

# Finding the source

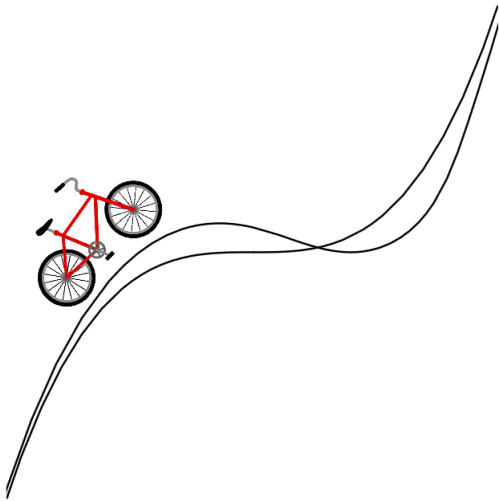
Jacob Richey

joint with: Chris Hoffman, Miki Racz

University of British Columbia

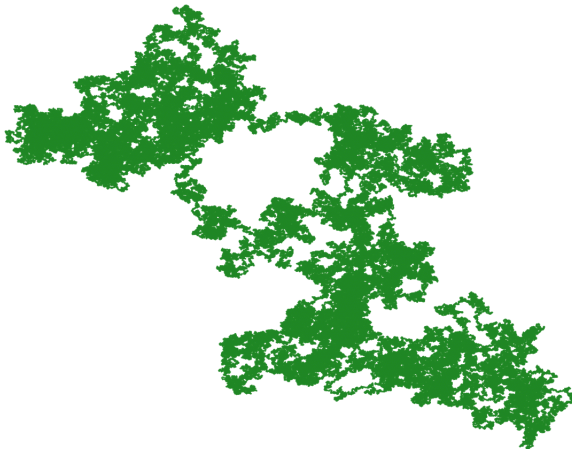
*[jfrichey@math.ubc.edu](mailto:jfrichey@math.ubc.edu)*

Dartmouth, November 2020



Which way did the bicycle go?

Simple random walk on  $\mathbb{Z}^2$ , run for  $5 \cdot 10^6$  steps.



**Q:** Given a snapshot of a (random) process, what can be determined?

**Q:** Given a snapshot of a (random) process, what can be determined?

- Starting/ending point?
- Most/least visited points?
- Step distribution/generator?
- Properties of the underlying graph?

**Q:** Given a snapshot of a (random) process, what can be determined?

- Starting/ending point?
- Most/least visited points?
- Step distribution/generator?
- Properties of the underlying graph?

Warmup: simple random walk on  $\mathbb{Z}$ .

**Problem:** Given the set of visited sites, guess the starting point.

Warmup: simple random walk on  $\mathbb{Z}$ .

**Problem:** Given the set of visited sites, guess the starting point.





Warmup: simple random walk on  $\mathbb{Z}$ .

**Problem:** Given the set of visited sites, guess the starting point.



Given that the range is an interval of length  $N$ , what's the most likely starting point? Purple, red, or green?

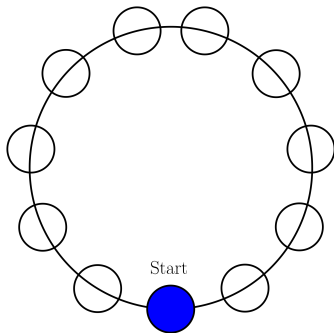
**A:** They're all equally likely!

Proof sketch: think of the range as a 'coin switching' markov chain, compute transition probabilities recursively.

**A:** They're all equally likely!

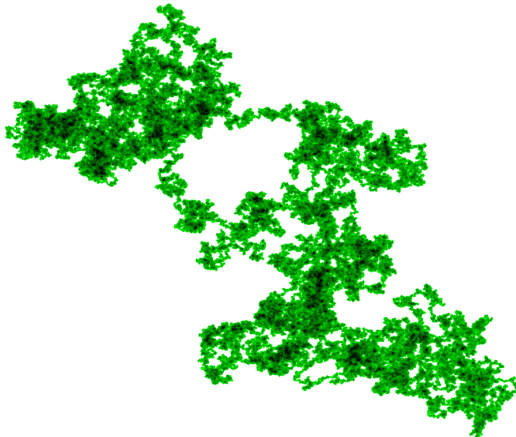
Proof sketch: think of the range as a 'coin switching' markov chain, compute transition probabilities recursively.

