

# Prise en compte d'aspects sémantiques dans la contruction d'un réseau pair-à-pair

Yann Busnel

► **To cite this version:**

Yann Busnel. Prise en compte d'aspects sémantiques dans la contruction d'un réseau pair-à-pair. MajecSTIC 2005: Manifestation des Jeunes Chercheurs francophones dans les domaines des STIC, IRISA – IETR – LTSI, Nov 2005, Rennes, pp.435-439. inria-00000836

**HAL Id: inria-00000836**

**<https://hal.inria.fr/inria-00000836>**

Submitted on 23 Nov 2005

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Prise en compte d'aspects sémantiques dans la construction d'un réseau pair-à-pair.

Yann Busnel

IRISA - Campus universitaire de Beaulieu - 35042 Rennes Cedex

ENS Cachan - Antenne de Bretagne - Campus de Ker Lann - Avenue Robert Schumann - 35170 Bruz

Yann.Busnel@irisa.fr

**Résumé :** Les systèmes pair-à-pair permettent de développer des systèmes distribués de grande taille. Ces systèmes construisent un réseau logique recouvrant le réseau physique existant. Par opposition au modèle classique dit de *client-serveur*, chaque nœud joue à la fois le rôle de client et de serveur et est connecté logiquement à un sous-ensemble du système. La plupart des systèmes existants optimisent le réseau recouvrant en minimisant la distance géographique (ou physique) entre chaque nœud. L'exploitation des intérêts communs entre pairs représente un axe de recherche récent dans l'optimisation de ces systèmes. Ceci consiste à ajouter des liens reflétant cette similarité d'intérêt en sus du réseau recouvrant. L'objectif de nos recherches dans ce domaine est de prendre en compte cette localité d'intérêt dans la construction même du réseau recouvrant. Cet article présente d'abord le paradigme de la communication épidémique, permettant de construire des réseaux pair-à-pair. Ensuite nous exposons des techniques permettant de prendre en compte une localité d'intérêt.

**Mots-clés :** Systèmes distribués, réseaux pair-à-pair, protocoles épidémiques, profil sémantique, proximité d'intérêt.

## 1 INTRODUCTION

De nombreuses approches permettant de construire un réseau reposant sur le paradigme pair-à-pair ont récemment émergés, formant des systèmes distribués de grande taille permettant de passer à l'échelle. Contrairement aux approches *client-serveur*, les nœuds sont connectés entre eux au-dessus du réseau physique et jouent à la fois les rôles de client et serveur. Ces systèmes permettent d'agréger des ressources et des services entre les différents nœuds du système. Ces derniers ne possèdent qu'une connaissance locale du système permettant ainsi de développer des applications à large échelle.

La majeure partie de ces réseaux est optimisée en fonction de la localité géographique des nœuds. Cela dit, derrière chaque machine se trouve un utilisateur. Les préférences des utilisateurs peuvent aussi être exploitées pour optimiser l'efficacité de l'application mise en œuvre sur un réseau pair-à-pair. Certaines études ont montrés la pertinence de la prise en compte de ces aspects [Handurukande, 2004b]. Dans le contexte de cette

étude, nous analysons comment prendre en compte ceux-ci via une proximité que nous qualifierons de sémantique. La recherche de nouvelles méthodes dans les systèmes pair-à-pair est actuellement un axe d'étude très actif.

Nous présentons dans la section 2, les différentes classes de réseaux pair-à-pair. Nous nous concentrons sur l'une d'elles ainsi que sur une classe de protocoles permettant de construire de tels réseaux pour ces systèmes : les protocoles épidémiques. Nous introduisons ensuite le principe de proximité sémantique dans la section 3, ainsi que des résultats préliminaires sur les différents impacts de la prise en compte de cette dernière. Enfin, nous concluons cet article dans la section 4.

## 2 LES RÉSEAUX PAIR-À-PAIR NON-STRUCTURÉS

Les systèmes pair-à-pair se sont popularisés ces dernières années avec les systèmes de partage de fichiers sur Internet. Par exemple, les systèmes populaires comme Gnutella [Gnutella, WWW] et KaZaA [KaZaA, WWW] supportent des millions d'utilisateurs partageant des pétaoctets de données sur Internet. Ces systèmes sont dits non-structurés car les applications n'imposent pas de structure au réseau sous-jacent. Ils sont simples et mettent en œuvre des techniques de recherche reposant sur l'inondation du réseau. Cela conduit cependant à des problèmes de performance en terme d'encombrement du réseau. Il existe également des systèmes pair-à-pair, dits hiérarchiques, tels eMule [eMule, WWW], eDonkey [eDonkey, WWW] ou Napster [Napster, WWW] reposant sur l'existence de *super-pairs* – sorte de nœuds contrôleurs. Certaines études se sont penchées sur des systèmes hybrides tel que [Castro, 2003]. Dans la suite, nous ne considérerons que les réseaux pair-à-pair non-structurés, étant donné leur simplicité de mise en œuvre et leur application potentielle à la localité d'intérêt.

