# Faster maximal clique enumeration in large real-world link streams

Alexis BAUDIN

January 30, 2024

IRIF – Graph seminar
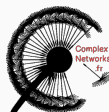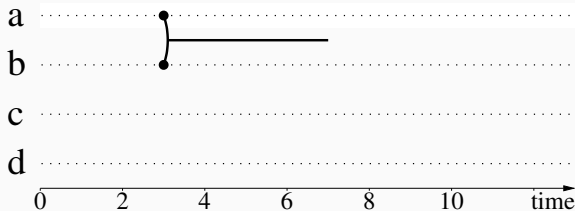
## Table of contents

Alexis BAUDIN    Faster maximal clique enumeration in link streams

# 1 - Maximal clique enumeration in link streams
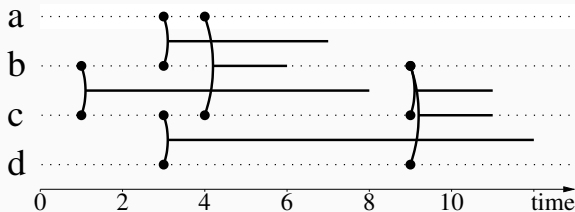
## Link stream

- Vertices : $a$, $b$, $c$ et $d$
- Time period : [0,12]
- Interaction : temporal links
  - $a$, $b$ linked over $[3, 7]$

## Advantages

- deals directly with the stream of interactions
- no arbitrary choice of time scale
- time is continuous
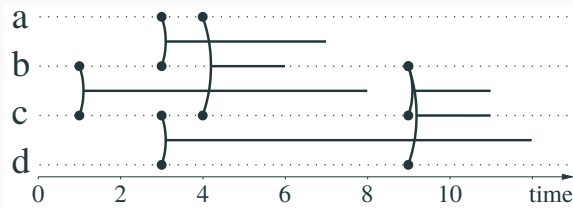
### Link stream

- Vertices : *a*, *b*, *c* et *d*
- Time period : [0,12]
- Interaction : temporal links
    - *a*, *b* linked over [3, 7]
    - ...

### Advantages

- deals directly with the stream of interactions
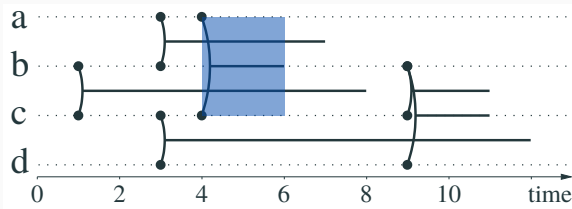- no arbitrary choice of time scale
- time is continuous

$(\{a, b, c\}, [4, 6])$ **is a clique**

$(\{a, b, c\}, [4, 6])$ **is a clique**

$\rightarrow$ $(\{a, b, c\}, \underline{[4, 5]})$ is not time-maximal

$(\{a, b, c\}, [4, 6])$ **is a clique**

$\rightarrow$ $(\{a, b, c\}, \underline{[4, 5]})$ is not time-maximal

$\rightarrow$ $(\underline{\{a, b\}}, [4, 6])$ is not vertex-maximal

$(\{a, b, c\}, [4, 6])$ **is a clique**

$\rightarrow$ $(\{a, b, c\}, [4, 5])$ is not time-maximal

$\rightarrow$ $(\{a, b\}, [4, 6])$ is not vertex-maximal

**maximal clique** = *time-maximal* and *vertex-maximal*

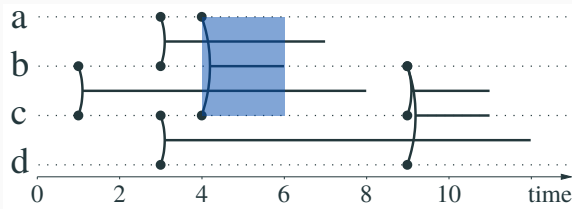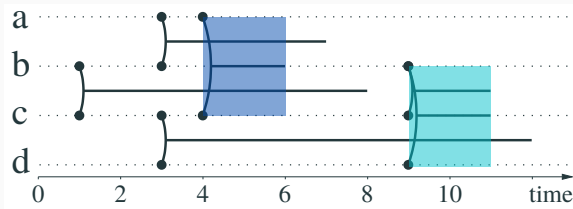$(\{a, b, c\}, [4, 6])$ **is a clique (maximal).**

