

Surveillance passive dans l'Internet

C. Chaudet¹, E. Fleury², I. Guérin-Lassous² et H. Rivano³, M.-E. Voge³

¹ ENST - 46, rue Barrault - 75634 Paris Cedex 13, France - Claude.Chaudet@enst.fr

² INRIA Ares - CITI, INSA de Lyon, 69621 Villeurbanne, France - Prénom.Nom@insa-lyon.fr

³ I3S/INRIA Mascotte - INRIA Sophia Antipolis, BP 93, 06902 Sophia-Antipolis - Prénom.Nom@sophia.inria.fr

Afin d'obtenir les informations nécessaires à une bonne gestion des ressources de leur réseau, les opérateurs placent des sondes passives sur les liens de leurs points de présence. Dans cet article, nous donnons des écritures en programmes linéaires mixtes des problèmes de placement de sondes simples ou avec échantillonnage, et donnons une stratégie pour la maintenance de la surveillance partielle de trafics dynamiques dans un point de présence. Ces formulations améliorent les résultats de deux articles récents de la littérature.

Keywords: surveillance passive, positionnement de sondes, métrologie, problème de flot, simulations

1 Introduction

Le problème de la surveillance passive dans l'Internet consiste à placer sur les liens d'un réseau des équipements spécifiques (appelés *sondes*) analysant le trafic transitant sur les liens. Cette surveillance permet d'évaluer dynamiquement la sécurité, la connectivité et les différentes utilisations faites du réseau, et donc de gérer efficacement l'infrastructure et les ressources du réseau. Le positionnement des sondes dans le réseau est un problème clé pour la surveillance dans l'Internet, et suscite beaucoup d'intérêt dans la communauté "Réseaux". Ces équipements sont très coûteux et présentent de nombreuses contraintes physiques comme l'espace mémoire, l'alimentation, etc. Des articles très récents se sont intéressés au positionnement optimal de sondes [1, 2]. Les approches et problèmes résolus dans ces deux articles sont très similaires. Dans [2], les auteurs introduisent des notions de coûts de déploiement et de maintenance pour les sondes, alors que dans [1] les sondes ont un coût unitaire. En revanche, dans [1], les auteurs s'intéressent à la résolution optimale du positionnement par le biais de la programmation linéaire en nombres entiers, tandis que dans [2], les auteurs utilisent des algorithmes d'approximation basés sur une approche gloutonne. Enfin, dans [2], les sondes sont capables d'échantillonner le trafic, *i.e.* de contrôler le nombre de paquets à analyser.

Dans cet article, nous présentons le problème ainsi que le cadre de l'étude en section 2, puis nous améliorons et étendons les résultats donnés dans [1] et [2]. Notamment, nous présentons dans la section 3 une amélioration de la formulation de [1] : nous modifions le programme linéaire initial en un autre programme linéaire répondant toujours à notre problème mais nécessitant un temps de résolution bien plus court. Enfin, dans la section 4, nous étendons cette modélisation pour prendre en compte la notion d'échantillonnage introduite dans [2], et obtenir un algorithme rapide de gestion de la surveillance du réseau lorsque le trafic est dynamique.

2 Présentation du problème

Un point de présence (POP) permet de diriger les connexions des utilisateurs, connectés aux *routeurs d'accès*, vers le cœur du réseau de l'opérateur. La surveillance du trafic passant à travers un POP est un enjeu majeur, notamment pour donner aux opérateurs les moyens de négocier au mieux les accords de service avec d'autres opérateurs. À ces fins, un opérateur installe des sondes sur les interfaces des routeurs du

réseau. Chacune de ces sondes surveille une partie du trafic, fixe ou adaptable (sonde avec *échantillonnage*). L'objectif d'un opérateur n'est pas forcément de surveiller la totalité du trafic, mais une certaine proportion k . En effet l'expérience montre que le surcoût nécessaire à la couverture des derniers pourcentages de trafic est très important, alors que le gain pour la maintenance du réseau ne l'est pas.

