

Identification of logical websites

Pierre Senellart



1. Introduction

✓ What is a website?

2. Flow Simulation
3. Seed Extension
4. GEMO website
5. Conclusion

What is a logical website?

⑥ Simple idea: website = **webserver**

⑥ But:

- △ Some websites **span over** several webserver (e.g. `www-rocq.inria.fr`, `osage.inria.fr...`).
- △ Some webserver **contain** different websites (e.g. `perso.wanadoo.fr`).
- △ Some websites may be **split into** different websites (e.g. INRIA website).

⑥ Limits of a website: a **subjective** notion

1. Introduction

✓ Outline

2. Flow Simulation
3. Seed Extension
4. GEMO website
5. Conclusion

Talk outline

1. Introduction

2. Flow Simulation

- ✓ Maximal Flow / Minimal Cut
- ✓ Preflow-Push algorithm
- ✓ Adaptation to the Web

3. Seed Extension

- ✓ Markov CLustering algorithm
- ✓ Flow simulation from MCL

4. GEMO website

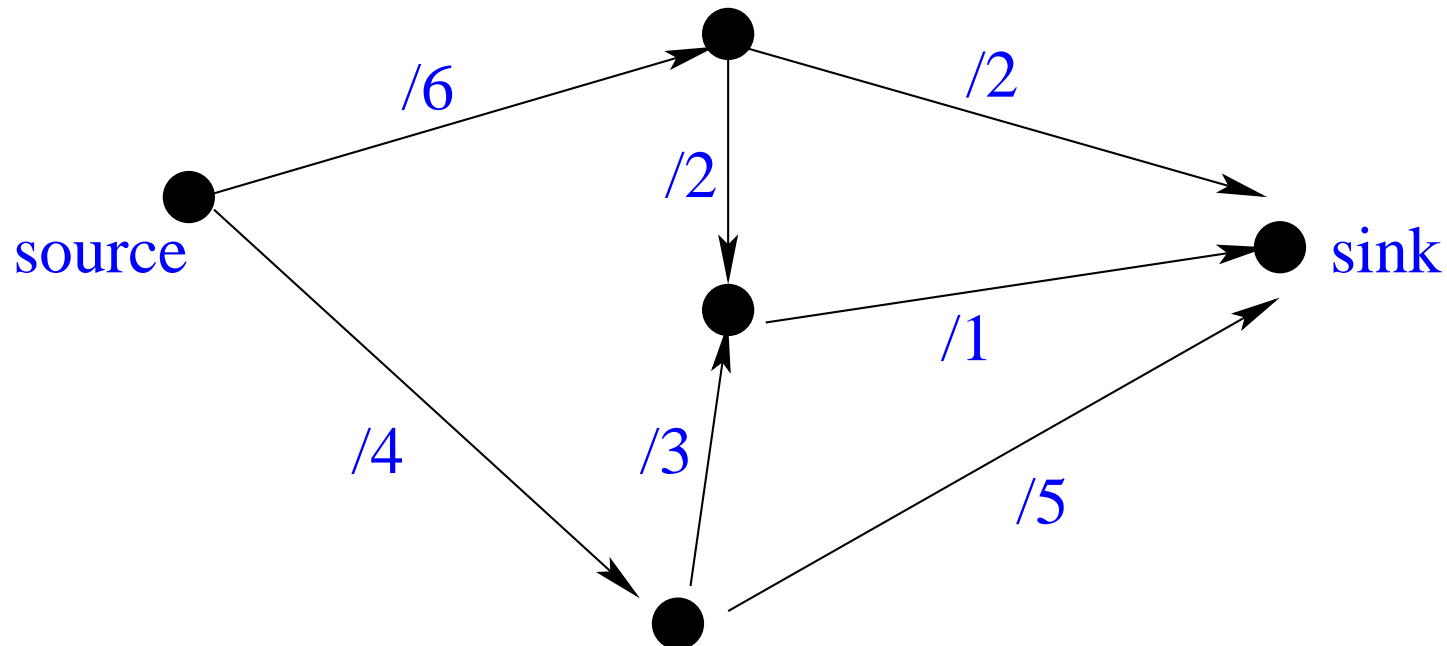
- ✓ Experiments
- ✓ Results

5. Conclusion

1. Introduction
2. Flow Simulation
✓ MaxFlow / MinCut
3. Seed Extension
4. GEMO website
5. Conclusion

Maximal Flow / Minimal Cut

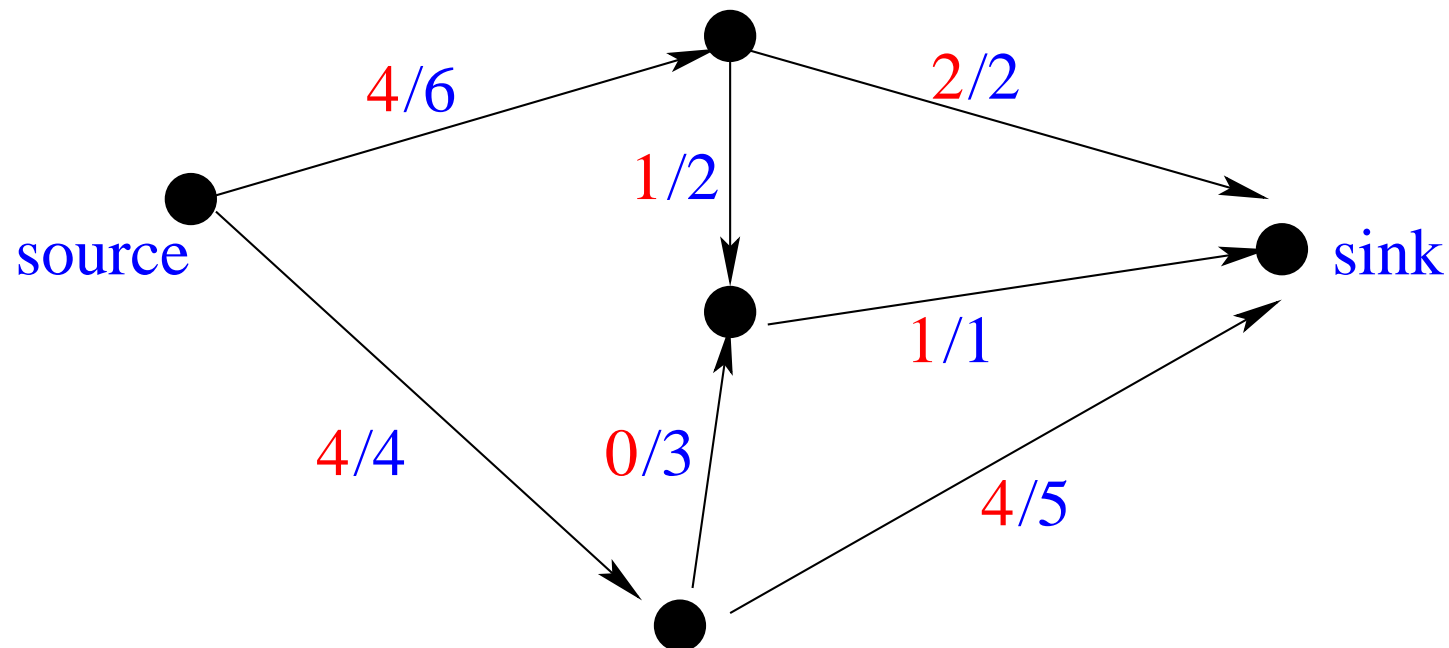
- ⑥ Transport network
- ⑥ Maximum flow \equiv Minimal cut



