

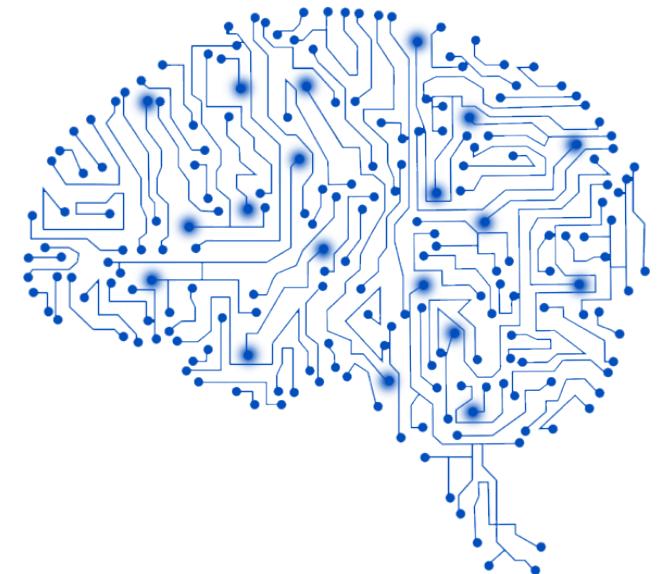
Master's Thesis

Motif Discovery for Connectomics: Finding Computational Units in the Brain's Wiring Diagram

June 26th, 2019

Presenter : Alleon Antoine, CSE – Master 2

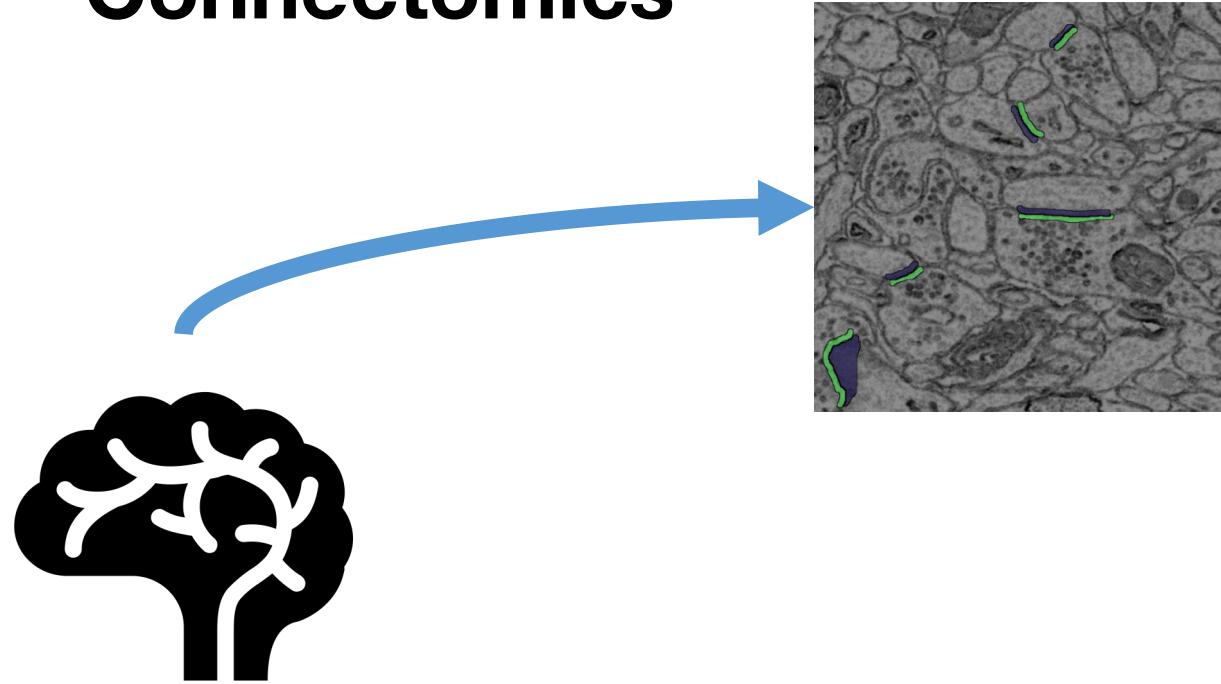
École Polytechnique Fédérale de Lausanne – EPFL
Visual Computing Group, Harvard SEAS



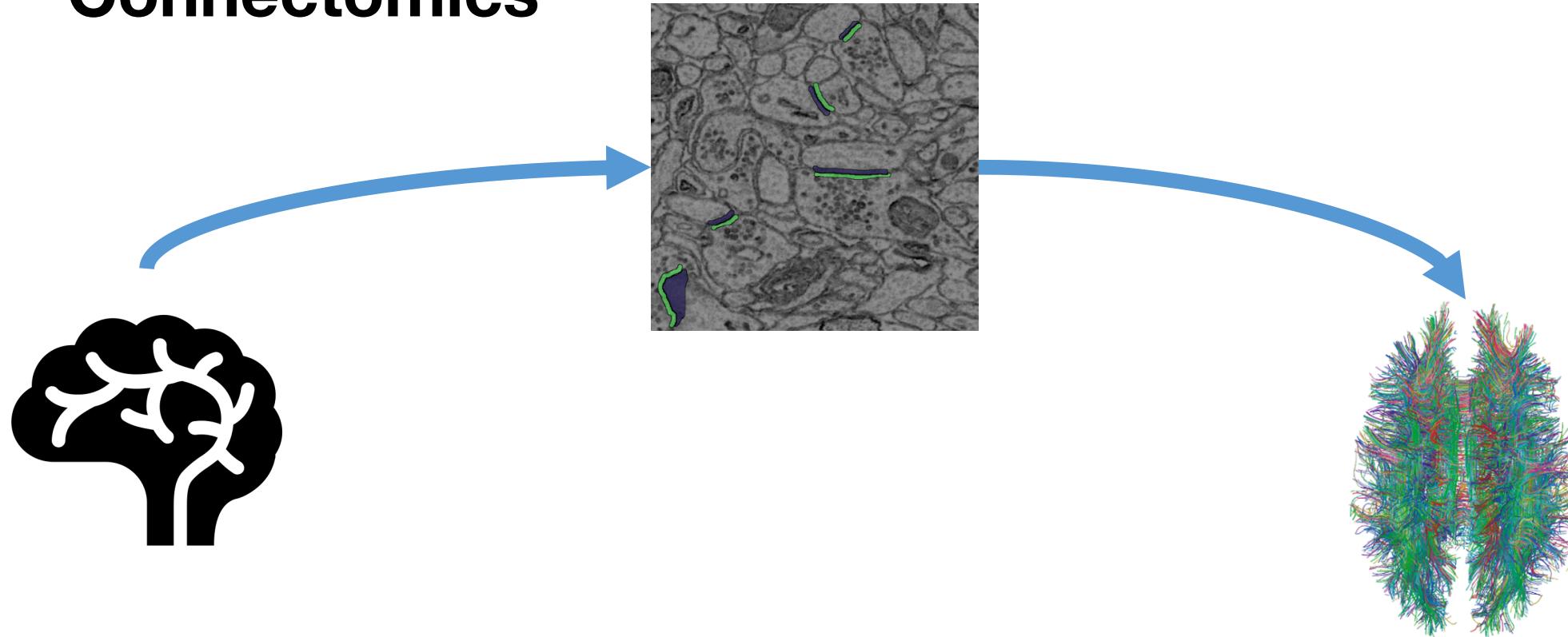
Connectomics



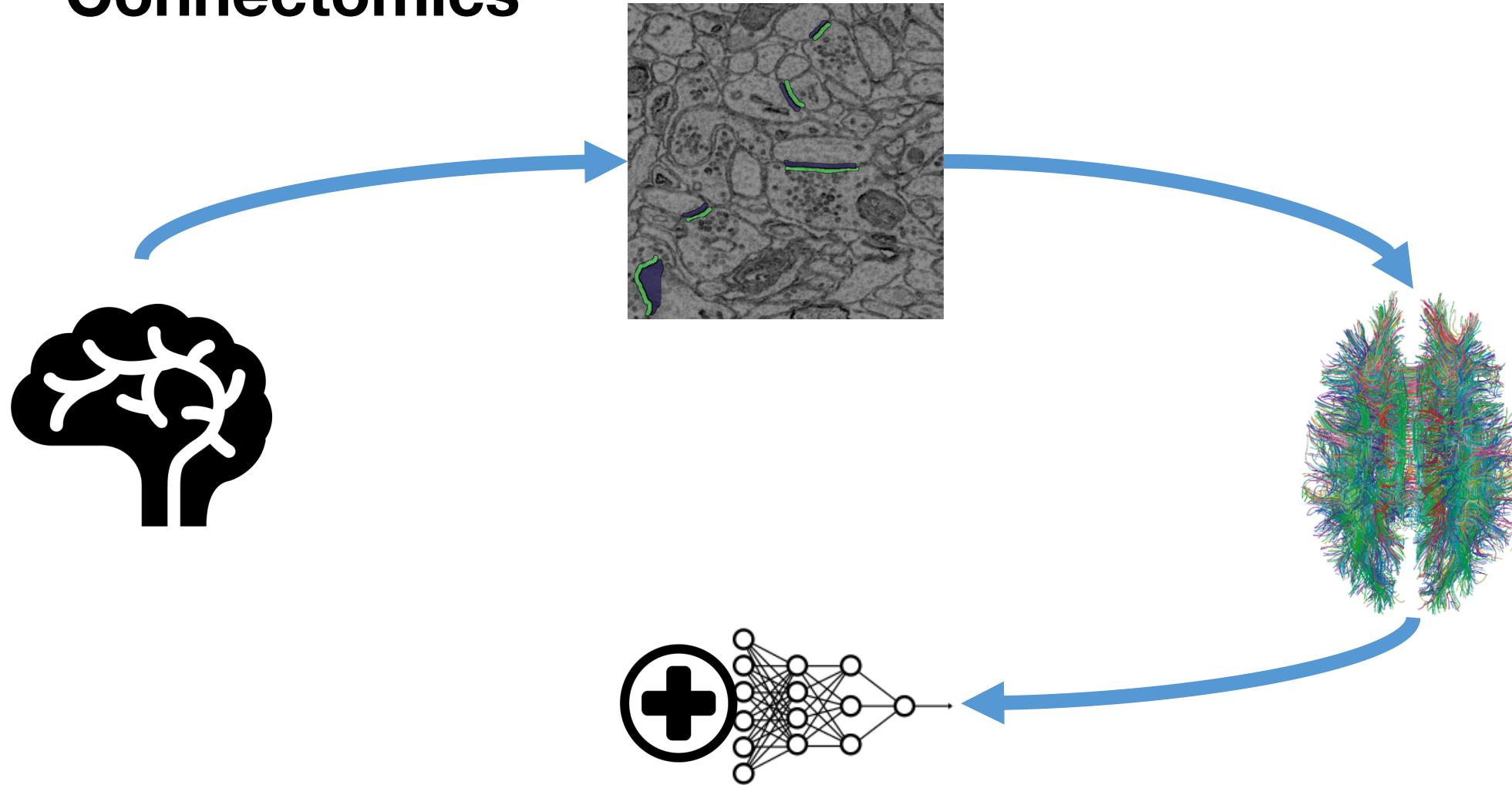
Connectomics



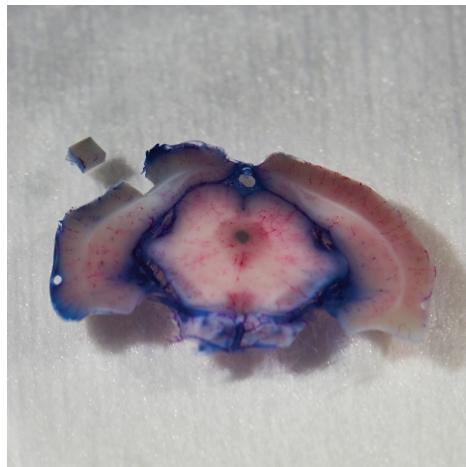
Connectomics



Connectomics



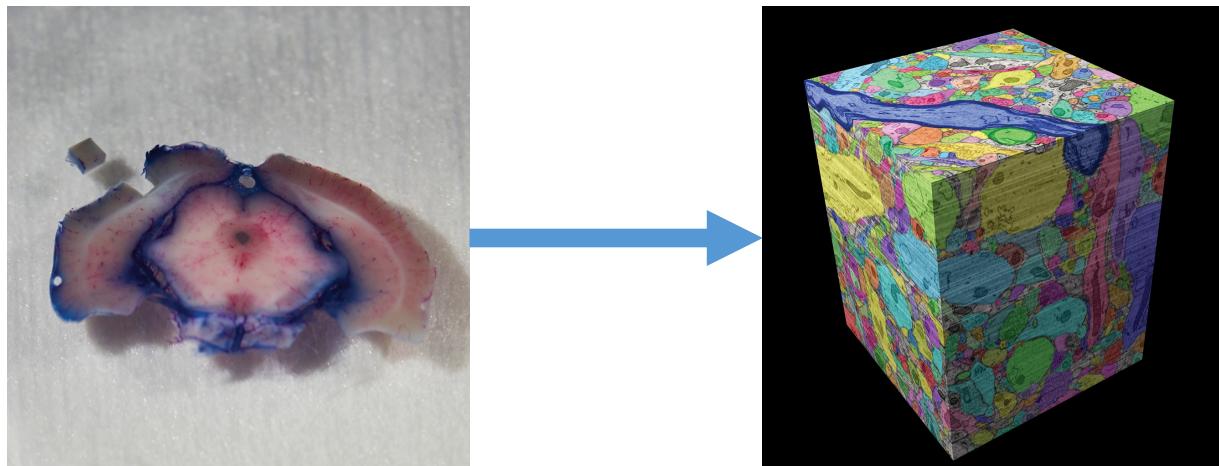
Connectomics



Bio-engineering:

- Brain cutting
- Brain imaging

Connectomics



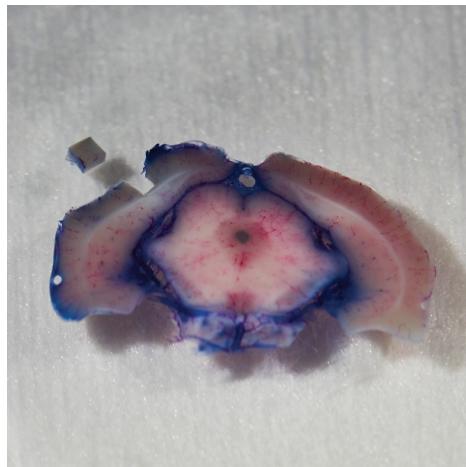
Bio-engineering:

- Brain cutting
- Brain imaging

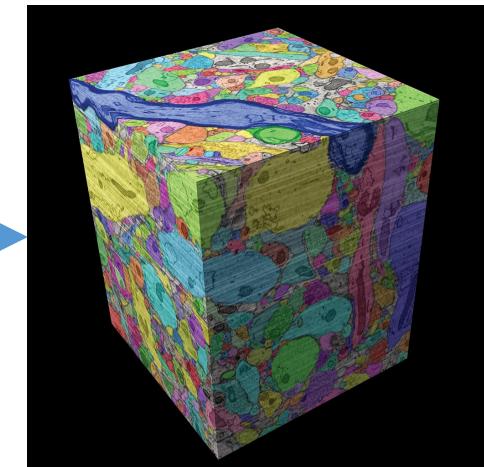
Computer Sciences:

- Alignment of images
- Segmentation

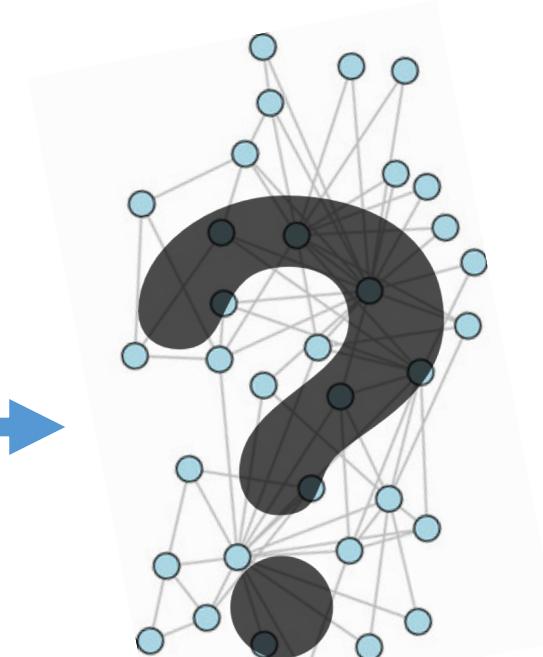
Connectomics



- Bio-engineering:
- Brain cutting
 - Brain imaging



- Computer Sciences:
- Alignment of images
 - Segmentation

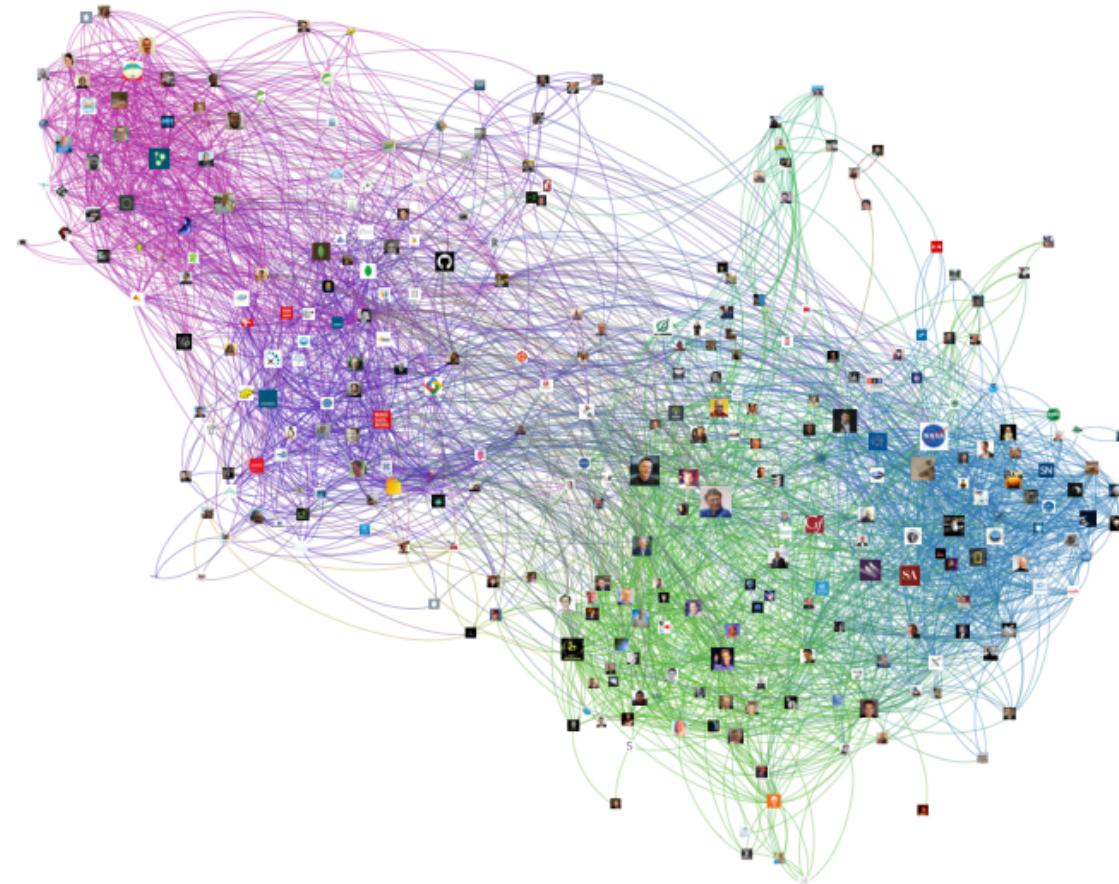
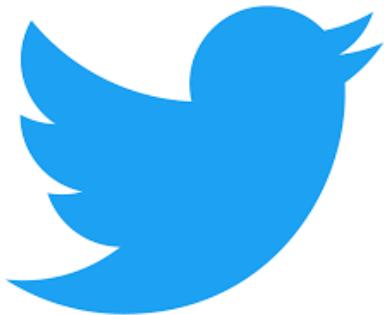


