

From trees to seeds:  
on the inference of the seed from large random trees

Joint work with  
Sébastien Bubeck, Ronen Eldan, and Elchanan Mossel

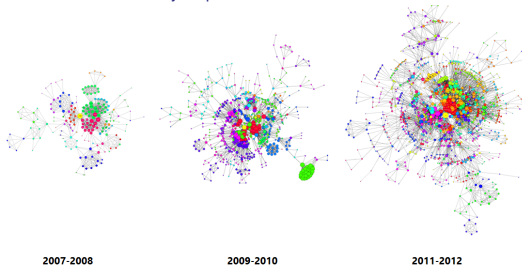
Miklós Z. Rácz

UC Berkeley

University of Bristol  
Probability and Statistics Seminar  
March 27, 2015.

# Statistical inference in non-equilibrium networks

Apple's inventor network over a 6-year period. Source: Kenedict.



## Given the current state of a network, what can we say about a previous state?

### Inferring network mechanisms: The *Drosophila melanogaster* protein interaction network

Manuel Middendorff<sup>1</sup>, Eytan Ziv<sup>2</sup>, and Chris H. Wiggins<sup>3\*</sup>

<sup>1</sup>Department of Physics, <sup>2</sup>College of Physicians and Surgeons, <sup>3</sup>Department of Applied Physics and Applied Mathematics, and <sup>4</sup>Center for Computational Biology and Bioinformatics, Columbia University, New York, NY, 10027

Communicated by Barry H. Honig, Columbia University, New York, NY, December 20, 2004 (invited for review September 7, 2006)

**Naturally occurring networks exhibit quantitative features revealing underlying growth mechanisms. Numerous network mechanisms assess significance of given subgraphs relative to an assumed null model, generated by Monte Carlo sampling of**

### Recovering time-varying networks of dependencies in social and biological studies

Amer Ahmed and Eric P. Xing<sup>1\*</sup>

<sup>1</sup>Language Technology Institute and Machine Learning Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213

Edited by Douglas E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, and approved April 26, 2006 (invited for review February 23, 2006)

**A plausible representation of the relational information among entities in dynamic systems such as a living cell or a social community is a stochastic network that is topologically rewiring and semantically evolving over time. Although there is a rich literature underlying this phenomenon is the unavailability of social snapshots of the rewiring network during the unfolding and progression of the**

OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

## Not All Scale-Free Networks Are Born Equal: The Role of the Seed Graph in PPI Network Evolution

Fereydoun Hormozdizadeh<sup>1</sup>, Petra Berenbrink<sup>2</sup>, Nataša Pržulj<sup>3</sup>, S. Cenk Sahinalp<sup>1\*</sup>

<sup>1</sup> School of Computing Science, Simon Fraser University, Burnaby, British Columbia, Canada, <sup>2</sup> Department of Computer Science, University of California Irvine, California, United States of America

OPEN ACCESS Freely available online

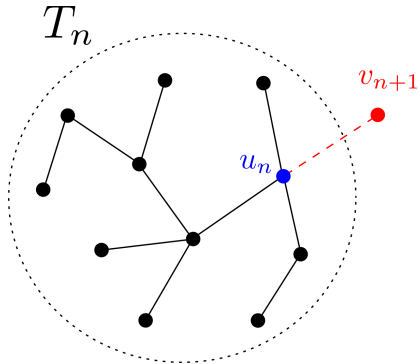
PLOS COMPUTATIONAL BIOLOGY

## Network Archaeology: Uncovering Ancient Networks from Present-Day Interactions

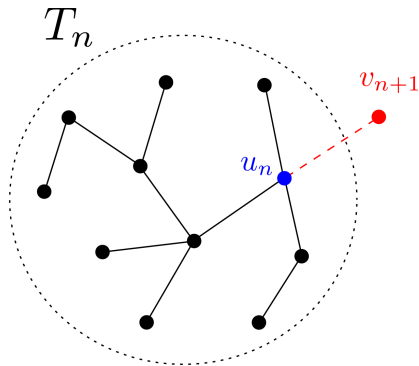
Saket Navlakha, Carl Kingsford<sup>1\*</sup>

<sup>1</sup> Department of Computer Science and Center for Bioinformatics and Computational Biology, University of Maryland, College Park, Maryland, United States of America

