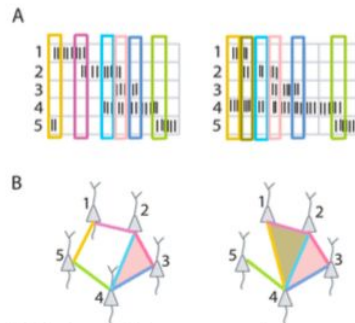**IUNI** INDIANA UNIVERSITY
**NETWORK SCIENCE INSTITUTE**

# Computational tools for handling simplicial complexes in real datasets
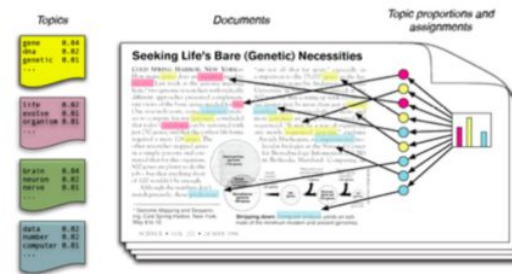
Alice Patania
Network Science Institute (IUNI), Indiana University

# Why?



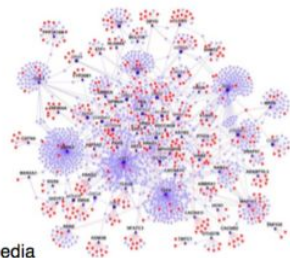Curto, Carina, and Vladimir Itskov. "Cell groups reveal structure of stimulus space." PLoS Comput Biol 4, no. 10 (2008): e1000205.
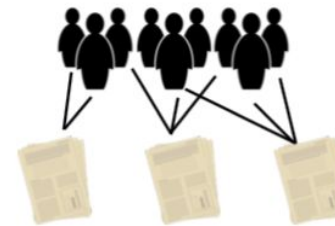
**brain networks**



Rhody, Lisa. "Topic modeling and figurative language." Journal of Digital Humanities 2, no. 1 (2012): 19-35.

**topic modelling**



Wikipedia
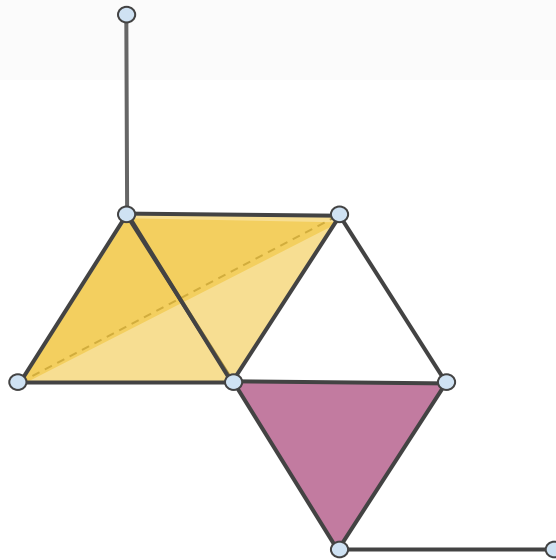
**protein interaction networks**



**collaboration networks**
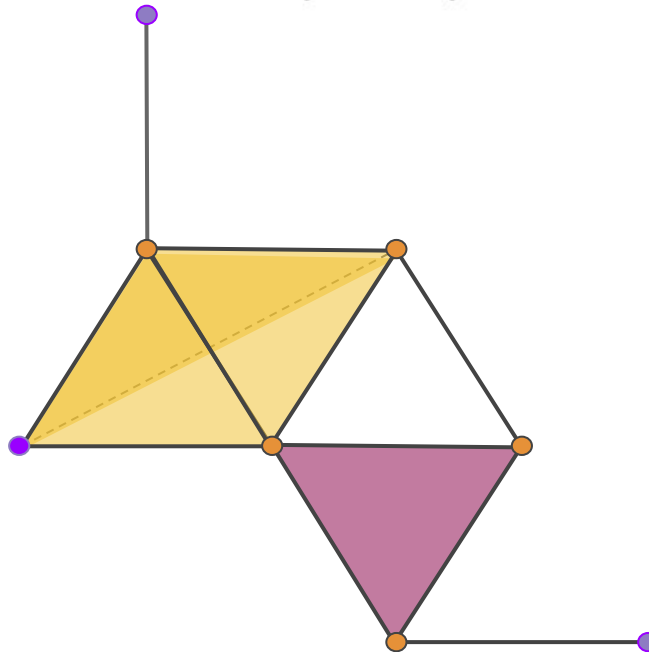
# Simplicial Complexes

A **simplicial complex** X is a collection of simplices such that:

- $\forall \sigma \in X$ its faces are still in X,
- $\forall \sigma, \tau \in X$, $\sigma \cap \tau$ is either the empty set or a face of both $\sigma$ and $\tau$.

# Simplicial Complexes

The **simplicial degree** of a node in a simplicial complex is the number
if maximal simplices under inclusion (**facets**) incident on the node.

# Simplicial Complexes 101

The **simplicial degree** of a node in a simplicial complex is the number if maximal simplices under inclusion (**facets**) incident on the node.
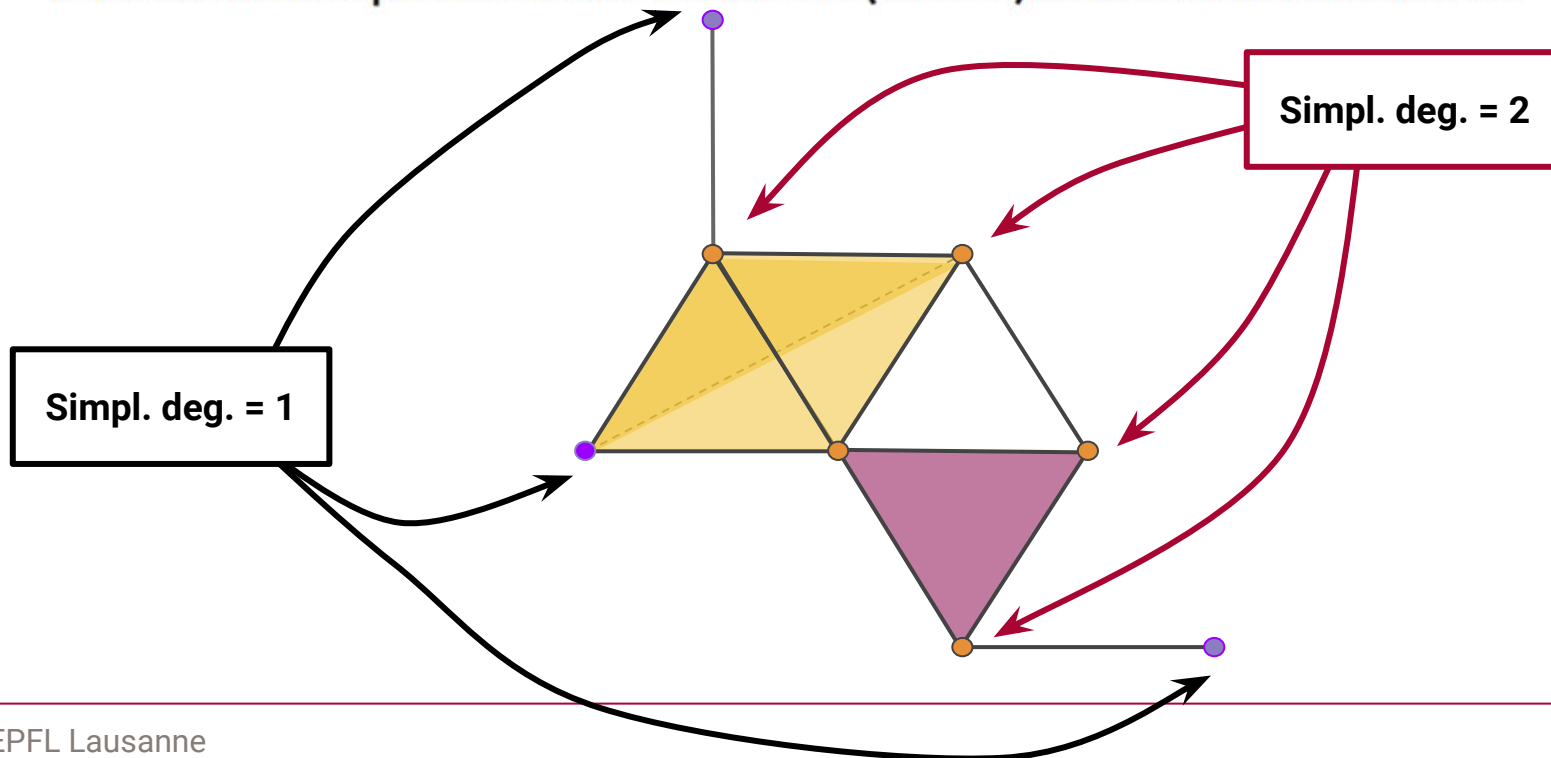


Simpl. deg. = 2

Simpl. deg. = 1
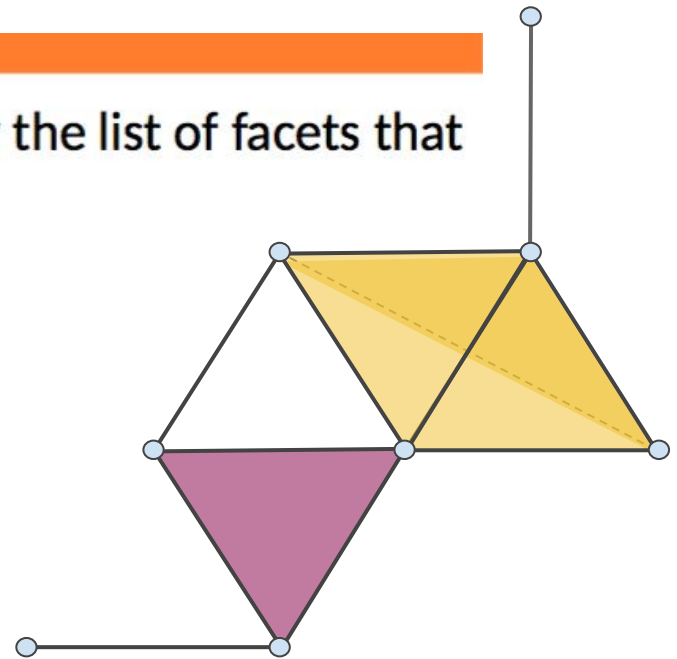
# Simplicial Complexes

The **simplicial degree** of a node in a simplicial complex is the number if maximal simplices under inclusion (**facets**) incident on the node.

A simplicial complex is completely described by the list of facets that belong to it.
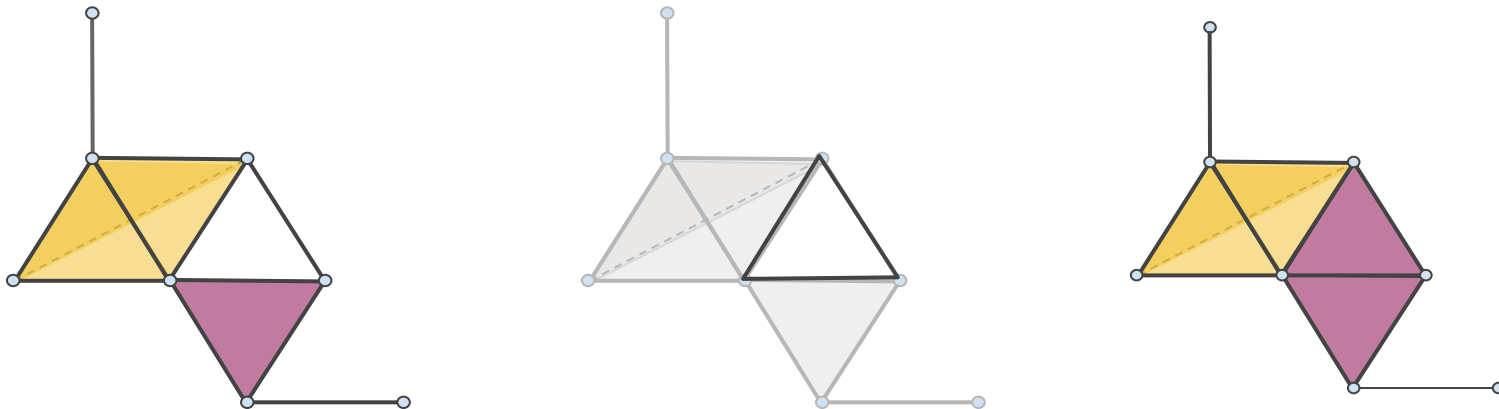
# Simplicial Complexes

## Simplicial homology

### What hinders application of homology to data?

**Representation:** it's difficult to find an optimal representative

**Memory and efficiency:** the algorithm for computing homology grows with the number of simplices in the complex

**Null model:** There is a lack of samplers that could easily be used in practice.

# Structure

- **Random Simplicial Complexes**
  - Simplicial Configuration Model

- **Reducing the complexity of homology computation**

- **1D-Homology and network communities**
  - arXiv case study

# Random Simplicial Complexes

# Simplicial Complexes
## Sampling

**Erdos-Renyi inspired:**

**Random pure simplicial complexes** [Linial-Meshulam (2006)]
**Random simplicial complexes** [Kahle (2009)]
**Multi-parameter random simplicial complexes** [Costa-Farber (2015)]
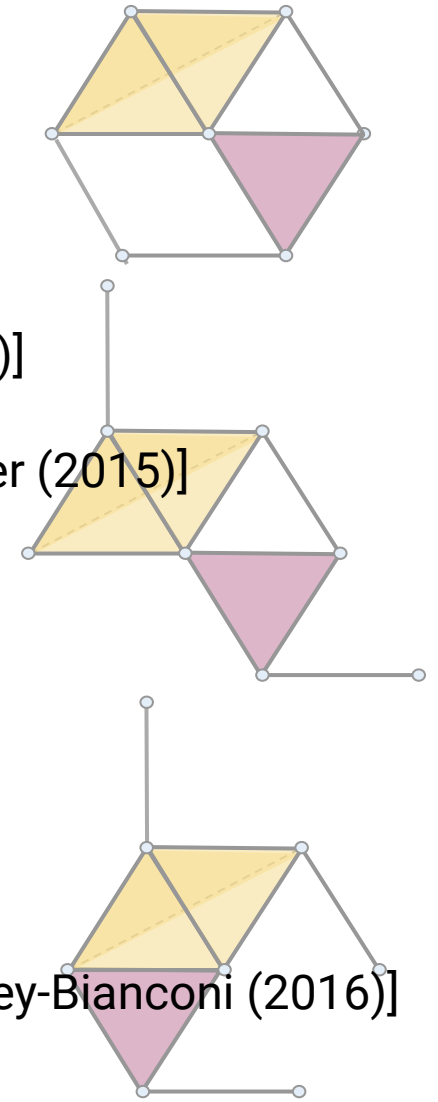
**Exponential random graphs inspired:**

**Exponential random simplicial complex** [Zuev et al. (2015)]

**Preferential attachment for simplicial complexes:**

**Network geometry with flavor** [Bianconi-Rahmede (2016)]

**Configuration model:**

**Configuration model for pure simplicial complexes** [Courtney-Bianconi (2016)]

# Configuration model

The **configuration model** is a generative model that creates a random graph with a fixed degree sequence.

**It implies the following are fixed?**

the number of nodes $n$

the number of edges in the network $m = \frac{1}{2} \sum_i k_i$

# Configuration model

The **configuration model** is a generative model that creates a random graph with a fixed degree sequence.