Alternatively: last vertex visited by SRW on the ring is uniform.



# Previous work

Same SRW as before, with occupation times



## Previous work

Assume *partial* information about the occupation measure.

## Previous work

Assume *partial* information about the occupation measure.

(Warren, Yor '98) Brownian burgler: BM conditioned on local times

## Previous work

Assume *partial* information about the occupation measure.

(Warren, Yor '98) Brownian burgler: BM conditioned on local times

Theorem (Pemantle, Peres, Pitman, Yor '00)

Let  $d \geq 3$ , and consider Brownian motion in  $\mathbb{R}^d$  run for time 1.

Given the *occupation measure* of the path projected onto the sphere, you can recover the *range* and the *endpoint* with probability 1.

## Previous work

Assume *partial* information about the occupation measure.

(Warren, Yor '98) Brownian burglar: BM conditioned on local times

### Theorem (Pemantle, Peres, Pitman, Yor '00)

*Let  $d \geq 3$ , and consider Brownian motion in  $\mathbb{R}^d$  run for time 1.*

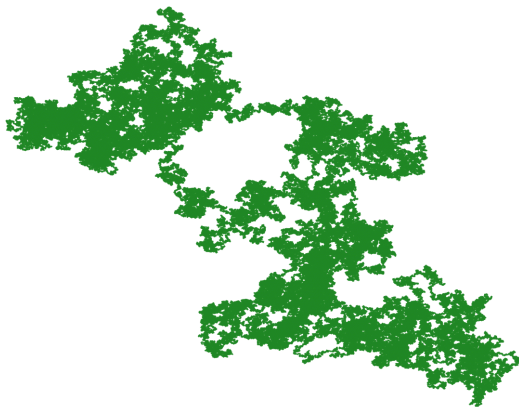
*Given the **occupation measure** of the path projected onto the sphere, you can recover the **range** and the **endpoint** with probability 1.*

### Conjecture (PPPY '00)

*In dimension  $d = 2$ , the range cannot be recovered with probability 1.*

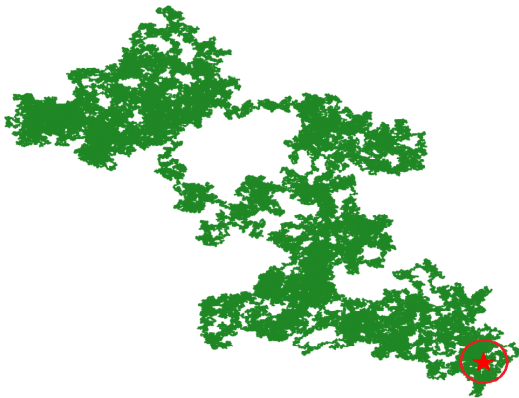


SRW in  $\mathbb{Z}^d$



**Q:** where is the starting point?

SRW in  $\mathbb{Z}^d$



$R_t$  = range of SRW up to time  $t$ .

### Definition

An **estimator**  $\hat{v}$  for the start point is a random element of  $R_t$ .

$R_t$  = range of SRW up to time  $t$ .

### Definition

An **estimator**  $\hat{v}$  for the start point is a random element of  $R_t$ .

Example:  $\widehat{v_{CM}}$  = closest point to the center of mass of  $R_t$ .

$R_t$  = range of SRW up to time  $t$ .

### Definition

An **estimator**  $\hat{v}$  for the start point is a random element of  $R_t$ .

Example:  $\widehat{v}_{CM}$  = closest point to the center of mass of  $R_t$ .