$\rightarrow$ $(\{a, b, c\}, [4, 5])$ is not time-maximal

$\rightarrow$ $(\{a, b\}, [4, 6])$ is not vertex-maximal

**maximal clique** = *time-maximal* and *vertex-maximal*

$(\{a, b, c\}, [4, 6])$ **is a clique (<u>maximal</u>).**

$\rightarrow$ $(\{a, b, c\}, \underline{[4, 5]})$ is not time-maximal

$\rightarrow$ $(\underline{\{a, b\}}, [4, 6])$ is not vertex-maximal

**maximal clique** = *time-maximal* and *vertex-maximal*

$(\{a, b, c\}, [4, 6])$ **is a clique (maximal).**
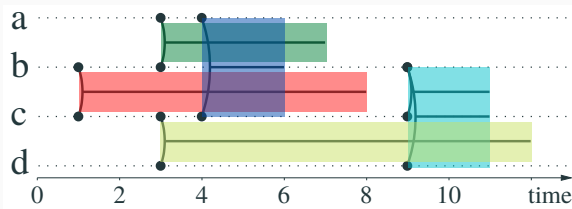
$\rightarrow$ $(\{a, b, c\}, [4, 5])$ is not time-maximal

$\rightarrow$ $(\{a, b\}, [4, 6])$ is not vertex-maximal

**maximal clique** = *time-maximal* and *vertex-maximal*

> *How to efficiently enumerate maximal cliques in massive real-world link streams?*

**Input:** A link stream.

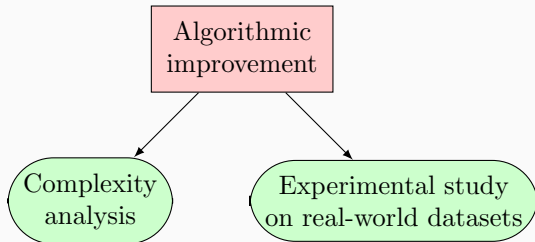**Output:** List of all its maximal cliques.

*How to efficiently enumerate maximal cliques in massive real-world link streams?*

**Input:** A link stream.

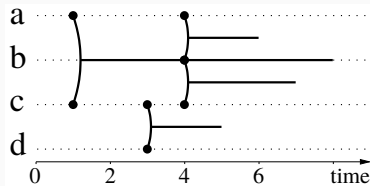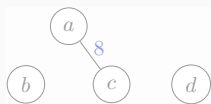**Output:** List of all its maximal cliques.

# 2 - New algorithm

**List of temporal links sorted chronologically**



$\downarrow$

**Evolving instantaneous graph $G_t$, with end dates**



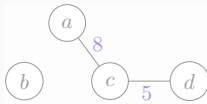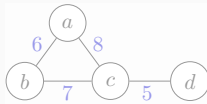$G_1 : t = 1$    $G_3 : t = 3$    $G_4 : t = 4$    $G_5 : t = 5$

**List of temporal links sorted chronologically**



↓

**Evolving instantaneous graph $G_t$, with end dates**



$G_1 : t = 1$  $G_3 : t = 3$  $G_4 : t = 4$  $G_5 : t = 5$

**List of temporal links sorted chronologically**



↓

**Evolving instantaneous graph $G_t$, with end dates**



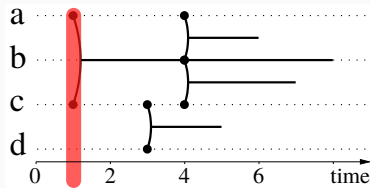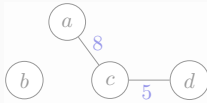$G_1 : t = 1$          $G_3 : t = 3$          $G_4 : t = 4$          $G_5 : t = 5$

