

Identification de sites Web logiques

Pierre Senellart



1. Introduction

✓ Qu'est-ce qu'un site Web ?

2. Propagation de flots
3. Extension de la graine
4. Site Web de GEMO
5. Conclusion

Qu'est-ce qu'un site Web logique ?

⑥ Idée simple : site Web = **serveur Web**.

⑥ Mais :

- △ Certains sites Web **couvrent plusieurs serveurs** (ex. `www-rocq.inria.fr`, `osage.inria.fr...`).
- △ Certains serveurs Web **contiennent des sites différents** (ex. `perso.wanadoo.fr`).
- △ Qu'est-ce qu'un serveur Web en présence d'**hôtes virtuels**, de **miroirs** ?
- △ Certains sites Web peuvent être **décomposés en différents sous-sites** (ex. le site Web de l'INRIA).

1. Introduction

✓ Plan

2. Propagation de flots
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Plan de l'exposé

1. Introduction

2. Propagation de flots

- ✓ Flot maximal / Coupe minimale
- ✓ Algorithme de préflots
- ✓ Adaptation au Web

3. Extension de la graine

- ✓ Markov CLustering algorithm
- ✓ Propagation de flots à partir de MCL

4. Site de GEMO

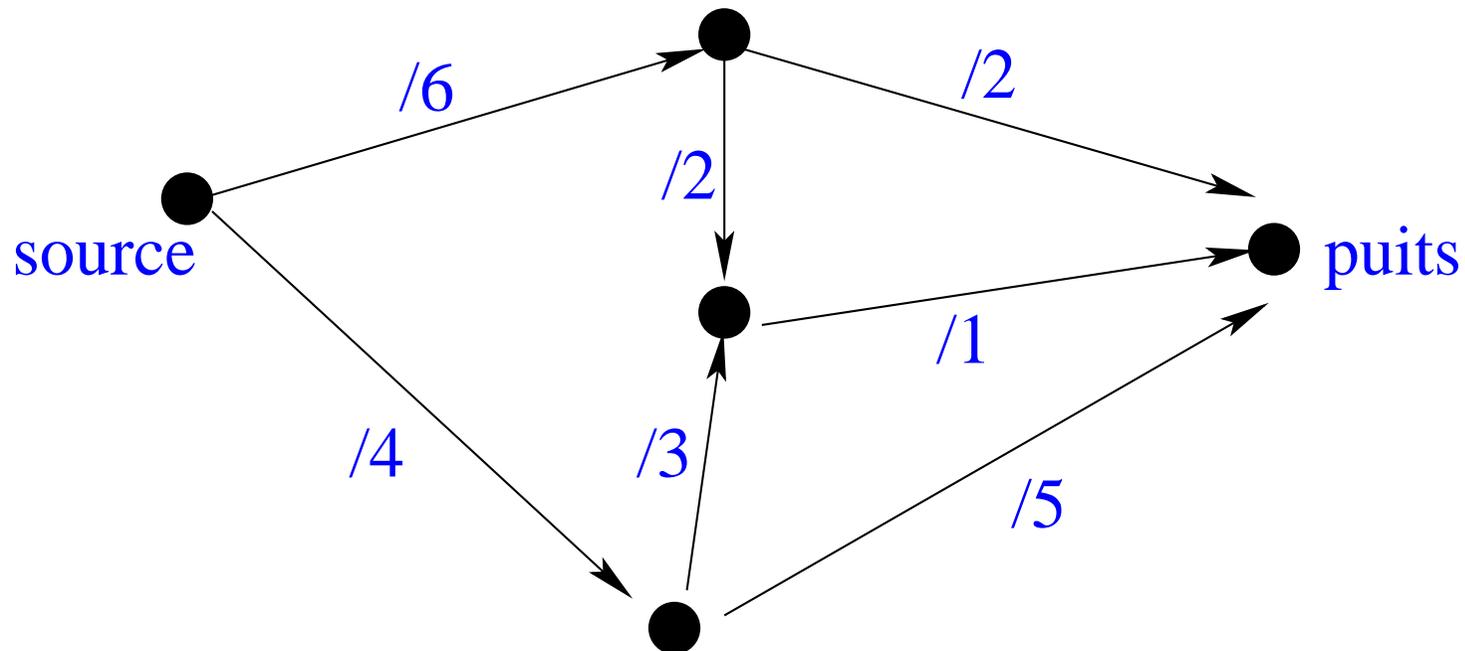
- ✓ Expériences
- ✓ Résultats

5. Conclusion

1. Introduction
2. Propagation de flots
✓ Flot max. / Coupe min.
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Flot maximal / Coupe minimale

