

# Sparsity in Practice

*a bit of introspection*



Felix Reidl  
Dagstuhl 2021

# *Part I*

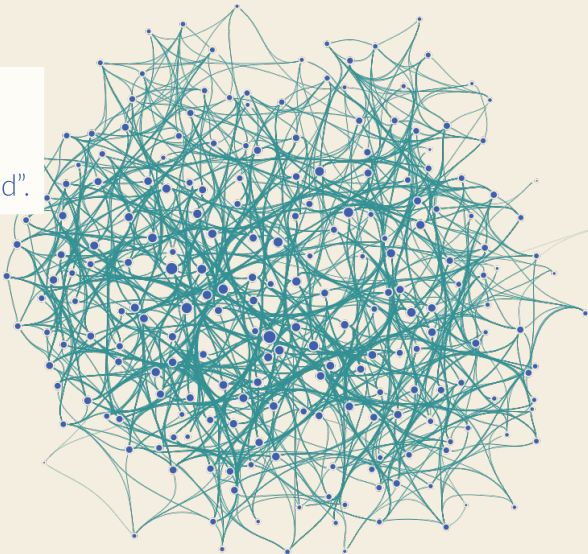
The work so far



# Residence hall

- Student in ANU Hall
- Friendship

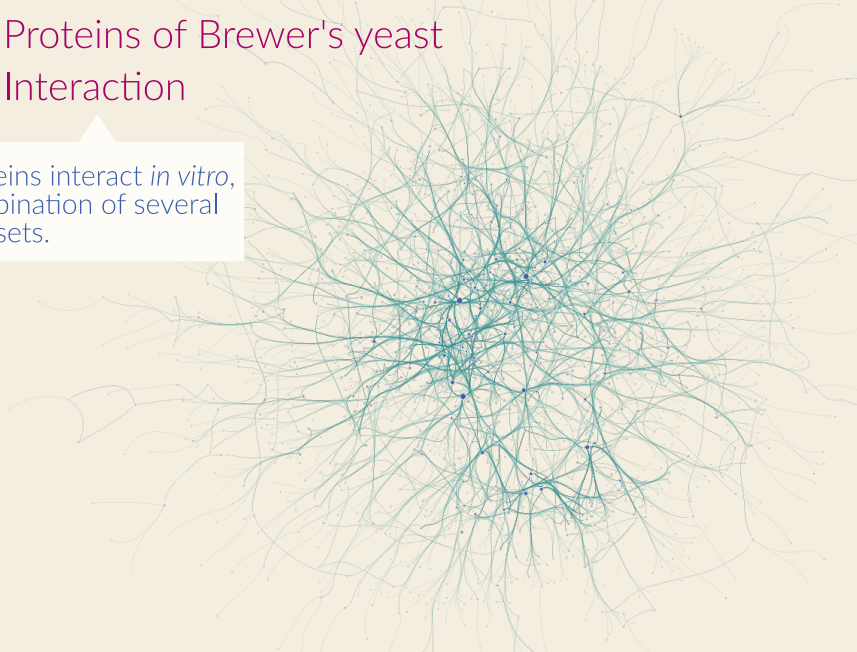
Collected via interviews by Cynthia Webster, ranked as “best friend”, “close friend”, and “friend”.



# Y2H union (yeast)

- Proteins of Brewer's yeast
- Interaction

Proteins interact *in vitro*,  
combination of several  
datasets.

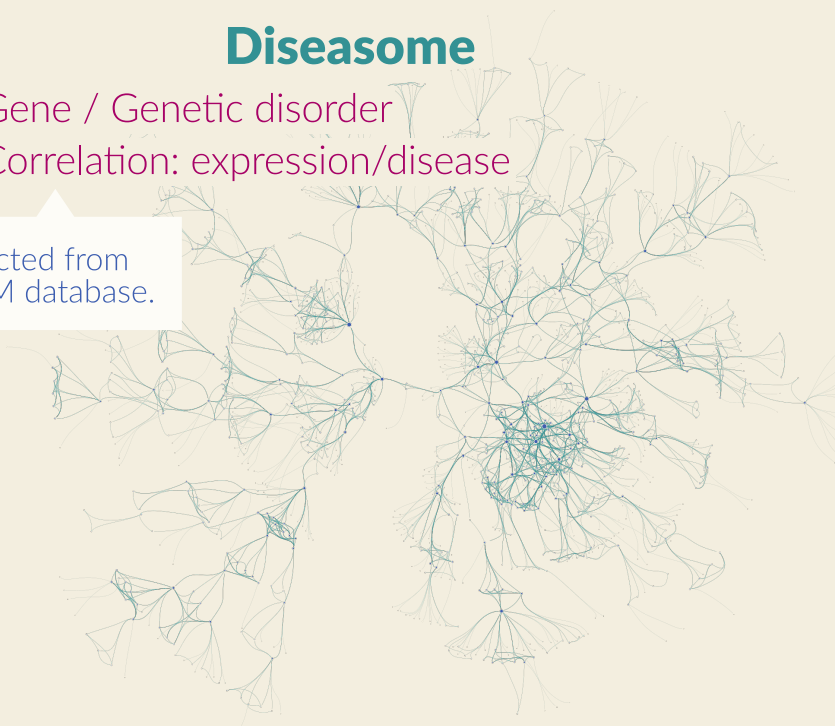




# Diseaseome

- Gene / Genetic disorder
- Correlation: expression/disease

Extracted from  
OMIM database.

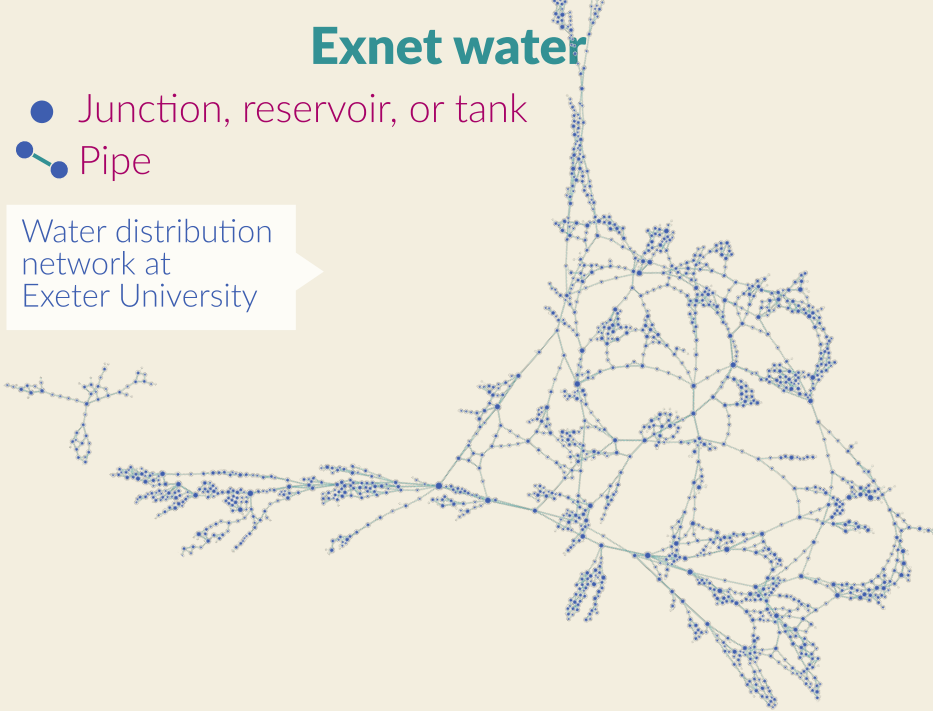


# Exnet water

● Junction, reservoir, or tank

● Pipe

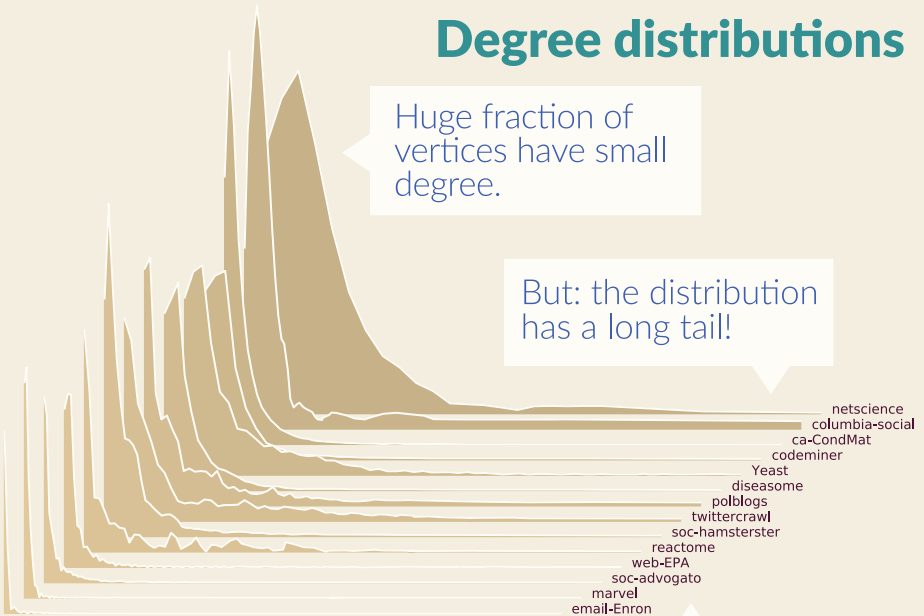
Water distribution network at Exeter University



# Degree distributions

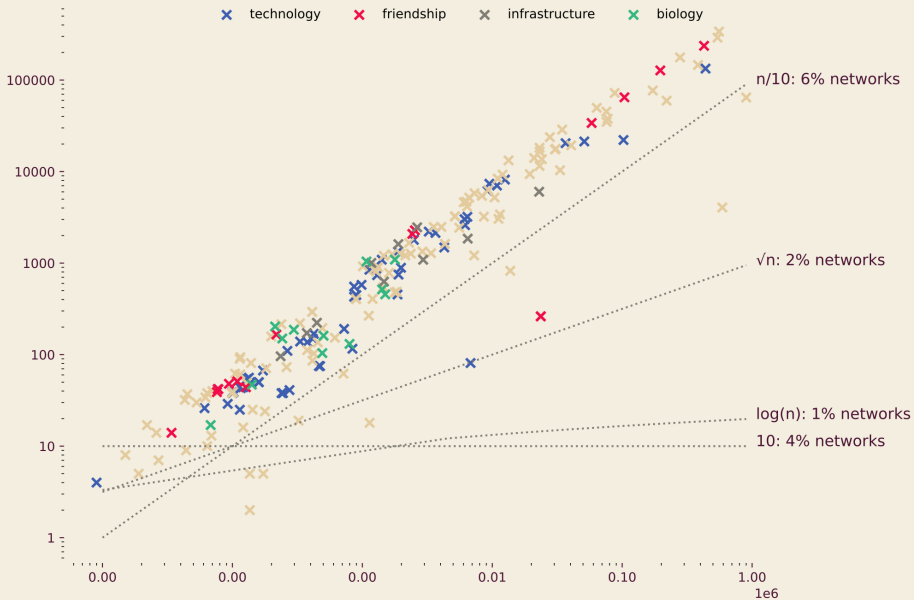
Huge fraction of vertices have small degree.