# Randomly growing trees



# Randomly growing trees



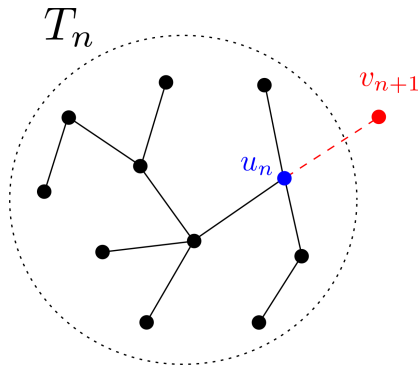
Preferential attachment:

$$\mathbb{P}(u_n = u) = \frac{d_{T_n}(u)}{2n-2}$$

Uniform attachment:

$$\mathbb{P}(u_n = u) = \frac{1}{n}$$

# Randomly growing trees



Preferential attachment:

$$\mathbb{P}(u_n = u) = \frac{d_{T_n}(u)}{2n-2}$$

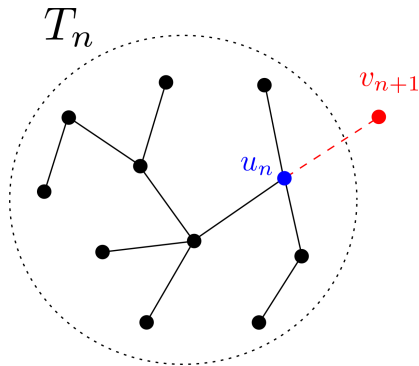
Uniform attachment:

$$\mathbb{P}(u_n = u) = \frac{1}{n}$$

In general:

$$\mathbb{P}(u_n = u) = \frac{(d_{T_n}(u))^\alpha}{Z}$$

# Randomly growing trees



Preferential attachment:

$$\mathbb{P}(u_n = u) = \frac{d_{T_n}(u)}{2n - 2}$$

Uniform attachment:

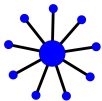
$$\mathbb{P}(u_n = u) = \frac{1}{n}$$

In general:

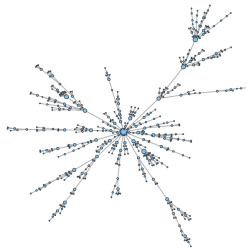
$$\mathbb{P}(u_n = u) = \frac{(d_{T_n}(u))^\alpha}{Z}$$

Many other tree growth models...

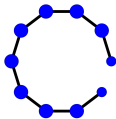
# The influence of the seed — preferential attachment



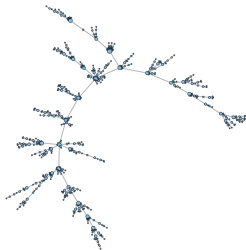
seed  $S_{10}$



PA ( $n = 500, S_{10}$ )

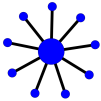


seed  $P_{10}$

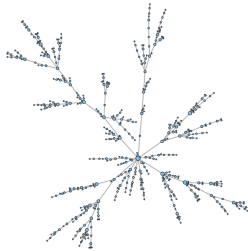


PA ( $n = 500, P_{10}$ )

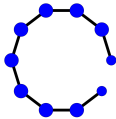
# The influence of the seed — uniform attachment



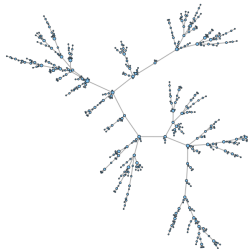
seed  $S_{10}$



UA ( $n = 500, S_{10}$ )



seed  $P_{10}$



UA ( $n = 500, P_{10}$ )



# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree

# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree
  - ↪ seed has no influence for PA or UA

# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree

↪ seed has no influence for PA or UA

But for superlinear attachment ( $\alpha > 1$ ),  
see Oliveira, Spencer (2005)

# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree

↪ seed has no influence for PA or UA

But for superlinear attachment ( $\alpha > 1$ ),  
see Oliveira, Spencer (2005)

- ▶ A finer measure: weak local limit (Benjamini-Schramm)

# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree

↪ seed has no influence for PA or UA

But for superlinear attachment ( $\alpha > 1$ ),  
see Oliveira, Spencer (2005)

- ▶ A finer measure: weak local limit (Benjamini-Schramm)

↪ seed has no influence for PA or UA

# Measuring the influence of the seed

- ▶ A crude measure: limit as a countably infinite tree

↪ seed has no influence for PA or UA

But for superlinear attachment ( $\alpha > 1$ ),  
see Oliveira, Spencer (2005)

- ▶ A finer measure: weak local limit (Benjamini-Schramm)

↪ seed has no influence for PA or UA

See Rudas, Tóth, Valkó (2007) (PA trees)  
and Berger, Borgs, Chayes, Saberi (2014) (in general)  
for weak local limits.

