



Recherche de motifs dans des graphes colorés

Florian Sikora (encadré par Guillaume Blin et Stéphane Vialette)

Université Paris-Est, LIGM - UMR CNRS 8049

Séminaire Combi LINA - 06/2010

Plan

Introduction

Complexité paramétrée et motifs avec topologie

Motifs sans topologie

Le problème GRAPH MOTIF Des logiciels pour GRAPH MOTIF

Conclusion

Plan

Introduction

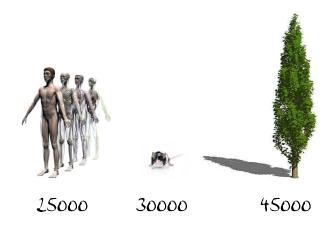
Complexité paramétrée et motifs avec topologie

Motifs sans topologie

Le problème GRAPH MOTIF Des logiciels pour GRAPH MOTIF

Conclusion

Motivations



- ► Complexité de l'homme # de gènes ?
- ▶ Complexité de l'homme⇔ protéines ?

Les protéines..

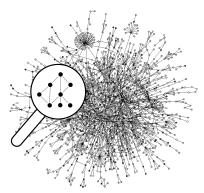
- Nouveaux interêts concernant les protéines...
- ... et sur leurs interactions: Protein-Protein Interactions (PPI)
- ▶ Obtenues biologiquement... avec beaucoup de bruit!

Réseau de protéines

Les protéines peuvent interagir avec d'autres protéines



Réseau de protéines



- Modélisation par un graphe (éventuellement pondéré)
 - ▶ Les protéines sont représentées par les nœuds
 - Les interactions sont représentées par les arêtes
 - Les arêtes peuvent être pondérées par la probabilité de l'interaction

Motivations

- Nouvelles techniques : l'information augmente très rapidement [Sharan & Ideker 2006]
 - ▶ 2001: quelques centaines d'interactions
 - 2006: plusieurs milliers
- ▶ Beaucoup de BDD

Motivations

- Nouvelles techniques : l'information augmente très rapidement [Sharan & Ideker 2006]
 - 2001: quelques centaines d'interactions
 - 2006: plusieurs milliers
- ▶ Beaucoup de BDD
- Chercher des motifs pour retrouver des fonctions connues
- Déduire les informations d'espèces peu connues depuis des espèces bien connues

Plan

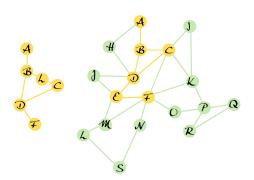
Complexité paramétrée et motifs avec topologie

Motifs sans topologie

Le problème GRAPH MOTIF Des logiciels pour GRAPH MOTIF

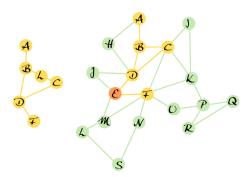
Rechercher des motifs

- Rechercher des motifs (ici, ensemble de protéines avec une topologie) dans un réseau PPI
- ▶ Une protéine est dite homologue à une autre protéine selon une analyse de séquences (avec BLASTp)



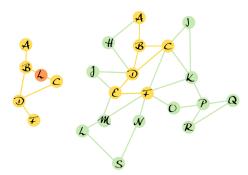
Rechercher des motifs

- Rechercher des motifs (ici, ensemble de protéines avec une topologie) dans un réseau PPI
- Une protéine est dite homologue à une autre protéine selon une analyse de séquences (avec BLASTp)
- Un nombre borné d'insertions et deletions peut-être autorisé



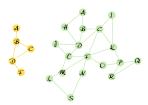
Rechercher des motifs (ici, ensemble de protéines avec une topologie) dans un réseau PPI

- Une protéine est dite homologue à une autre protéine selon une analyse de séquences (avec BLASTp)
- Un nombre borné d'insertions et deletions peut-être autorisé



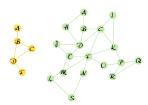
Recherche avec topologie dans un réseau PPI

- ► Le motif peut être un chemin, un arbre, un graphe
- Problèmes NP-complets donc une solution exacte entraine une complexité exponentielle

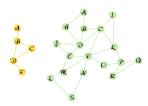


Recherche avec topologie dans un réseau PPI

- ► Le motif peut être un chemin, un arbre, un graphe
- Problèmes NP-complets donc une solution exacte entraine une complexité exponentielle
- Mais tout n'est pas perdu!



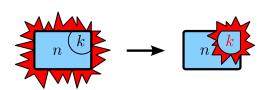
- ▶ Idée: exploiter le fait que les **motifs sont plus petits** $(\sim 5 15)$ que le réseau (e.g. ~ 5.000 pour la levure)
- ▶ Restreindre la partie exponentielle à k (taille motif) au lieu de n (taille du réseau): complexité paramétrée



- ▶ Beaucoup de problèmes (i.e. problèmes paramétrés) sont de la forme :
 - ▶ Input : Un objet X, |X| = n, un entier k
 - Question : Est-ce que X a une propriété dépendant de k ?

- Beaucoup de problèmes (i.e. problèmes paramétrés) sont de la forme :
 - ▶ Input : Un objet X, |X| = n, un entier k
 - Question : Est-ce que X a une propriété dépendant de k ?
- Exemples :
 - ▶ (VERTEX COVER) Le graphe G = (V, E) contient-il un **sous-ensemble** V' de sommets **de taille** k t.q. pour chaque arête (u, v) de G, soit u soit v est dans V'?
 - ► (LONGEST COMMON SUBSEQUENCE) Existe t-il une string de taille au moins k qui soit une sous-séquence de n strings?
 - ► (SET COVER) Etant donné un ensemble *S* de *n* ensembles, existe-il un sous-ensemble $S' \subseteq S$ de k ensembles t.g. chaque élément dans les ensembles de S est présent dans un ensemble de S'?