But: the distribution has a long tail!

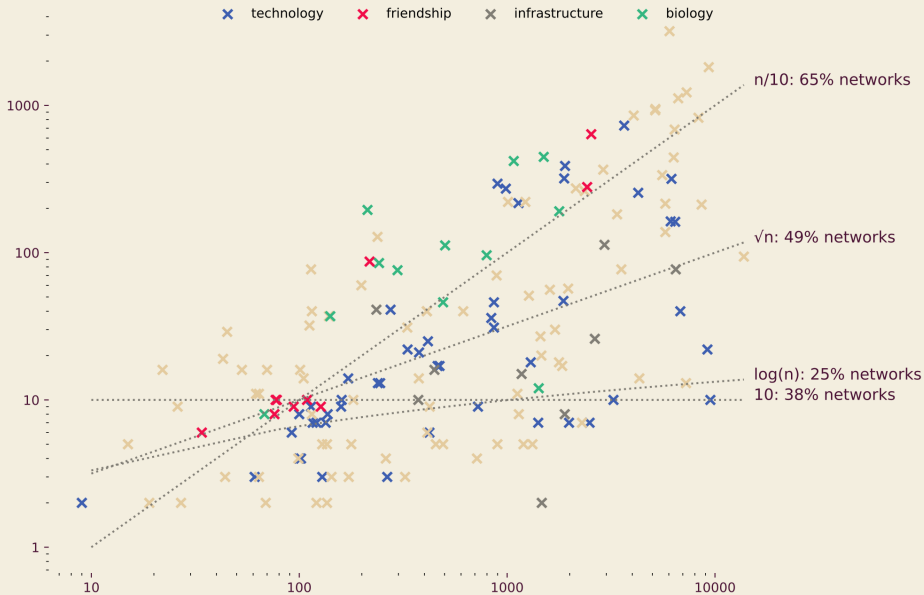


More extreme in large networks!

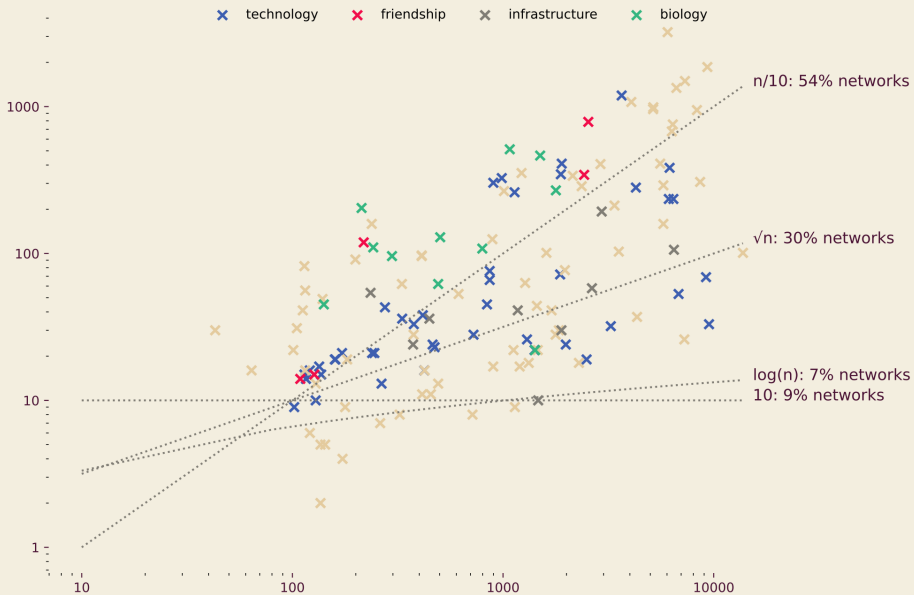
# Real vertex cover



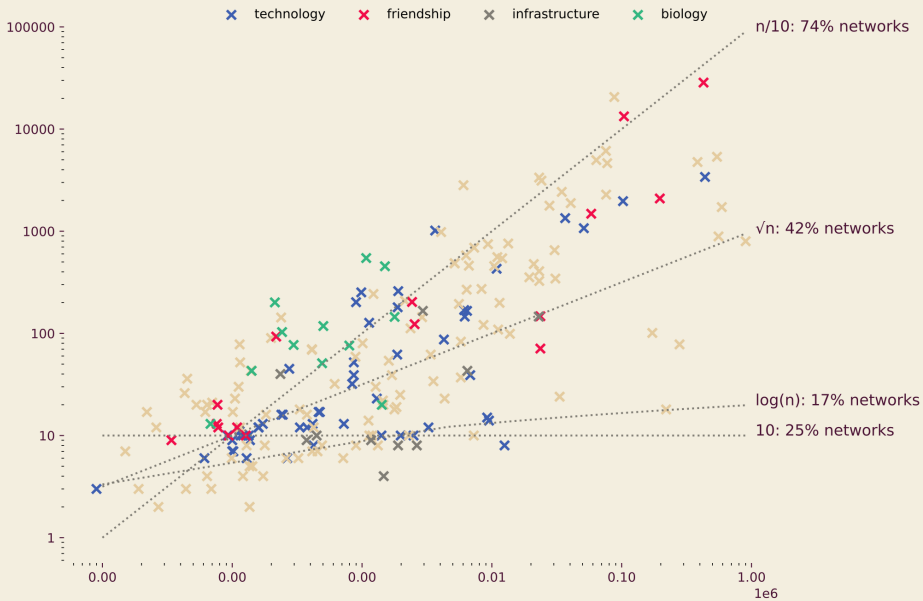
# Real treewidth



# Real treedepth

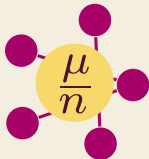


# Real $wcol_3$



# Random model sparsity

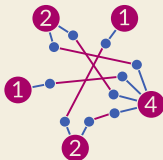
$$G(n, \frac{\mu}{n})$$



$$G^{CL}(D_n)$$



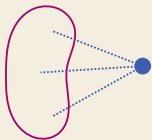
$$G^{CF}(D_n)$$



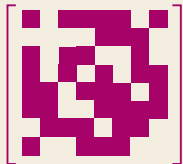
$$G^{RIG}(n, \alpha, \beta, \gamma)$$



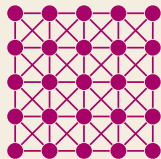
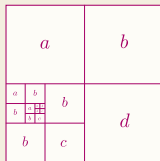
$$G^{BA}(n, n_0, k)$$



$$G^{SGK}(k, \alpha, \dots, \gamma)$$

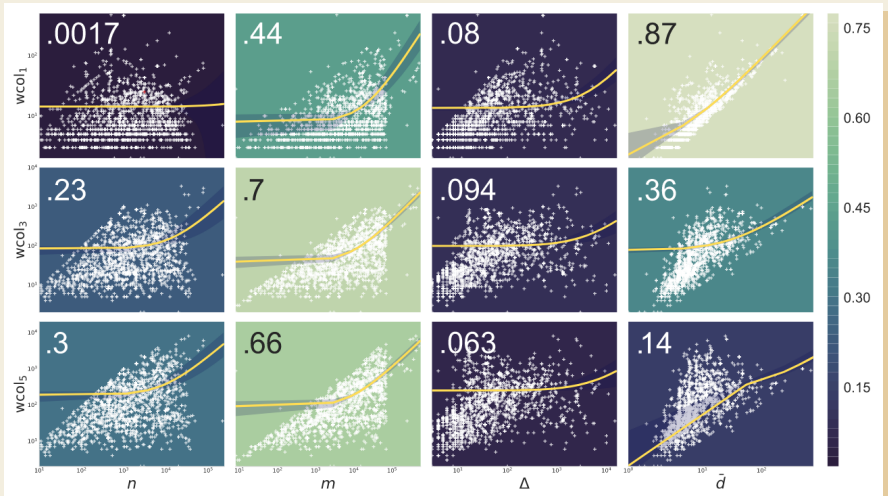


$$G^{RMAT}(k, m, a, b, c) \quad G^{KL}(n, p, q, \gamma)$$





# Real structural sparseness



W. Nadara, M. Pilipczuk, R. Rabinovich, FR, S. Siebertz:  
Empirical Evaluation of Approximation Algorithms for  
Generalized Graph Coloring and Uniform Quasi-Widness.  
SEA 2018: 14:1-14:16

# A hard-learnt lesson



Nešetřil J, de Mendez PO, Wood DR.  
**Characterisations and examples of graph classes with bounded expansion.**  
European Journal of Combinatorics. 2012 Apr 30;33(3):350-73.

Demaine ED, FR, Rossmanith P, Sánchez Villamil F, Sikdar S, Sullivan BD.  
**Structural sparsity of complex networks: Bounded expansion in random models and real-world graphs.**  
Journal of Computer and System Sciences. 2019 May 24.

Farrell M, Goodrich TD, Lemons N, FR, Sánchez Villamil F, Sullivan BD.  
**Hyperbolicity, degeneracy, and expansion of random intersection graphs.**  
In International Workshop on Algorithms and Models for the Web-Graph 2015  
Dec 10 (pp. 29-41). Springer, Cham.

