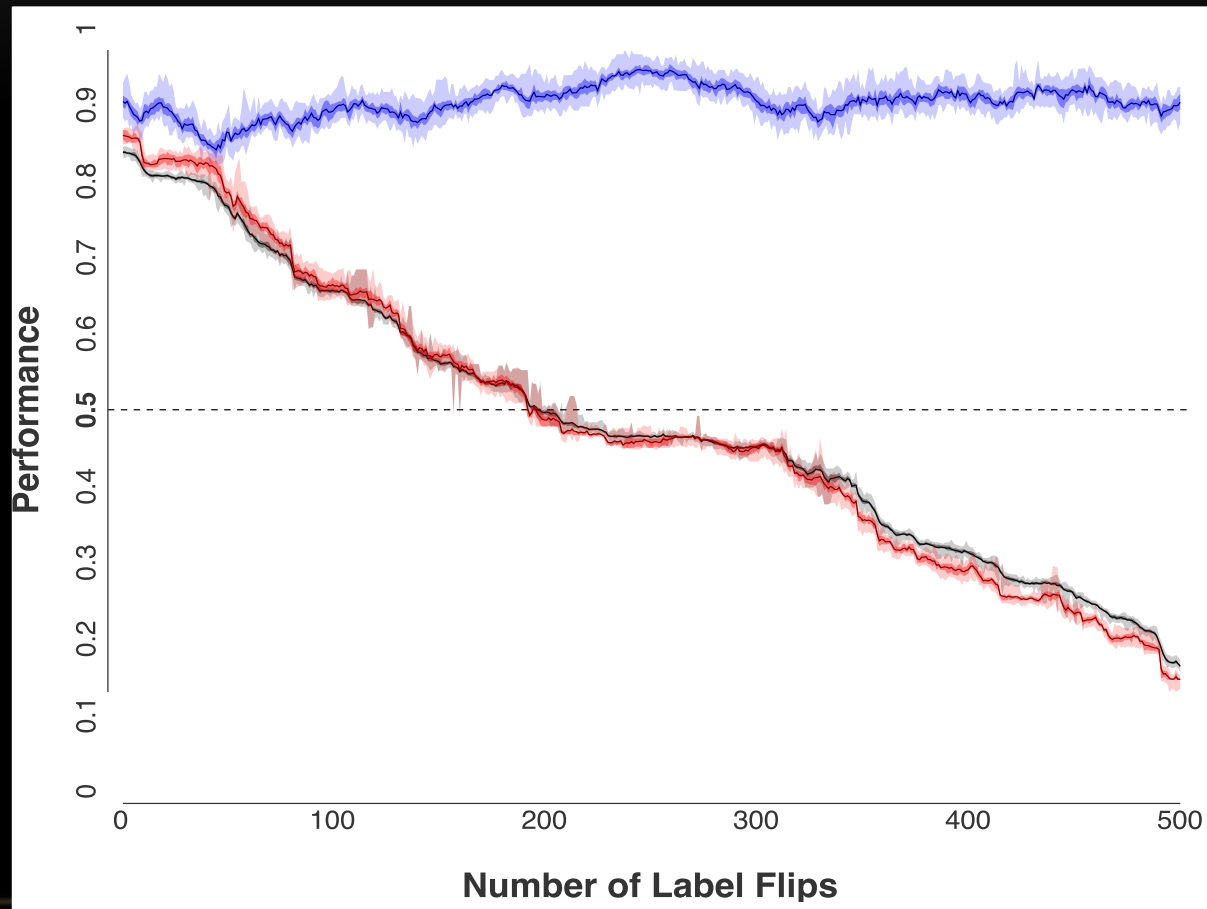
# MEASURING SUPERVISED LEARNING

- Accuracy = Num Correct / Num Possible

  - Requires knowing correct labels

  - However, the defender no longer has correct labels!