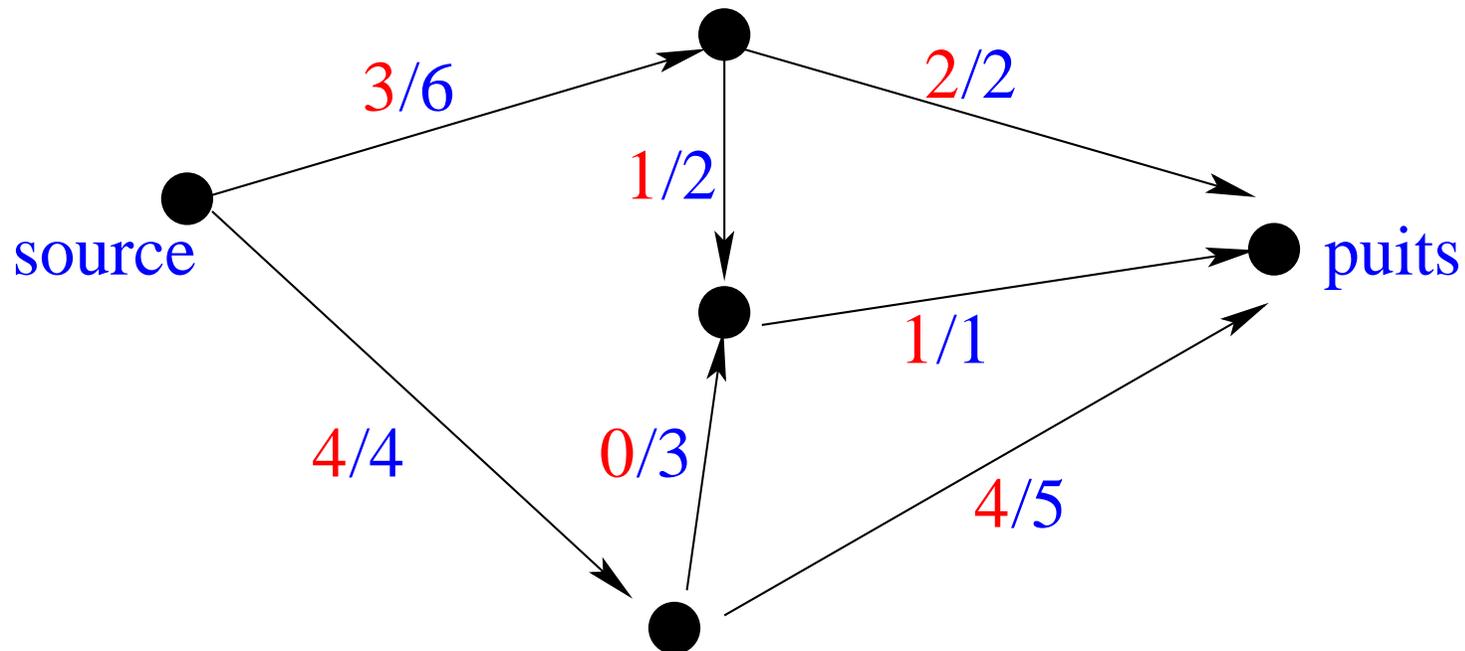
- ⑥ Réseau de transport.
- ⑥ Flot maximal \equiv Coupe minimale.



1. Introduction
2. Propagation de flots
✓ Flot max. / Coupe min.
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Flot maximal / Coupe minimale

