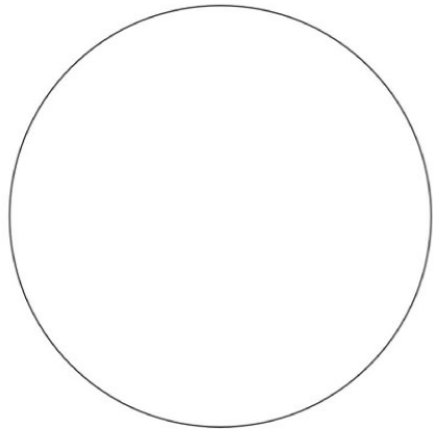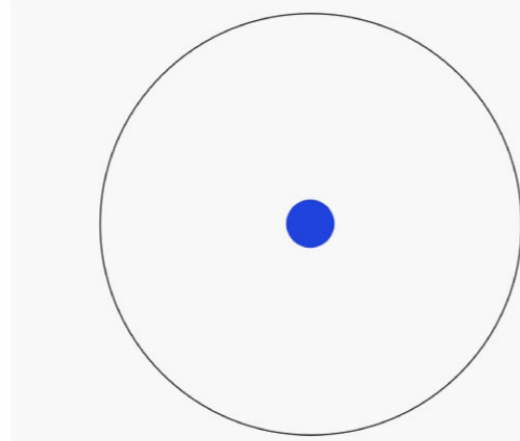# A RESEARCH QUESTION

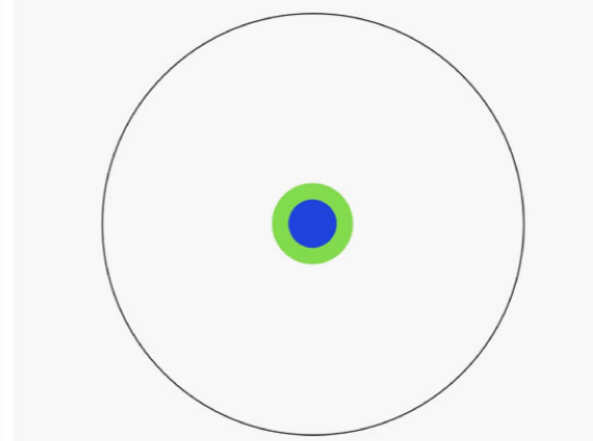Communities in degenerate link streams

Imagine a circle that contains all of human knowledge:
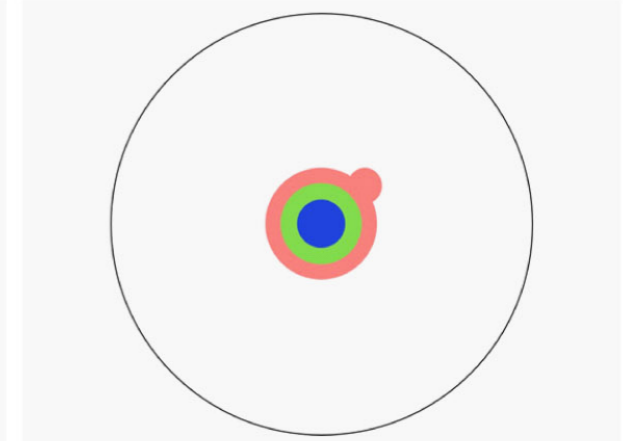
By the time you finish elementary school, you know a little:

By the time you finish high school, you know a bit more:

With a bachelor's degree, you gain a specialty:

A master's degree deepens that specialty:

Reading research papers takes you to the edge of human knowledge:

Once you're at the boundary, you focus:

You push at the boundary for a few years:

Until one day, the boundary gives way:

And, that dent you've made is called a Ph.D.:

Ph.D.

Of course, the world looks different to you now:

So, don't forget the bigger picture:

Ph.D.

Keep pushing.

# DYNAMIC NETWORKS

- Most real world networks are dynamic
  - ‣ Facebook friendship
    - People joining/leaving
    - Friend/Unfriend
  - ‣ Twitter mention network
    - Each mention has a timestamp
    - Aggregated every day/month/year => still dynamic
  - ‣ World Wide Web
  - ‣ Urban network
  - ‣ …

# DYNAMIC NETWORKS

- Most real world networks are dynamic
  - ‣ Nodes can appear/disappear
  - ‣ Edges can appear/disappear
  - ‣ Nature of relations can change

- How to represent those changes?

- How to manipulate dynamic networks?

# SEVERAL FORMALISMS

# TEMPORAL NETWORK

Collected dataset, for instance in (t,u,v) format

| Time | u | v |
|------|------|------|
| 1353304100 | 1148 | 1644 |
| 1353304100 | 1613 | 1672 |
| 1353304100 | 656 | 682 |
| 1353304100 | 1632 | 1671 |
| | | |
| 1353304120 | 1492 | 1613 |
| 1353304120 | 656 | 682 |
| 1353304120 | 1632 | 1671 |
| | | |
| 1353304140 | 1148 | 1644 |
| | | |
| 1353304160 | 656 | 682 |
| 1353304160 | 1108 | 1601 |
| 1353304160 | 1632 | 1671 |
| 1353304160 | 626 | 698 |

Examples:
-SocioPatterns
-Enron

-…

# TEMPORAL NETWORK

## Snapshots

1353304100   1148  1644
1353304100   1613  1672
1353304100   656   682
1353304100   1632  1671

1353304120   1492  1613
1353304120   656   682
1353304120   1632  1671

1353304140   1148  1644

1353304160   656   682
1353304160   1108  1601
1353304160   1632  1671
1353304160   626   698

## Link Stream

1353304100   1148  1644
1353304100   1613  1672
1353304100   656   682
1353304100   1632  1671

1353304120   1492  1613
1353304120   656   682
1353304120   1632  1671

1353304140   1148  1644

1353304160   656   682
1353304160   1108  1601
1353304160   1632  1671
1353304160   626   698

## Interval Graph

1353304100   1148  1644
1353304100   1613  1672
1353304100   656   682
1353304100   1632  1671

1353304120   1492  1613
1353304120   656   682
1353304120   1632  1671

1353304140   1148  1644

1353304160   656   682
1353304160   1108  1601
1353304160   1632  1671
1353304160   626   698

# SLOWLY EVOLVING NETWORKS (SEN)

# SLOWLY EVOLVING NETWORKS

- Edges change (relatively) slowly

- The network is well defined at any t
  ‣ Nodes/edges described by (long lasting) intervals
  ‣ Enough snapshots to track nodes

- A static analysis at every (relevant) t gives a dynamic vision

- No formal distinction with previous case (higher observation frequency)

# SLOWLY EVOLVING NETWORKS

- Visualization
  - ‣ Problem of stability of node positions

# SLOWLY EVOLVING NETWORKS

- Global graph properties

Leskovec, Jure, Jon Kleinberg, and Christos Faloutsos. "Graph evolution: Densification and shrinking diameters." *ACM Transactions on Knowledge Discovery from Data (TKDD)* 1.1 (2007): 2.

# UNSTABLE/DEGENERATE TEMPORAL NETWORKS

Matthieu Latapy, Tiphaine Viard, and Clémence Magnien. "Stream graphs and link streams for the modeling of inter- actions over time". In: *Social Network Analysis and Mining* 8.1 (2018), p. 61.

# UNSTABLE TEMPORAL NETWORK

- The network at a given $t$ is not meaningful

- How to analyze such a network?

# UNSTABLE TEMPORAL NETWORK

# UNSTABLE TEMPORAL NETWORK

- Common solution: transform into SEN using aggregation/ sliding windows
  - ‣ Information loss
  - ‣ How to chose a proper aggregation window size?