Suppose to have $n$ vertices with fixed degrees $k_i$ for $i = 1, \ldots, n$, the random graph is constructed in the following way.

1. Each vertex $i$ is provided with $k_i$ edge 'stubs', there are therefore $\sum_i k_i = 2m$ stubs.

2. Uniformly at random two stubs are chosen and an edge is created connecting the two of them, until no free stubs are left in the graph.

# Bipartite graphs and simplicial complexes

## Theorem

Let G be a bipartite graph with vertex sets $\{F, V\}$, $G_V$ its one-mode projections onto the vertex set V.

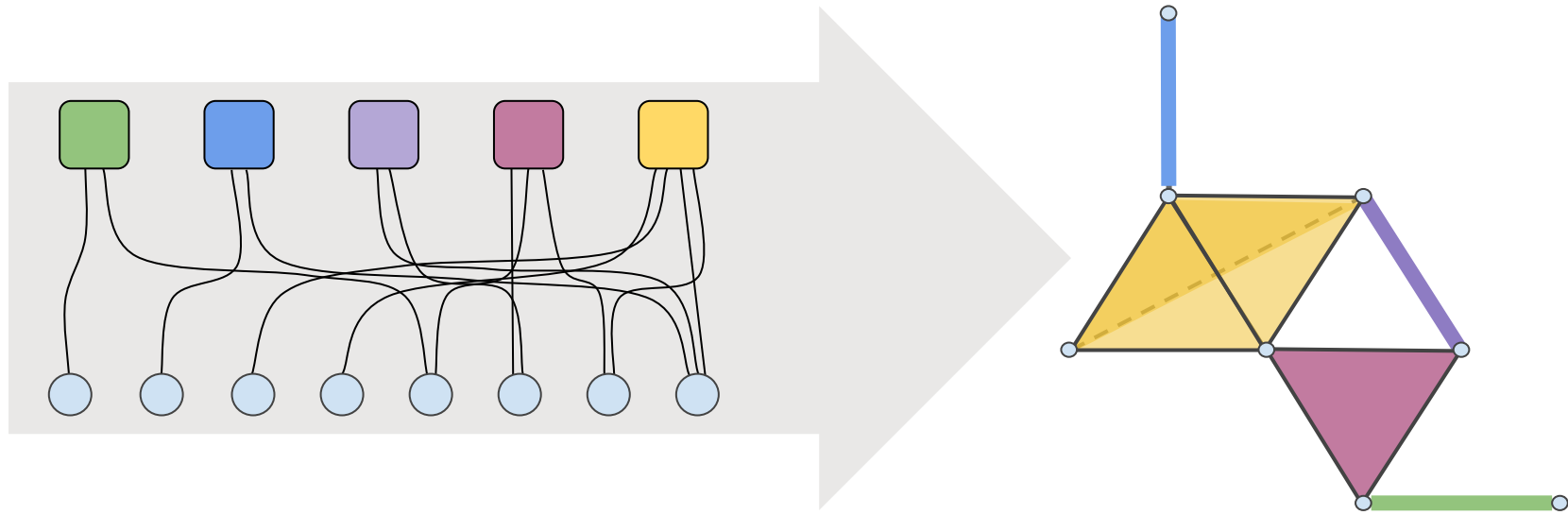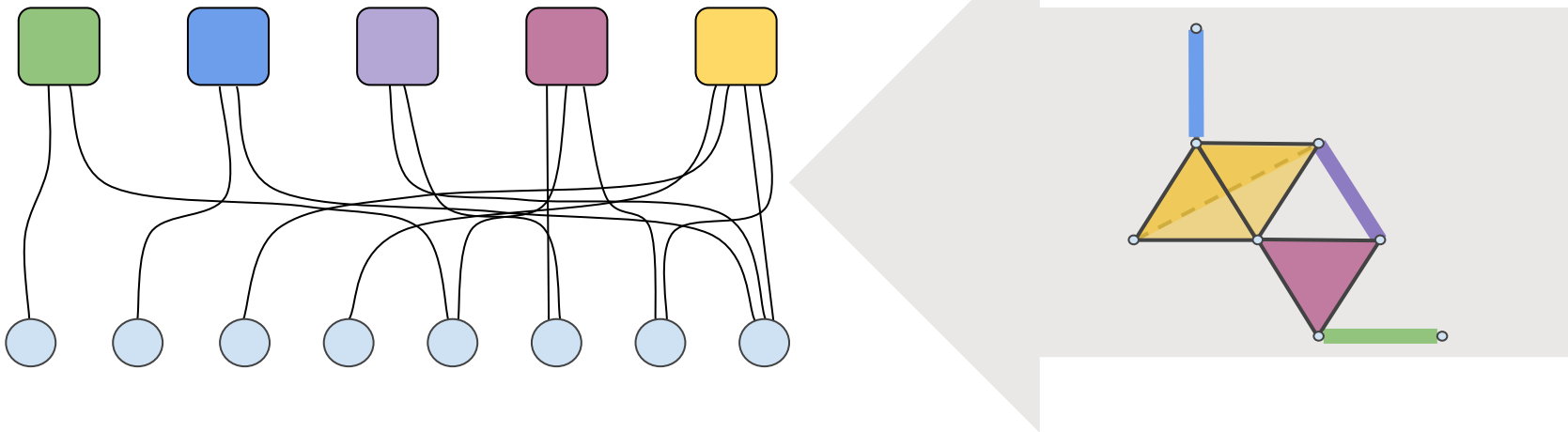Then it exists a simplicial complex $\Sigma$ whose underlying graph is $G_V$.

The neighbours $\mathcal{N}(f_i)$ of $f_i$ are the vertices that form the maximal simplex $f_i$, for each $i$, or equivalently, the neighbours $\mathcal{N}(v_i)$ of vertex $v_i$ are the facets in which node $v_i$ appears.

# Bipartite graphs and simplicial complexes

## Theorem

Let G be a bipartite graph with vertex sets $\{F, V\}$, $G_V$ its one-mode projections onto the vertex set V.

Then it exists a simplicial complex $\Sigma$ whose underlying graph is $G_V$.

# Bipartite graphs and simplicial complexes
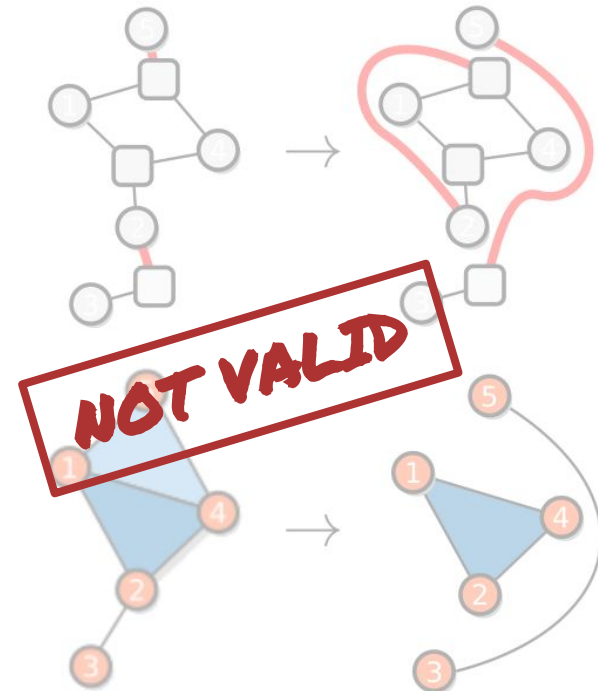
# Bipartite graphs and simplicial complexes

**Idea:**

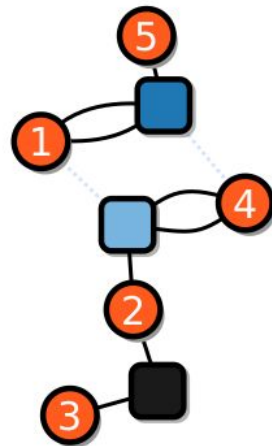Use the configuration model for bipartite graphs and the maps to construct a sampling method for SCM.

# Bipartite graphs and simplicial complexes

**Idea:**

Use the configuration model for bipartite graphs and the maps to construct a sampling method for SCM.
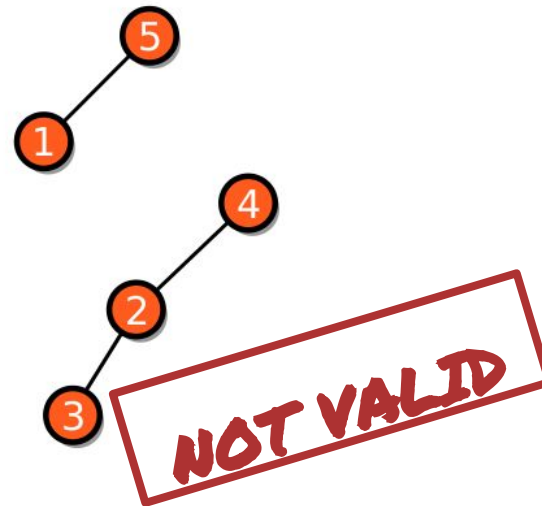
# Adding constraints

First constraint: **No multi-edges**

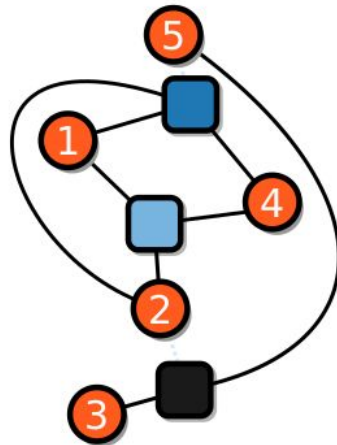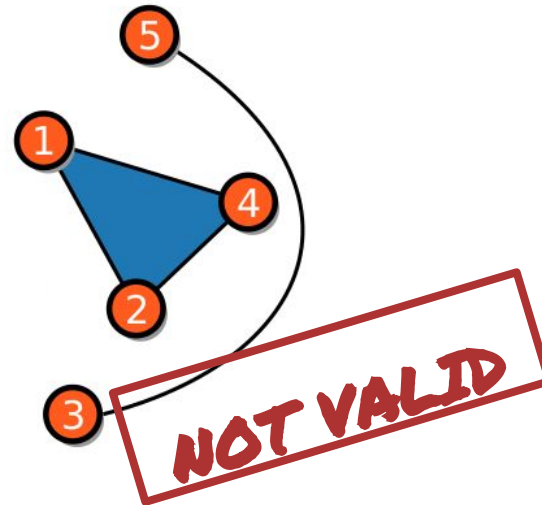Multi-edges decrease the size of the maximal simplices.

# Adding constraints

Second constraint: **No included neighborhoods**

Included neighborhoods violate the maximality assumption of the facets.
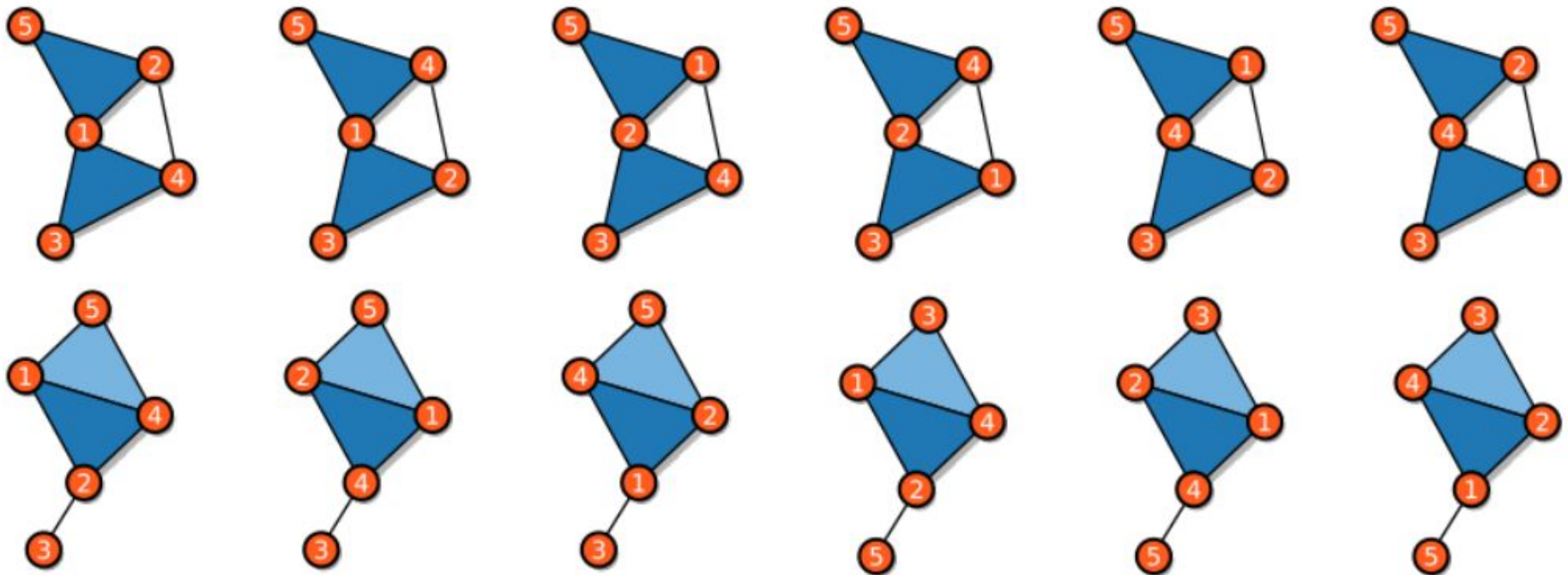


