

# Studying large networks via local weak limit theory

Venkat Anantharam

EECS Department  
University of California, Berkeley

March 10, 2017

Advanced Networks Colloquium

University of Maryland, College Park

(Joint work with Justin Salez and Payam Delgosha )

# Outline

- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs
- 5 Graph indexed data
- 6 Universal compression of graphical data



Justin Salez



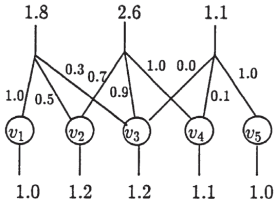
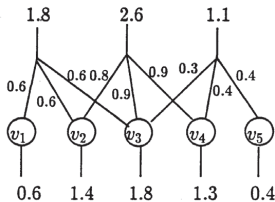
Payam Delgosha

# Outline

- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs
- 5 Graph indexed data
- 6 Universal compression of graphical data

# Resource allocation

Consumers above, Resources below



# Balanced resource allocation

- Let  $f$  be a convex function on the nonnegative reals.
- Over assignments  $\theta$ , the objective is to minimize

$$J(\theta) := \sum_{i=1}^M f(\theta(i)) .$$

where  $\theta(i)$  is the load at resource  $i$  and  $M$  is the number of resources.

# Balanced resource allocation

- Let  $f$  be a convex function on the nonnegative reals.
- Over assignments  $\theta$ , the objective is to minimize

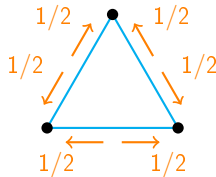
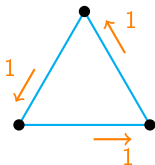
$$J(\theta) := \sum_{i=1}^M f(\partial\theta(i)) .$$

where  $\partial\theta(i)$  is the load at resource  $i$  and  $M$  is the number of resources.

- **Theorem** ( *Hajek*): The assignment  $\theta$  minimizes  $J(\theta)$  iff for all pairs of resources  $i, i'$  available to consumer  $u$  we have  $\theta_u(i) = 0$  whenever  $\partial\theta(i) > \partial\theta(i')$ .
- Note that the condition for an assignment to be balanced does not depend on  $f$ .

# Uniqueness of the balanced loads

The assignment  $\theta$  need not be unique, but  $\partial\theta(i)$  is unique





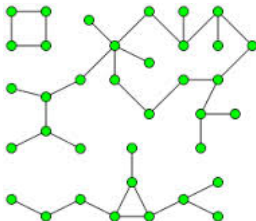
# Many consumers and resources

- We want to understand the local environment of a typical agent (consumer, resource) in resource allocation problem with many agents.
- We will first describe how this can be done for the basic load balancing problem in the case of **large sparse graphs**.

# Outline

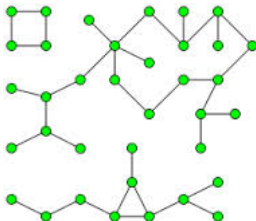
- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs**
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs
- 5 Graph indexed data
- 6 Universal compression of graphical data

# Graphs



- A graph corresponds to a load balancing problem where each consumer has access to two resources.

# Graphs



- A graph corresponds to a load balancing problem where each consumer has access to two resources.
- Each edge is a consumer with one unit of load and has to decide how to distribute its load between the two vertices that define the edge.
- Multiple edges between a pair of vertices are okay.

# Load percolation

- Note that the local structure of the balanced allocation depends on the global structure of the graph, not just on its local structure.

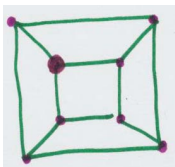


Figure: Graph A



Figure: Graph B

# Load percolation

- Note that the local structure of the balanced allocation depends on the global structure of the graph, not just on its local structure.

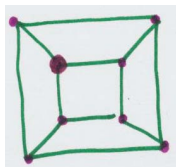


Figure: Graph A

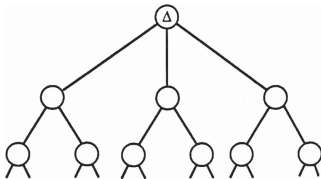


Figure: Graph B

- The marked vertex in graph A has the same depth-1 neighborhood as the root in graph B.
- However the induced balanced load is  $\frac{3}{2}$  at each vertex in graph A and is  $\frac{4}{5}$  in graph B.
- The phenomenon underlying this is called **load percolation** by Hajek.

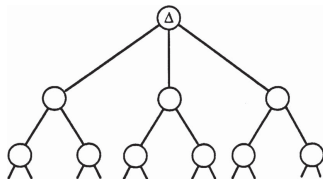
# Load percolation as nonuniqueness in the limit

- An infinite sparse graph can exhibit nonuniqueness in its balanced allocations.



# Load percolation as nonuniqueness in the limit

- An infinite sparse graph can exhibit nonuniqueness in its balanced allocations.

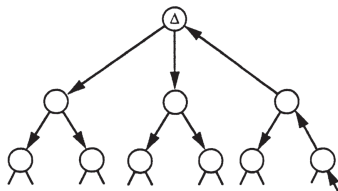


- In this infinite 3-regular tree, start by assigning the load of each edge to the vertex that is furthest from the marked vertex.
- This gives induced load 1 at all vertices except for the marked one, which has induced load 0.



# Nonuniqueness: an example due to Hajek

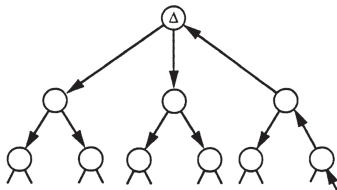
- Pick a path from infinity to the marked node and flip the allocations of edges along this path.



- This allocation is balanced. Each vertex has induced load 1.

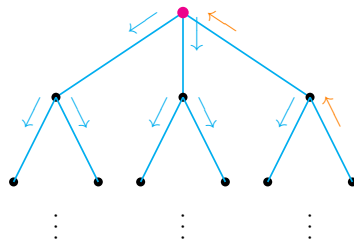
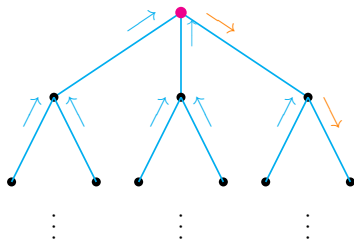
# Nonuniqueness: an example due to Hajek

- Pick a path from infinity to the marked node and flip the allocations of edges along this path.



- This allocation is balanced. Each vertex has induced load 1.
- Now flip the allocation of **each edge**.
- This is another balanced allocation **!!**. The induced load at each vertex is 2.
- These examples are due to Hajek.

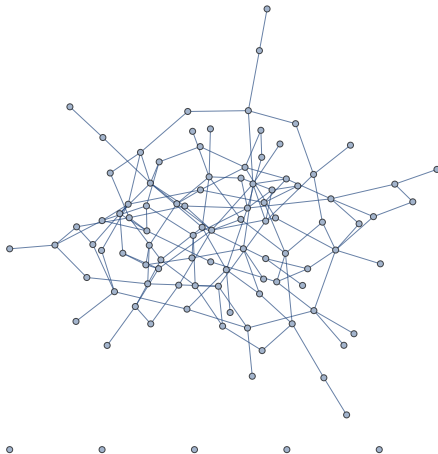
# Another look at the Hajek counterexample



# Hajek's conjectures

- To develop insight into the structure of the balanced load allocation in large graphs Hajek carried out simulations.
- He picked random graphs according to a sparse Erdős-Rényi model and studied the corresponding balanced allocations.

# A sparse Erdős-Rényi graph



# Numerics on Erdős-Rényi graphs (*Hajek*)

$\alpha M$  consumers and  $M$  resources; edges picked at random

SAMPLE LOAD DISTRIBUTION *BEFORE* BALANCING ( $\alpha = 2$ ,  $M = 10000$ )

$\tau$	Load $\leq \tau$	Load $\approx \tau$
0.0	201	201
0.5	921	720
1.0	2382	1461
1.5	4299	1917
2.0	6291	1992
2.5	7896	1605
3.0	8899	1003
3.5	9472	573
4.0	9778	306
4.5	9912	134
5.0	9962	50
5.5	9987	25
6.0	10000	13

SAMPLE LOAD DISTRIBUTION *AFTER* BALANCING ( $\alpha = 2$ ,  $M = 10000$ )

$\tau$	Load $\leq \tau$	Load $\approx \tau$	Product
0.00000000	201	201	0
0.50000000	223	22	11
1.00000000	992	769	769
1.25000000	996	4	5
1.33333333	1023	27	36
1.50000000	1239	216	324
1.60000000	1244	5	8
1.66666667	1313	69	115
1.75000000	1353	40	70
1.77777778	1362	9	16
1.80000000	1392	30	54
1.83333333	1398	6	11
1.85714286	1405	7	13
1.92307692	1418	13	25
2.00000000	3316	1898	3796
2.07692308	3329	13	27
2.11111111	3338	9	19
2.12500000	3362	24	51
2.14285714	3404	42	90
2.16666667	3440	36	78
2.18181818	3462	22	48
2.20000000	3562	100	220
2.20782852	10000	6438	14214

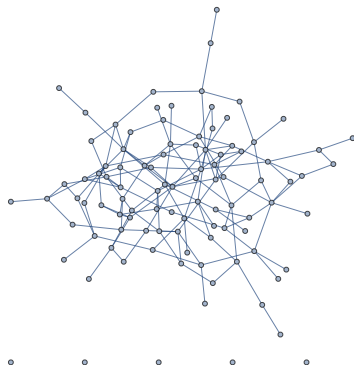
# Numerics on Erdős-Rényi graphs (*Hajek*) (cont'd)