### Definition

For  $x \in \mathbb{Z}^d$ , let  $R_t^x$  be the range of an independent SRW started from  $x$ .  
For any estimator  $\hat{v}$ , the **likelihood** of  $\hat{v}$  is

$$L(\hat{v}) = \mathbb{P}(R_t^{\hat{v}} = R_t | R_t).$$

We want to find an estimator with large likelihood. Canonical best guess is the **maximum likelihood estimator**:

$$\hat{v}_{MLE} = \arg \max_{x \in R_t} L(x)$$

We want to find an estimator with large likelihood. Canonical best guess is the **maximum likelihood estimator**:

$$\hat{v}_{MLE} = \arg \max_{x \in R_t} L(x)$$

**Goal:** Find a convenient estimator  $\hat{v}$  with

$$L(\hat{v}) \approx L(\hat{v}_{MLE})$$

## Theorem (Hoffman, R.)

The following hold for SRW in  $\mathbb{Z}^d$  as  $t \rightarrow \infty$ .

i. If  $d = 2$ ,

$$\frac{L(\hat{v}_{MLE})}{\sum_{w \in R_T} L(w)} \rightarrow_p 0$$

ii. If  $d \in \{3, 4, 5, 6\}$ , there exists an estimator  $\hat{v}$  such that

$$\mathbb{P}(\hat{v} = 0) \geq \Theta(t^{-c_d})$$

for some constant  $c_d \in (0, 1)$ .

iii. If  $d \geq 7$ , there exists an estimator  $\hat{u}$  such that

$$\mathbb{P}(\hat{u} = 0) = \Theta(1).$$



## Conjecture

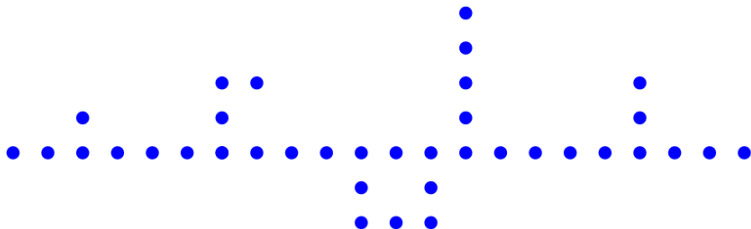
$$\mathbb{P}(\hat{v}_{MLE} = 0) = \begin{cases} o(1), & d = 2 \\ \Theta(1), & d \geq 5 \end{cases}$$

Further Q's:

- SRW on  $d$ -regular tree, biased RW
- Performance of 'longest path' estimator for transient RW's?
- Good estimator for  $\mathbb{Z}^2$ ?

Proof ideas:

- ① Get rid of the 'middle' of the range, using [transience](#).
- ② Infer chronological info using 'cut points.'



Proof sketch:

- ① Get rid of the 'middle' of the range, using [transience](#).
- ② Infer chronological info using 'cut points.'



Ingredients:

- ① Long cycles: return probabilities / self-intersection exponents (Lawler)

Ingredients:

- ① Long cycles: return probabilities / self-intersection exponents (Lawler)
- ② A **cut time** for  $X$  is a time  $s \in [0, t]$  such that

$$X_{[0,s)} \cap X_{(s,t]} = \emptyset$$

If  $s$  is a cut time,  $X_s$  is called a **cut point**.

**Theorem (James, Peres, '96)**

*In dimension  $d \geq 3$ , there are infinitely many cut times. In dimension  $d \geq 5$ , cut times have positive density.*

Cutpoints are totally ordered (by their cut times).

Given all the cut points, find the 'first' and 'last' ones, pick uniformly from their small components.

Cutpoints are totally ordered (by their cut times).

Given all the cut points, find the 'first' and 'last' ones, pick uniformly from their small components.

**Problem:** not all 'divider' points are cut points!

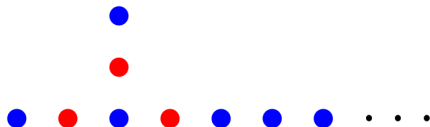


Figure: The three red 'divider' points can't all be cut points.

Cutpoints are totally ordered (by their cut times).

Given all the cut points, find the 'first' and 'last' ones, pick uniformly from their small components.

**Problem:** not all 'divider' points are cut points!

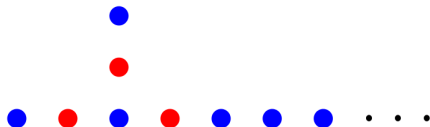


Figure: The three red 'divider' points can't all be cut points.

Need more information about how cutpoints are distributed.



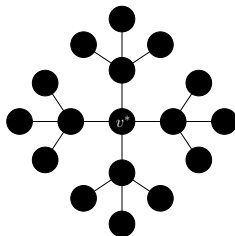
Consider a rumor spreading through a network.