W. Nadara, M. Pilipczuk, R. Rabinovich, FR, S. Siebertz:  
**Empirical Evaluation of Approximation Algorithms for Generalized Graph Coloring and Uniform Quasi-Widness.**  
SEA 2018: 14:1-14:16

# A hard-learnt lesson



Nešetřil J, de Mendez PO, Wood DR.  
**Characterisations and examples of graph classes with bounded expansion.**  
European Journal of Combinatorics. 2012 Apr 30;33(3):350-73.

Demaine ED, FR, Rossmanith P, Sánchez Villalamil F, Sikdar S, Sullivan BD.  
**Structural sparsity of complex networks: Bounded expansion in random models and real-world graphs.**  
Journal of Computer and System Sciences. 2019 May 24.

Farrell M, Goodrich TD, Lemons N, FR, Sánchez Villalamil F, Sullivan BD.  
**Hyperbolicity, degeneracy, and expansion of random intersection graphs.**  
In International Workshop on Algorithms and Models for the Web-Graph 2015 Dec 10 (pp. 29-41). Springer, Cham.

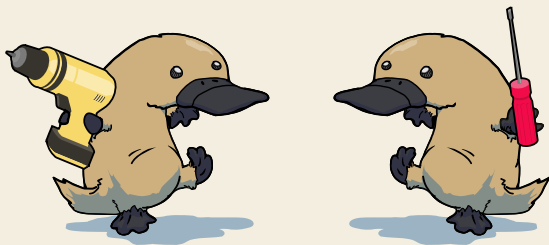
W. Nadara, M. Pilipczuk, R. Rabinovich, FR, S. Siebertz:  
**Empirical Evaluation of Approximation Algorithms for Generalized Graph Coloring and Uniform Quasi-Wideness.**  
SEA 2018: 14:1-14:16

**No one<sup>\*</sup>  
cares.**

<sup>\*</sup>not enough people

## *Part II*

# Reflection



# A false dichotomy



# Reality\*



**ALGORITHMS**



**NETWORK SCIENCE**



**DATA SCIENCE**



**AI**



**DATABASES**



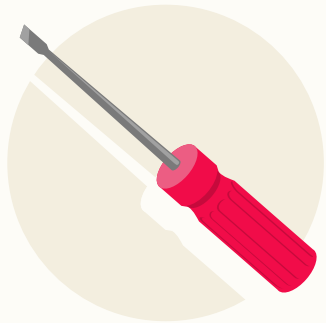
**BIOINFORMATICS**



**INDUSTRIES**

\*Highly subjective opinion

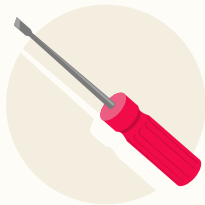
# Solvers **vs** Solutions



# Solvers vs Solutions



Communities



Individuals



# Solvers vs Solutions

Diffuse problems

Competing views

Often sceptical of new approaches

Concrete problems

Single perspective

Appreciate new approaches



Must be very user-friendly

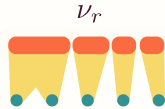
High maintenance

Usable for stakeholder

Maintenance by stakeholder



# The hammer scale



$$\Delta^-(\vec{G}_r)$$



FO model  
checking



r-Dominating  
Set approx.



splitter games

Usability  
in *practice*



Usability  
in *theory*





**Big hammers don't implement**



# Medium hammers don't scale

O'Brien MP, Sullivan BD.

Experimental evaluation of counting subgraph isomorphisms  
in classes of bounded expansion.

arXiv preprint arXiv:1712.06690. 2017 Dec 18.



## Small hammers might just work!

Nadara W, Pilipczuk M, Rabinovich R, Reidl F, Siebertz S.  
**Empirical evaluation of approximation algorithms for generalized graph coloring and uniform quasi-wideness.**  
Journal of Experimental Algorithmics (JEA). 2019 Dec 10;24:1-34.

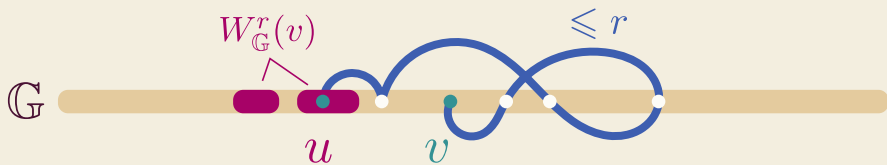
Brown CT, Moritz D, O'Brien MP, Reidl F, Reiter T, Sullivan BD.  
**Exploring neighborhoods in large metagenome assembly graphs using spacegraphcats reveals hidden sequence diversity.** Genome biology. 2020 Dec;21(1):1-6.

[github.com/  
spacegraphcats/  
spacegraphcats](https://github.com/spacegraphcats/spacegraphcats)

Reidl F, Sullivan BD. **A color-avoiding approach to subgraph counting in bounded expansion classes.**  
arXiv preprint arXiv:2001.05236. 2020 Jan 15.

[github.com/  
theoryinpractice/  
mandoline](https://github.com/theoryinpractice/mandoline)

# Weak colouring & bounded expansion



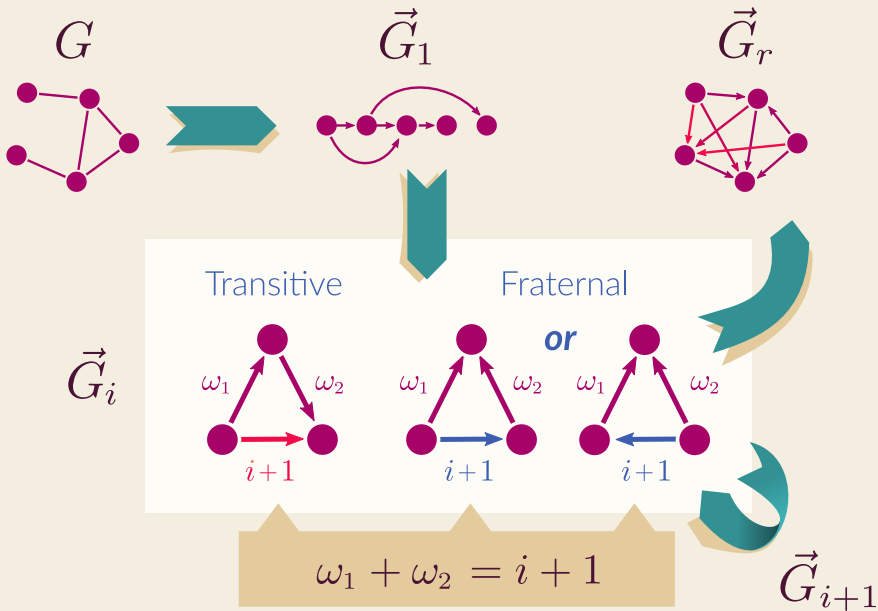
$u$  is **weakly  $r$ -reachable** from  $v$  if there exists a path from  $v$  to  $u$  of length at most  $r$  such that  $u$  is the path's leftmost vertex.

$$\text{wcol}_r(G) := \min_{G \in \Pi(G)} \max_{v \in G} |W_G^r(v)|$$



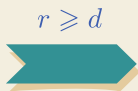
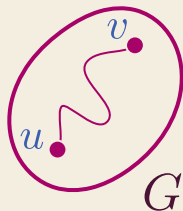
A graph class has bounded expansion iff it is  $\text{wcol}_r$ -bounded.

# dtf-augmentations



# Distances under dtf-augmentations

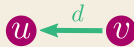
Let  $u$  and  $v$  be at distance  $d$  in  $G$  :



1

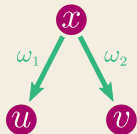


2



$\vec{G}_r$

3



$$\omega_1 + \omega_2 = d$$

Pairs at distance at most  $r$  in the original graph have distance at most two in the  $r^{\text{th}}$  augmentation.



## B.E. & dtf-augmentations

There exist two (horrible) polynomials  $P$  and  $Q$  such that:

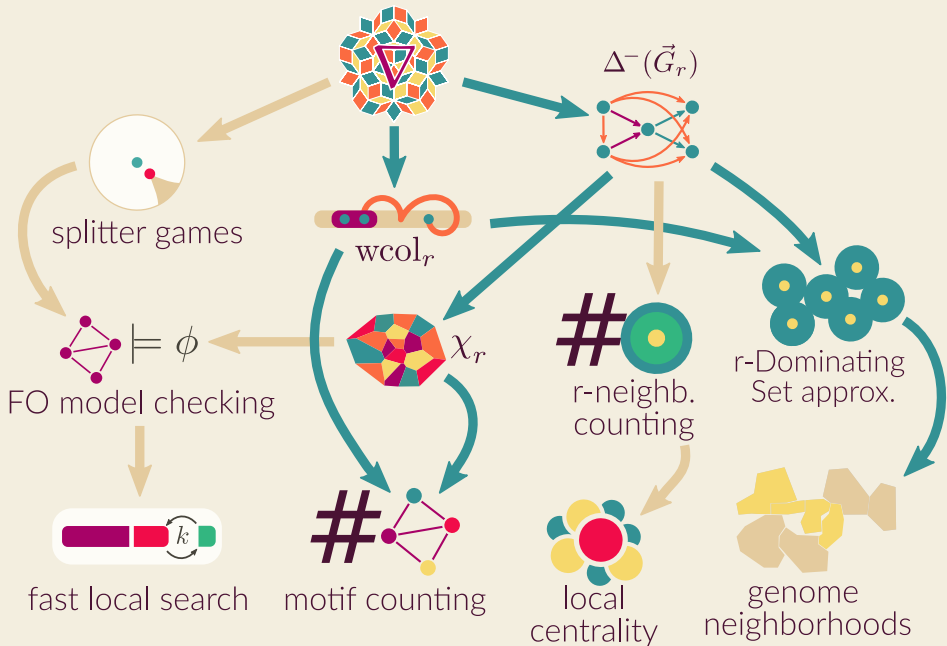
$$\chi_r(G) \leq P(\tilde{\nabla}_{(2 \log r)^r}(G))$$
$$\Delta^-(\vec{G}_r) \leq Q(\tilde{\nabla}_r(G) \Delta^-(\vec{G}_1))$$



A graph class has bounded expansion iff it is  $\Delta^-(\vec{G}_r)$ -bounded.