# Measuring the influence of the seed

- ▶ A much finer measure: total variation distance

$$\delta_{\text{PA}}(\mathcal{S}, \mathcal{T}) := \lim_{n \rightarrow \infty} \text{TV}(\text{PA}(n, \mathcal{S}), \text{PA}(n, \mathcal{T}))$$

$$\delta_{\text{PA}}(\text{circle}, \text{star}) = \lim_{n \rightarrow \infty} \text{TV}(\text{fractal circle}, \text{fractal star})$$

# Measuring the influence of the seed

- ▶ A much finer measure: total variation distance

$$\delta_{\text{PA}}(\mathcal{S}, \mathcal{T}) := \lim_{n \rightarrow \infty} \text{TV}(\text{PA}(n, \mathcal{S}), \text{PA}(n, \mathcal{T}))$$

$$\delta_{\text{PA}}(\text{circle}, \text{star}) = \lim_{n \rightarrow \infty} \text{TV}(\text{fractal circle}, \text{fractal star})$$

Hypothesis testing question:

$$H_0 : R \sim \text{PA}(n, \mathcal{S}), \quad H_1 : R \sim \text{PA}(n, \mathcal{T})$$

Q: test with asymptotically (in  $n$ ) non-negligible power?



# Main results

## Preferential attachment:

**Theorem (Bubeck-Mossel-R., arXiv:1401.4849v3, March 2014)**

If the degree profiles of  $S$  and  $T$  are different, and both have at least 3 vertices, then

$$\delta_{\text{PA}}(S, T) > 0.$$

**Theorem (Curien-Duquesne-Kortchemski-Manolescu, June 2014)**

If  $S$  and  $T$  are non-isomorphic and both have at least 3 vertices, then

$$\delta_{\text{PA}}(S, T) > 0.$$

## Uniform attachment:

**Theorem (Bubeck-Eldan-Mossel-R., arXiv:1409.7685, Sept. 2014)**

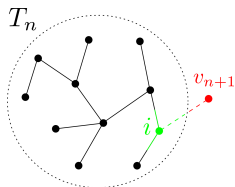
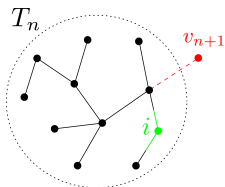
If  $S$  and  $T$  are non-isomorphic and both have at least 3 vertices, then

$$\delta_{\text{UA}}(S, T) > 0.$$

# PA heuristics: maximum degree

Degree evolution governed by Pólya urns

$$(2n - 2 - d_{\text{PA}(n,S)}(i), d_{\text{PA}(n,S)}(i))$$

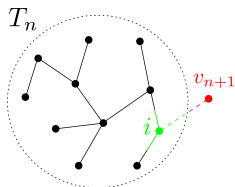
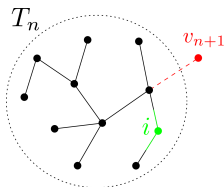


- ▶ Replacement matrix:  $\begin{pmatrix} 2 & 0 \\ 1 & 1 \end{pmatrix}$ ; **initial condition:**
  - ▶ If  $i \in S$  then  $(2|S| - 2 - d_S(i), d_S(i))$ ;
  - ▶ If  $i \notin S$  then  $(2i - 3, 1)$ .

# PA heuristics: maximum degree

Degree evolution governed by **Pólya urns**

$$(2n - 2 - d_{\text{PA}(n,S)}(i), d_{\text{PA}(n,S)}(i))$$



- ▶ Replacement matrix:  $\begin{pmatrix} 2 & 0 \\ 1 & 1 \end{pmatrix}$ ; **initial condition:**
  - ▶ If  $i \in S$  then  $(2|S| - 2 - d_S(i), d_S(i))$ ;
  - ▶ If  $i \notin S$  then  $(2i - 3, 1)$ .