- The rumor starts from a 'source' vertex
- At integer times, nodes can pass the rumor to neighbors
- Observer sees which nodes have the rumor at some time

Consider a rumor spreading through a network.

- The rumor starts from a 'source' vertex
- At integer times, nodes can pass the rumor to neighbors
- Observer sees which nodes have the rumor at some time

For today, focus on the  $d$ -regular tree.



The rumor is spread by a random algorithm known to the observer.

Goals for the rumor spreader:

- *Spreading*: spread the rumor to many nodes
- *Obfuscation*: minimize the probability that the observer guesses the source correctly
- *Multiple observations*: obfuscate the source even when the observer has access to multiple independent rumors

The rumor is spread by a random algorithm known to the observer.

Goals for the rumor spreader:

- *Spreading*: spread the rumor to many nodes
- *Obfuscation*: minimize the probability that the observer guesses the source correctly
- *Multiple observations*: obfuscate the source even when the observer has access to multiple independent rumors
- *Local spreading* (new): ensure that all vertices close to the source learn the rumor quickly

Anonymous messaging platforms, e.g. Secret, Yik Yak, Whisper

Obfuscating the source  $\leftrightarrow$  protecting user data

Anonymous messaging platforms, e.g. Secret, Yik Yak, Whisper

Obfuscating the source  $\leftrightarrow$  protecting user data

Contact tracing / finding patient zero

Previous work: SI/SIR. MLE well understood. Rumor centrality

New algorithm: **adaptive diffusion**

- $G = d$ -regular tree
- $G_t =$  set of nodes that know the rumor at time  $t$
- $vs_t =$  virtual source at time  $t$
- $G_t$  is a ball of radius  $t/2$  centered at  $vs_t$  at even times  $t$
- Defined by transition probabilities  $\alpha(t, h)$  for the virtual source

Virtual source evolves according to:

- Start with  $vs_0 = v^*$



Virtual source evolves according to:

- Start with  $vs_0 = v^*$
- $vs_2$  is a uniform neighbor of  $v^*$ .

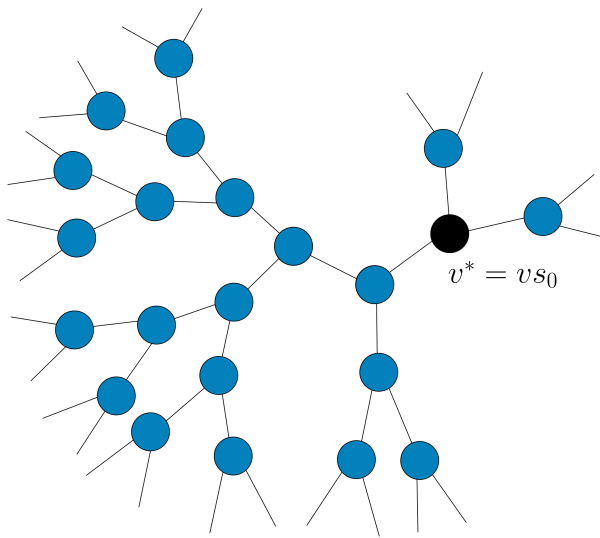
Virtual source evolves according to:

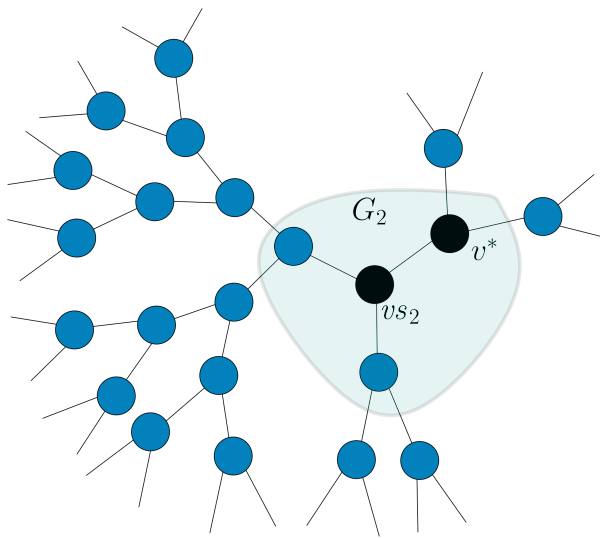
- Start with  $vs_0 = v^*$
- $vs_2$  is a uniform neighbor of  $v^*$ .
- Let  $h = \text{dist}(vs_t, v^*)$
- Probability  $\alpha(t, h)$ :  $vs_{t+2} =$  uniform neighbor of  $vs_t$  excluding previous virtual sources
- Probability  $1 - \alpha(t, h)$ :  $vs_{t+2} = vs_t$