Bipartite graph    Simplicial complex

NOT VALID

# Adding constraints

Constraints:
**No multi-edges**
**No included neighborhoods**

Then the acceptable configurations for the toy example are the following:
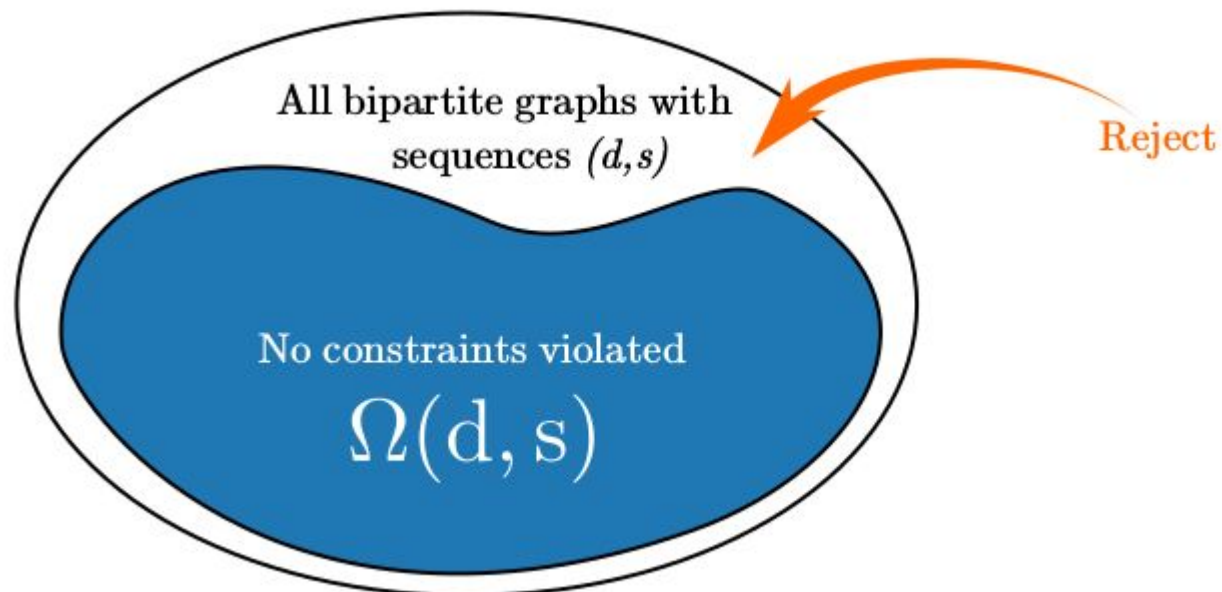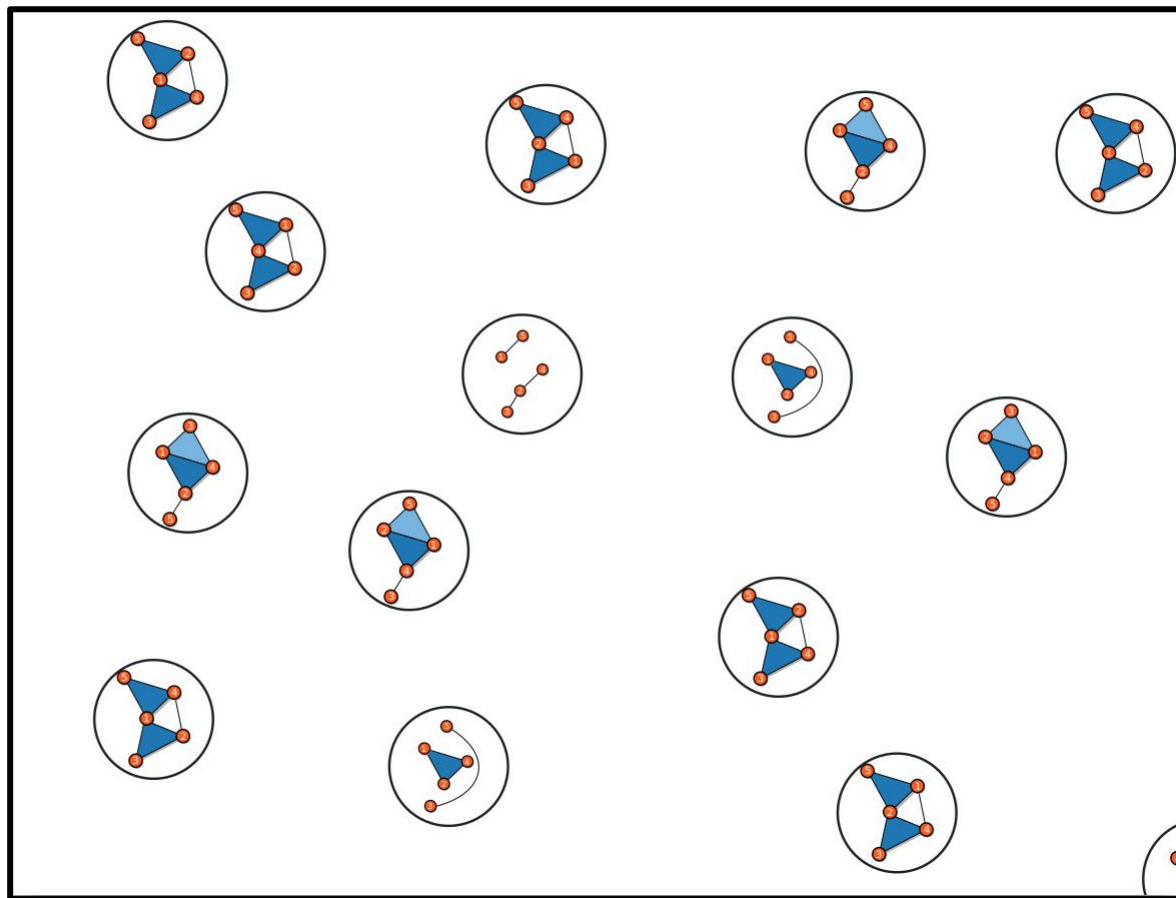
# Adding constraints

**!** **Problem with rejection sampling**: Far too many rejections!
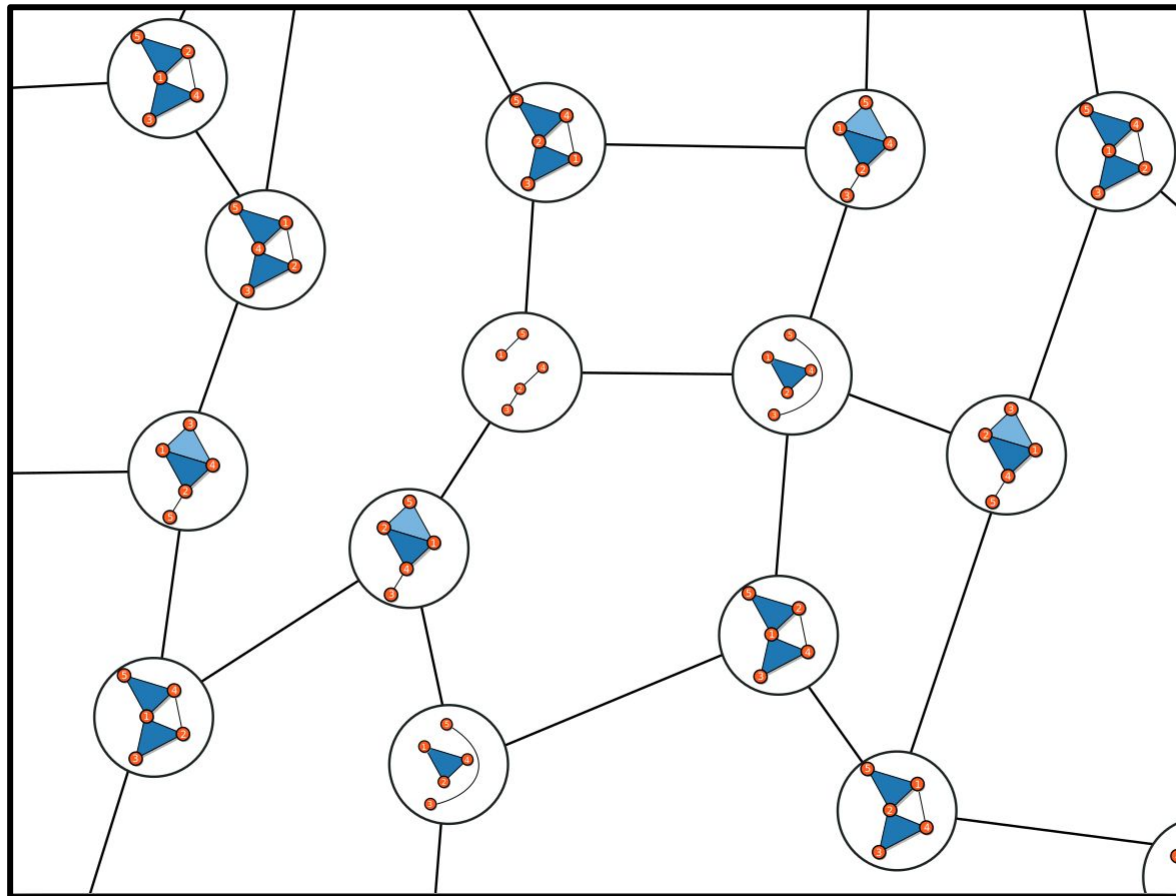
Loose upper bound :

$$\Pr[\text{reject}] > \exp[-0.5(\langle d^2\rangle/\langle d\rangle - 1)(\langle s^2\rangle/\langle s\rangle - 1)]$$
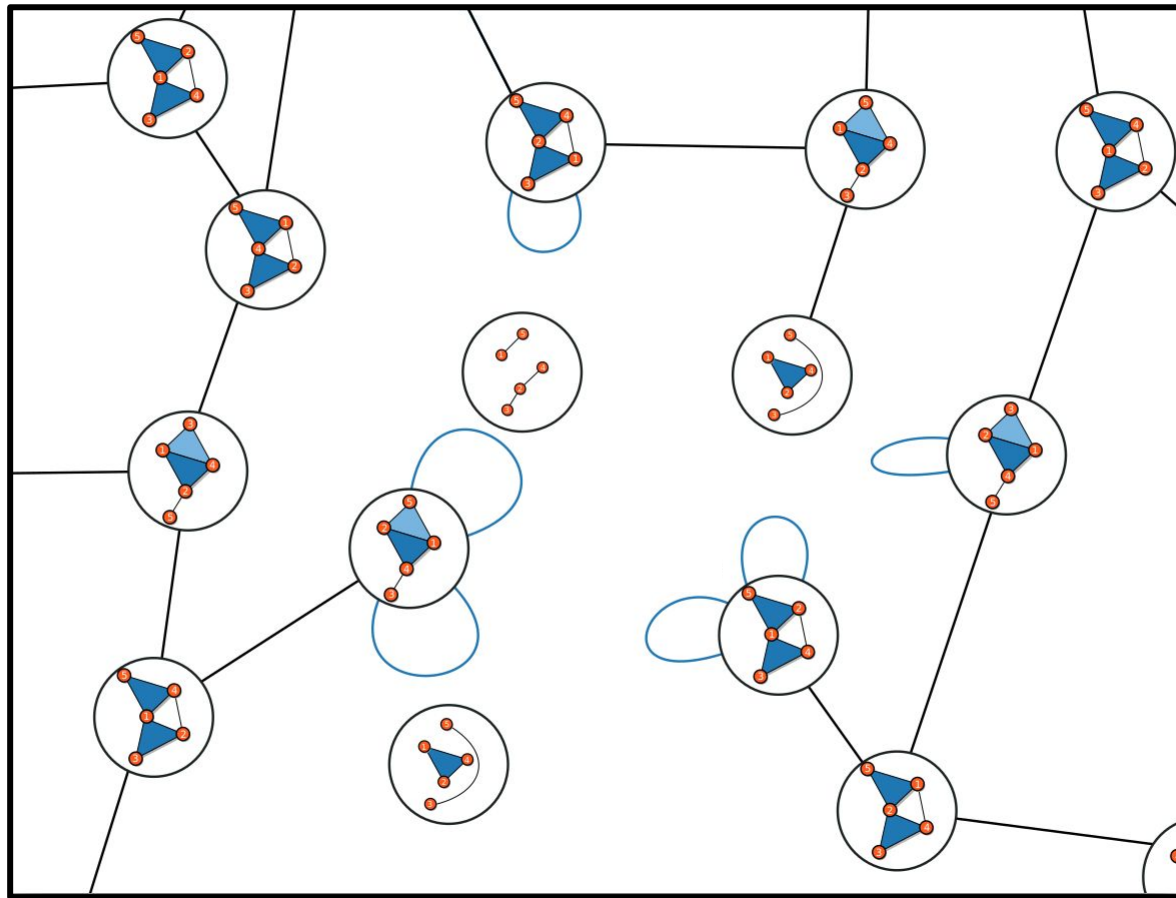


All bipartite graphs with
sequences $(d,s)$

Reject

No constraints violated

$\Omega(d, s)$

# Markov Chain Monte Carlo sampling

# Markov Chain Monte Carlo sampling

# Markov Chain Monte Carlo sampling

# MCMC sampling: The details

**Move set**

1. Pick L~P random edges in bipartite graph
   P can be arbitrary, we use $Pr[L = l] = \exp[\lambda l]/Z$

2. Rewire edges. If multi-edge or included neighbors, reject.

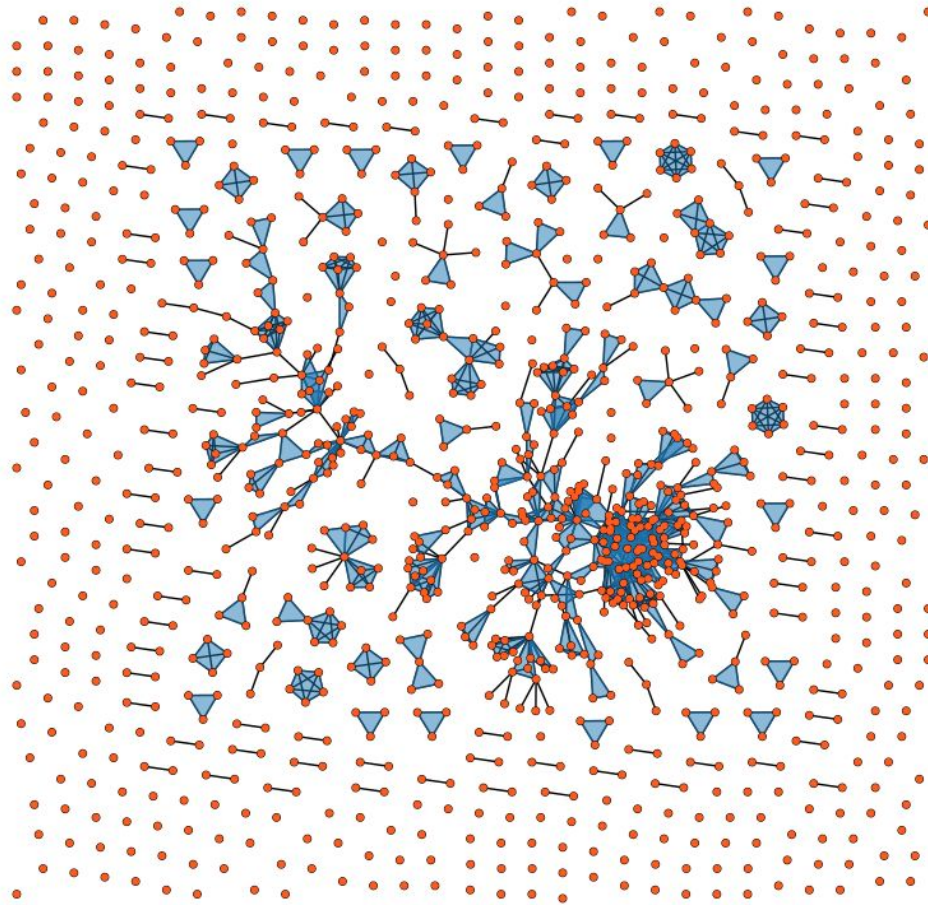Similar to [Miklós−Erdős−Soukup, Electron. J. Combin., 20, (2013)]

- MCMC is uniform over $\Omega(d, s)$
- Move set yields aperiodic chain
- Move set connects the space

