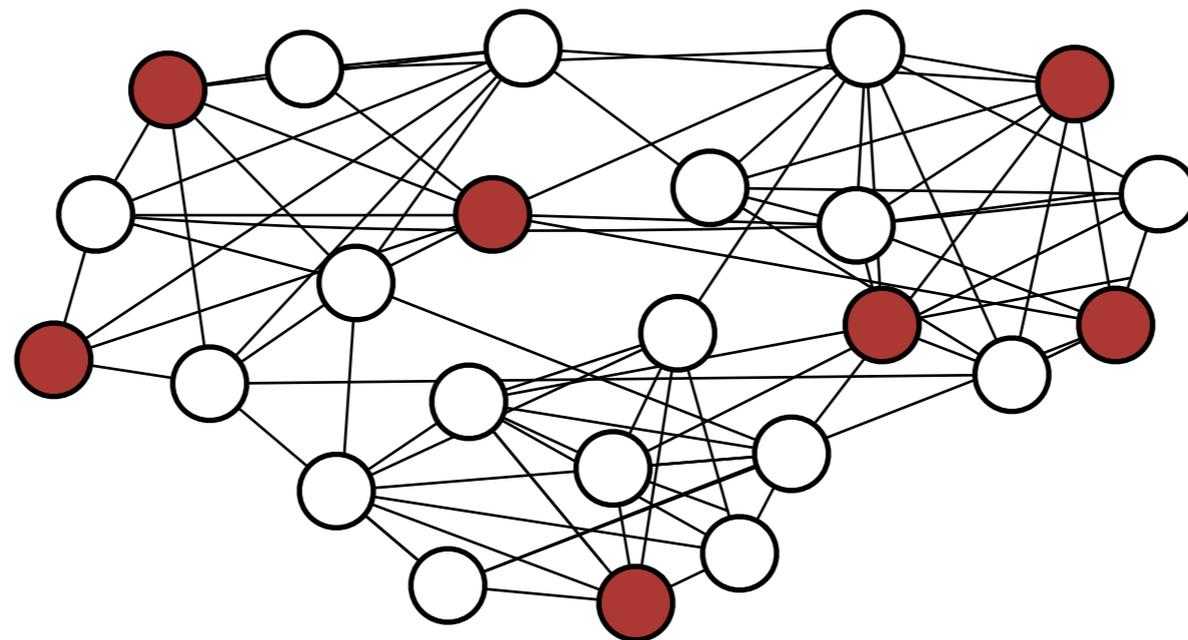
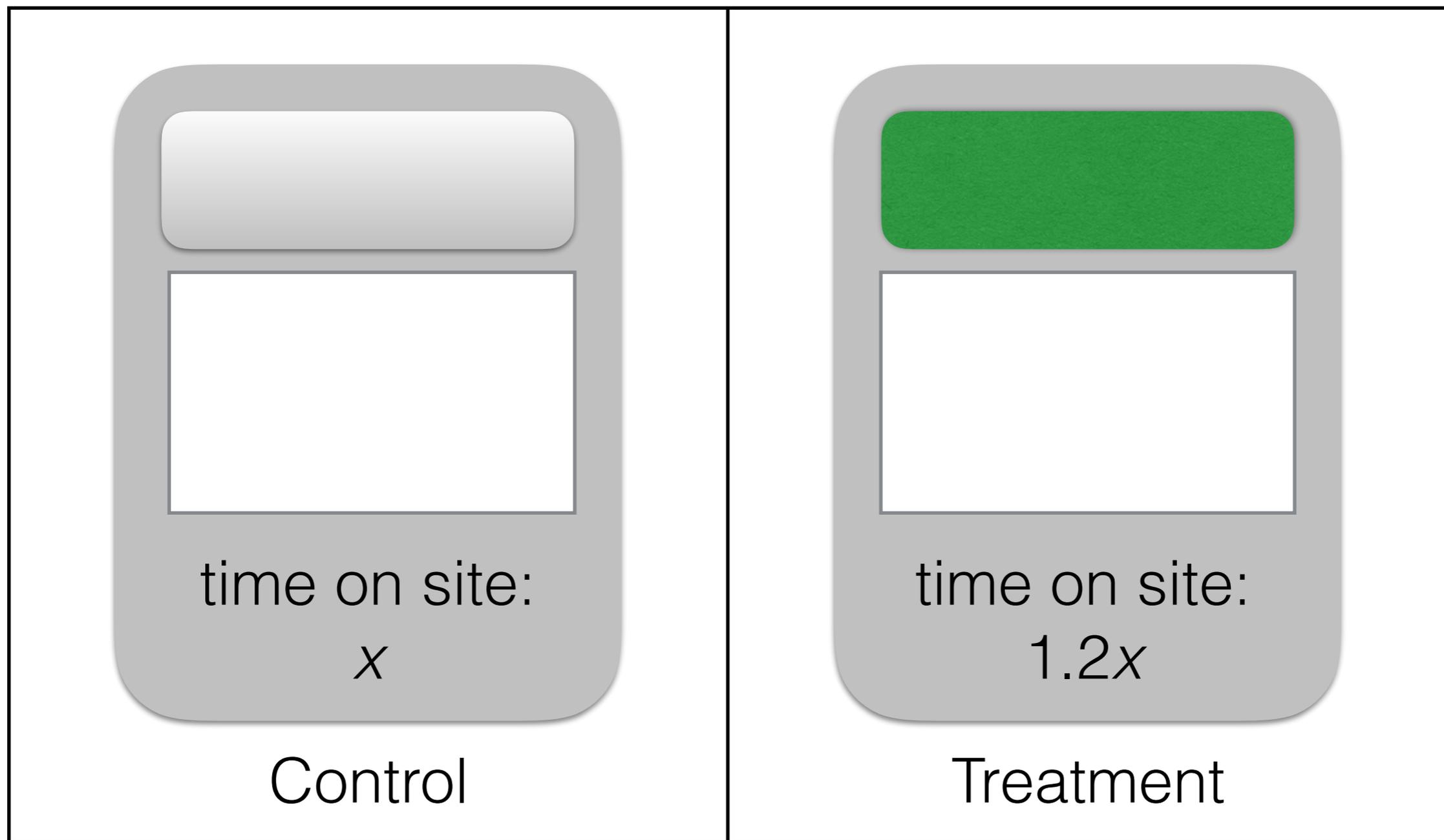


# A/B Testing in Networks with Adversarial Members

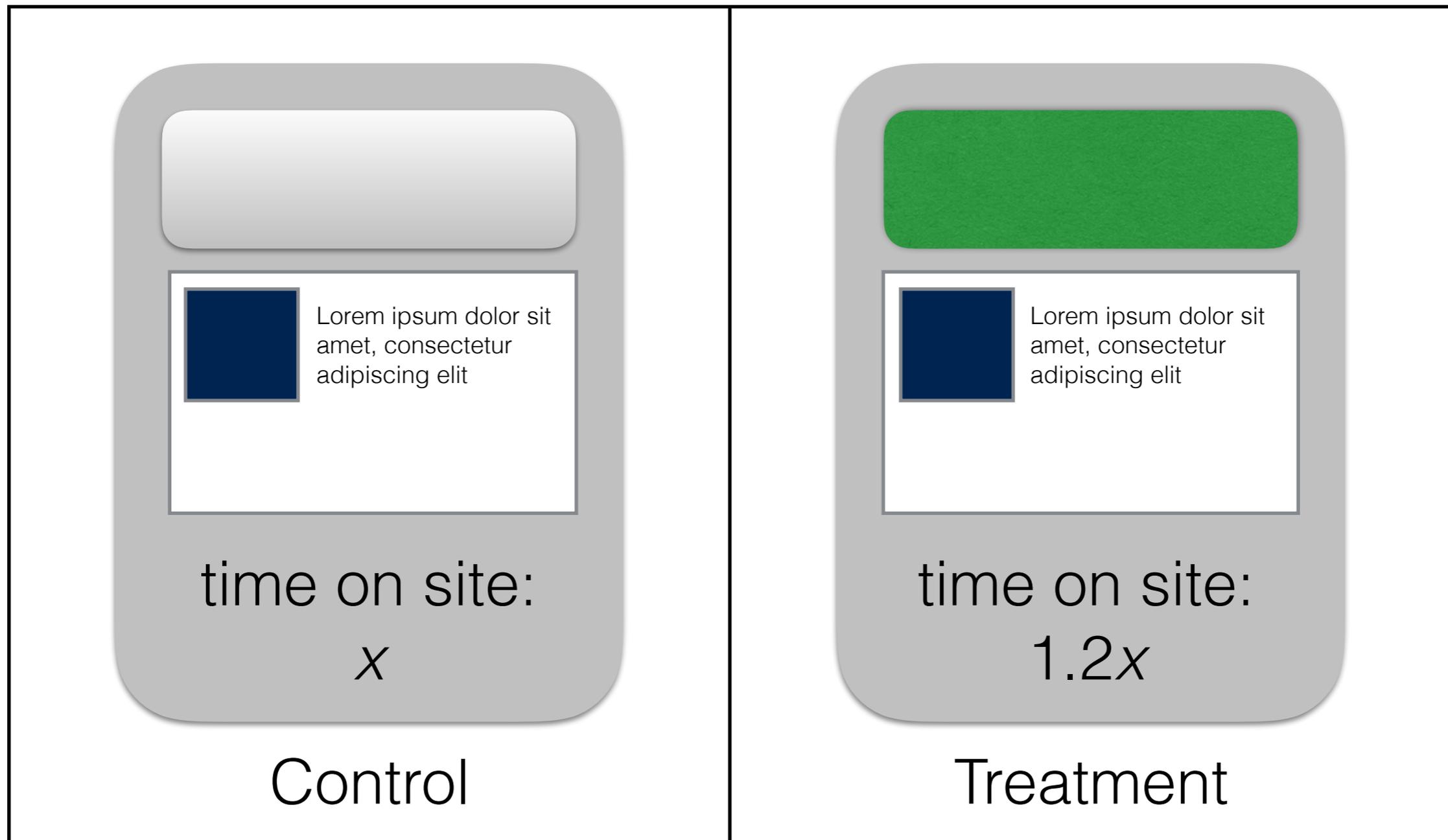
MLG Workshop 2017  
**Kaleigh Clary** and David Jensen  
University of Massachusetts Amherst  
August 14, 2017