1. Introduction
2. Flow Simulation
✓ MaxFlow / MinCut
3. Seed Extension
4. GEMO website
5. Conclusion

Maximal Flow / Minimal Cut

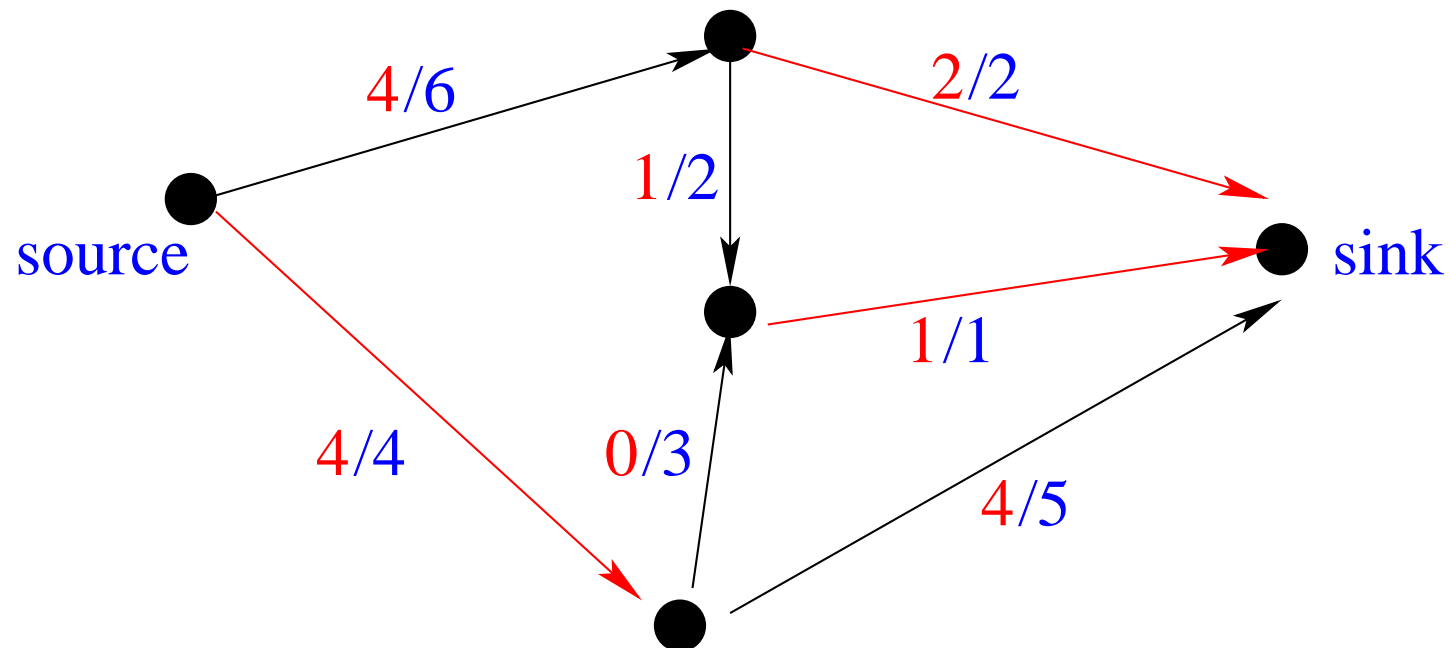
- ⑥ Transport network
- ⑥ Maximum flow \equiv Minimal cut



1. Introduction
2. Flow Simulation
✓ MaxFlow / MinCut
3. Seed Extension
4. GEMO website
5. Conclusion

Maximal Flow / Minimal Cut

- ⑥ Transport network
- ⑥ Maximum flow \equiv Minimal cut



1. Introduction
2. Flow Simulation
✓ Preflow-Push
3. Seed Extension
4. GEMO website
5. Conclusion

Preflow-Push algorithm

1. All nodes are assigned a **height** h : $h(\text{source}) = N$,
 $\forall k \neq \text{source}, h(k) = 0$ (N is the number of nodes)
2. Nodes with an **overflow** are visited, in any order.
 - ⌚ If possible, the flow is **pushed** toward a lower node.
Capacities of edges are respected.
 - ⌚ Otherwise, the node is **heightened**.

Theorem

*The process **converges**, whatever the sequence of visited nodes may be. The **maximal flow** is obtained at the limit.*

1. Introduction
2. Flow Simulation
✓ Adaptation to the Web
3. Seed Extension
4. GEMO website
5. Conclusion

Adaptation to the Web

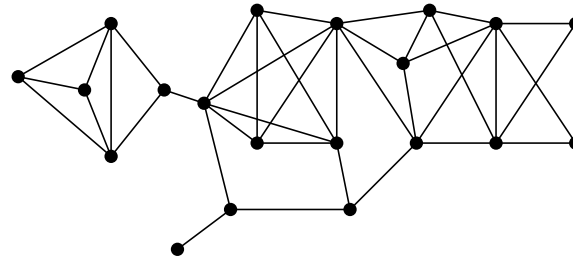
Website: nodes of a flow network delimited by a MaxFlow / MinCut.