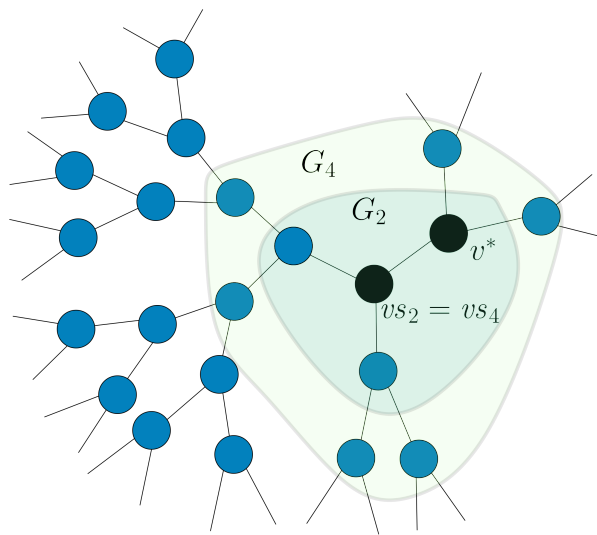
Virtual source evolves according to:

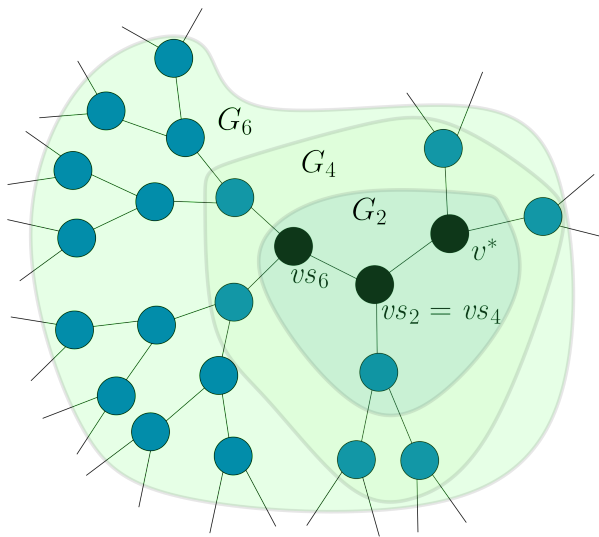
- Start with  $vs_0 = v^*$
- $vs_2$  is a uniform neighbor of  $v^*$ .
- Let  $h = \text{dist}(vs_t, v^*)$
- Probability  $\alpha(t, h)$ :  $vs_{t+2} =$  uniform neighbor of  $vs_t$  excluding previous virtual sources
- Probability  $1 - \alpha(t, h)$ :  $vs_{t+2} = vs_t$

When the virtual source moves, it always moves in a uniform direction away from  $v^*$ .









Equivalently, work with  $p(t, h) = \mathbb{P}(\text{dist}(vs_t, v^*) = h)$ .



Equivalently, work with  $p(t, h) = \mathbb{P}(\text{dist}(vs_t, v^*) = h)$ .

MLE for the source vertex:

$$\hat{v}_{MLE} = \arg \max_{v \in G_t} L(v),$$

where  $L(v) = \mathbb{P}(G_t^v = G_t | G_t)$ .

Equivalently, work with  $p(t, h) = \mathbb{P}(\text{dist}(vs_t, v^*) = h)$ .

MLE for the source vertex:

$$\hat{v}_{MLE} = \arg \max_{v \in G_t} L(v),$$

where  $L(v) = \mathbb{P}(G_t^v = G_t | G_t)$ .