- New theoretical tools developed to deal with such networks

- **Link Streams** & **Stream Graphs** (Latapy, Viard, and Magnien 2018)

- **Temporal Networks**, **Contact Sequences** and **Interval Graphs** (Holme and Saramäki 2012)

- **Time Varying Graphs** (Casteigts et al. 2012)

# CENTRALITIES
# &
# NETWORK PROPERTIES
# IN STREAM GRAPHS

# STREAM GRAPHS

## Stream Graph (SG)- Definition

Stream Graphs have been proposed in[a] as a generic formalism – it can represent any type of dynamic networks, continuous, discrete, with or without duration, with the objective or redefining typical notions of graphs on dynamic networks, including degenerate ones.

Let's define a Stream Graph

$$S = (T, V, W, E)$$

| | |
|---|---|
| $T$ | **Set of Possible times** (Discrete or Time intervals) |
| $V$ | **Set of Nodes** |
| $W$ | **Vertices presence time** $V \times T$ |
| $E$ | **Edges presence time** $V \times V \times T$ |

_____

[a]Latapy, Viard, and Magnien 2018.

# STREAM GRAPHS

## SG - Time-Entity designation

It is useful to work with Stream Graphs to introduce some new notions mixing entities (nodes, edges) and time:

| | |
|---|---|
| $V_t$ | **Nodes At Time**: set of nodes present at time $t$ |
| $E_t$ | **Edges At Time**: set of edges present at time $t$ |
| $G_t$ | **Snapshot**: Graph at time $t$, $G_t = (V_t, E_t)$ |
| $v_t$ | **Node-time**: $v_t$ exists if node $v$ is present at time $t$ |
| $(u, v)_t$ | **Edge-time**: $(u, v)_t$ exists if edge $(u, v)$ is present at time $t$ |
| $T_u$ | **Times Of Node**: the set of times during which $u$ is present |
| $T_{uv}$ | **Times Of Edge**: the set of times during which edge $(u, v)$ is present |

# STREAM GRAPHS

| | |
|---|---|
| $N_u$ | **Node presence**: The fraction of the total time during which $u$ is present in the network $\frac{|T_u|}{|T|}$ |
| $L_{uv}$ | **Edge presence**: The fraction of the total time during which $(u, v)$ is present in the network $\frac{|T_{uv}|}{|T|}$ |

# STREAM GRAPHS

## SG - Redefining Graph notions

The general idea of redefining static network properties on Stream Graphs is that if the network stays unchanged along time, then properties computed on the stream graph should yield the same values as the same property computed on the aggregated graph.

# STREAM GRAPHS

## SG - $N$ & $L$

The number/quantity of nodes in a stream graph is defined as the total presence time of nodes divided by the dataset duration. In general, it isn't an integer.
More formally:

$$N = \sum_{v \in V} N_v = \frac{|W|}{|T|}$$

For instance, $N = 2$ if there are 4 nodes present half the time, or two nodes present all the time.

Figure 2: **Two strea**
**densities:** Left: $\delta = 0$

In addition, $\delta(L)$ is eq

$$\frac{1}{|T| \cdot |V \otimes V|} \int_t |E_t| \, \mathrm{d}t = \frac{\int_t}{\int_t |W|}$$

Finally, if we consid
of the corresponding g

# STREAM GRAPHS



$$N = 2$$

# STREAM GRAPHS

of nodes $V$ and the same

induced by a subset $V'$ of

$i.e. \ldots = (T', (T \times V') \cap W,$

induced by a subset $T'$ of

$V) \cap W, (T' \times V \otimes V) \cap E$

For the example in Fig

is $([6,9], \{a, b, c\}, [6, 9] \times$

**SG -** $L$

The number of edges is defined as the total presence of nodes divided by the total dataset duration.
More formally:

$$L = \sum_{(u,v), u,v \in V} L_{uv} = \frac{|E|}{|T|}$$

For instance, $L = 2$ if there are 4 edges present half the time, or two edges present all the time.

# 7 Cliques

A clique of graph $G$ is a

involved in $C$ are linked t

$C'$ such that $C \subset C'$.

# STREAM GRAPHS



$$L = 1$$

# 7 Cliques

*A clique of graph G is a*
*involved in C are linked t*
*C' such that C ⊂ C'.*

We define a **clique** of
all pairs of nodes involved
C is maximal if there is n
We say that a clique
(resp. uniform). It is the
set) meaning that all pair

**SG - Edge domain - $L_{\max}$**

In Stream Graphs, several possible definitions of $L_{\max}$ could exist.

- Ignoring nodes duration: $L_{\max}\,1 = |V|^2$
- Ignoring nodes co-presence $L_{\max}\,2 = N$
- Taking nodes co-presence into account:
  $L_{\max}{}^3 = \sum_{(u,v),u,v \in V} |T_u \bigcap T_v|$

betw
deno

Figu
disp
indu
$\{ac\}$
of $G$
set o
insta
$V_t$
at ti
E
$T^C_{ab} =$

(resp. uniform). It is then fu

set) meaning that all pairs o

$T^C_{ab} = [6, 8$

The density in static networks can be understood as the fraction of existing edges among all possible edges,

$$d = \frac{L}{L_{\max}}.$$

In the following, we will use $L_{max}$, as in Latapy et al.

Figure 4: **Examples of ma**

pact cliques involving three

and $[7, 8] \times \{b, c, d\}$. Its oth

$[2, 5] \times \{a, c\}, [1, 8] \times \{b, c\},$

For instance, in Figure 4

a

b

a

b

# STREAM GRAPHS

$$N = 2 \qquad\qquad L = 1$$



$$d = \frac{3}{6} = \frac{1}{2} \qquad\qquad d = \frac{3}{4} \qquad\qquad d = \frac{3}{3} = 1$$

# STREAM GRAPHS

## SG - Clusters & Substreams

In static networks, a cluster is a set of nodes, and we have defined an (induced) subgraph of this cluster as a graph composed of the nodes of the cluster and the edges existing between those nodes. In Stream Graphs, a clusters $C$ is as subset of $W$, and the corresponding (induced) substream $S(C) = (T, V, C, E(C))$, with $E(C) = \{(t, (u, v)) \in E, (t, u), (t, v) \in C\}$.



Example of subgraph (red,left) and induced substream (right).

# STREAM GRAPHS

## SG - Cliques

Having defined substreams and density, we can now naturally define a clique by analogy with static networks as a substream of density 1. A clique is said to be a **maximal clique** if it is not included in any other clique.



Red and Grey are the two maximal cliques of size three in this Stream Graph.

**1. Examples of maximal compact cliques.** We display the two maximal cliques involving three nodes of the link stream $L$ of Figure 1 (right): $[2,4] \times \{$ $8] \times \{b,c,d\}$. Its other maximal compact cliques are $[0,4] \times \{a,b\}$, $[6,9] \times \{a,c\}$, $[1,8] \times \{b,c\}$, $[7,10] \times \{b,d\}$, $[6,9] \times \{c,d\}$ (involving two nodes each