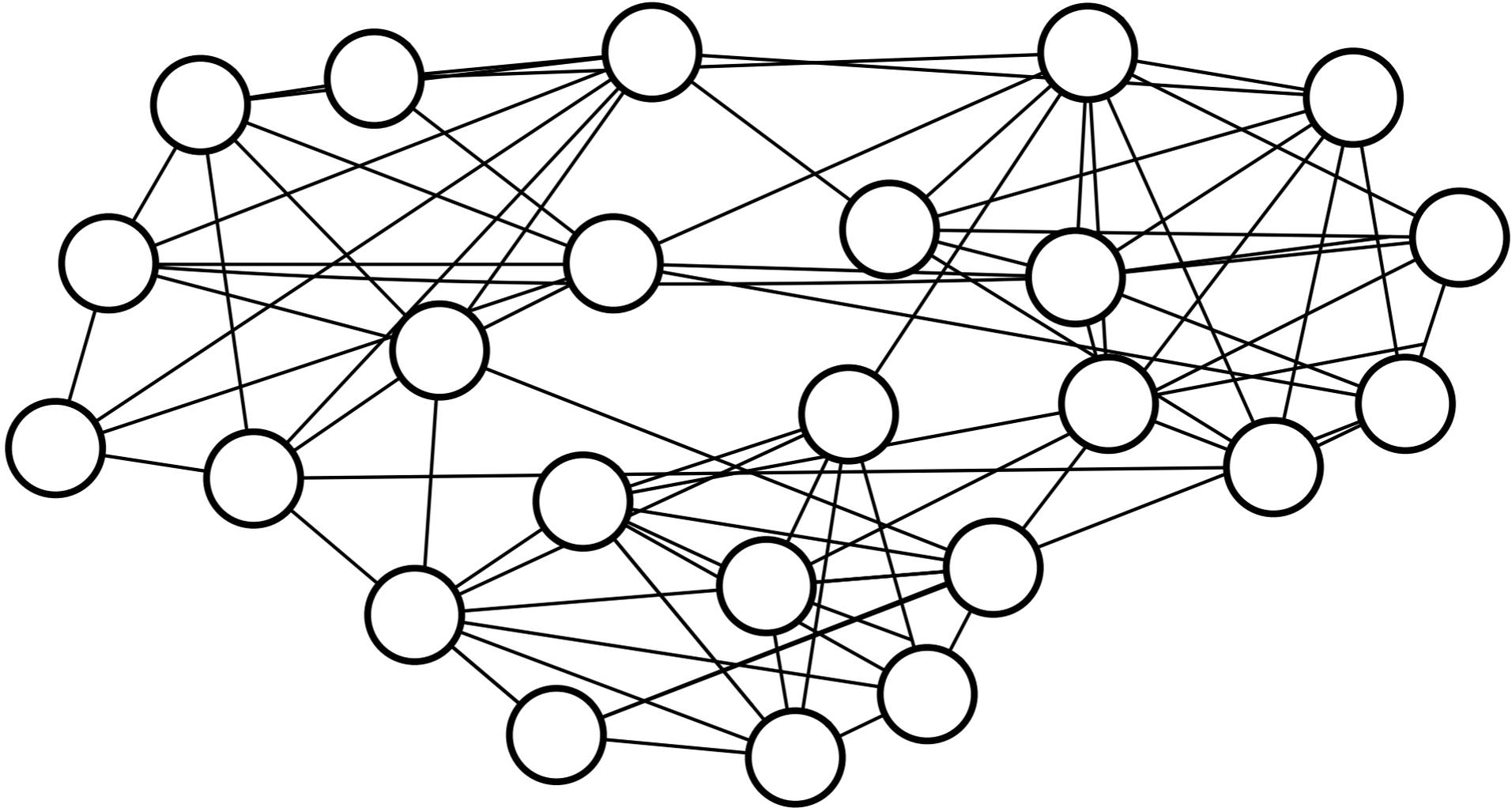


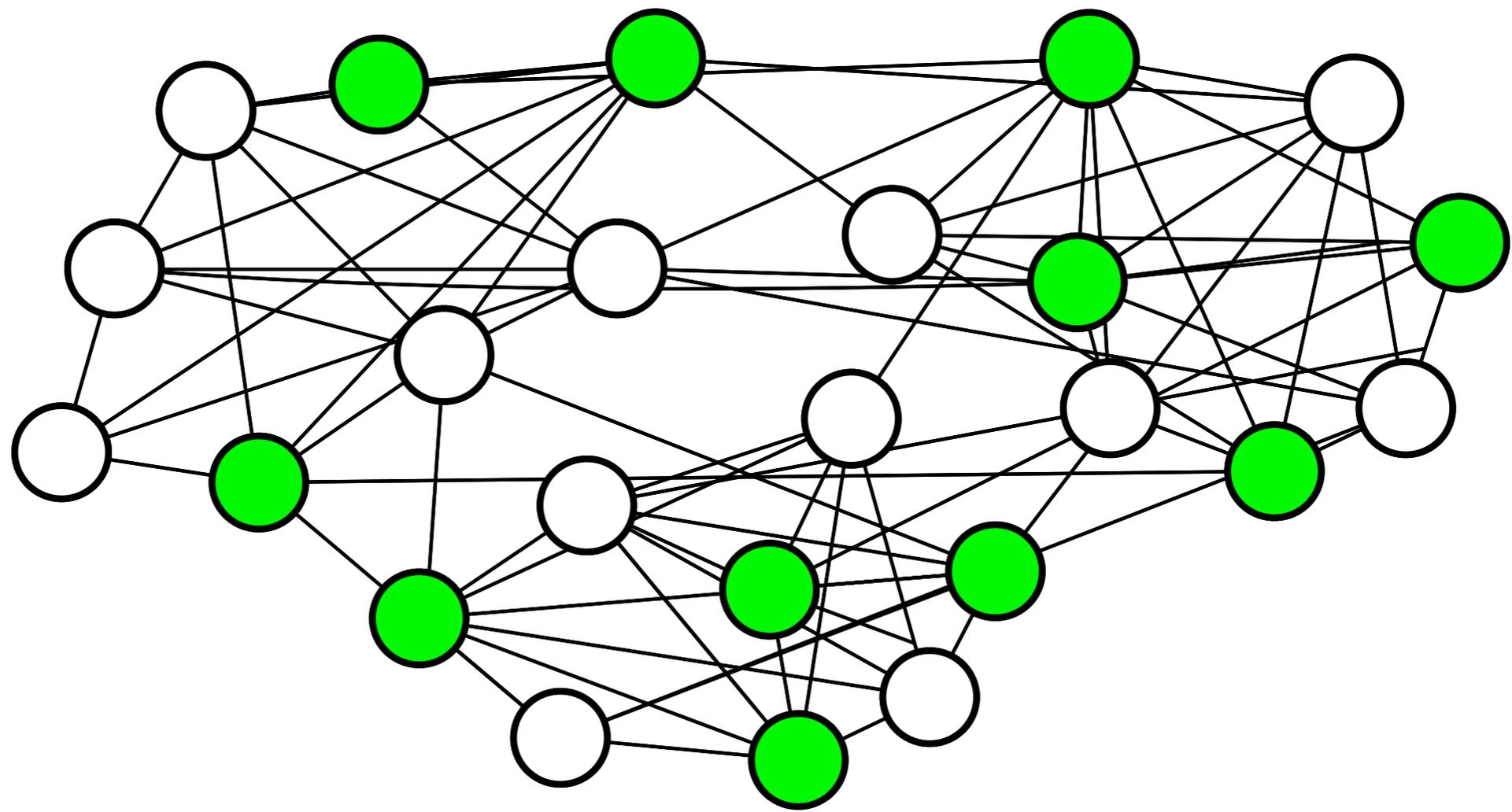
# A/B Testing

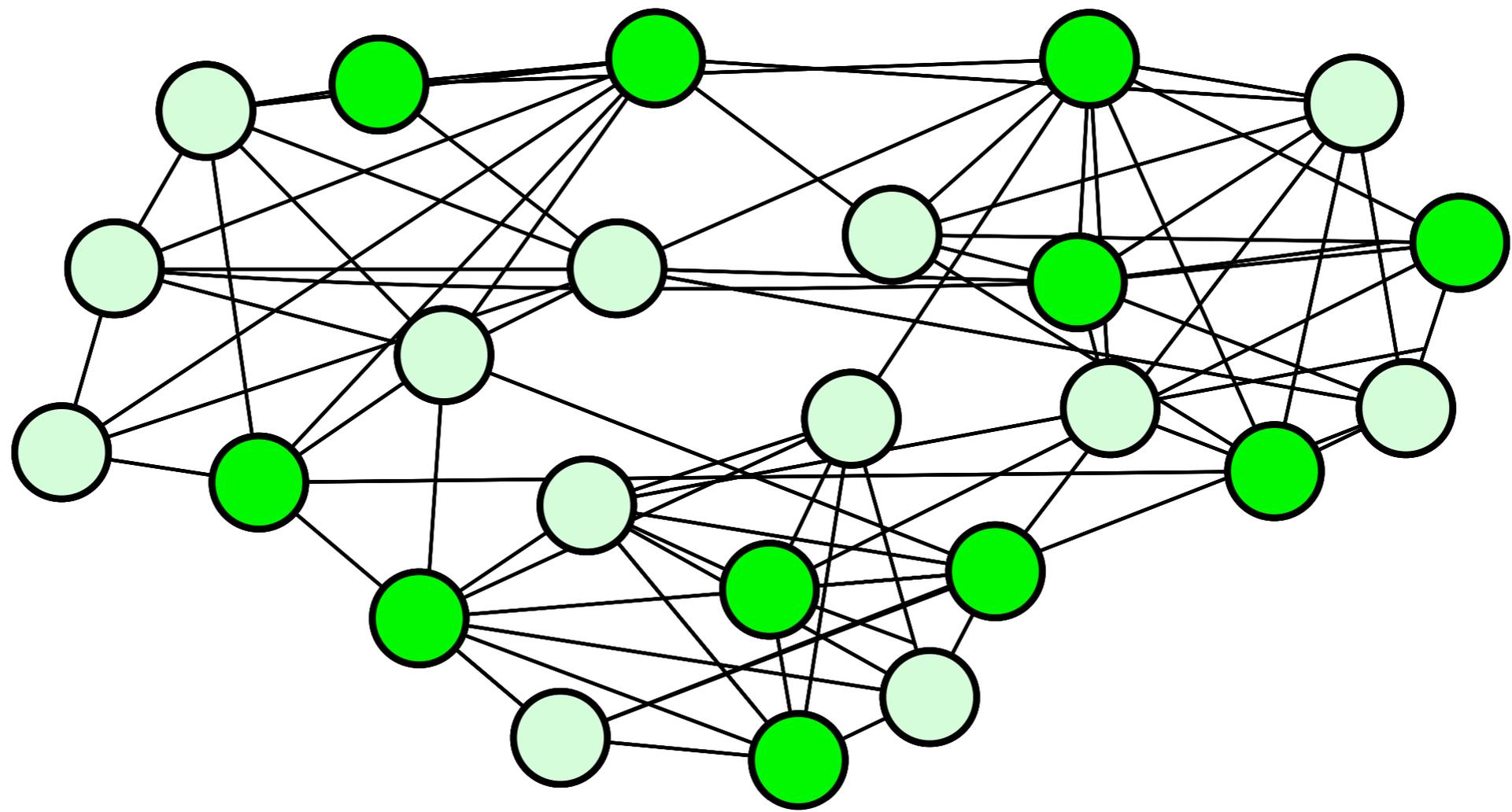