**List of temporal links sorted chronologically**
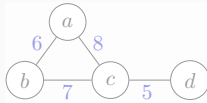


↓

**Evolving instantaneous graph $G_t$, with end dates**



$G_1 : t = 1$     $G_3 : t = 3$     $G_4 : t = 4$     $G_5 : t = 5$

**List of temporal links sorted chronologically**



↓

**Evolving instantaneous graph $G_t$, with end dates**



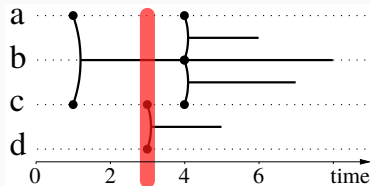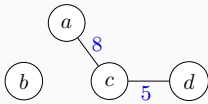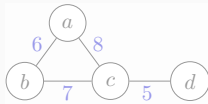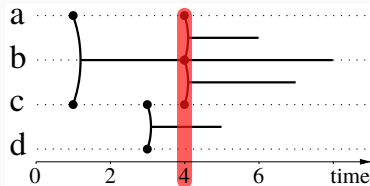$G_1 : t = 1$          $G_3 : t = 3$          $G_4 : t = 4$          $G_5 : t = 5$

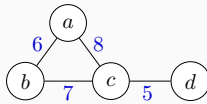Alexis Baudin    Faster maximal clique enumeration in link streams

**Instantaneous graph** $G_4$

Alexis BAUDIN    Faster maximal clique enumeration in link streams

## 2 - New algorithm
### > Maximal cliques that begin at each time $t$



**Instantaneous graph** $G_4$

## $>$ Maximal cliques that begin at each time $t$



**Instantaneous graph $G_4$**

$(\{b, c\}, \quad ), (\{a, b, c\}, \quad ), (\{a, b\}, \quad )$

$\rightarrow$ cliques of $G_4$ containing a new edge

Alexis Baudin    Faster maximal clique enumeration in link streams

**Instantaneous graph** $G_4$

$(\{b, c\}, [4, \quad), (\{a, b, c\}, [4, \quad), (\{a, b\}, [4, \quad)$

$\rightarrow$ cliques of $G_4$ containing a new edge

$\rightarrow$ starting time: $t = 4$

**Instantaneous graph $G_4$**

$(\{b, c\}, [4, 7]),\ (\{a, b, c\}, [4, 6]),\ (\{a, b\}, [4, 6])$

$\rightarrow$ cliques of $G_4$ containing a new edge

$\rightarrow$ starting time: $t = 4$

$\rightarrow$ end time: min of edge end times

**Instantaneous graph** $G_4$

$(\{b, c\}, [4, 7]), (\{a, b, c\}, [4, 6]),$ ~~$(\{a, b\}, [4, 6])$~~

$\rightarrow$ cliques of $G_4$ containing a new edge

$\rightarrow$ starting time: $t = 4$

$\rightarrow$ end time: min of edge end times

$\rightarrow$ filter vertex-maximal cliques

**Instantaneous graph $G_4$**

$(\{b, c\}, [4, 7])$, $(\{a, b, c\}, [4, 6])$, $(\{a, b\}, [4, 6])$

$\rightarrow$ cliques of $G_4$ containing a new edge

$\rightarrow$ starting time: $t = 4$

$\rightarrow$ end time: min of edge end times

$\rightarrow$ filter vertex-maximal cliques

**Back to a graph problem: clique enumeration in $G_t$.**

Clique enumeration
  in some $G_t$

Clique enumeration
in some $G_t$



$c$ chosen as pivot at first step

Clique enumeration
in some $G_t$



Example 1: do not prune $\{a; b; e\}$





$c$ chosen as pivot at first step

Clique enumeration
  in some $G_t$





$c$ chosen as pivot at first step

_Example 1:_ do not prune $\{a; b; e\}$



_Example 2:_ prune $\{a; b; d\}$

**Instantaneous graph** $G_3$

**Instantaneous graph** $G_3$

Alexis Baudin          Faster maximal clique enumeration in link streams

**Instantaneous graph** $G_3$

**Summary of the new algorithm**

$\rightarrow$ cliques not stored in memory $\checkmark$

$\rightarrow$ interactions reduced at each time step $\checkmark$

**Summary of the new algorithm**

$\rightarrow$ cliques not stored in memory $\checkmark$

$\rightarrow$ interactions reduced at each time step $\checkmark$

**State of the art : four main works**

- Viard *et al.* 2016
- Viard *et al.* 2018

  Store all cliques
  $\Rightarrow$ too costly in memory

- Himmel *et al.* 2017
- Bentert *et al.* 2019

  Need all past and future interactions
  when processing a vertex.