- Applications:
- Graph analysis
 - Visualization

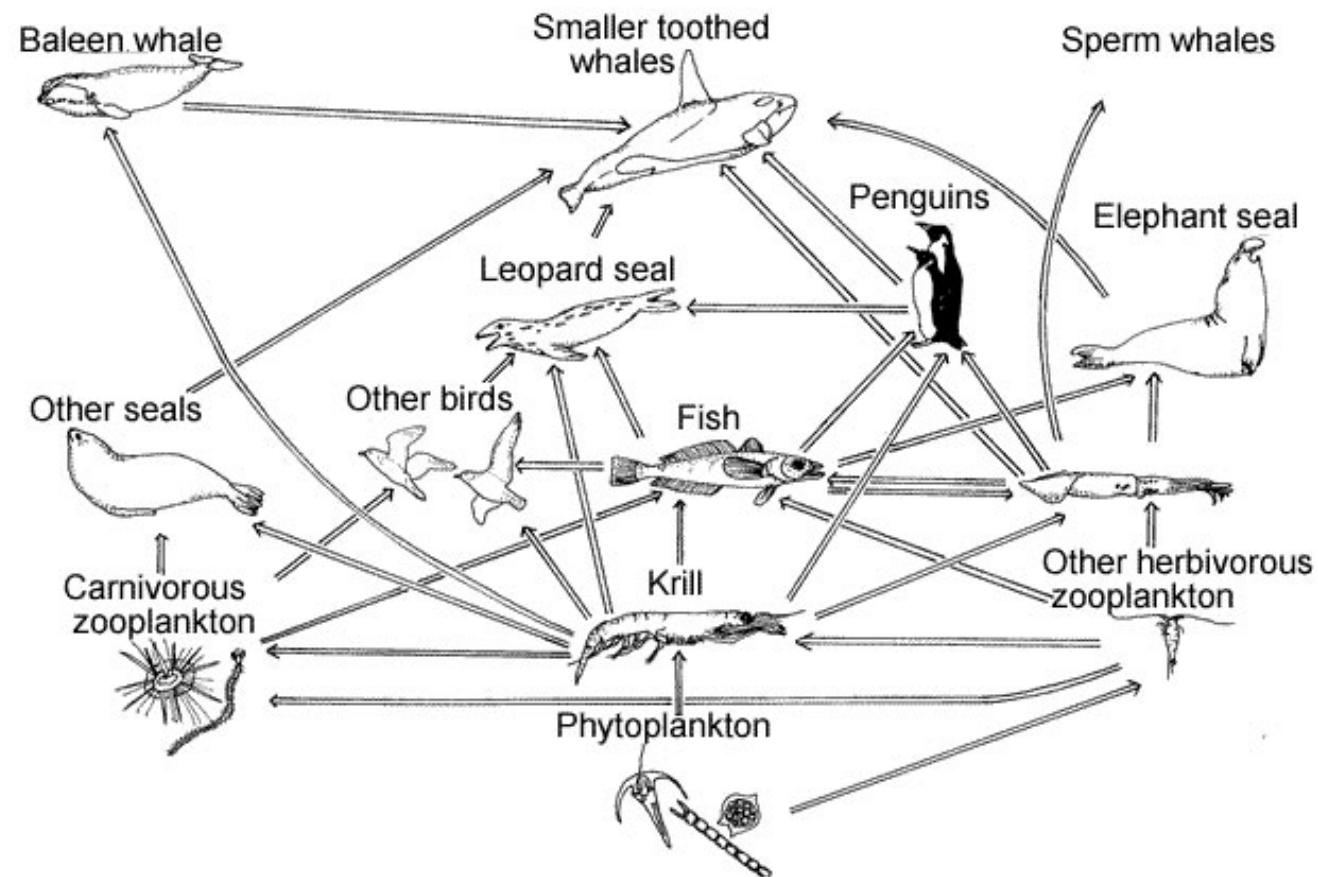
Motif Discovery



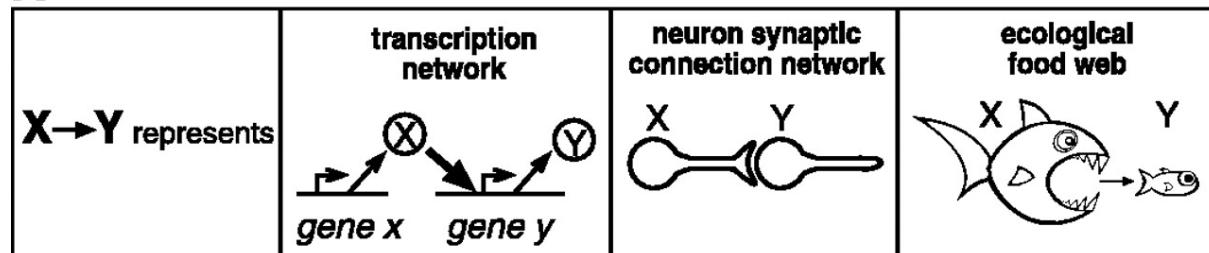
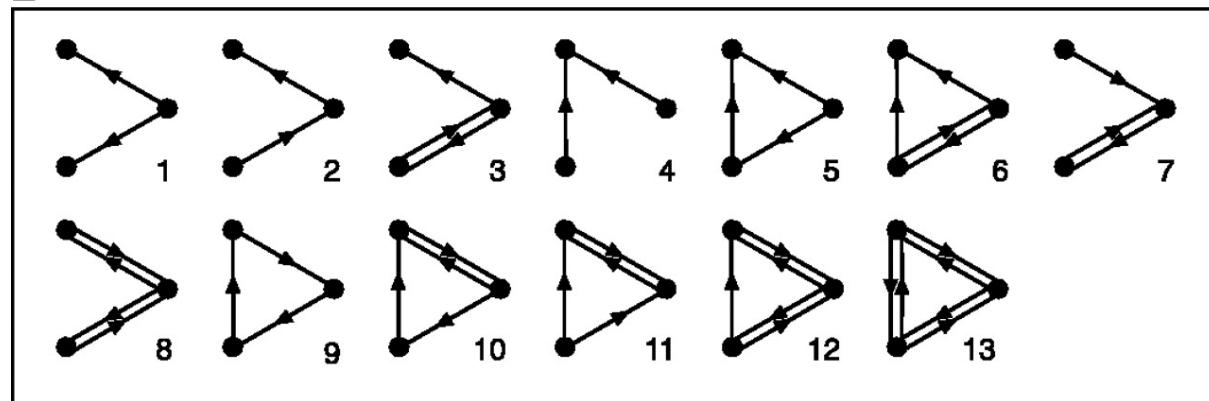
Motif Discovery



Motif Discovery



Motif Discovery

A**B**

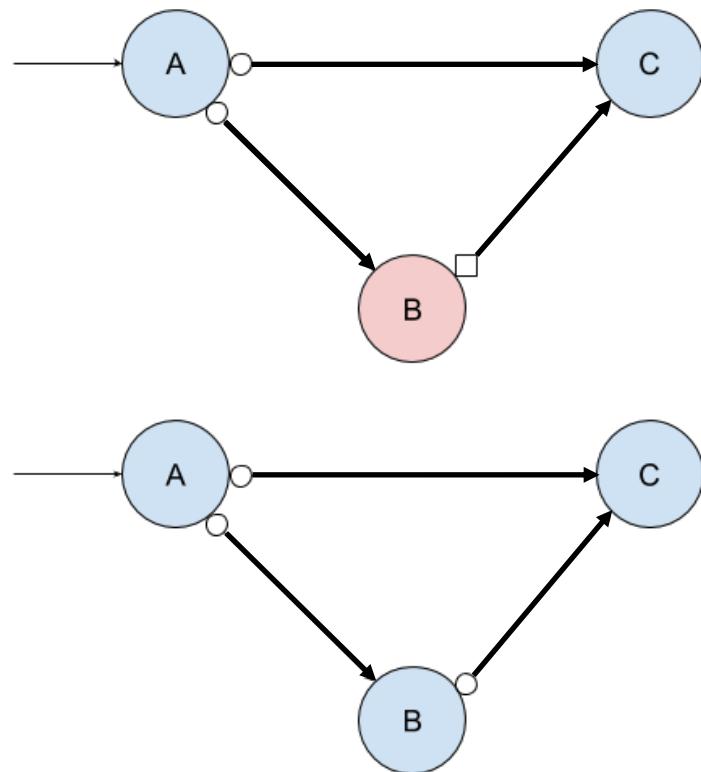
Gene regulation

Metabolic system

Voter networks

Neuronal connections

Motif Discovery

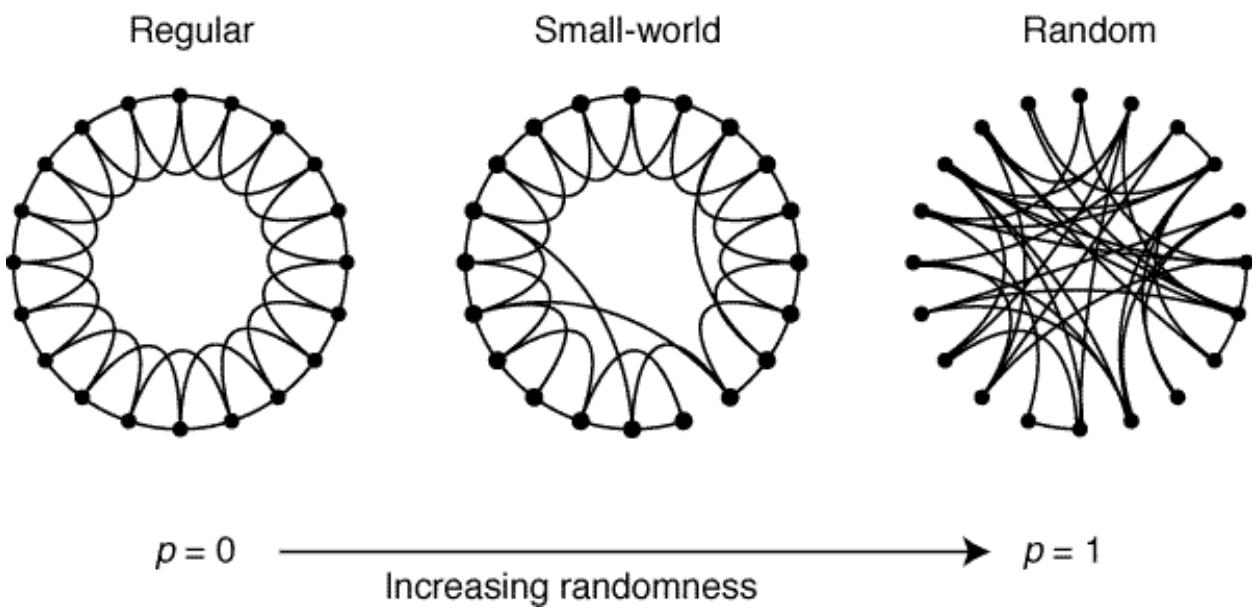


Feed-Forward loop motif

Two types of neurons

- (Light Blue Circle) Excitatory neuron
- (Pink Circle) Inhibitory neuron

Motif Discovery



Small-world properties of neural networks

Highly clustered

Small path length

Agenda

Survey & Application: Small motifs (3 to 5 nodes)

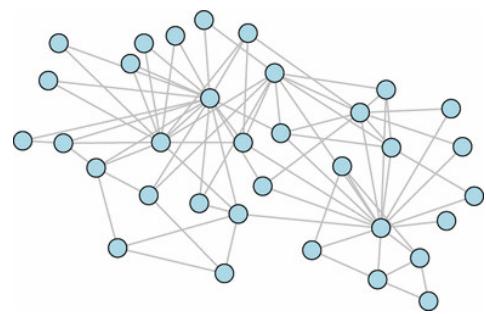
- Brain wiring diagram
- Computational performance

Exploration: Large motifs (6 and more nodes)

- Improve efficiency on large motif detection

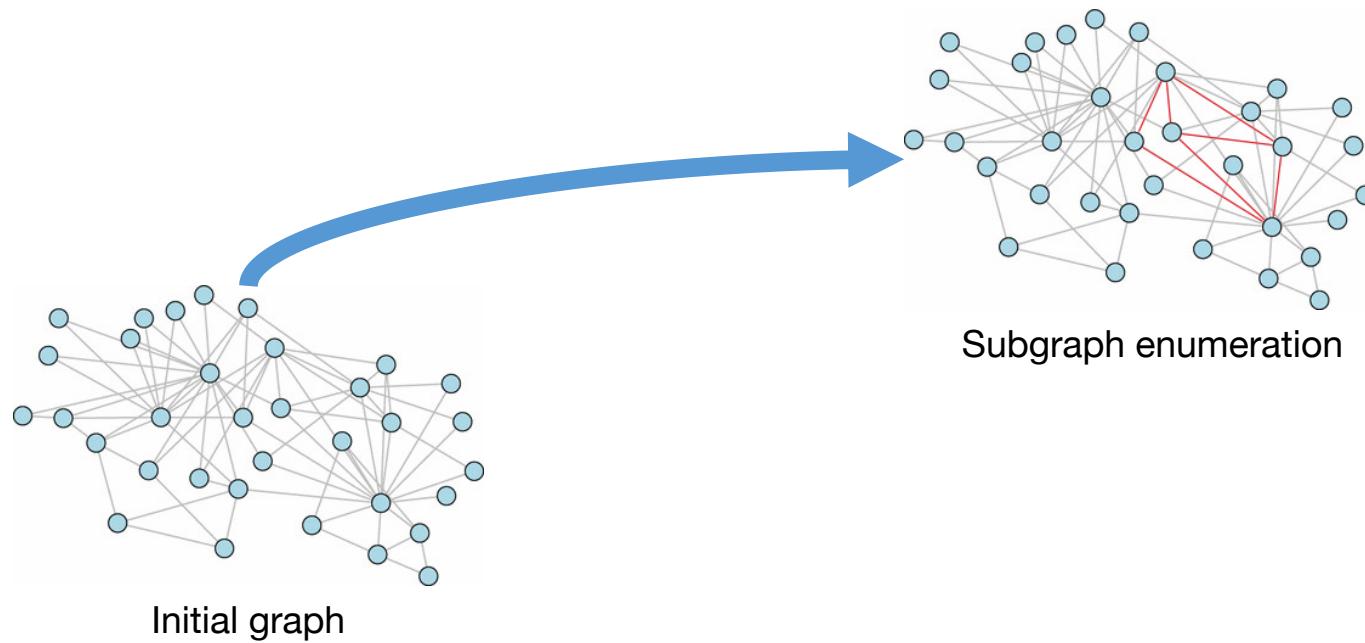
Small Subgraph enumeration: Survey and applications