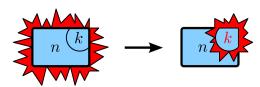
- Un algorithme FPT [DOWNEY & FELLOWS 1999]: algorithme exact exponentiel seulement en son paramètre k (et pas en la taille de l'entrée n)
- ► $f(k).n^c$, avec c une constante, et f n'importe quelle fonction $(e.g. \mathcal{O}(2^k.n^2))$
- ► L'algorithme devient souvent "praticable", car même si *f* est exponentiel, *k* est petit



▶ Un algorithme FPT [Downey & Fellows 1999]: algorithme exact **exponential** seulement en son

paramètre k (et pas en la taille de l'entrée n)

- $ightharpoonup f(k).n^c$, avec c une constante, et f n'importe quelle fonction (e.g. $\mathcal{O}(2^k.n^2)$)
- ▶ L'algorithme devient souvent "praticable", car même si *f* est exponentiel, k est petit
- ► Attention, 2^{22^{22²}} .n est FPT mais rédibitoire même pour k = 1



- Si un problème est W[1]-difficile, il n'y a probablement pas d'algorithme FPT possible pour ce problème (plus ou moins l'analogue de la classe NP)
- Se prouve avec des réductions paramétrées (préservation de la taille de l'instance et du paramètre)

- ► Probleme : Trouver un chemin de taille *k* sans cycle dans un graphe avec *n* noeuds
- ▶ NP-Complet (Chemin hamiltonien où k = n)
- ▶ Algorithme naïf : n^k (tous les chemins de taille k)

- (Randomized) Color-coding :
 - Choisir k couleurs differentes
 - Choisir aléatoirement une couleur pour chaque noeud de G





- (Randomized) Color-coding :
 - Choisir k couleurs differentes
 - Choisir aléatoirement une couleur pour chaque noeud de G
 - ▶ Chercher un chemin "colorful" sur les k couleurs \Rightarrow chemin sans cycle de taille k
 - ▶ S'effectue en seulement $\mathcal{O}(2^k)$ (1 bit pour chaque couleur)



Couleurs introduites



 $D(v, S) = \bigvee_{u} D(u, S \setminus \{col(u)\})$ t.q. u voisin de v et $col(u) \in S$

- (Randomized) Color-coding:
 - Choisir k couleurs differentes
 - ► Choisir aléatoirement une couleur pour chaque noeud de *G*
 - ▶ Chercher un chemin "colorful" sur les k couleurs \Rightarrow chemin sans cycle de taille k
 - ▶ S'effectue en seulement $\mathcal{O}(2^k)$ (1 bit pour chaque couleur)
 - ▶ Un "bon chemin" est colorful avec la probabilité $\frac{k!}{kk} \geq e^{-k}$



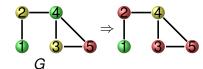


- (Randomized) Color-coding:
 - Choisir k couleurs differentes
 - ► Choisir aléatoirement une couleur pour chaque noeud de *G*
 - ▶ Chercher un chemin "colorful" sur les k couleurs \Rightarrow chemin sans cycle de taille k
 - ▶ S'effectue en seulement $\mathcal{O}(2^k)$ (1 bit pour chaque couleur)
 - ▶ Un "bon chemin" est colorful avec la probabilité $\frac{k!}{kk} \geq e^{-k}$
 - Donc, relancer la procédure e^k fois pour avoir une probabilité d'erreur faible





Couleurs introduites



- (Randomized) Color-coding :
 - ▶ Choisir *k* couleurs differentes
 - Choisir aléatoirement une couleur pour chaque noeud de G
 - ► Chercher un chemin "colorful" sur les k couleurs \Rightarrow chemin sans cycle de taille k
 - ▶ S'effectue en seulement $\mathcal{O}(2^k)$ (1 bit pour chaque couleur)
 - ▶ Un "bon chemin" est colorful avec la probabilité $\frac{k!}{k^k} \geq e^{-k}$
 - Donc, relancer la procédure e^k fois pour avoir une probabilité d'erreur faible
 - ► Exponentiel en k seulement!

Recherche de motifs

- Motif est un chemin : QPath, O(2^k|E|) [SHLOMI ET AL. 2006]
 - Application directe du color-coding
- ▶ Motif est un arbre : QNet, $\mathcal{O}^*(2^{\mathcal{O}(k)})$ [Dost et al. 2007]
 - Extension du color-coding (programmation dynamique sur l'arbre)
- Motif est un graphe : W[1]-dur [Downey & Fellows 1999] (pas d'algo FPT possible)
 - ► Mais FPT si le motif est un graphe à treewidth bornée, ou à Vertex Feedback Set borné [DOST ET AL. 2007, BLIN ET AL. 2009]

Le problème GRAPH MOTIF

Plan

Complexité paramétrée et motifs avec topologie

Motifs sans topologie

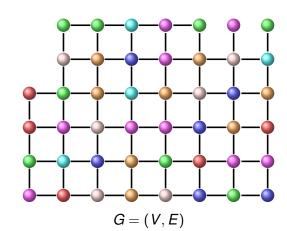
Le problème GRAPH MOTIF Des logiciels pour GRAPH MOTIF

- Constat : les données biologiques sont très bruitées
 - Manque des informations, faux négatifs. Estimé à 50% [GAVIN ET AL. 2002]
 - ► Informations erronées, faux positifs. Estimé à 65% [REGULY ET AL. 20061
- La topologie du motif peut ne pas être connue a priori
- La topologie du motif peut ne pas être pertinente

- Chaque nœud du réseau est coloré par sa "fonction"
- ► Le motif sera juste un ensemble (ou multi ensemble) de couleurs à rechercher dans un réseau coloré [Lacroix et al. 20061

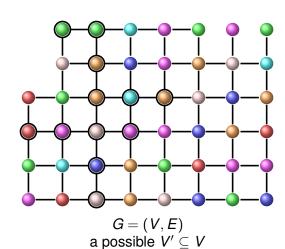
- S'applique à différent type de réseaux biologiques
 - Dans les réseaux PPI, chaque protéine du motif recoit une couleur
 - Le réseau est coloré selon les homologies avec les protéines du motif
- ► Selon [Betzler et al. 2008], peut être utilisé pour réseaux sociaux ou autres réseaux complexes...