Routing Le routage des flux obéit à des règles propres à l'opérateur, pouvant aller d'un calcul des plus courts chemins type OSPF à des algorithmes d'équilibrage de charge. Nous supposons donc que les chemins peuvent être arbitraires, mais qu'il n'y en a qu'un "petit" nombre par flux. Par exemple, le protocole de routage OSPF garde moins d'une dizaine de chemins différents pour chaque flux.

Modélisation Les routeurs et liens de communication d'un POP seront représentés par un graphe $G = (V, E)$. Les flux sont transportés sur des chemins $p \in \mathcal{P}$. On agrège dans un même trafic tous les flux passant sur un ensemble $\mathcal{P}_{s,t}$ de chemins allant d'un client/POP source s à un client/POP destination t . Chaque chemin p transporte une quantité de trafic v_p . Le problème qui nous intéresse est donc de sélectionner certaines arêtes du graphe et, le cas échéant, de régler la fréquence d'échantillonnage sur chaque arête, afin de mesurer au moins $k \cdot \sum_{p \in \mathcal{P}} v_p$ unités de trafic ($k \in [0, 1]$).

Dans le cas, théorique, de sondes sans échantillonnage mesurant l'intégralité du trafic qu'elles voient passer, si l'on désire mesurer 100% du trafic, ce problème s'écrit comme un *Minimum Set Cover* dans le graphe d'incidence des chemins sur les arêtes[†].

Pour toute valeur de $k \in [0, 1]$, le problème se modélise assez directement par un *Minimum Edge Cost Flow* (MECF) dans un graphe auxiliaire. Le MECF est un problème de flot dans lequel le coût de certaines arêtes est binaire, nul quand aucun flot ne passe par l'arête et 1 quand une quantité quelconque passe l'arête. Ce problème est \mathcal{NP} -difficile et non approximable [3]. Étant donnée la structure du graphe auxiliaire, il est possible de combiner les équations de flot pour obtenir les formulations en programmes linéaires mixtes compactes présentées dans les sections suivantes.

3 Amélioration de la modélisation initiale

Programme linéaire 1 (PPM(k)).

$$\begin{aligned}
 & \text{Minimiser} && \sum_{e \in E} x(e) && \text{le nombre de sondes posées} \\
 & \text{t.q.} && \sum_{e \in p} x(e) \geq \delta(p) && \forall p \in \mathcal{P} \\
 & && \sum_{p \in \mathcal{P}} \delta(p) \cdot v_p \geq k \cdot \sum_{p \in \mathcal{P}} v_p \\
 & && \delta(p) \in [0, 1] && \forall p \in \mathcal{P} \\
 & && x(e) \in \{0, 1\} && \forall e \in E
 \end{aligned}$$

Dans le programme linéaire PPM(k), $x(e)$ indique si une sonde est placée sur l'arc e (si oui $x(e)$ est alors égal à 1) et $\delta(p)$ indique si le trafic passant par p est surveillé ou non. Dans la modélisation initiale de [1], $\delta(p)$ est entier et vaut 0 ou 1. La première contrainte détermine les trafics qui seront surveillés par les sondes, alors que la deuxième permet de s'assurer que suffisamment de trafic est surveillé. Il est possible de considérer $\delta(p)$ comme une variable décimale dans $[0, 1]$, ce qui ne change rien à la solution et accélère

[†] Le graphe d'incidence des chemins sur les arêtes est le graphe biparti avec un sommet par chemin dans une partie, un sommet par arête dans l'autre, et un arc entre deux sommets si le chemin correspondant au premier sommet passe sur l'arête correspondante au deuxième. Il est à noter que cette réduction prouve la \mathcal{NP} -difficulté et la non-approximabilité du problème : c'est le cas de *Minimum Set Cover* dans les graphes bipartis, et que tout graphe biparti est le graphe d'incidence de chemins sur les arcs d'une grille.

grandement sa résolution. En effet, dès qu'on a une solution avec $0 < \delta(p) < 1$, on peut arrondir $\delta(p)$ à 1, les $x(e)$ restent inchangés, et la valeur de $\sum_{p \in \mathcal{P}} \delta(p) v_p$ augmente, comme recherché.