SAMPLE LOAD DISTRIBUTION *AFTER* BALANCING ( $\alpha = 10$ ,  $M = 10000$ )

$\tau$	Load $\leq \tau$	Load $= \tau$	Product
6.00000000	2	2	12
7.00000000	6	4	28
8.00000000	17	11	88
9.00000000	51	34	306
9.33333333	54	3	28
9.50000000	56	2	19
10.00000000	114	58	580
10.00799110	10000	9886	98939

# Large Erdős Rényi graphs

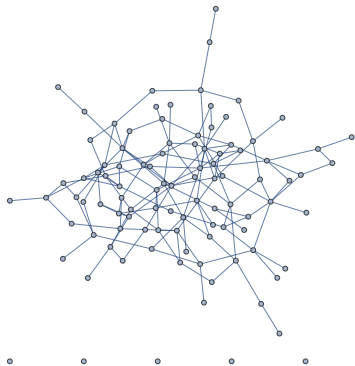
$$\mathcal{G}(n, \alpha/n)$$



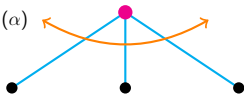


# Large Erdős Rényi graphs

$$\mathcal{G}(n, \alpha/n)$$

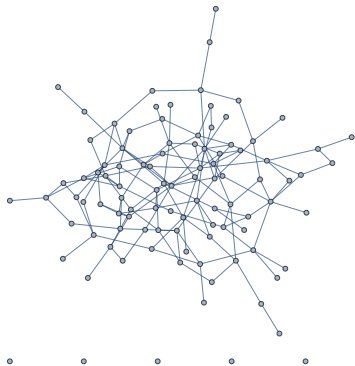


$$(n-1)\text{Ber}(\alpha/n) \approx \text{Poi}(\alpha)$$

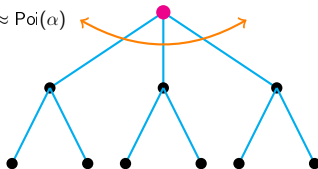


# Large Erdős Rényi graphs

$$\mathcal{G}(n, \alpha/n)$$

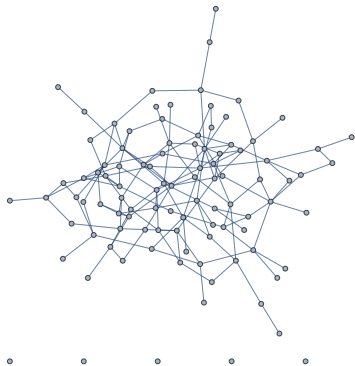


$$(n-1)\text{Ber}(\alpha/n) \approx \text{Poi}(\alpha)$$

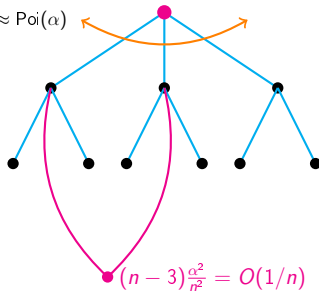


# Large Erdős Rényi graphs

$$\mathcal{G}(n, \alpha/n)$$



$$(n-1)\text{Ber}(\alpha/n) \approx \text{Poi}(\alpha)$$



# The Poisson Galton-Watson tree

## Poisson Galton-Watson tree :

- Start with a root.
- Pick a Poisson ( $\lambda$ ) number of neighbors (at depth 1).
- For each of these, independently pick a Poisson ( $\lambda$ ) number of neighbors (at depth 2).
- $\vdots$
- Etc.

# The Poisson Galton-Watson tree

## Poisson Galton-Watson tree :

- Start with a root.
- Pick a Poisson ( $\lambda$ ) number of neighbors (at depth 1).
- For each of these, independently pick a Poisson ( $\lambda$ ) number of neighbors (at depth 2).

⋮

Etc.

- The local environment of a typical vertex in an Erdős - Rényi graph converges to a Poisson Galton-Watson tree as  $M \rightarrow \infty$ .

# A recursive distributional equation

- The numerics suggest that there should be a well defined limiting distribution ( $M \rightarrow \infty$ ) for the induced load (in a balanced allocation) at a typical vertex.
- Natural guess: the limiting induced load distribution obeys a fixed point equation (a **recursive distributional equation**).
- This was conjectured by Hajek.

# Our contribution

- We verify this conjecture of Hajek as a special case of a broader result.
- Our results are in the language of **local weak convergence** of sequences of graphs, also called **the objective method**.
- In this theory graphs are viewed through the lens of probability distributions on rooted graphs.

# What we prove (with Justin Salez)

- There is a uniquely defined balanced allocation associated to any probability distribution on infinite rooted graphs that can arise as a local weak limit of a sequence of finite graphs.
- The unique balanced allocation on the finite graphs converges to the corresponding unique balanced allocation on its local weak limit.
- The induced load distribution at the root in the infinite limit rooted graph obeys the expected recursive distributional equation.



# Outline

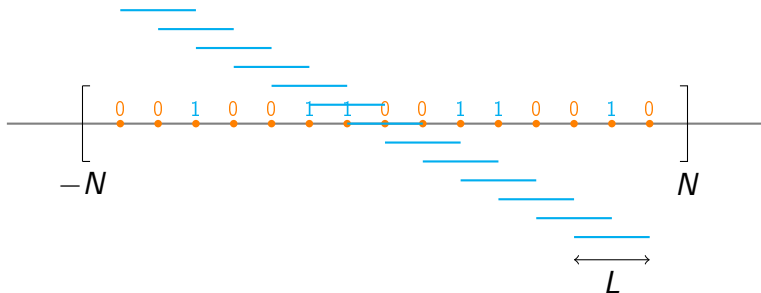
- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence**
- 4 Load balancing on hypergraphs
- 5 Graph indexed data
- 6 Universal compression of graphical data

# Stochastic processes as a model for data samples

A stochastic process is a model for the structure of data samples.

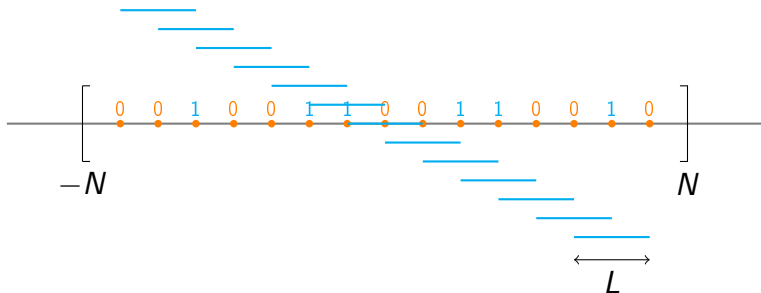
# Stochastic processes as a model for data samples

A stochastic process is a model for the structure of data samples.



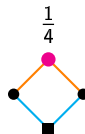
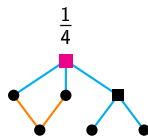
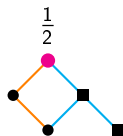
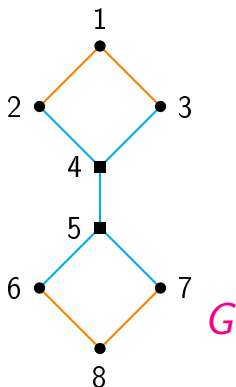
# Stochastic processes as a model for data samples

A stochastic process is a model for the structure of data samples.



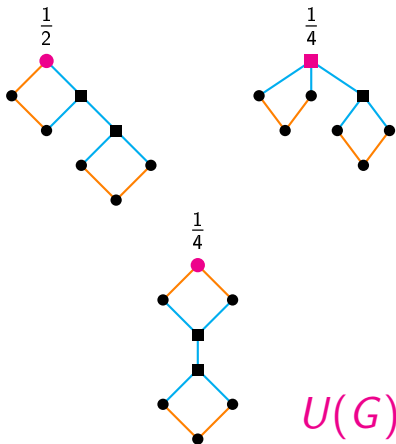
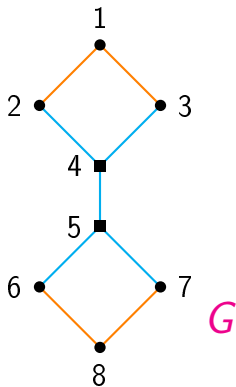
$$\frac{1}{2(N+1) - L} \sum_{i=-N}^{N-L+1} \delta_{x_i, \dots, x_{i+L-1}} \Rightarrow P_{x_0, \dots, x_{L-1}}.$$

# "Empirical distribution" of a marked graph

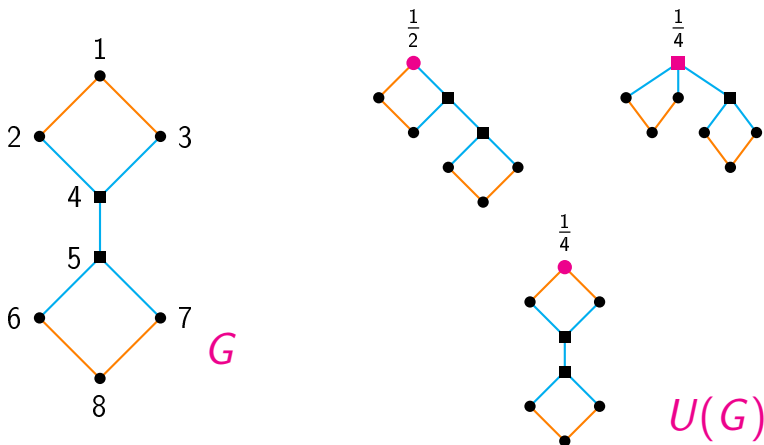


$U_2(G)$

# Rooted marked graph process from a marked graph

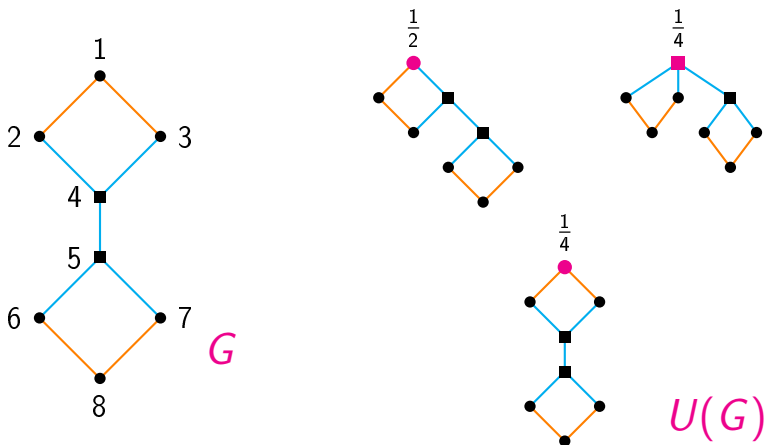


# Rooted marked graph process from a marked graph



- $\mathcal{G}_*$ : space of unlabelled marked rooted graphs

# Rooted marked graph process from a marked graph



- $\mathcal{G}_*$ : space of unlabelled marked rooted graphs
- A process with values in rooted marked graphs:  $\mu \in \mathcal{P}(\mathcal{G}_*)$
- We will first consider the unmarked case.