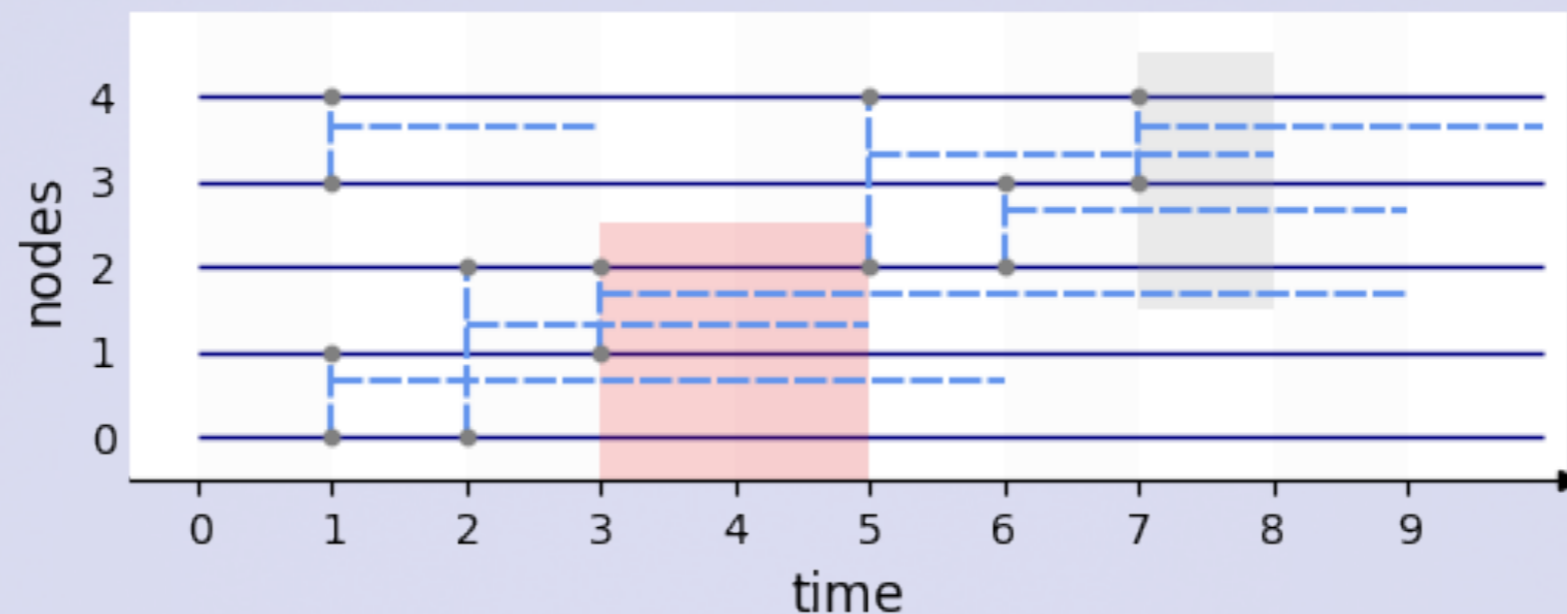instance, in Figure ... ct clique. However ximal, as it is included in $[0,4] \times \{a,b\}$ which is a maximal compact clique intersects another maximal compact clique, $[2,4] \times \{a,b,c\}$. There is a unique al compact clique involving three nodes, $[8,9] \times \{b,c,d\}$. The maximal co $[0,4] \times \{a,b\}$ is not a maximal clique because it is for instance included in the $\{a,b\} \cup [6,9] \times \{c$ ... not maximal eit instance included in the clique $[0,4] \times \{a,b\} \cup [6,9] \times \{c,d\} \cup [5,6] \times \{d\}$.

$G_t := \frac{1}{T} \int_T \delta(G_t) dt = \frac{1}{T}$ ... clique in $S$ does not in general induce a clique in $G(S)$: for instance, $[0,1] \times \{$ $= V$ for all $t$, then its density is equal to the density $\{c,d\}$ is a clique for the example in Figure 4 but $\{a,b,c,d\}$ is not a clique graph. Instead, we define the den ... $[b,e] \times X$ is a compact node $v$ in $V$, and the density

$\delta(t) = \frac{|E_t|}{|V_t \otimes V_t|}$

$\frac{}{|T_v|}$ and $\delta(t) = \frac{|E_t|}{|V_t \otimes V_t|}$.

0, respectively, then we define $\delta(uv)$,
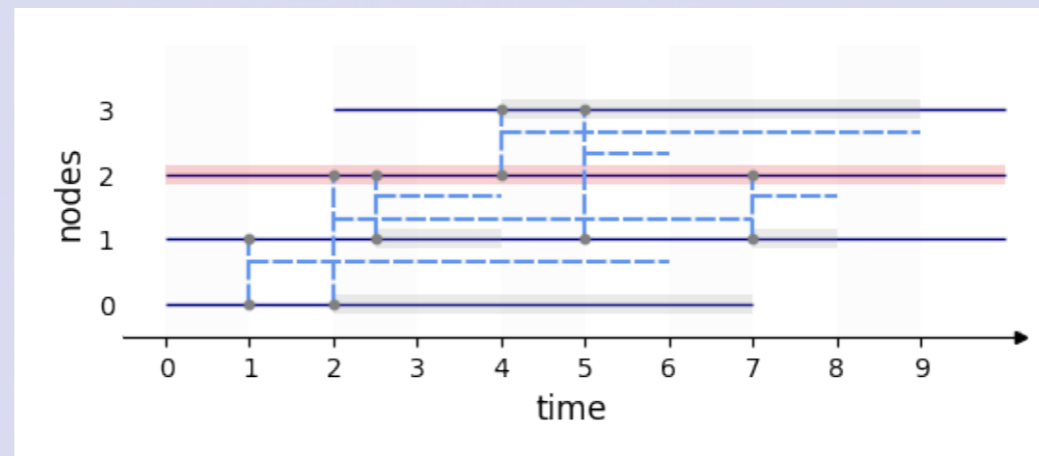


**SG - Neighborhood $N(u)$**

The neighborhood $N(u)$ of node $u$ is defined as the cluster composed of node-times such as an edge-time exists between it and a node-time of $u$, i.e.,

$$N(u) = \{v_t, (u,v)_t \in E\}$$

**SG - Degree $k(u)$**

The degree $k(u)$ of node $u$ is defined as the quantity of node in the Neighborhood of node $u$, i.e.

$$k(u) = |N(u)|$$

Example, the neighborhood of node 2 is highlighted in grey.
$k(c) = \frac{5 + 2.5 + 5}{10} = 1.25$.

8 time

0 2 4 6 8 time

$(t_i, u_i v_i) \in E$, $[\alpha, t_0] \times \{$

This sequence is similar

**f neighborhoods and degrees of nodes.** We display

node under concern, and in grey the other links. Left:

not necessarily have $t_0 \geq$

is symmetric: if $(\alpha, u)$--

$4.5, 7.5] \times \{c\}$ is in blue, leading to $d(a) = \frac{3}{10} + \frac{3}{\overline{W}} = 0.6$.

$(9, d)$--$(3, g)$ through t

We say that $S$ is **wea**

$] \times \{b\} \cup [6, 9] \times \{d\}$ is in blue, leading to $d(c) = \frac{13}{\overline{W}} = 1.3$.

We say that a cluster $C$

**node de**gree of $S$ as follows.

connected. It is a weakly

cluster of $S$. Intuitively,

$$\frac{1}{n} \cdot \sum_{v \in V} n_v \cdot d(v) = \sum_{v \in V} \frac{|T_v|}{|W|} \cdot d(v)$$

Figure 14 for an illustrat

## SG - Ego-network

The Ego network $G_u$ of node $u$ is defined as the substream induced by its neighborhood, i.e., $G_u = (T, V, N(u), E(S(u)))$.