Motif discovery - Intuition

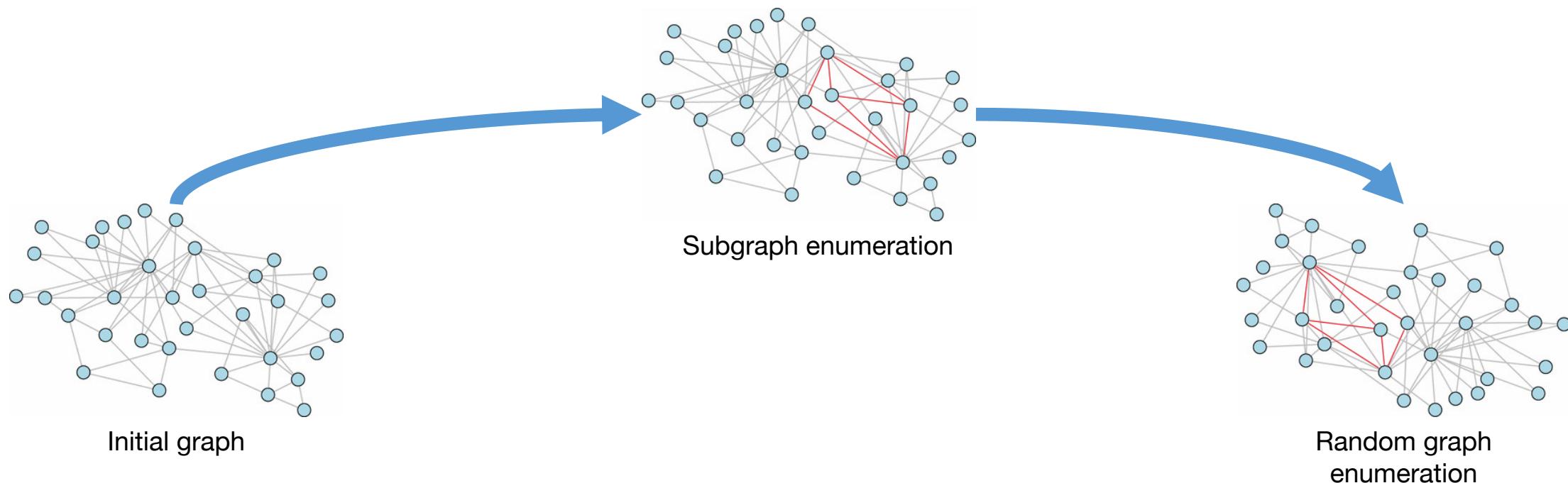


Initial graph

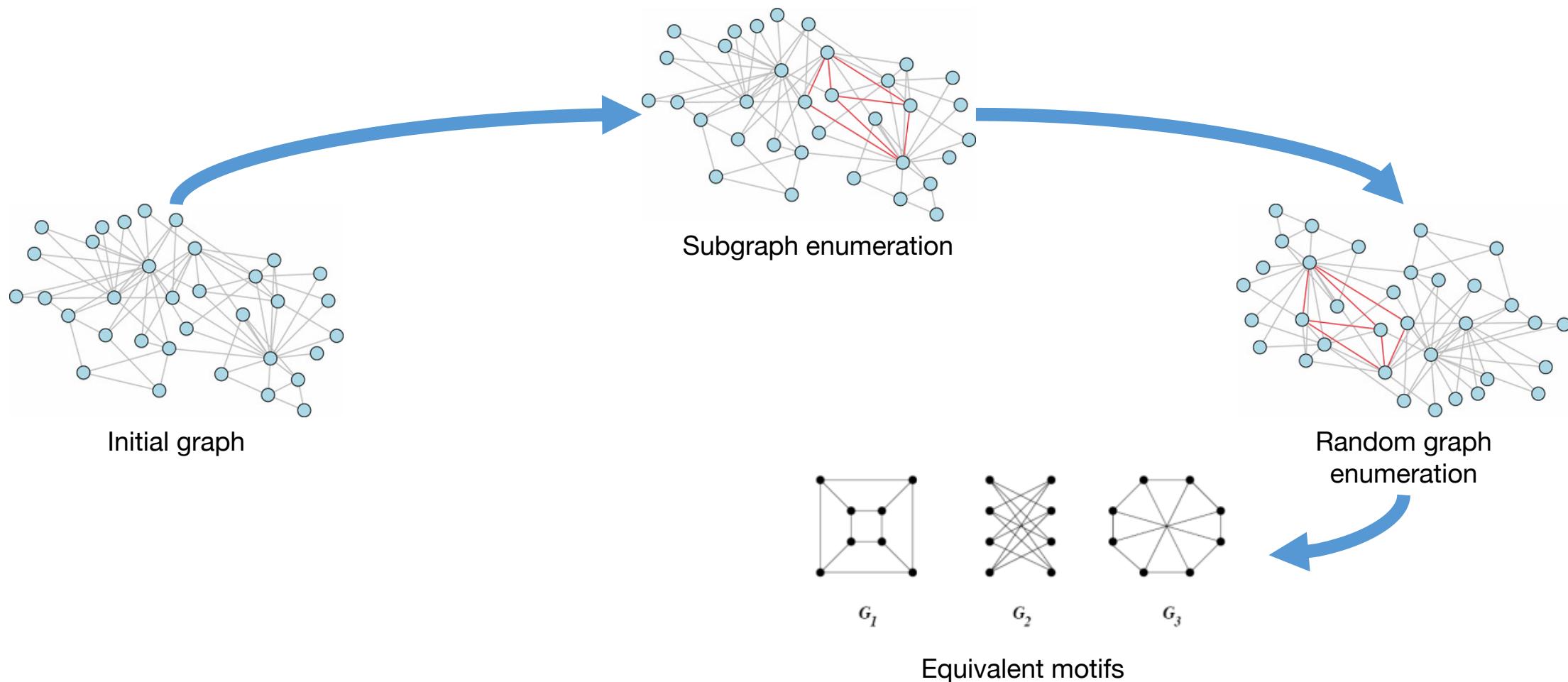
Motif discovery - Intuition



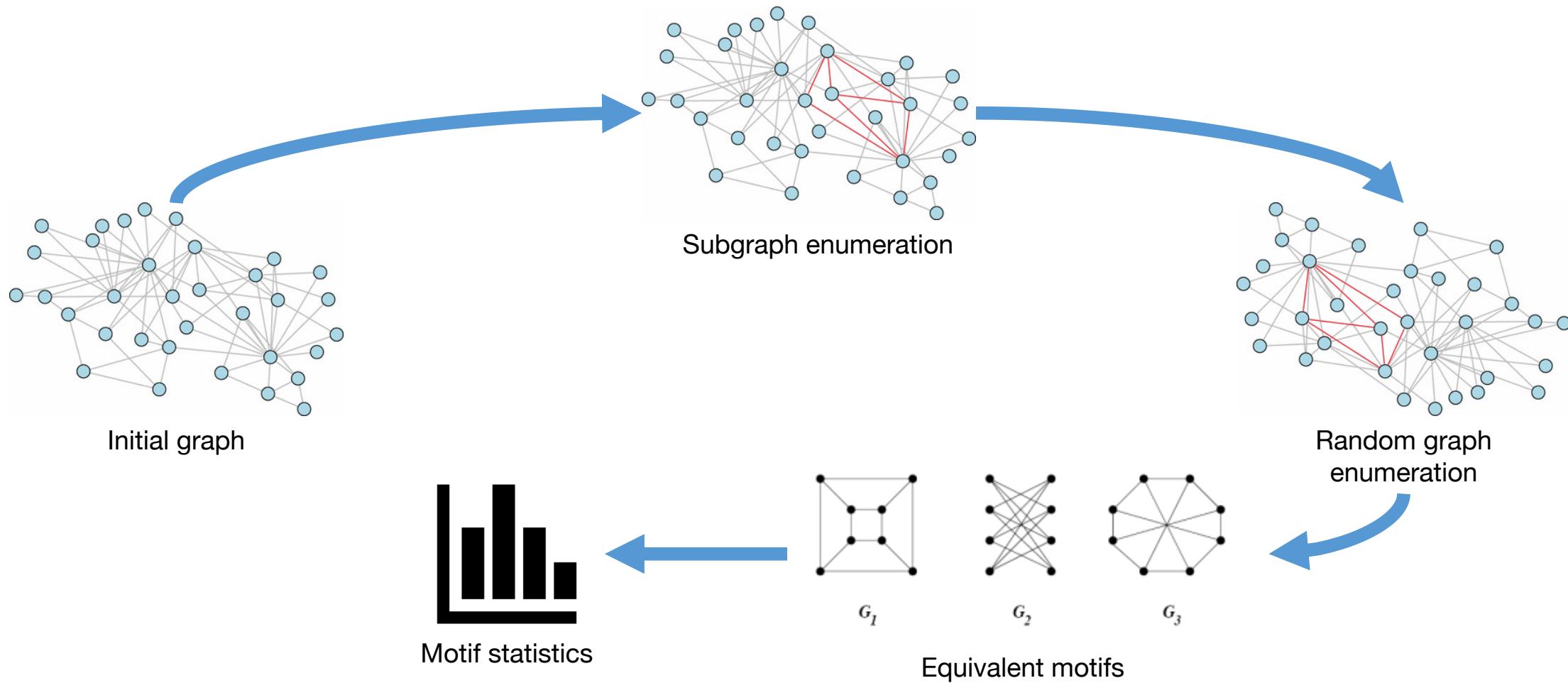
Motif discovery - Intuition



Motif discovery - Intuition



Motif discovery - Intuition



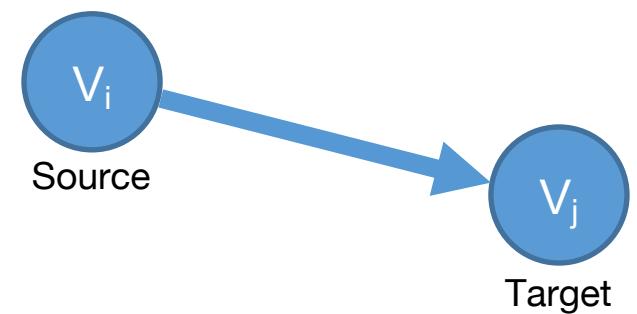
Methods – Graph Notation

Graph - $G = \{V, E\}$

Methods – Graph Notation

Graph - $G = \{V, E\}$

Edge - $E_{ij} = (V_i, V_j)$ with $V_i, V_j \in V, i \neq j$

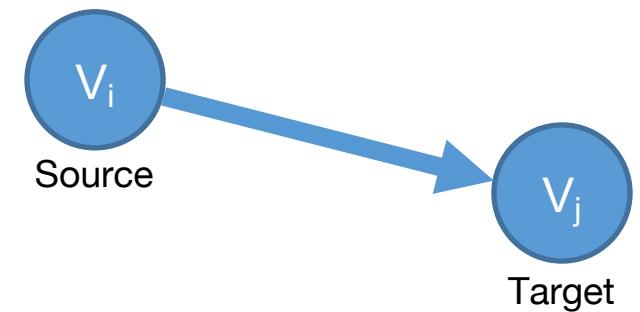


Methods – Graph Notation

Graph - $G = \{V, E\}$

Edge - $E_{ij} = (V_i, V_j)$ with $V_i, V_j \in V, i \neq j$

Subgraph - $G' = \{V', E'\}$ with $V' \subseteq V; E' \subseteq (V', V') \cap E$

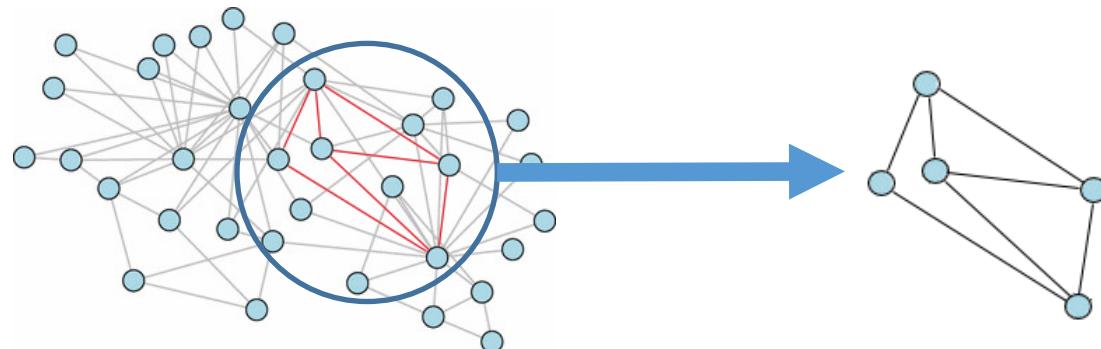
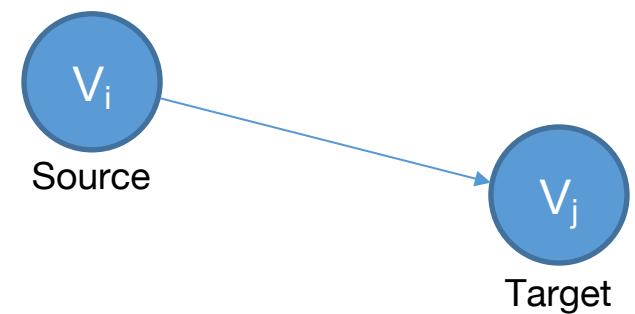


Methods – Graph Notation

Graph - $G = \{V, E\}$

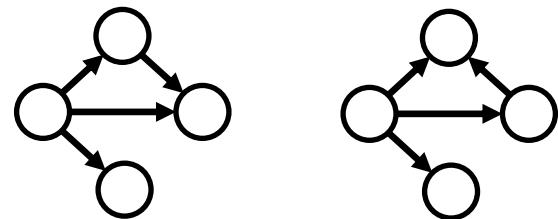
Edge - $E_{ij} = (V_i, V_j)$ with $V_i, V_j \in V, i \neq j$

Subgraph - $G' = \{V', E'\}$ with $V' \subseteq V; E' \subseteq (V', V') \cap E$



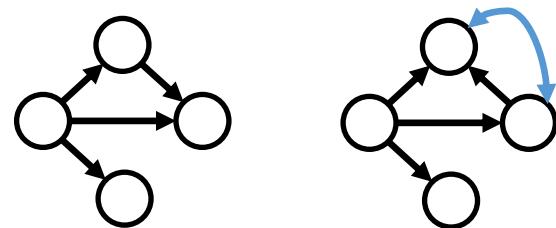
Methods – Graph Notation

Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$



Methods – Graph Notation

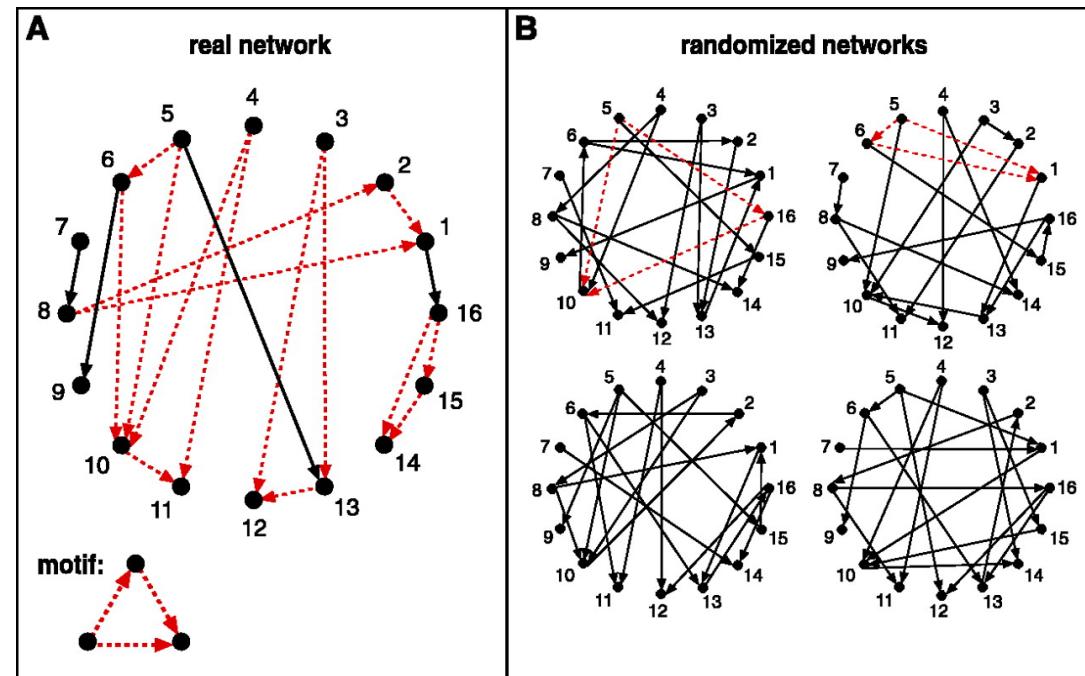
Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$



Methods – Graph Notation

Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$

Z-Score - $Z_m = (f_m - \mu_m)/\sigma_m$

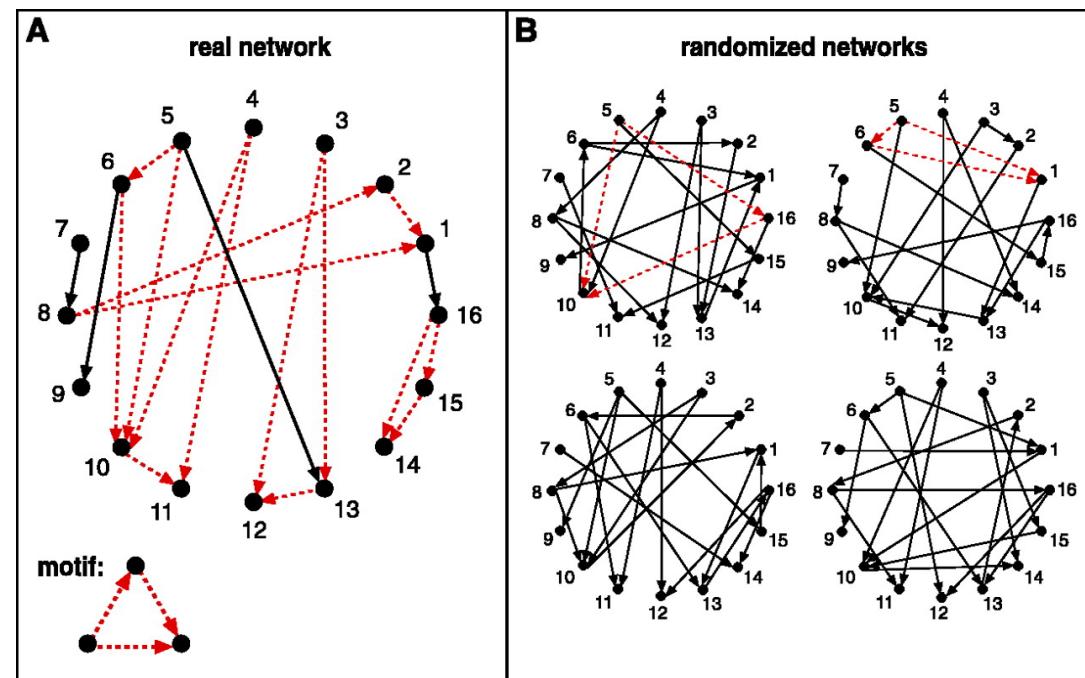


Methods – Graph Notation

Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$

Z-Score - $Z_m = (f_m - \mu_m)/\sigma_m$

$$\begin{aligned} F_0 &= 5 \\ F_1 &= 1 \\ F_2 &= 1 \\ F_3 &= 0 \\ F_4 &= 0 \end{aligned}$$

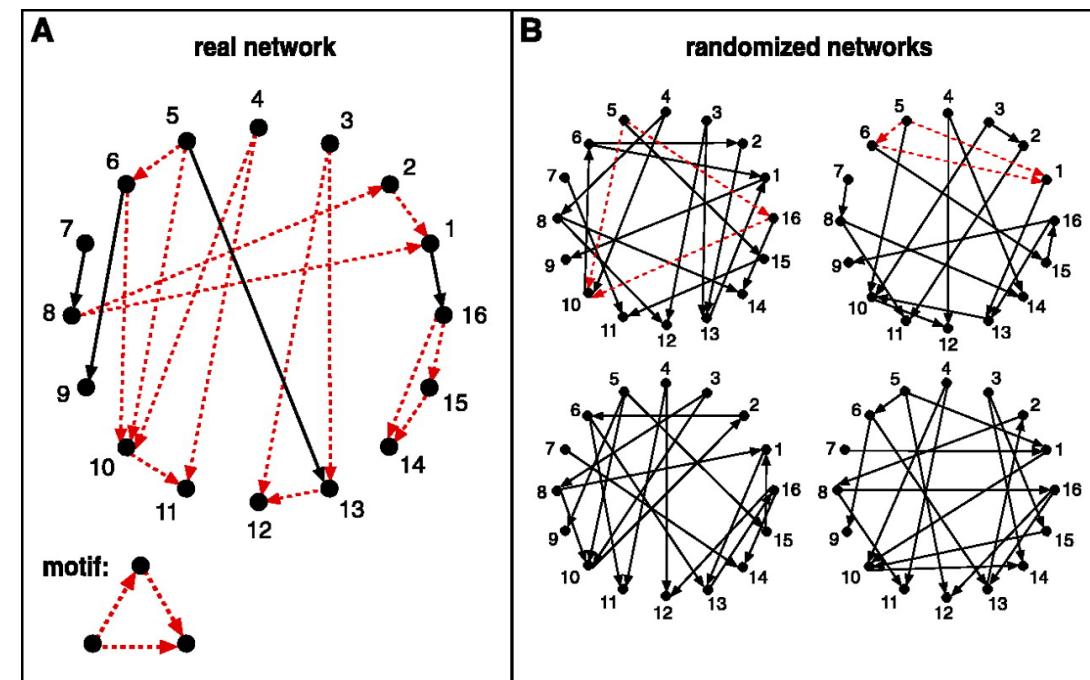


Methods – Graph Notation

Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$

Z-Score - $Z_m = (f_m - \mu_m)/\sigma_m$

$$\begin{aligned} F_0 &= 5 \\ F_1 &= 1 \\ F_2 &= 1 \\ F_3 &= 0 \\ F_4 &= 0 \end{aligned} \quad \xrightarrow{\hspace{1cm}} \quad \begin{aligned} f &= 5 \\ \mu &= 0.5 \\ \sigma &= 0.5 \end{aligned}$$

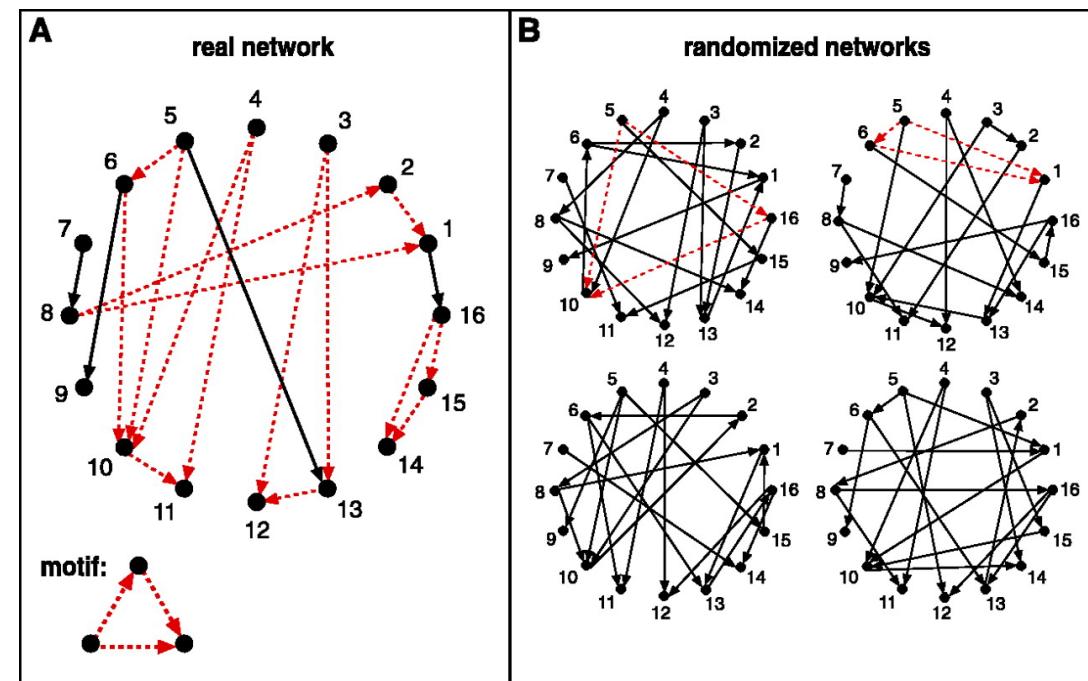


Methods – Graph Notation

Isomorphic group - $G_1 \approx G_2: \exists F; F(V_1) = V_2, E_2 = \tilde{F}(E_1) = (F(V_{1i}), F(V_{1j}))$