We can compute dtf-augmentations in linear time (in bounded expansion classes)

# Applications & Algorithms



# Exhibit A

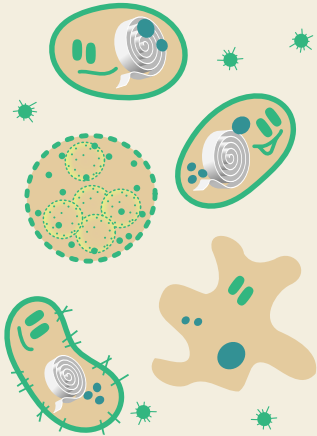
# CATLAS



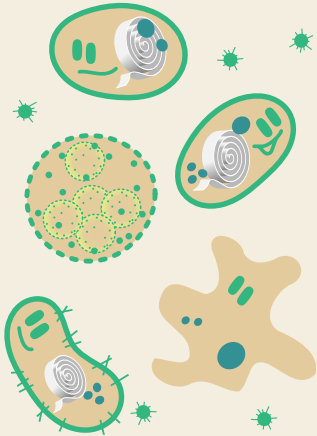
Metagenome exploration  
using hierarchical domination  
of de-Bruijn graphs

Joint work with C. Titus Brown, Dominik Moritz,  
Michael P. O'Brien, Taylor Reiter, Blair D. Sullivan

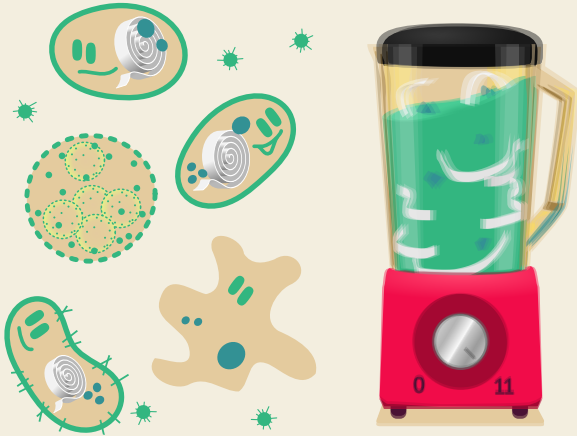
# Metagenomics



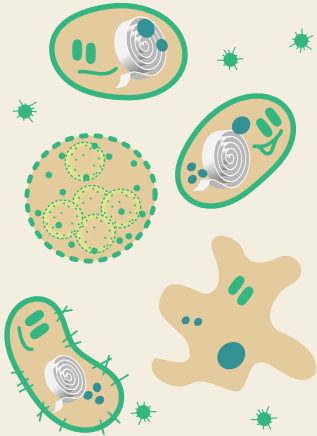
# Metagenomics



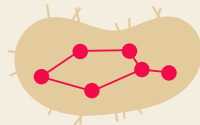
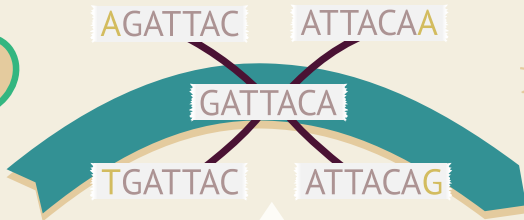
# Metagenomics



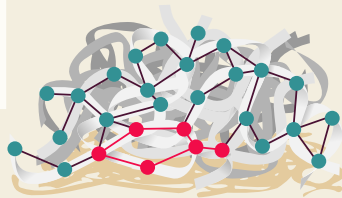
# Metagenomics



# De-Bruijn graphs



Bounded  
degree

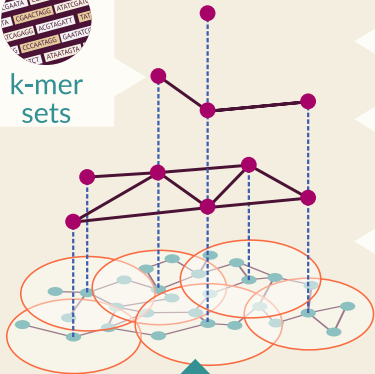




# CATLAS Overview



k-mer  
sets



Domset



r-Domset



r-DTFAs



Dvořák's  
Algorithm\*



de-Bruijn  
graphs

# Reminder: Dvořák's algorithm



$wcol_{2r}$



$\vec{G}_{2r}$

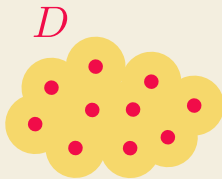
while ● undominated:

$W^{2r}(\bullet)$

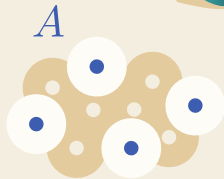


$N_{2r}^-(\bullet)$

$$|D| \leq c_r |A|$$

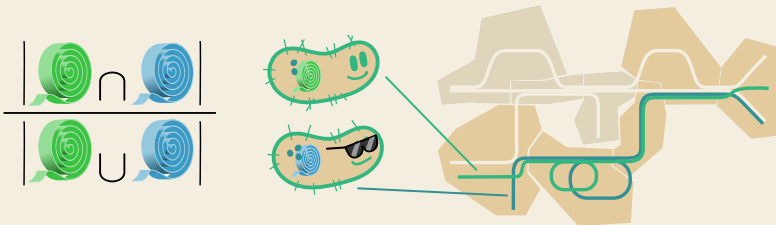
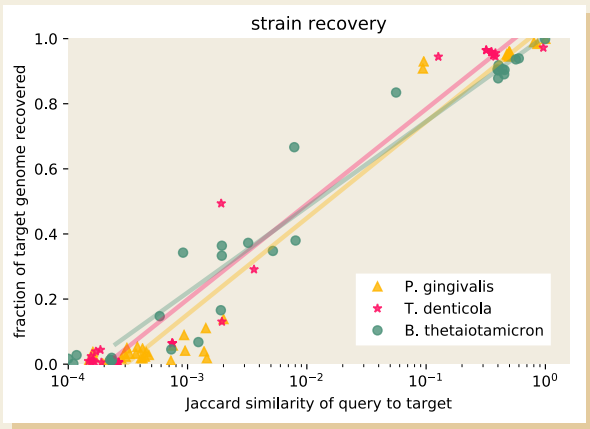


r-dominating set



r-scattered set

# CATLAS-1 Results



# What really matters



“The efficient computation of dominating sets will open up a whole new range of possibilities in bioinformatics.”

C. Titus Brown, Associate Professor at UC Davis

“spacegraphcats will transform the way biologists interact with genome assemblies.

It allows us to access previously discarded sequencing information thereby allowing more robust functional characterization.”



Taylor Reiter, his much more eloquent post-doc

# Engineering: efficiency

**Lesson:** Stick to the 'sparsity methodology'.

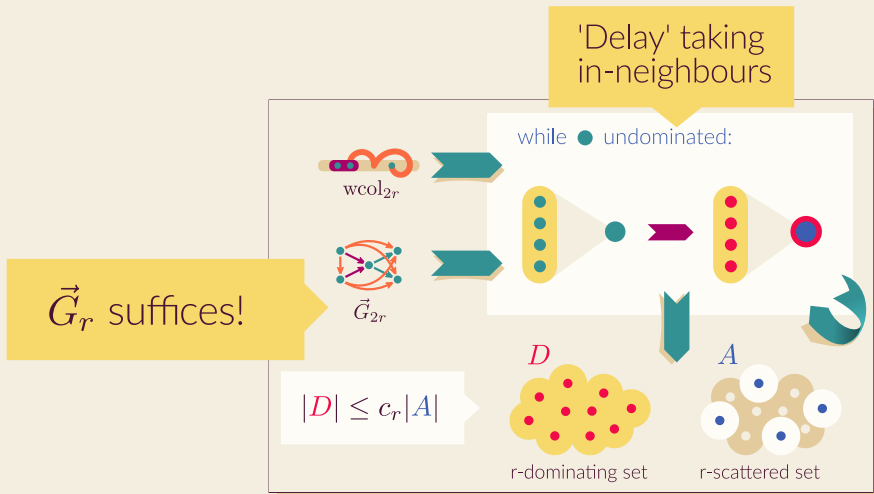
Example: computing partition from r-domset



We can do this in linear time using the already computed dtf-augmentation

# Engineering: efficiency

**Lesson:** Practical issues inspire theoretical questions.



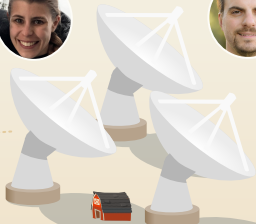
# Collaboration

**Lesson:** Expect roles to shift over time.

Then



Now



# Next steps

- 1 Faster language  
Believe it or not, so far this is all done in Python!
- 2 Improve network partition  
E.g. consider additional constraints
- 3 Whatever our collaborators need!



## ? Contracted de-Bruijn graphs

These graphs have bounded degree, but probably much more structure that we could exploit!