# The space of rooted graphs

- $\mathcal{G}_*$  denotes the set of locally finite connected rooted graphs considered up to rooted isomorphism.
- The distance between two elements of  $\mathcal{G}_*$  is  $\frac{1}{1+r}$ , where  $r$  is the largest depth of a neighborhood around the root up to which they agree.
- This distance makes  $\mathcal{G}_*$  into a complete separable metric space.

# Local weak limit of a sequence of graphs

- A fixed finite graph  $G$  corresponds to a probability distribution on  $\mathcal{G}_*$  by picking the root at random from the vertices of  $G$ .
- A sequence of finite graphs is said to converge in the sense of local weak convergence if the corresponding probability distributions on  $\mathcal{G}_*$  converge weakly.

The definitions extend naturally to **marked graphs**, i.e. graphs where each edge and each vertex carries an element of some other separable metric space.

# The space of (edge, vertex) rooted graphs

- $\mathcal{G}_{**}$  denotes the set of locally finite connected graphs with a distinguished oriented edge, considered up to isomorphism (preserving the distinguished oriented edge).
- $\mathcal{G}_{**}$  can be metrized to give a complete separable metric space, just as for  $\mathcal{G}_*$ .

# Moving between $\mathcal{G}_*$ and $\mathcal{G}_{**}$

- A function  $f : \mathcal{G}_{**} \mapsto \mathbb{R}$  gives rise to a function  $\partial f : \mathcal{G}_* \mapsto \mathbb{R}$  via

$$\partial f(G, o) = \sum_{i \sim o} f(G, i, o) .$$

- A probability distribution  $\mu$  on  $\mathcal{G}_*$  gives rise to a measure  $\vec{\mu}$  on  $\mathcal{G}_{**}$  via

$$\int_{\mathcal{G}_{**}} f d\vec{\mu} = \int_{\mathcal{G}_*} \partial f d\mu , \quad \text{for all bounded continuous } f .$$

- Note that  $\vec{\mu}(\mathcal{G}_{**}) = \deg(\mu) := \int_{\mathcal{G}_*} \deg(\text{root}) d\mu$  .

# Unimodularity

- Given  $f : \mathcal{G}_{**} \mapsto \mathbb{R}$ , define  $f^* : \mathcal{G}_{**} \mapsto \mathbb{R}$  via

$$f^*(G, i, o) = f(G, o, i) .$$

- A probability distribution  $\mu$  on  $\mathcal{G}_*$  is called **unimodular** if

$$\int_{\mathcal{G}_{**}} f d\vec{\mu} = \int_{\mathcal{G}_{**}} f^* d\vec{\mu} , \text{ for all bounded continuous } f .$$

- It is known that the local weak limit of any sequence of finite graphs is unimodular (*Aldous and Lyons*).

# Asymptotic notion of a balanced allocation

- A function  $\Theta : \mathcal{G}_{**} \mapsto [0, 1]$  is called an **allocation** if  $\Theta + \Theta^* = 1$ .
- An allocation  $\Theta$  is called a **balanced allocation** for a given unimodular  $\mu$  if for  $\vec{\mu}$  almost all  $(G, i, o)$  it holds that

$$\partial\Theta(G, i) < \partial\Theta(G, o) \implies \Theta(G, i, o) = 0 .$$

# Formal statement of the main results

- We prove that for any unimodular  $\mu$  with  $\deg(\mu) < \infty$  there is a  $\Theta_0$  that is a balanced allocation for  $\mu$  with the property that it simultaneously minimizes  $\int_{\mathcal{G}_*} f(\partial\Theta) d\mu$  over allocations  $\Theta$  for **every** convex real valued function  $f$  on  $\mathbb{R}_+$ .
- Further,  $\Theta_0$  is  $\mu$ -almost surely unique.
- For any sequence of finite graphs with local weak limit  $\mu$ , the empirical distribution of the induced load in the unique balanced allocation on these graphs converges weakly to the law of  $\partial\Theta_0$  (for the  $\Theta_0$  of the limit).

# Variational characterization of the limit

- Given unimodular  $\mu$  on  $\mathcal{G}_*$  with  $\deg(\mu) < \infty$ , define, for each  $t \geq 0$ ,

$$\Phi_\mu(t) := \int_{\mathcal{G}_*} (\partial\Theta_0 - t)^+ d\mu .$$

- $t \mapsto \Phi_\mu(t)$  is the **mean-excess function** of the almost surely unique balanced allocation associated to  $\mu$ .
- We have the variational characterization

$$\Phi_\mu(t) = \max_{f : \mathcal{G}_* \rightarrow [0,1], \text{ Borel}} \left\{ \frac{1}{2} \int_{\mathcal{G}_{**}} \hat{f} d\vec{\mu} - t \int_{\mathcal{G}_*} f d\mu \right\} ,$$

for each  $t$ , where

$$\hat{f}(G, i, o) := f(G, i) \wedge f(G, o) .$$



# Intuition behind the variational characterization

- The optimizing function is  $f = 1(\partial\Theta_0 > t)$ .
- To check this, observe that

$$\begin{aligned}\frac{1}{2} \int_{\mathcal{G}_{**}} \hat{f} d\vec{\mu} &= \frac{1}{2} \int_{\mathcal{G}_*} (\partial\hat{f}) d\mu \\ &= \frac{1}{2} \int_{\mathcal{G}_*} \sum_{i \sim o} 1(\partial\Theta_0(G, i) > t \text{ and } \partial\Theta_0(G, o) > t) d\mu\end{aligned}$$

- Thus

$$\int_{\mathcal{G}_*} (\partial\Theta_0 - t)^+ d\mu = \frac{1}{2} \int_{\mathcal{G}_{**}} \hat{f} d\vec{\mu} - t \int_{\mathcal{G}_*} f d\mu ,$$

for this choice of  $f$ .

# Unimodular Galton-Watson trees

- Given a probability distribution  $\{\pi(i) , \ i \geq 0\}$  on the nonnegative integers, with finite mean  $\sum_i i\pi(i)$ , define

$$\hat{\pi}(i) := \frac{(i+1)\pi(i+1)}{\sum_i i\pi(i)} , \ i \geq 0 .$$

$\{\hat{\pi}(i) , \ i \geq 0\}$  is also a probability distribution.

- The **unimodular Galton-Watson tree**, **UGWT**( $\pi$ ) is the random tree constructed as follows: Start with a root and give it a random number of children (at depth 1) with the number of children distributed as  $\pi$ . For each child, give it a random number of children (at depth 2), the number distributed as  $\hat{\pi}$ , independently. Repeat (using  $\hat{\pi}$  from now on).
- Many standard sequences of bipartite graph models, such as the pairing model based on half edges and fixed degree distributions which shows up in the theory of LDPC codes, have a unimodular Galton-Watson tree as their local weak limit.

# Recursive distributional equation characterization of the limit on unimodular Galton-Watson trees

- If  $\mu$  is the law of  $\text{UGWT}(\pi)$ , then for every  $t$ , we have

$$\Phi_\mu(t) = \max_{Q=F_{\pi,t}(Q)} \left\{ \frac{E[D]}{2} P(\xi_1 + \xi_2 > 1) - tP(\xi_1 + \dots + \xi_D > t) \right\},$$

where  $F_{\pi,t}(Q)$  is the law of  $[1 - t + \xi_1 + \dots + \xi_{\hat{D}}]_0^1$ .

- Here  $[a]_0^1$  equals 0 if  $a < 0$ , 1 if  $a > 1$  and  $a$  otherwise. Also,  $\hat{D}$  has the law  $\hat{\pi}$ ,  $D$  has the law  $\pi$ , and the  $\xi_i$  are i.i.d. with law  $Q$ .
- Recall that

$$t \mapsto \Phi_\mu(t) := \int_{\mathcal{G}_*} (\partial\Theta_0 - t)^+ d\mu,$$

characterizes the limiting distribution of the induced load at the root.

- The above recursive distributional characterization of is in effect the one conjectured by Hajek.