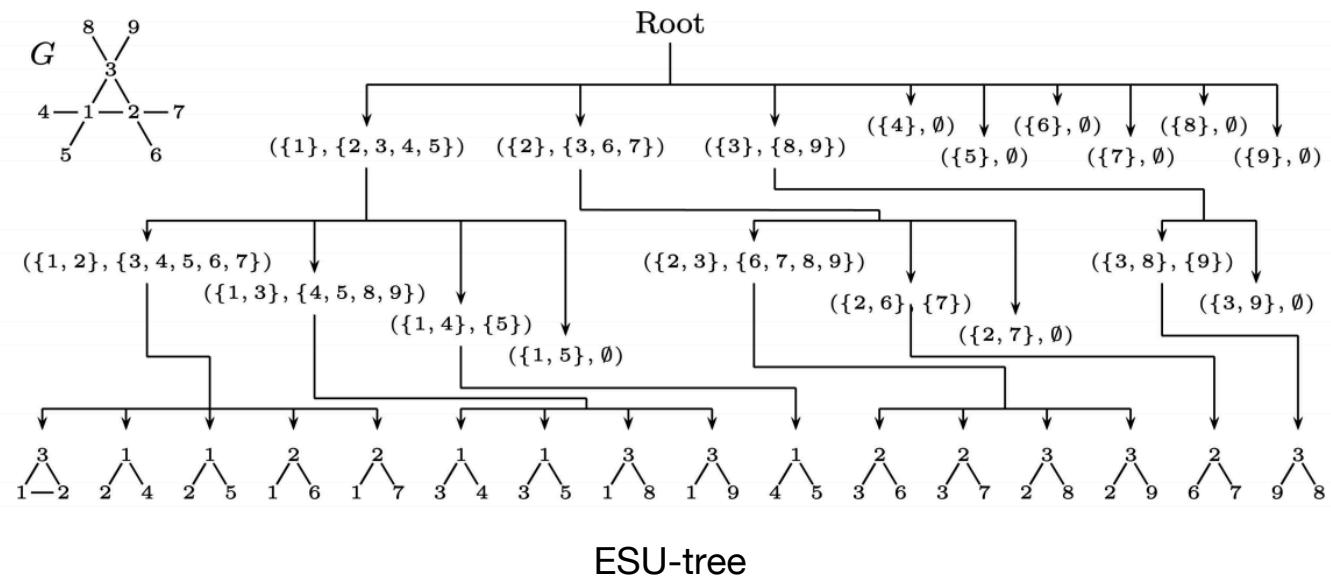
Z-Score - $Z_m = (f_m - \mu_m)/\sigma_m$

$$\begin{array}{l} F_0 = 5 \\ F_1 = 1 \\ F_2 = 1 \\ F_3 = 0 \\ F_4 = 0 \end{array} \xrightarrow{\hspace{2cm}} f = 5 \quad \begin{array}{l} \mu = 0.5 \\ \sigma = 0.5 \end{array} \xrightarrow{\hspace{2cm}} Z = 9$$



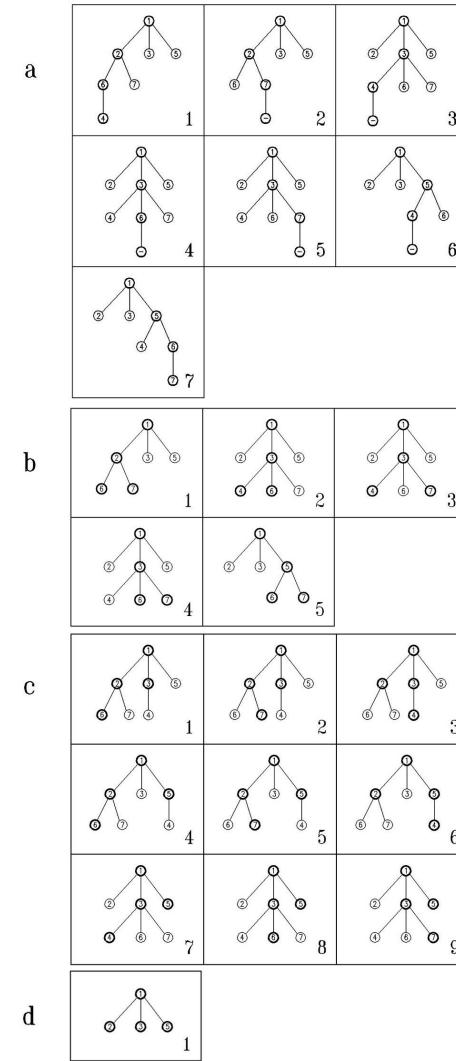
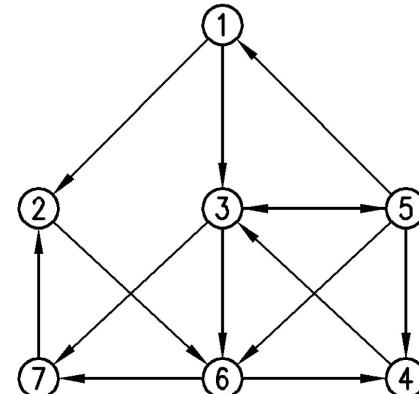
Motif discovery - History

- *Nauty* – 1981
- *Mfinder* – 2002, 2004
 - Brute force method
- *ESU/RAND-ESU* – 2006
 - implements *Nauty*
 - subgraph size restricted to 8

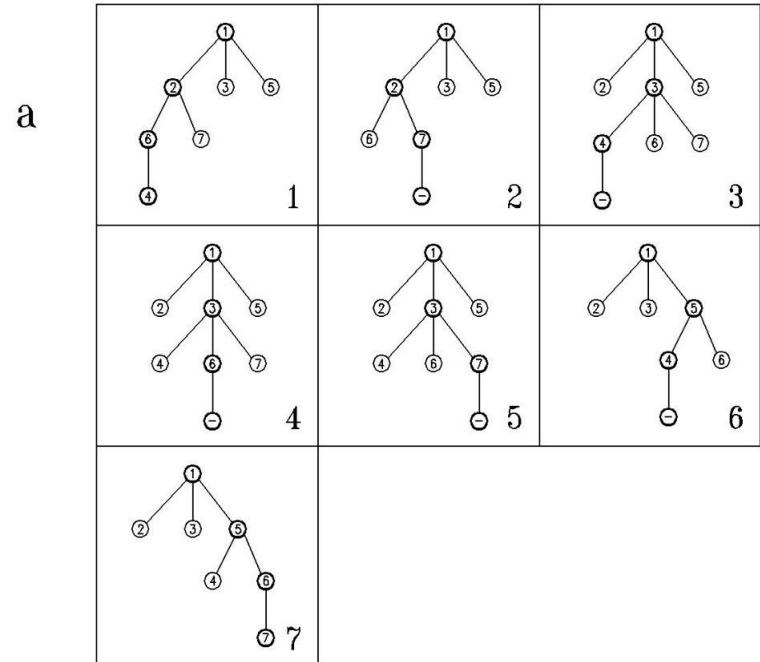
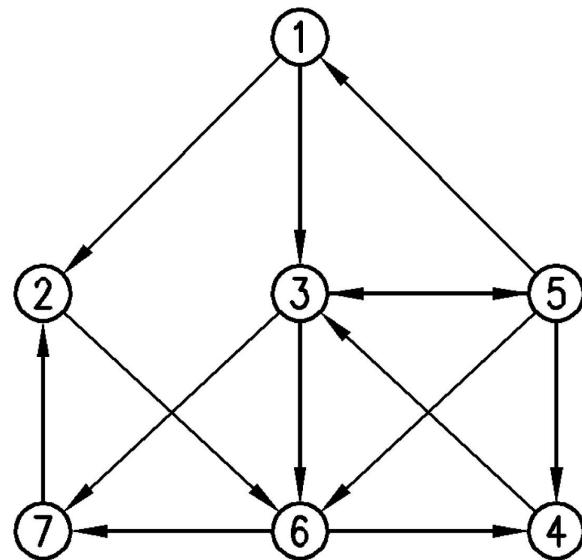


Motif discovery – Kavosh - 2009

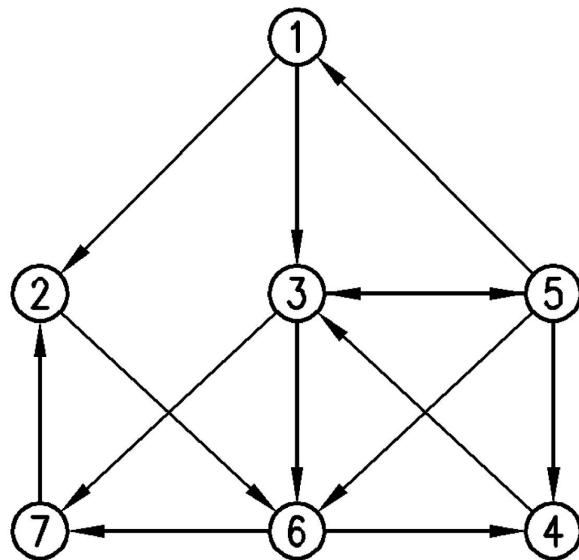
- Iterates on each node.
- Build a tree with max depth k .
- Looks at all compositions of $k-1$
- Isomorphic test with *Nauty*



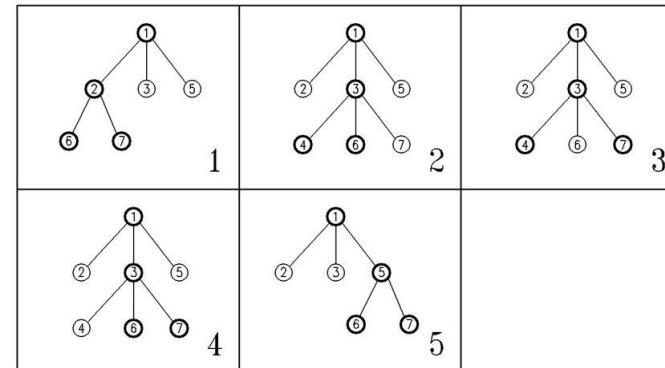
Motif discovery - Kavosh



Motif discovery - Kavosh

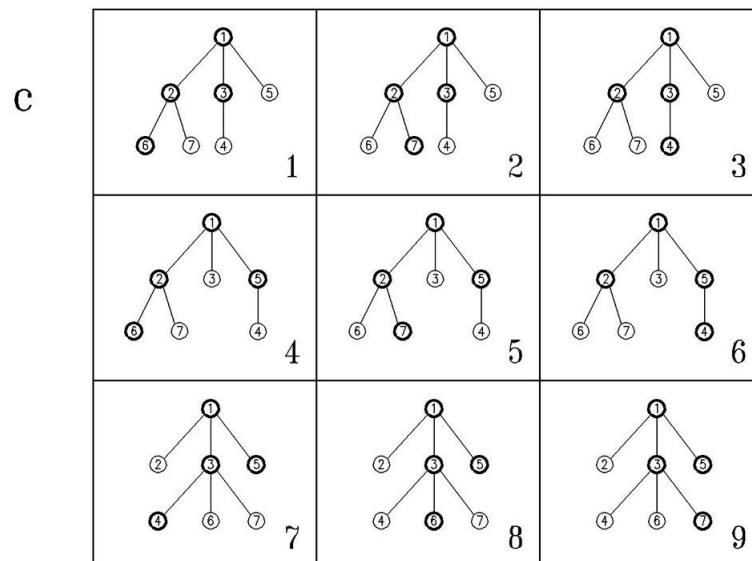
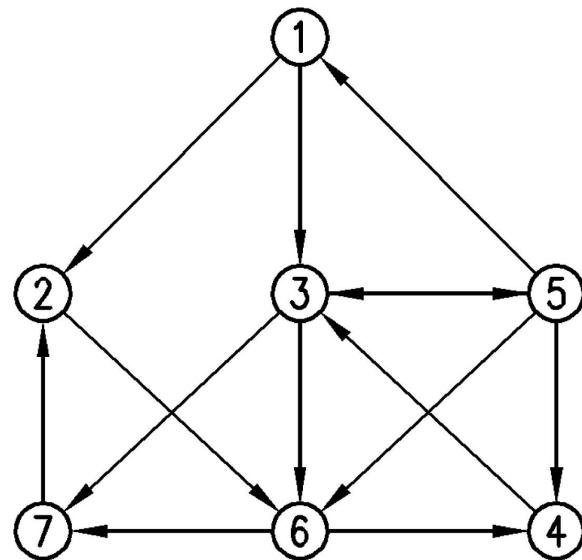


b



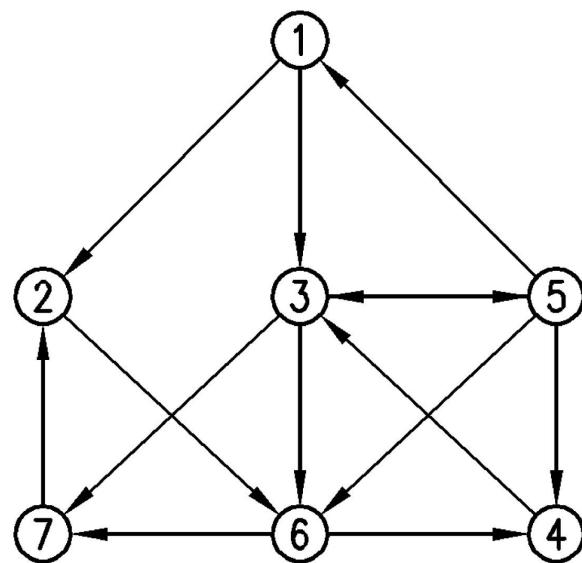
{ 1 , 1 , 2 , 0 }

Motif discovery - Kavosh

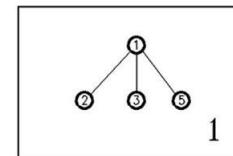


{ 1 , 2 , 1 , 0 }

Motif discovery - Kavosh



d



{ 1 , 3 , 0 , 0 }

Motif discovery - Kavosh

- Pros:

- Optimized memory usage

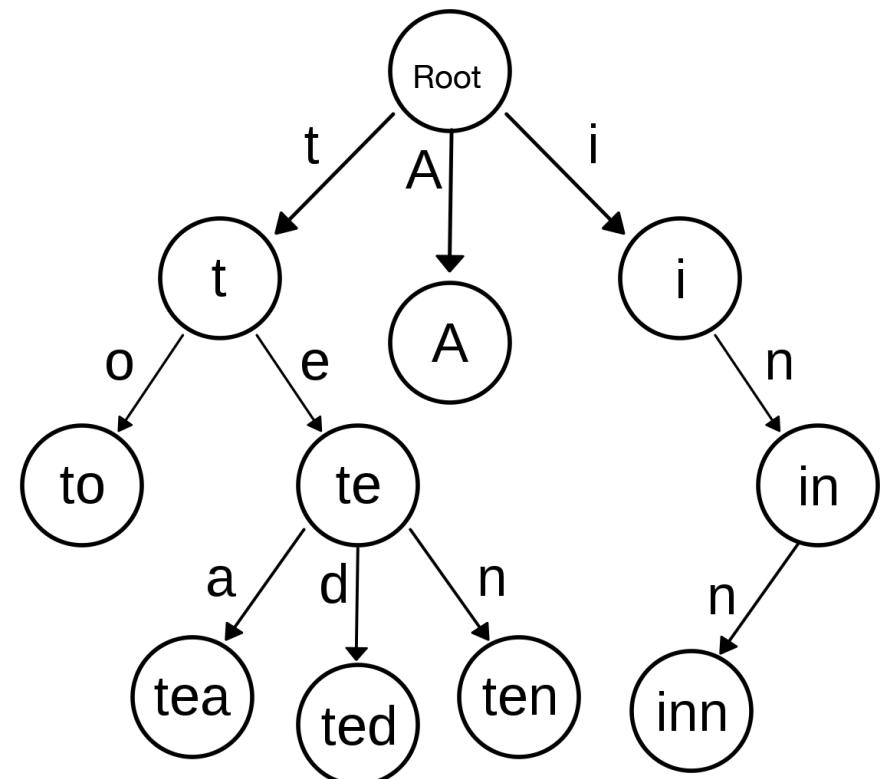
- No subgraph size limit

- Cons:

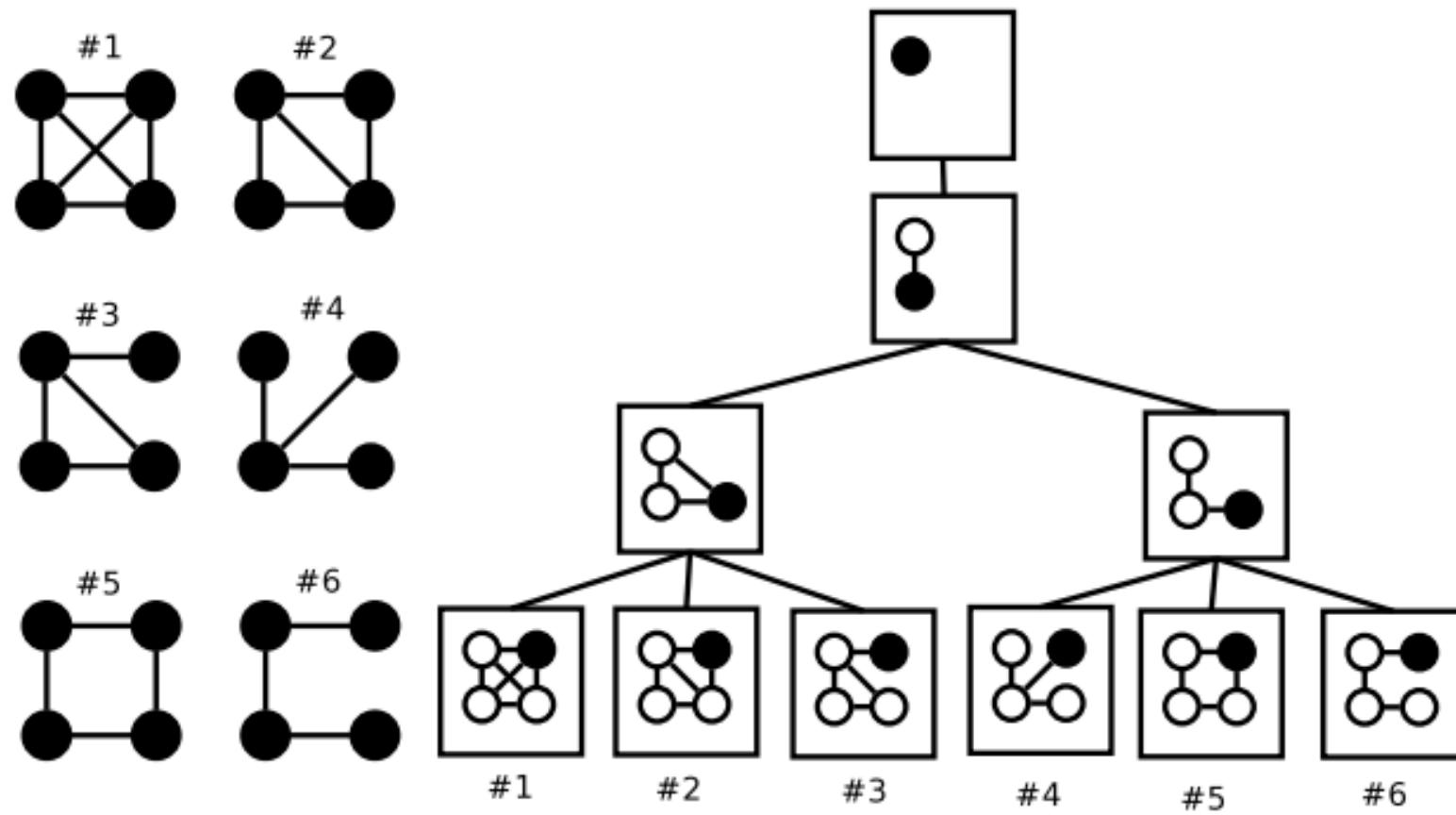
- Not optimized for random network subgraph enumeration

Motif discovery – GtrieScanner - 2010

- Gtrie : Prefix tree structure
- First enumeration to build Trie with ESU
- Canonically labelled subgraph adjacency matrices



Motif discovery - GtrieScanner



Gtrie for undirected motifs of size 4

Motif discovery - GtrieScanner

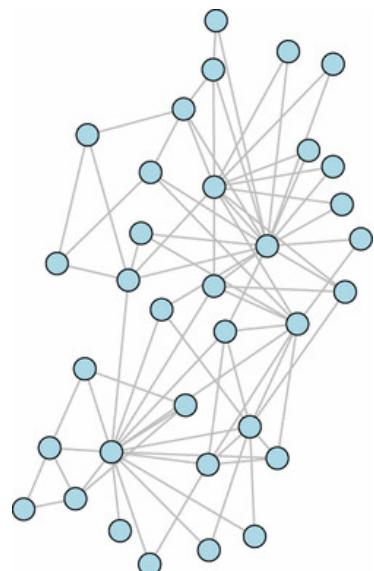
- Pros:

- Optimized for random network subgraph enumeration

- Cons:

- Need to build Gtrie

Motif discovery – Random network generation



Source list Target list

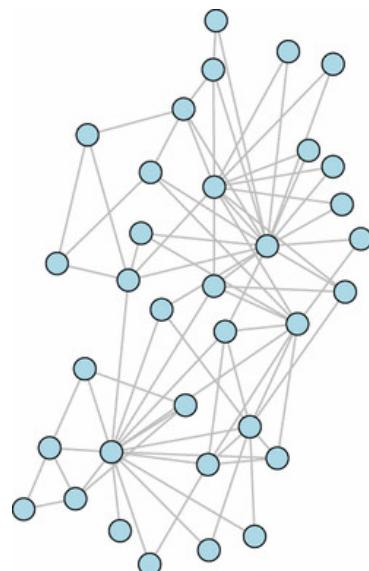
1	→ 12
1	→ 4
2	→ 19
3	→ 8
4	→ 1
4	→ 11
8	→ 5
.	.
.	.