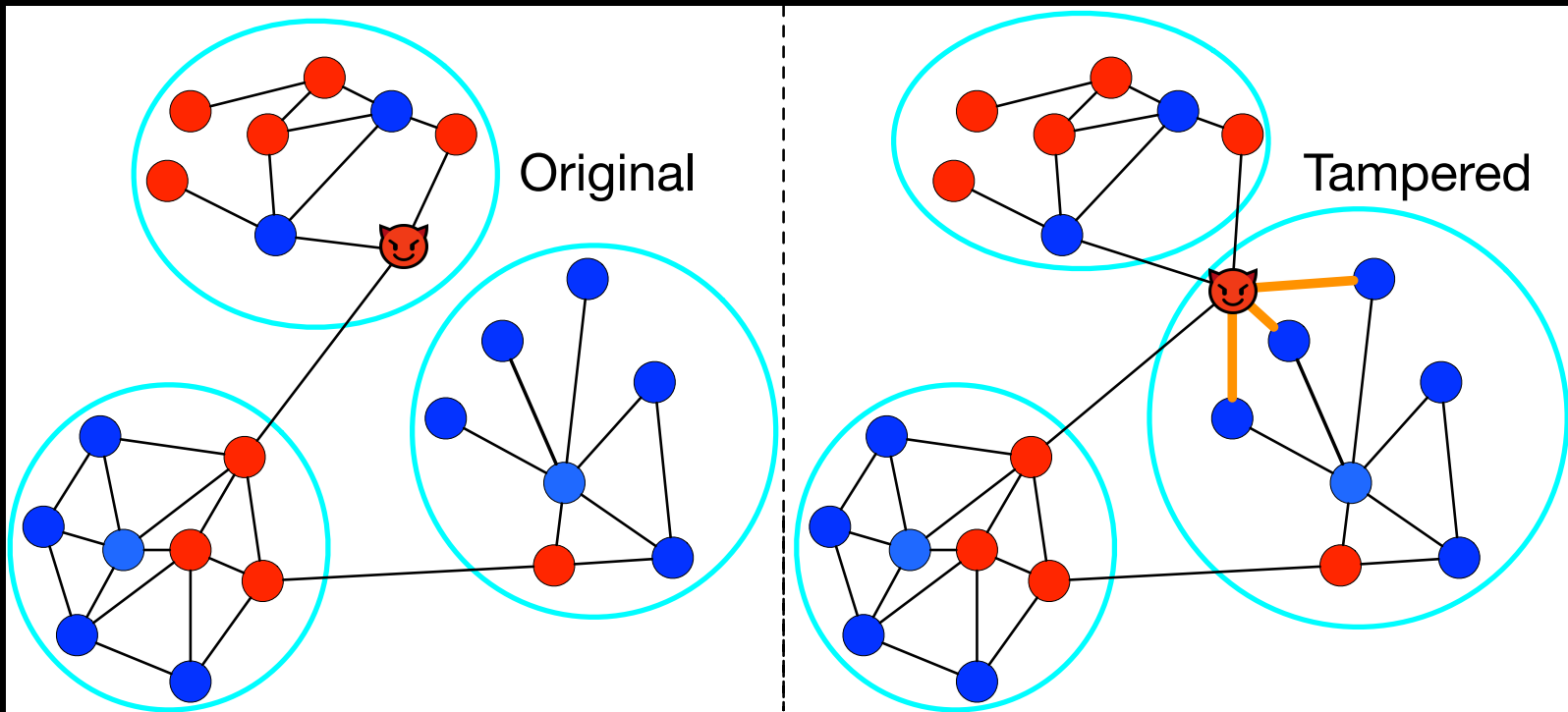- Defender can only measure his accuracy against the altered labels