# Results - True sustems

## Disease regulation dataset
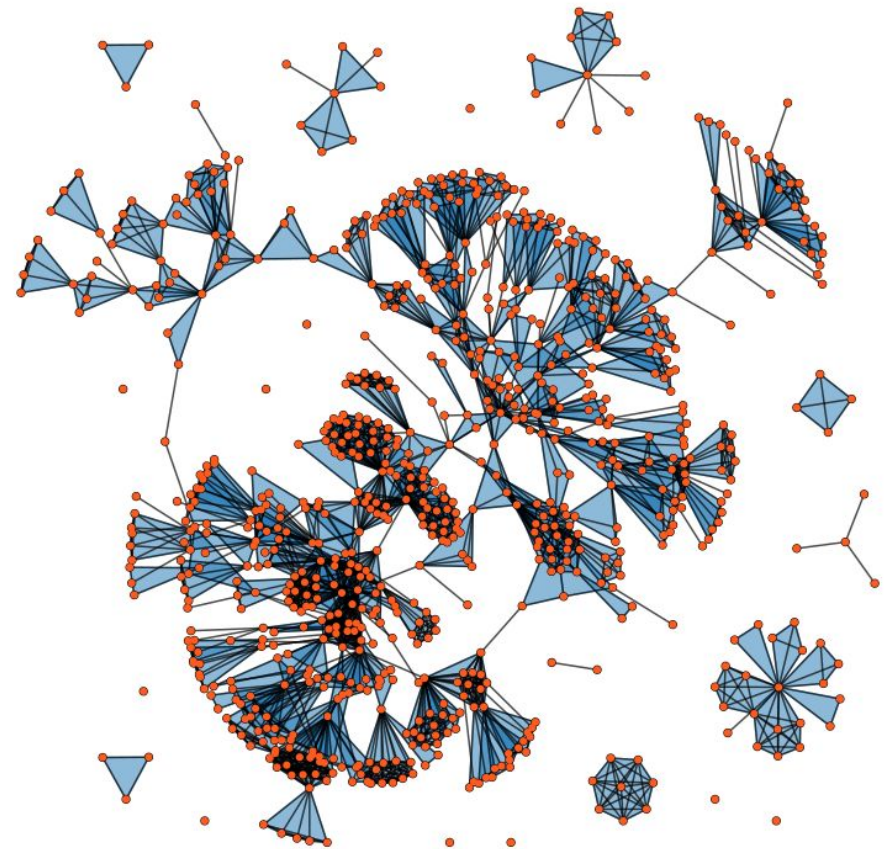(facets : *genes*, nodes : *human diseases*)
[Goh et al., PNAS, **104**, (2007)]

## Crimes in St-Louis (true system)
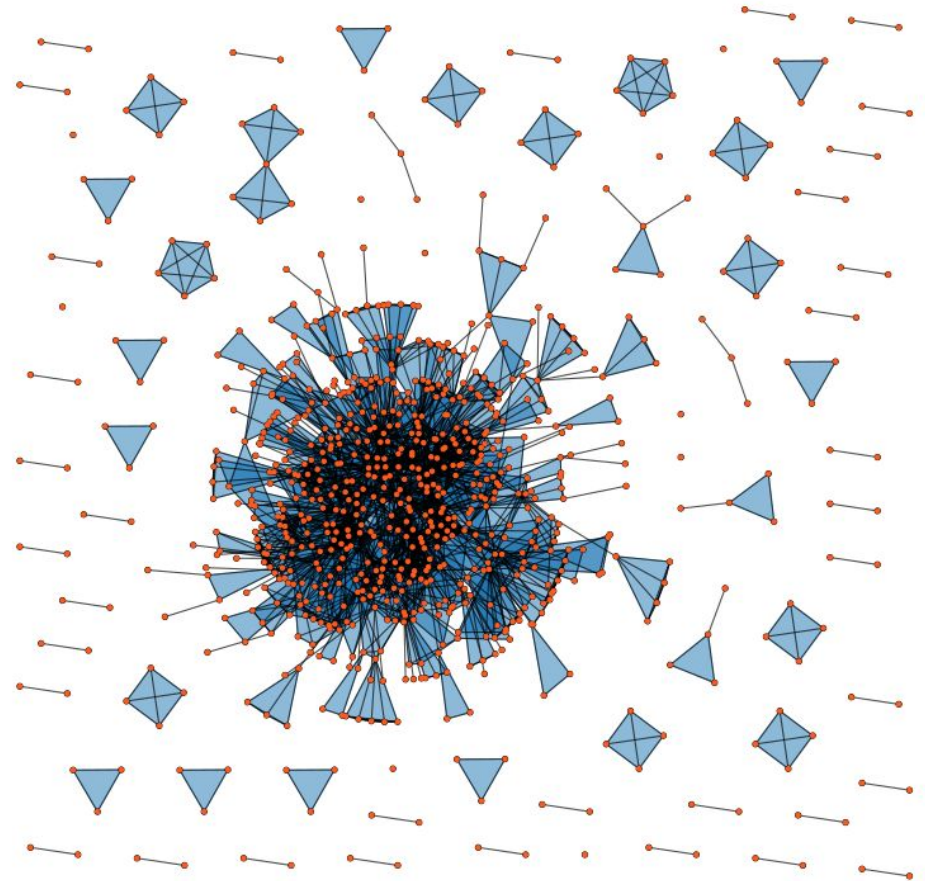(facets : *people*, nodes : *crimes*)
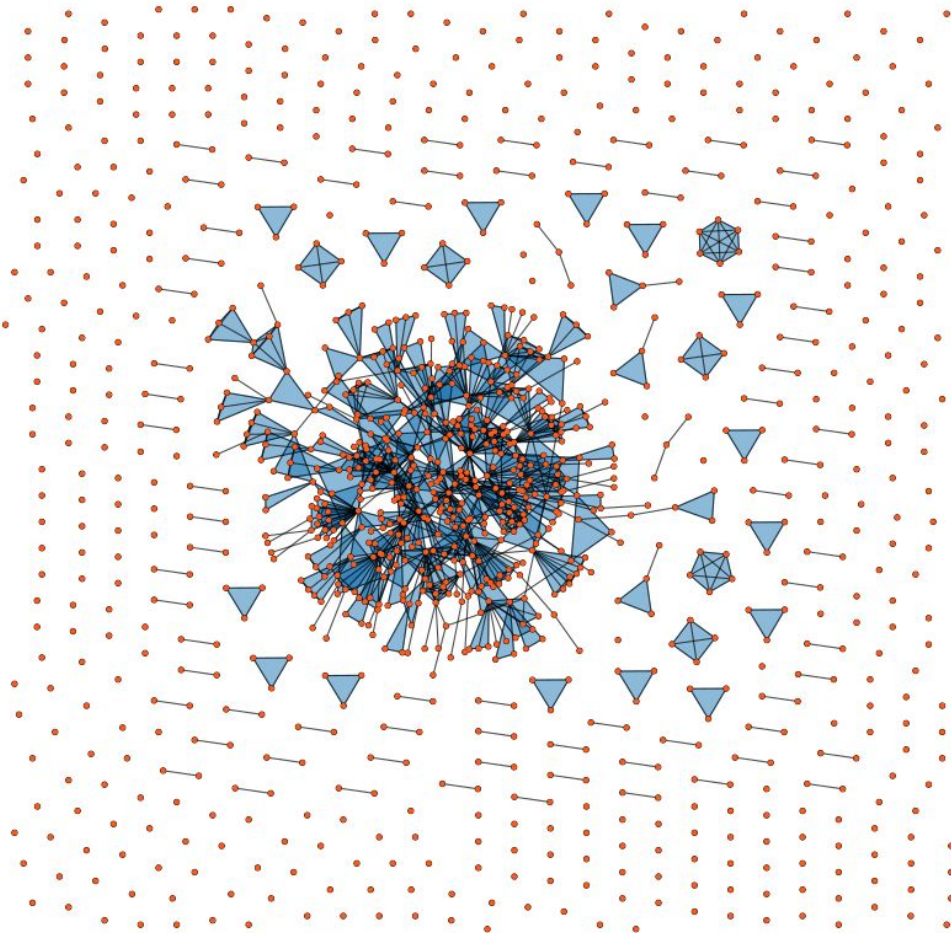[Rosenfeld et al., (1991)]

# Results - Random instances

Disease regulation dataset (random instance)
(facets : *genes*, nodes : *human diseases*)
[Goh et al., PNAS, **104**, (2007)]
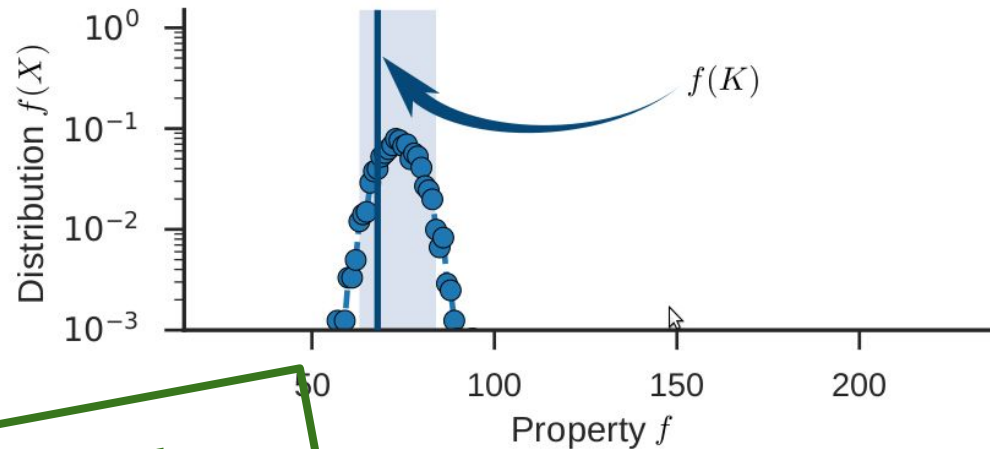
Crimes in St-Louis (random instance)
(facets : *people*, nodes : *crimes*)
[Rosenfeld et al., (1991)]

# Concept for a null model

**Null model**

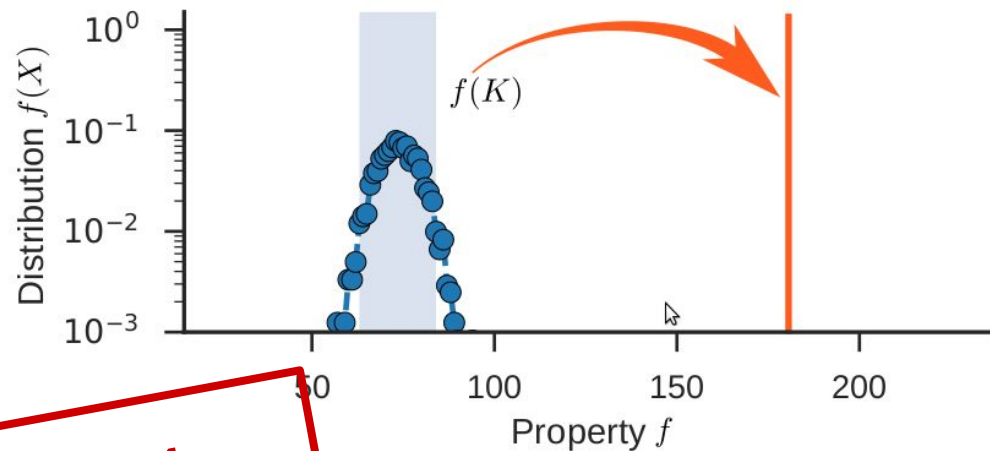Is the quantity **f(X)** close to **f(K)** for random simplicial complexes
**X~SCM[d(K), s(K)]** ?



$$\text{Pr}[|f(K) - f(X)| < 1] \approx 1$$

**K** is typical, the local quantities (**d**,**s**) explain **f**.
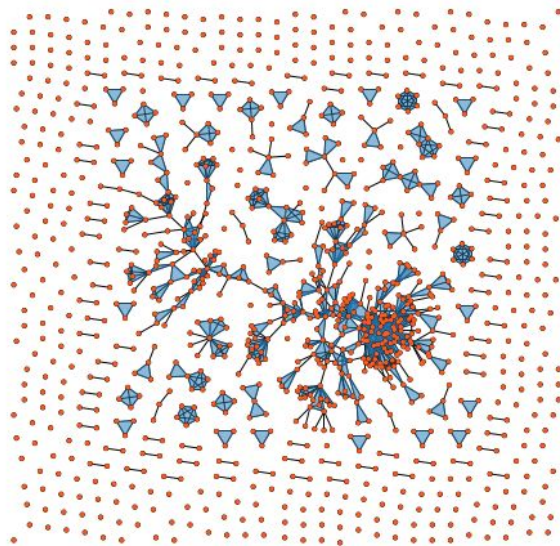
# Concept for a null model

**Null model**

Is the quantity f(X) close to f(K) for random simplicial complexes
X~SCM[d(K), s(K)] ?



$$Pr[|f(K) - f(X)| < ] \ll 1$$

EPFL Lausanne                    K is atypical, K is organized beyond the local scale.
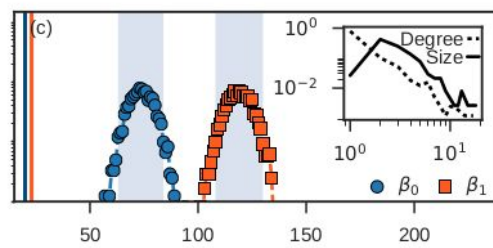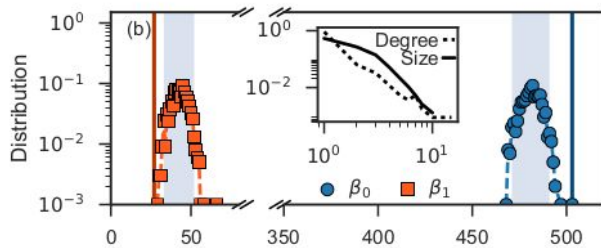
# Results on Betti numbers of real data sets



Diseases        Crime        Pollinators

EPFL Lausanne
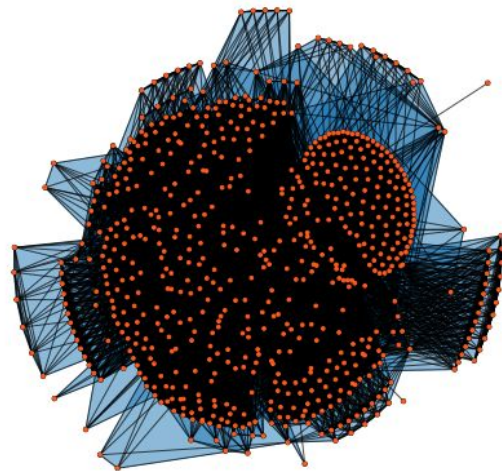
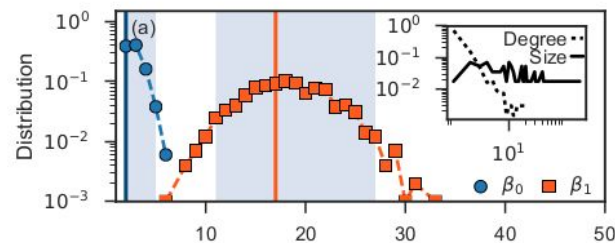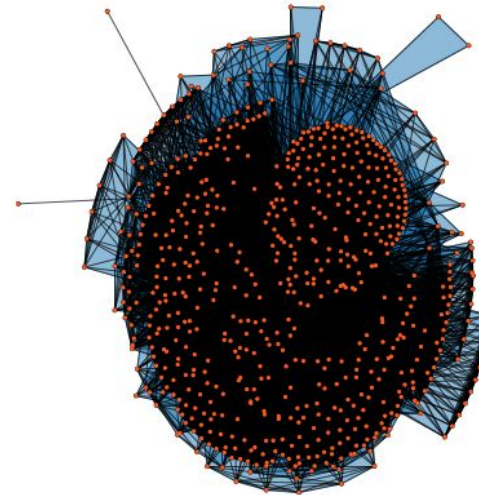# Results on Betti numbers of real data sets



Pollinators (*real*)  Pollinators (*random*)

# Reducing the algorithmic complexity

# Simplicial Complexes
## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{max(s)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.