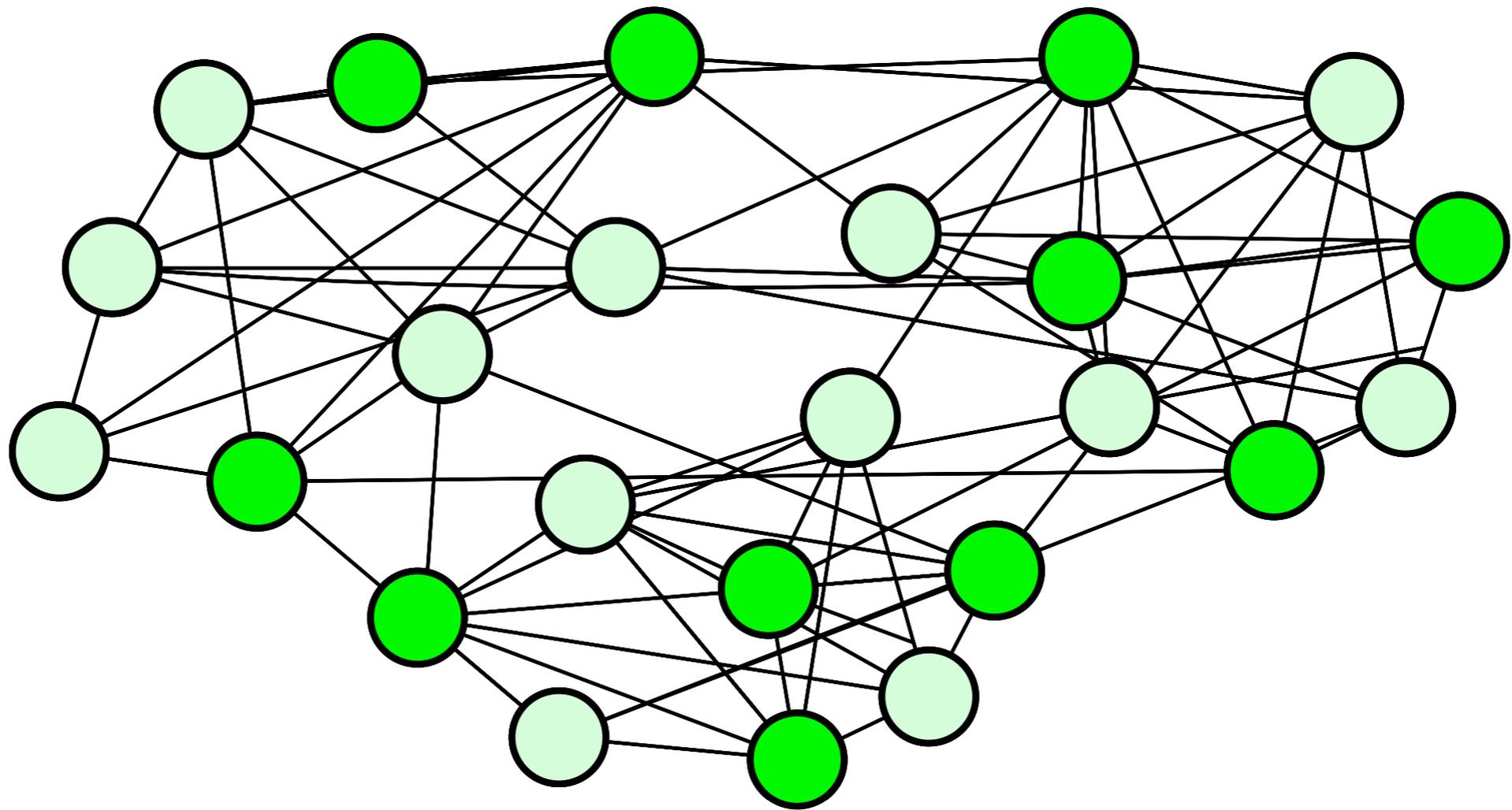
# A/B Testing





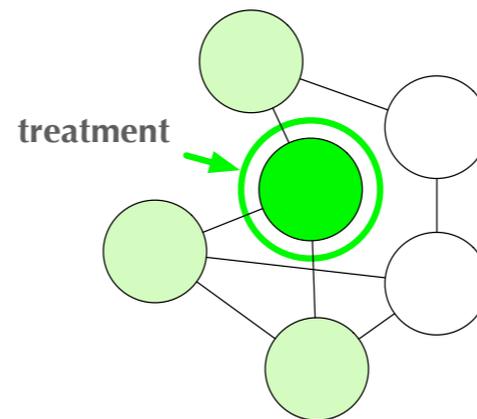
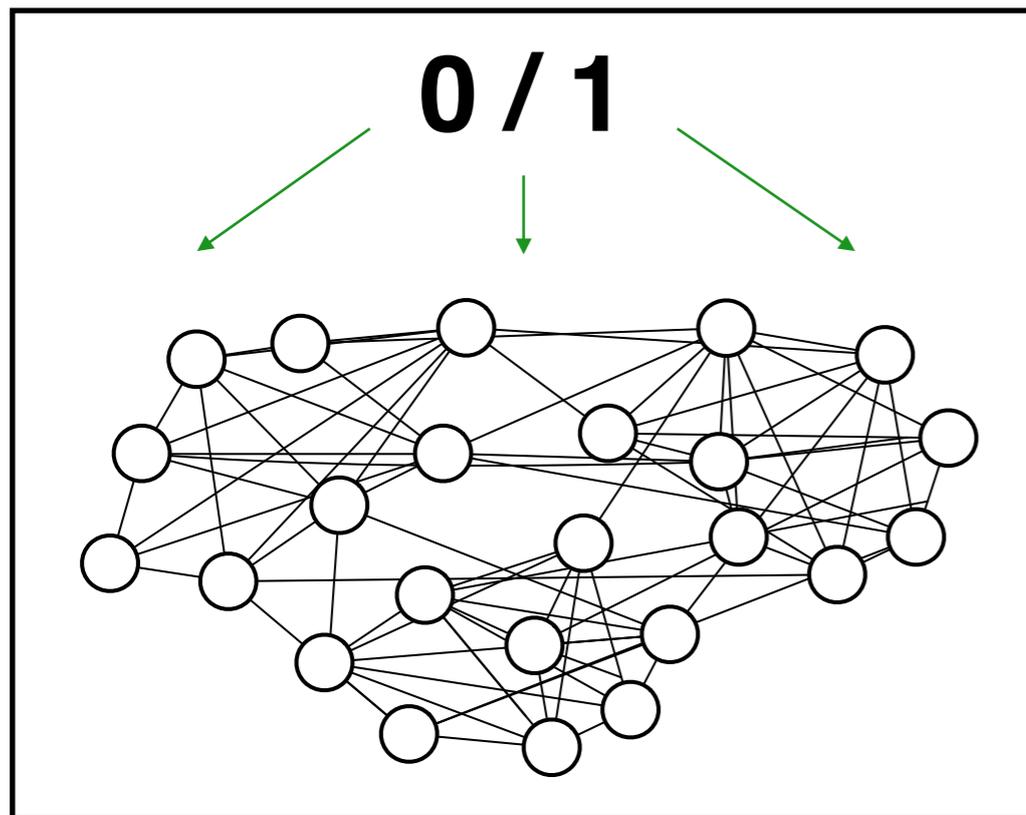