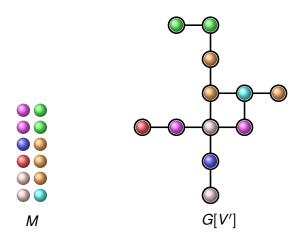
GRAPH MOTIF – Un exemple





Le problème GRAPH MOTIF

GRAPH MOTIF – Un exemple



GRAPH MOTIF - NP-C

- ▶ Le problème est NP-Complet, même si
 - ▶ Le réseau est un arbre [Lacroix et al. 2006]
 - Cet arbre est de degré max 3 et le motif est un ensemble simple [Fellows et al. 2008]
 - Le motif n'est constitué que de 2 couleurs et le réseau est un graphe biparti de degré max 4 [Fellows et al. 2008]

GRAPH MOTIF – Complexité paramétré

▶ Le problème est FPT par la taille du motif [Lacroix et al. 2006]

GRAPH MOTIF – Complexité paramétré

- ► Le problème est FPT par la taille du motif [Lacroix et al. 20061
- ► Mais est W[1]-difficile si paramétré par le nombre de couleurs [Fellows et al. 2008] (pas d'algo FPT possible)

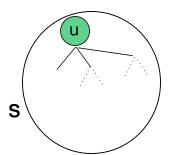
GRAPH MOTIF - FPT - Motif colorful

▶ Complexité en $\mathcal{O}^*(3^k)$ si le motif est un ensemble simple (colorful), avec k taille du motif [Betzler et al. 2008]

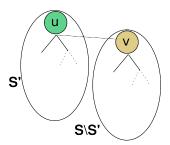
GRAPH MOTIF - FPT - Motif colorful

- ▶ Complexité en $\mathcal{O}^*(3^k)$ si le motif est un ensemble simple (colorful), avec k taille du motif [Betzler et al. 2008]
- ▶ Idée :
 - Solution connecté : on cherche un arbre
 - Motif colorful: on cherche donc un arbre colorful
 - 3. k couleurs différentes \Rightarrow k nœuds différents

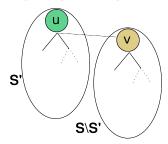
▶ But : trouver un arbre, enraciné en *u*, colorful sur *S*. D(u, S) = True?



- ▶ But : trouver un arbre, enraciné en *u*, colorful sur *S*. D(u, S) = True?
- ▶ Oui si *u* à un voisin *v* t.q.
 - 1. *u* enracine un arbre, colorful sur S' couleurs
 - **2.** v enracine un arbre, colorful sur $S \setminus S'$ couleurs



- ▶ But : trouver un arbre, enraciné en *u*, colorful sur *S*. D(u, S) = True?
- ▶ Oui si *u* à un voisin *v* t.q.
 - 1. u enracine un arbre, colorful sur S' couleurs
 - **2.** v enracine un arbre, colorful sur $S \setminus S'$ couleurs
- ▶ Regarder pour chaque voisin v et chaque $S' \subset S$
- \blacktriangleright $D(u, S) = D(u, S') \land D(v, S \setminus S')$

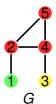


- ▶ But : trouver un arbre, enraciné en *u*, colorful sur *S*. D(u, S) = True?
- ▶ Oui si *u* à un voisin *v* t.q.
 - 1. u enracine un arbre, colorful sur S' couleurs
 - **2.** v enracine un arbre, colorful sur $S \setminus S'$ couleurs
- ▶ Regarder pour chaque voisin v et chaque $S' \subset S$
- \blacktriangleright $D(u,S) = D(u,S') \land D(v,S \setminus S')$
- Recursivement
- ightharpoonup ...jusqu'à S = col(u)

- Programmation dynamique valable que si le motif est un simple ensemble, car S et S' doivent être distincts (sinon plus assuré d'avoir k nœuds différents)
- On peut ramener le cas multi-ensemble au cas colorful avec le color-coding, en $\mathcal{O}^*(4.32^k)$ [Betzler et al. 2008]

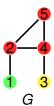
- ▶ Pour chaque couleur c du motif dont $occ_M(c) \ge 2$
- ▶ Créer $occ_M(c)$ nouvelles couleurs et remplacer c dans le motif par ces nouvelles couleurs





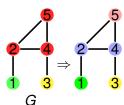
- ▶ Pour chaque couleur c du motif dont $occ_M(c) \ge 2$
- ▶ Créer $occ_M(c)$ nouvelles couleurs et remplacer c dans le motif par ces nouvelles couleurs





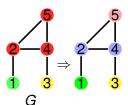
- ▶ Pour chaque couleur c du motif dont $occ_M(c) \ge 2$
- ▶ Créer $occ_M(c)$ nouvelles couleurs et remplacer c dans le motif par ces nouvelles couleurs
- Recolorer aléatoirement les nœuds du réseau portant la couleur c par une des nouvelles couleurs





- ▶ Pour chaque couleur c du motif dont $occ_M(c) \ge 2$
- ► Créer $occ_M(c)$ nouvelles couleurs et remplacer c dans le motif par ces nouvelles couleurs
- ► Recolorer aléatoirement les nœuds du réseau portant la couleur c par une des nouvelles couleurs

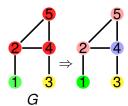




► Finalement, le motif est colorful : en chercher une occurence dans le réseau modifié avec l'algo précédent