Nous avons comparé les temps de calcul obtenus avec la modélisation initiale et celle proposée ici. Les POPs et les trafics considérés sont les mêmes que ceux présentés dans [1], en revanche les programmes linéaires ont été résolus par la bibliothèque CPLEX[4]. La table 1 compare les différents temps de calcul pour des POPs de 10, 15 et 29 nœuds et un taux de surveillance de 80%.

$ V - \mathcal{P} $	version modifiée	version initiale
10 – 132	0,02 s	0,05 s
15 – 1980	1,9 s	30,6 s
29 – 11130	34,2 s	873 s

TAB. 1: Temps de calcul moyen pour 40 exécutions sur des POPs de 10, 15 et 29 nœuds

4 Surveillance avec échantillonnage

Certaines sondes sont en mesure de ne capturer qu'une partie du trafic, grâce à un échantillonnage des paquets. Sur des liens très haut débit atteignant des débits de plusieurs Gb/s (OC-48, OC-192, OC-255), ne pas avoir à stocker et analyser tous les paquets, mais seulement une certaine proportion, apporte un gain significatif dans les coûts d'exploitation des sondes. Cette proportion de trafic échantillonné va dépendre du coût d'exploitation des sondes placées, et donc du coût d'échantillonnage d'un paquet sur une sonde, ce coût pouvant varier d'une sonde à l'autre en fonction de la vitesse du lien surveillé. On souhaite donc minimiser le coût d'exploitation des sondes tout en garantissant toujours une proportion globale k de trafic surveillé. Si l'opérateur du POP souhaite avoir une vision de tous les trafics, sans pour autant devoir surveiller tous les chemins, on garantit une couverture minimum h sur chaque trafic. On notera que $h \leq k$ étant donné qu'il s'agit d'un minimum de couverture de chaque trafic alors que k est un minimum sur la totalité des trafics.

4.1 Écriture en programmation linéaire mixte

Nous allons modéliser le coût d'installation d'une sonde sur un lien e par $cost_i(e)$ et son coût d'exploitation par $cost_e(e)$. Ces fonctions de coût peuvent être quelconques sans que cela n'ait d'impact sur la formulation du problème linéaire 2, présenté ci-dessous. Néanmoins, le coût d'exploitation est souvent une fonction croissante concave [2], ce qui permet de prendre en compte le facteur d'échelle.

Programme linéaire 2 (PPME(h,k)).

$$\begin{aligned}
 \text{Minimiser} \quad & \sum_{e \in E} (cost_i(e) \cdot x(e) + cost_e(e) \cdot f(e)) && \text{le coût d'installation et d'exploitation} \\
 \text{t.q.} \quad & \sum_{e \in p} f(e) \geq \delta(p) \quad \forall p \\
 & x(e) \geq f(e) \quad \forall e \in E \\
 & \sum_{p \in \mathcal{P}_{s,t}} \delta(p) \cdot v_p \geq h \cdot \sum_{p \in \mathcal{P}_{s,t}} v_p \quad \text{pour tout trafic } (s,t) \\
 & \sum_{p \in \mathcal{P}} \delta(p) \cdot v_p \geq k \cdot \sum_{p \in \mathcal{P}} v_p \\
 & \delta(p), f(e) \in [0, 1] \quad \forall p \in \mathcal{P}, \forall e \in E \\
 & x(e) \in \{0, 1\} \quad \forall e \in E
 \end{aligned}$$

Comme pour le programme linéaire 1, $x(e)$ indique si une sonde est placée sur l'arc e . Par contre $\delta(p)$ indique quelle proportion du trafic passant par p est surveillée. On introduit ici les variables $f(e)$ qui

modélisent le taux d'échantillonnage de la sonde placée sur e , et les contraintes qui indiquent tout naturellement que pour échantillonner du trafic sur un arc e , il est nécessaire d'y avoir installé une sonde au préalable. Les contraintes suivantes indiquent qu'une proportion minimale de h de chaque trafic doit être surveillé et que l'on doit aussi avoir une proportion de couverture de trafic globale d'au moins k .