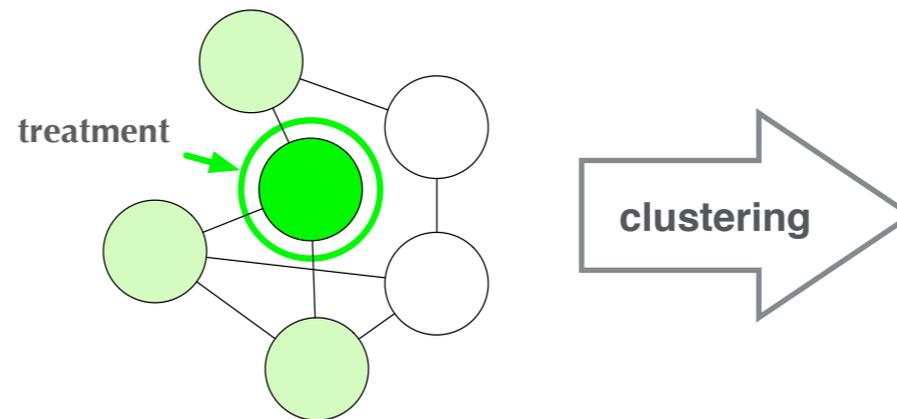
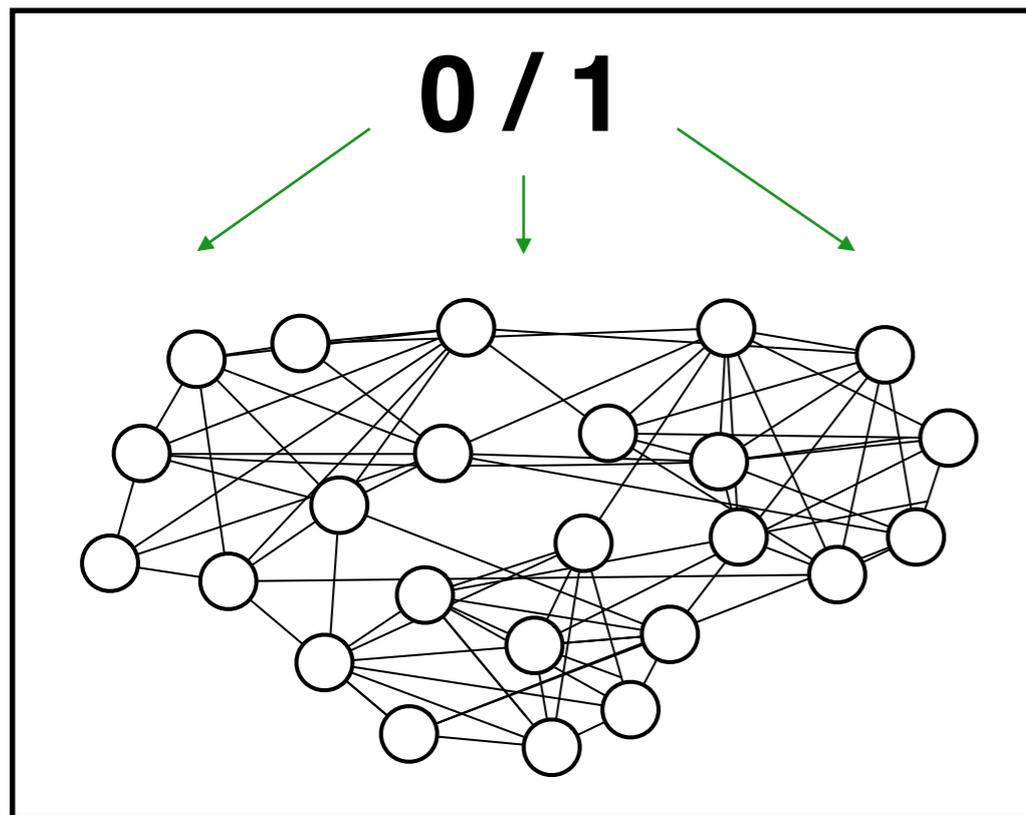


violation of SUTVA

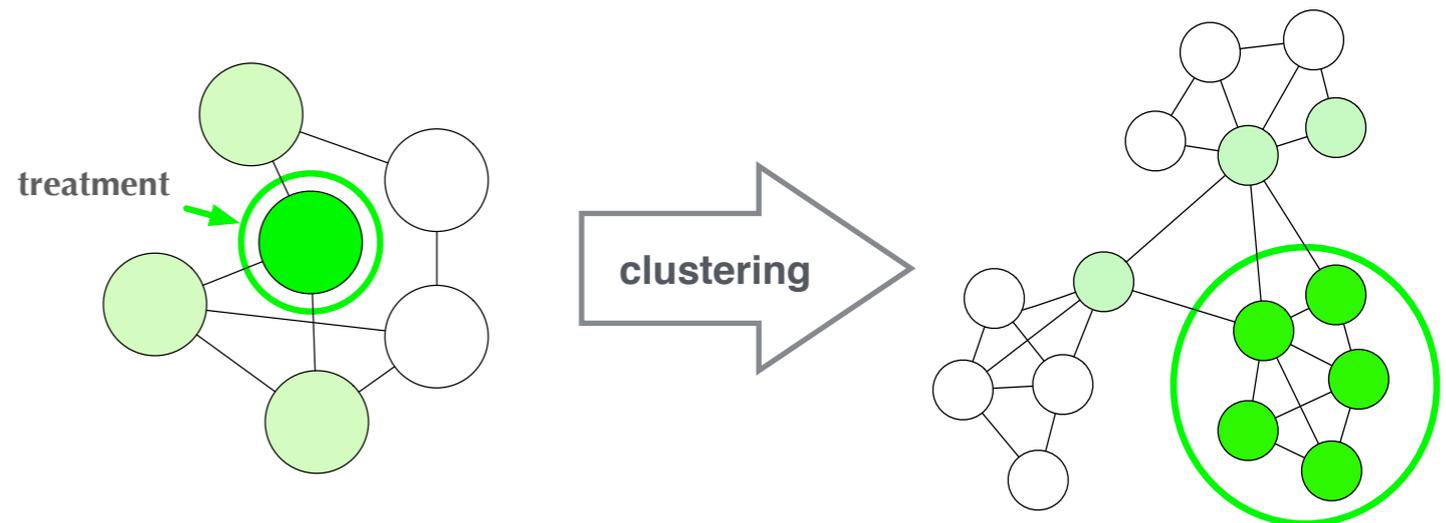
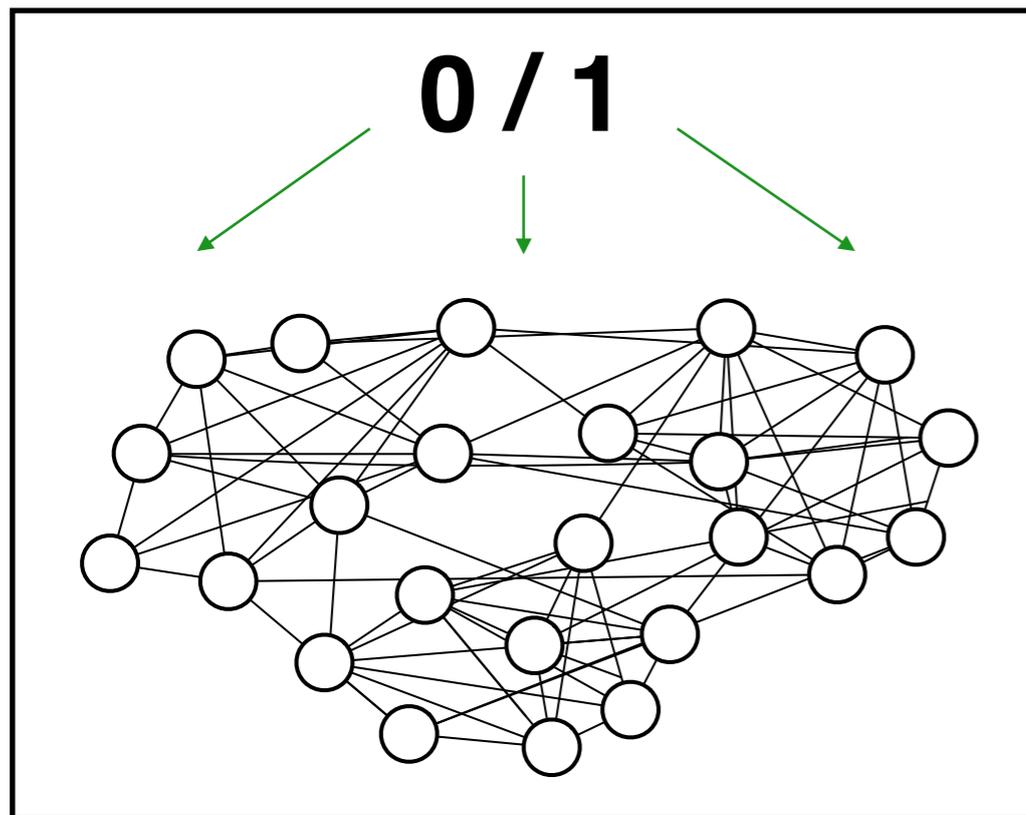
# A/B Testing in Networks: Graph Cluster Randomization