## Fact

$\hat{v}_{MLE}$  is uniform over all vertices at distance  $h^*$  from  $vs_t$ , where

$$h^* = \arg \max_{h \in \{1, 2, \dots, t/2\}} \frac{p(t, h)}{(d-1)^h}.$$

## Spreading

For adaptive diffusion,

$$|G_t| = N_t = \frac{1}{d-2}(d-1)^{t/2}.$$

deterministically at even times  $t$ . (Order-optimal spreading)

## Obfuscation

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = \begin{cases} \Theta(N_t^{-1}) & \text{(perfect obfuscation)} \\ \Theta(N_t^{-\gamma}) & \text{(polynomial obfuscation)} \\ o(1) & \text{(weak obfuscation)} \end{cases}$$

SI/SIR: good spreading, weak obfuscation; not even weak obfuscation under multiple observations. [Shah, Zaman, Dong, Tan, Wang, Zhang]

### Adaptive diffusion (Fanti, Kairouz, Oh, Viswanath '15)

Let  $G = d$ -regular tree. There exists an adaptive diffusion algorithm that achieves **perfect obfuscation**:

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = \Theta(N_t^{-1})$$

SI/SIR: good spreading, weak obfuscation; not even weak obfuscation under multiple observations. [Shah, Zaman, Dong, Tan, Wang, Zhang]

### Adaptive diffusion (Fanti, Kairouz, Oh, Viswanath '15)

Let  $G = d$ -regular tree. There exists an adaptive diffusion algorithm that achieves **perfect obfuscation**:

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = \Theta(N_t^{-1})$$

*Proof sketch:* choose  $p(t, h) \sim (d - 1)^h$ , so the MLE picks a uniform random vertex in  $G_t$ . Show this is realizable for some values  $\alpha(t, h)$ .

**Q:** Does it have good local spreading?

**Q:** Does it have good local spreading?

### Definition

The *local spread*  $R_t$  is the radius of the largest ball centered at  $v^*$  and contained in  $G_t$ .

For adaptive diffusion,  $R_t = \#$  of times the virtual source has *not* moved:

$$R_t = t/2 - \text{dist}(v^*, v_{s_t}).$$

**Q:** Does it have good local spreading?

### Definition

The *local spread*  $R_t$  is the radius of the largest ball centered at  $v^*$  and contained in  $G_t$ .

For adaptive diffusion,  $R_t = \#$  of times the virtual source has *not* moved:

$$R_t = t/2 - \text{dist}(v^*, v_{s_t}).$$

The algorithm from the theorem doesn't even achieve weak local spread!

$$p(t, h) \sim (d-1)^h \implies \text{dist}(v_{s_t}, v^*) \approx t/2 - O(1).$$



## Spreading/obfuscation trade-off [Racz, Richey '18]

Consider any adaptive diffusion with **polynomial obfuscation** of order  $\gamma \in (0, 1)$ , i.e.

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = O(N_t^{-\gamma}).$$

Then the **local spreading** is bounded from above:

$$\mathbb{E}[R_t] \leq (1 - \gamma) \frac{t}{2} + O(\log t).$$

## Spreading/obfuscation trade-off [Racz, Richey '18]

Consider any adaptive diffusion with **polynomial obfuscation** of order  $\gamma \in (0, 1)$ , i.e.

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = O(N_t^{-\gamma}).$$

Then the **local spreading** is bounded from above:

$$\mathbb{E}[R_t] \leq (1 - \gamma) \frac{t}{2} + O(\log t).$$

Obfuscation and local spreading are **inversely linked** in this case.

The trade-off is essentially tight:

### Spreading/obfuscation trade-off [Racz, Richey '18]

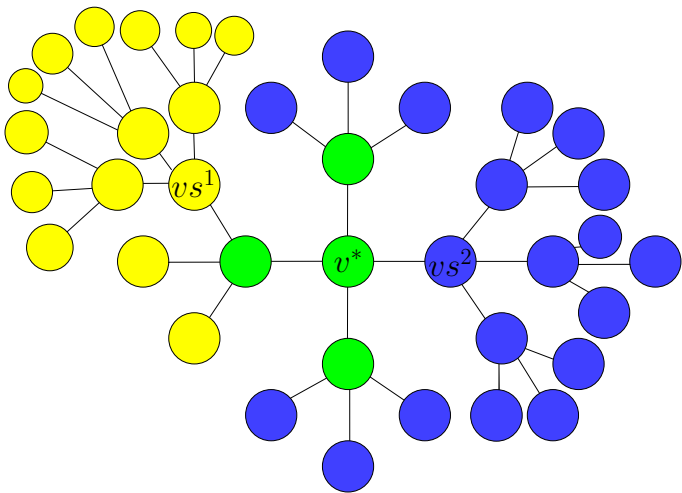
For every  $\gamma \in (0, 1)$ , there exists an adaptive diffusion with both **polynomial obfuscation** of order  $\gamma$ ,