# Intuition behind the RDE

- We consider the RDE  $Q = F_{\pi,t}(Q)$ , where  $F_{\pi,t}(Q)$  is the law of  $[1 - t + \xi_1 + \dots + \xi_{\hat{D}}]_0^1$ , where  $\xi_1, \xi_2, \dots$  are i.i.d with the law  $Q$ .
- Consider an edge  $(i, o)$ . We are “solving for the load that passes in the direction from  $o$  to  $i$ .”
- For  $1 \leq k \leq \hat{D}$ ,  $1 - \xi_k$  has the meaning of the amount of load that can be absorbed by the  $k$ -th child of  $o$  (think of  $i$  as the parent of  $o$  and not as a child), this child of course supporting its own subtree of children, such as to make the net load at that child equal to  $t$ .
- The number  $[1 - (t - \xi_1 - \dots - \xi_{\hat{D}})]_0^1$  is then the amount that would be presented in the direction from node  $o$  to node  $i$  in order to maintain a total load of  $t$  at node  $o$ .

# Convergence of the maximum load

- Under a mild additional on the degree distributions the **maximum load** also converges to the maximum of the limit.
- This verifies the conjecture of Hajek regarding the limit of the maximum load.
- One must exclude "local pockets of high edge density" in the graph.
- Assume that for some  $\lambda > 0$  we have

$$\sup_{n \geq 1} \left\{ \frac{1}{n} \sum_{i=1}^n e^{\lambda d_n(i)} \right\} < \infty .$$

- Let  $Z_{\delta,t}^{(n)}$  denote the number of subsets  $S$  of  $\{1, \dots, n\}$  of size  $|S| \leq \delta n$  with edge count  $|E(S)| \geq t|S|$  in the given random pairing model. Then we can show that

$$P(Z_{\delta,t}^{(n)} > 0) \rightarrow 0 , \quad \text{as } n \rightarrow \infty .$$

This suffices.

# Sketch of the proof of the main result

- The key idea is to consider so-called  $\epsilon$ -balanced allocations, i.e. allocations  $\theta$  on a locally finite graph  $G$  that satisfy

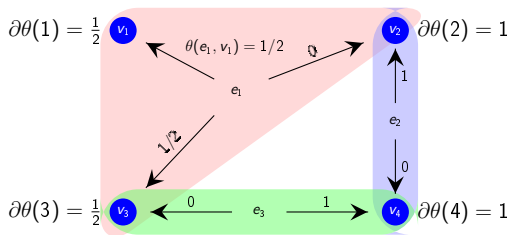
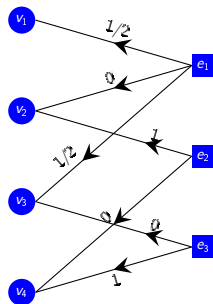
$$\theta(i,j) = \left[ \frac{1}{2} + \frac{1}{2\epsilon}(\partial\theta(i) - \partial\theta(j)) \right]_0^1.$$

- There is a built-in contractivity in this definition for bounded degree graphs, which allows one to establish the uniqueness of  $\epsilon$ -balanced allocations for such graphs.
- The case of locally finite graphs can be handled by a truncation argument.
- The claimed  $\Theta_0$  can then be shown to exist as a limit in  $L^2$  of the  $\epsilon$ -balanced allocations as  $\epsilon \rightarrow 0$ .
- The  $\epsilon$ -relaxation can be roughly thought of as analogous to working at finite temperature (versus zero temperature) in statistical mechanics.

# Outline

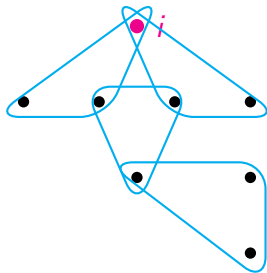
- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs**
- 5 Graph indexed data
- 6 Universal compression of graphical data

# Load Balancing on a hypergraph



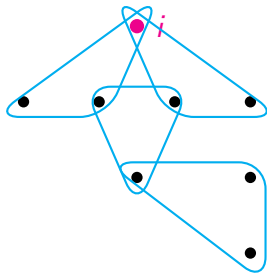


$\mathcal{H}_*$  and  $\mathcal{H}_{**}$

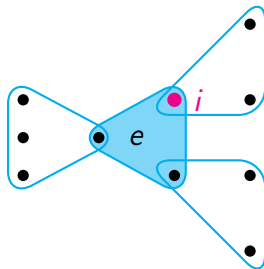


$$\mathcal{H}_* = \{[H, i]\}$$

$\mathcal{H}_*$  and  $\mathcal{H}_{**}$

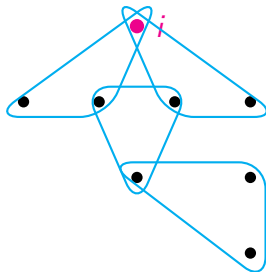


$$\mathcal{H}_* = \{[H, i]\}$$

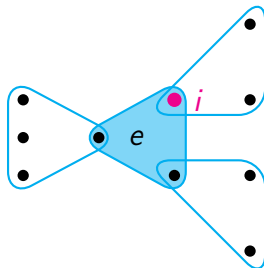


$$\mathcal{H}_{**} = \{[H, e, i]\}$$

$\mathcal{H}_*$  and  $\mathcal{H}_{**}$



$$\mathcal{H}_* = \{[H, i]\}$$



$$\mathcal{H}_{**} = \{[H, e, i]\}$$

Simple, connected, finite edges, locally finite

# Unimodularity

- Finite  $H_n$

$$U(H) = \frac{1}{|V(H)|} \sum_{i \in V(H)} \delta_{[H,i]} \in \mathcal{P}(\mathcal{H}_*)$$

# Unimodularity

- Finite  $H_n$

$$U(H) = \frac{1}{|V(H)|} \sum_{i \in V(H)} \delta_{[H, i]} \in \mathcal{P}(\mathcal{H}_*)$$

- $H_n \xrightarrow{lwc} \mu$  when  $U(H_n) \Rightarrow \mu$

# Unimodularity

- Finite  $H_n$

$$U(H) = \frac{1}{|V(H)|} \sum_{i \in V(H)} \delta_{[H,i]} \in \mathcal{P}(\mathcal{H}_*)$$

- $H_n \xrightarrow{lwc} \mu$  when  $U(H_n) \Rightarrow \mu$
- Not all  $\mu$  can be local weak limits of finite hypergraphs

# Unimodularity

- Finite  $H_n$

$$U(H) = \frac{1}{|V(H)|} \sum_{i \in V(H)} \delta_{[H, i]} \in \mathcal{P}(\mathcal{H}_*)$$

- $H_n \xrightarrow{lwc} \mu$  when  $U(H_n) \Rightarrow \mu$
- Not all  $\mu$  can be local weak limits of finite hypergraphs
- For  $f : \mathcal{H}_{**} \rightarrow \mathbb{R}$ , let

$$\partial f : \mathcal{H}_* \rightarrow \mathbb{R} \quad \partial f(H, i) = \sum_{e \ni i} f(H, e, i)$$

## Unimodularity (cont'd)

- For  $\mu \in \mathcal{P}(\mathcal{H}_*)$ , define  $\vec{\mu} \in \mathcal{M}(\mathcal{H}_{**})$  as

$$\int f d\vec{\mu} = \int \partial f d\mu$$

for all Borel function  $f$  on  $\mathcal{H}_{**}$ .



# Unimodularity (cont'd)

- For  $\mu \in \mathcal{P}(\mathcal{H}_*)$ , define  $\vec{\mu} \in \mathcal{M}(\mathcal{H}_{**})$  as

$$\int f d\vec{\mu} = \int \partial f d\mu$$

for all Borel function  $f$  on  $\mathcal{H}_{**}$ .

- For  $f : \mathcal{H}_{**} \rightarrow \mathbb{R}$ , let

$$\nabla f : \mathcal{H}_{**} \rightarrow \mathbb{R} \quad \nabla f(H, e, i) = \frac{1}{|e|} \sum_{j \in e} f(H, e, j).$$

# Unimodularity (cont'd)

- For  $\mu \in \mathcal{P}(\mathcal{H}_*)$ , define  $\vec{\mu} \in \mathcal{M}(\mathcal{H}_{**})$  as

$$\int f d\vec{\mu} = \int \partial f d\mu$$

for all Borel function  $f$  on  $\mathcal{H}_{**}$ .

- For  $f : \mathcal{H}_{**} \rightarrow \mathbb{R}$ , let

$$\nabla f : \mathcal{H}_{**} \rightarrow \mathbb{R} \quad \nabla f(H, e, i) = \frac{1}{|e|} \sum_{j \in e} f(H, e, j).$$

- $\mu \in \mathcal{P}(\mathcal{H}_*)$  is called unimodular if

$$\int f d\vec{\mu} = \int \nabla f d\vec{\mu}$$

# Unimodularity (cont'd)

- For  $\mu \in \mathcal{P}(\mathcal{H}_*)$ , define  $\vec{\mu} \in \mathcal{M}(\mathcal{H}_{**})$  as

$$\int fd\vec{\mu} = \int \partial fd\mu$$

for all Borel function  $f$  on  $\mathcal{H}_{**}$ .