### 2.1 Présentation générale

De manière générale, les réseaux pair-à-pair possèdent des capacités d'auto-organisation et de tolérance aux défaillances. En effet, à chaque arrivée et départ de nœud le système est capable de réorganiser les liens entre les nœuds restant, en fonction des contraintes imposées à celui-ci. De plus, la symétrie entre les nœuds, l'équilibrage de charge et la connaissance uniquement lo-

cale du système fournissent aux réseaux pair-à-pair une capacité inhérente de passage à l'échelle.

Les applications pour les systèmes pair-à-pair ne se limitent pas au partage de fichiers. Par exemple, de nombreuses applications ont été développées telles que l'archivage de données, la gestion de cache, la diffusion, etc. Plusieurs domaines de recherche sont en cours d'exploration tels que la distribution de données ou le monitoring. Ces domaines donnent lieu à de nombreux travaux portant sur l'organisation des nœuds, la sécurité, les algorithmes de routage et de recherche, etc. Une optimisation courante consiste à essayer de connecter les pairs voisins géographiquement entre eux.

Les réseaux pair-à-pair non-structurés reposent sur une construction aléatoire du graphe de connexions. Un nœud se joint au réseau par l'intermédiaire d'un autre nœud déjà connecté. Chaque nœud tient à jour une liste de voisins. Cette liste, que l'on appellera *vue* par la suite, représente les liens logiques entre les différents nœuds du système, formant ainsi le réseau recouvrant. Une fois inséré, le nœud sonde de manière périodique son voisinage afin de maintenir et découvrir un certain nombre de connexions, de ressources ou de services. La recherche de ressources (physique ou logique) dans un tel réseau peut se faire selon une technique d'inondation : un nœud désirant localiser une ressource  $r$  demande à ses voisins, inclus dans sa vue, s'ils la connaissent. À leur tour, ces voisins demandent à leurs voisins s'ils ont connaissance de cette ressource  $r$  et ce, jusqu'à une profondeur fixée par le système. Un nœud possédant la ressource  $r$  avertit l'auteur de la requête en lui envoyant une réponse qui parcourt le chemin initial en sens inverse. Nous présenterons uniquement une classe de protocoles permettant de construire, de réorganiser ou de maintenir le réseau recouvrant selon une heuristique particulière.

## 2.2 Les protocoles épidémiques

Cette classe regroupe les protocoles dits épidémiques. Chaque nœud met périodiquement à jour sa liste de voisins en échangeant de l'information avec l'un d'entre eux choisi selon une méthode particulière. Ces protocoles épidémiques peuvent être modélisés selon le modèle présenté dans [Jelasity, 2004a]. Celui-ci repose sur trois appels de fonctions, spécifiés et analysés dans [Jelasity, 2004a], permettant de paramétrer le protocole en fonction du graphe requis. Leur efficacité est évaluée dans ce même article. Ces trois fonctions se déclinent de la manière suivante :

**Sélection du voisin** Elle détermine l'heuristique du choix du nœud participant à l'échange (aléatoirement, tête de la liste de voisins, ...)

**Propagation de sa vue** L'information peut se propager de trois façons différentes : l'envoi de sa vue au nœud sélectionné (*push*), la demande de la vue du nœud sélectionné (*pull*) ou l'échange des vues entre ces deux nœuds (*push/pull*).

**Sélection de sa vue** La liste est ordonnée suivant des critères donnés (proximité géographique, ordre d'arrivée). Le choix des nœuds qui composent la liste