4.2 Trafic dynamique et adaptation de l'échantillonnage

Le trafic transporté par un POP est intrinsèquement dynamique et fluctue selon les périodes d'activité de la journée. Un changement des débits des flux et /ou des routes peut dégrader fortement la pertinence d'une surveillance avec échantillonnage. S'il n'est pas envisageable pour un opérateur de modifier la position des sondes à chaque changement de trafic, il reste tout à fait possible d'adapter les taux d'échantillonnage des sondes déjà installées. Il s'agit alors de trouver une solution au programme linéaire $PPME(h, k)$ en fixant a priori les $x(e)$, puisque les sondes sont déjà mises en œuvre. On note $PPME^*(x, h, k)$ ce problème.

$PPME^*(x, h, k)$ s'écrit comme le programme linéaire 2 en considérant les $x(e)$ comme des constantes. Il n'y a donc plus de variables binaires, et il est alors possible de trouver une solution optimale en temps polynomial. En fait, il s'agit même d'un calcul de flot de coût minimum qui peut s'effectuer rapidement, sans recours à la programmation linéaire.

Si un opérateur se donne un niveau de surveillance minimal par trafic h , un niveau global k et un seuil de tolérance $T < k$, une stratégie de maintien de la surveillance dans un POP pourrait être :

1. Tant que $\sum_{p \in \mathcal{P}} \delta(p) \cdot v_p \geq T \cdot \sum_{p \in \mathcal{P}} v_p$, attendre ;
2. Dès que $\sum_{p \in \mathcal{P}} \delta(p) \cdot v_p < T \cdot \sum_{p \in \mathcal{P}} v_p$, calculer $PPME^*(x, h, k)$, mettre à jour les taux d'échantillonnage ;
3. Goto 1.

La résolution de $PPME$ devient donc une phase initiale de la conception du POP pour laquelle la complexité n'est pas cruciale. Par contre, lors de l'adaptation de l'échantillonnage au trafic, le temps de calcul est très important car il s'agit de réaction à des situations potentiellement fortement dynamiques. Le calcul de $PPME^*$ est rapide et, s'agissant d'un calcul de flot, nécessite peu de ressources.

5 Conclusion

Cet article donne une modélisation des divers problèmes de placement de sondes dans un réseau afin de mettre en œuvre une surveillance passive. Il donne des solutions efficaces et améliore les résultats précédents grâce à une modélisation par flot qui permet aussi de déduire les différents résultats de \mathcal{NCP} -complétude et de non-approximabilité.

Les extensions possibles de cette thématique concernent la prise en compte totale des multi-chemins, problème difficile. Le but est d'arriver à une modélisation qui s'affranchit du facteur multiplicatif entre le nombre de chemins employés et le nombre de flux considérés.

Pour finir, on peut s'interroger sur l'impact que peuvent avoir la topologie des POPs, qui est le plus souvent multi-étages, et les protocoles de routage de plus court chemin (OSPF) employés, sur le placement des sondes au sein d'un POP.

Références

- [1] C. Chaudet, E. Fleury, and I. Guérin Lassous. Positionnement optimal de sondes pour la surveillance active et passive de réseaux. In *CFIP 2005*, Bordeaux, France, March 2005.
- [2] K. Suh, Y. Guo, J. Kurose, and D. Towsley. Locating network monitors : complexity, heuristics, and coverage. In *Infocom 2005*, Miami, FL, USA, March 2005.
- [3] G. Even, G. Kortsarz, and W. Slany. On network design problems : fixed cost flows and the Covering Steiner Problem. *Transactions on Algorithms*. à paraître.
- [4] ILOG CPLEX. <http://www.ilog.com/products/cplex/index.cfm>. (v7.5).