# Simplicial Complexes
## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{\max(s)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.
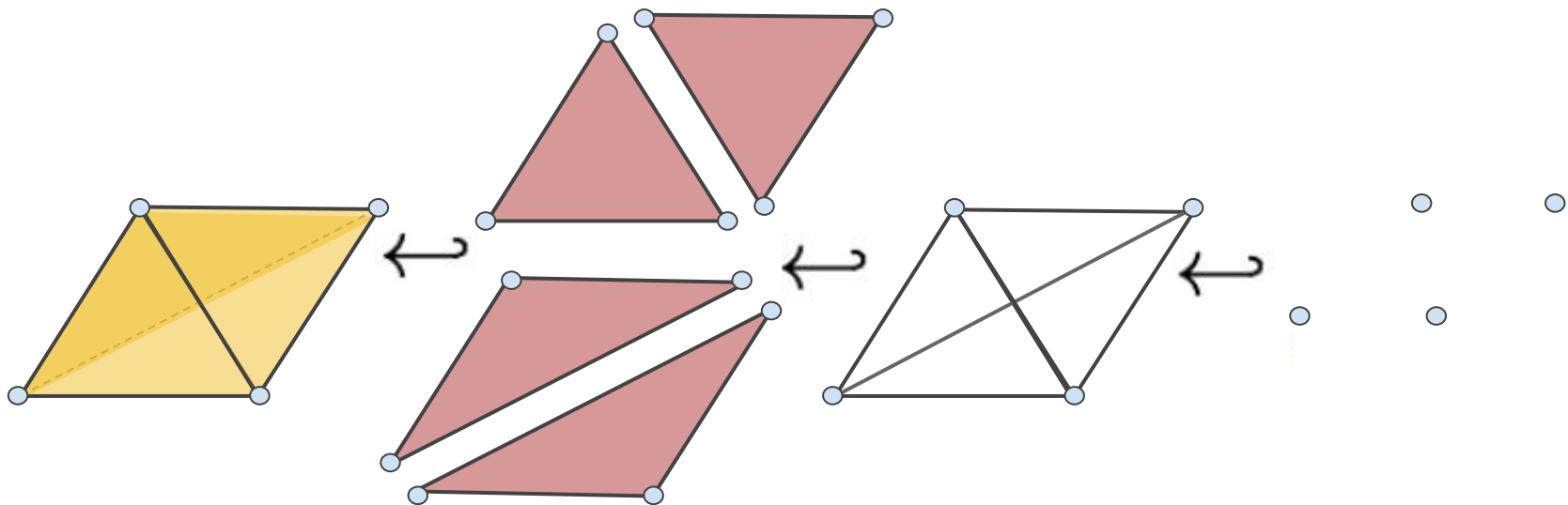


REDUCTION

# Simplicial Complexes
## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{max(s)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.



EPFL Lausanne

# Simplicial Complexes
## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{\max(s)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.
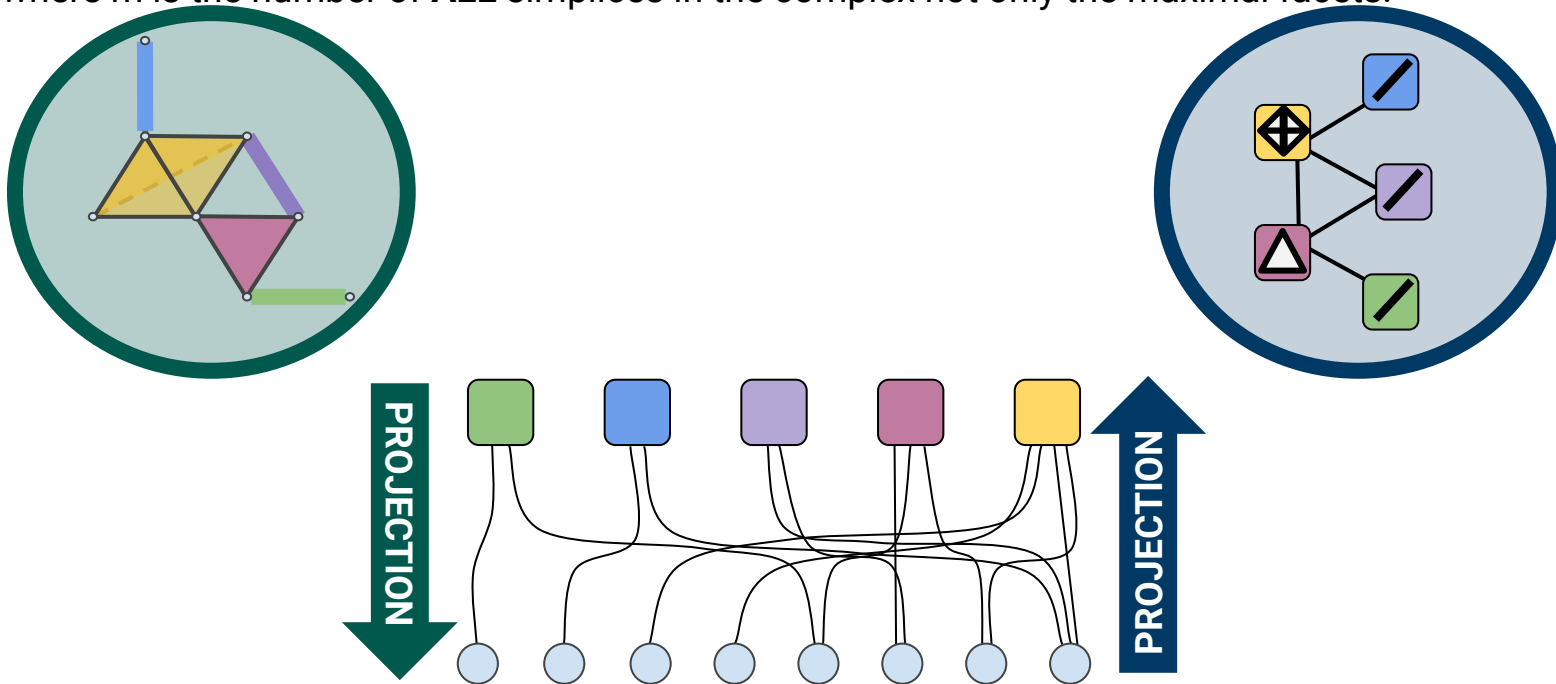
This method does **NOT** guarantee automatically that the new complex will have fewer simplices.

PROJECTION

PROJECTION

# Simplicial Complexes
## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{\max(d)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.
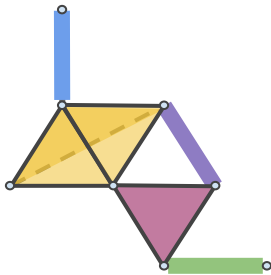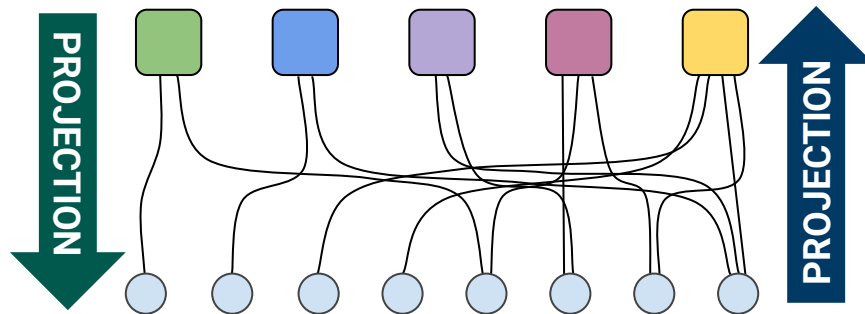
This method does **NOT** guarantee automatically that the new complex will have fewer simplices.

PROJECTION

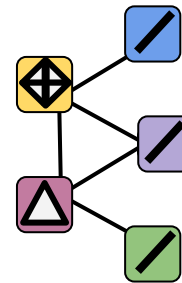PROJECTION

# Simplicial Complexes

## homological hiccups

The computational complexity of homology is $O(m^3) > O([2^{\max(d)}]^3)$

where m is the number of **ALL** simplices in the complex not only the maximal facets.
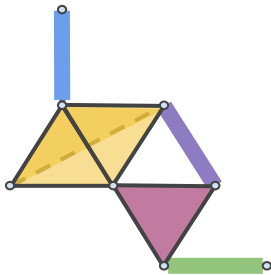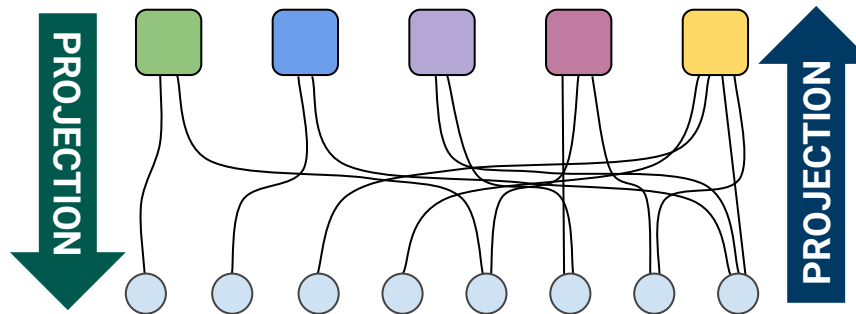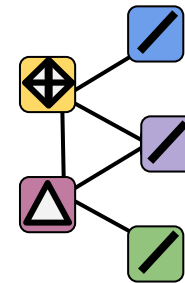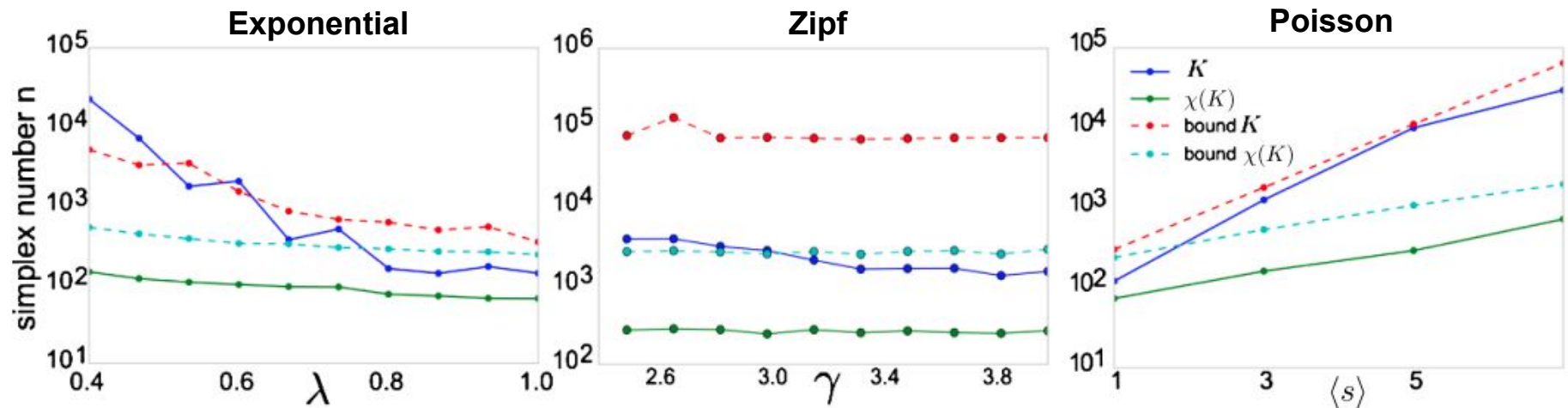


Simplicial Configuration Model, J-G Young et. al. 2017 PRE

# 1D Homology and communities

## Empirical proof of concept

# The data set

The data span 9 years, from 2007 to 2016, and are split according to the 18 major categories of arXiv.

This major categories correspond to different thematic areas and thus can be used as rough representative of different scientific fields.

Notice: Due to arXiv's history, there is a bias toward mathematical and physical topics.

Cornell University Library

arXiv.org

Set of authors

List of categories

Date of publication

# The data set

The data span 9 years, from 2007 to 2016, and are split according to the 18 major categories of arXiv.

This major categories correspond to different thematic areas and thus can be used as rough representative of different scientific fields.

Notice: Due to arXiv's history, there is a bias toward mathematical and physical topics.



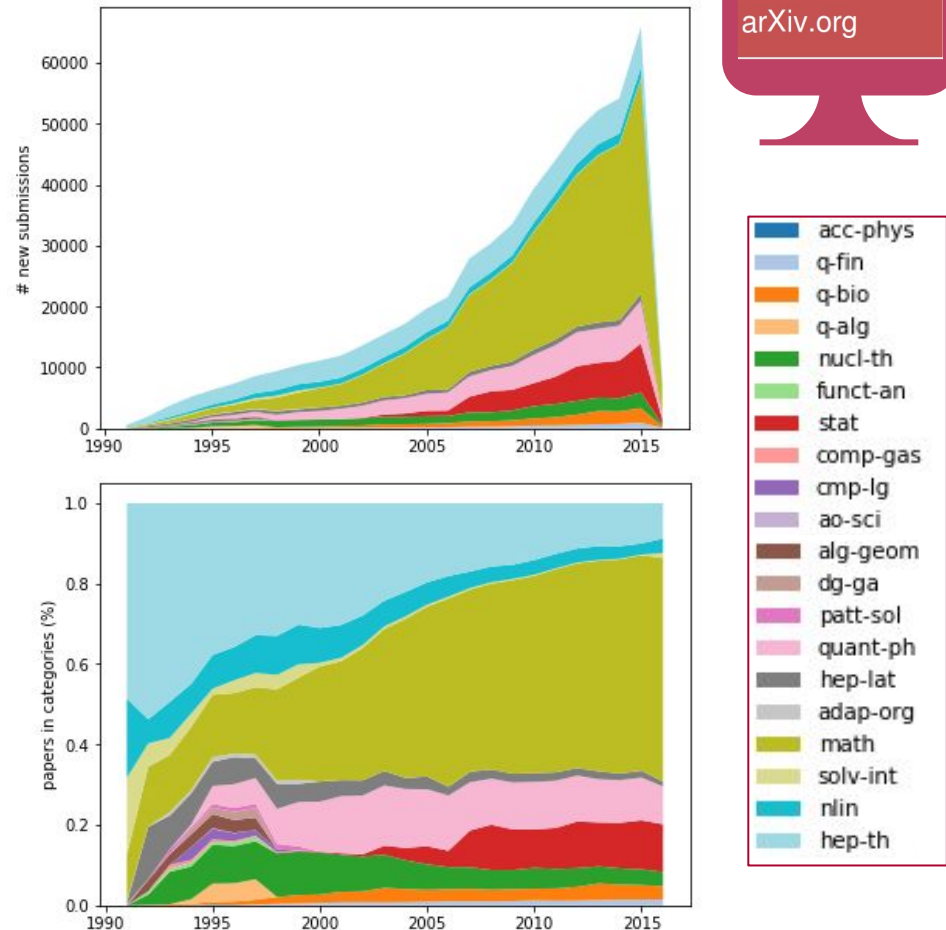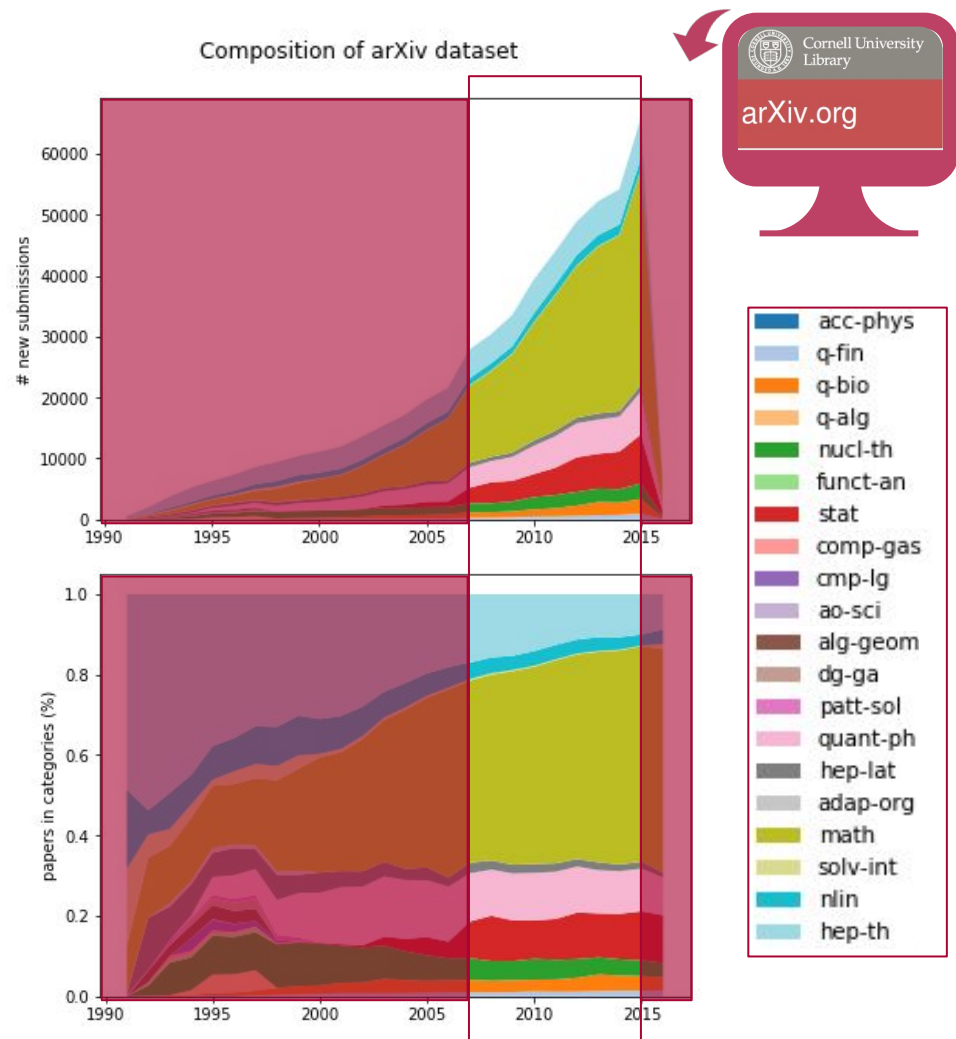Composition of arXiv dataset