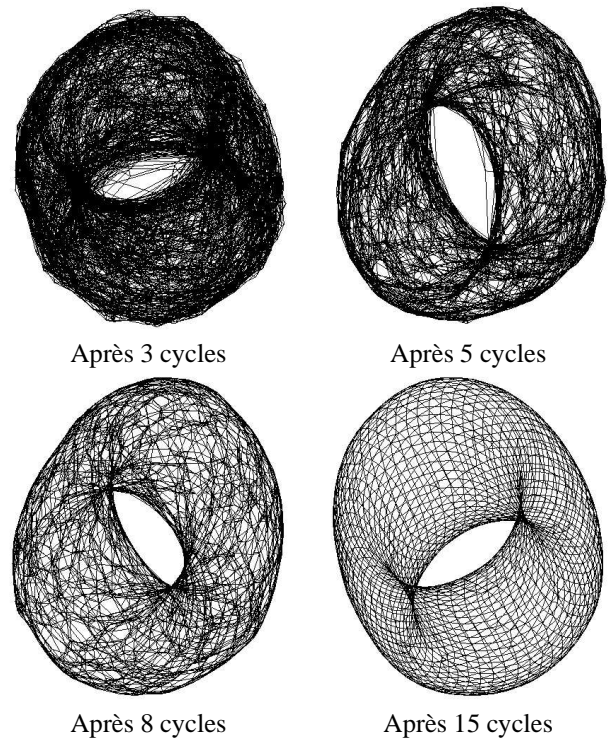


FIG. 1 – T-MAN : EXEMPLE ILLUSTRANT LA CONSTRUCTION D'UN TORE SUR  $50 * 50 = 2500$  NŒUDS RELIÉS À 4 VOISINS

des voisins peut se faire soit aléatoirement (sélection de  $c$  éléments de la vue), soit en prenant les  $c$  premiers nœuds apparaissant en tête de sa liste de voisins, soit en prenant les  $c$  derniers nœuds apparaissant en queue de celle-ci.

Afin d'illustrer l'efficacité de ce type de protocoles, nous présentons dans la figure 1 les résultats du protocole de réorganisation épidémique appelé T-Man [Jelasity, 2004b]. Dans cet article, les auteurs optimisent selon la localité géographique des nœuds du système. La vue de chaque nœud est ordonnée selon un critère de proximité géographique. Le nœud participant à l'échange est choisi aléatoirement dans la première moitié de la liste de voisins. La propagation est de type *push/pull* et la vue mise à jour correspond aux  $c$  nœuds les plus proches géographiquement du nœud en question. L'expérience modélisée sur la figure 1 illustre la rapidité de convergence de cette classe de protocole. En effet, 2500 nœuds ont été placés sur un tore de dimension (50, 50) et chaque nœud s'est vu attribuer une vue aléatoire. En moins de 15 cycles, plus de 99% des liens sont optimaux. Nous nous intéressons à ces protocoles étant donné la possibilité d'y appliquer d'autres critères, notamment sémantiques.

## 3 LA SÉMANTIQUE DANS LES RÉSEAUX PAIR-À-PAIR

De nombreux travaux se penchent aujourd'hui sur des techniques similaires pour capturer non plus une localité

géographique mais une localité d'intérêt (ou sémantique) dans les réseaux pair-à-pair. Ces systèmes exploitent des similitudes d'intérêts entre les utilisateurs d'une application donnée. De manière générale, une couche additionnelle est développée au dessus du réseau recouvrant prenant en compte ce type de localité. Des évaluations sont conduites sur ces systèmes [Handurukande, 2004a, Voulgaris, 2004]. Les résultats de cette étude offrent une comparaison entre l'efficacité de la recherche avec une liste de voisins aléatoires et l'efficacité au sein d'un groupement de nœuds (*Clustering*) basée sur des intérêts communs.

### 3.1 Description

Les aspects sémantiques peuvent être exposés ou non au sein de l'application. Afin d'éviter un investissement important de la part de l'utilisateur, nous nous intéressons aux systèmes utilisant une sémantique dite *implicite* ou non-intrusive. Ici, les applications cherchent à détecter les intérêts et les préférences de chacun des utilisateurs. Souvent, ces systèmes se basent sur l'utilisation d'un historique récent afin de capturer la proximité d'intérêt des nœuds du réseau pair-à-pair. En effet, différentes techniques implicites de groupement de nœuds axées sur la similarité des intérêts de chacun des participants ont été proposées et évaluées [Handurukande, 2004a]. Ainsi, une recherche exploitant cette sémantique améliore les chances de succès d'une requête dans les systèmes de partage de fichiers. En effet, pour une recherche sur 20 nœuds, alors que la stratégie aléatoire possède un taux de succès de 1%, les méthodes prenant en compte la sémantique via l'analyse d'un historique récent obtiennent des taux autour de 45%.

Un choix de conception de ces mises en œuvre provient de leur localisation. Se situant dans une couche supérieure au réseau recouvrant, elles n'interviennent pas, ou très peu, dans la construction de celui-ci. Nous proposons une alternative : remplacer les liens existants par des liens sémantiques. Afin de prendre en compte les aspects sémantiques directement au sein du réseau recouvrant, nous nous sommes appuyés sur les excellents résultats des protocoles épidémiques.