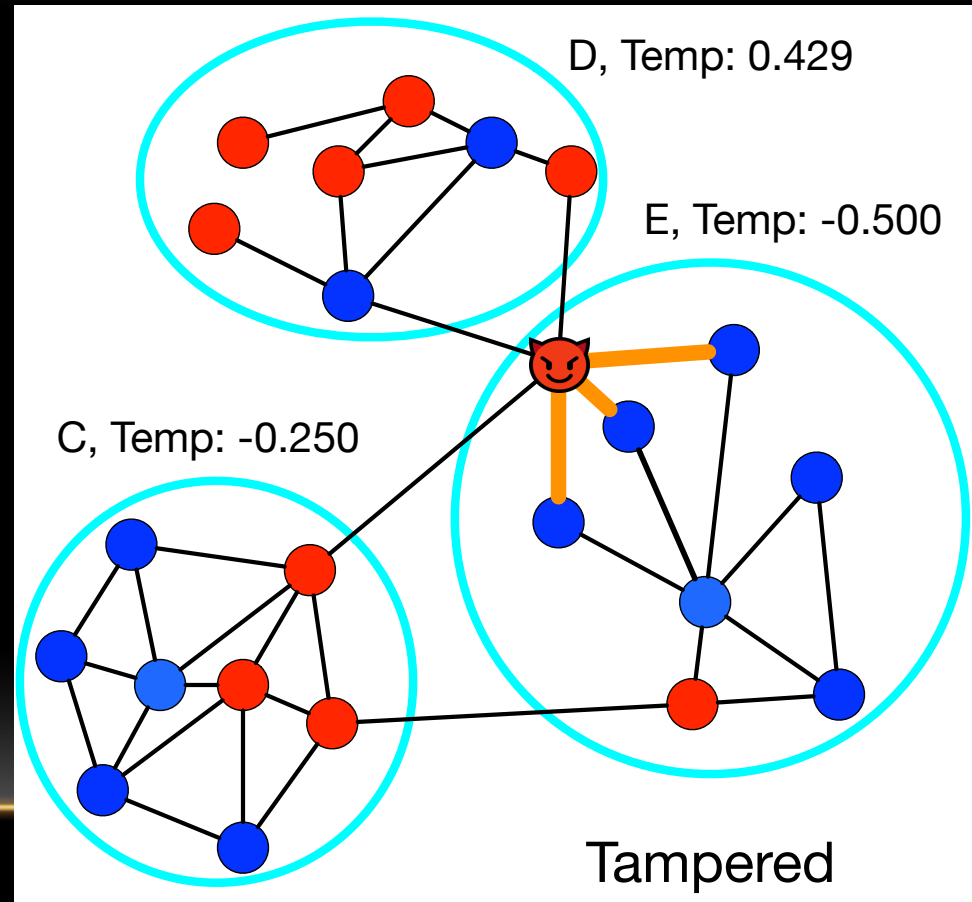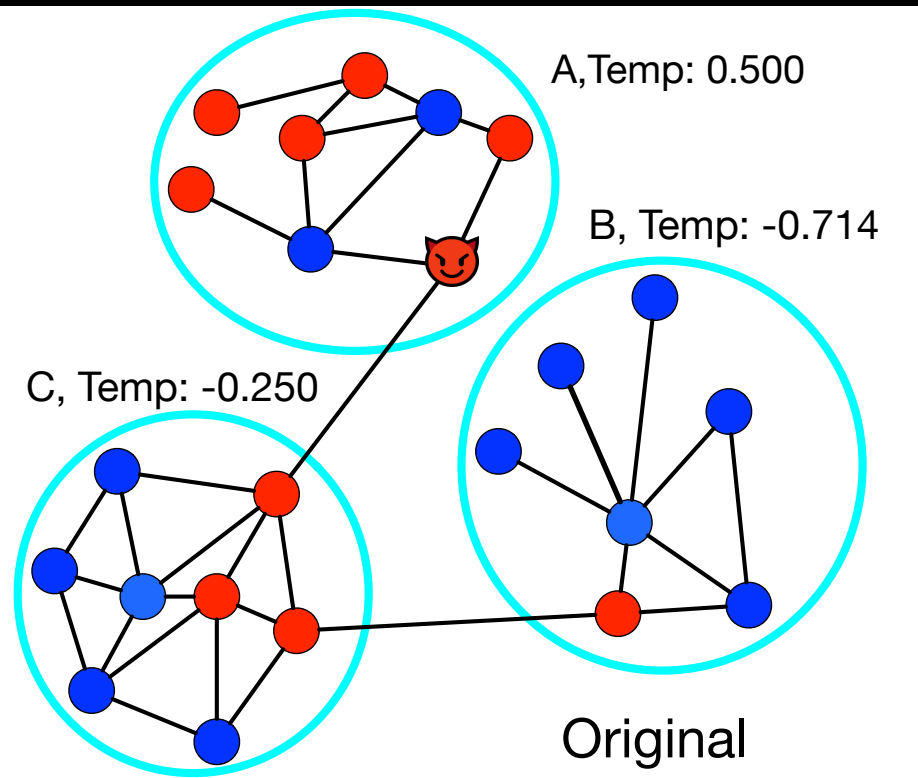
# SUBVERT CONFIDENCE