## SG - Clustering coefficient

The clustering coefficient $C(u)$ of node $u$ is defined as the density of the ego-network of $u$, i.e.,

*(hence $v_{i-1} = u_i = v_j = u_{j+1}$) then $P' = (u_0, v_0), \ldots, (u_{i-1}, v_{i-1}), (u_{j+1}, v_{j+1}), \ldots, (u_k,$*

$$C(u) = d(N(u))$$

*also is a path from $u$ to $v$. If one iteratively removes the cycles of $P$ in this way, eventually obtains a simple path from $u$ to $v$.*

*The path $P$ is a shortest path from $u$ to $v$ if there is no path in $G$ of length lower th $k$. Then, $k$ is called the distance between $u$ and $v$ and it is denoted by $\partial(u, v)$. If ther no path between $u$ and $v$ then their distance is infinite. The diameter of $G$ is the lar finite distance between two nodes in $V$.*

Figure 14: **Weakly co**

# RANDOM MODELS FOR DYNAMIC NETWORKS

Laetitia Gauvin et al. "Randomized reference models for temporal networks". In: *SIAM Review* 64.4 (Nov. 2022)
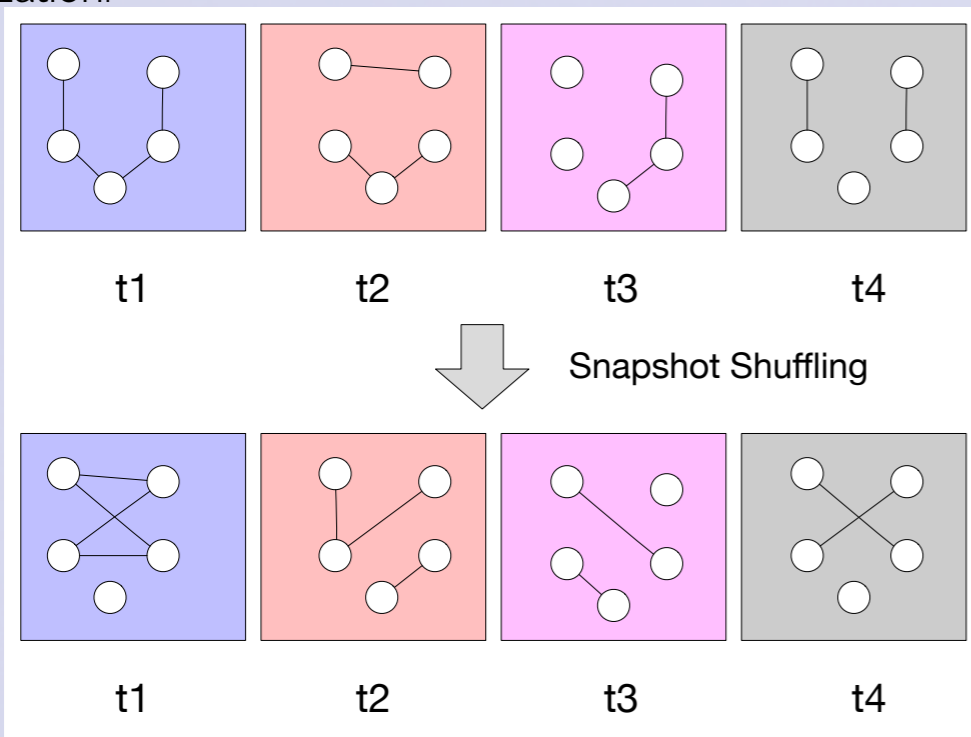
# RANDOM MODELS

- In many cases, in network analysis, useful to compare a network to a randomized version of it
  - ‣ Clustering coefficient, assortativity, modularity, …

- In a static graph, 2 main choices:
  - ‣ Keep only the number of edges (ER model)
  - ‣ Keep the number of edges and the degree of nodes (Configuration model)

- In dynamic networks, it is more complex…

# RANDOM M



## Snapshot Shuffling

**Snapshot Shuffling** keeps the order of snapshots, randomize edges inside snapshots. Any random model for static network can be used, such as ER random graphs or a degree preserving randomization.



## Sequence Shuffling

**Sequence Shuffling** keeps each snapshot identical, switch randomly their order.

# RANDOM MODELS

## Link Shuffling

**Link Shuffling** keeps activation time per node pairs, randomize the aggregated graph. For instance, a simple way to achieve this is to pick two node pairs at random (connected or not) of the aggregated graph, and to exchange activation time of these node pa...



## Timeline Shuffling

**Timeline Shuffling** keeps the aggregated graph, randomize edges activation time. For instance, a simple way to achieve this is to redistribute randomly activation period among all edges, e.g.:

# DYNAMIC COMMUNITY DETECTION

Rossetti, G., & Cazabet, R. (2018). Community discovery in dynamic networks: a survey. *ACM Computing Surveys (CSUR)*, *51*(2), 1-37.

Cazabet, R., Boudebza, S., & Rossetti, G. (2020). Evaluating community detection algorithms for progressively evolving graphs. *Journal Of Complex Networks*

# COMMUNITY DETECTION

## Static networks

Clusters: Sets of nodes



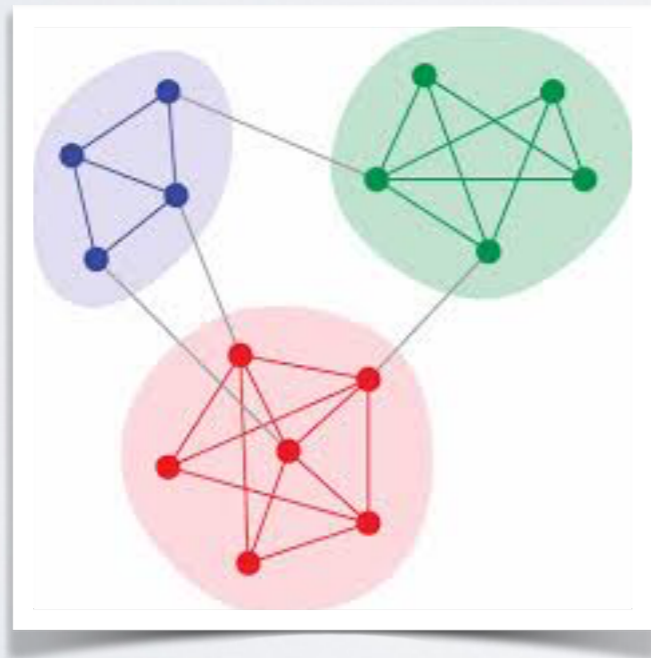## Dynamic Networks

Clusters: Sets of time-nodes, i.e., pairs (node,time)