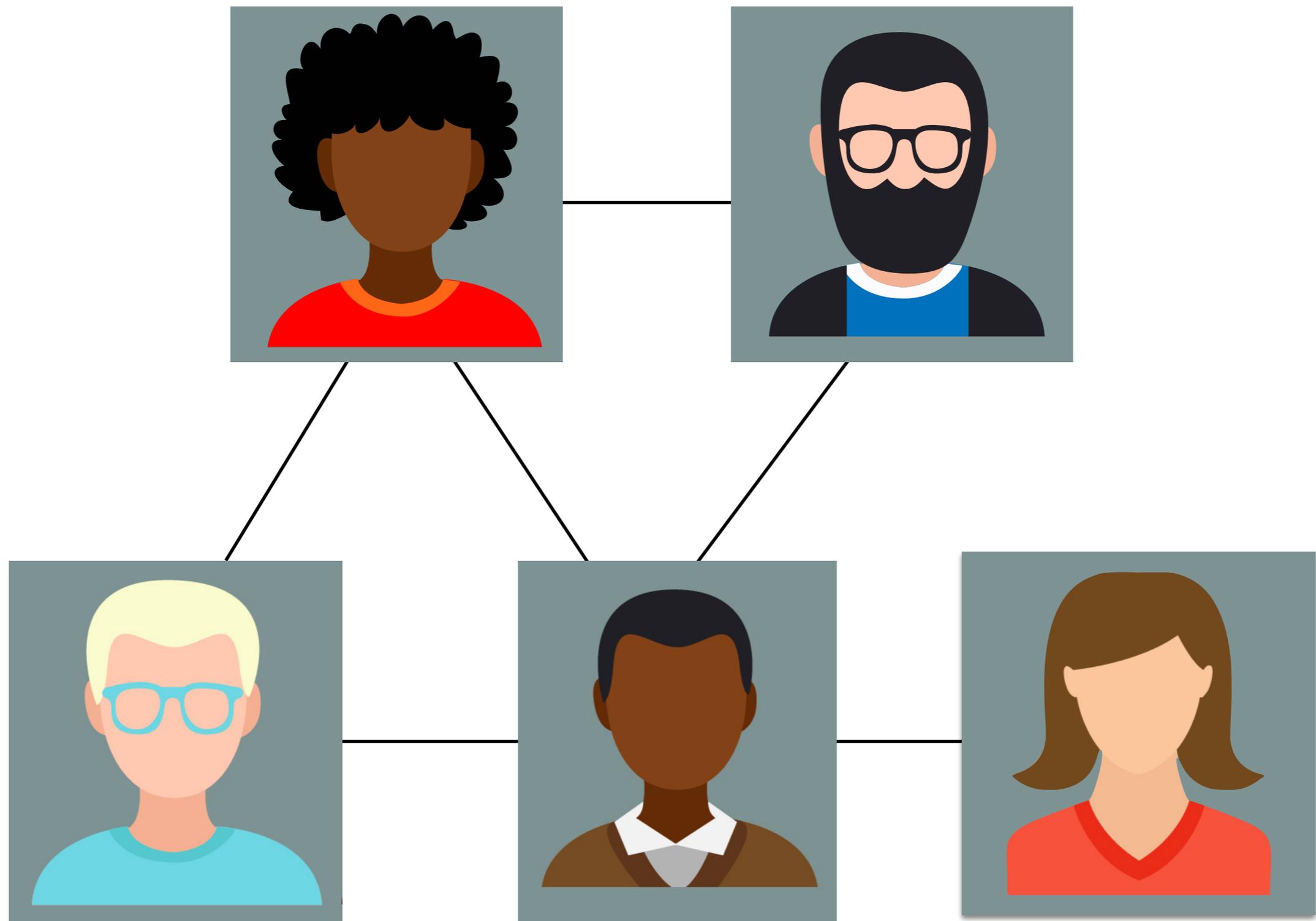


# A/B Testing in Networks: Graph Cluster Randomization

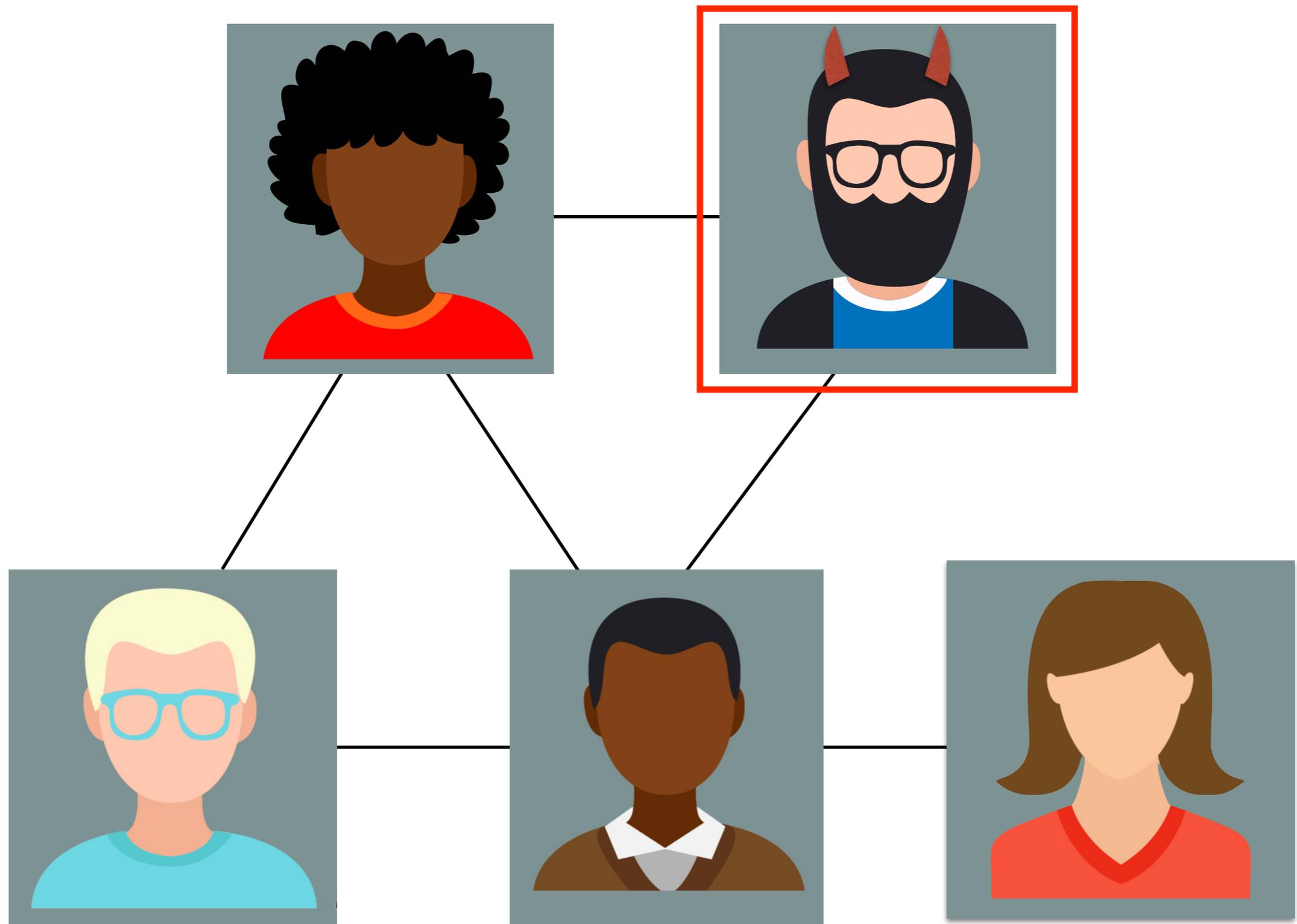


# A/B Testing in Networks: Graph Cluster Randomization

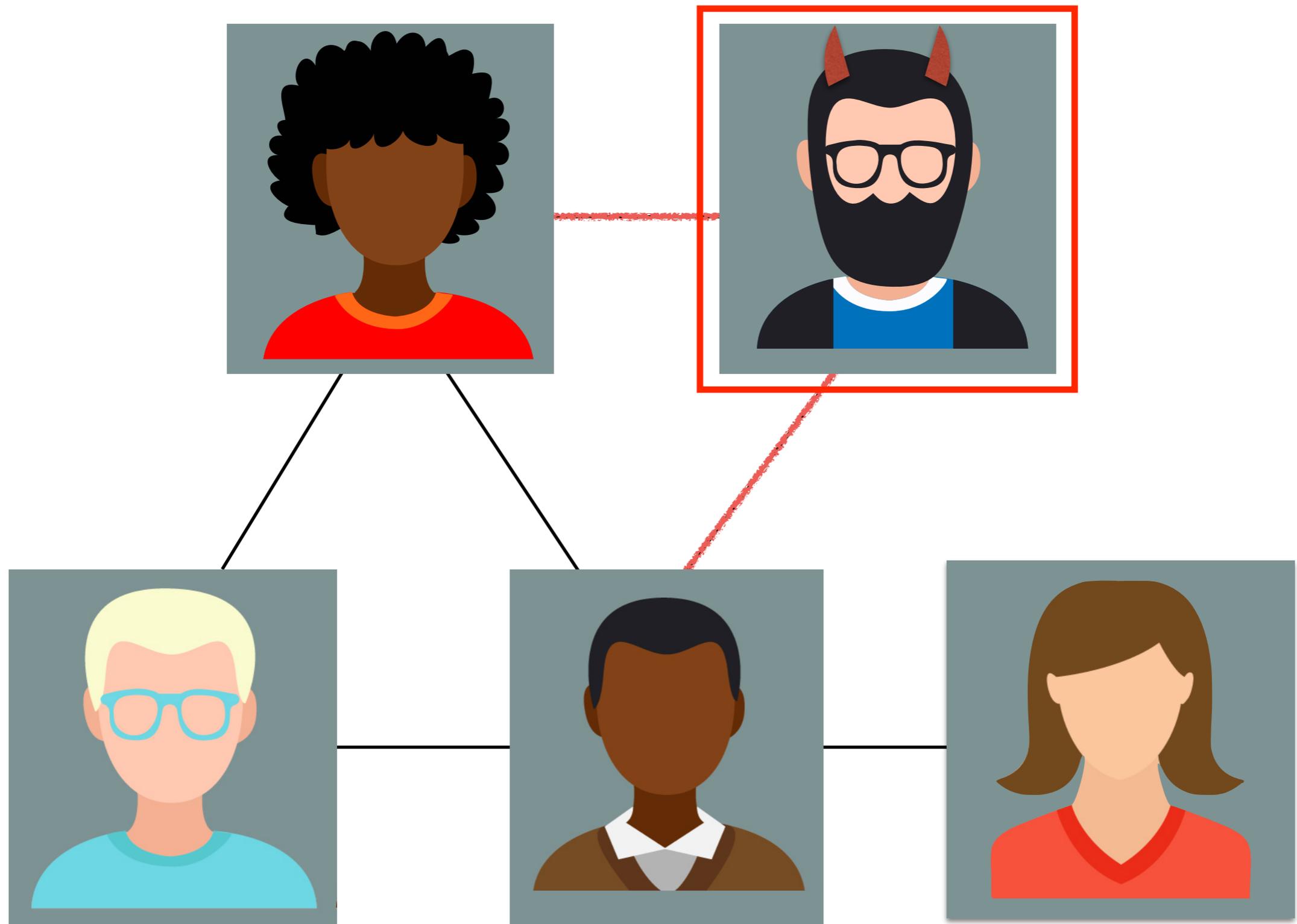




Avatar vector images designed by Freepik



Avatar vector images designed by Freepik



Avatar vector images designed by Freepik

# Adversaries

- Participants in the experiment who would like to influence the estimate of interest, e.g. ATE
  - Increasing or decreasing the estimate of interest
  - Introducing random behavior to the estimate

# Adversaries

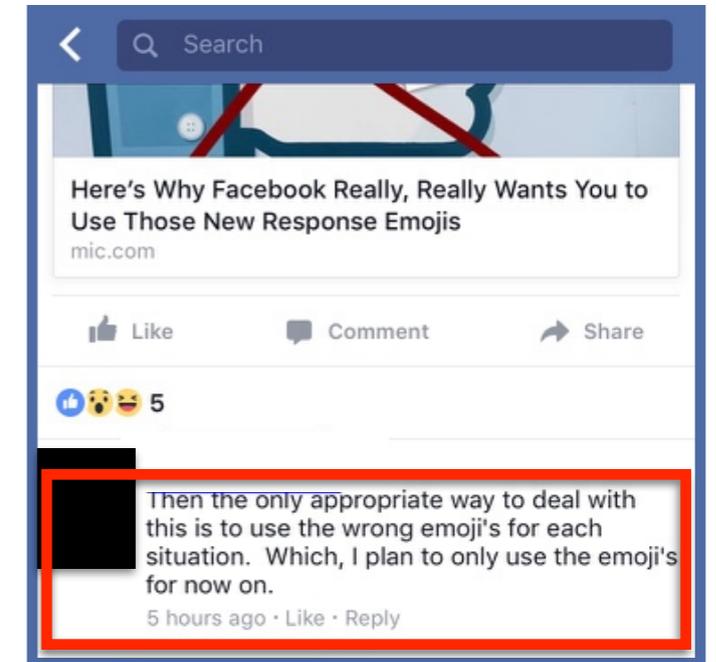
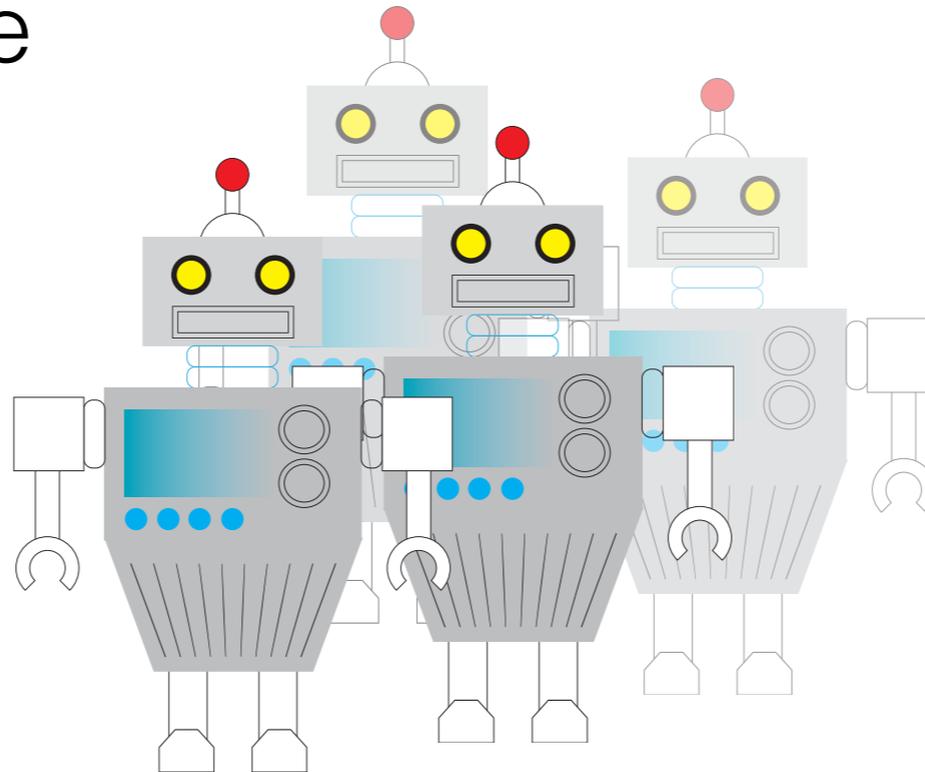