---

Rescaled degrees converge almost surely:

$$\begin{aligned} d_{\text{PA}(n,S)}(i) / \sqrt{n} &\xrightarrow{n \rightarrow \infty} D_i(S) \\ \Delta(\text{PA}(n,S)) / \sqrt{n} &\xrightarrow{n \rightarrow \infty} D_{\max}(S) \\ D_{\max}(S) &= \max_{i \geq 1} D_i(S) \end{aligned}$$

See Móri (2005), Janson (2006), Peköz, Röllin, Ross (2013, 2014).

# Influence of the seed on the maximum degree

## Lemma (Tail behavior of the maximum degree)

Let  $\mathcal{S}$  be a finite tree and let  $m := |\{i \in \{1, \dots, |\mathcal{S}|\} : d_{\mathcal{S}}(i) = \Delta(\mathcal{S})\}|$ .  
Then

$$\mathbb{P}(D_{\max}(\mathcal{S}) > t) \sim m \times c(|\mathcal{S}|, \Delta(\mathcal{S})) t^{1-2|\mathcal{S}|+2\Delta(\mathcal{S})} \exp(-t^2/4)$$

as  $t \rightarrow \infty$ , where the constant  $c$  is explicit.

# Influence of the seed on the maximum degree

## Lemma (Tail behavior of the maximum degree)

Let  $\mathcal{S}$  be a finite tree and let  $m := |\{i \in \{1, \dots, |\mathcal{S}|\} : d_{\mathcal{S}}(i) = \Delta(\mathcal{S})\}|$ .  
Then

$$\mathbb{P}(D_{\max}(\mathcal{S}) > t) \sim m \times c(|\mathcal{S}|, \Delta(\mathcal{S})) t^{1-2|\mathcal{S}|+2\Delta(\mathcal{S})} \exp(-t^2/4)$$

as  $t \rightarrow \infty$ , where the constant  $c$  is explicit.

↪ the seed influences the polynomial factor!

# Influence of the seed on the maximum degree

## Lemma (Tail behavior of the maximum degree)

Let  $S$  be a finite tree and let  $m := |\{i \in \{1, \dots, |S|\} : d_S(i) = \Delta(S)\}|$ .  
Then

$$\mathbb{P}(D_{\max}(S) > t) \sim m \times c(|S|, \Delta(S)) t^{1-2|S|+2\Delta(S)} \exp(-t^2/4)$$

as  $t \rightarrow \infty$ , where the constant  $c$  is explicit.

↪ the seed influences the polynomial factor!

## Corollary (Distinguishing seeds)

If  $|S| - \Delta(S) \neq |T| - \Delta(T)$ , then

$$\delta_{\text{PA}}(S, T) > 0.$$

# Influence of the seed on the maximum degree

## Lemma (Tail behavior of the maximum degree)

Let  $\mathcal{S}$  be a finite tree and let  $m := |\{i \in \{1, \dots, |\mathcal{S}|\} : d_{\mathcal{S}}(i) = \Delta(\mathcal{S})\}|$ .  
Then

$$\mathbb{P}(D_{\max}(\mathcal{S}) > t) \sim m \times c(|\mathcal{S}|, \Delta(\mathcal{S})) t^{1-2|\mathcal{S}|+2\Delta(\mathcal{S})} \exp(-t^2/4)$$

as  $t \rightarrow \infty$ , where the constant  $c$  is explicit.

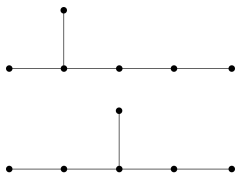
↪ the seed influences the polynomial factor!

## Corollary (Distinguishing seeds)

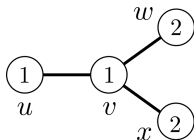
If  $|\mathcal{S}| - \Delta(\mathcal{S}) \neq |\mathcal{T}| - \Delta(\mathcal{T})$ , then

$$\delta_{\text{PA}}(\mathcal{S}, \mathcal{T}) > 0.$$

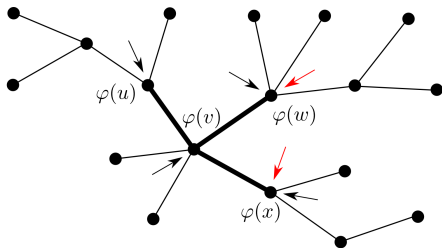
Two trees with the same  
degree profile:



# The approach of Curien et al.



$$\underline{\mathcal{T}} = (\mathcal{T}, \ell)$$



$$D_{\underline{\mathcal{T}}}(T) := \sum_{\varphi} \prod_{u \in \underline{\mathcal{T}}} [d_{\mathcal{T}}(\varphi(u))]_{\ell(u)}$$

Combinatorial interpretation:  $D_{\underline{\mathcal{T}}}(T) = \#$  decorated embeddings

Heuristic:

- ▶ large degree nodes contribute the most;
- ▶ captures geometric structure of large degree nodes.