- For  $f : \mathcal{H}_{**} \rightarrow \mathbb{R}$ , let

$$\nabla f : \mathcal{H}_{**} \rightarrow \mathbb{R} \quad \nabla f(H, e, i) = \frac{1}{|e|} \sum_{j \in e} f(H, e, j).$$

- $\mu \in \mathcal{P}(\mathcal{H}_*)$  is called unimodular if

$$\int fd\vec{\mu} = \int \nabla fd\vec{\mu}$$

- If  $H_n \xrightarrow{lwc} \mu$ ,  $\mu$  is unimodular

# Borel Allocations and Balancedness

- $\Theta : \mathcal{H}_{**} \rightarrow [0, 1]$  is called a Borel allocation if

$$\sum_{j \in e} \Theta(H, e, j) = 1 \quad \forall [H, e, i] \in \mathcal{H}_{**}$$

# Borel Allocations and Balancedness

- $\Theta : \mathcal{H}_{**} \rightarrow [0, 1]$  is called a Borel allocation if

$$\sum_{j \in e} \Theta(H, e, j) = 1 \quad \forall [H, e, i] \in \mathcal{H}_{**}$$

- $\Theta$  is balanced w.r.t.  $\mu \in \mathcal{P}(\mathcal{H}_*)$  if for  $\vec{\mu}$ -almost all  $[H, e, i] \in \mathcal{H}_{**}$

$$j \in e \quad \partial\Theta(H, i) > \partial\Theta(H, j) \quad \Rightarrow \quad \Theta(H, e, i) = 0.$$

# Main results (with Payam Delgosha)

## Theorem

*Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then*

# Main results (with Payam Delgosha)

## Theorem

Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then

- 1 (existence)  $\exists$  a balanced allocation  $\Theta_0$

# Main results (with Payam Delgosha)

## Theorem

Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then

- 1 (existence)  $\exists$  a balanced allocation  $\Theta_0$
- 2 (uniqueness)  $\Theta_1, \Theta_2$  two balanced allocations, then  $\partial\Theta_1 = \partial\Theta_2$ ,  $\mu$ -a.s.



# Main results (with Payam Delgosha)

## Theorem

Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then

- 1 (existence)  $\exists$  a balanced allocation  $\Theta_0$
- 2 (uniqueness)  $\Theta_1, \Theta_2$  two balanced allocations, then  $\partial\Theta_1 = \partial\Theta_2$ ,  $\mu$ -a.s.
- 3 (continuity)  $H_n \xrightarrow{lwc} \mu$  then  $\mathcal{L}_n \Rightarrow \mathcal{L}$

# Main results (with Payam Delgosha)

## Theorem

Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then

- 1 (existence)  $\exists$  a balanced allocation  $\Theta_0$
- 2 (uniqueness)  $\Theta_1, \Theta_2$  two balanced allocations, then  $\partial\Theta_1 = \partial\Theta_2$ ,  $\mu$ -a.s.
- 3 (continuity)  $H_n \xrightarrow{\text{lwc}} \mu$  then  $\mathcal{L}_n \Rightarrow \mathcal{L}$
- 4 (optimality)  $\Theta$  is balanced iff it minimizes  $\int f(\partial\Theta)d\mu$  for strictly convex  $f : [0, \infty) \rightarrow \mathbb{R}$ .

# Main results (with Payam Delgosha)

## Theorem

Take  $\mu \in \mathcal{P}(\mathcal{H}_*)$  unimodular,  $\deg(\mu), \text{Var}(\mu) < \infty$ , then

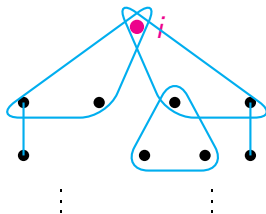
- 1 (existence)  $\exists$  a balanced allocation  $\Theta_0$
- 2 (uniqueness)  $\Theta_1, \Theta_2$  two balanced allocations, then  $\partial\Theta_1 = \partial\Theta_2$ ,  $\mu$ -a.s.
- 3 (continuity)  $H_n \xrightarrow{lwc} \mu$  then  $\mathcal{L}_n \Rightarrow \mathcal{L}$
- 4 (optimality)  $\Theta$  is balanced iff it minimizes  $\int f(\partial\Theta)d\mu$  for strictly convex  $f : [0, \infty) \rightarrow \mathbb{R}$ .
- 5 (variational characterization)  $t \in \mathbb{R}$  and  $\Theta$  balanced, then

$$\int (\partial\Theta - t)^+ d\mu = \max_{f \in \mathcal{H}_* \xrightarrow{\text{Borel}} [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

where  $\tilde{f}_{\min}(H, e, i) = \frac{1}{|e|} \min_{j \in e} f(H, j)$ .

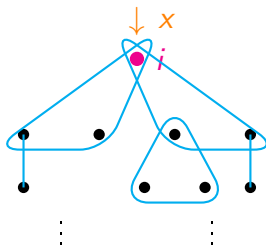
# Response Function

$\rho_{T,i}(x)$  = total load at  $i$  with baseload  $x$



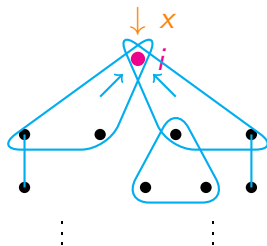
# Response Function

$\rho_{T,i}(x)$  = total load at  $i$  with baseload  $x$

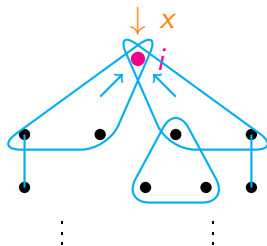


# Response Function

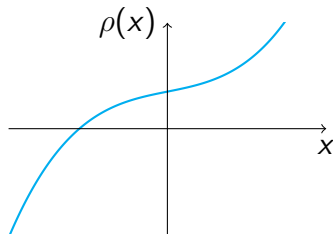
$\rho_{T,i}(x)$  = total load at  $i$  with baseload  $x$



# Response Function

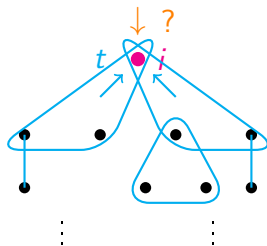
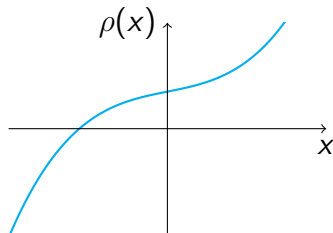
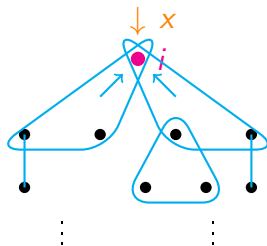


$\rho_{T,i}(x) = \text{total load at } i \text{ with baseload } x$



# Response Function

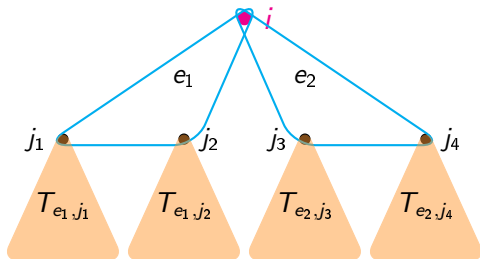
$\rho_{T,i}(x)$  = total load at  $i$  with baseload  $x$



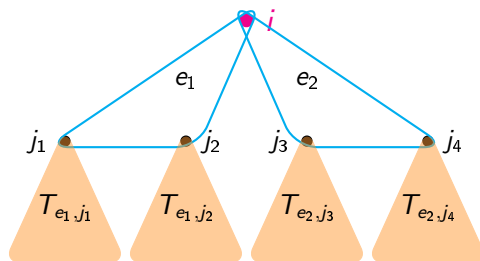
$\rho_{T,i}^{-1}(t)$ : the amount of extra load so that the total load becomes  $t$



# Recursion of Response Function

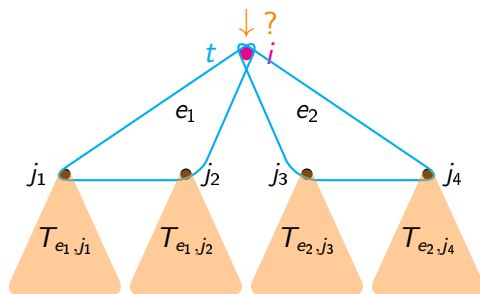


# Recursion of Response Function



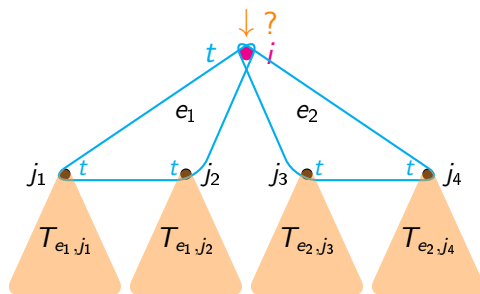
$$\rho_{T,i}^{-1}(t) = ?$$

# Recursion of Response Function



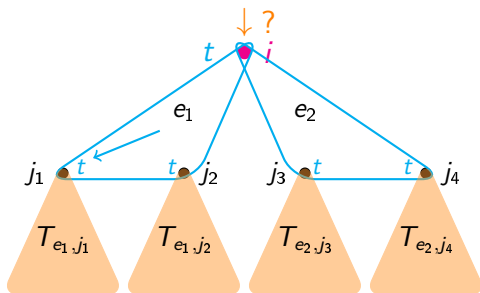
$$\rho_{T,i}^{-1}(t) = ?$$

# Recursion of Response Function



$$\rho_{T,i}^{-1}(t) = ?$$

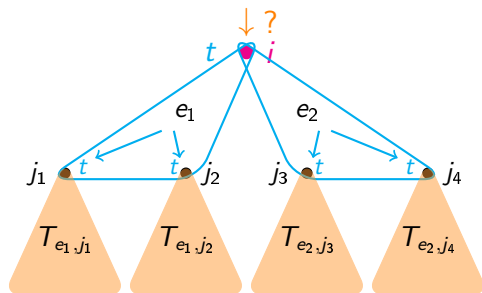
## Recursion of Response Function



$$\rho_{T,i}^{-1}(t) = ?$$

$$\theta(e_1, j_1) = \rho_{T_{e_1, j_1}}^{-1}(t)$$

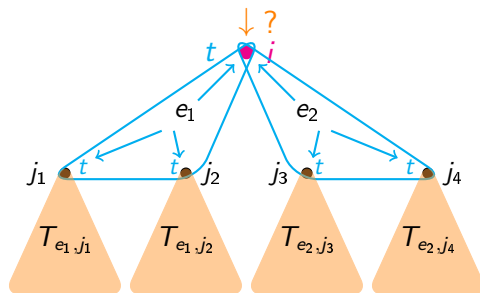
# Recursion of Response Function



$$\rho_{T,i}^{-1}(t) = ?$$

$$\theta(e_1, j_1) = \rho_{T_{e_1,j_1}}^{-1}(t)$$

# Recursion of Response Function

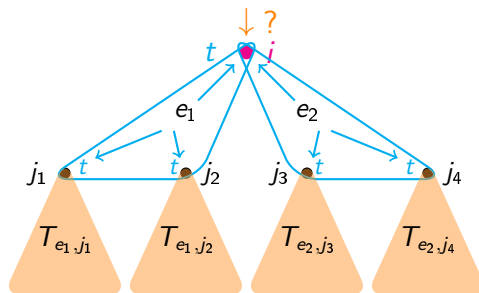


$$\rho_{T,i}^{-1}(t) = ?$$

$$\theta(e_1, j_1) = \rho_{T_{e_1,j_1}}^{-1}(t)$$

$$\theta(e_1, i) = 1 - \sum_{\substack{j \in e_1 \\ j \neq i}} \rho_{T_{e_1,j}}^{-1}(t)$$