Gaumont, N., Viard, T., Fournier-S'Niehotta, R., Wang, Q., & Latapy, M. (2016). Analysis of the temporal and structural features of threads in a mailing-list. In *Complex Networks VII*
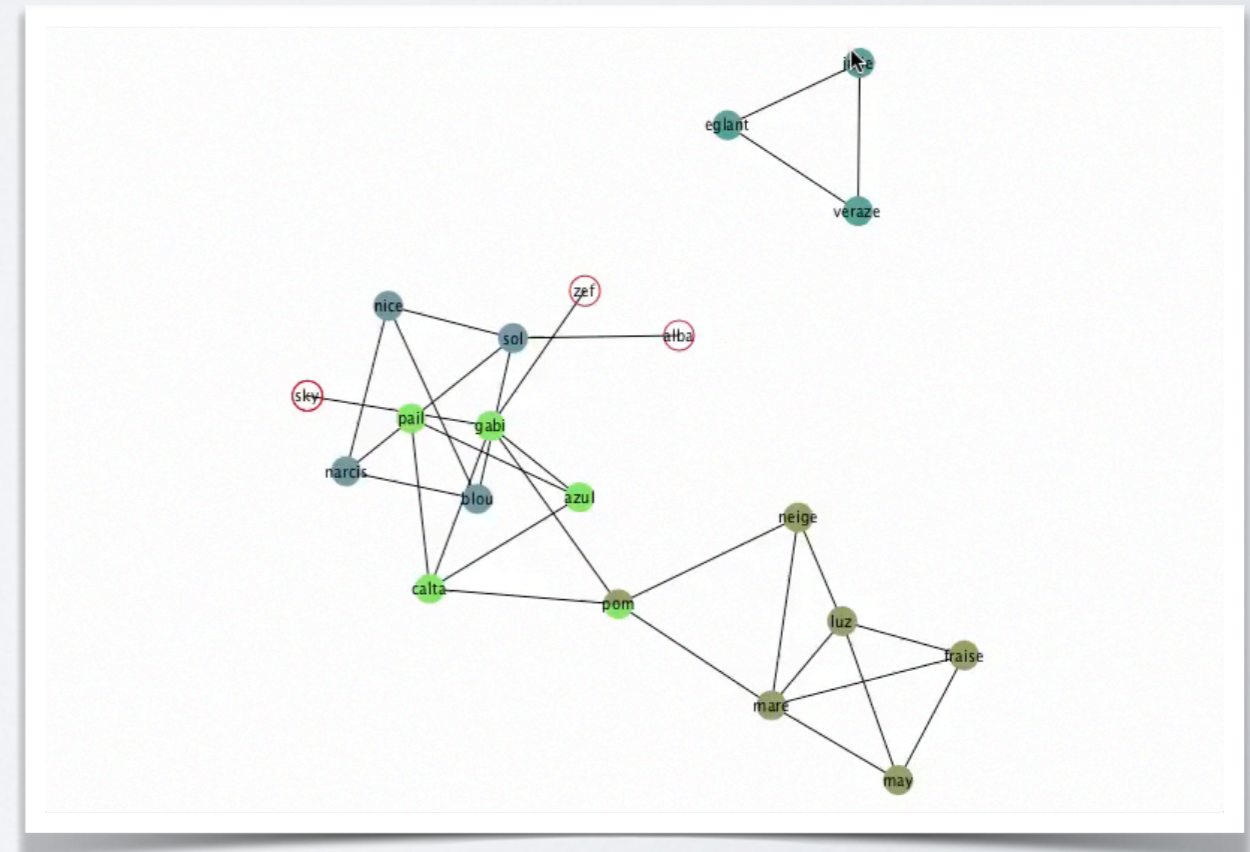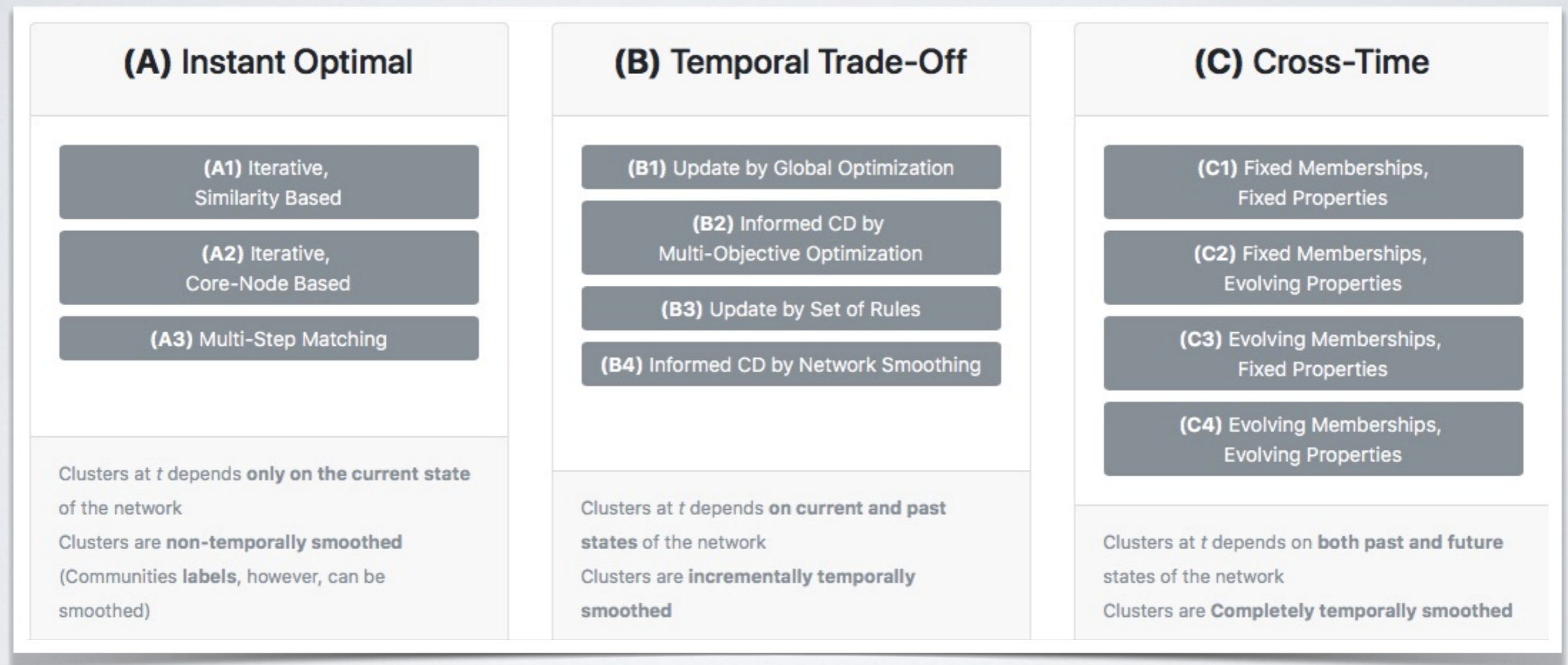
# COMMUNITY DETECTION

Static networks

Dynamic Networks

Clusters: Sets of nodes

Clusters: Sets of time-nodes, i.e., pairs (node,time)

# APPROACHES TO DCD

# DYNAMIC COMMUNITIES ?

## More than 50 methods published, broad categories

| (A) Instant Optimal | (B) Temporal Trade-Off | (C) Cross-Time |
|---|---|---|
| **(A1)** Iterative, Similarity Based | **(B1)** Update by Global Optimization | **(C1)** Fixed Memberships, Fixed Properties |
| **(A2)** Iterative, Core-Node Based | **(B2)** Informed CD by Multi-Objective Optimization | **(C2)** Fixed Memberships, Evolving Properties |
| **(A3)** Multi-Step Matching | **(B3)** Update by Set of Rules | **(C3)** Evolving Memberships, Fixed Properties |
| | **(B4)** Informed CD by Network Smoothing | **(C4)** Evolving Memberships, Evolving Properties |
| Clusters at t depends **only on the current state** of the network. Clusters are **non-temporally smoothed** (Communities **labels**, however, can be smoothed) | Clusters at t depends **on current and past states** of the network. Clusters are **incrementally temporally smoothed** | Clusters at t depends on **both past and future** states of the network. Clusters are **Completely temporally smoothed** |

Rossetti, G., & Cazabet, R. (2018). Community discovery in dynamic networks: a survey. *ACM Computing Surveys (CSUR), 51*(2), 1-37.

# CATEGORIES

- Instant optimal:
  - ‣ Allows reusing static algorithms
  - ‣ No partition smoothing
  - ‣ Labels can be smoothed
  - ‣ Simple to parallelize

# CATEGORIES

- Temporal trade-off
  - ‣ Cannot be parallelized (iterative)
  - ‣ => Best suited for real-time analysis / tasks