# The approach of Curien et al.

## General framework:

- ▶ Construct a **family of martingales** using decorated embeddings:

$$M_{\underline{\tau}}^{(S)}(n) = \sum_{\underline{\tau}' \preceq \underline{\tau}} c_n(\underline{\tau}, \underline{\tau}') D_{\underline{\tau}'}(\text{PA}(n, S)).$$

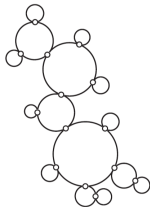
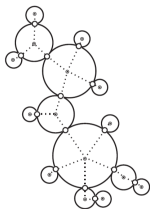
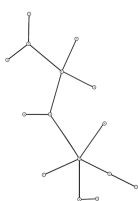
- ▶ For any  $S$  and  $T$ , there exists  $\underline{\tau}$  and  $n$  such that

$$\mathbb{E} \left[ M_{\underline{\tau}}^{(S)}(n) \right] \neq \mathbb{E} \left[ M_{\underline{\tau}}^{(T)}(n) \right].$$

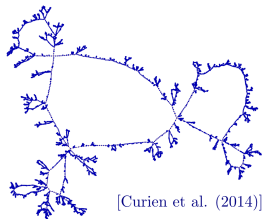
- ▶ Prove that the martingales are **bounded in  $L^2$** .
- ▶ Conclude using the **Paley-Zygmund inequality** that

$$\delta_{\text{PA}}(S, T) > 0.$$

# The Brownian looptree



[Curien et al. (2014)]



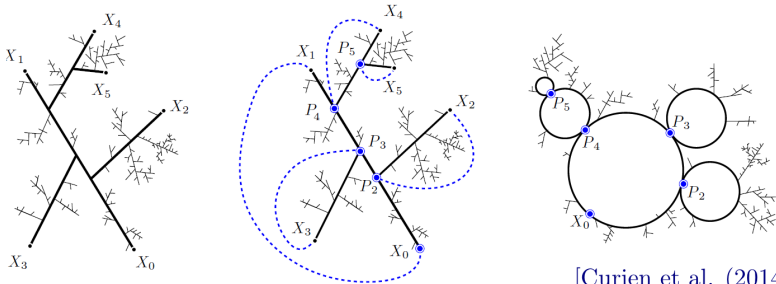
[Curien et al. (2014)]

## Theorem (Curien-Duquesne-Kortchemski-Manolescu, June 2014)

For any  $S$  there exists a random compact metric space  $\mathcal{L}^{(S)}$  such that the following convergence holds a.s. in the Gromov-Hausdorff topology:

$$n^{-1/2} \cdot \text{Loop}(\text{PA}(n, S)) \xrightarrow{n \rightarrow \infty} 2\sqrt{2} \cdot \mathcal{L}^{(S)}.$$

# The Brownian looptree



[Curien et al. (2014)]

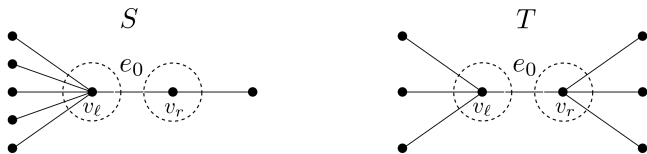
The metric space  $\mathcal{L}$  is constructed as a quotient of  
Aldous's **Brownian Continuum Random Tree**.