# Recursion of Response Function



$$\rho_{T,i}^{-1}(t) = ?$$

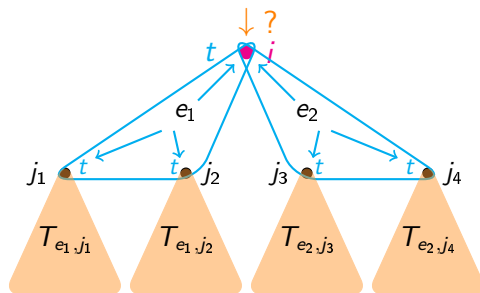
$$\theta(e_1, j_1) = \rho_{T_{e_1,j_1}}^{-1}(t)$$

$$\theta(e_1, i) = 1 - \sum_{\substack{j \in e_1 \\ j \neq i}} \rho_{T_{e_1,j}}^{-1}(t)$$

$$\rho_{T,i}^{-1}(t) = t - \sum_{e \ni i} \left( 1 - \sum_{\substack{j \in e \\ j \neq i}} \rho_{T_{e,j}}^{-1}(t) \right)$$



# Recursion of Response Function



$$\rho_{T,i}^{-1}(t) = ?$$

$$\theta(e_1, j_1) = \rho_{T_{e_1,j_1}}^{-1}(t)$$

$$\theta(e_1, i) = 1 - \sum_{\substack{j \in e_1 \\ j \neq i}} \rho_{T_{e_1,j}}^{-1}(t)$$

$$\rho_{T,i}^{-1}(t) = t - \sum_{e \ni i} \left( 1 - \sum_{\substack{j \in e \\ j \neq i}} \rho_{T_{e,j}}^{-1}(t) \right)$$

$$\rho_{T,i}^{-1}(t) = t - \sum_{e \ni i} \left[ 1 - \sum_{\substack{j \in e \\ j \neq i}} \rho_{T_{e,j}}^{-1}(t) \right]_0^1$$

# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)

# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers

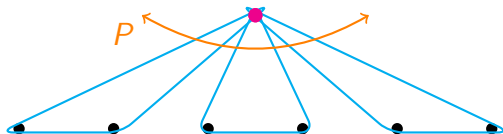
# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers



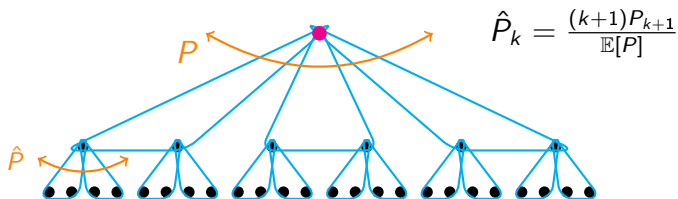
# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers



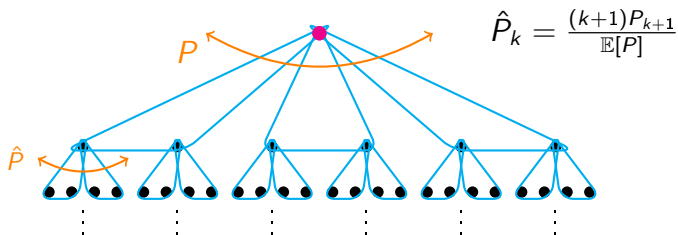
# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers



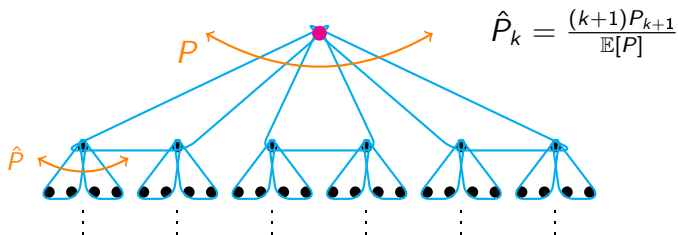
# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers



# Unimodular Galton Watson Hypertrees

- All the hyperedges have size  $c$  (say 3)
- distribution  $P$  on non-negative integers



- $\text{UGWT}_c(P) \in \mathcal{P}(\mathcal{H}_*)$  is unimodular



# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f:\mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

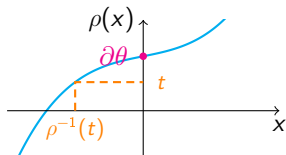
optimal  $f = \mathbb{1}[\partial\Theta > t]$

# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

optimal  $f = \mathbb{1}[\partial\Theta > t]$

$$\partial\Theta > t \Leftrightarrow \rho(0) > t \Leftrightarrow \rho^{-1}(t) < 0$$

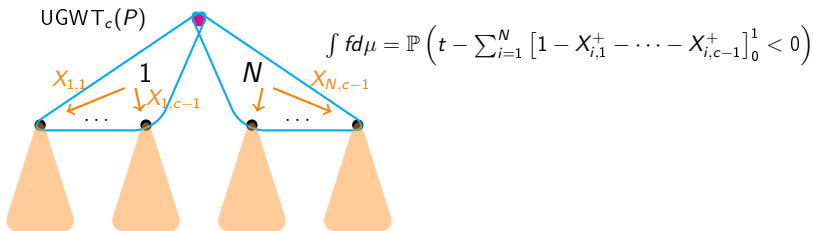
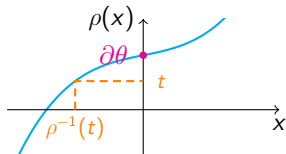


# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

optimal  $f = \mathbb{1}[\partial\Theta > t]$

$$\partial\Theta > t \Leftrightarrow \rho(0) > t \Leftrightarrow \rho^{-1}(t) < 0$$

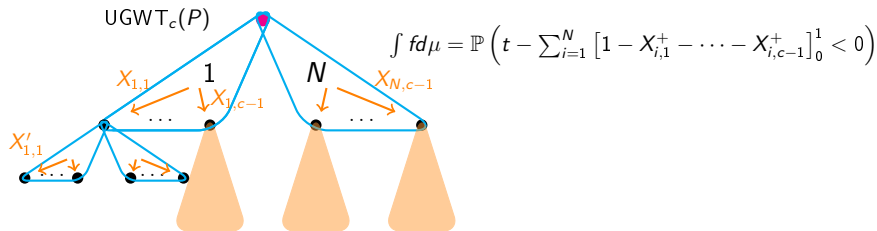
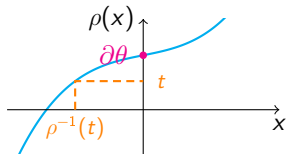


# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

optimal  $f = \mathbb{1}[\partial\Theta > t]$

$$\partial\Theta > t \Leftrightarrow \rho(0) > t \Leftrightarrow \rho^{-1}(t) < 0$$

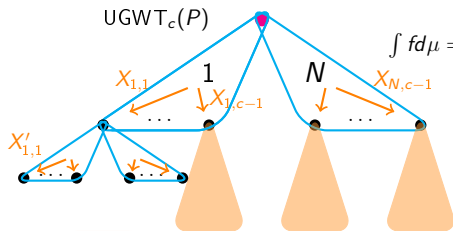
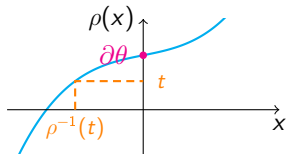


# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

optimal  $f = \mathbb{1}[\partial\Theta > t]$

$$\partial\Theta > t \Leftrightarrow \rho(0) > t \Leftrightarrow \rho^{-1}(t) < 0$$



$$\int f d\mu = \mathbb{P}\left(t - \sum_{i=1}^N [1 - X_{i,1}^+ - \dots - X_{i,c-1}^+]_0^1 < 0\right)$$

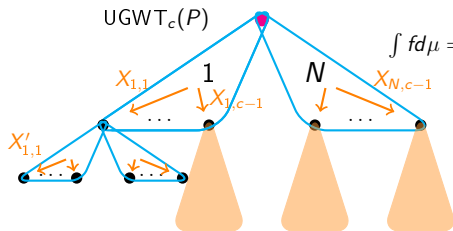
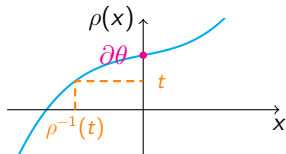
$$X_{1,1} = t - \sum_{j=1}^{\hat{N}} [1 - X_{j,1}'^+ - \dots - X_{j,c-1}'^+]_0^1$$