- Cross-Time
  - ‣ Requires to know the whole evolution in advance
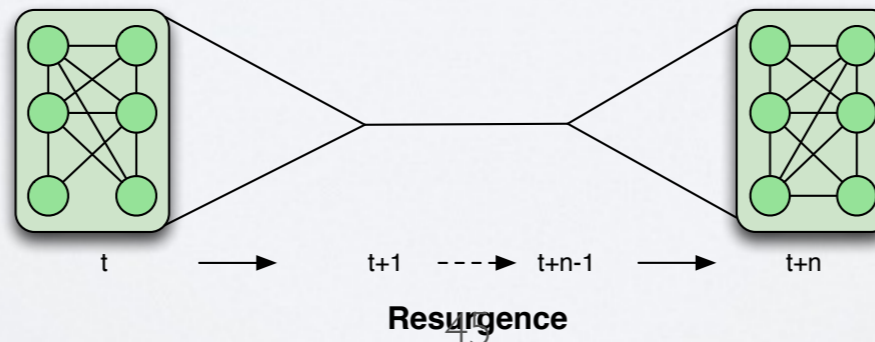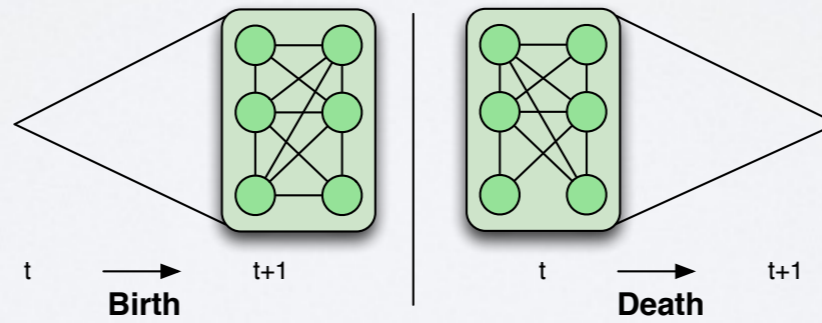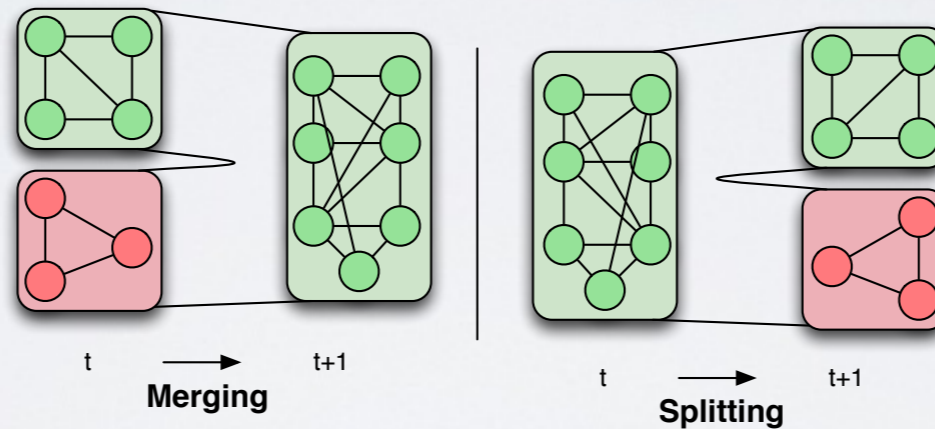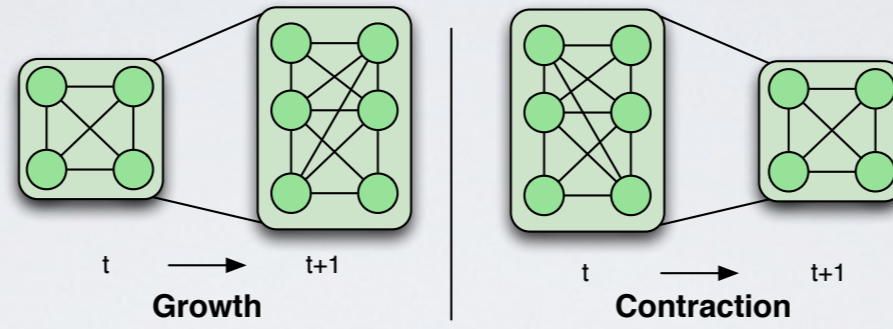  - ‣ => Not suited for real-time analysis, potentially the best smoothed (a posteriori interpretation)

# WHAT MAKES DCD INTERESTING

NARRATIVES ?

# SMOOTHNESS / STABILITY

- No Smoothness: Partition at **t** should be the same as found by a static algorithm.

- Smoothness: Partition at **t** is a trade-off between "good" communities for the graph at **t** and similarity with partitions at different times

# COMMUNITY EVENTS



t ⟶ t+1
**Growth**

t ⟶ t+1
**Contraction**

t ⟶ t+1
**Merging**

t ⟶ t+1
**Splitting**

t ⟶ t+1
**Birth**

t ⟶ t+1
**Death**

t ⟶ t+1 ⇢ t+n-1 ⟶ t+n
**Resurgence**

# PROGRESSIVE EVOLUTION



2 communities

??

Intermediate state

1 community

How to *track* communities, giving a *coherent* dynamic structure ?

# ONGOING WORK 1

How to adapt modularity for link streams?

# MODULARITY

- Most popular approach in static networks: modularity optimization

- Fraction of edges inside communities **-** Fraction of edges expected inside communities according to a **null model**

# MODULARITY

- How to adapt for link streams ?

# LONGITUDINAL MODULARITY



- Fraction inside communities=> OK

- Fraction expected?
  - ‣ We need to choose a null model
    - Keep node degrees (whole period or locally ?) => Timeline Shuffling => Globally
    - How many links are expected between two nodes on an interval $[t_{start}, t_{end}]$ ?

$$\approx \frac{k_u k_v}{2m} \frac{[t_{start}, t_{end}]}{T}$$
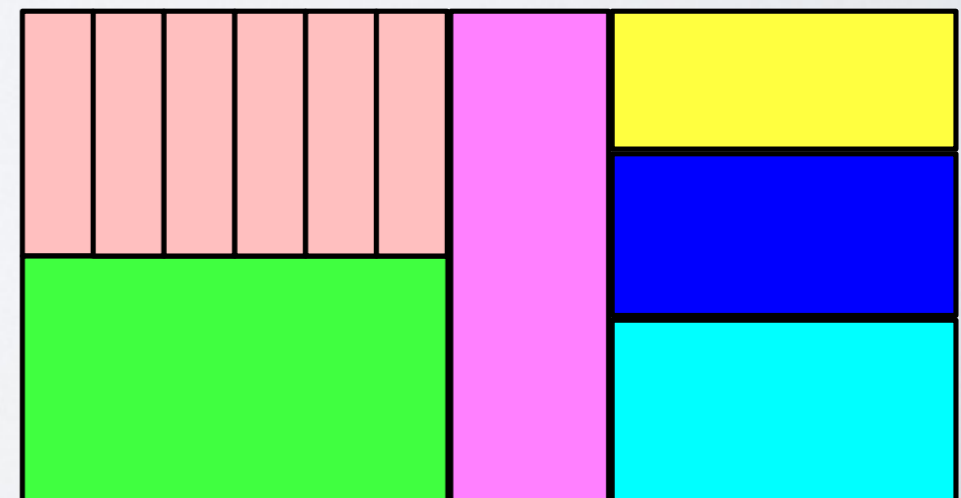
# LONGITUDINAL MODULARITY

Intuitively correct: find communities where
More links than expected by chance

# LONGITUDINAL MODULARITY

- But… it does not work as expected

- Consider the red communities as a single community or separate ones
  - Number of edges inside does not change
  - Number of expected edges does not change
  - $$\approx \frac{k_u k_v}{2m} \frac{[t_{start}, t_{end}]}{T}$$

# LONGITUDINAL MODULARITY

- The problem is a "smoothness" problem
  - No interest/gain to make a community last longer if it makes it less attractive at some point in time

- Solution Proposed: Work with edge repetitions
  - Modularity: Fraction of edges inside communities **-** Fraction of edges expected inside communities according to a null model
  - **Lmodularity:** Fraction of edges **repetitions** inside communities **-** Fraction of edges **repetitions** expected inside communities according to a null model

# LONGITUDINAL MODULARITY

- Modularity works thanks to a trade-off between:
  - ‣ Each node added to a community allows to add some edges inside
    - Large gain for each edge (1)
    - Linear gain
  - ‣ Each node added to a community increases the potential number of edges
    - Small penalty for each edge (<<1)
    - **Quadratic** penalty (square of nodes inside the community)