- ⑥ **Nodes:** webpages, progressively crawled
- ⑥ **Edges:** hyperlinks
- ⑥ **Capacities:** edition distance between URLs
- ⑥ A **virtual source**, pointing to a **seed** of pages with infinite capacity edges
- ⑥ A **virtual sink**, pointed by all nodes with very low capacity edges

1. Introduction
2. Flow Simulation
3. Seed Extension
✓ MCL
4. GEMO website
5. Conclusion

Markov CLustering algorithm

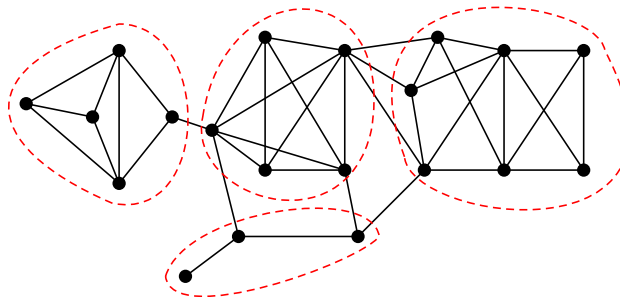
- ⑥ An **off-line** Graph Clustering Algorithm based on flow simulation
- ⑥ Alternation between **expansion** (multiplication) and **inflation** (rescaling) of a stochastic matrix
- ⑥ Work (almost) only on **undirected** graphs
- ⑥ **Complexity:** $O(n)$ (but with a high multiplying factor)



1. Introduction
2. Flow Simulation
3. Seed Extension
✓ MCL
4. GEMO website
5. Conclusion

Markov CLustering algorithm

- ⑥ An **off-line** Graph Clustering Algorithm based on flow simulation
- ⑥ Alternation between **expansion** (multiplication) and **inflation** (rescaling) of a stochastic matrix
- ⑥ Work (almost) only on **undirected** graphs
- ⑥ **Complexity:** $O(n)$ (but with a high multiplying factor)



1. Introduction
2. Flow Simulation
3. Seed Extension
✓ F.S. from MCL
4. GEMO website
5. Conclusion

Flow Simulation from MCL clusters

Process:

1. **MCL Clustering** of a large relevant portion of the graph
2. Identification of the **most relevant cluster**
3. **Flow Simulation** starting from this cluster

Advantages over pure MCL:

- ⑥ **Dynamic** discovery of clusters
- ⑥ Use the fact that the graph is **directed**

1. Introduction
2. Flow Simulation
3. Seed Extension
4. GEMO website
✓ Experiments
5. Conclusion

Experiments

- | | |
|--|-----------------|
| ⑥ Crawl of a large part of <code>*.inria.fr/*</code> | <i>72 hours</i> |
| ⑥ MCL clustering of the obtained graph | <i>24 hours</i> |
| ⑥ Identification of the GEMO cluster | <i>< 1s</i> |
| ⑥ Flow Simulation from this cluster | <i>3 hours</i> |

87,140 webpages in the graph.
276 clusters.

1. Introduction
2. Flow Simulation
3. Seed Extension Results
4. GEMO website
✓ Results
5. Conclusion

Results

	Size	Precision	Recall	THESUS
Flow Simulation	8	87.5	1.3	xml web
MCL	320	99.7	33.0	gemo report
MCL + Flow Simulation	788	90.4	86.4	gemo report
<code>www-rocq.inria.fr/ verso/*</code>	221	100.0	44.4	diapositive texte
<code>{www-rocq,osage}. inria.fr/verso/*</code>	683	100.0	68.6	report diapositive

1. Introduction
2. Flow Simulation
3. Seed Extension
4. GEMO website
5. Conclusion
✓ In Brief

In Brief

- ⑥ Website: a **subjective, not obvious** notion
- ⑥ Flow Simulation used to **discover the limits** of a website
- ⑥ Best results given by a **combination** of an **off-line graph clustering** and an **on-line flow simulation**

1. Introduction
2. Flow Simulation
3. Seed Extension
4. GEMO website
5. Conclusion
✓ Future works

Future works

- ⑥ Optimization, scaling
- ⑥ On-line computing of MCL?
- ⑥ Hierarchical clustering
- ⑥ Application of the Preflow-Push algorithm to peer-to-peer networks

1. Introduction
2. Flow Simulation
3. Seed Extension
4. GEMO website
5. Conclusion
✓ Future works

Future works

- ⑥ Optimization, scaling
- ⑥ On-line computing of MCL?
- ⑥ Hierarchical clustering
- ⑥ Application of the Preflow-Push algorithm to peer-to-peer networks

THANK YOU FOR YOUR ATTENTION!