- ▶ Pour chaque couleur c du motif dont $occ_M(c) \ge 2$
- ► Créer $occ_M(c)$ nouvelles couleurs et remplacer c dans le motif par ces nouvelles couleurs
- ► Recolorer aléatoirement les nœuds du réseau portant la couleur c par une des nouvelles couleurs





- ► Finalement, le motif est colorful : en chercher une occurence dans le réseau modifié avec l'algo précédent
- ▶ Il faut répeter $|\ln(\epsilon)|e^k$ fois cette coloration aléatoire pour avoir une probabilité $1 - \epsilon$ de succès

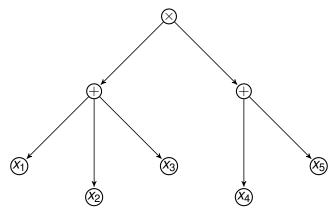
► On peut diminuer la complexité de GRAPH MOTIF en utilisant un framework dû à Koutis et Williams [ICALP 08.09]

- ► On peut diminuer la complexité de GRAPH MOTIF en utilisant un framework dû à Koutis et Williams [ICALP 08,09]
- Résultat clef : On peut déterminer par un algorithme randomisé, en temps $\mathcal{O}^*(2^k)$ et en espace polynomial, si un polynôme représenté par un circuit arithmétique contient un monome multilineaire de degré k

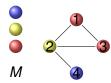
- Un monôme est multilineaire si chaque variable du monôme n'apparait qu'une seule fois (pas de carré)
- ▶ Par définition, le degré d'un monôme multilinéaire est son nombre de variables
- Exemple : $P(X) = (x_1^2 x_3 x_5 + x_1 x_2 x_4 x_6)$
 - x₁x₂x₄x₆ est un monôme multilineaire de degré 4
 - \rightarrow $x_1^2 x_3 x_5$ n'est pas un monôme multilineaire

GRAPH MOTIF et polynômes

- ▶ Un circuit arithmétique sur un ensemble de variables X est un DAG t.g.:
 - ▶ les noeuds internes sont les opérations × ou +
 - les feuilles sont des éléments de X
- Exemple pour $P(X) = (x_1 + x_2 + x_3)(x_4 + x_5)$



- Pour résoudre GRAPH MOTIF avec motif colorful :
 - On introduit des variables correspondant aux couleurs
 - On construit un circuit correspondant aux couleurs de tous les sous arbres de taille k dans G
 - ▶ Un monôme multilineaire de degré *k* correspond alors à un sous arbre avec *k* couleurs différentes (motif colorful)

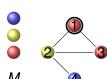


Construction récursive.

Les monômes de degré k=3 contenant le noeud 1 sont la somme des produits des monômes:

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

Le problème GRAPH MOTIF



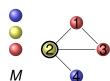
Construction récursive.

Les monômes de degré k=3 contenant le noeud 1 sont la somme des produits des monômes:

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

Pour
$$k' = 1$$
:

 X_R



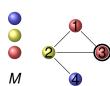
Construction récursive

Les monômes de degré k = 3 contenant le noeud 1 sont la somme des produits des monômes:

Le problème GRAPH MOTIF

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

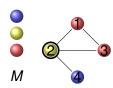
Pour
$$k' = 1$$
: $x_R .(P_{2,2})$



Construction récursive

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

Pour
$$k' = 1$$
:
 $x_R \cdot (P_{2,2} + P_{2,3})$



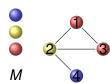
Construction récursive.

Les monômes de degré k=3 contenant le noeud 1 sont la somme des produits des monômes :

Le problème GRAPH MOTIF

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

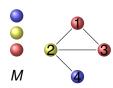
Pour
$$k' = 1$$
:
 $x_R \cdot (P_{2,2} + P_{2,3})$
 $= x_R \cdot (x_Y \cdot (P_{1,1} + P_{1,3} + P_{1,4}) + P_{2,3})$



Construction récursive.

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

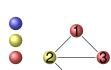
Pour
$$k' = 1$$
:
 $x_R \cdot (P_{2,2} + P_{2,3})$
 $= x_R \cdot (x_Y \cdot (P_{1,1} + P_{1,3} + P_{1,4}) + P_{2,3})$
 $= x_R \cdot (x_Y \cdot (x_R + x_R + x_R) + P_{2,3})$



Construction récursive.

- **1.** de degré k' < k contenant le noeud 1,
- **2.** de degré k k' contenant un voisin de 1.

Pour
$$k' = 1$$
:
 $x_R \cdot (P_{2,2} + P_{2,3})$
 $= x_R \cdot (x_Y \cdot (P_{1,1} + P_{1,3} + P_{1,4}) + P_{2,3})$
 $= x_R \cdot (x_Y \cdot (x_R + x_R + x_B) + P_{2,3})$
 $= x_R \cdot (x_Y \cdot x_R + x_Y \cdot x_R + x_Y \cdot x_R + P_{2,3})$



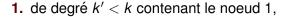
Μ

Construction récursive.

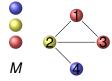
- **1.** de degré k' < k contenant le noeud 1.
- **2.** de degré k k' contenant un voisin de 1.

Pour
$$k' = 1$$
:
 x_R . $(P_{2,2}+P_{2,3})$
= x_R . $(x_Y$. $(P_{1,1}+P_{1,3}+P_{1,4})+P_{2,3})$
= x_R . $(x_Y$. $(x_R+x_R+x_B)+P_{2,3})$
= x_R . $(x_Y$. x_R+x_Y . x_R+x_Y . x_R+x_Y . $x_B+P_{2,3})$
= x_R . x_Y . x_R+x_Y . x

Construction récursive.