# Mean Excess Characterization for $\text{UGWT}_c(P)$

$$\int (\partial\Theta - t)^+ d\mu = \sup_{f: \mathcal{H}_* \rightarrow [0,1]} \int \tilde{f}_{\min} d\vec{\mu} - t \int f d\mu$$

optimal  $f = \mathbb{1}[\partial\Theta > t]$

$$\partial\Theta > t \Leftrightarrow \rho(0) > t \Leftrightarrow \rho^{-1}(t) < 0$$



$$\int f d\mu = \mathbb{P}\left(t - \sum_{i=1}^N [1 - X_{i,1}^+ - \dots - X_{i,c-1}^+]_0^1 < 0\right)$$

$$X_{1,1} = t - \sum_{j=1}^{\hat{N}} [1 - X'_{j,1}^+ - \dots - X'_{j,c-1}^+]_0^1$$

$$Q = F_{P,t}^c(Q)$$

# Mean Excess Characterization for $UGWT_c(P)$

(cont'd)

## Theorem

Assume  $P$  is a distribution on nonnegative integers with finite variance and  $\mu = UGWT^c(P)$ . Then, we have

$$\int (\partial\Theta - t)^+ d\mu = \max_{Q: F_{P,t}^c(Q)=Q} \frac{\mathbb{E}[N]}{c} \mathbb{P}(X_1^+ + \dots + X_c^+ < 1) \\ - t\mathbb{P}(Y_1 + \dots + Y_N > t),$$

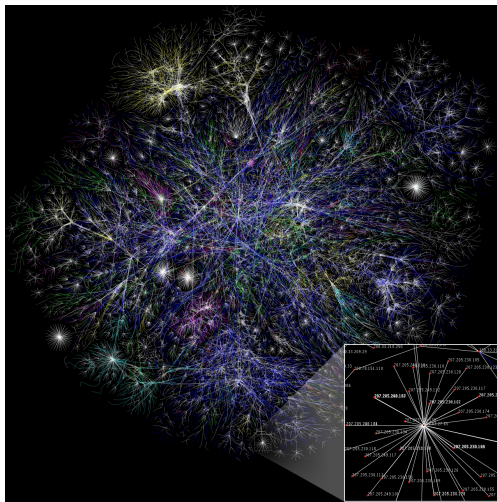
where  $N$  has distribution  $P$ ,  $X_i$ 's are independent and have distribution  $Q$  and  $Y_i$ 's are independent and each have distribution of  $[1 - (X_1^+ + \dots + X_{c-1}^+)]_0^1$  where  $X_i$ 's are i.i.d. from  $Q$ .



# Outline

- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs
- 5 Graph indexed data**
- 6 Universal compression of graphical data

# Sources of big graphical data: The web



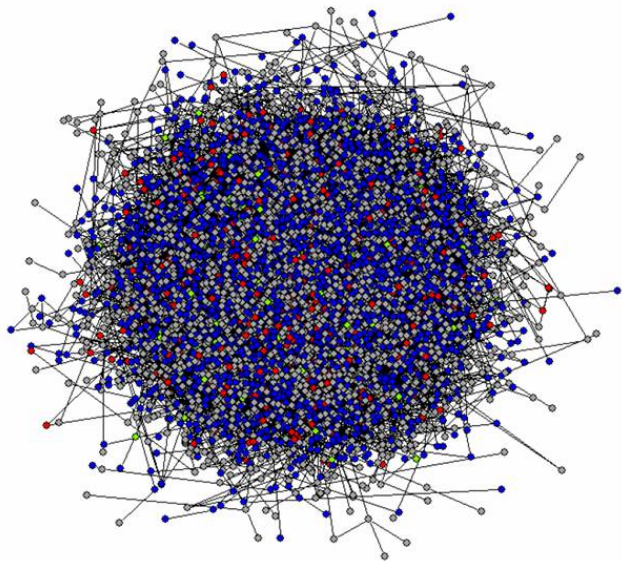
$\approx 47$  billion webpages

# Sources of big graphical data: Social networks



$\approx$  1.8 billion active users on Facebook

# Sources of big graphical data: Biological networks



0.25 million - 1 million estimated human proteins

# Universal compression of marked graphical data

- We want to compress down to the “entropy” of the data.
- Universality means that the scheme should work irrespective of the underlying “statistics” of the data.
- Ideally, the compressed representation should enable analysis and querying in the compressed form

# Universal compression of marked graphical data

- We want to compress down to the “entropy” of the data.
- Universality means that the scheme should work irrespective of the underlying “statistics” of the data.
- Ideally, the compressed representation should enable analysis and querying in the compressed form

The local weak limit theory allows one to precisely formulate the universal compression problem and to provide a solution.

# The BC entropy: counting typical graphs

- $\Xi$ : edge marks,  $\Theta$ : vertex marks, both finite
- $\mathcal{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)}$ : set of graphs on  $n$  vertices with  $\mathbf{m}_n(x)$  many edges with mark  $x \in \Xi$  and  $\mathbf{u}_n(t)$  many vertices with mark  $t \in \Theta$ .
- $\mathcal{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)}(\mu, \epsilon) = \{G \in \mathcal{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)} : U(G) \in B(\mu, \epsilon)\}$ .
- For  $\mu \in \mathcal{P}(\mathcal{G}_*)$  and  $x \in \Xi$ ,  $\deg_x(\mu)$ : expected number of edges connected to the root with mark  $x$ ,
- $t \in \Theta$ ,  $\Pi_t(\mu)$ : probability of root having mark  $t$ .

# The BC entropy: counting typical graphs

- Fix sequences  $\mathbf{m}_n, \mathbf{u}_n$  such that  $\mathbf{m}_n(x)/n \rightarrow \deg_x(\mu)/2$  and  $\mathbf{u}_n(t)/n \rightarrow \Pi_t(\mu)$  for all  $x \in \Xi, t \in \Theta$ .
- $\log |\mathcal{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)}| = \|\mathbf{m}_n\|_1 \log n + cn + o(n)$  where  $\|\mathbf{m}_n\|_1 = \sum_{x \in \Xi} \mathbf{m}_n(x)$ .
- 

$$\bar{\Sigma}(\mu) := \lim_{\epsilon \downarrow 0} \limsup_{n \rightarrow \infty} \frac{\log |\mathbf{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)}(\mu, \epsilon)| - \|\mathbf{m}_n\|_1 \log n}{n}$$

$$\underline{\Sigma}(\mu) := \lim_{\epsilon \downarrow 0} \liminf_{n \rightarrow \infty} \frac{\log |\mathbf{G}_{\mathbf{m}_n, \mathbf{u}_n}^{(n)}(\mu, \epsilon)| - \|\mathbf{m}_n\|_1 \log n}{n}$$

- If they are equal, define the common value as  $\Sigma(\mu)$   
(Generalizing work of [Bordenave](#) and [Caputo](#))



# Outline

- 1 A resource allocation problem studied by Hajek
- 2 Load balancing on graphs
- 3 The framework of local weak convergence
- 4 Load balancing on hypergraphs
- 5 Graph indexed data
- 6 Universal compression of graphical data

# Our target for the graph regime

- **Goal:** design  $f_n : \mathcal{G}_n \rightarrow \{0, 1\}^*$  and  $g_n : \{0, 1\}^* \rightarrow \mathcal{G}_n$
- $g_n \circ f_n = \text{Id}$
- $\mu \in \mathcal{P}(\mathcal{G}_*)$  a process
- Target: typical graphs
- **Optimal** if  $G_n \xrightarrow{\text{lwc}} \mu$

$$\limsup_{n \rightarrow \infty} \frac{l(f_n(G_n)) - m_n \log n}{n} \leq \Sigma(\mu),$$

where  $m_n$  is the total number of edges in  $G_n$ .

# A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \deg \leq \Delta_n\}$$

# A First Step Coding Scheme : Example

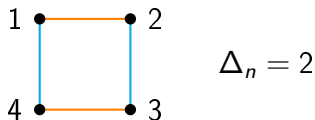
$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \text{deg} \leq \Delta_n\}$$

$$n = 4, k_n = 1$$

# A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \text{deg} \leq \Delta_n\}$$

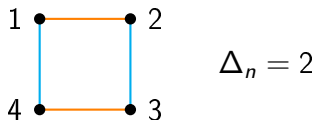
$$n = 4, k_n = 1$$



# A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \text{deg} \leq \Delta_n\}$$

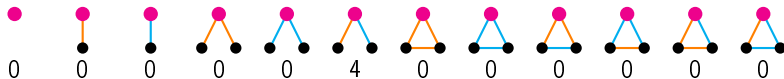
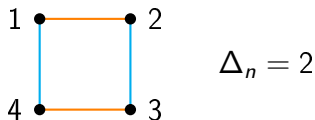
$$n = 4, k_n = 1$$



## A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \deg \leq \Delta_n\}$$

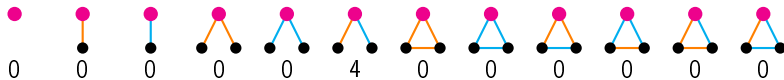
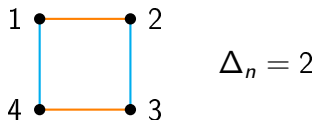
$$n = 4, \quad k_n = 1$$



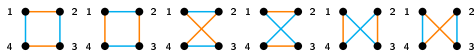
# A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \deg \leq \Delta_n\}$$

$$n = 4, k_n = 1$$



$W_n :=$  the set of graphs with the same sequence

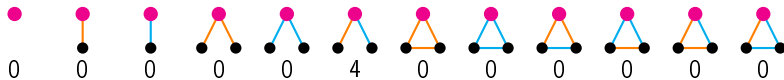
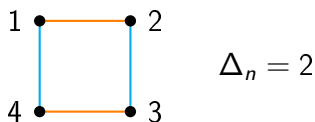