## Exhibit B

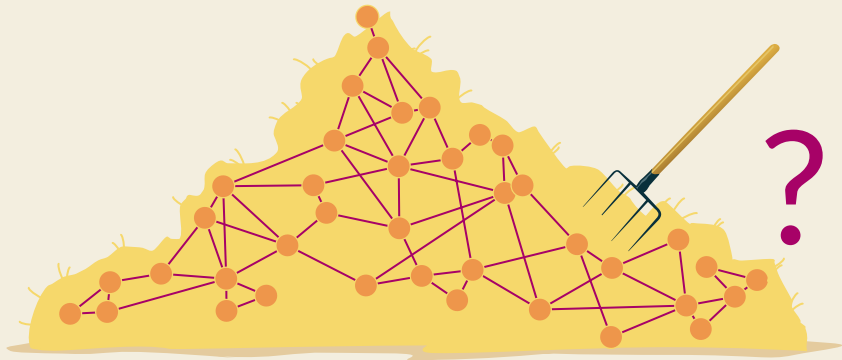
# *Mandoline*



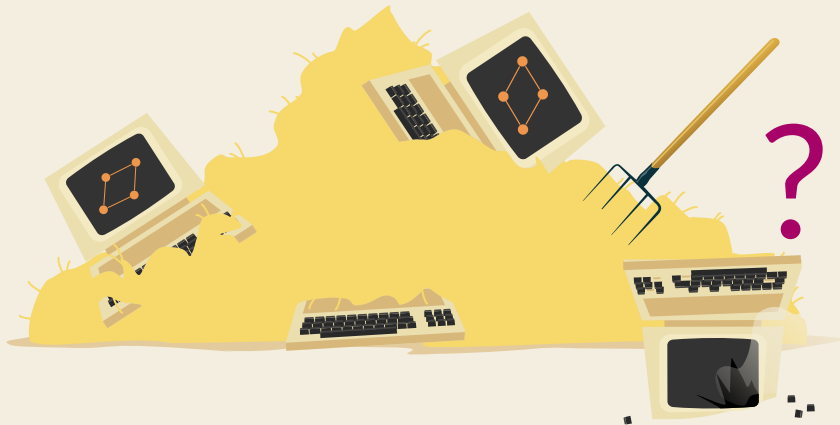
Motif counting using  
generalized colourings

Joint work with Blair D. Sullivan

How many  are in a

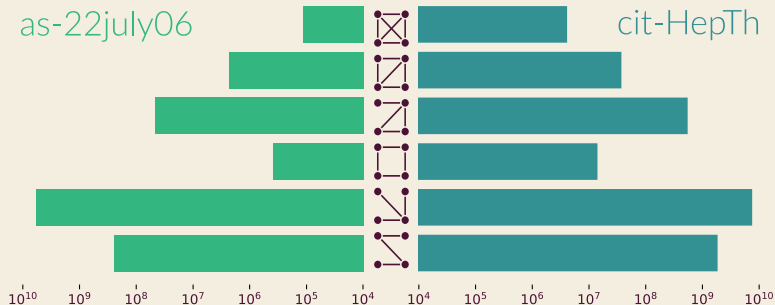


How many  are in a



# Graphlets

We want to count all (connected) induced subgraphs up to a given size.



The **graphlet degree distribution** or the **graphlet degree** can be used to compare networks.

Pržulj N, Corneil DG, Jurisica I.

**Modeling interactome: scale-free or geometric?**

Bioinformatics. 2004 Jul 29;20(18):3508-15

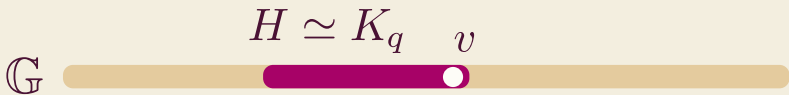
Pržulj N. **Biological network comparison using graphlet degree distribution.**

Bioinformatics. 2007 Jan 15;23(2):e177-83.

# Let's start with something easy!

We count cliques in a  $d$ -degenerate graph.

**Observation:** every clique is contained in the left-neighbourhood of its *last* vertex.

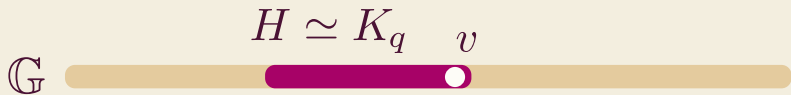


$$V(H) \subseteq N_G^-(v)$$

# Let's start with something easy!

We count cliques in a  $d$ -degenerate graph.

**Observation:** every clique is contained in the left-neighbourhood of its *last* vertex.



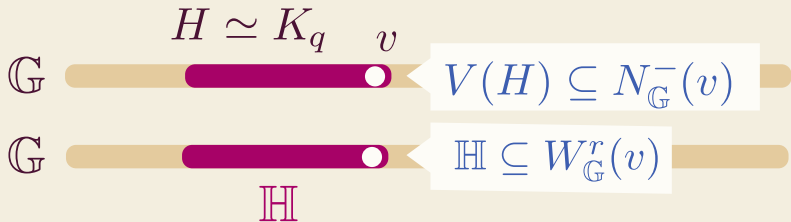
$$V(H) \subseteq N_G^-(v)$$

Therefore we can enumerate all cliques by enumerating all cliques in  $N^-(v)$  for all  $v \in G$ !

$$O(2^d n) \text{ time!}$$

## Does it blend?

Can we 'lift' this algorithm to wcol?

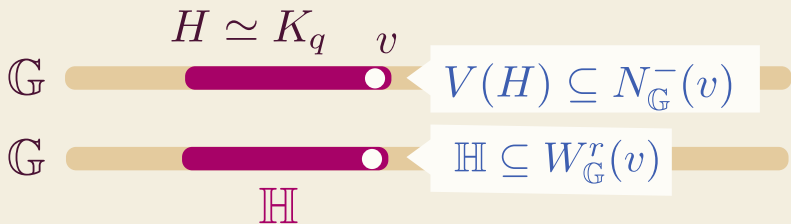


1 What is the 'last' vertex of  $H$ ?

Enumerate all orderings  $\mathbb{H}$  of  $H$ .

# Does it blend?

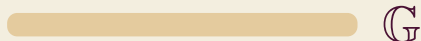
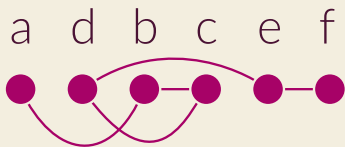
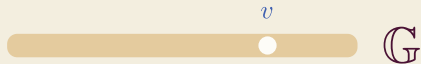
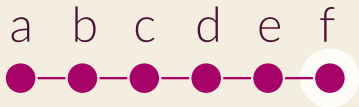
Can we 'lift' this algorithm to wcol?



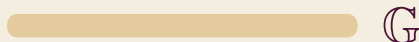
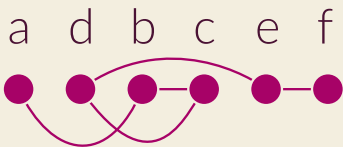
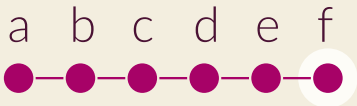
- 1 What is the 'last' vertex of  $H$ ?  
Enumerate all orderings  $\mathbb{H}$  of  $H$ .
- 2 Does  $\mathbb{H} \subseteq W_{\mathbb{G}}^r(v)$  actually hold?  
Only sometimes!



# Two ways to order a $P_6$



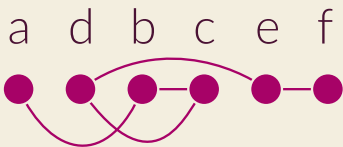
# Two ways to order a $P_6$



# Two ways to order a $P_6$



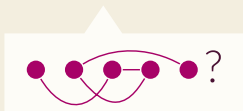
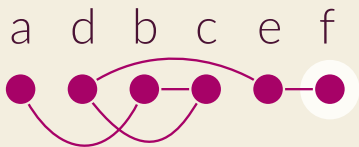
$$W_{P_6}^5(f)!$$



# Two ways to order a $P_6$



$$W_{P_6}^5(f)!$$



# Two ways to order a $P_6$



$$W_{P_6}^5(f)!$$



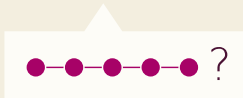
$$W_{P'_6}^5(f) \dots$$



# Two ways to order a $P_6$



$$W_{P_6}^5(f)!$$

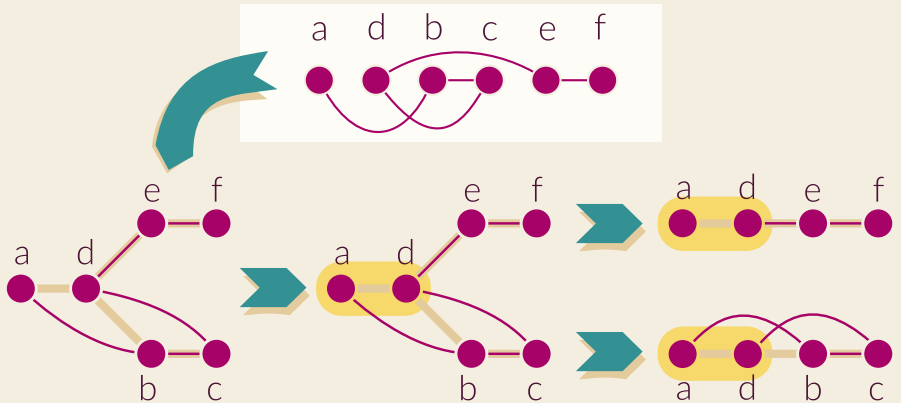


$$W_{P_6}^5(f) \dots$$



Is there a nice formalization of this property?

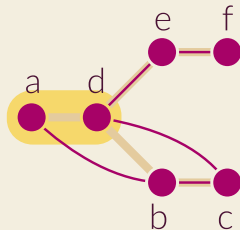
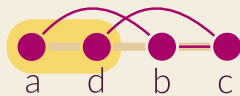
# Decomposition!



We can count linear pieces!

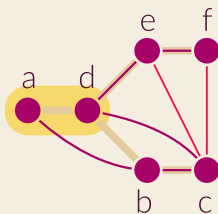
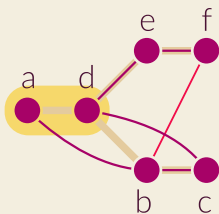
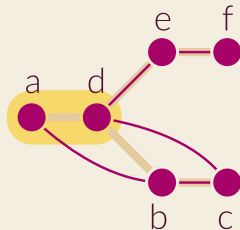
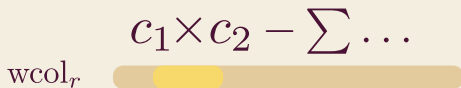
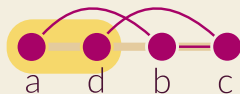
Progress! These pieces are linear!