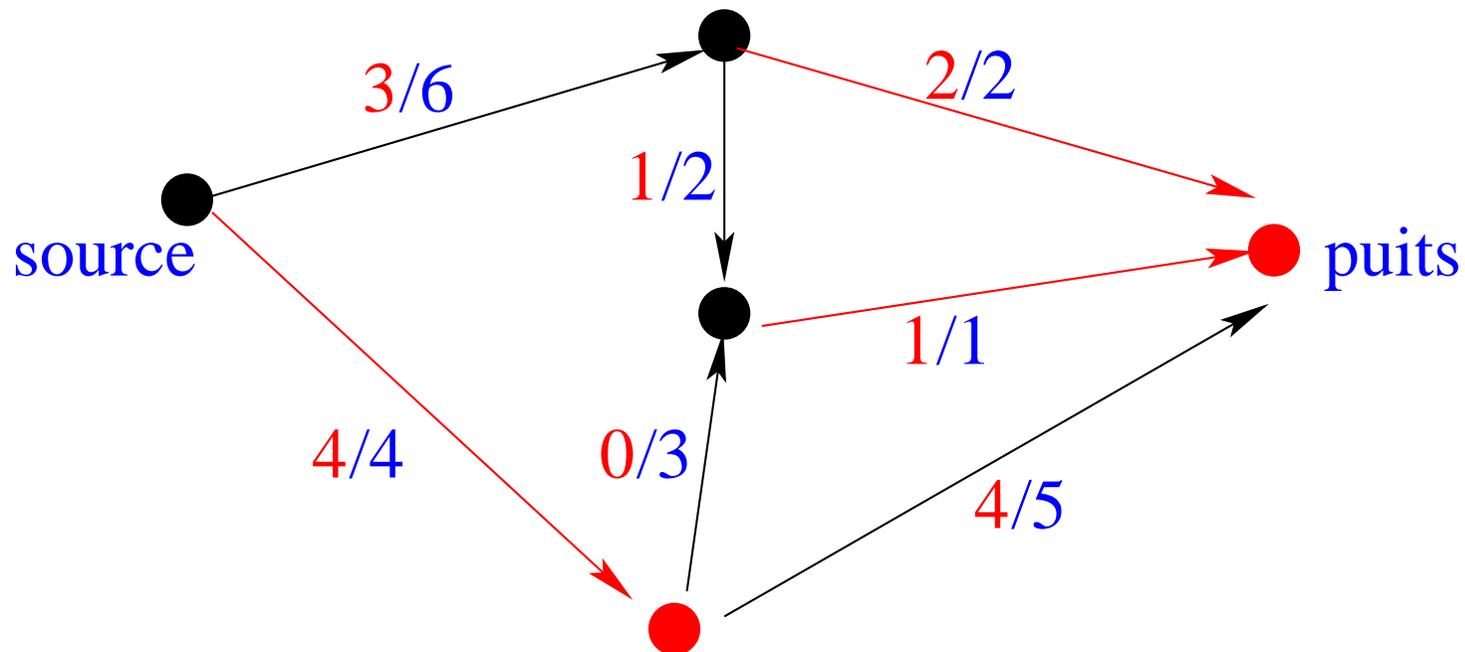
- ⑥ Réseau de transport.
- ⑥ Flot maximal \equiv Coupe minimale.



1. Introduction
2. Propagation de flots
✓ Flot max. / Coupe min.
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Flot maximal / Coupe minimale

- ⑥ Réseau de transport.
- ⑥ Flot maximal \equiv Coupe minimale.



1. Introduction
2. Propagation de flots
✓ Préflots
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Algorithme de préflots

1. On affecte une hauteur initiale h à tous les nœuds :
 $h(\text{source}) = N$ et $\forall k \neq \text{source}, h(k) = 0$
(N est le nombre de nœuds).
2. Les nœuds avec excédent de flot sont visités, dans un ordre quelconque.
 - ⑥ Si possible, le flot est poussé vers un nœud plus bas. Les capacités des arêtes sont respectées.
 - ⑥ Sinon, on élève le nœud.

Théorème

Le processus converge, quel que soit l'ordre de visite des nœuds. Un flot maximal est obtenu à la limite.

1. Introduction
2. Propagation de flots
✓ Adaptation au Web
3. Extension de la graine
4. Site de GEMO
5. Conclusion

Adaptation au Web

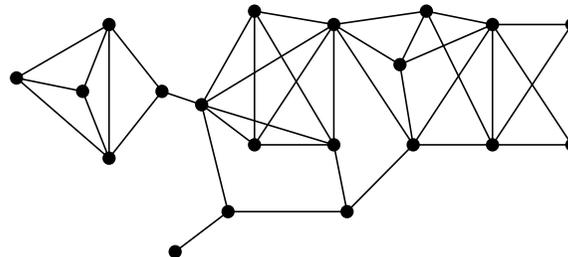
Site Web : Nœuds d'un réseau de transport délimité par un Flot max. / Coupe min.

- ⑥ **Nœuds** : pages Web, parcourues progressivement.
- ⑥ **Arêtes** : hyperliens.
- ⑥ **Capacité** : fonction décroissante des distances d'édition entre les URLs.
- ⑥ Une source virtuelle, pointant vers une graine de pages avec des arêtes de capacité infinie.
- ⑥ Un puits virtuel, pointé par tous les nœuds avec des arêtes de très faible capacité.