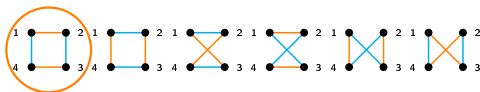
# A First Step Coding Scheme : Example

$$\mathcal{A}_{k_n, \Delta_n} = \{[G, o] \in \mathcal{G}_* : \text{depth} \leq k_n, \max \deg \leq \Delta_n\}$$

$$n = 4, k_n = 1$$



$W_n :=$  the set of graphs with the same sequence



# Analysis Outline

- $I(f_n(G_n))$ , the total number of bits we use:
  - ▶  $\log n$  bits for  $\Delta_n$ ,
  - ▶  $|\mathcal{A}_{k_n, \Delta_n}| \log n$  bits for specifying how many times each pattern appears in the graph
  - ▶  $\log |W_n|$  bits to specify the input graph among the graphs with the same pattern counts.
- We need to show that if  $G_n \xrightarrow{IWC} \mu$ ,

$$\frac{I(f_n(G_n)) - m_n \log n}{n} \leq \bar{\Sigma}(\mu).$$

- If  $|\mathcal{A}_{k_n, \Delta_n}| = o(n/\log n)$ , we only need to consider the  $\log |W_n|$  term.
- Graphs in  $W_n$  are typical  $\Rightarrow$  yields  $\bar{\Sigma}(\mu)$  as an upper bound.

# First step algorithm: Main Result

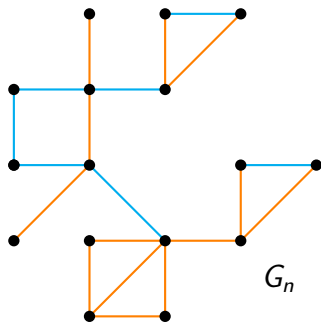
## Proposition

*If parameters  $k_n$  and  $\Delta_n$  are such that  $|\mathcal{A}_{k_n, \Delta_n}| = o(\frac{n}{\log n})$  and  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$ , for any sequence  $G_n$  with maximum degree no more than  $\Delta_n$  and local weak limit  $\mu \in \mathcal{P}(\mathcal{G}_*)$  such that  $\bar{\Sigma}(\mu) > -\infty$  we have*

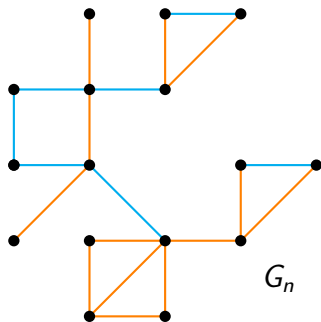
$$\limsup_{n \rightarrow \infty} \frac{I(f_n(G_n)) - m_n \log n}{n} \leq \bar{\Sigma}(\mu), \quad (1)$$

*where  $m_n$  is the number of edges in  $G_n$ .*

# General Algorithm

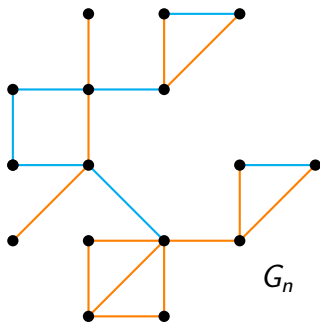


# General Algorithm

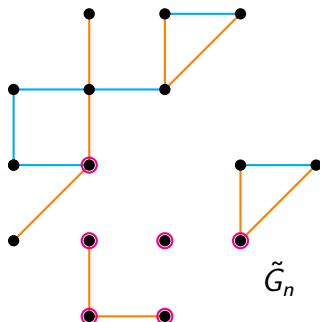


$$\Delta_n \xrightarrow{=} 5$$

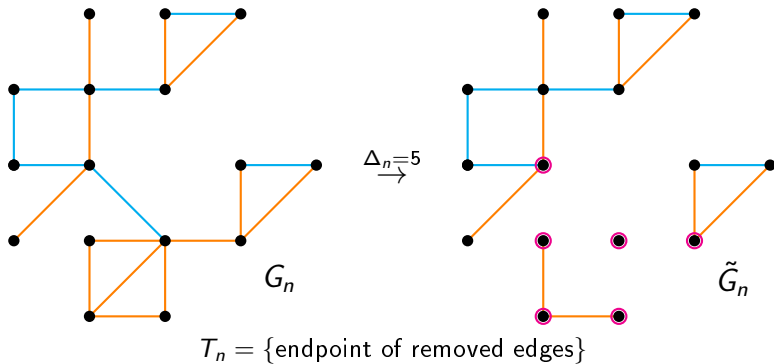
# General Algorithm



$\Delta_n \stackrel{=}{\rightarrow} 5$



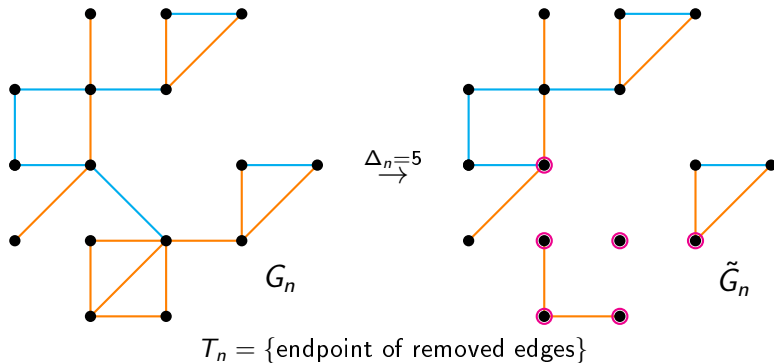
# General Algorithm





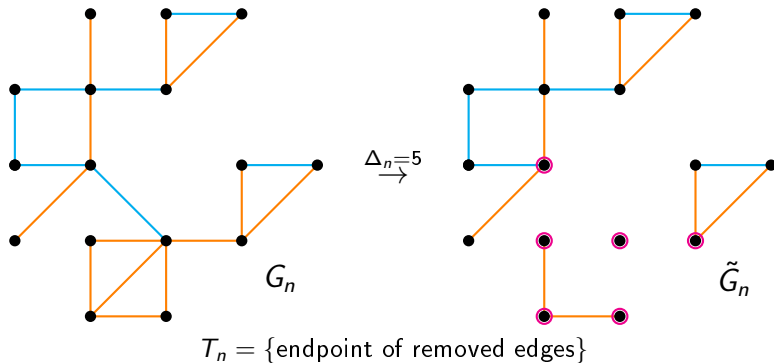


# General Algorithm



Compress  $\tilde{G}_n$  using the first step scheme, then compress removed edges  $\Delta_n = \log \log n$   $k_n = \sqrt{\log \log n}$

# General Algorithm



Compress  $\tilde{G}_n$  using the first step scheme, then compress removed edges  $\Delta_n = \log \log n$   $k_n = \sqrt{\log \log n}$

$$|T_n|/n \rightarrow 0 \quad |\mathcal{A}_{k_n, \Delta_n}| = o(n/\log n) \quad G_n \xrightarrow{lwc} \mu \Rightarrow \tilde{G}_n \xrightarrow{lwc} \mu$$

# Result: Achievability

## Theorem

Assume  $\mu \in \mathcal{G}_*$  with  $\deg_x(\mu) < \infty$  for all  $x$  and  $\bar{\Sigma}(\mu) > -\infty$ . If  $G_n$  is a sequence of marked graphs with local weak limit  $\mu$ , we have

$$\limsup_{n \rightarrow \infty} \frac{I(f_n(G_n)) - m_n \log n}{n} \leq \bar{\Sigma}(\mu),$$

where  $m_n$  is the number of edges in  $G_n$ .

# Result: Converse

## Theorem

*Assume  $\mu \in \mathcal{P}(\mathcal{G}_*)$  with  $\underline{\Sigma}(\mu) > -\infty$  and  $\deg_x(\mu) < \infty$  for all  $x \in \Xi$ . Then there exists a sequence of graph ensembles  $G_n$  converging to  $\mu$  such that with probability one for any sequence of compression schemes  $f_n$  we have*

$$\liminf_{n \rightarrow \infty} \frac{I(f_n(G_n)) - m_n \log n}{n} \geq \underline{\Sigma}(\mu),$$

*where  $m_n$  is the number of edges in  $G_n$ .*

# Concluding remarks

- Just like a stochastic process is a model for the statistics of a long string of data, a local weak limit of marked graphs is a model for the statistics of data that lives on large graphs.

# Concluding remarks

- Just like a stochastic process is a model for the statistics of a long string of data, a local weak limit of marked graphs is a model for the statistics of data that lives on large graphs.
- This provides a methodology to address networking problems and data centric problems arising in networks that parallels how stochastic processes are used in the study of time series.

# Concluding remarks

- Just like a stochastic process is a model for the statistics of a long string of data, a local weak limit of marked graphs is a model for the statistics of data that lives on large graphs.
- This provides a methodology to address networking problems and data centric problems arising in networks that parallels how stochastic processes are used in the study of time series.
- This was illustrated with two kinds of applications: **resource allocation** in graphs and hypergraphs, and **universal lossless compression** of graph-structured data.

# Concluding remarks

- Just like a stochastic process is a model for the statistics of a long string of data, a local weak limit of marked graphs is a model for the statistics of data that lives on large graphs.
- This provides a methodology to address networking problems and data centric problems arising in networks that parallels how stochastic processes are used in the study of time series.
- This was illustrated with two kinds of applications: **resource allocation** in graphs and hypergraphs, and **universal lossless compression** of graph-structured data.
- A world of other applications awaits.



# The End