# 3 - Complexity analysis

### Input characteristics

$d$: maximal instantaneous degree

$m$: number of links

---

**Input characteristics**

  $d$: maximal instantaneous degree

  $m$: number of links

**Algorithm:** for each link ($u \xrightarrow{[t_0, t_1]} v$):

$\rightarrow$ List and process cliques in $G_{t_0}$ that contain $\{u, v\}$

- Number of those cliques: $\mathcal{O}\left(2^d\right)$
- Cost of computing and processing: $\mathcal{O}\left(d^2\right)$

---

**Input characteristics**

$d$: maximal instantaneous degree

$m$: number of links

**Algorithm:** for each link ($u \xrightarrow{[t_0, t_1]} v$):

$\rightarrow$ List and process cliques in $G_{t_0}$ that contain $\{u, v\}$

- Number of those cliques: $\mathcal{O}\left(2^d\right)$
- Cost of computing and processing: $\mathcal{O}\left(d^2\right)$

$\Rightarrow$ **Complexity:** $\boxed{\mathcal{O}\left(m \cdot d^2 \cdot 2^d\right)}$

**Instantaneous graph** $G_3$

**Instantaneous graph** $G_3$

**Instantaneous graph** $G_3$



**Ratio "good" leaves**

$$r = \frac{\text{nb maximal clique leaves}}{\text{nb leaves}}$$

Here: $r = \frac{1}{4}$; optimal: $r = 1$.

**Output characteristics**

$\alpha$: number of maximal cliques

$q$: maximal number of vertices in a clique

### Output characteristics

$\alpha$: number of maximal cliques

$q$: maximal number of vertices in a clique

**Complexity:** $\mathcal{O}\left((\text{nb nodes in the trees}) \cdot (\text{cost of a node})\right)$

**Output characteristics**

$\alpha$: number of maximal cliques

$q$: maximal number of vertices in a clique

**Complexity:** $\mathcal{O}\left((\text{nb nodes in the trees}) \cdot (\text{cost of a node})\right)$

$\rightarrow$ Nb nodes: $\mathcal{O}\left((\text{max depth}) \cdot (\text{nb leaves})\right)$

- max depth $= q$
- nb leaves $= \frac{1}{r} \cdot (\text{nb maximal clique leaves}) = \mathcal{O}\left(\frac{1}{r} \cdot \alpha\right)$

**Output characteristics**

  $\alpha$: number of maximal cliques

  $q$: maximal number of vertices in a clique

**Complexity:** $\mathcal{O}\left((\text{nb nodes in the trees}) \cdot (\text{cost of a node})\right)$

  $\rightarrow$ Nb nodes: $\mathcal{O}\left((\text{max depth}) \cdot (\text{nb leaves})\right)$
  - max depth $= q$
  - nb leaves $= \frac{1}{r} \cdot (\text{nb maximal clique leaves}) = \mathcal{O}\left(\frac{1}{r} \cdot \alpha\right)$
  $\rightarrow$ Cost of a node: $\mathcal{O}\left(d^2\right)$

**Output characteristics**

   $\alpha$: number of maximal cliques

   $q$: maximal number of vertices in a clique

**Complexity:** $\mathcal{O}\left((\text{nb nodes in the trees}) \cdot (\text{cost of a node})\right)$

$\rightarrow$ Nb nodes: $\mathcal{O}\left((\text{max depth}) \cdot (\text{nb leaves})\right)$
- max depth $= q$
- nb leaves $= \frac{1}{r} \cdot (\text{nb maximal clique leaves}) = \mathcal{O}\left(\frac{1}{r} \cdot \alpha\right)$