1. Introduction
2. Propagation de flots
3. Extension de la graine
✓ MCL
4. Site de GEMO
5. Conclusion

Markov CLustering algorithm

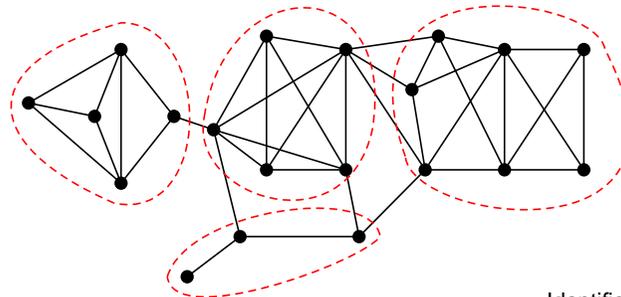
- ⑥ Un algorithme **hors-ligne** de segmentation de graphe, basé sur la propagation de flots.
- ⑥ Alternativement, **expansion** (multiplication) et **inflation** (renormalisation) d'une matrice stochastique.
- ⑥ Fonctionne (presque) seulement sur des graphes **non dirigés**.
- ⑥ **Complexité** : $O(n)$ (mais avec un important facteur).



1. Introduction
2. Propagation de flots
3. Extension de la graine
✓ MCL
4. Site de GEMO
5. Conclusion

Markov CLustering algorithm

- ⑥ Un algorithme **hors-ligne** de segmentation de graphe, basé sur la propagation de flots.
- ⑥ Alternativement, **expansion** (multiplication) et **inflation** (renormalisation) d'une matrice stochastique.
- ⑥ Fonctionne (presque) seulement sur des graphes **non dirigés**.
- ⑥ **Complexité** : $O(n)$ (mais avec un important facteur).



1. Introduction
2. Propagation de flots
3. Extension de la graine
✓ Propagation de flots à partir de MCL
4. Site de GEMO
5. Conclusion

Propagation de flots à partir des classes de MCL

Processus :

1. **Segmentation MCL** d'une grande partie (appropriée) du graphe.
2. Identification de la **classe la plus pertinente**.
3. **Propagation de flots** à partir de cette classe.

Avantages par rapport à MCL seul :

- ⑥ Découverte **dynamique** de classes.
- ⑥ Utilise le fait que le graphe est **dirigé**.

1. Introduction
2. Propagation de flots
3. Extension de la graine
4. Site de GEMO
✓ Expériences
5. Conclusion

Expériences

- ⑥ Parcours d'une large partie de `*.inria.fr/*` 72 heures
- ⑥ Segmentation par MCL du graphe obtenu 24 heures
- ⑥ Identification de la classe GEMO < 1s
- ⑥ Propagation de flots à partir de cette classe 3 heures

87 140 pages Web dans le graphe.

276 classes.

(5 247 en segmentant à nouveau la classe la plus grande).

1. Introduction
2. Propagation de flots
3. Extension de la graine Results
4. Site de GEMO
✓ Résultats
5. Conclusion

Résultats

	Taille	Précision	Rappel
Propagation de flots	8	87.5%	1.3%
MCL	320	99.7%	33.0%
MCL + Propagation de flots	788	90.4%	86.4%
<code>www-rocq.inria.fr/verso/*</code>	221	100.0%	44.4%
<code>{www-rocq,osage}.inria.fr/verso/*</code>	683	100.0%	68.6%

1. Introduction
2. Propagation de flots
3. Extension de la graine
4. Site de GEMO
5. Conclusion
✓ En Bref

En bref

- ⑥ Site Web : une notion **subjective**, **non évidente**
- ⑥ Propagation de flots utilisée pour **découvrir les limites** d'un site Web
- ⑥ Meilleurs résultats obtenus par la **combinaison** d'une segmentation de graphe hors-ligne et d'une propagation de flots en-ligne

1. Introduction
2. Propagation de flots
3. Extension de la graine
4. Site de GEMO
5. Conclusion
✓ Perspectives

Perspectives

- ⑥ Optimisation, passage à l'échelle.
- ⑥ Calcul en-ligne de MCL ?
- ⑥ Choix d'une stratégie de parcours efficace.
- ⑥ Combinaison avec une méthode sémantique.

1. Introduction
2. Propagation de flots
3. Extension de la graine
4. Site de GEMO
5. Conclusion
✓ Perspectives

Perspectives

- ⑥ Optimisation, passage à l'échelle.
- ⑥ Calcul en-ligne de MCL ?
- ⑥ Choix d'une stratégie de parcours efficace.
- ⑥ Combinaison avec une méthode sémantique.

MERCI POUR VOTRE ATTENTION !