Afin de pouvoir utiliser ceux-ci, il est nécessaire de définir une relation d'ordre. En effet, chaque nœud doit pouvoir ordonner sa vue en fonction de l'heuristique choisie. Dans T-Man [Jelasity, 2004a], la relation d'ordre est induite par une mesure de distance géographique entre deux nœuds du système. Afin de pouvoir exploiter cette classe de protocoles, nous devons donc déterminer une mesure de distance permettant d'évaluer une proximité d'intérêt. Ainsi, nous pourrions comparer le profil sémantique de différents utilisateurs de l'application, et classer les nœuds d'une vue selon un critère sémantique (et non plus géographique).

Considérons un utilisateur d'un système de partage de fichiers. Celui-ci possède un ensemble de fichiers auxquels tous les utilisateurs de l'application ont accès. Nous appellerons cet ensemble la *cache* d'un nœud. Par la suite, nous prenons pour hypothèse qu'il n'y a qu'un utilisateur

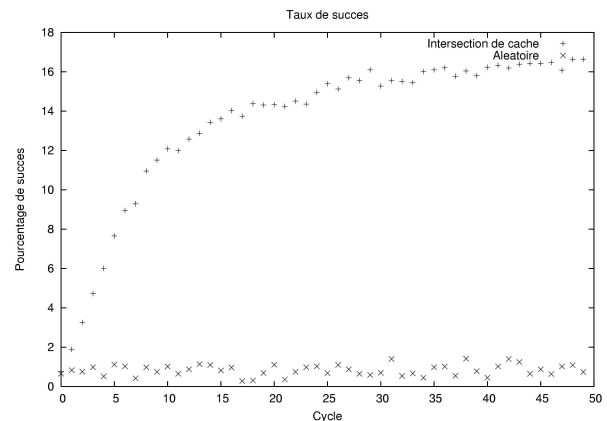


FIG. 2 – TAUX DE SUCCÈS MOYEN D'UNE REQUÊTE EN FONCTION DU CYCLE DU PROTOCOLE ÉPIDÉMIQUE

par nœud. Afin de déterminer la proximité sémantique entre deux nœuds, il est nécessaire de connaître les centres d'intérêts des utilisateurs. Aussi, étant donné l'orientation non-intrusive de l'approche dans la prise en compte de la sémantique, ces travaux sont basés sur l'étude du contenu des caches. Deux nœuds ayant beaucoup de fichiers en commun dans leur cache sont plus proches sémantiquement que deux nœuds en ayant très peu. Ainsi, la mesure de proximité sémantique entre deux nœuds du système est déterminée en fonction du cardinal de l'intersection des deux caches. La figure 2 illustre l'efficacité de cette mesure par rapport à une sélection de voisins aléatoires. En effet, dès les premiers cycles du protocole, le taux de succès d'une requête aux voisins d'un nœud est bien supérieur comparativement au taux dans un système ayant une structure aléatoire. Nous obtenons ces résultats par simulation sur une trace du logiciel eDonkey [eDonkey, WWW] datant de novembre 2003 [Handurukande, 2004a].

### 3.2 Effets de bords et problèmes ouverts

Malgré les bons résultats en terme de taux de succès d'une requête, cette méthode induit divers effets de bords qui peuvent biaiser les résultats et donc, réduire son efficacité. Nous présentons ci-dessous les différents points potentiellement à prendre en compte à terme dans la mesure de proximité.

**Générosité d'un nœud** La générosité d'un nœud biaise les résultats obtenus par la mesure introduite dans la partie 3.1. En effet, un nœud très généreux a une probabilité plus importante de posséder dans son cache un grand nombre des fichiers partagés par un autre nœud. Un nœud généreux possède donc une forte probabilité d'être choisi comme voisin sémantique – à tort. Ce biais se révèle en effet un biais positif mais induit potentiellement une forte charge sur le nœud généreux. En effet, la charge de ce nœud a une forte probabilité d'être nettement supérieure à la moyenne. La nécessité de limiter le choix d'un nœud généreux par rapport à un autre moins généreux,

mais plus proche sémantiquement, doit être intégrée dans notre mesure de proximité.

**Popularité des fichiers** De même, la popularité des fichiers entraîne un biais sur notre proximité. Un nœud possédant dans son cache une majorité de fichiers populaires risque d'être plus souvent choisi comme voisin sémantique. Cependant, les fichiers rares définissent plus précisément le profil sémantique d'un utilisateur que les fichiers populaires. La mesure de proximité doit donc limiter les liens vers des nœuds possédant beaucoup de fichiers populaires, afin de conserver une valuation de la distance sémantique correcte. Cependant, la popularité d'un fichier est une information globale et le problème consistant à obtenir une information globale de manière locale est encore ouvert.