Pour
$$k' = 1$$
:
 x_R . $(P_{2,2}+P_{2,3})$
= x_R . $(x_Y$. $(P_{1,1}+P_{1,3}+P_{1,4})+P_{2,3})$
= x_R . $(x_Y$. $(x_R+x_R+x_B)+P_{2,3})$
= x_R . $(x_Y$. x_R+x_Y . x_R+x_Y . $x_B+P_{2,3})$
= $x_Rx_Yx_R+x_Rx_Yx_R+x_Rx_Yx_B+\dots$
Il y a un monôme multilineaire (donc une solution)

GRAPH MOTIF et polynômes

- ► Si le **motif est un multiensemble**, un monôme non multilinéaire sur les couleurs est une solution
- La construction précédente n'est plus suffisante
- On introduit des variables pour les noeuds du graphe

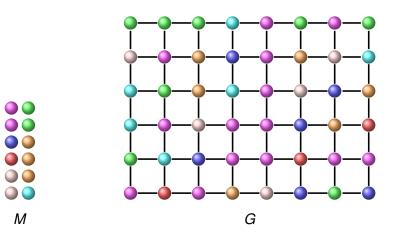
Une variante de GRAPH MOTIF: MAX MOTIF

- Comme données bruitées, chercher une occurence exacte peut-être impossible
- Autoriser insertions et deletions

Une variante de Graph Motif: Max Motif

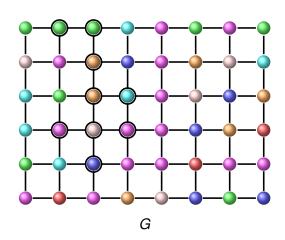
- ► Comme données bruitées, chercher une occurence exacte peut-être impossible
- Autoriser insertions et deletions
- ▶ Un exemple : MAX MOTIF
- Trouver une occurence qui matche "le plus possible" de couleurs du motif

Exemple MAX MOTIF



Exemple MAX MOTIF

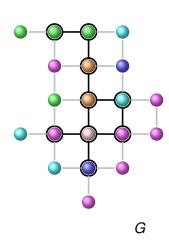




Le problème GRAPH MOTIF

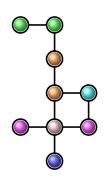
Exemple MAX MOTIF





Exemple MAX MOTIF





G

- (Malheuresement,)MAX MOTIF est aussi dur à approximer que INDEPENDANT SET, même si
 - Le réseau est un arbre
 - Le motif est colorful
 - Chaque couleur apparait au plus 2 fois dans le réseau
- ▶ INDEPENDANT SET : Dans un graphe G = (V, E), trouver le plus grand $V' \subseteq V$ t.q. aucun nœud dans V' est relié à un autre nœud de V'

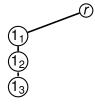
▶ A partir d'un graphe $G = (V_G, E_G)$, on construit un arbre



$$G = (V_G, E_G)$$

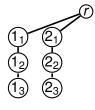


$$G = (V_G, E_G)$$



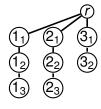


$$G = (V_G, E_G)$$



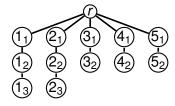


$$G = (V_G, E_G)$$





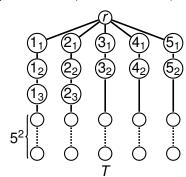
$$G = (V_G, E_G)$$



► Ajout de $|V_G|^2$ nœuds à chaque chemin ("boites noires")



$$G = (V_G, E_G)$$



Le problème GRAPH MOTIF

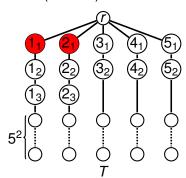
Coloration (2 occurences...) + motif (colorful)



$$G = (V_G, E_G)$$

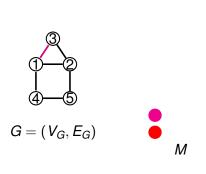


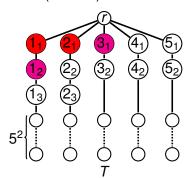
Μ



MAX MOTIF – Innaproximabilité

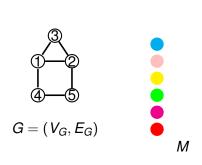
Coloration (2 occurences...) + motif (colorful)

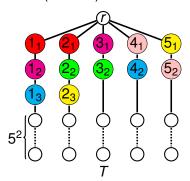




MAX MOTIF – Innaproximabilité

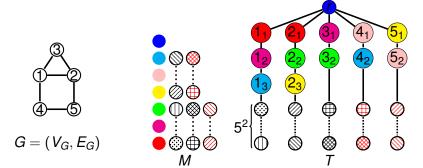
Coloration (2 occurences...) + motif (colorful)



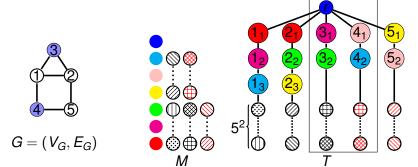


Des couleurs différentes pour la racine et les boites noires

Le problème GRAPH MOTIF

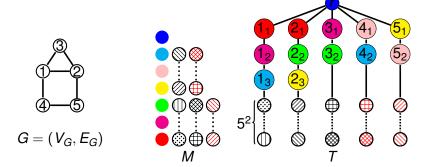


► INDEPENDANT SET ⇒ MAX MOTIF colorful



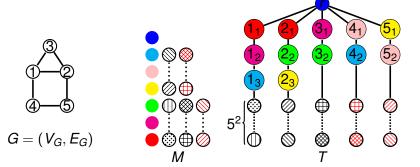
MAX MOTIF – Innaproximabilité

 On prouve que s'il existait un algo d'approx. pour MAX MOTIF, il y aurait un algo d'approx. pour INDEPENDANT SET



▶ Idées :

- La racine doit être prise pour rendre la solution connexe
- Les boites noires sont données "gratuitement" à la solution si le chemin depuis la racine est dans la solution.
- Une boite noire est de plus grande taille que l'ensemble des autres noeuds



- ► Torque [Bruckner et al. 2009] : un web service, qui ne gère que les motifs colorfuls
- ► GraMoFoNe [BLIN ET AL. 2010]: un plugin cytoscape

GraMoFoNe – Cytoscape (2002)

- Plateforme java open-source, gratuite, permettant
 - l'importation / l'exportation nombreux formats, BDD,...
 - la visualisation et l'analyse de réseaux d'interactions
 - l'integration d'annotations à ces réseaux
- Largement utilisé ("des centaines" d'articles l'utilisent pour de l'analyse)
- Maintenu...