$\rightarrow$ Cost of a node: $\mathcal{O}\left(d^2\right)$

$\Rightarrow$ Complexity of the algorithm: $\boxed{\mathcal{O}\left(\frac{1}{r} \cdot \alpha \cdot q \cdot d^2\right)}$

**From input characteristics**

$$\mathcal{O}\left(m \cdot d^2 \cdot 2^d\right)$$

**From output characteristics**

$$\mathcal{O}\left(\frac{1}{r} \cdot d^2 \cdot q \cdot \alpha\right)$$

**From input characteristics**

$$\mathcal{O}\left(m \cdot d^2 \cdot 2^d\right)$$

- Viard *et al.* 2018: $\mathcal{O}\left(n^3 \cdot m^2 \cdot 2^n\right)$
- Himmel *et al.* 2018: $\mathcal{O}\left(m \cdot n \cdot |T| \cdot 3^{c/3} \cdot 2^c\right)$
- Bentert *et al.* 2019: $\mathcal{O}\left(n^4 \cdot |T|^2 \cdot 2^c\right)$

  ($c \leq d$ degeneracy, $|T|$ number of time steps)

**From output characteristics**

$$\mathcal{O}\left(\frac{1}{r} \cdot d^2 \cdot q \cdot \alpha\right)$$

**From input characteristics**

$$\mathcal{O}\left(m \cdot d^2 \cdot 2^d\right)$$

- Viard *et al.* 2018: $\mathcal{O}\left(n^3 \cdot m^2 \cdot 2^n\right)$
- Himmel *et al.* 2018: $\mathcal{O}\left(m \cdot n \cdot |T| \cdot 3^{c/3} \cdot 2^c\right)$
- Bentert *et al.* 2019: $\mathcal{O}\left(n^4 \cdot |T|^2 \cdot 2^c\right)$

  ($c \leq d$ degeneracy, $|T|$ number of time steps)

**From output characteristics**

$$\mathcal{O}\left(\tfrac{1}{r} \cdot d^2 \cdot q \cdot \alpha\right)$$

- output size $= \mathcal{O}\left(q \cdot \alpha\right) \Rightarrow$ factor $\frac{1}{r} \cdot d^2$ from output size

**From input characteristics**

$$\mathcal{O}\left(m \cdot d^2 \cdot 2^d\right)$$

- Viard *et al.* 2018: $\mathcal{O}\left(n^3 \cdot m^2 \cdot 2^n\right)$
- Himmel *et al.* 2018: $\mathcal{O}\left(m \cdot n \cdot |T| \cdot 3^{c/3} \cdot 2^c\right)$
- Bentert *et al.* 2019: $\mathcal{O}\left(n^4 \cdot |T|^2 \cdot 2^c\right)$

  ($c \leq d$ degeneracy, $|T|$ number of time steps)

**From output characteristics**

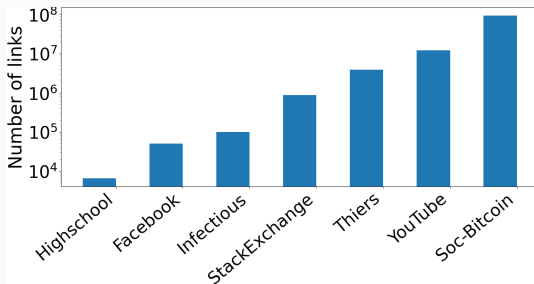$$\mathcal{O}\left(\tfrac{1}{r} \cdot d^2 \cdot q \cdot \alpha\right)$$

- output size $= \mathcal{O}\left(q \cdot \alpha\right) \Rightarrow$ factor $\tfrac{1}{r} \cdot d^2$ from output size
- $1 \leq \tfrac{1}{r} \leq 2^q$ but $\tfrac{1}{r} \approx 1.1$ in practice (experiments)

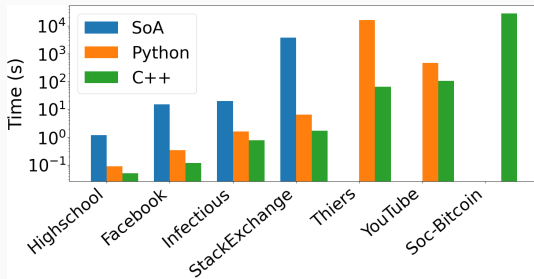**4 - Experimental study: performance gains**