# Count & combine!



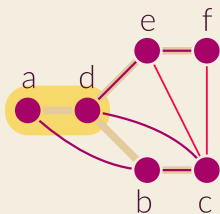
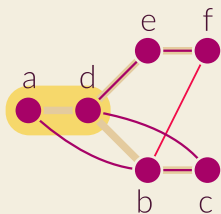
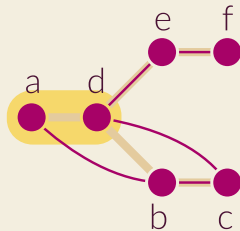
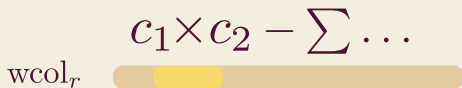
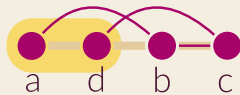


# Count & combine!



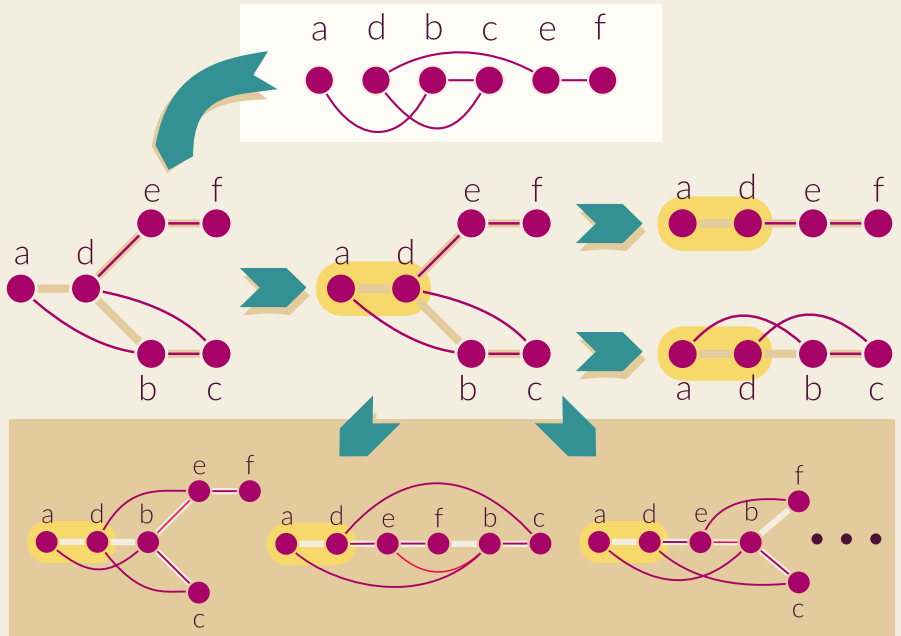
...

# Count & combine!



How do we count these graphs?

# Decomposition!

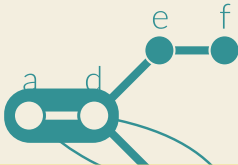
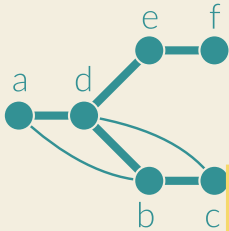


# Decomposition!

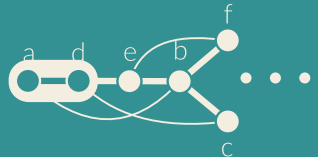
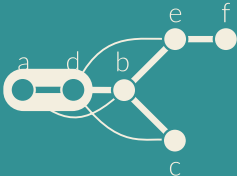
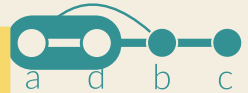
a d b c e f



Less branches



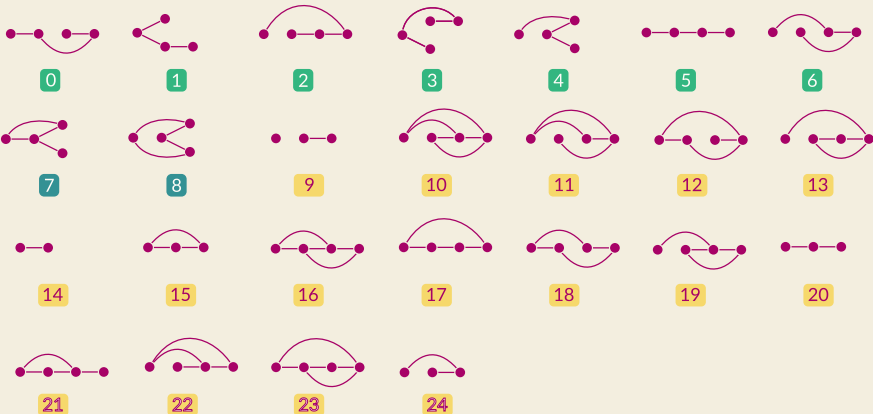
More edges =  
longer decomposition



# Counting $P_4$ s using $wcol_3$

**Lemma 6.** Let  $\mathbf{H} \in \mathcal{H}$  be a (non-linear) pattern relaxation and let  $\mathbf{H}_1 \oplus_{\bar{x}} \mathbf{H}_2 = \mathbf{H}$ . Fix an ordered vertex set  $\bar{y} \in \mathbb{G}$  such that  $\mathbf{H}[\bar{x}] \simeq \mathbb{G}[\bar{y}]$ . Then

$$\#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}, \mathbb{G}) = \#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}_1, \mathbb{G}) \#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}_2, \mathbb{G}) - \sum_{\mathbf{D} \in \mathcal{D}(\mathbf{H}_1, \mathbf{H}_2)} \frac{\#_{\bar{x} \mapsto \bar{x}}(\mathbf{H}, \mathbf{D} \mid \mathbf{H}_1, \mathbf{H}_2) \#_{\bar{x} \mapsto \bar{y}}(\mathbf{D}, \mathbb{G})}{\#_{\bar{x} \mapsto \bar{x}}(\mathbf{D}, \mathbf{D})}.$$

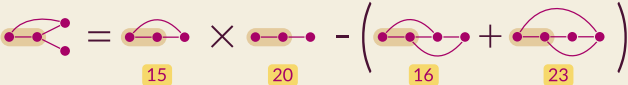


# Counting $P_4$ s using $wcol_3$

**Lemma 6.** Let  $\mathbf{H} \in \mathcal{H}$  be a (non-linear) pattern relaxation and let  $\mathbf{H}_1 \oplus_{\bar{x}} \mathbf{H}_2 = \mathbf{H}$ . Fix an ordered vertex set  $\bar{y} \in \mathbb{G}$  such that  $\mathbf{H}[\bar{x}] \simeq \mathbb{G}[\bar{y}]$ . Then

$$\#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}, \mathbb{G}) = \#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}_1, \mathbb{G}) \#_{\bar{x} \mapsto \bar{y}}(\mathbf{H}_2, \mathbb{G}) - \sum_{\mathbf{D} \in \mathcal{D}(\mathbf{H}_1, \mathbf{H}_2)} \frac{\#_{\bar{x} \mapsto \bar{x}}(\mathbf{H}, \mathbf{D} \mid \mathbf{H}_1, \mathbf{H}_2) \#_{\bar{x} \mapsto \bar{y}}(\mathbf{D}, \mathbb{G})}{\#_{\bar{x} \mapsto \bar{x}}(\mathbf{D}, \mathbf{D})}.$$

4 

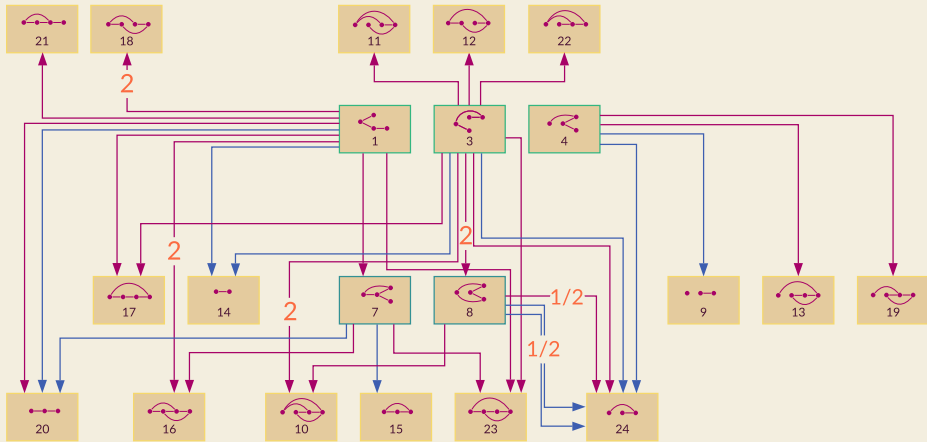
7 

8 

Not shown:

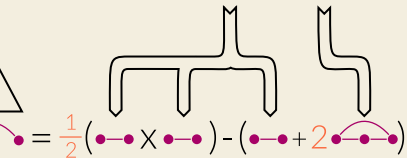
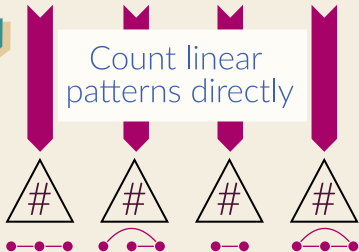
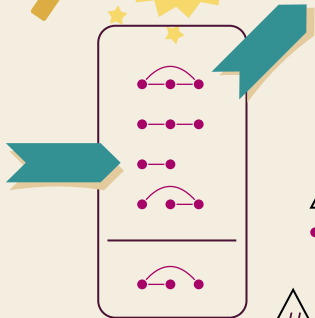
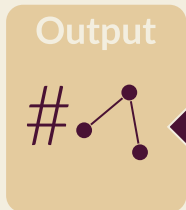
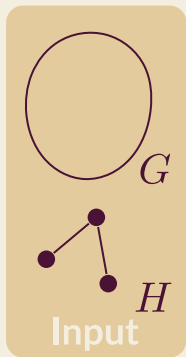


# Counting $P_4$ s using $wcol_3$



This 'counting-DAG' has bounded depth

# Subgraph counting using $wcol_r$



Aggregate

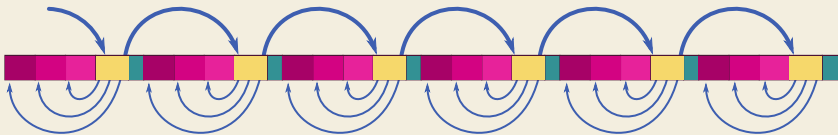
Compute composite pattern counts



# Engineering: memory

**Lesson:** We cannot ignore memory locality.  
But this is really hard when working with graphs.

## 1 Flatten everything



# Engineering: memory

**Lesson:** We cannot ignore memory locality.  
But this is really hard when working with graphs.

2 Be as specific as possible

Vertex ids

$n$

(int16)