Preserve in and out degree

Source list Target list

1	12
1	4
2	19
3	8
4	1
4	11
8	5
.	.
.	.

Motif discovery – Random network generation



Source list	Target list
1	12
1	4
2	19
3	8
4	1
4	11
8	5
.	.
.	.

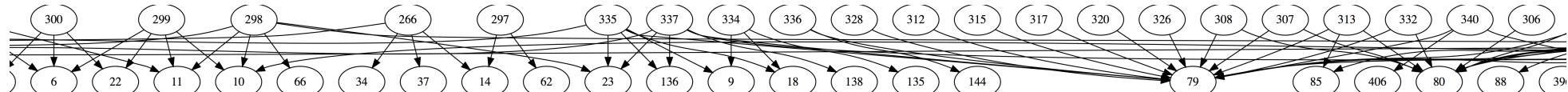
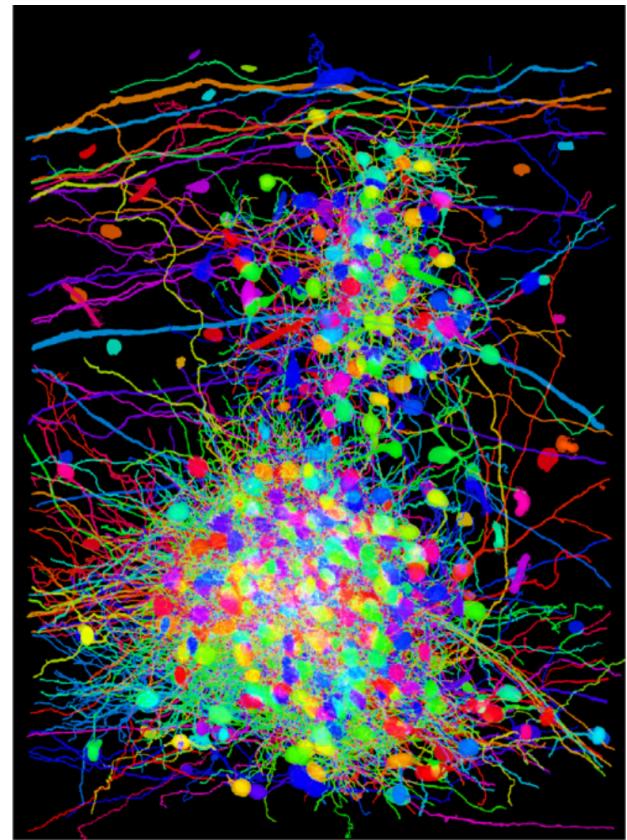
Preserve in and out degree

Assign random source and target

Source list	Target list
1	12
1	4
2	19
3	8
4	1
4	11
8	5
.	.
.	.

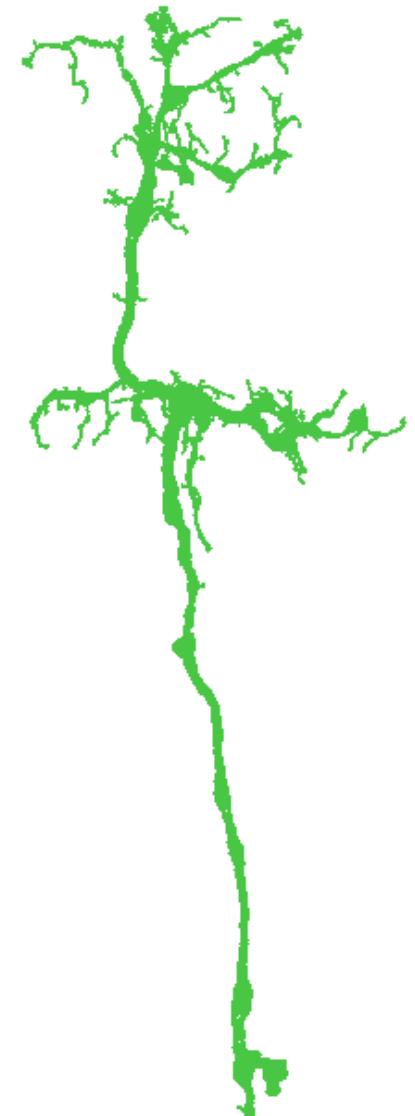
Experiments – LGN dataset

- Mouse dorsal Lateral Geniculate Nucleus
- $400 \times 600 \times 280 \mu\text{m}^3$
- 412 nodes
- 820 directed edges



Experiments – Fib25 dataset

- Medulla Oblongata of fruit fly ($64 \times 66 \times 80 \mu\text{m}^3$)
- Initially 798 neuron segments and 31000 synapses
- Removed edges with less than 10 connections
- 279 nodes
- 441 directed edges



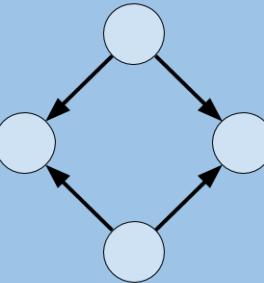
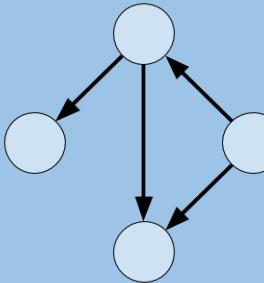
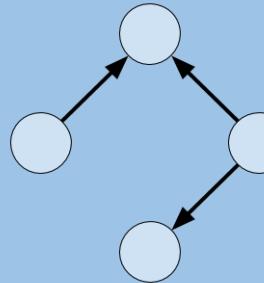
Survey and application

- Small subgraph enumeration
 - Motif relevance
 - Computing performance

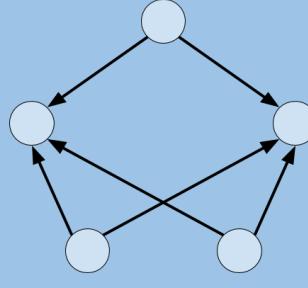
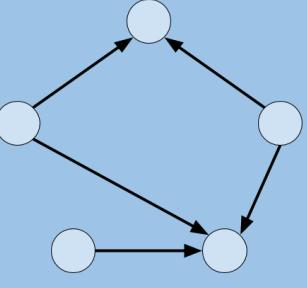
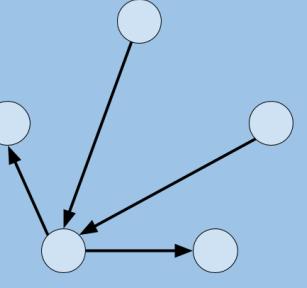
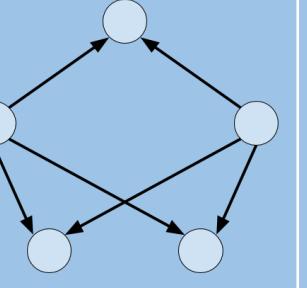
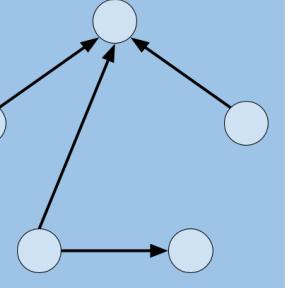
Results – Small motif relevance

Motif Size	LGN		Fib-25	
	Number of Subgraphs	Number of iso. classes	Number of Subgraphs	Number of iso. classes
4	89 641	18	17 239	91
5	1 289 357	85	157 160	666

Results – Small motif relevance

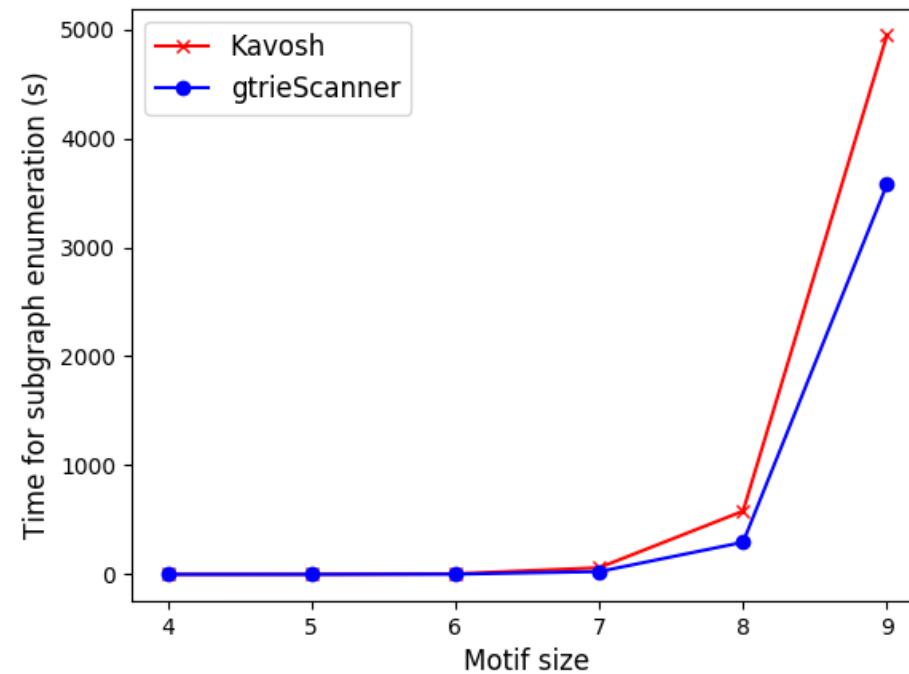
Motif ID	4.1	4.2	4.3	
Motif				
LGN	Normalized frequency	1.8%	0.15%	42%
	Z-Score	7.7	4.0	-10.9
Fib-25	Normalized frequency	0.68%	1.6%	18%
	Z-Score	14.1	5.2	-3.9

Results – Small motif relevance

Motif ID	5.1	5.2	5.3	5.4	5.5
Motif					
LGN	Normalized frequency	0.15%	3.0%	0.04%	0.08%
	Z-Score	1.4	1.1	-2.3	14.6
Fib-25	Normalized frequency	0.06%	0.79%	6.3%	0.01%
	Z-Score	52.1	10.7	-0.7	9.4

Results – Small motif computing performance

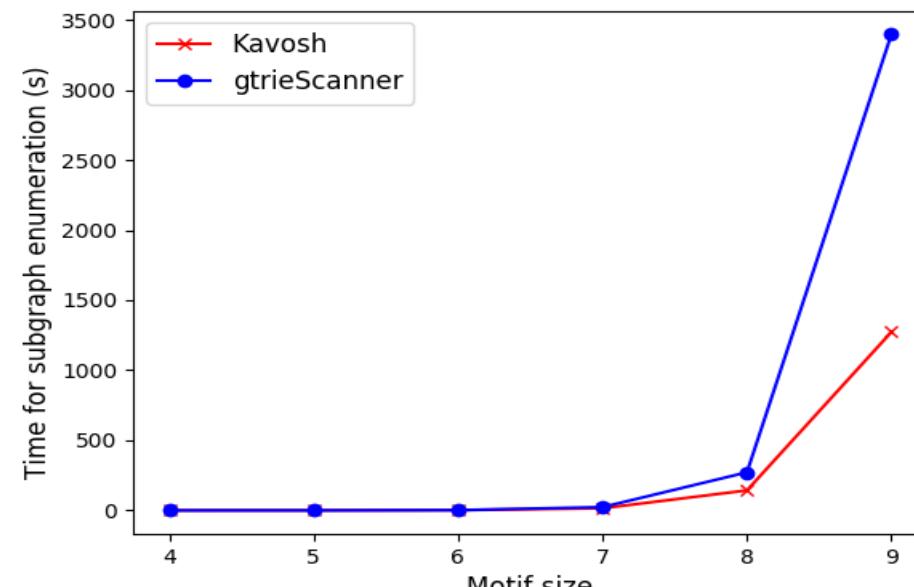
Exponential increase of runtime for both algorithms



Total running time on Fib-25 with 1000 random networks

Results – Small motif computing performance

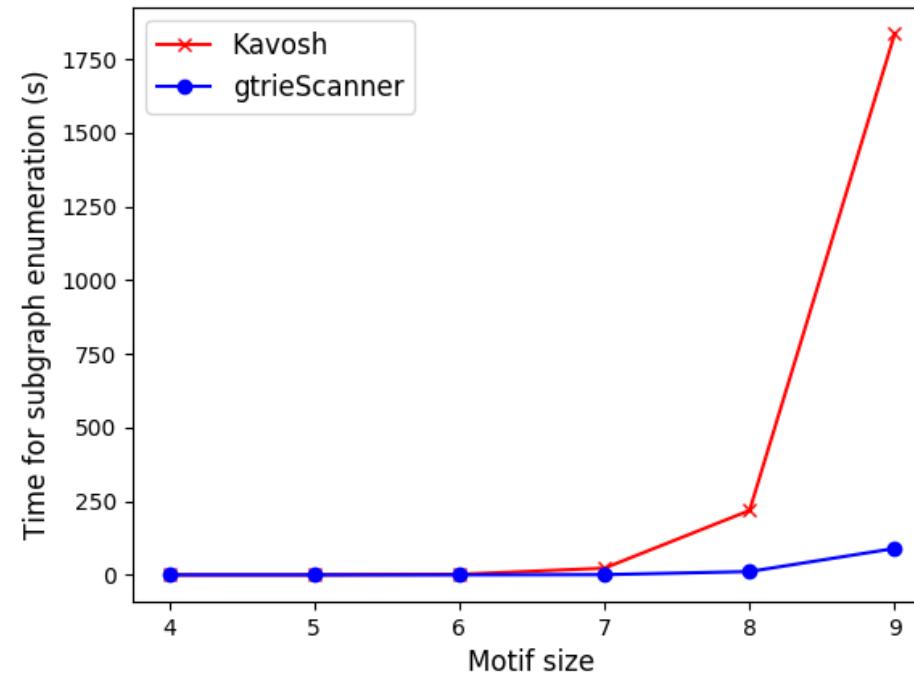
Kavosh is much faster for one graph enumeration