GraMoFoNe – Cytoscape

- Principal avantage : plugins
- Ajout de fonctions :
 - Analyses
 - Nouveaux layouts
 - Support de fichiers
 - Connection avec des bases de données

GraMoFoNe – PB

- On exprime GRAPH MOTIF avec de la programmation linéaire Pseudo-Booléenne
- ▶ i.e. de la programmation linéaire avec des variables booléennes

GraMoFoNe – PB

- Trouver un assignement de variables satisfaisant les contraintes et maximisant l'objectif
- Un exemple simple :
 - ▶ Variables : $x_i \in \{0, 1\}, \forall i = 1, 2, 3$
 - ▶ **Objectif**: $\max x_1 + 2x_2 x_3$
 - Contraintes :
 - 1. $x_1 2x_2 + 3x_3 > 1$
 - 2. $X_1 + X_2 + X_3 = 1$
 - 3. $2x_1 + x_2 + x_3 < 3$
- ► Solution : $x_1 = 1, x_2 = 0, x_3 = 0$

GraMoFoNe - PB

- Donne une solution exacte
- Beaucoup de solvers efficaces existent
- Dont beaucoup de gratuits (\neq CPLEX)

- Donne une solution exacte
- Beaucoup de solvers efficaces existent
- Dont beaucoup de gratuits (\neq CPLEX)
- ► On utilise SAT4JPseudo [Le Berre, Parrain 2007]
 - Gratuit
 - Maintenu
 - Bonnes performances lors des compétitions
 - En java

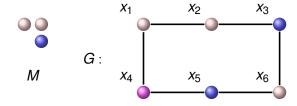


GraMoFoNe – PB

- On utilise 23 contraintes et 9 domaines de variables.
- ▶ On cherche une occurence d'un motif *M* dans un graphe *G*
- A respecter :
 - La taille de la solution.
 - 2. La coloration de la solution par rapport au motif
 - 3. La connexité de la solution (partie difficile)

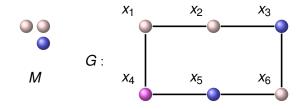
GraMoFoNe – Variables

▶ Une variable x pour chaque nœud



GraMoFoNe – Variables

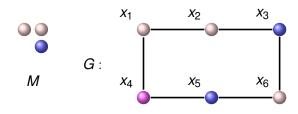
▶ Une variable x pour chaque nœud



- ► Contrainte "taille solution" : $\sum_{v \in V} x_v = |M|$:
 - $X_1 + X_2 + X_3 + X_4 + X_5 + X_6 = 3$

GraMoFoNe – Variables

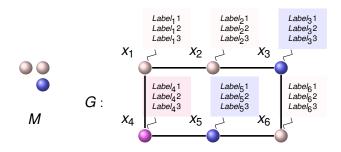
Une variable x pour chaque nœud



- ► Contrainte "taille solution" : $\sum_{v \in V} x_v = |M|$:
 - $X_1 + X_2 + X_3 + X_4 + X_5 + X_6 = 3$
- ► Contraintes "coloration" $\sum_{\substack{v \in V \\ c \in col(v)}} x_v = occ_M(c)$:
 - $x_1 + x_2 + x_6 = 2$ (gris)
 - $x_3 + x_5 = 1$ (bleu)
 - $x_4 = 0$ (rose)

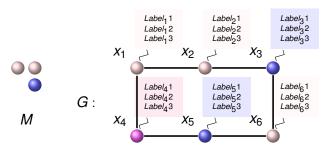
GraMoFoNe – Variables

► Connexité : |M| variables Label_v pour chaque nœud v



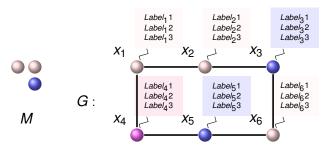
GraMoFoNe – Variables

▶ Connexité : |M| variables Label_v pour chaque nœud v



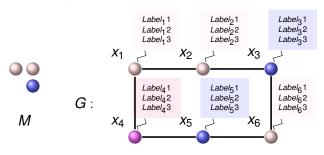
- ▶ |V| contraintes "un seul label par nœud dans la solution" :
 - ▶ Pour chaque $v, x_v \Rightarrow (\sum_{i=1}^{|M|} Label_v i = 1)$

▶ Connexité : |M| variables Label_v pour chaque nœud v



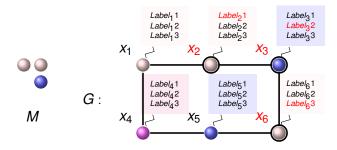
- ▶ |V| contraintes "un seul label par nœud dans la solution" :
 - ▶ Pour chaque $v, x_v \Rightarrow (\sum_{i=1}^{|M|} Label_v i = 1)$
- ► |M| contraintes "un seul nœud par label donné"
 - ▶ Pour un label *i* donné, $\sum_{v \in V} Label_v i = 1$