$$\mathbb{P}(\hat{v}_{MLE} = v^*) = O(N_t^{-\gamma}),$$

and **order optimal local spreading**

$$\mathbb{E}[R_t] \geq (1 - \gamma) \frac{t}{2}.$$

Suppose the observer has access to  $k > 1$  independent snapshots  $\{G_t^i\}_{i=1}^k$  of the diffusion started from the same source  $v^*$ .



## Two independent observations (Racz, Richey '18)

Suppose the observer has access to two independent observations  $G_t^1$  and  $G_t^2$  started from a fixed source  $v^*$ . There exists a nice estimator  $\hat{v}$ , not depending on the spreading algorithm, such that for any  $t$ ,

$$\mathbb{P}(\hat{v} = v^*) \geq \frac{d-1}{d} \cdot \frac{2}{t}.$$

Moreover, there exists a protocol such that for any  $t$ ,

$$\mathbb{P}(\hat{v}_{ML} = v^*) \leq \frac{d-1}{d} \cdot \frac{7}{t}.$$

Only **weak obfuscation** now!

It gets worse:

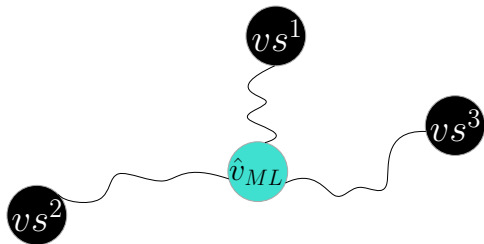
### Three or more independent observations (Racz, Richey '18)

Suppose the observer has access to  $k \geq 3$  independent observations  $G_t^i$ ,  $i \in [k]$ . There exists a nice estimator  $\hat{w}$ , not depending on the spreading algorithm, such that for any  $t$ ,

$$\mathbb{P}(\hat{w} = v^*) \geq 1 - d \exp\left(-\frac{(d-2)^2}{2d^2}k\right).$$

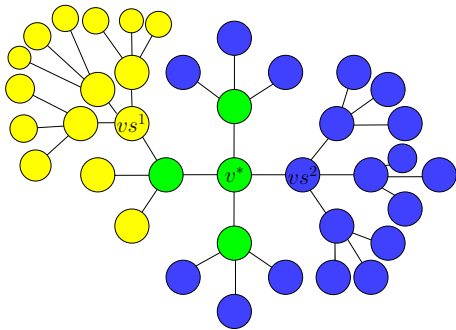
Not even **weak obfuscation**!

*Proof:* Pick any three virtual sources and draw the paths between them.



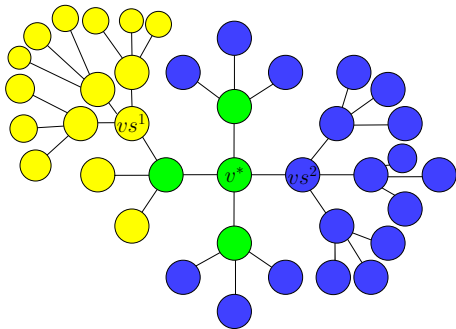
When the three virtual sources lie in different sub-trees away from the root, there will be a unique intersection point  $\hat{w}$ .

Simple estimator: guess a green vertex!





Simple estimator: guess a green vertex!



Obfuscation and local spreading are **positively linked** in this case:

$$\mathbb{P}(\hat{v}_{MLE} = v^*) \geq \mathbb{E} \left[ \left| \bigcap_{i=1}^k G_t^i \right|^{-1} \right]$$

## Question

Does there exist a spreading algorithm that achieves **order-optimal spreading** and **polynomial obfuscation** under  $\geq 2$  observations?

Should look at algorithms that have **order-optimal local spreading**.

Also, need more randomness: adaptive diffusion is indexed by a single vertex. Too simple!