Running time on Fib-25 for the original network

Results – Small motif computing performance

GtrieScanner is very quick to enumerate random graphs

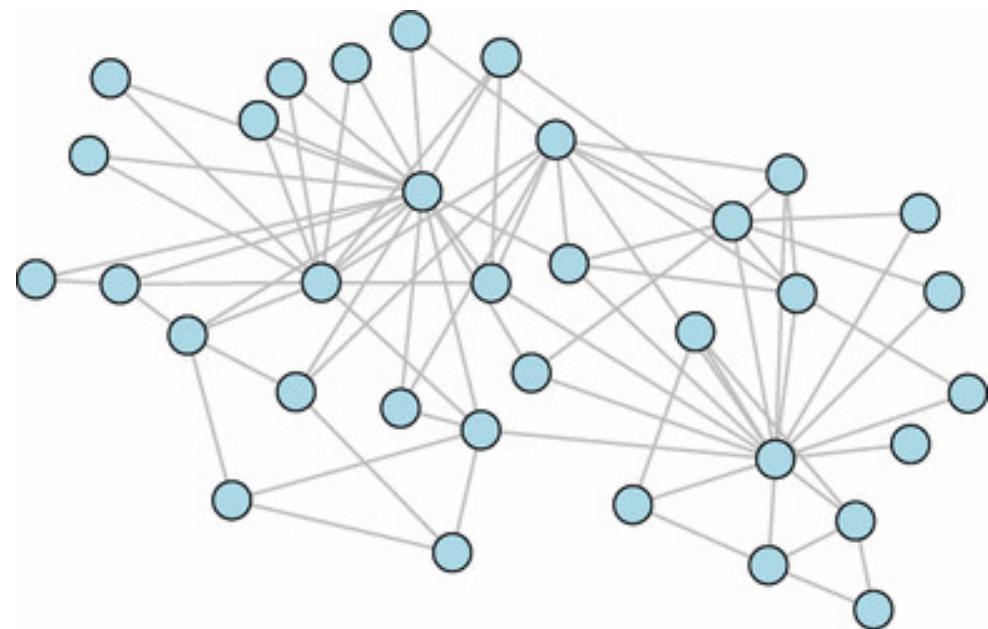


Running time on Fib-25 for 1000 random networks

Large Subgraph enumeration: Exploration

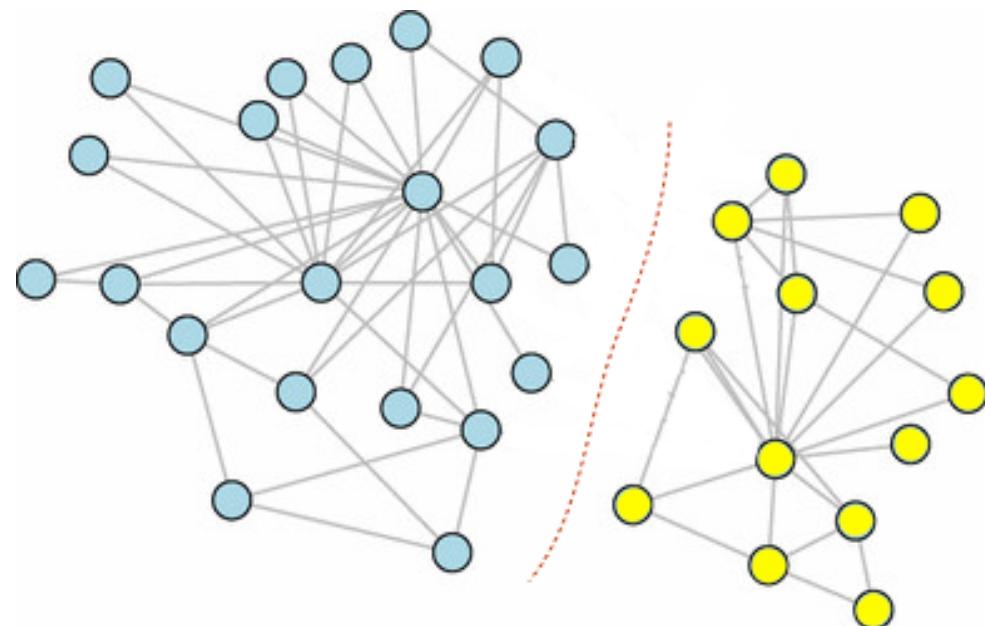
Methods – Graph theory

Graph cut - $C(V) \in 1, \dots, n_c$



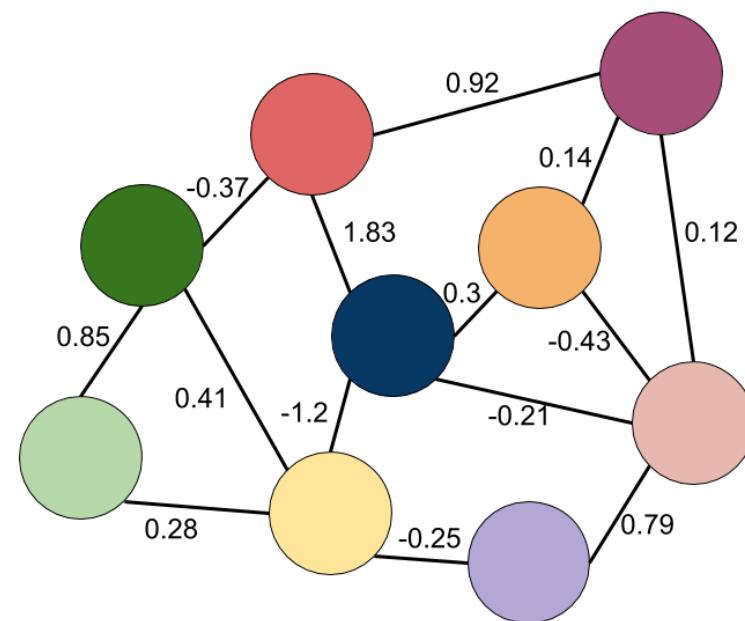
Methods – Graph theory

Graph cut - $C(V) \in 1, \dots, n_c$



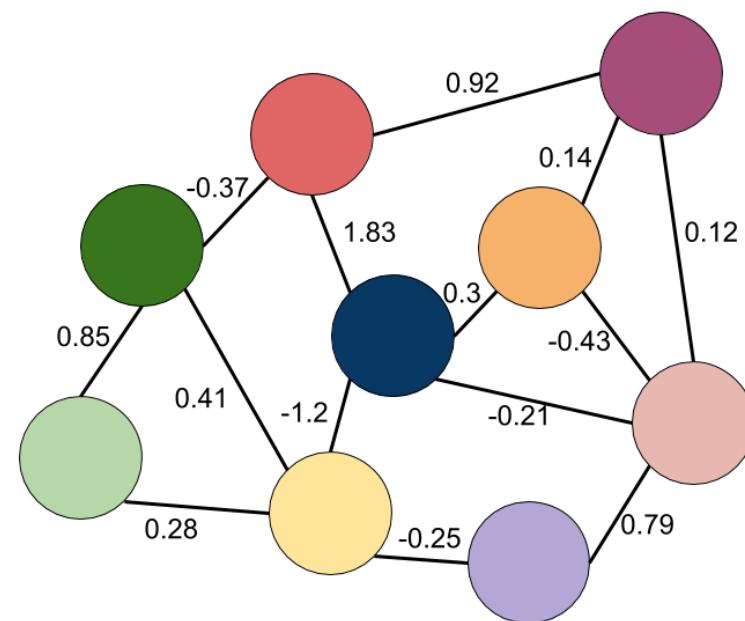
Motif discovery – Multicut clustering

- Minimum Cost Multicut Problem
- NP hard problem



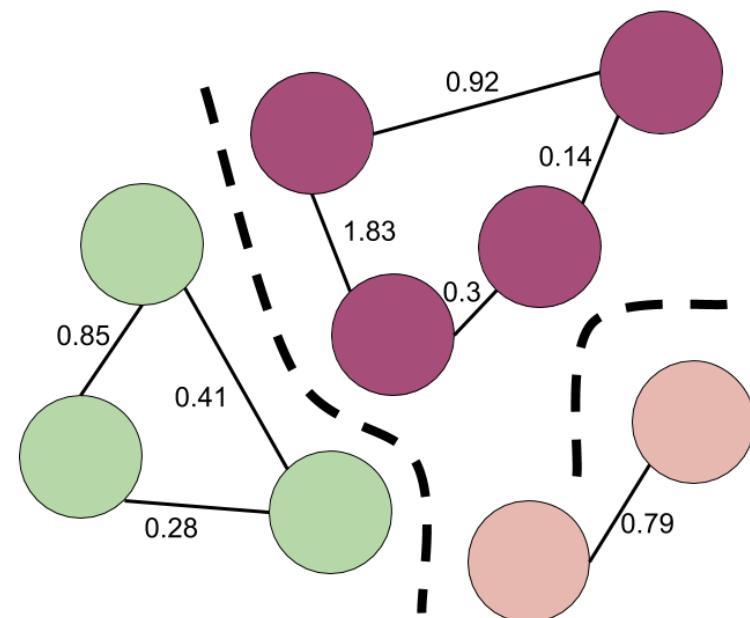
Motif discovery – Multicut clustering

- Minimum Cost Multicut Problem
- NP hard problem
- Modified Kernighan-Lin algorithm
- Iterative algorithm to minimize cost function

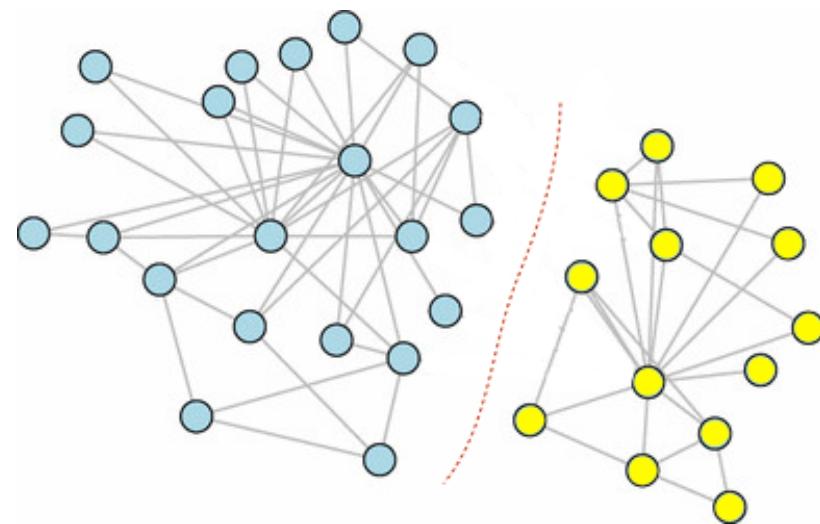


Motif discovery – Multicut clustering

- Minimum Cost Multicut Problem
- NP hard problem
- Modified Kernighan-Lin algorithm
- Iterative algorithm to minimize cost function
- No need to know number of clusters
- No need for sources

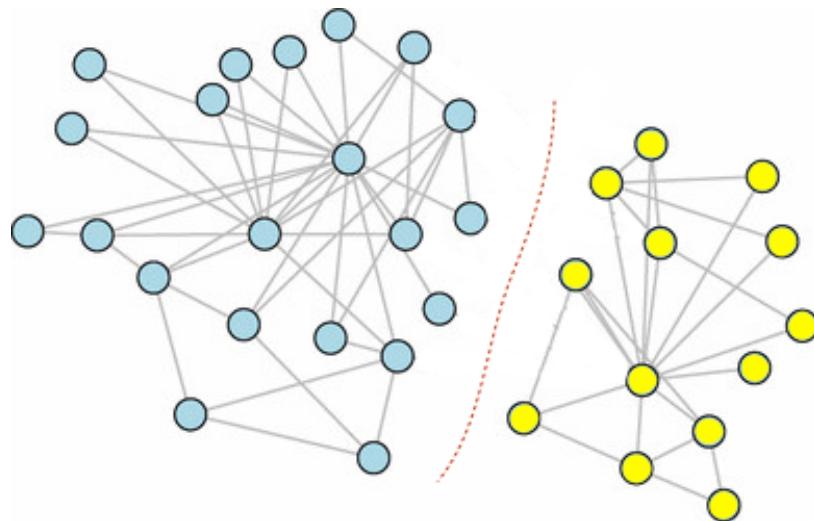


OUR PROPOSAL

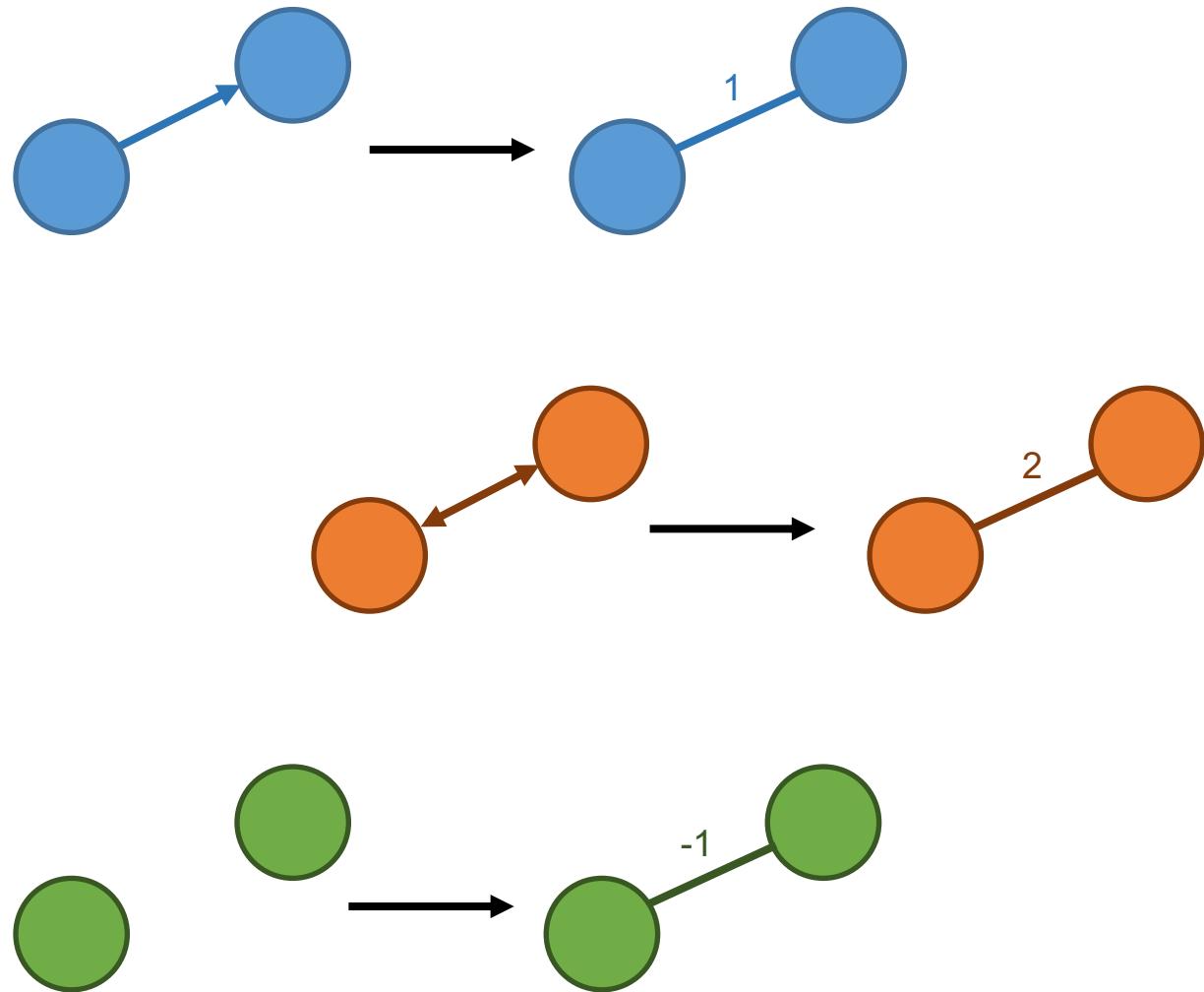


1. Cluster graph

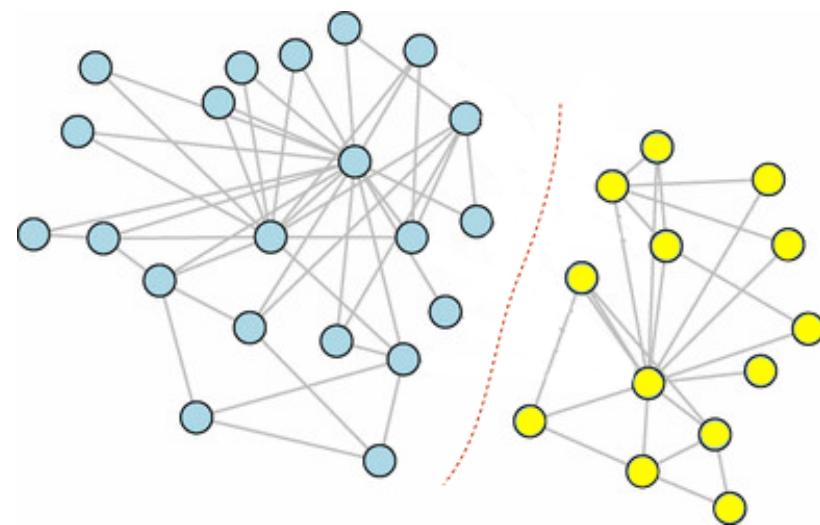
OUR PROPOSAL



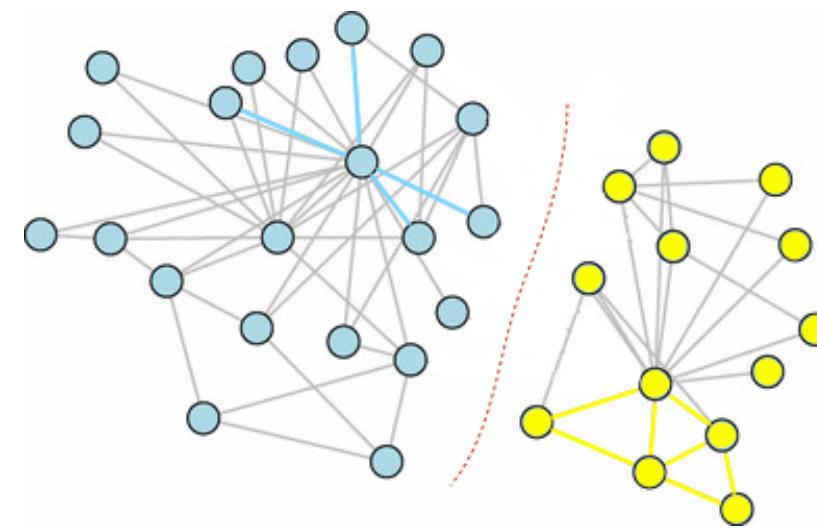
1. Cluster graph



OUR PROPOSAL



1. Cluster graph



2. Enumerate subgraphs

Exploration

- Large subgraph enumeration
 - Subgraph frequency
 - Computing performance

Results – Large motif clustering on *LGN*

Cluster ID	Number of nodes	Number of edges
0	200	345
1	38	45
2	20	19
3	5	4
5	10	7
8	10	9
10	36	43
15	4	3
16	78	70

Results – Large motif clustering on *LGN*

Cluster ID	Number of nodes	Number of edges
0	200	345
1	38	45
2	20	19
3	5	4
5	10	7
8	10	9
10	36	43
15	4	3
16	78	70

Results – Large motif relevance on *LGN*

- 3 622 227 359 subgraphs occurrences in original graph
- 1 633 037 684 subgraphs occurrences in clustered graph
- 55% search space reduction

Results – Large motif relevance on *LGN*

Motif ID	1	2	3	4	5	6	7	8
Z-Score 	5100	307	115	83.3	72.8	56.7	39.3	22.9
Original frequency 	510	2 385	1 945	32 235	8 415	10 055	49 182	186 620
Clusters frequency	510	2 385	1 945	32 235	8 415	10 055	465	186 620

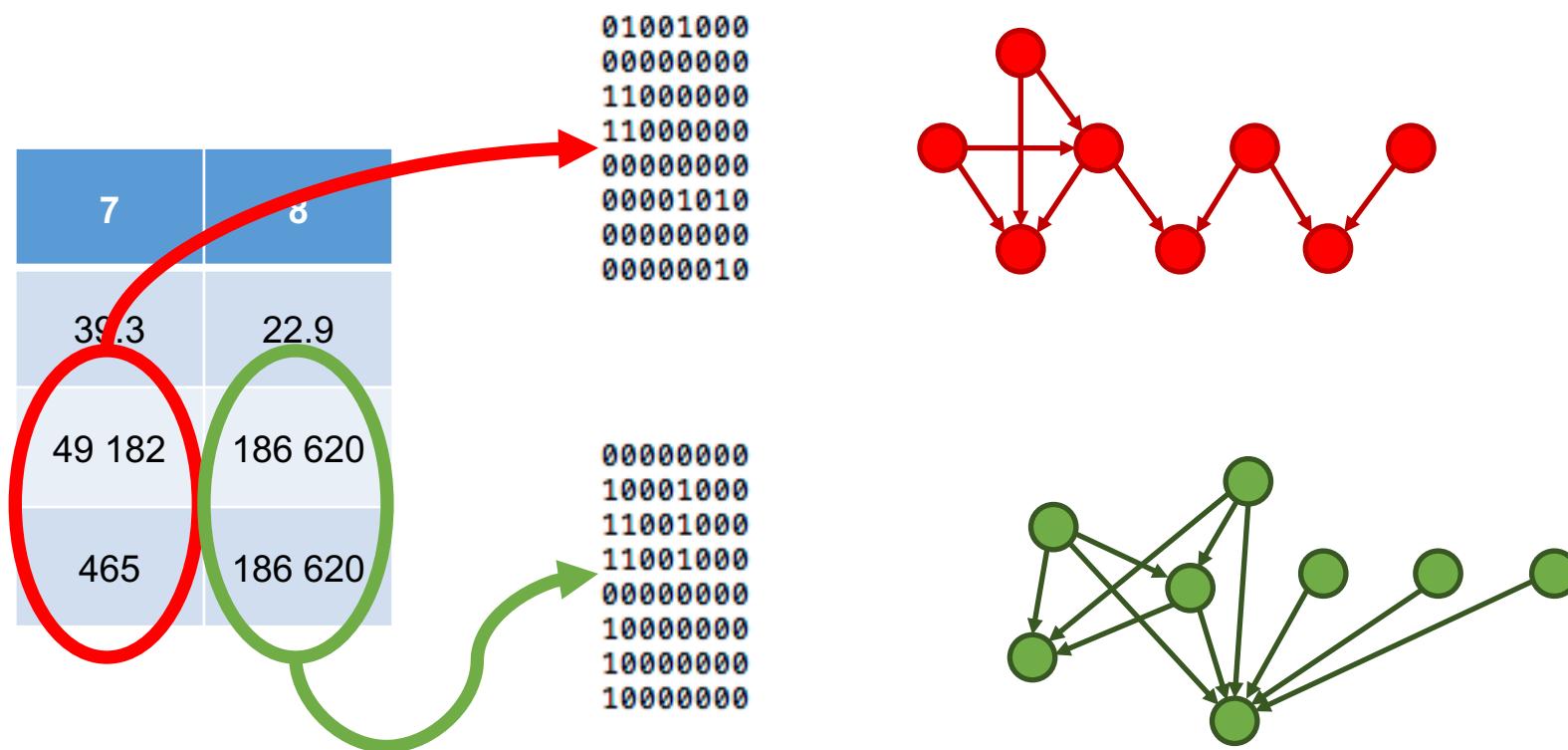
Results – Large motif relevance on *LGN*

Motif ID	1	2	3	4	5	6	7	8
Z-Score ↑	5100	307	115	83.3	72.8	56.7	39.3	22.9
Original frequency ↑	510	2 385	1 945	32 235	8 415	10 055	49 182	186 620
Clusters frequency	510	2 385	1 945	32 235	8 415	10 055	465	186 620

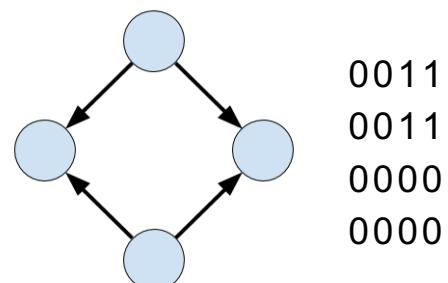
Results – Large motif relevance on *LGN*

Motif ID	1	2	3	4	5	6	7	8
Z-Score ↑	5100	307	115	83.3	72.8	56.7	39.3	22.9
Original frequency ↑	510	2 385	1 945	32 235	8 415	10 055	49 182	186 620
Clusters frequency	510	2 385	1 945	32 235	8 415	10 055	465	186 620

Results – Large motif relevance on *LGN*



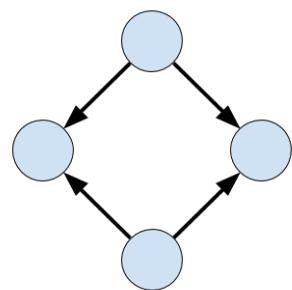
Results – Large motif relevance on *LGN*



Motif 4.1

0011
0011
0000
0000

Results – Large motif relevance on *LGN*



0011
0011
0000
0000

Motif 4.1

00110000	00110000	00110000	00110000	00000000	00110001	01001000	00000000
10110010	10110000	10110000	10110000	10001000	10110000	00000000	10001000
00000000	00000000	00000001	00000000	11001100	00000000	11000000	11001000
00000000	00000000	00000000	00000000	11001000	00000000	11000000	11001000
10110000	10110000	10110000	10110000	00000000	10110000	00000000	00000000
10010010	10110000	10110000	10110000	10000000	10110000	00001010	10000000
00100000	10110000	00100000	00100000	10000000	00100000	00000000	10000000
00100000	00100000	00000000	00100000	10000000	00000000	00000010	10000000