# The data set

The data span 9 years, from 2007 to 2016, and are split according to the 18 major categories of arXiv.

This major categories correspond to different thematic areas and thus can be used as rough representative of different scientific fields.

Notice: Due to arXiv's history, there is a bias toward mathematical and physical topics.

# Facets size and simplicial degree

Assess commonalities in the statistical properties of the different categories
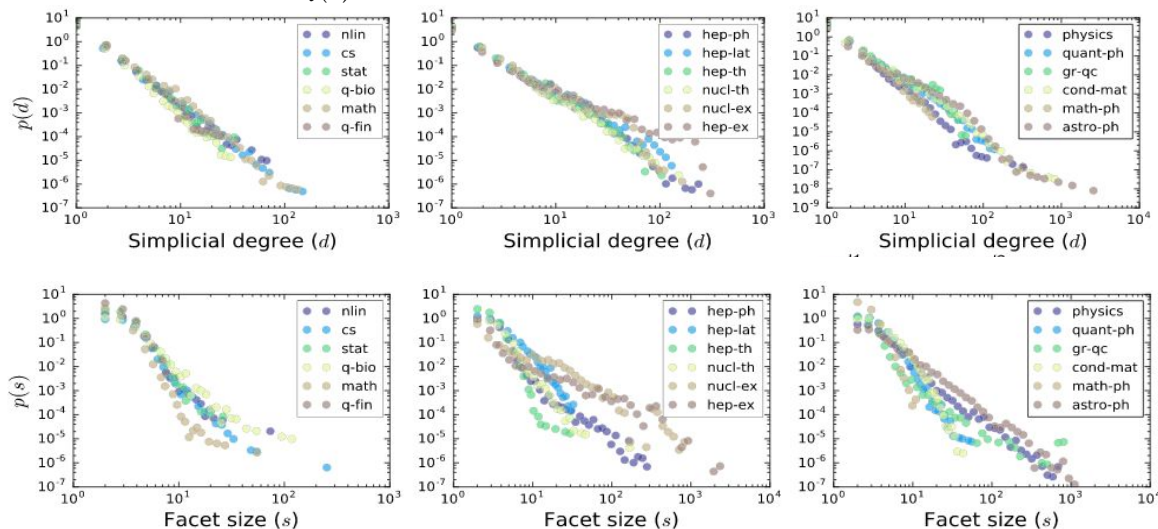
**Jensen-Shannon Divergence**

$$JSD(P,Q) = \tfrac{1}{2}D_{KL}(P \mid M) + \tfrac{1}{2}D_{KL}(Q \mid M)$$

where:

$$M = P + Q$$

$$D_{KL}(P \mid Q) = -\sum_x P(x) \log \frac{P(x)}{Q(x)}$$
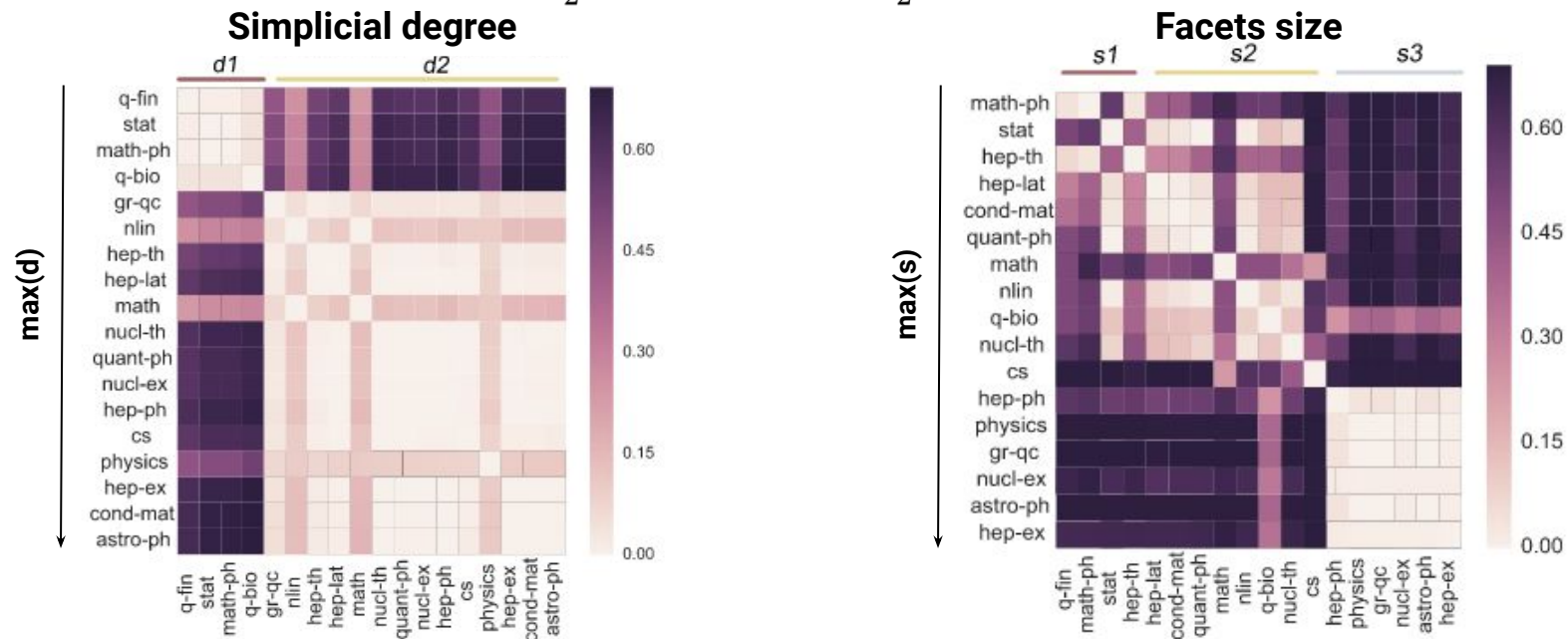


EPFL Lausanne

# Facets size and simplicial degree

Assess commonalities in the statistical properties of the different categories

**Jensen-Shannon Divergence**

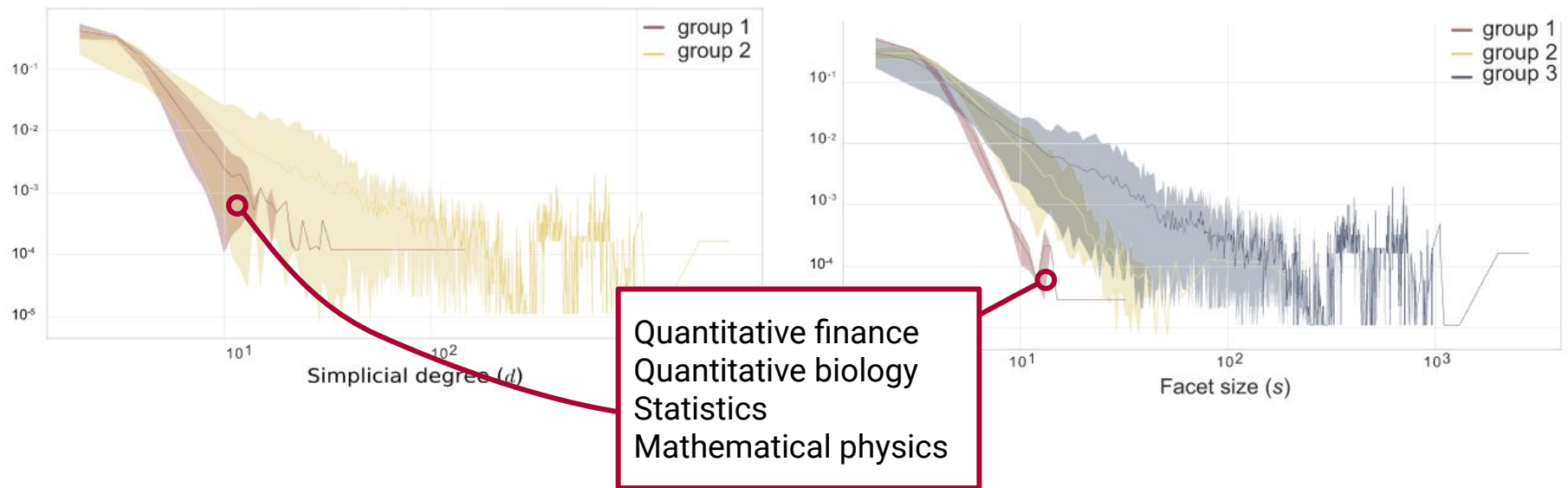$$JSD(P, Q) = \tfrac{1}{2}D_{KL}(P \mid M) + \tfrac{1}{2}D_{KL}(Q \mid M)$$

**Simplicial degree**                                              **Facets size**



EPFL Lausanne

# Facets size and simplicial degree

Examples for the biggest connected component for each group.