# SUBVERT WITHOUT DETECTION

# ATTACKING COMMUNITY DETECTION



Original

Tampered

Kegelmeyer, Wendt, and Pinar, 2018

# COMMUNITY TEMPERATURE

Hot (red) = 1; Cold (blue) = -1; Unknown = 0
Community temp is average of nodes' temps
Higher temperature is more suspicious
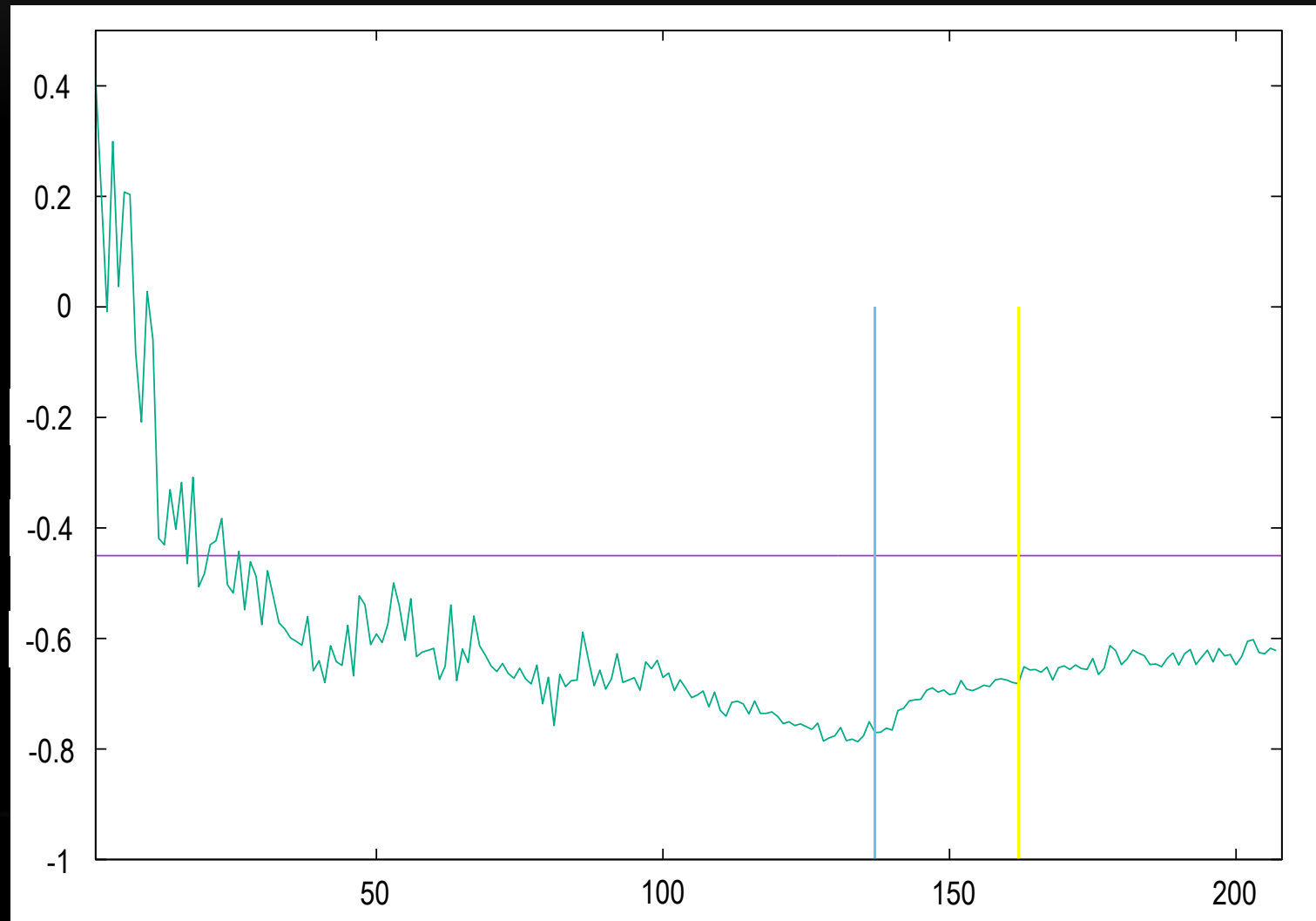


Original



Tampered

# ATTACK HEURISTICS

- Attacker wants to be in cold community

  - He can add edges between his node and others – which to add?