$\approx 60K$

int32

$\approx 4.3Mio$

int64

All we could  
ever need

Node ids

$k$

3 bit

4 bit

char

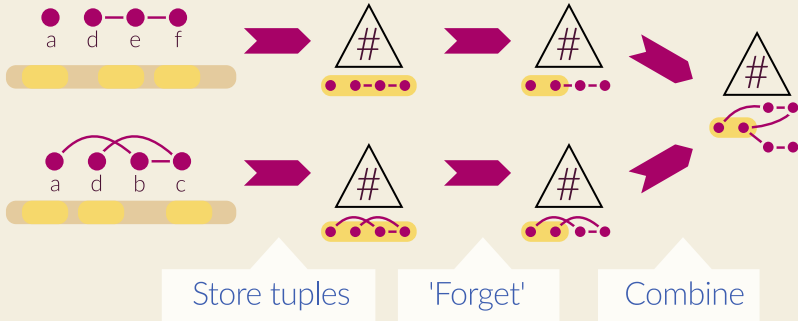
All we could  
ever need

Tough design  
choices

# Engineering: data structures

**Lesson:** Common data structures not enough

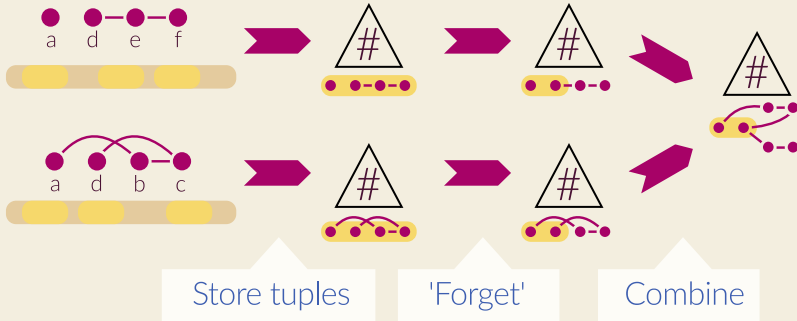
## 3 Design special-purpose data structures



# Engineering: data structures

**Lesson:** Common data structures not enough

## 3 Design special-purpose data structures



Our solution: specialized trie. Better options?

# Next steps

## 1 More features

Count coloured graphs, sum-of-weights, other types of embeddings (homomorphisms, non-induced, etc.)

## 2 Implement for $\text{col}_r(\mathbb{G})$

We know how to do this in theory and  $\text{col}_r(\mathbb{G})$  seems to be smaller in practice.

## ? Preprocessing!

Simple preprocessing rules (based e.g. min-degree of pattern) help a lot. Can we generate more elaborate, pattern-dependent rules?



# Solvers vs Solutions



Mandoline

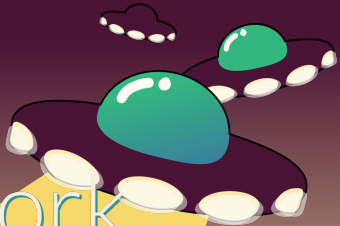


CATLAS

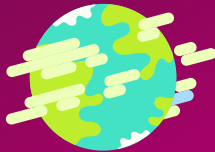
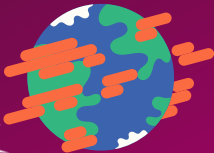
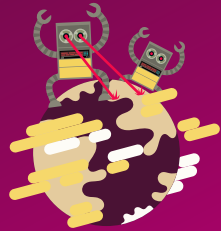
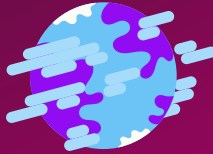


Part III

Future work



# The big open question

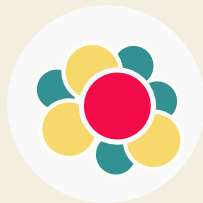




# THANKS



# Exhibit X



Local centrality measures  
via neighbourhood surveys  
(Secret slides)

# Close-to-Closeness Centralities

 $C(v)$ 

r-Local version

Closeness

$$\left( \sum_{u \in G} \text{dist}(u, v) \right)^{-1} \quad \left( \sum_{u \in N^r[v]} \text{dist}(v, u) \right)^{-1}$$

Harmonic

$$\sum_{u \in G} \text{dist}(u, v)^{-1} \quad \sum_{u \in N^r[v]} \text{dist}(v, u)^{-1}$$

Lin's index

$$\frac{|\{u \mid \text{dist}(u, v) < \infty\}|^2}{\sum_{\text{dist}(u, v) < \infty} \text{dist}(u, v)} \quad \frac{|N^r[v]|^2}{\sum_{u \in N^r[v]} \text{dist}(v, u)}$$

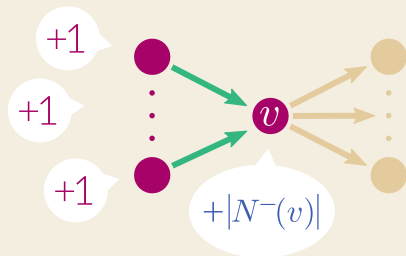
All three measures can be computed quickly if we know  $|N^d(v)|$  for  $1 \leq d \leq r$ .

Can we compute this quickly in sparse graphs?

# Warm-up: Counting with degeneracy

Let  $G$  be  $(d-1)$ -degenerate.

- 1 Compute orientation  $\vec{G}$  with  $\Delta^-(\vec{G}) \leq d$  in linear time.
- 2 Initialize counter  $C[v] = 0$  for all  $v \in G$ .
- 3 For every  $v \in G$ , increment  $C[v]$  and  $C[u]$  for every in-neighbour  $u \in N^-(v)$ .



# Degeneracy to dtf-augmentations

**Thm.** Given a graph  $G$  and an integer  $r$ , we can compute the size of  $|N^d(v)|$  for all  $v \in G$  and  $1 \leq d \leq r$  in total time  $O(2^{\Delta^-(\vec{G}_r)} n)$ .



# Counting using dtf-augmentations

We compute the size of the  $r^{\text{th}}$  nbhds:

- 1 Compute dtf-augm.  $\vec{G}_r$  with small  $\Delta^-(\vec{G}_r)$  in linear time.

# Counting using dtf-augmentations

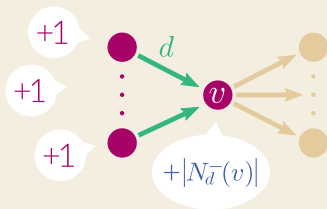
We compute the size of the  $r^{\text{th}}$  nbhds:

- 1 Compute dtf-augm.  $\vec{G}_r$  with small  $\Delta^-(\vec{G}_r)$  in linear time.
- 2 Initialize counter  $C[v][d] = 0$  for all  $v \in G$  and  $d \leq r$ .

# Counting using dtf-augmentations

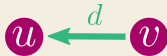
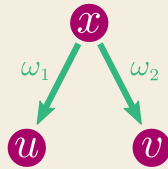
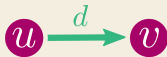
We compute the size of the  $r^{\text{th}}$  nbhds:

- 1 Compute dtf-augm.  $\vec{G}_r$  with small  $\Delta^-(\vec{G}_r)$  in linear time.
- 2 Initialize counter  $C[v][d] = 0$  for all  $v \in G$  and  $d \leq r$ .
- 3 For every  $v \in G$ , increment  $C[v][d]$  and  $C[u][d]$  for every in-neighbour  $u \in N_d^-(v)$ .





# Counting using dtf-augmentations



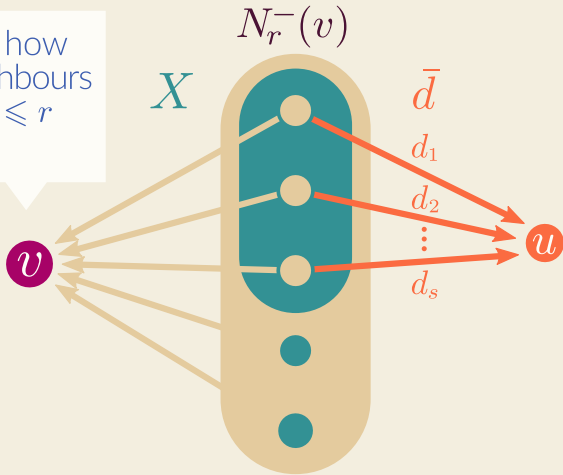
The counting so far takes care of the first two cases, but what about the *indirect* neighbours?

This is where the algorithm becomes **interesting**.

# Counting using dtf-augmentations

$v$  needs to know how many indirect neighbours at distance  $2 \leq d \leq r$  there are.

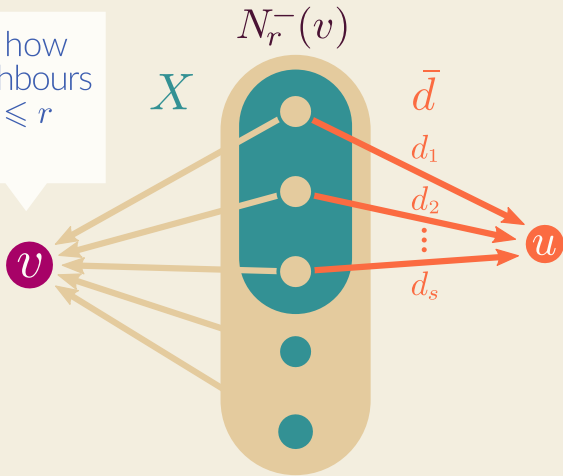
Indirect neighbours connect via  $N_r^-(v)$ .



# Counting using dtf-augmentations

$v$  needs to know how many indirect neighbours at distance  $2 \leq d \leq r$  there are.