# Facets size and simplicial degree
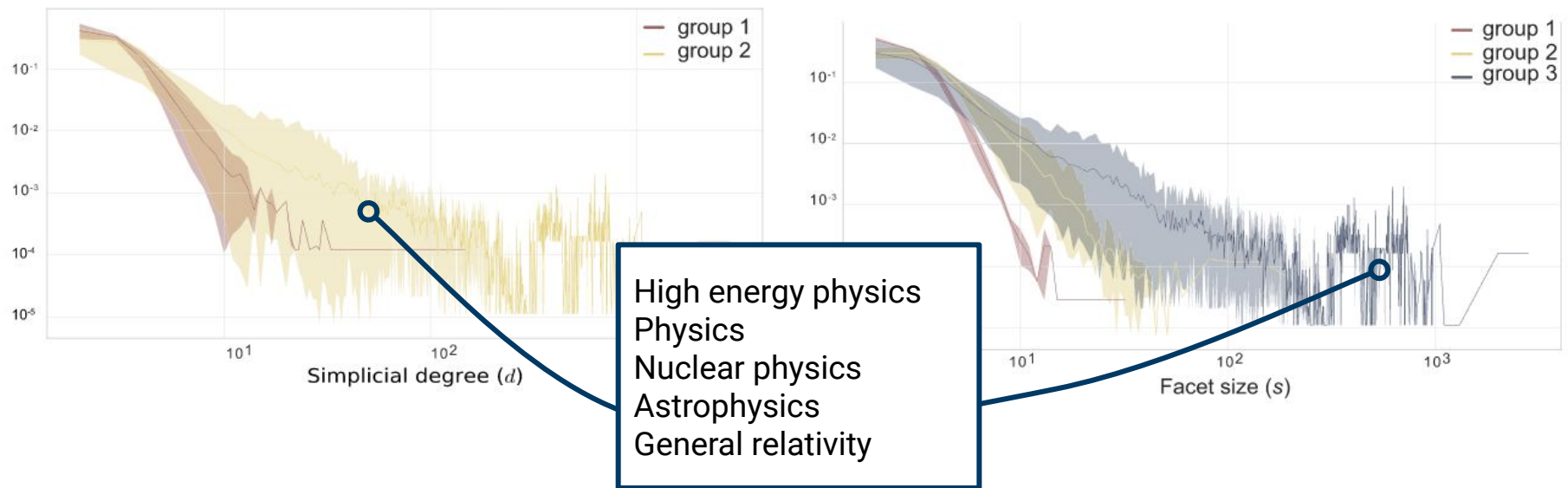
**Simplicial degree**

**Facets size**



Quantitative finance
Quantitative biology
Statistics
Mathematical physics

# Facets size and simplicial degree

**Simplicial degree**

**Facets size**



High energy physics
Physics
Nuclear physics
Astrophysics
General relativity

EPFL Lausanne
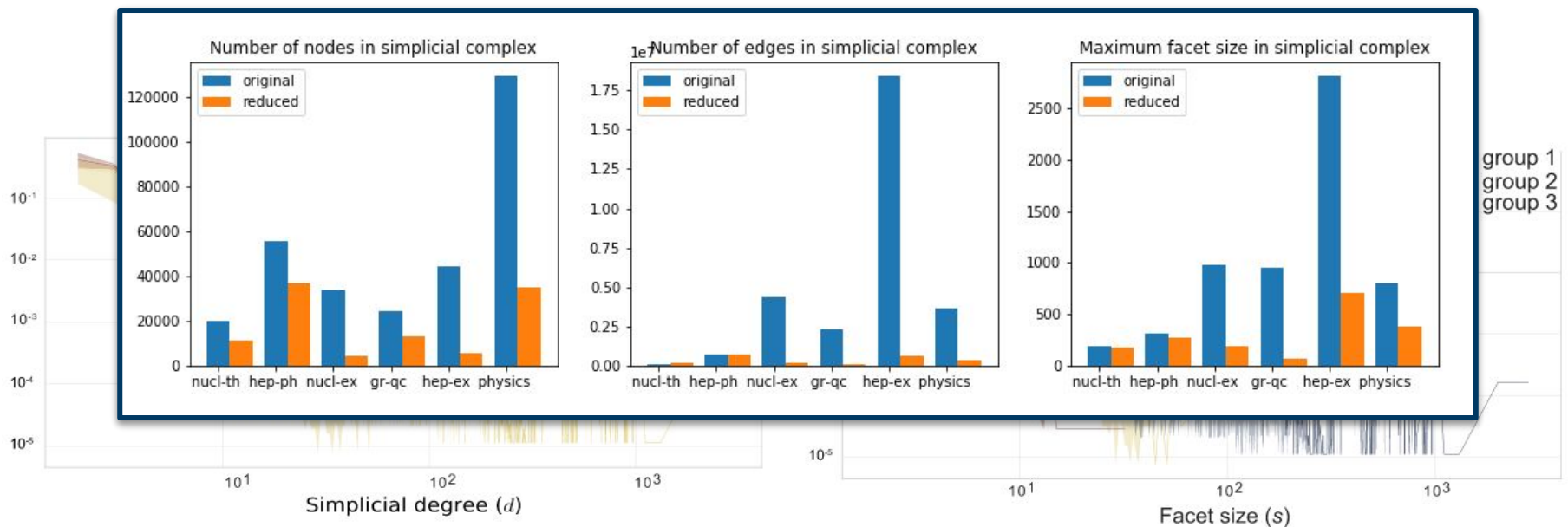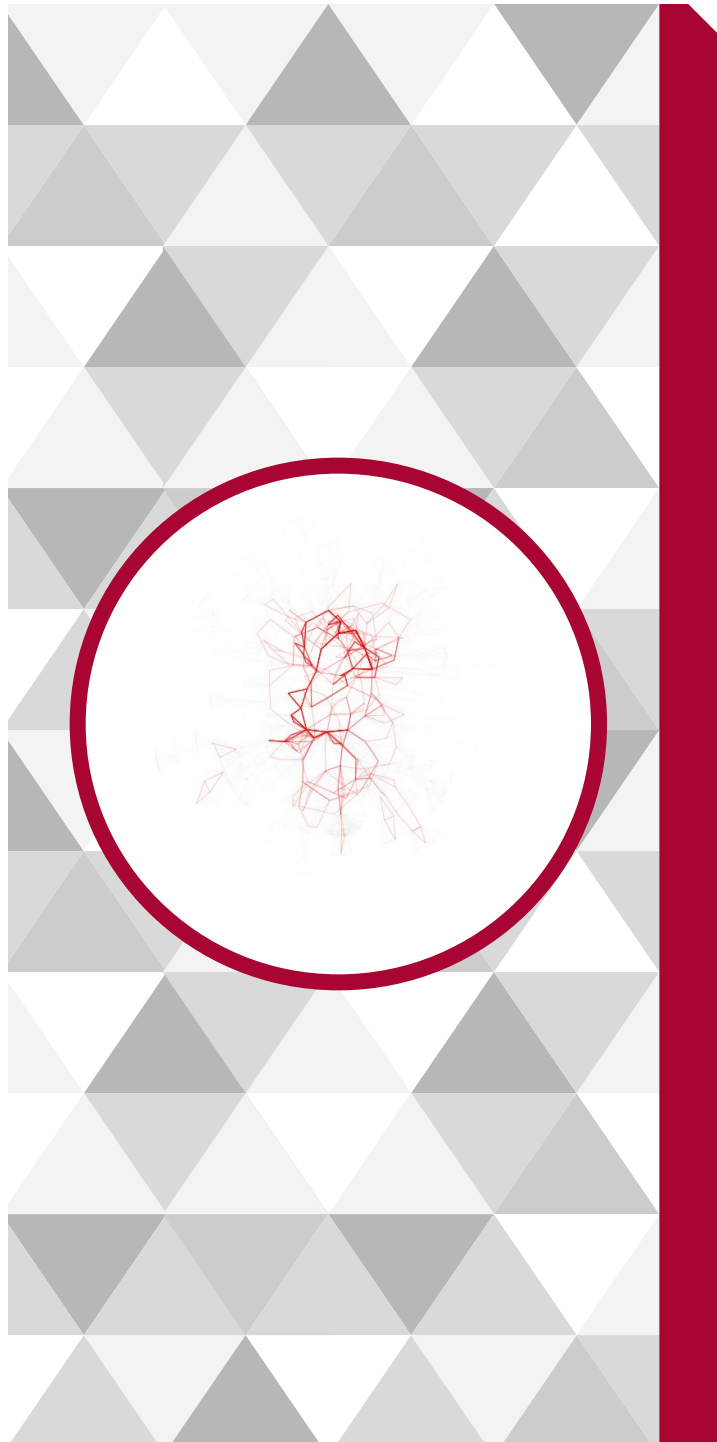
# Facets size and simplicial degree
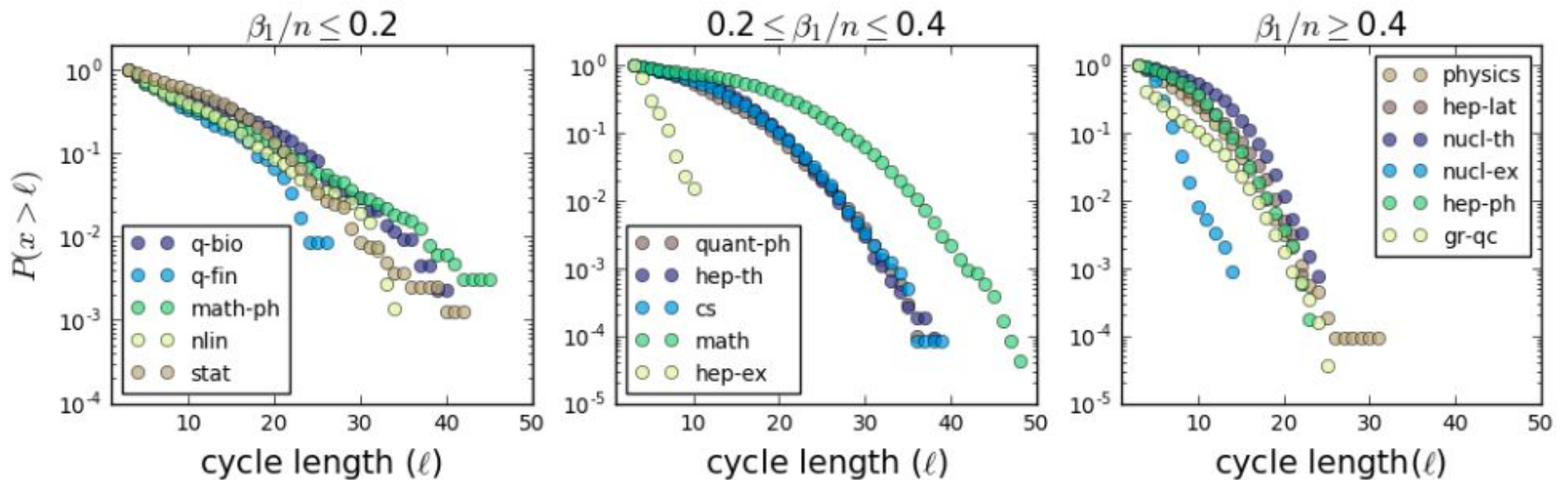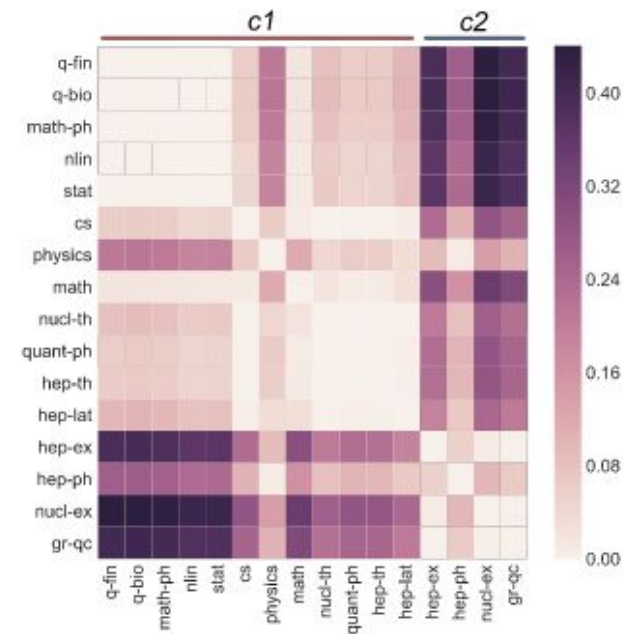## Homology equivalent simplicial complex
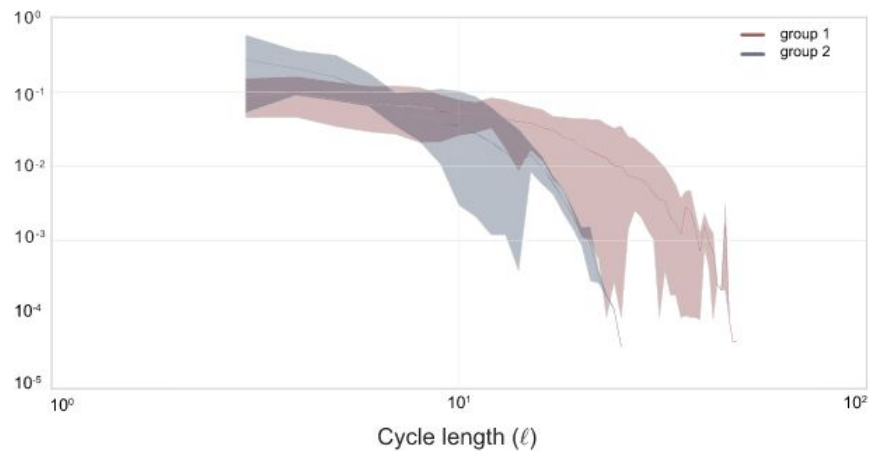
# Homological Results

# Homology

We introduce a new quantity **$\beta_1/n$** the ratio of the number of cycles over the number of nodes in the simplicial complex
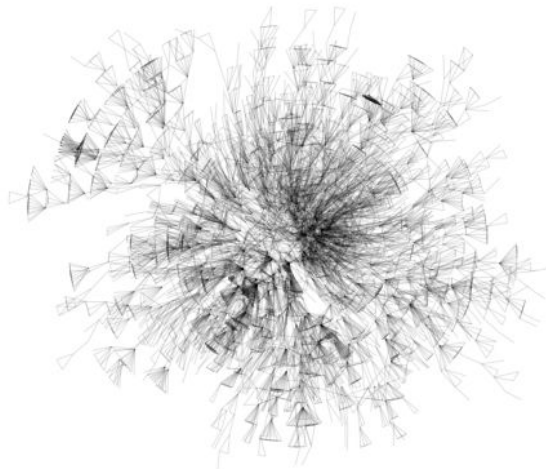
# Homology

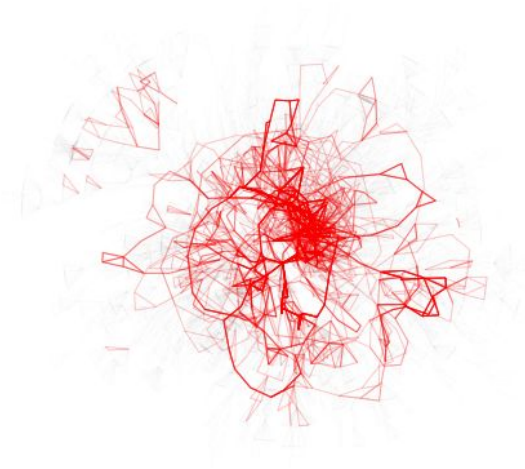**Cycle length distribution**
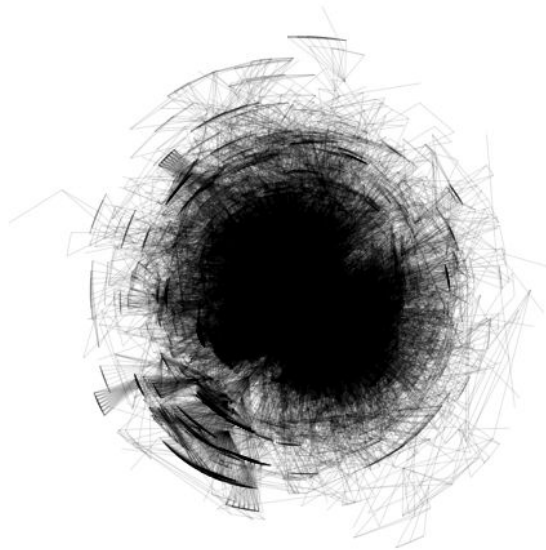
# Homology

## Cycle length distribution

math-ph

math-ph

Group 1

# Cycles length distribution

hep-lat

hep-lat

Group 2

EPFL Lausanne

# Community detection

**Assumption**:

Homological cycles act as bridges between communities of the underlying graph

SIMPLICIAL COMPLEX

Homological cycles

COMMUNITY detection
in underlying graph

# Community detection



**Infomap**



| i Cut-based perspective | ii Clustering perspective | iii Stochastically equivalent nodes | iv Dynamical perspective |

SCHAUB, Michael T., et al. The many facets of community detection in complex networks. *App. Net. Sc*, 2017.

EPFL Lausanne

# Homology and Communities

If cycles do not act as bridges between communities, then we expect them to go in and out randomly.

But we can clearly see that as cycles get longer the go through a larger number of communities.



math-ph - cycle len. = 3

# Homology bridges communities?

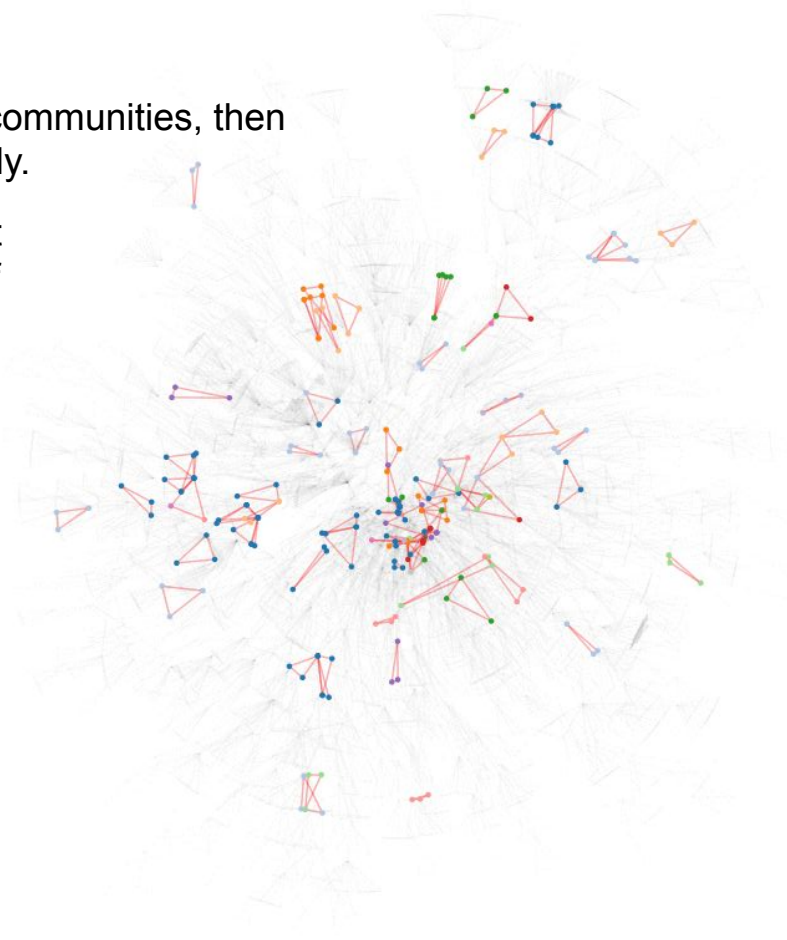If cycles do not act as bridges between communities, then we expect them to go in and out randomly.