▶ Connexité : |M| variables Label_v pour chaque nœud v

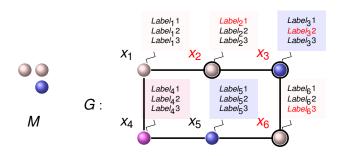


- ▶ |V| contraintes "un seul label par nœud dans la solution" :
 - ▶ Pour chaque $v, x_v \Rightarrow (\sum_{i=1}^{|M|} Label_v i = 1)$
- ► |M| contraintes "un seul nœud par label donné"
 - ▶ Pour un label *i* donné, $\sum_{v \in V} Label_v i = 1$
- ► |V|.|M| contraintes "un nœud avec un label a un voisin avec un label supérieur" (sauf le dernier)
 - ▶ Label_v $i \Rightarrow (\sum_{u \in N(v)} \sum_{i>i} Label_u j \ge 1)$

GraMoFoNe – Une solution

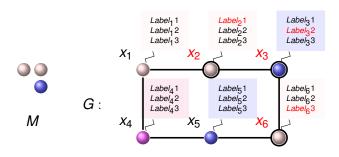


GraMoFoNe – Une solution



- "taille": $x_1 + x_2 + x_3 + x_4 + x_5 + x_6 = 3$
- "coloration":
 - $x_1 + x_2 + x_6 = 2$ (gris)
 - $x_3 + x_5 = 1$ (bleu)
 - $\rightarrow x_4 = 0$ (rose)

GraMoFoNe – Une solution



- "un label par nœud" : $\forall v, x_v \Rightarrow (\sum_{i=1}^{|M|} Label_v i = 1)$
- "un nœud par label" : $\sum_{v \in V} Label_v i = 1$
- "voisin avec label supérieur" : $Label_v i \Rightarrow (\sum_{u \in N(v)} \sum_{i>i} Label_u j \geq 1)$

GraMoFoNe - PB

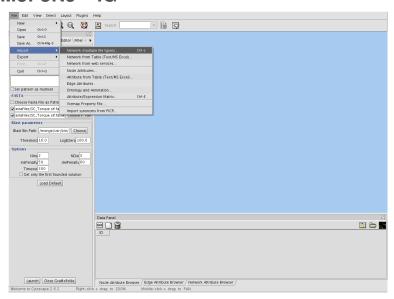
- Gère GRAPH MOTIF "classique" (colorful et multi-ensemble)
- PB permet d'avoir toutes les solutions possibles...

GraMoFoNe – PB

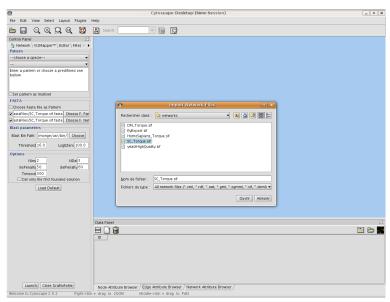
- Gère GRAPH MOTIF "classique" (colorful et multi-ensemble)
- PB permet d'avoir toutes les solutions possibles...
- Avec variables et contraintes supplémentaires, on peut gérer
 - Les insertions
 - Les deletions
 - Une liste de couleurs associée à chaque nœud du graphe

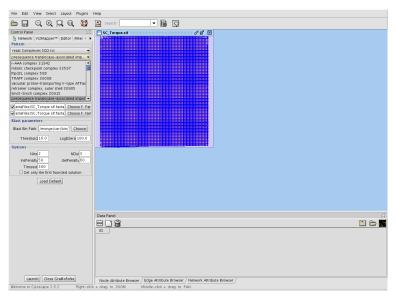
- Gère GRAPH MOTIF "classique" (colorful et multi-ensemble)
- PB permet d'avoir toutes les solutions possibles...
- Avec variables et contraintes supplémentaires, on peut gérer
 - Les insertions
 - Les deletions
 - Une liste de couleurs associée à chaque nœud du graphe
- Devient complexe (compensations "une deletion compense une insertion", garder la bijection entre motif et graphe si plusieurs couleurs dans le graphe,...)