Motif 1

Motif 2

Motif 3

Motif 4

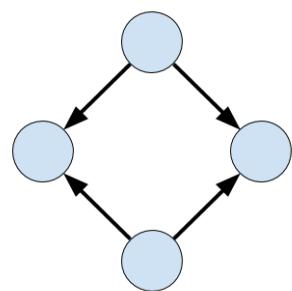
Motif 5

Motif 6

Motif 7

Motif 8

Results – Large motif relevance on *LGN*



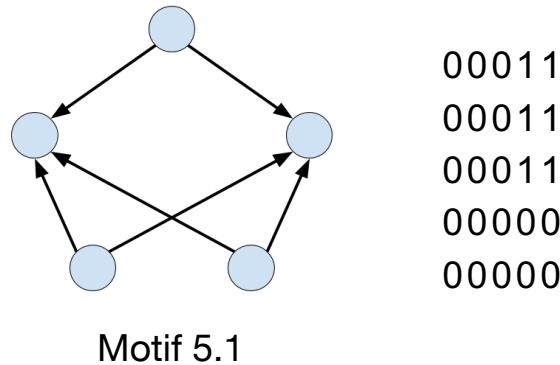
0011
0011
0000
0000

Motif 4.1

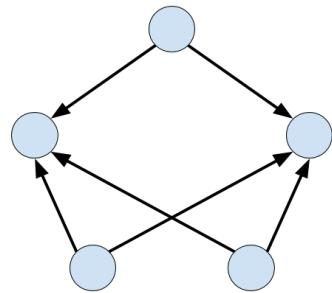
Identified in all size 8 motifs of LGN dataset

Motif 1	Motif 2	Motif 3	Motif 4	Motif 5	Motif 6	Motif 7	Motif 8
0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10010010 00100000 00100000	0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10010010 00100000 00100000	0 1100 0 1 1100 0 0 0000 1 0 0000 0 1 1100 0 10010010 00100000 00100000	0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10010010 00100000 00100000	00000000 10001000 11001100 11000000 00000000 10000000 10000000 10000000	0 1100 1 1 1100 0 0 0000 0 0 0000 0 1 1100 0 11001100 00000000 00000000	01001000 00000000 11000000 11000000 00000000 11000000 00000010 00000000	00000000 10001000 11001100 11000000 00000000 00000000 10000000 10000000

Results – Large motif relevance on *LGN*



Results – Large motif relevance on *LGN*



Motif 5.1

00011
00011
00011
00000
00000

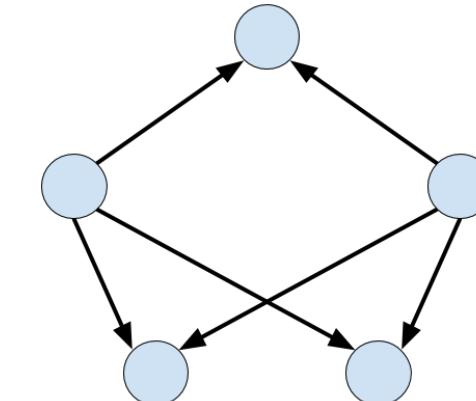
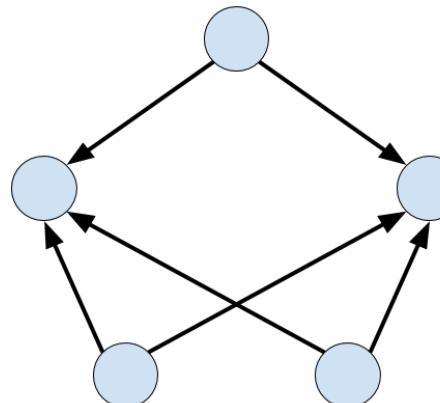
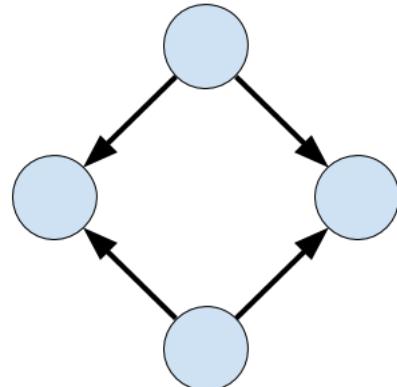
Motif 1	Motif 2	Motif 3	Motif 4	Motif 6
0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10010010 00100000 00100000	0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10110000 10110000 00100000	0 1100 0 1 1100 0 0 0000 1 0 0000 0 1 1100 0 10110000 10110000 00100000	0 1100 0 1 1100 0 0 0000 0 0 0000 0 1 1100 0 10110000 00100000 00100000	0 1100 1 1 1100 0 0 0000 0 0 0000 0 0 0000 0 1 1100 0 1 1100 0 00100000 00000000

Results – Large motif performance on *LGN*

	GtrieScanner	OUR PROPOSAL
Clustering time (s)	0	4
Subgraph enumeration time (s)	15 860	4 821
Total time (s)	15 860 (~4h30m)	4 825 (~1h20m)

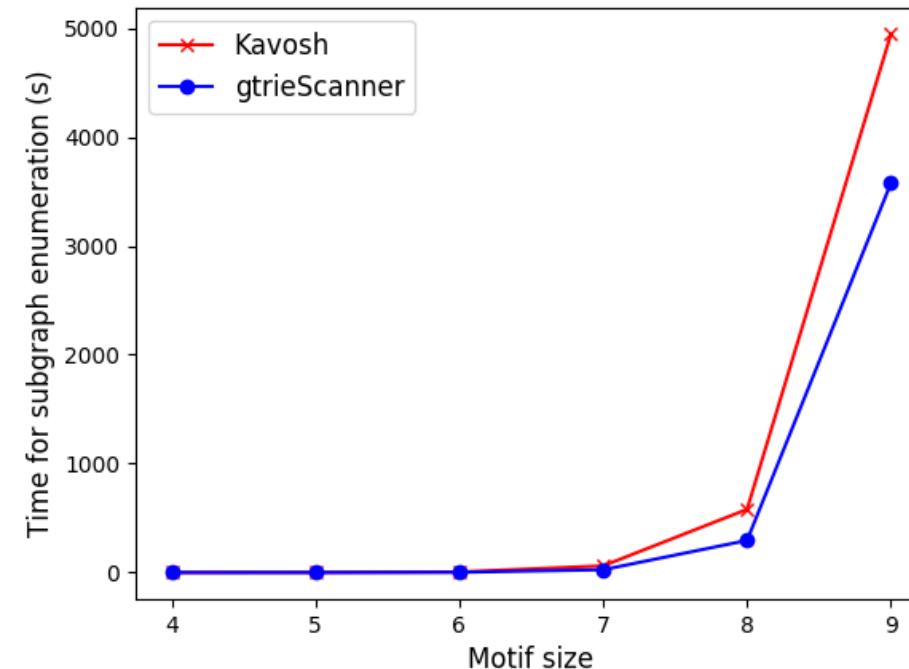
Discussion

- From the survey, same motifs have been identified.
- Two real connectomes (mouse and fruit fly).



Discussion

- Kavosh is very efficient to enumerate a single graph
- GtrieScanner is more time efficient when computing subgraph statistics



Discussion

- We have proposed a solution to speed up subgraph enumeration

3x faster

- Drawbacks:
 - Clusters are mostly uneven.

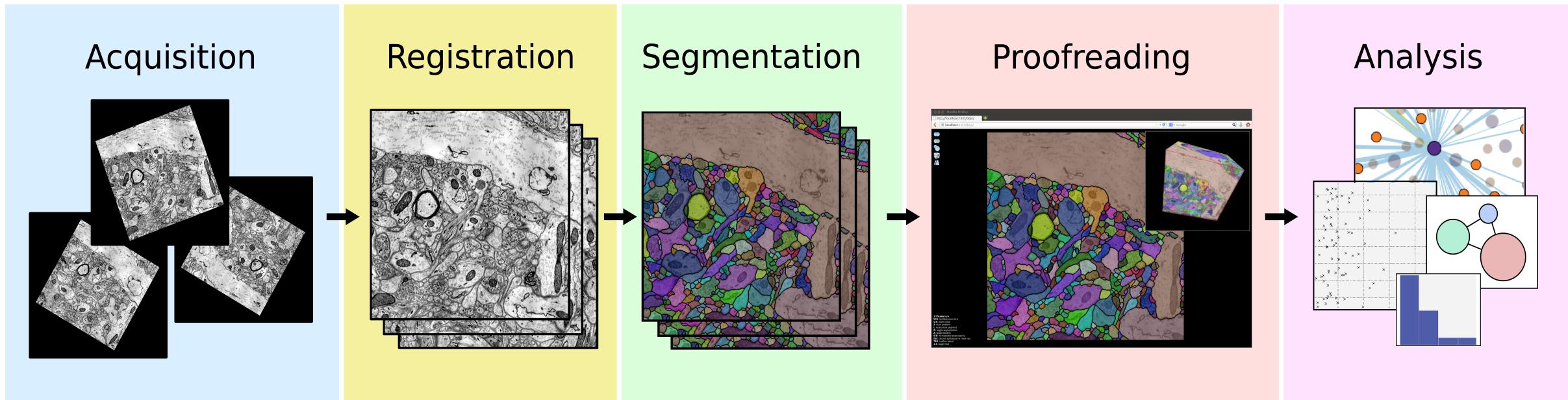
Future Works

- Improve the clustering on our algorithm.
- Apply small motif discovery to more brain wiring diagram.

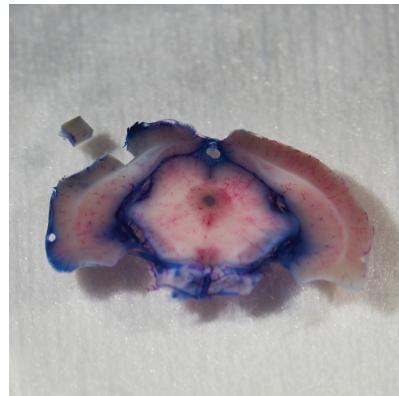
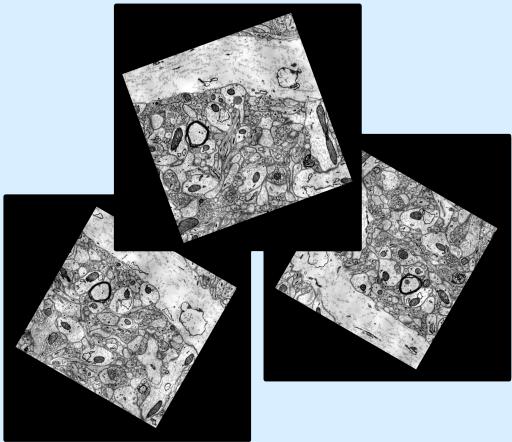
Thank you to the entire VCG lab and to Prof. Kathryn Hess!

Thank you

Connectomics

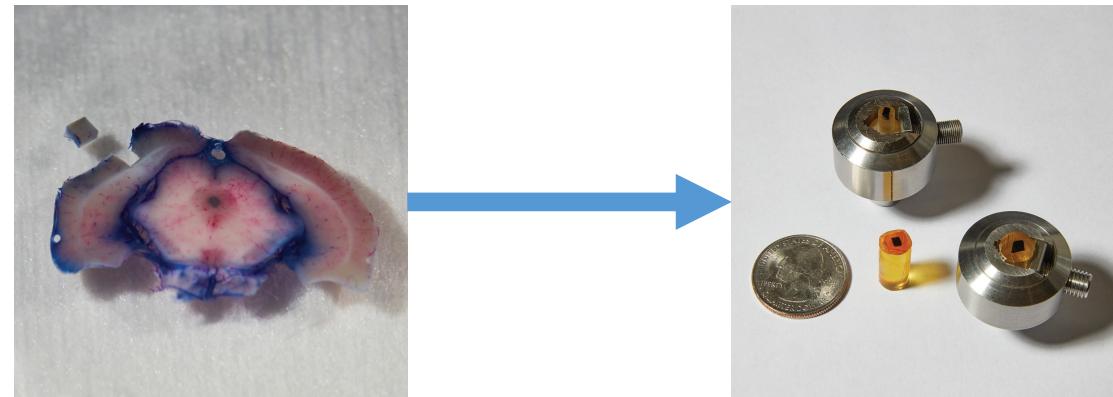
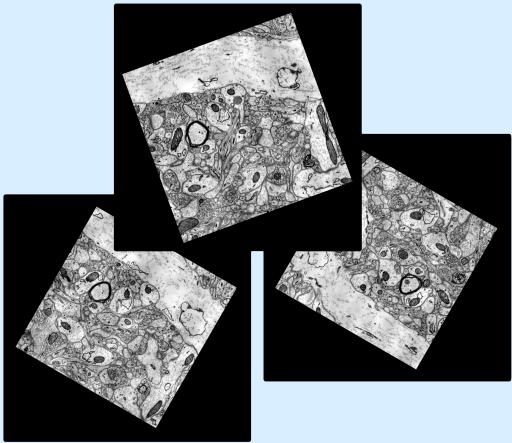


Acquisition

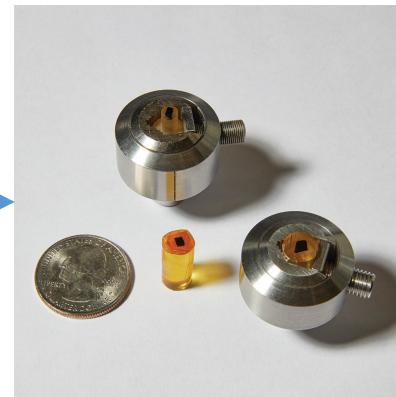


Cubic millimeter of
brain tissue.

Acquisition

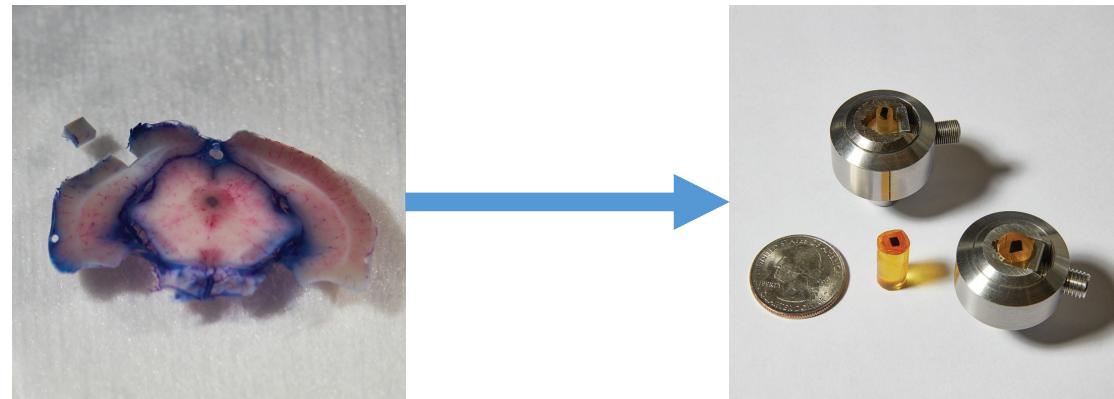
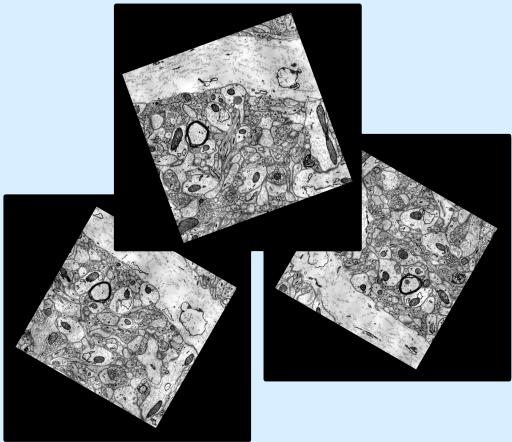


Cubic millimeter of
brain tissue.

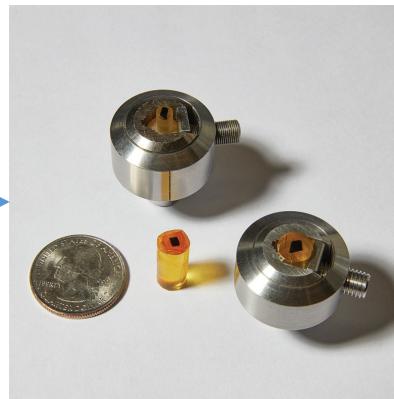


Thinly sliced in an
acrylic encasing.

Acquisition



Cubic millimeter of
brain tissue.

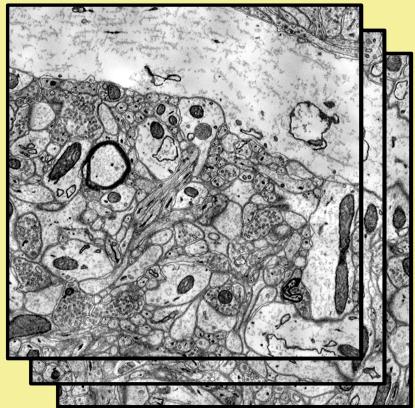


Thinly sliced in an
acrylic encasing.



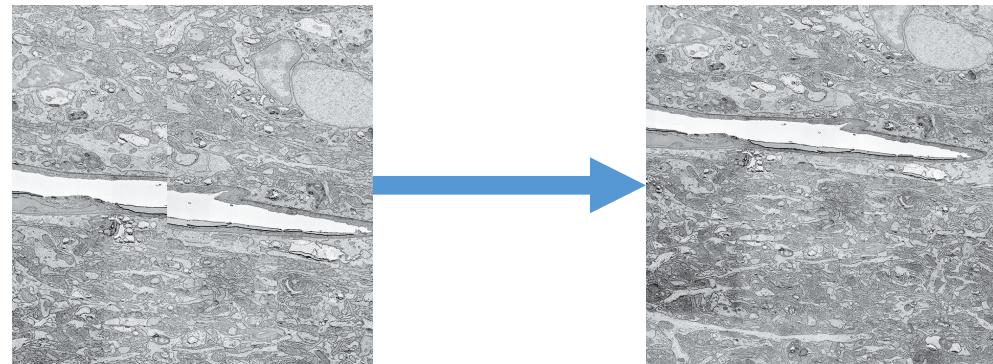
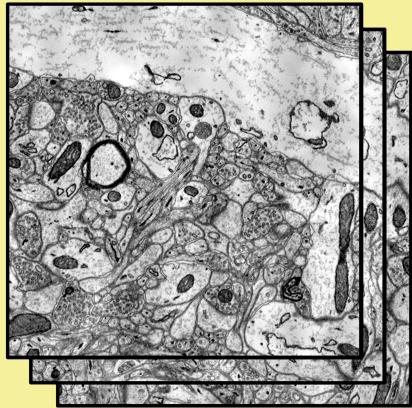
EM imaging