**Datasets: state of the art + massive link streams**

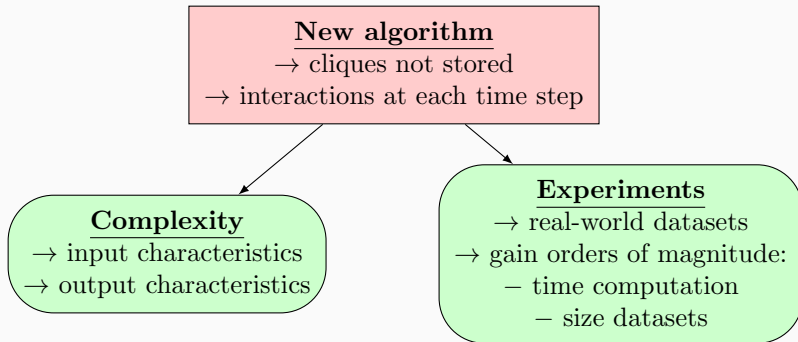- communication networks
- human interactions

**Experimental protocol**

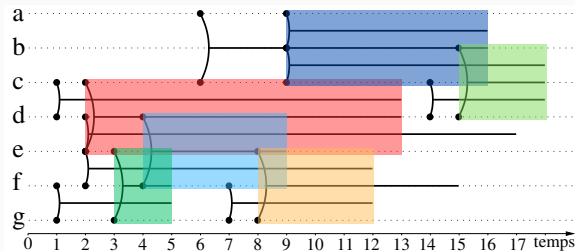- Code in `Python` and `C++`
- Maximum 24h and 390Gb RAM.

# 5 - Conclusion and perspectives

**Communities in link streams by clique percolation**

[Baudin et al. 2023]



$\rightarrow$ temporal data analysis
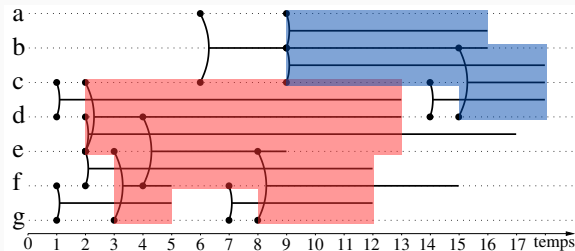
$\rightarrow$ anomaly detection

**Communities in link streams by clique percolation**

[Baudin et al. 2023]



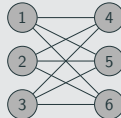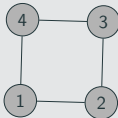$\rightarrow$ temporal data analysis

$\rightarrow$ anomaly detection

**Improve enumeration by ordering nodes**

Ordering the nodes of each instantaneous graph $G_t$.

[Eppstein *et al.* 2010]: core ordering of vertices.

**Improve enumeration by ordering nodes**

Ordering the nodes of each instantaneous graph $G_t$.

[Eppstein *et al.* 2010]: core ordering of vertices.

**Other temporal enumeration using a same framework**

Enumerate other motifs in each instantaneous graph $G_t$.
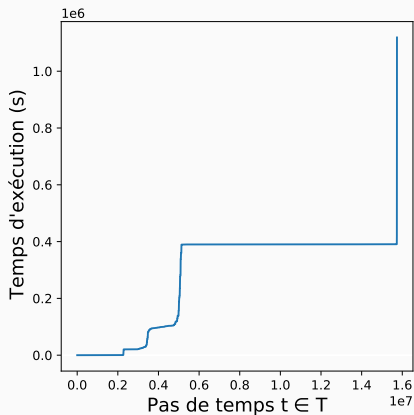
**Thank you for your attention!**

*Code in open-access:*
https://gitlab.lip6.fr/baudin/maxcliques-linkstream

*Contact:* alexis.baudin@lip6.fr

# Appendix

## Study the limits of computation

## Study the limits of computation