**Non-connexité du réseau recouvrant** Une précédente étude a montré que la prise en compte des intérêts des utilisateurs dans le partage de fichier induit un regroupement des nœuds (*clustering*) [Handurukande, 2004a]. En effet, les nœuds proches sémantiquement vont se choisir mutuellement en tant que voisins. Ainsi, le risque de perdre la connexité du réseau recouvrant devient important. En effet, certains groupes peuvent ne plus posséder aucun lien avec le reste du réseau. D'après le fonctionnement des protocoles épidémiques, ces groupes ne peuvent alors plus recréer de liens vers le système et restent donc exclus. Afin d'éviter de partitionner le réseau recouvrant, il est donc indispensable que chaque nœud conserve quelques liens vers d'autres nœuds choisis aléatoirement.

#### **Pertinence de la mesure de proximité dans les requêtes**

L'observation des traces disponibles met en évidence une différence entre le profil sémantique reflété par l'étude du cache d'un nœud et celui reflété par l'étude des requêtes vers les autres nœuds. De manière générale, ce comportement est assez difficile à prendre en compte directement dans la mesure de proximité et pourra faire l'objet de prochaines études.

## **4 CONCLUSION**

Dans cet article, nous avons commencé par introduire les réseaux pair-à-pair non-structurés. Cette classe de systèmes distribués permet de développer des applications à très grande échelle dans lesquelles chaque nœud ne possède qu'une vision locale du réseau recouvrant. Après avoir introduit les protocoles épidémiques, nous avons présenté comment adapter leur efficacité à notre problème. En effet, nous cherchons à intégrer des aspects sémantiques directement dans la couche du réseau recouvrant. Ces protocoles permettant de construire et de maintenir ce type de réseau selon une proximité donnée, nous avons présenté une mesure de distance sémantique simple mais efficace.

Nos travaux de recherche actuels portent sur l'amélioration et l'optimisation de cette mesure afin

de prendre en compte les divers biais et effets de bords induits par la prise en compte d'aspects sémantiques et non géographiques. Notre objectif principal est d'établir une mesure de proximité efficace, prenant en compte la majorité des points présentés dans la section 3.2 et ne nécessitant qu'une connaissance restreinte et locale du système.

## **BIBLIOGRAPHIE**

- [Gnutella, WWW] The Gnutella project, <http://www.gnutella.com/>
- [Napster, WWW] The Napster project, <http://www.napster.com/>
- [eDonkey, WWW] The E-Donkey project, <http://www.edonkey2000.com/>
- [eMule, WWW] The E-mule project, <http://www.emule-project.net/>
- [KaZaA, WWW] The KaZaA project, <http://www.kazaa.com/>
- [Castro, 2003] Castro M., Costa M., and Rowstron A. : "Should we build gnutella on a structured overlay?". In the 2nd Workshop on Hot Topics in Networks (HotNets-II), Cambridge, MA, USA.
- [Jelasity, 2004a] Jelasity M., Guerraoui R., Kermarrec A.-M., Van Steen M. : "The Peer Sampling Service : Experimental Evaluation of Unstructured Gossip-Based Implementations". In the 5th International Middleware Conference ACM/IFIP/USENIX (Middleware'04), Toronto, Ontario, Canada.
- [Jelasity, 2004b] Jelasity M., Babaoglu O. : "T-Man : Fast Gossip-based Construction of Large-Scale Overlay Topologies". Technical Report, University of Bologna, Department of Computer Science, UBLCS-2004-7, Bologna, Italy.
- [Handurukande, 2004a] Handurukande S., Kermarrec A.-M., Le Fessant F., Massoulié L. : "Exploiting Semantic Clustering in the eDonkey P2P network". In the 11th ACM SIGOPS European Workshop (SIGOPS'04), Leuven, Belgium.
- [Handurukande, 2004b] Handurukande S., Kermarrec A.-M., Le Fessant F., Massoulié L. : "Clustering in peer-to-peer file sharing workloads". In the 3rd International Workshop on Peer-to-Peer Systems (IPTPS'04), San Diego, CA, USA.
- [Voulgaris, 2004] Voulgaris S., Kermarrec A.-M., Massoulié L., and Van Steen M. : "Exploiting semantic proximity in peer-to-peer content searching." In the 10th International Workshop on Future Trends in Distributed Computing Systems (FTDCS'04), Suzhou, China.
- [Sripanidkulchai, 2003] Sripanidkulchai K., Maggs B., and Zhang H. : "Efficient content location using interest-based locality in peer-to-peer systems." In the 22nd Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE INFOCOM'03), San Francisco, CA, USA.