Registration



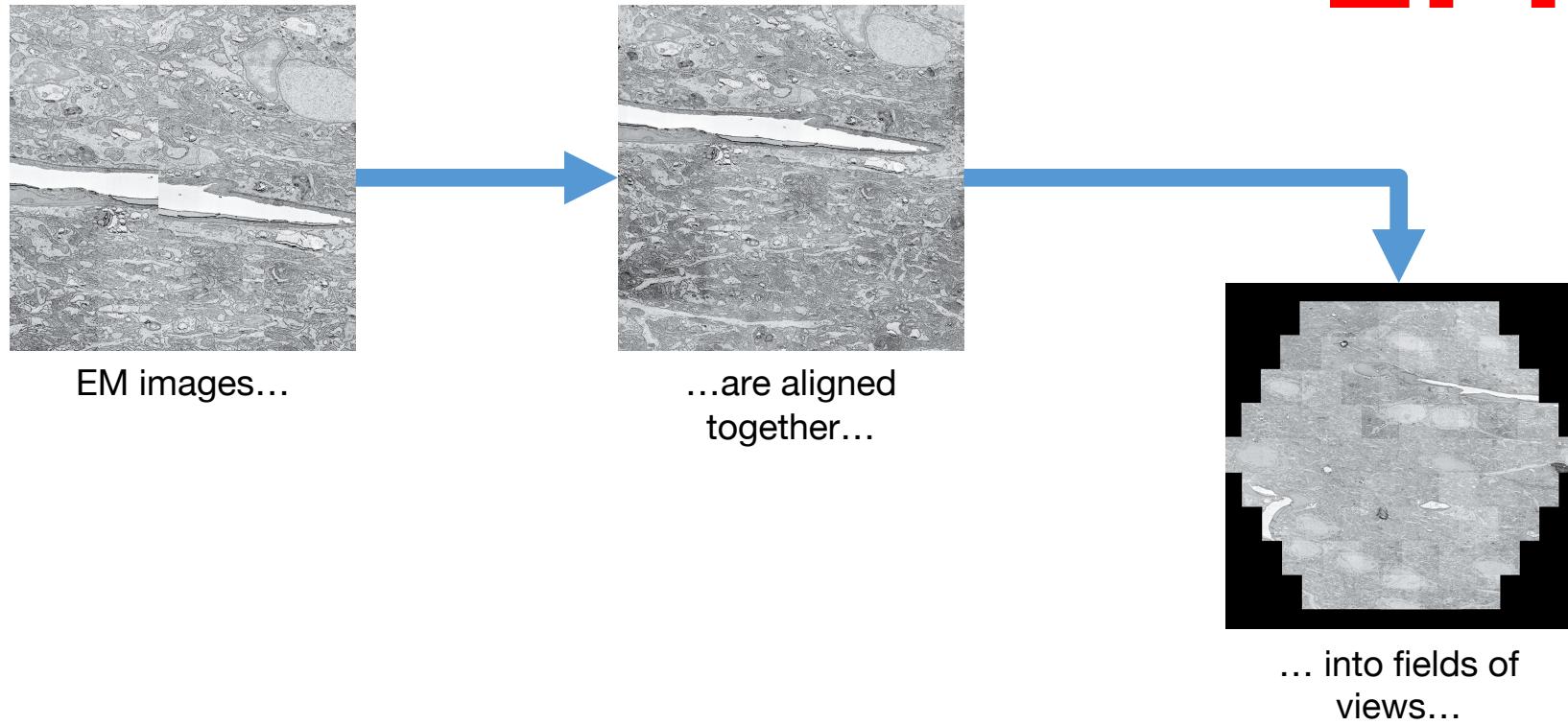
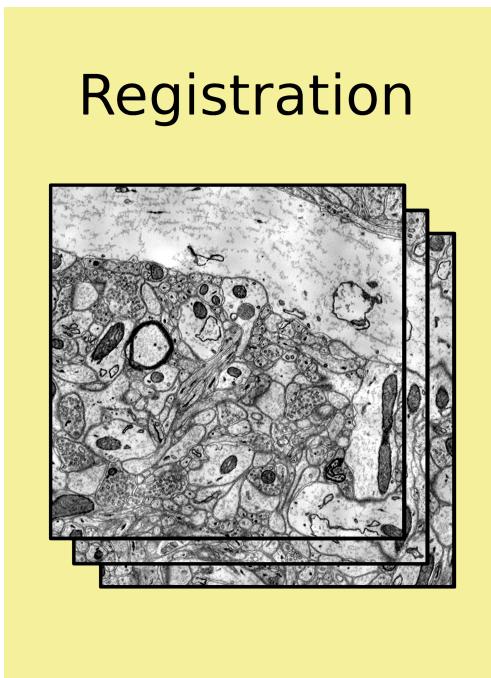
EM images...

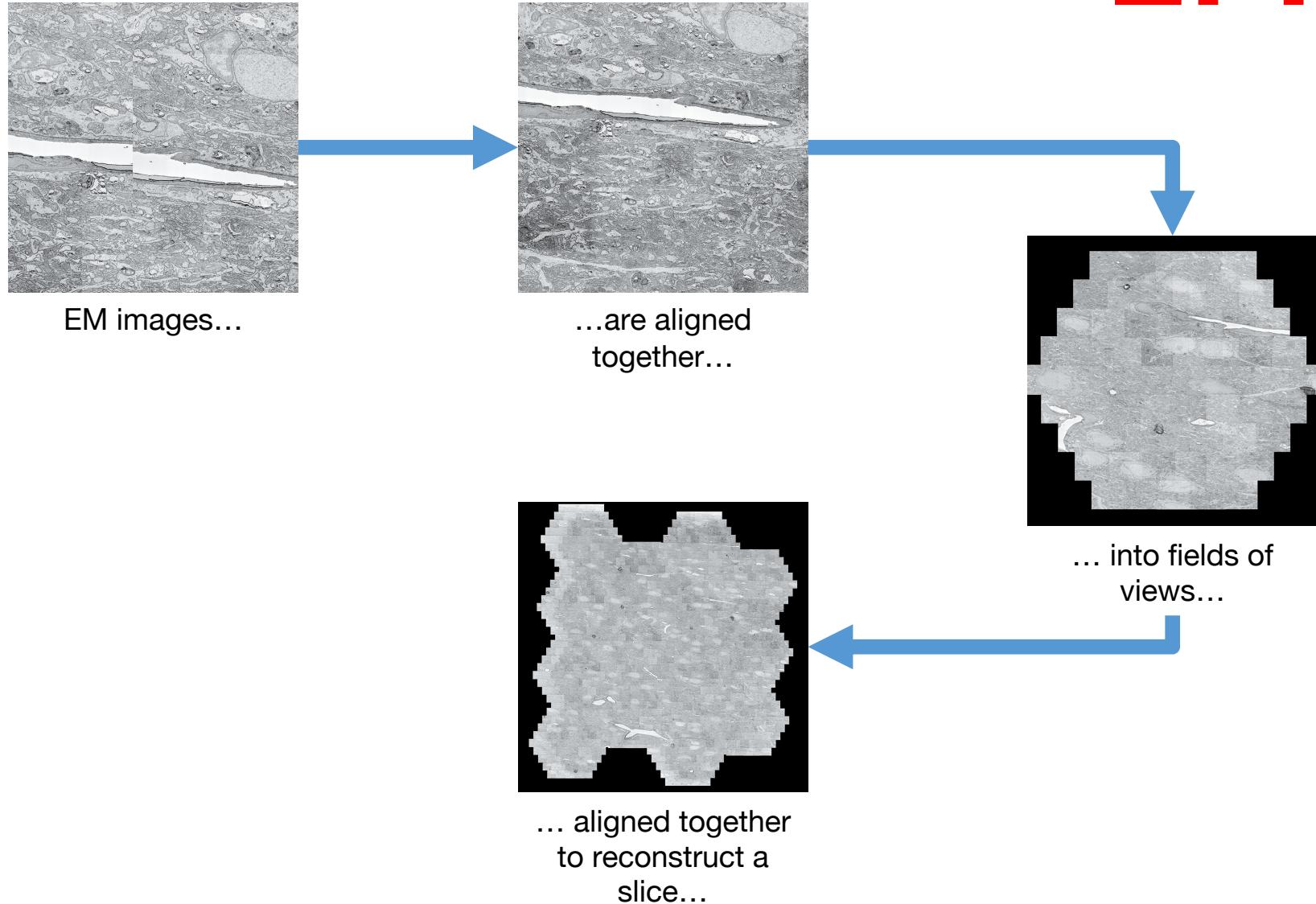
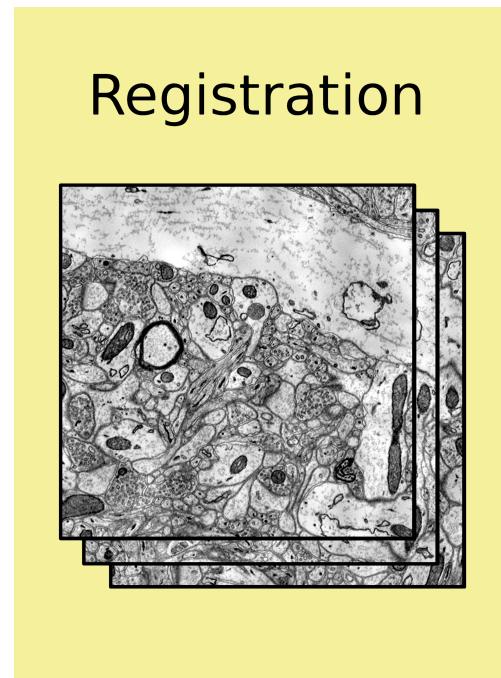
Registration



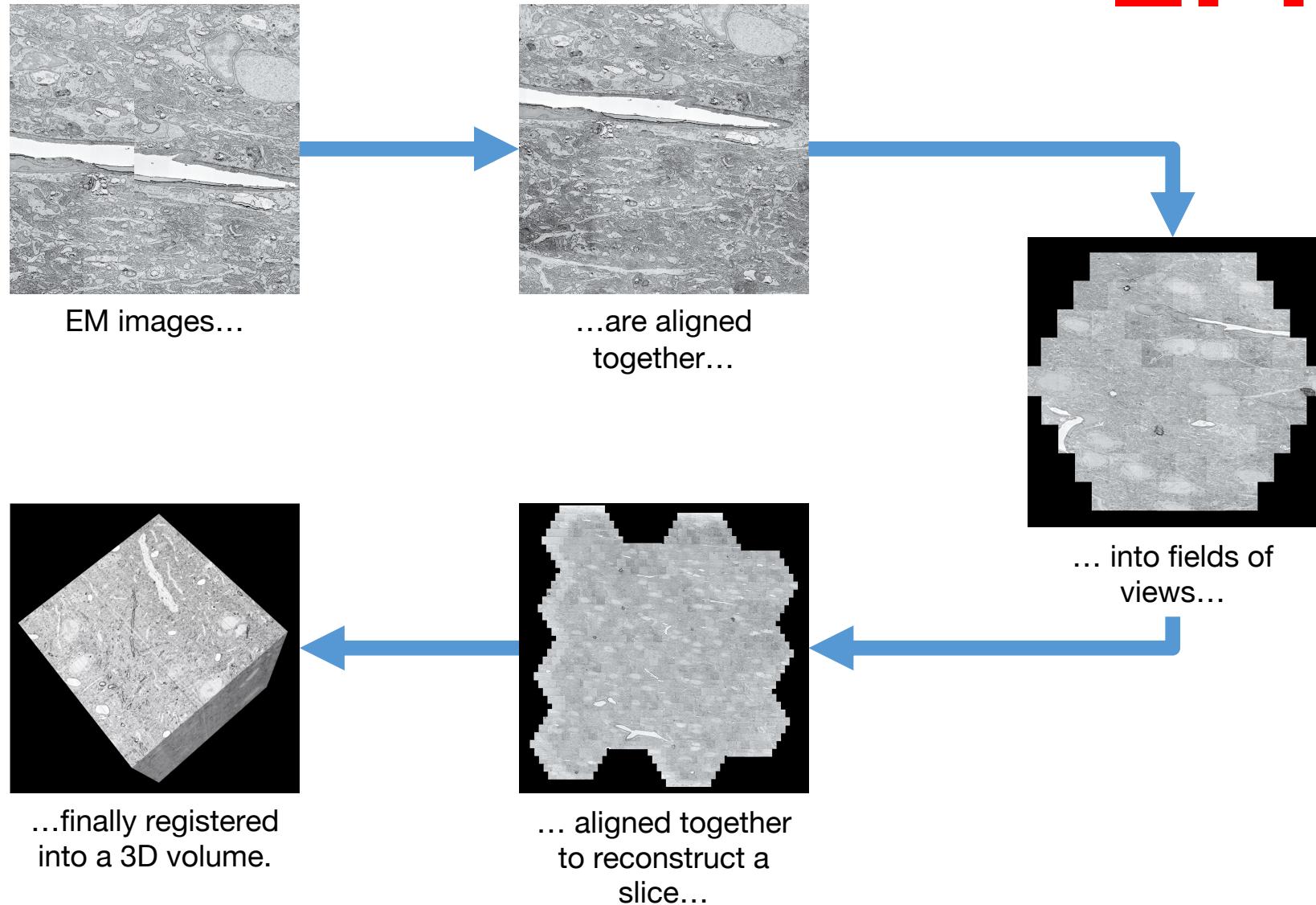
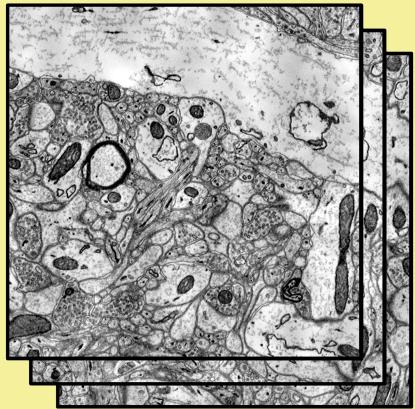
EM images...

...are aligned
together...

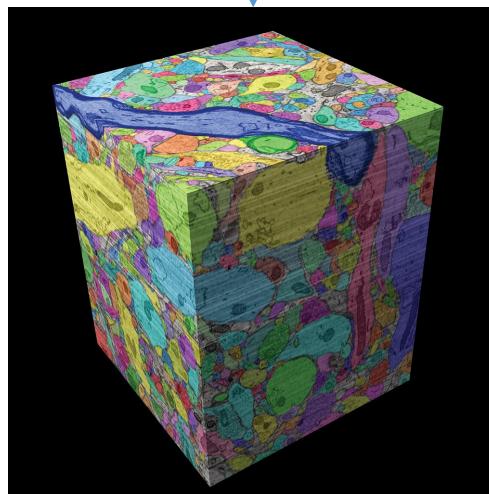
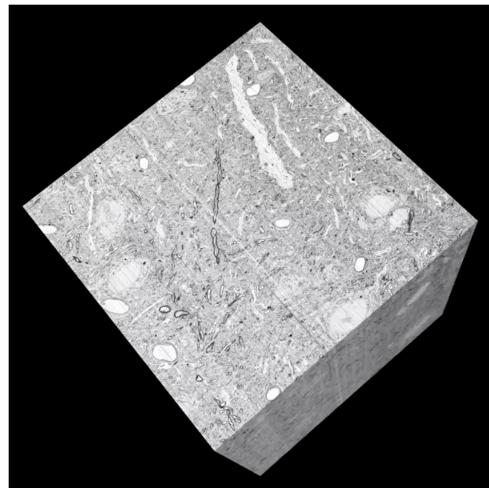
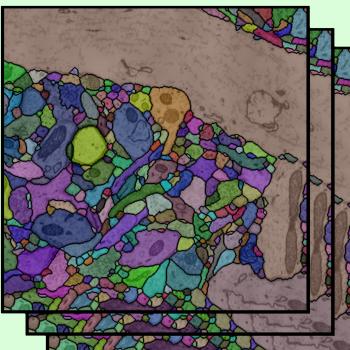




Registration



Segmentation

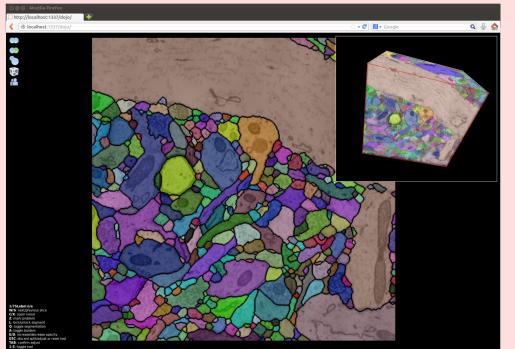


Apply automatic 3D segmentation to registered slices.

Segment neurons and synapses.

Computationally very costly.

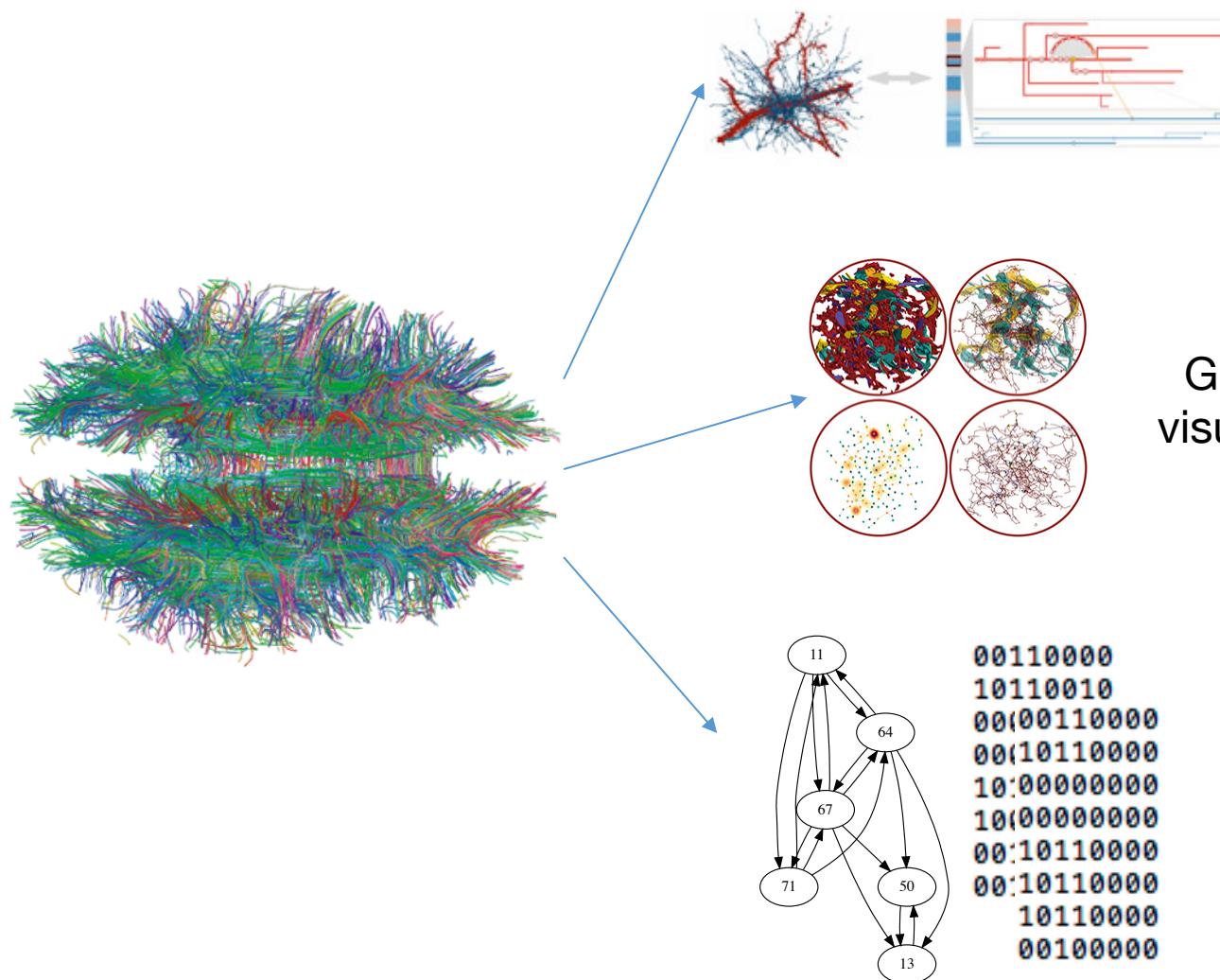
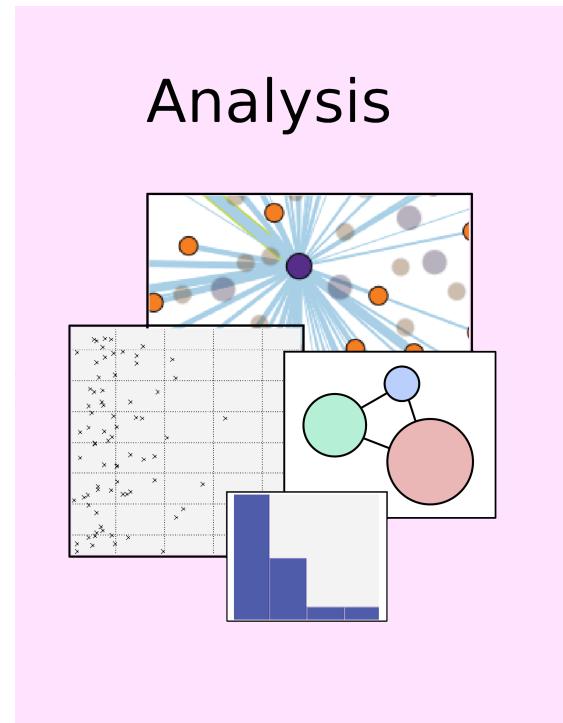
Proofreading



Proofreading is a necessary step in connectomics

Focuses on streamlining human involvement





Connectome
visualization

Glial cells
visualization

Motif Discovery