**Conjecture (Curien-Duquesne-Kortchemski-Manolescu, June 2014)**

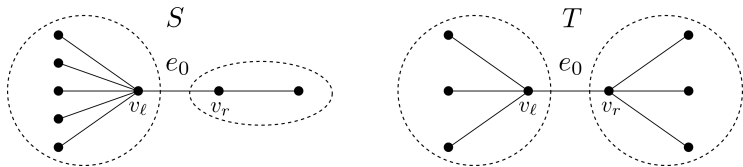
For any pair of seeds  $S$  and  $T$ ,

$$\delta_{\text{PA}}(S, T) = \text{TV} \left( \mathcal{L}^{(S)}, \mathcal{L}^{(T)} \right).$$

# Uniform attachment

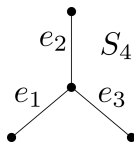
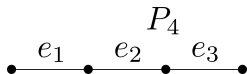


**Preferential attachment:** the **degrees** of  $v_l$  and  $v_r$  are unbalanced in  $S$  but balanced in  $T$ , and this likely remains so throughout the process.



**Uniform attachment:** the **subtree sizes** under  $v_l$  and  $v_r$  are unbalanced in  $S$  but balanced in  $T$ , and this likely remains so throughout the process.

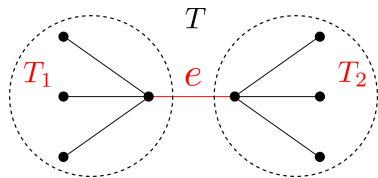
# An example: distinguishing $P_4$ and $S_4$



Measuring balancedness:

$$g(T, e) := \frac{|T_1|^2 |T_2|^2}{|T|^4}$$

$$G(T) := \sum_e g(T, e)$$



In order to show that  $\delta_{\text{UA}}(P_4, S_4) > 0$ , it suffices to show that

$$\liminf_{n \rightarrow \infty} |\mathbb{E}[G(\text{UA}(n, P))] - \mathbb{E}[G(\text{UA}(n, S))]| > 0$$

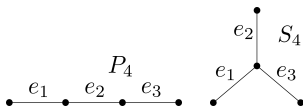
$$\limsup_{n \rightarrow \infty} (\text{Var}[G(\text{UA}(n, P))] + \text{Var}[G(\text{UA}(n, S))]) < \infty$$

## An example: distinguishing $P_4$ and $S_4$

Let  $\{e_j^P\}$  and  $\{e_j^S\}$  denote the edges.

For every  $j \geq 4$ :

$$g(\text{UA}(n, P), e_j^P) \stackrel{d}{=} g(\text{UA}(n, S), e_j^S).$$



We also have this for  $j = 1$  and  $j = 3$ , so

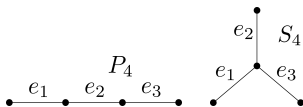
$$\begin{aligned} \mathbb{E}[G(\text{UA}(n, P))] - \mathbb{E}[G(\text{UA}(n, S))] &= \mathbb{E}[g(\text{UA}(n, P), e_2^P)] - \mathbb{E}[g(\text{UA}(n, S), e_2^S)] \\ &= \frac{2n^3 + 5n^2 + 8n + 5}{140n^3} \rightarrow \frac{1}{70}. \end{aligned}$$

## An example: distinguishing $P_4$ and $S_4$

Let  $\{e_j^P\}$  and  $\{e_j^S\}$  denote the edges.

For every  $j \geq 4$ :

$$g\left(\text{UA}(n, P), e_j^P\right) \stackrel{d}{=} g\left(\text{UA}(n, S), e_j^S\right).$$



We also have this for  $j = 1$  and  $j = 3$ , so

$$\begin{aligned}\mathbb{E}[G(\text{UA}(n, P))] - \mathbb{E}[G(\text{UA}(n, S))] &= \mathbb{E}[g(\text{UA}(n, P), e_2^P)] - \mathbb{E}[g(\text{UA}(n, S), e_2^S)] \\ &= \frac{2n^3 + 5n^2 + 8n + 5}{140n^3} \rightarrow \frac{1}{70}.\end{aligned}$$