- LM works in a similar way:
  - ‣ Making communities last longer allows to have more edge repetitions
    - Large gain for each edge
  - ‣ Making communities last longer increases **quadratically** the expectation of the number of repetitions
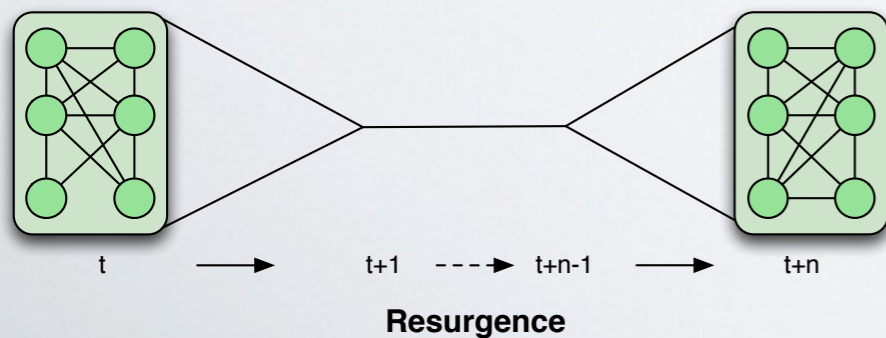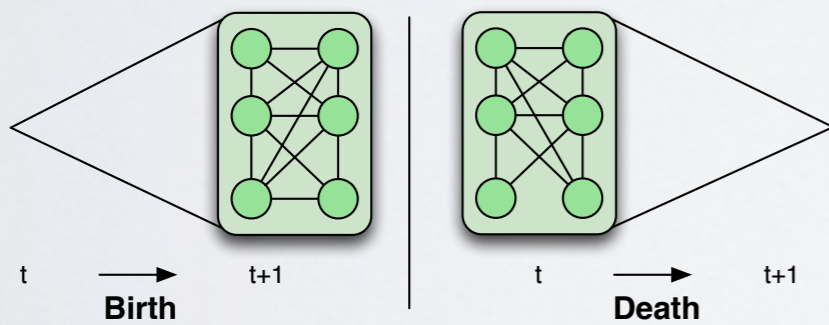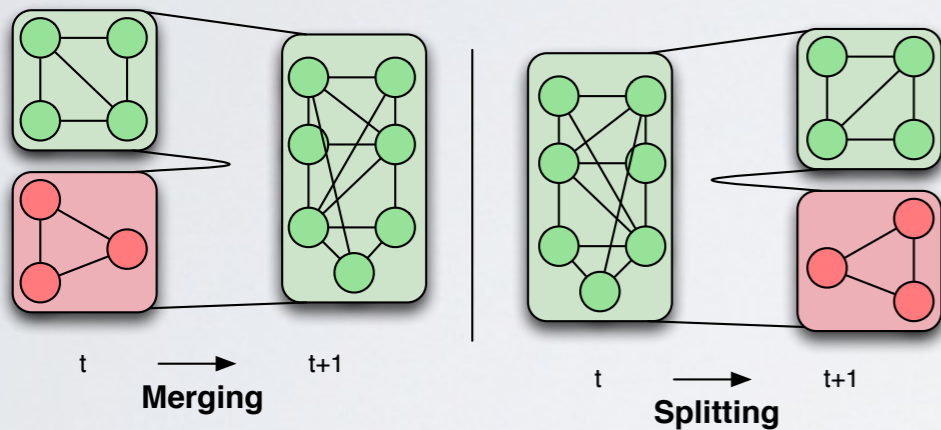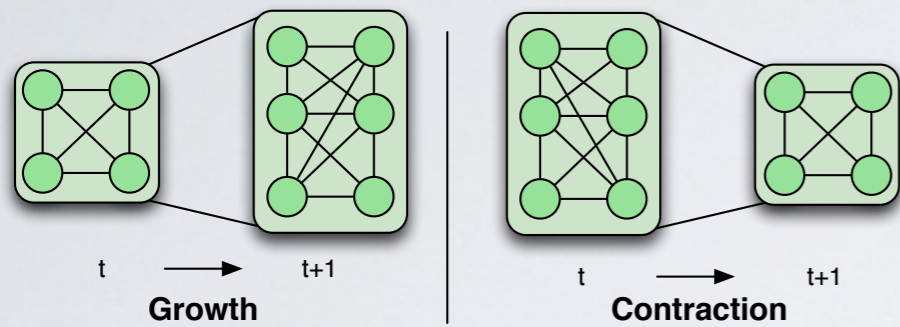
# LONGITUDINAL MODULARITY

**Static:** $Q(A, \mathcal{C}) = \displaystyle\sum_{C \in \mathcal{C}} \sum_{i,j \in C^2} \left[ \frac{A_{ij}}{2m} - \frac{k_i k_j}{4m^2} \right]$

**Longitudinal:** $Q_{\mathcal{L}}(L, \mathcal{C}) = \displaystyle\sum_{C \in \mathcal{C}} \sum_{u,v \in V^2} \left[ \frac{L^2_{uv \in C}}{2\mu} - \frac{\kappa_u \kappa_v}{4\mu^2} \frac{|T_{u \in C} \cap T_{v \in C}|^2}{|T|^2} \right]$
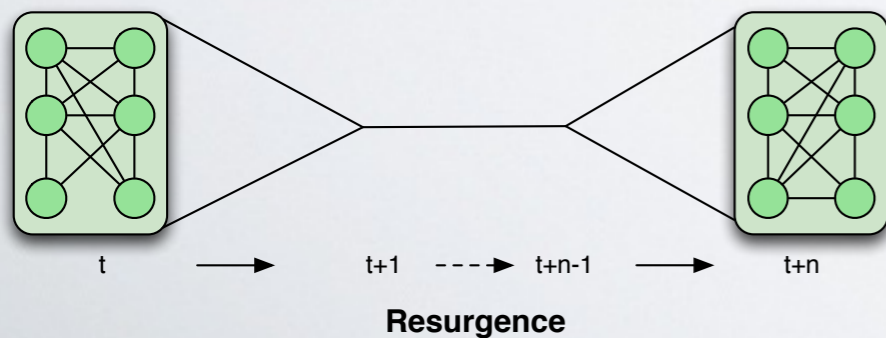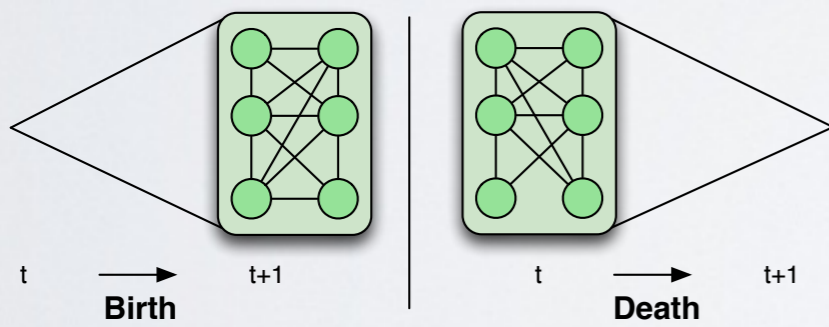
# ONGOING WORK 2

How to define community events quantitatively?
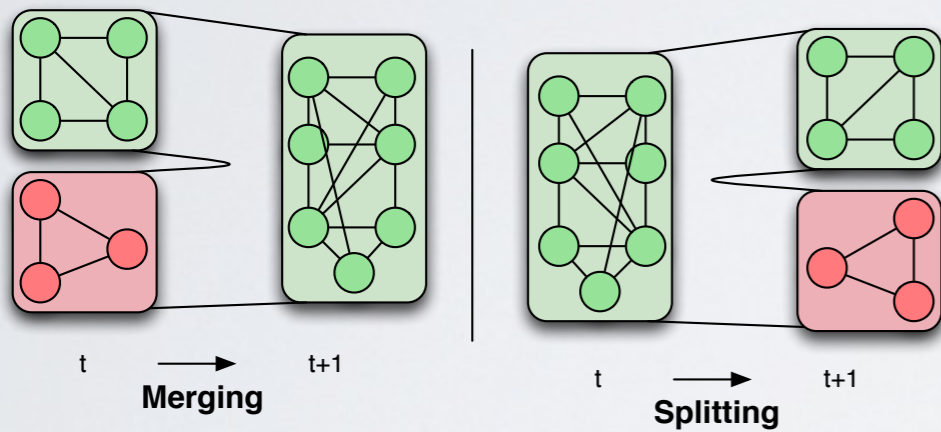
# COMMUNITY EVENTS



**The beautiful theory**

# COMMUNITY EVENTS



Growth

Contraction

Resurgence

**The ugly truth**

# COMMUNITY EVENTS

- Given the successive partitions of a dynamic graph
  - ‣ How to decide what events take place
  - ‣ For instance to study quantitatively
    - Do I have many merges?
    - Many splits?
    - Are large communities splitting more than small ones?
    - etc.