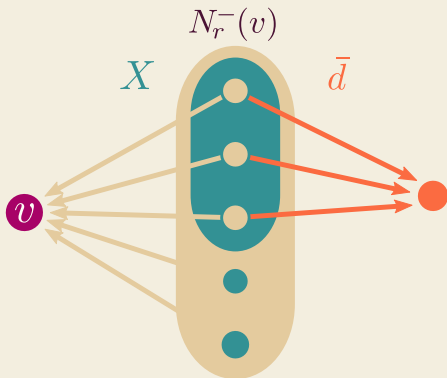
Indirect neighbours connect via  $N_r^-(v)$ .



We compute the distance between  $v, u$  as follows:

$$\text{dist}(u, v) = \min(\text{dist}(v, X) + \text{dist}(u, X))$$

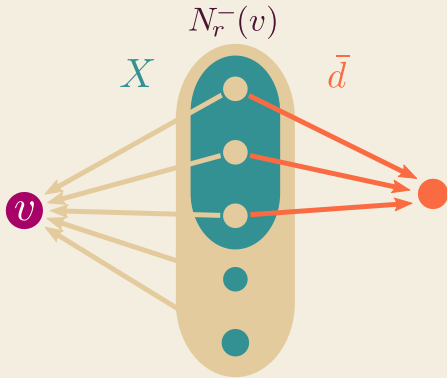
# Counting using dtf-augmentations



We need to compute for every set  $X \subseteq N_r^-(v)$  and every possible dist.-vector  $\bar{d} \in [r]^{|X|}$  the number of vertices  $u$  such that:

- 1  $N_r^-(u) \cap N_r^-(v) = X$
- 2  $\text{dist}(u, X) = \bar{d}$

# Counting using dtf-augmentations



We need to compute for every set  $X \subseteq N_r^-(v)$  and every possible dist.-vector  $\bar{d} \in [r]^{|X|}$  the number of vertices  $u$  such that:

- 1  $N_r^-(u) \cap N_r^-(v) = X$
- 2  $\text{dist}(u, X) = \bar{d}$

Let us call this number  $c(v, X, \bar{d})$ . Our first goal is to compute it for every vertex.

## A data structure for $c(v, X, \bar{d})$

- 1 For every  $v \in \vec{G}_r$ ,  $X \subseteq N_r^-(v)$  and  $\bar{d} \in [r]^{|X|}$ , initialize  $R[X][\bar{d}] = 0$ .

## A data structure for $c(v, X, \bar{d})$

- 1 For every  $v \in \vec{G}_r$ ,  $X \subseteq N_r^-(v)$  and  $\bar{d} \in [r]^{|X|}$ , initialize  $R[X][\bar{d}] = 0$ .
- 2 For every  $v \in \vec{G}_r$ ,  $X \subseteq N_r^-(v)$ , increment  $R[X][\text{dist}(v, X)]$  by one.

## A data structure for $c(v, X, \bar{d})$

- 1 For every  $v \in \vec{G}_r$ ,  $X \subseteq N_r^-(v)$  and  $\bar{d} \in [r]^{|X|}$ , initialize  $R[X][\bar{d}] = 0$ .
- 2 For every  $v \in \vec{G}_r$ ,  $X \subseteq N_r^-(v)$ , increment  $R[X][\text{dist}(v, X)]$  by one.

**Claim.**

$$c(v, X, \bar{d}) = \sum_{X \subseteq Y \subseteq N_r^-(v)} (-1)^{|Y \setminus X|} \sum_{\bar{d}' : \bar{d}'|_X = \bar{d}} R[Y][\bar{d}'].$$



# Counting using dtf-augmentations

Given  $c(v, \bullet, \bullet)$  we can now count the number of indirect neighbours of  $v$ . For every subset  $X \subseteq N_r^-(v)$  and distance-vector  $\bar{d} \in [r]^{|X|}$ , apply the update:

$$C[v][\min(\bar{d} + \text{dist}(v, X))] += c(v, X, \bar{d})$$

Since the above counts  $v$  as a neighbour of itself, we apply the following correction:

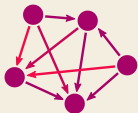
$$C[v][\min(\text{dist}(v, X) + \text{dist}(v, X))] -= 1$$

There are a few more corrections concerning direct neighbours, see paper.

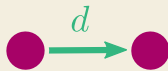
# Counting using dtf-augmentations



$G, r$   
Input



$\vec{G}_r$



Populate  $C[\bullet][\bullet]$   
for direct neighbours

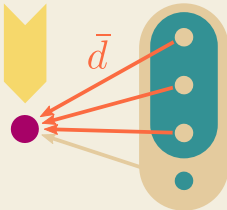
$C$   
 $|N^d(v)|$   
 $= C[v][d]$

Output

$c(\bullet, \bullet, \bullet)$

$R$

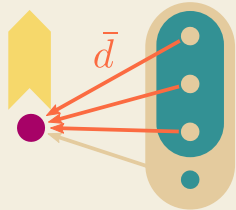
$\bar{d}$



Update  $C[\bullet][\bullet]$   
using  $R$  to count  
indirect neighbours

$R$

$\bar{d}$

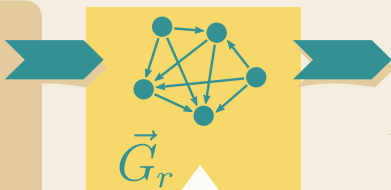


Populate  
 $R[\bullet][\bullet]$

# Counting using dtf-augmentations



$G, r$   
Input



$\vec{G}_r$

$O(\Delta^-(\vec{G}_r)^2 n)$



Populate  $C[\bullet][\bullet]$   
for direct neighbours



$C$   
 $|N^d(v)|$   
 $= C[v][d]$

Output

$c(\bullet, \bullet, \bullet)$

$\bar{d}$

Update  $C[\bullet][\bullet]$   
using  $R$  to count  
indirect neighbours

$R$

$\bar{d}$

Populate  
 $R[\bullet][\bullet]$

# Counting using dtf-augmentations



$G, r$   
Input



$\vec{G}_r$



Populate  $C[\bullet][\bullet]$   
for direct neighbours

$O(\Delta^-(\vec{G}_r)n)$

$R$

$c(\bullet, \bullet, \bullet)$

$\bar{d}$

$d$

$C$   
 $|N^d(v)|$   
 $= C[v][d]$

Output

Update  $C[\bullet][\bullet]$   
using  $R$  to count  
indirect neighbours

Populate  
 $R[\bullet][\bullet]$

# Counting using dtf-augmentations



$G, r$   
Input



$\vec{G}_r$

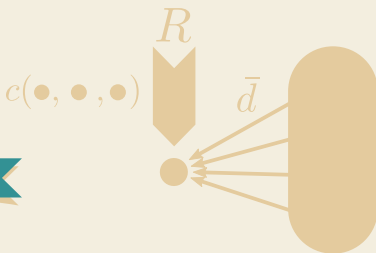


$$O(2^{\Delta^-(\vec{G}_r)} n)$$

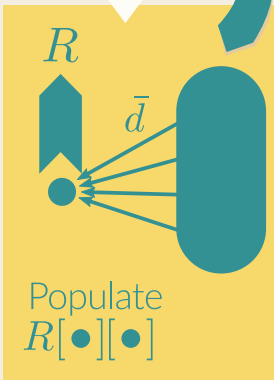
$C$

$$|N^d(v)| = C[v][d]$$

Output



Update  $C[\bullet][\bullet]$   
using  $R$  to count  
indirect neighbours



# Counting using dtf-augmentations



$G, r$   
Input



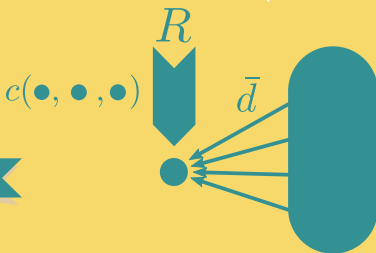
$\bar{G}$   $O(2^{\Delta^-}(\bar{G}_r)n)$



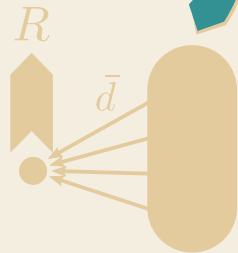
Populate  $C[\bullet][\bullet]$   
direct neighbours

$C$   
 $|N^d(v)|$   
 $= C[v][d]$

Output



Update  $C[\bullet][\bullet]$   
using  $R$  to count  
indirect neighbours

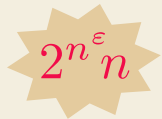


Populate  
 $R[\bullet][\bullet]$

# Counting using dtf-augmentations

**Thm.** Given a graph  $G$  and an integer  $r$ , we can compute the size of  $|N^d(v)|$  for all  $v \in G$  and  $1 \leq d \leq r$  in total time  $O(2^{\Delta^-(\vec{G}_r)} n)$ .

- Exponential vs quadratic?
- Does not scale to nowhere dense graphs!



Can we do *better*?

# Can we do better?

## CLOSED 2-NEIGHBOURHOOD SIZES

**Input:** A graph  $G$ .

**Output:**  $|N^2[v]|$  for every  $v \in G$ .

**Thm.** Unless SETH fails, 2-CNBS cannot be solved in time

①  $O(|G|^{2-\varepsilon})$

②  $O(2^{o(\Delta^-(\vec{G}_2))} n^{2-\varepsilon})$

Gutin G, Mertzios GB, Reidl F.

**Lower and Upper Bound for Computing the Size of All Second Neighbourhoods.**

arXiv preprint arXiv:1805.01684. 2018 May 4



# Engineering: compromises

**Lesson:** We cannot be too idealistic.  
Some vertices will have large in-neighbourhoods.

## 1 Mix algorithms

Conduct regular bfs for vertices with large in-neighbourhood, use dtf-magic for everything else.

Luckily, these two approaches mix!

# Next steps

## 1 More features

We can easily count the weights of neighbourhoods or neighbourhoods restricted to certain colours.

## 2 Optimize for small $r$

Most real-world cases will be for small distances. In particular for  $r=2$  we can simplify the algorithm.

## ? Nowhere dense?

Combine with neighbourhood complexity, handle large intersections differently.

## ? Approximate?

The Möbius inversion is the bottleneck. Sacrifice precision for speed?