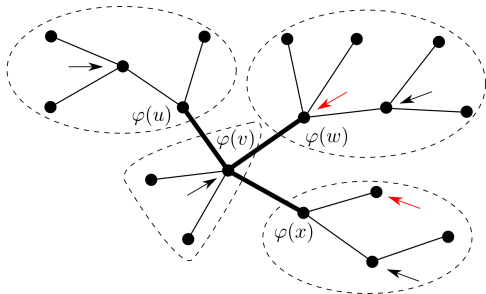
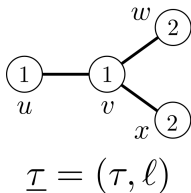
For the variance we use Cauchy-Schwarz:

$$\text{Var}[G(\text{UA}(n, S))] \leq \left( \sum_{j=1}^{n-1} \sqrt{\text{Var}[g(\text{UA}(n, S), e_j)]} \right)^2,$$

and estimates on moments of the beta-binomial distribution to give

$$\mathbb{E}[g(\text{UA}(n, S), e_j)^2] \leq C/j^4.$$

# General statistics



$$F_{\underline{\mathcal{T}}}(T) := \sum_{\varphi} \prod_{u \in \underline{\mathcal{T}}} [f_{\varphi(u)}(T)]_{\ell(u)}$$

Combinatorial interpretation:  $F_{\underline{\mathcal{T}}}(T) = \#$  decorated embeddings

Heuristic:

- ▶ embeddings that are “central” contribute the most;
- ▶ captures global balancedness properties of the tree.



# General framework

- ▶ Construct a **family of martingales** using decorated embeddings:

$$M_{\underline{T}}^{(S)}(n) = \sum_{\underline{T}' \preceq \underline{T}} c_n(\underline{T}, \underline{T}') F_{\underline{T}'}(\text{UA}(n, S)).$$

- ▶ For any  $S$  and  $T$ , there exists  $\underline{T}$  and  $n$  such that

$$\mathbb{E} \left[ M_{\underline{T}}^{(S)}(n) \right] \neq \mathbb{E} \left[ M_{\underline{T}}^{(T)}(n) \right].$$

- ▶ Prove that the martingales are **bounded in  $L^2$** .
- ▶ Conclude using the **Paley-Zygmund inequality** that

$$\delta_{\text{UA}}(S, T) > 0.$$

# Main technical issue: second moment

## Lemma (First moment)

Let  $\underline{\tau} \in \mathcal{D}_+$  be a decorated tree with positive labels and  $|\underline{\tau}| \geq 2$ , and let  $S$  be a seed tree. Then

$$n^{w(\underline{\tau})} \lesssim \mathbb{E} [F_{\underline{\tau}}(\text{UA}(n, S))] \lesssim n^{w(\underline{\tau})},$$

where  $w(\underline{\tau}) = \sum_{u \in \tau} \ell(u)$ .

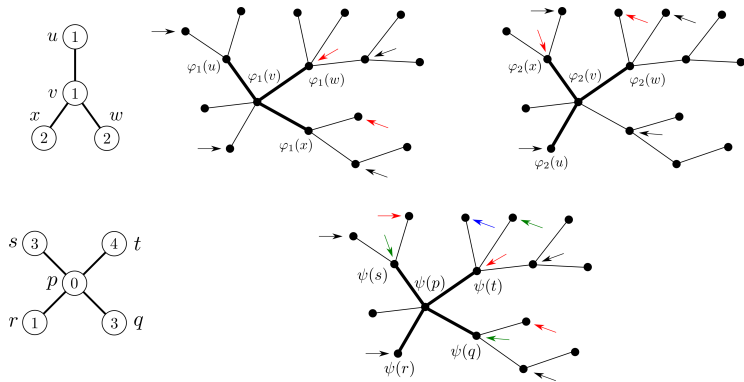
## Lemma (Second moment)

Let  $\underline{\tau} \in \mathcal{D}_+$  be a decorated tree with positive labels and  $|\underline{\tau}| \geq 2$ , and let  $S$  be a seed tree. Then

$$(a) \quad \mathbb{E} [F_{\underline{\tau}}(\text{UA}(n, S))^2] \lesssim n^{2w(\underline{\tau})},$$

$$(b) \quad \mathbb{E} [(F_{\underline{\tau}}(\text{UA}(n+1, S)) - F_{\underline{\tau}}(\text{UA}(n, S)))^2] \lesssim n^{2w(\underline{\tau})-2}.$$