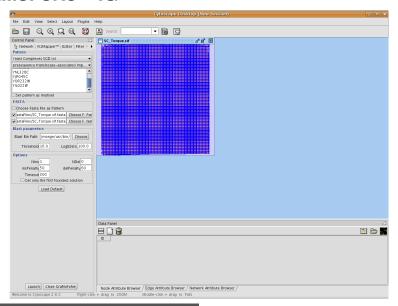
GraMoFoNe - IG

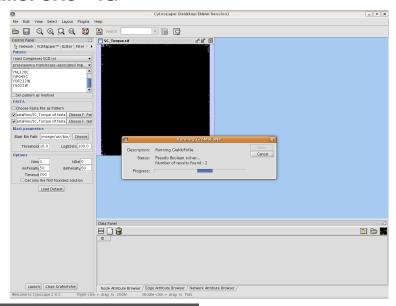


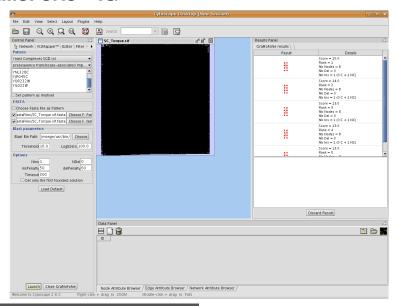
GraMoFoNe - IG

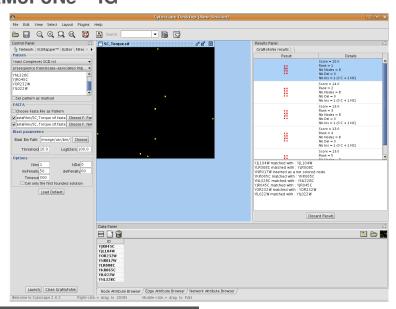


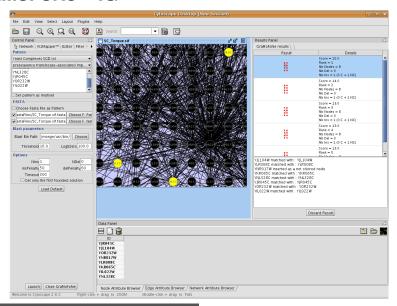


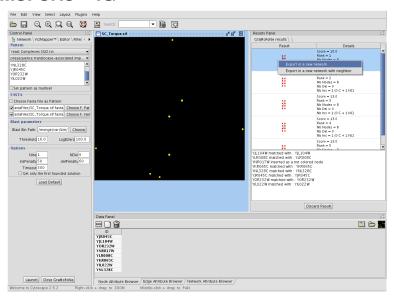


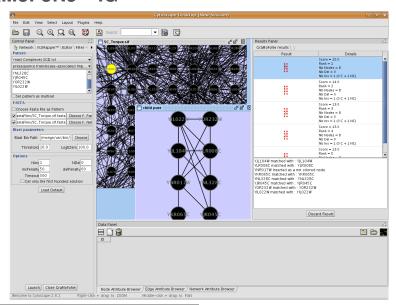








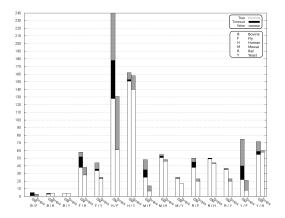




- Utilisé pour tests grande échelle
- Recherche de dizaines de complexes protéiques dans des réseaux d'autres espèces

- Utilisé pour tests grande échelle
- Recherche de dizaines de complexes protéigues dans des réseaux d'autres espèces
- Données :
 - Motifs de 6 espèces (levure, drosophile, homo sapiens, souris, boeuf, rat)
 - Réseaux de 3 espèces (levure, drosophile, homo sapiens)
 - De 2 à 4 indels autorisés selon la taille du motif

- Chaque motif est
 - Trouvé dans le temps imparti
 - 2. Marqué comme non trouvé dans le temps imparti
 - 3. Inconnu (temps imparti écoulé)



- ▶ + de bruit pour la mouche (faux negatifs déconnectent la solution, faux positifs donnent des "mauvaises" solutions)
- ▶ 5-20s (petits M), 40-60s (grands). Mais PB "boite noire"

Plan

Introduction

Complexité paramétrée et motifs avec topologie

Motifs sans topologie

Le problème GRAPH MOTIF Des logiciels pour GRAPH MOTIF

Conclusion

Conclusion

- Deux vues pour la recherche de motifs
 - 1. Motifs avec topologie
 - 2. Motifs sans topologie
- ► Problèmes difficiles donc :
 - Complexité paramétrée
 - Approximation

Conclusion

- ► Le bruit pousse à la recherche d'extensions plus souples
 - contrainte couleurs
 - contrainte connexité
 - Simple connexité : FPT
 - 2-connexité : W[1]-difficile
 - ► Modules ?
 - ► 7
 - √

Questions sur la recherche de motifs dans les graphes colorés ?

Florian Sikora (encadré par Guillaume Blin et Stéphane Vialette)

Université Paris-Est, LIGM - UMR CNRS 8049

Séminaire Combi LINA – 06/2010

► Si le **motif est un multiensemble**, un monôme non multilinéaire sur les couleurs est pourtant une solution

GRAPH MOTIF et polynômes

- Si le motif est un multiensemble, un monôme non multilinéaire sur les couleurs est pourtant une solution
- \triangleright On introduit des variables x_u pour les noeuds du graphe
- ▶ Pour chaque couleur c qui apparait m fois dans le motif, on introduit des variables y_{c,1}, y_{c,2},..., y_{c,m}
- On modifie la construction du circuit :

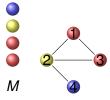
$$P_{u,1} = x_u \cdot (y_{c,1} + y_{c,2} + \cdots + y_{c,m})$$

- Les variables x assurent qu'un noeud n'est utilisé qu'une fois
- ► Les variables y donnent le bon nombre de couleurs demandés par le motif

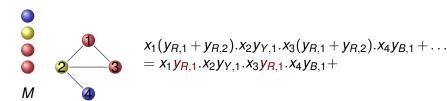
- Si le motif est un multiensemble, un monôme non multilinéaire sur les couleurs est pourtant une solution
- \triangleright On introduit des variables x_u pour les noeuds du graphe
- ▶ Pour chaque couleur c qui apparait m fois dans le motif, on introduit des variables y_{c,1}, y_{c,2},..., y_{c,m}
- On modifie la construction du circuit :

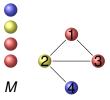
$$P_{u,1} = x_u \cdot (y_{c,1} + y_{c,2} + \cdots + y_{c,m})$$

- Les variables x assurent qu'un noeud n'est utilisé qu'une fois
- ► Les variables *y* donnent le bon nombre de couleurs demandés par le motif
- On cherche maintenant un monôme de degré 2k (noeuds + couleurs)



$$x_1(y_{R,1}+y_{R,2}).x_2y_{Y,1}.x_3(y_{R,1}+y_{R,2}).x_4y_{B,1}+\ldots$$

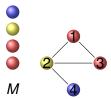




$$X_1(y_{R,1} + y_{R,2}).X_2y_{Y,1}.X_3(y_{R,1} + y_{R,2}).X_4y_{B,1} + \dots$$

$$= X_1y_{R,1}.X_2y_{Y,1}.X_3y_{R,1}.X_4y_{B,1} + \dots$$

$$X_1y_{R,1}.X_2y_{Y,1}.X_3y_{R,2}.X_4y_{B,1} + \dots$$



$$X_1(y_{R,1} + y_{R,2}).X_2y_{Y,1}.X_3(y_{R,1} + y_{R,2}).X_4y_{B,1} + \dots$$

$$= X_1y_{R,1}.X_2y_{Y,1}.X_3y_{R,1}.X_4y_{B,1} + \dots$$

$$X_1y_{R,1}.X_2y_{Y,1}.X_3y_{R,2}.X_4y_{B,1} + \dots$$

 Un monôme multilinéaire solution est de degré 8.