- Participants in the experiment who would like to influence the estimate of interest, e.g. ATE
    - Increasing or decreasing the estimate of interest
    - Introducing random behavior to the estimate
- 
- Introducing **bias**

# Adversaries

- Participants in the experiment who would like to influence the estimate of interest, e.g. ATE
  - Increasing or decreasing the estimate of interest  Introducing **bias**
  - Introducing random behavior to the estimate  Increasing **variance**

# Motivations for Adversaries

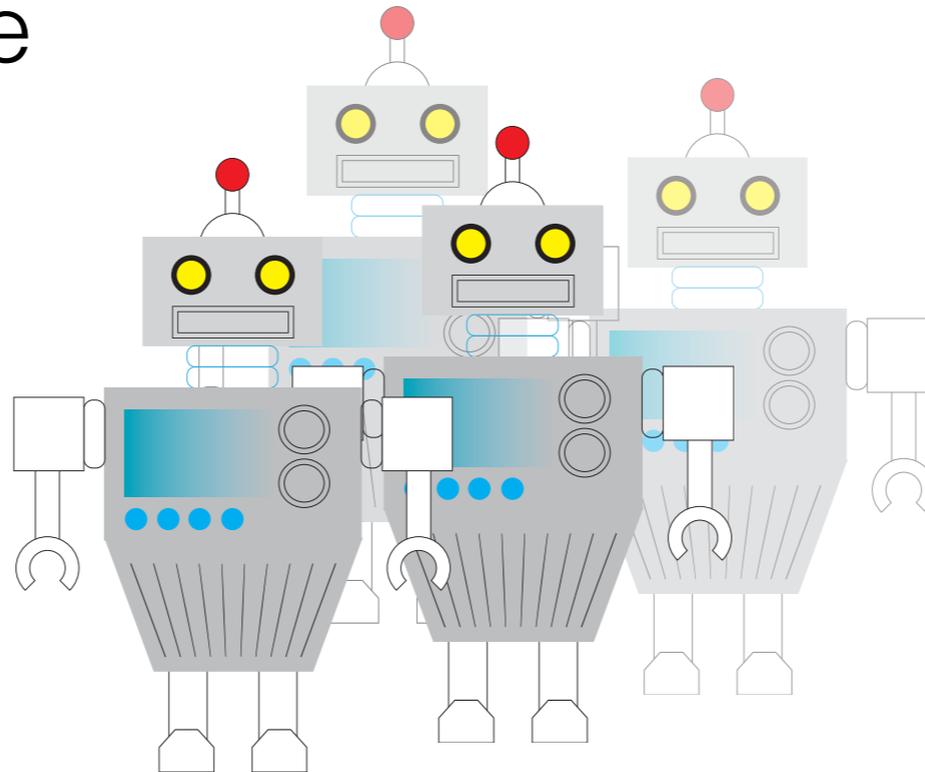
- Competition
- Noncompliance
- Privacy



Avatar vector image designed by Freepik  
Robot vector images designed by Vecteezy