# Main technical issue: second moment



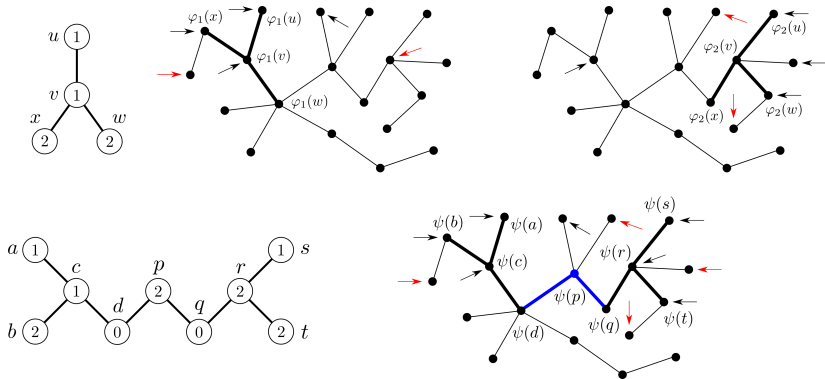
Top row: a decorated tree  $\underline{\tau}$  and two decorated embeddings,

$\underline{\varphi}_1$  and  $\underline{\varphi}_2$ , of it into a larger tree  $T$ .

Bottom row: an associated decorated tree  $\underline{\sigma}$  and the decorated embedding  $\underline{\psi}$  of it into  $T$ .

Note:  $w(\underline{\sigma}) \leq 2w(\underline{\tau})$ .

# Main technical issue: second moment



**Top row:** a decorated tree  $\tau$  and two decorated embeddings,

$\varphi_1$  and  $\varphi_2$ , of it into a larger tree  $T$ .

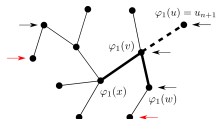
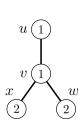
**Bottom row:** an associated decorated tree  $\sigma$  and the

decorated embedding  $\psi$  of it into  $T$ .

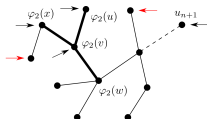
Note:  $w(\sigma) \leq 2w(\tau)$ , but no a priori bound on  $|\sigma|$ .

$\rightsquigarrow$  use the fact that  $\text{diam}(\text{UA}(n, S)) = O(\log n)$  whp.

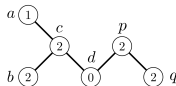
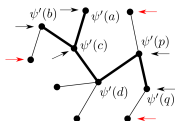
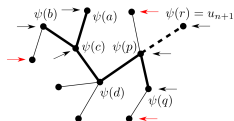
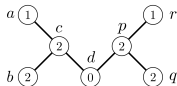
# Main technical issue: second moment



Type A



Type B



**Top row:** There are two types of decorated embeddings that use the new vertex.

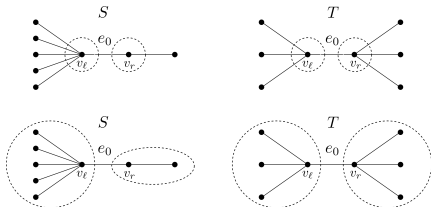
**Bottom row:** associated decorated trees and decorated embeddings.

Roughly speaking, the two arrows associated with the new vertex give the extra factor of  $n^{-2}$  required in the bound of (b).

# Summary and open questions

## Takeaways:

- ▶ Every seed has an influence, both in PA and in UA
- ▶ Degrees (PA) and balancedness (UA) are key statistics



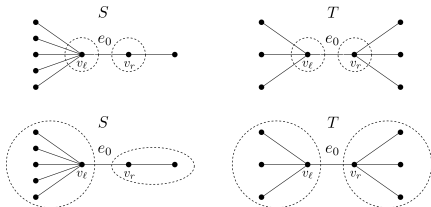
## Open questions:

- ▶ Multiple edges added at each time step?
- ▶ Is  $\delta_\alpha(S, T) > 0$  for  $\alpha \in (0, 1)$ ?  
Is it monotone in  $\alpha$ ? Is it convex?
- ▶ Other models of randomly growing graphs.
- ▶ Estimation. Finding the seed.
- ▶ The effect of extra information.
- ▶ Applications...

# Summary and open questions

## Takeaways:

- ▶ Every seed has an influence, both in PA and in UA
- ▶ Degrees (PA) and balancedness (UA) are key statistics



## Open questions:

- ▶ Multiple edges added at each time step?
- ▶ Is  $\delta_\alpha(S, T) > 0$  for  $\alpha \in (0, 1)$ ?  
Is it monotone in  $\alpha$ ? Is it convex?
- ▶ Other models of randomly growing graphs.
- ▶ Estimation. Finding the seed.
- ▶ The effect of extra information.
- ▶ Applications...

Thank you!