- *Random* – Add random edges; used as a baseline

- *Stratified Random* – Random to cold, then unknown, last hot

- *Cold and Lonely* – Stratified random, but lower degree first w/in temperature bands

- *Stable Structure* – Exploit graph structure; next slide

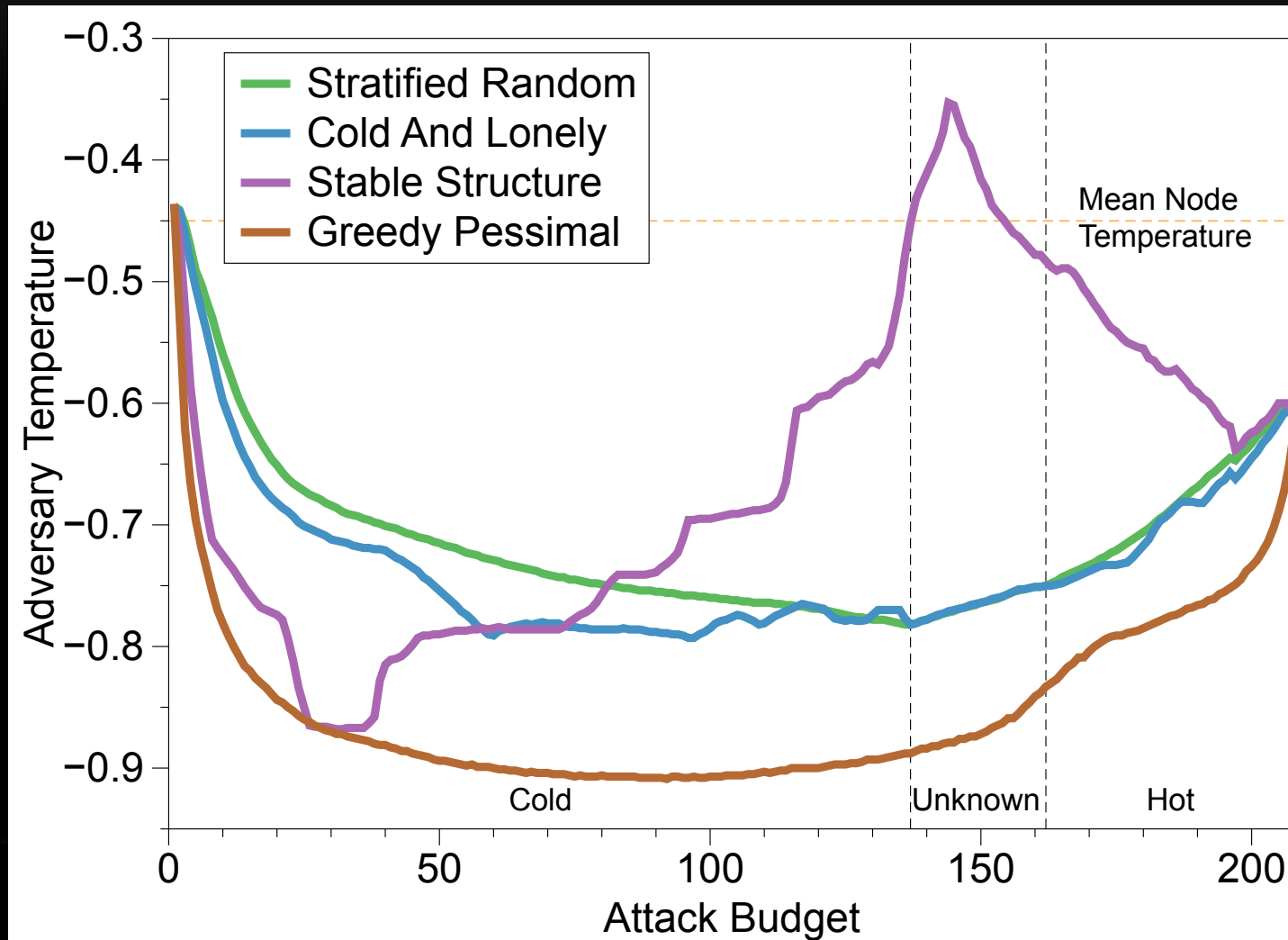- *Greedy Pessimal* – Exhaustively compute worst case; computationally infeasible in general, approximately worst case

# STABLE STRUCTURE

# CREATING ATTACK PLOTS

# DATA ANALYTICS NEED DEFENSES

- Computer Networks in 1990s → Data Analytics in 2019

# DYNAMITE



Alfred Nobel

# COMPETITION IS UBIQUITOUS