# Motivations for Adversaries

- Competition
- Noncompliance
- Privacy



★ Assume all adversaries have the same behavior model

Avatar vector image designed by Freepik  
Robot vector images designed by Vecteezy

# ATE Estimation in Networks

- Assume outcome is a linear additive function of treatment (Gui et al, 2015)

$$Y_i(Z) = \alpha + \beta z_i + \gamma A_i^T Z + \eta A_i^T Y / D_{ii}$$

$\beta$ : individual treatment effect  
 $\gamma$ : peer treatment effect  
 $\eta$ : peer outcome effect

- Estimate treatment effect from data

$$g(z_i, \sigma_i) = \alpha + \beta z_i + \gamma \sigma_i$$

$z_i$ : treatment assignment  
 $\sigma_i$ : treatment exposure

$$\widehat{ATE} = \hat{\beta} + \hat{\gamma}$$

# ATE Estimation in Networks

- Assume outcome is a linear additive function of treatment (Gui et al, 2015)

$$Y_i(Z) = \alpha + \beta z_i + \gamma A_i^T Z + \eta A_i^T Y / D_{ii}$$

$\beta$ : individual treatment effect  
 $\gamma$ : peer treatment effect  
 $\eta$ : peer outcome effect

- Estimate treatment effect from data

$$g(z_i, \sigma_i) = \alpha + \beta z_i + \gamma \sigma_i$$

$z_i$ : treatment assignment  
 $\sigma_i$ : treatment exposure

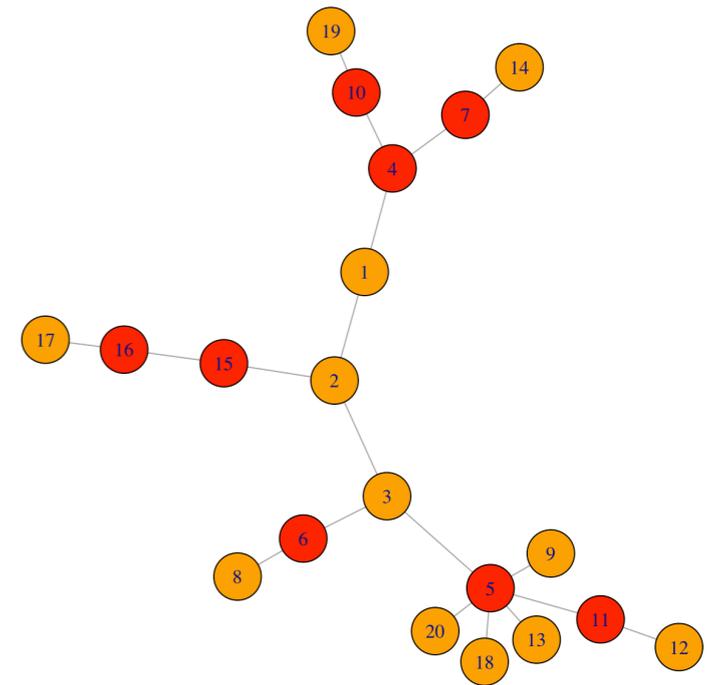
$$\hat{ATE} = \hat{\beta} + \hat{\gamma}$$

# Adversary Placement

(1) Random assignment over the graph

(2) Targeted adversary placement

- When adversaries form a **dominating set** over the graph, every vertex contains at least one adversary in their set of neighbors

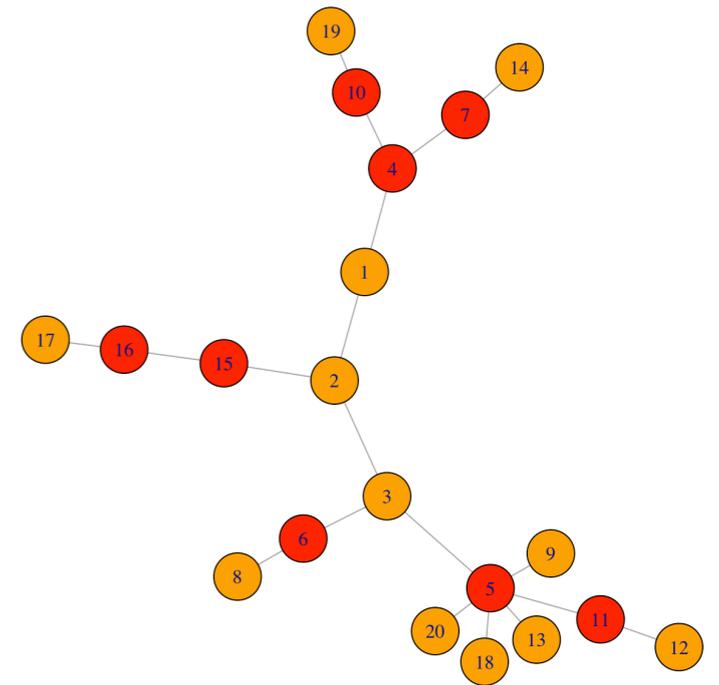


# Adversary Influence

- Adversary **influence** measures the sum of relative effects of a node on its neighbors' outcomes

$$\omega_i = D^{-1} A \mathbb{1}_i^N$$

column sum of transition probabilities

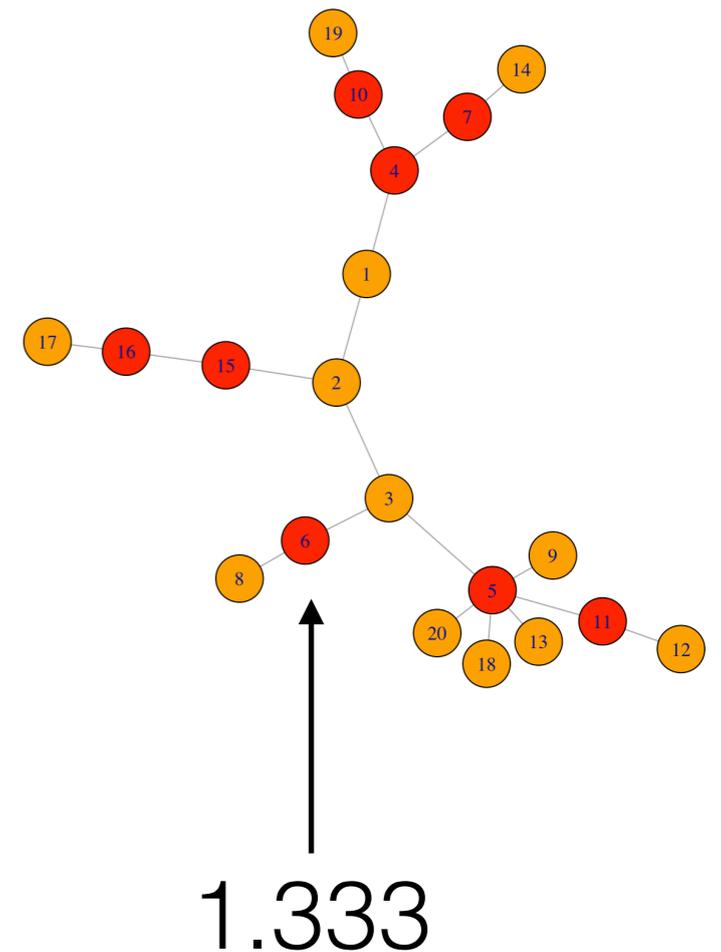


# Adversary Influence

- Adversary **influence** measures the sum of relative effects of a node on its neighbors' outcomes

$$\omega_i = D^{-1} A \mathbb{1}_i^N$$

column sum of transition probabilities

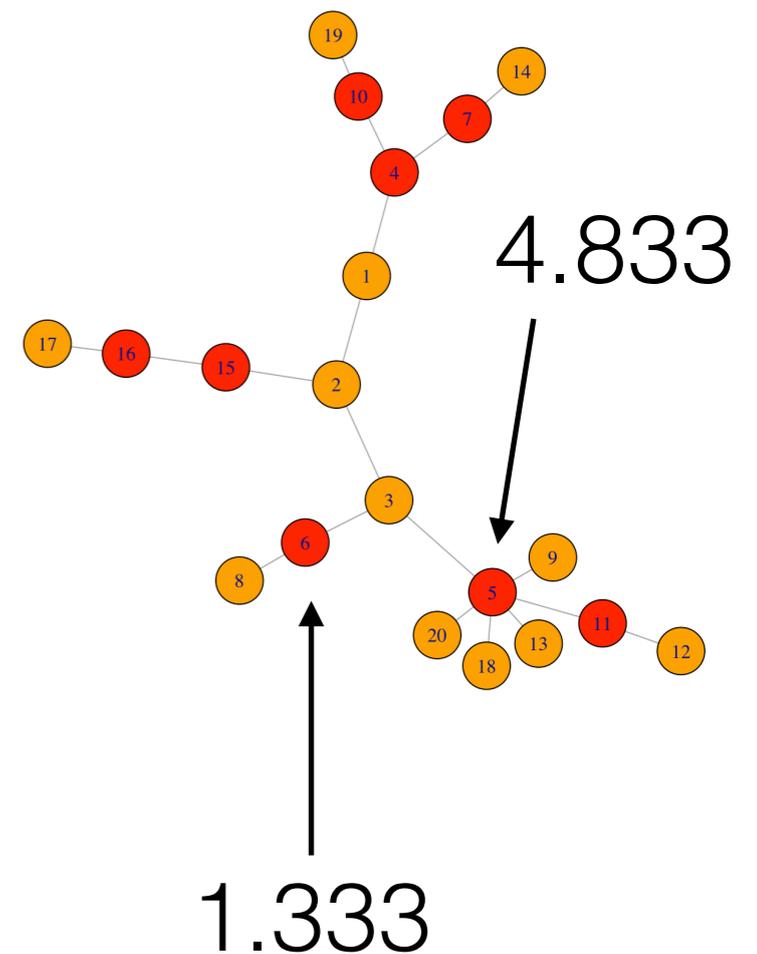


# Adversary Influence

- Adversary **influence** measures the sum of relative effects of a node on its neighbors' outcomes