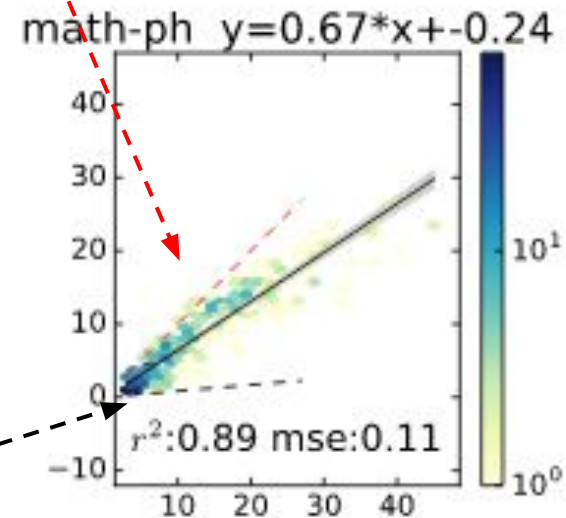But we can clearly see that as cycles get longer the go through a larger number of communities.

We decide to this as lower bound to assess if a cycle act as bridge in a category, and the length of the cycles as upper bound.
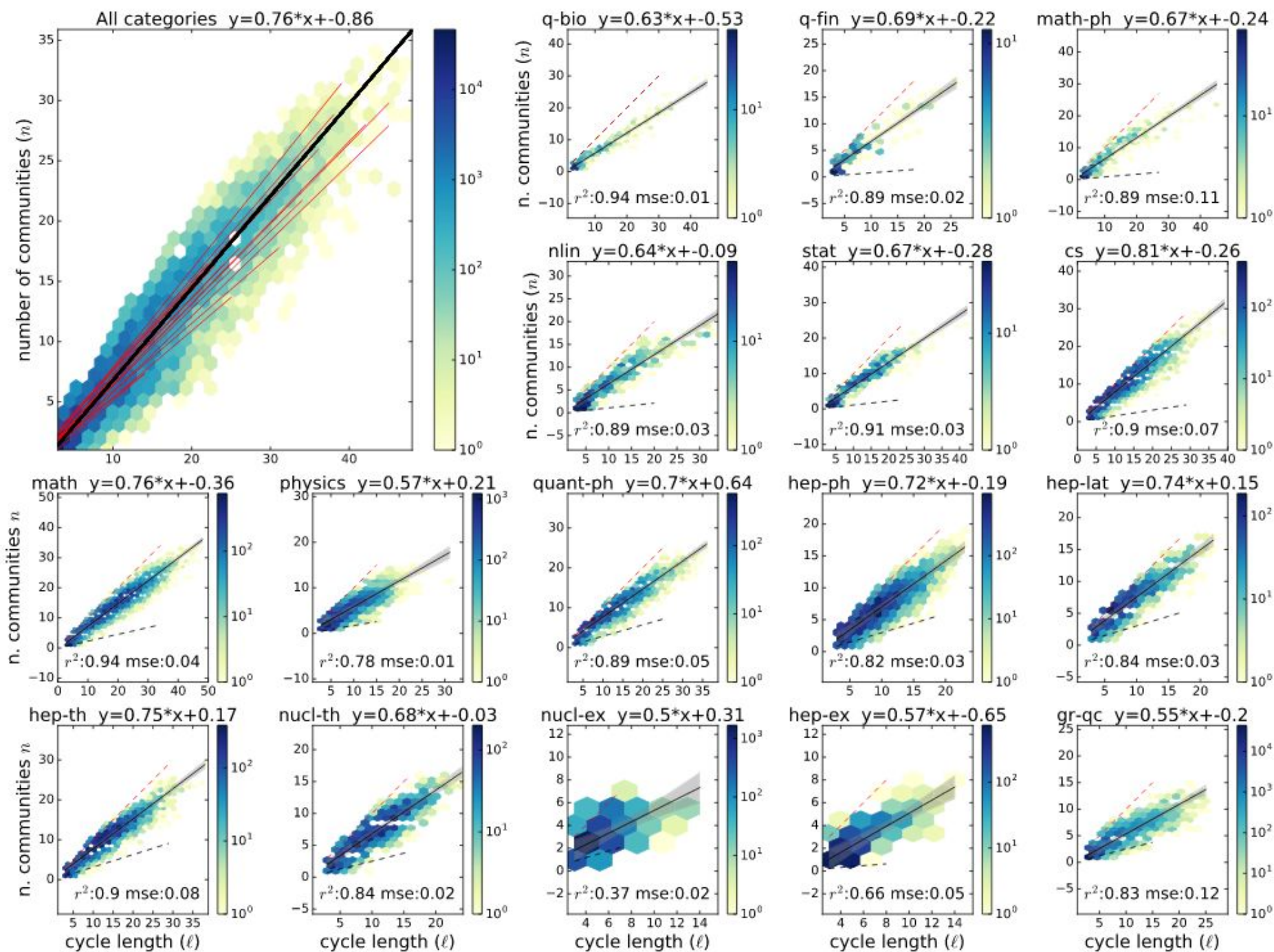
UPPER BOUND
**y = x**

**y=m*x**

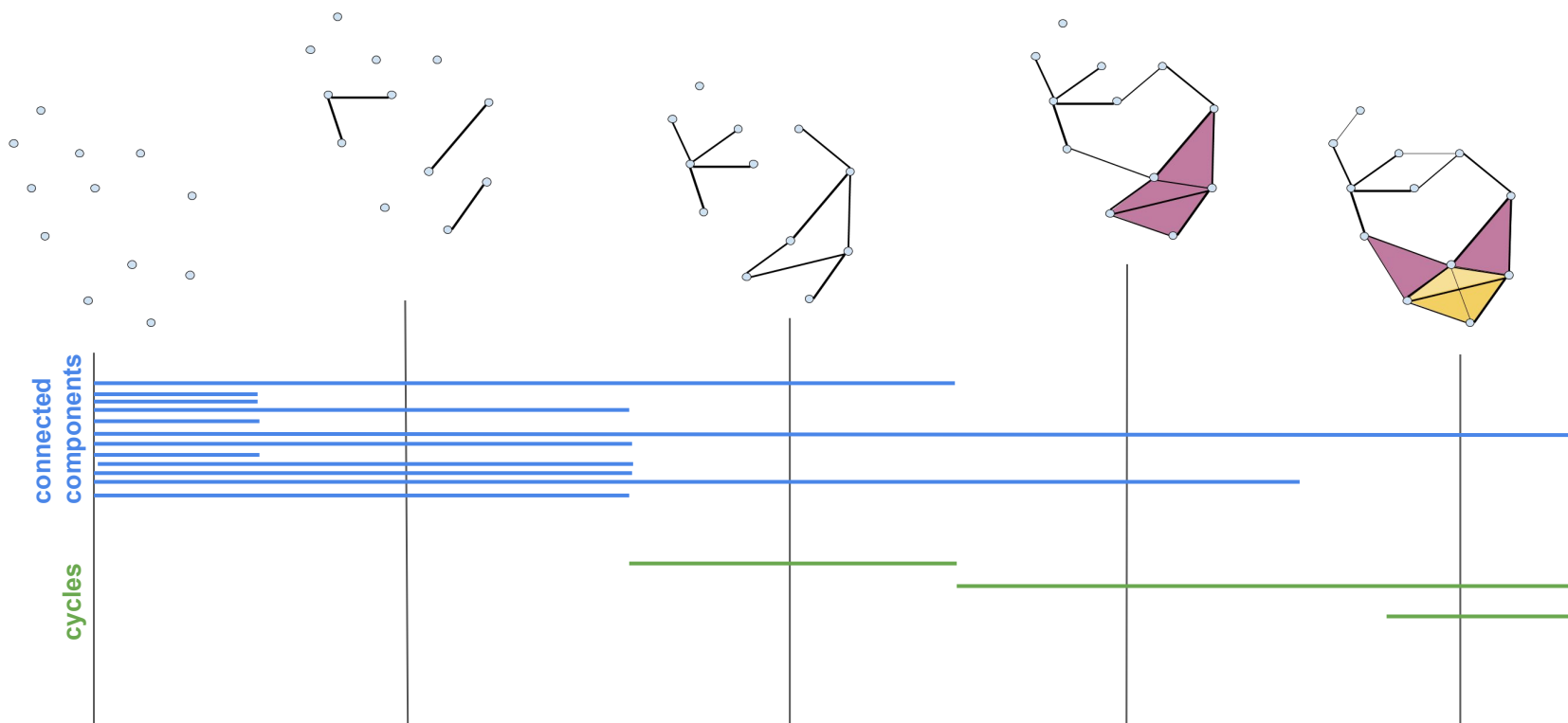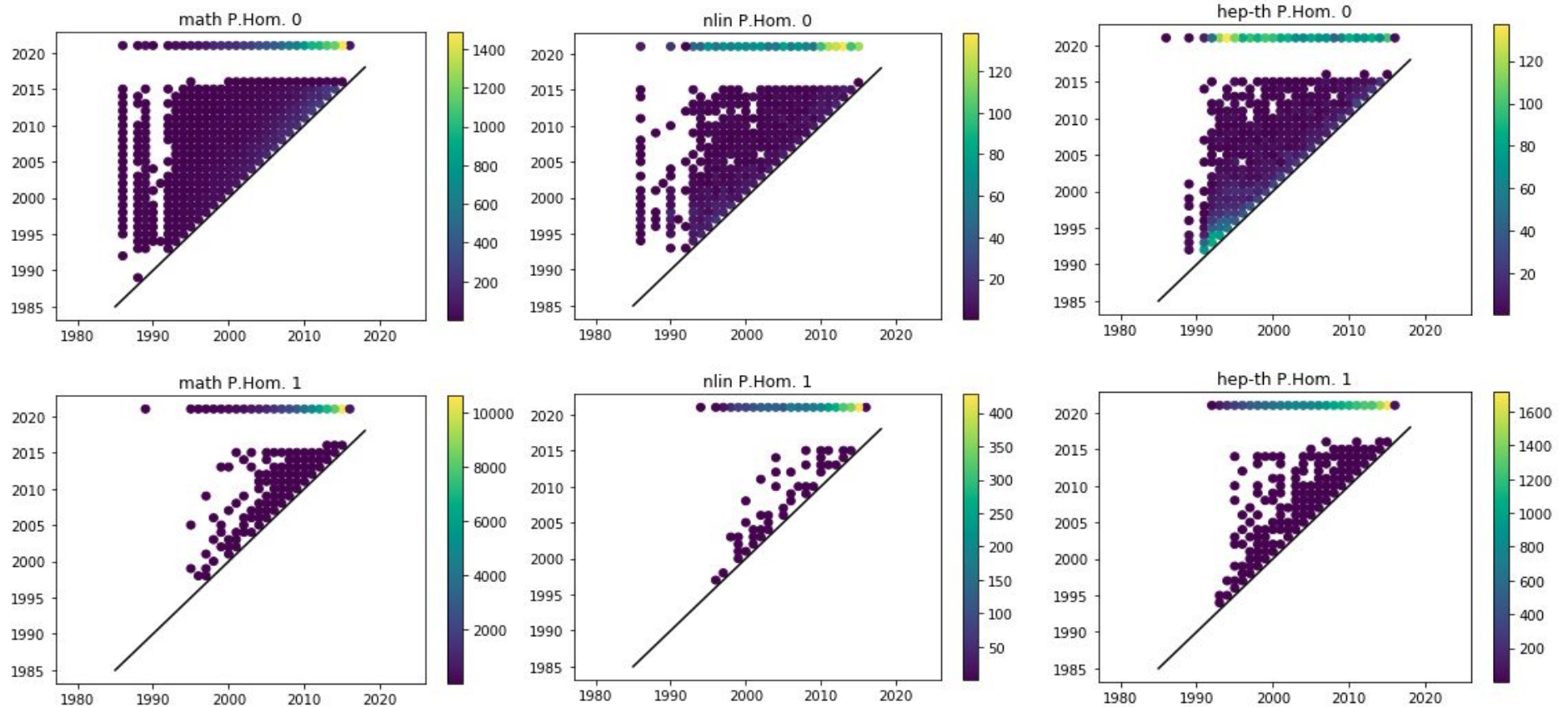Where **m** is the fraction of total edges between communities



math-ph  y=0.67*x+-0.24

$r^2$:0.89 mse:0.11

**All categories** y=0.76*x+-0.86

**q-bio** y=0.63*x+-0.53 — $r^2$:0.94 mse:0.01

**q-fin** y=0.69*x+-0.22 — $r^2$:0.89 mse:0.02

**math-ph** y=0.67*x+-0.24 — $r^2$:0.89 mse:0.11

**nlin** y=0.64*x+-0.09 — $r^2$:0.89 mse:0.03

**stat** y=0.67*x+-0.28 — $r^2$:0.91 mse:0.03

**cs** y=0.81*x+-0.26 — $r^2$:0.9 mse:0.07

**math** y=0.76*x+-0.36 — $r^2$:0.94 mse:0.04

**physics** y=0.57*x+0.21 — $r^2$:0.78 mse:0.01

**quant-ph** y=0.7*x+0.64 — $r^2$:0.89 mse:0.05

**hep-ph** y=0.72*x+-0.19 — $r^2$:0.82 mse:0.03

**hep-lat** y=0.74*x+0.15 — $r^2$:0.84 mse:0.03

**hep-th** y=0.75*x+0.17 — $r^2$:0.9 mse:0.08

**nucl-th** y=0.68*x+-0.03 — $r^2$:0.84 mse:0.02

**nucl-ex** y=0.5*x+0.31 — $r^2$:0.37 mse:0.02

**hep-ex** y=0.57*x+-0.65 — $r^2$:0.66 mse:0.05

**gr-qc** y=0.55*x+-0.2 — $r^2$:0.83 mse:0.12

number of communities ($n$) — cycle length ($\ell$)

# Future work

# Future work
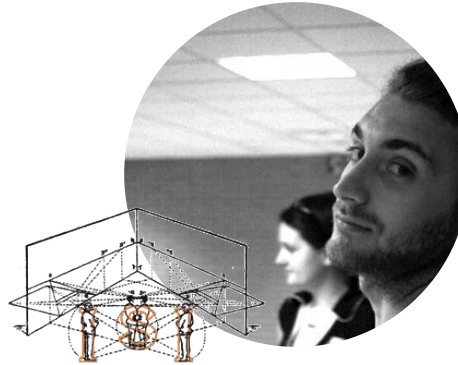
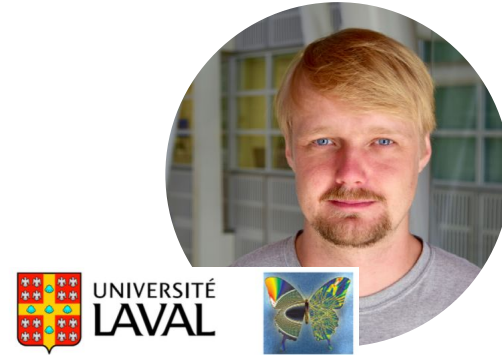## Extending to Persistent Homology

# Future work

Prof. F. Vaccarino

Dott. G. Petri

J.-G. Young

# Thank you
# for the attention

**Patania A., Petri G., and Vaccarino F.**
     "The shape of collaborations." EPJ Data Science 6.1 (2017): 18.

**Young J.-G., Petri G., Vaccarino F. and Patania A.**
     "Construction of and efficient sampling from the simplicial configuration model"
     PRE 96 (3), 032312 (2017)

Indiana University
**Network Science Institute**