# COMMUNITY EVENTS

- Given the successive partitions of a dynamic graph
  - ‣ How to decide what events take place
  - ‣ For instance to study quantitatively
    - Do I have many merges?
    - Many splits?
    - Are large communities splitting more than small ones?
    - etc.

- Actually, the problem is not specific to communities in networks
  - ‣ Dynamic clustering

# COMMUNITY EVENTS

**Definition 1** (Flow Entropy). *Let $X$ be the target set, and $\mathcal{R} = \{R_1, \ldots, R_{|\mathcal{R}|}\}$ be the reference set. Let $B = \bigcup_{R \in \mathcal{R}} X \cap R$ identify the flow of $X$, namely the subset of $X$ shared with any of the elements of $\mathcal{R}$. Let $\sigma(b)$ be a function that maps each $b \in B$ to a unique identifier indicating the set $R \in \mathcal{R}$ such that $b \in R$. The Flow Entropy is defined as:*

$$\mathcal{H} = \begin{cases} -\sum_{b \in B} \frac{p(\sigma(b)) \log_2 p(\sigma(b))}{\log_2 |\mathcal{R}|} & \textit{if } |\mathcal{R}| \geq 2 \\ 0 & \textit{o/w} \end{cases} \tag{1}$$

The flow entropy quantifies the extent to which the nodes in $X$ come from one or multiple sets at time $t-1$. The flow entropy is bounded in $[0, 1]$ due to the normalizing factor $\log_2 |\mathcal{R}|$. The more $\mathcal{H}$ approaches 0, the fewer sets contribute to $X$, and vice versa. However, two special cases need further discussion:

# COMMUNITY EVENTS

**Definition 2** (Contribution Factor)**.**

$$\mathcal{W} = \frac{1}{|X|} \sum_{R \in \mathcal{R}} |R \cap X| \frac{|R \cap X|}{|R|} \tag{2}$$

$\mathcal{W}$ measures the extent to which the target set is composed by the contributing sets from $t-1$ provide elements to $X$ (respectively, the extent to which $X$ provides elements

# COMMUNITY EVENTS

**Definition 3** (Difference Factor).

$$\mathcal{D} = \frac{|X - \bigcup_{R \in \mathcal{R}} R|}{|X|} \tag{3}$$

The difference factor quantifies the fraction of members in $X$ that are not observed in the previous timestamp, i.e., the new, never-before-seen elements. In the remainder of this work, we will refer to these new elements as "joining" elements, as opposed to the elements belonging to the target set's flow.
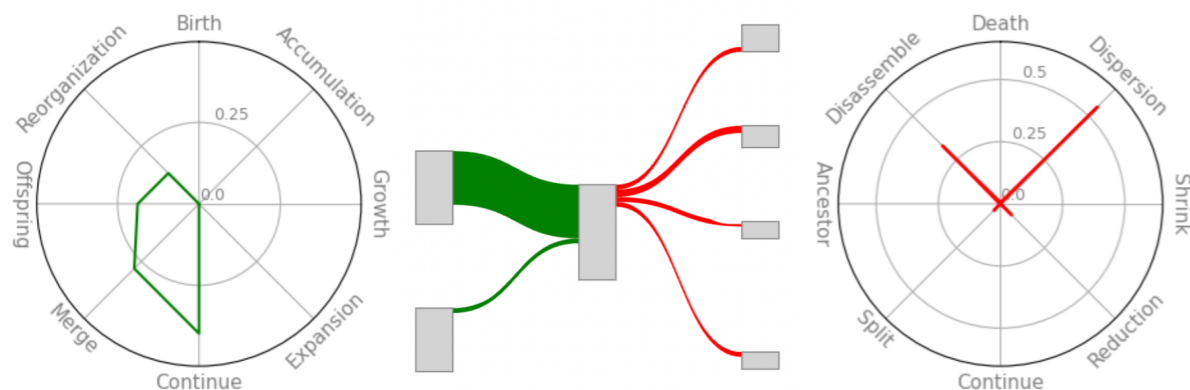
# COMMUNITY EVENTS

**Definition 5** (Backward Event Weights)**.** *Let $X$ be the target set and $\mathcal{R}$ be the reference set such that $X$ evolves from $\mathcal{R}$. Backward event weights quantify the extent to which $X$'s evolution from $\mathcal{R}$ approximates one of the following transformations:*

$$
\begin{aligned}
\text{BIRTH} &= (1 - \mathcal{H}) \cdot (1 - \mathcal{W}) \cdot \mathcal{D} \\
\text{ACCUMULATION} &= \mathcal{H} \cdot (1 - \mathcal{W}) \cdot \mathcal{D} \\
\text{CONTINUE} &= (1 - \mathcal{H}) \cdot \mathcal{W} \cdot (1 - \mathcal{D}) \\
\text{MERGE} &= \mathcal{H} \cdot \mathcal{W} \cdot (1 - \mathcal{D}) \\
\text{OFFSPRING} &= (1 - \mathcal{H}) \cdot (1 - \mathcal{W}) \cdot (1 - \mathcal{D}) \\
\text{REORGANIZATION} &= \mathcal{H} \cdot (1 - \mathcal{W}) \cdot (1 - \mathcal{D}) \\
\text{GROWTH} &= (1 - \mathcal{H}) \cdot \mathcal{W} \cdot \mathcal{D} \\
\text{EXPANSION} &= \mathcal{H} \cdot \mathcal{W} \cdot \mathcal{D}
\end{aligned}
$$

**Definition 6** (Forward Event Weights)**.** *Let $X$ be the target set and $\mathcal{R}$ be the reference set such that $X$ evolves into $\mathcal{R}$. Forward event weights quantify the extent to which $X$'s evolution into $\mathcal{R}$ approximates one of the following transformations:*

$$
\begin{aligned}
\text{DEATH} &= (1 - \mathcal{H}) \cdot (1 - \mathcal{W}) \cdot \mathcal{D} \\
\text{DISPERSION} &= \mathcal{H} \cdot (1 - \mathcal{W}) \cdot \mathcal{D} \\
\text{CONTINUE} &= (1 - \mathcal{H}) \cdot \mathcal{W} \cdot (1 - \mathcal{D}) \\
\text{SPLIT} &= \mathcal{H} \cdot \mathcal{W} \cdot (1 - \mathcal{D}) \\
\text{ANCESTOR} &= (1 - \mathcal{H}) \cdot (1 - \mathcal{W}) \cdot (1 - \mathcal{D}) \\
\text{DISASSEMBLE} &= \mathcal{H} \cdot (1 - \mathcal{W}) \cdot (1 - \mathcal{D}) \\
\text{SHRINK} &= (1 - \mathcal{H}) \cdot \mathcal{W} \cdot \mathcal{D} \\
\text{REDUCTION} &= \mathcal{H} \cdot \mathcal{W} \cdot \mathcal{D}
\end{aligned}
$$

# COMMUNITY EVENTS