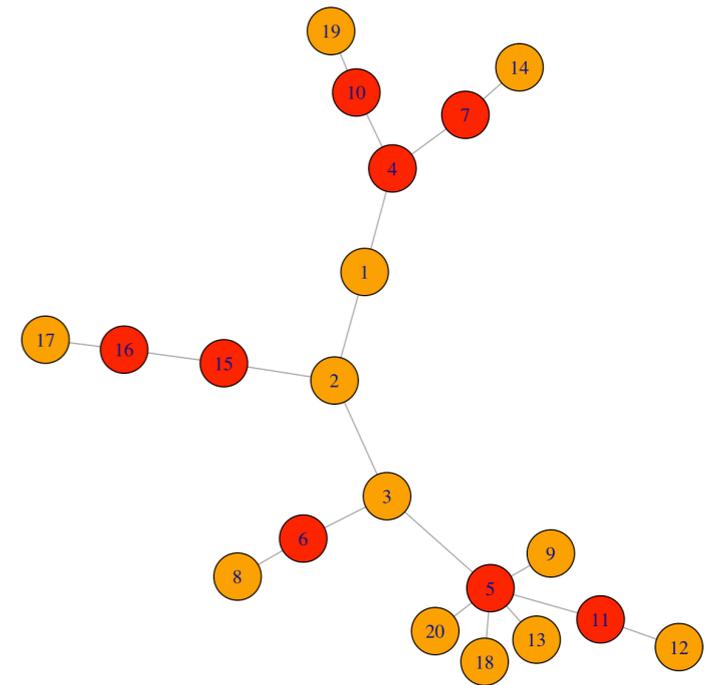
$$\omega_i = D^{-1} A \mathbb{1}_i^N$$

column sum of transition probabilities



# Bias in ATE Induced by Adversaries

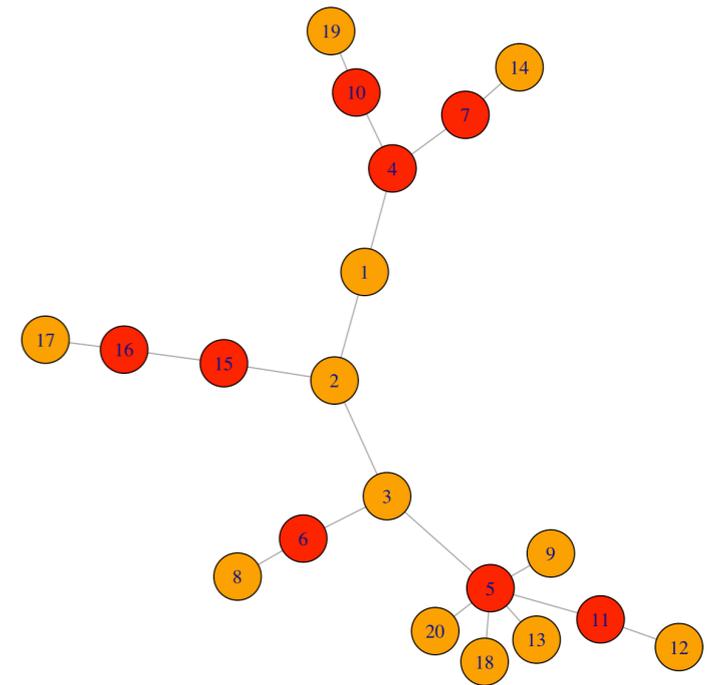
- Adversaries bias ATE estimates through:
  - (1) The value of their outcome
  - (2) The effect of their outcome on their neighbors' outcome
    - strength depends on true network effect



$$\hat{ATE}_{R_Y} = \sum_{j \in A_r} \frac{1}{d_j} (Y_r - \bar{Y}_{A_j \setminus r})$$

# Bias in ATE Induced by Adversaries

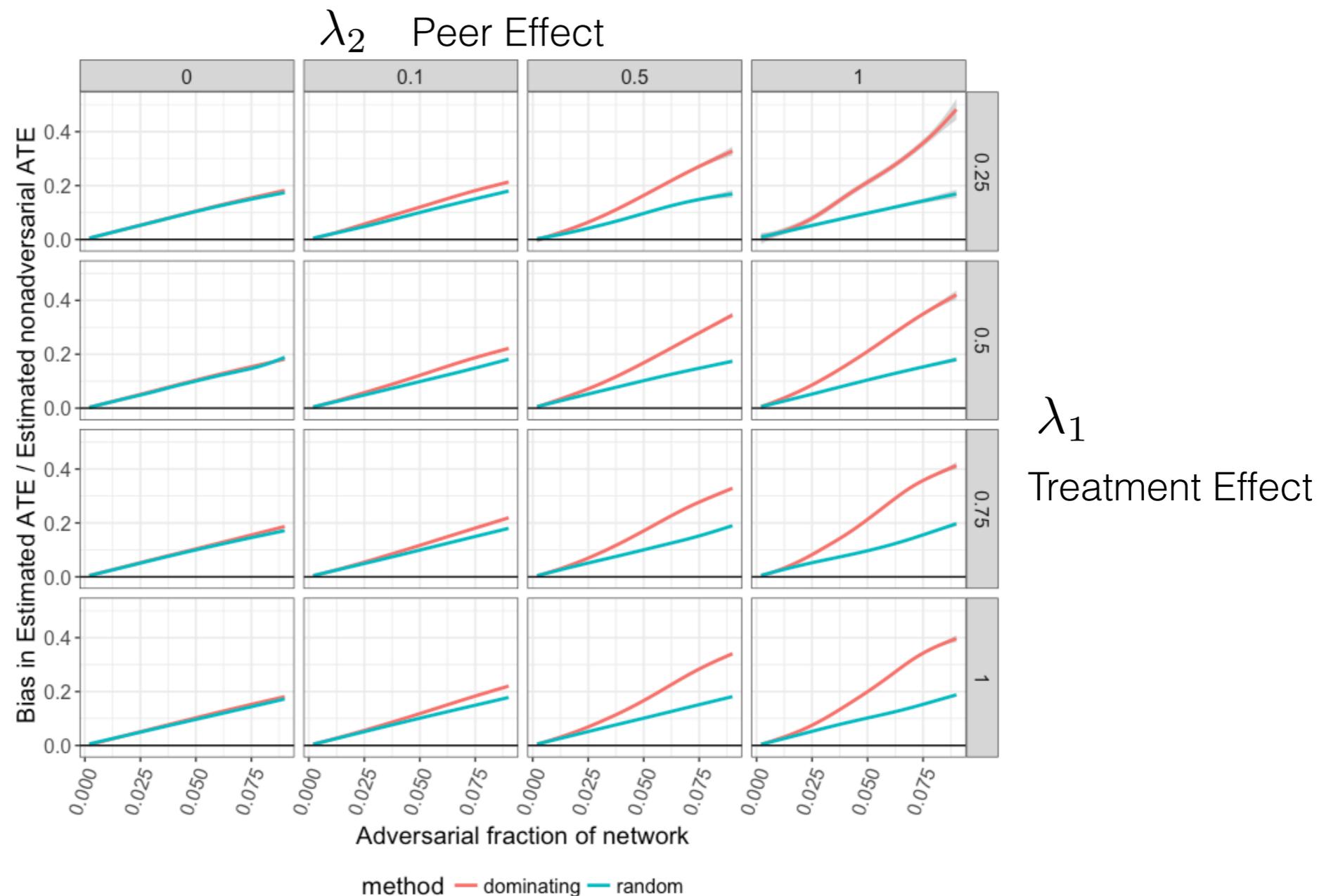
- Adversaries bias ATE estimates through:
  - (1) The value of their outcome
  - (2) The effect of their outcome on their neighbors' outcome
    - strength depends on true network effect



$$\begin{aligned} \hat{ATE}_{R_Y} &= \sum_{j \in A_r} \frac{1}{d_j} (Y_r - \bar{Y}_{A_j \setminus r}) \\ &\approx \omega_r (Y_r - \bar{Y}_{A_{2r}}) \end{aligned}$$

Approximate bias from effect of adversary outcome using influence,  $\omega$

# Experimental Results



Normal

$$Y_{i,t} = \lambda_0 + \lambda_1 z_i + \lambda_2 \frac{A_i Y_{t-1}}{D_{i,i}} + U_{i,t}$$

$$Y_r = \begin{cases} \lambda_0 & \text{if } z_r = 1, \\ \lambda_0 + \lambda_1 & \text{if } z_r = 0. \end{cases}$$

# Summary

- Derived expressions for the bias induced by adversary behavior
- Empirically demonstrated a vulnerability in network A/B testing to manipulation of ATE estimates from exploitation of peer effects
- Examined the difference between random and targeted placement of adversaries in the network
- *Future work*: Characterize the relationship between adversary detection and strength of adversarial response

Thank you