



HAL
open science

Exploitation d'approches système dans les réseaux sans fil

Frédéric Weis

► **To cite this version:**

Frédéric Weis. Exploitation d'approches système dans les réseaux sans fil. Informatique ubiquitaire. Université Rennes 1, 2012. tel-00790484v2

HAL Id: tel-00790484

<https://theses.hal.science/tel-00790484v2>

Submitted on 20 Feb 2013

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.

Habilitation à Diriger des Recherches

présentée devant l'Université de Rennes 1

Mention INFORMATIQUE

par

Frédéric Weis

Exploitation d'approches système dans les réseaux sans fil

Soutenue le 6 Juin 2012 devant le jury composé de

Michel Banâtre	Examineur
Yolande Berbers	Rapporteur
Andrzej Duda	Rapporteur
Olivier Festor	Rapporteur
Denis Rouffet	Examineur
David Simplot	Examineur
César Viho	Président

Table des matières

Introduction et fil conducteur de ma recherche	5
Exploitation des interactions sans fil courte portée « directes »	6
Exploitation des interactions sans fil courte portée dans un réseau étendu	7
Couplage système de deux réseaux mobiles	8
Structure du document	9
Partie 1 : Système d'Information Spontanés	13
1 Étude et mise en œuvre des Systèmes d'Information Spontanés	13
1.1 Principes généraux d'un S.I.S	14
1.1.1 Analyse des mécanismes sous-jacents	14
1.1.2 Définitions	15
1.2 Analyse de deux scénarios	18
1.3 Synthèse des problèmes de recherche identifiés	18
1.4 Découverte et construction du voisinage physique	20
1.4.1 Définition d'une nouvelle relation de voisinage	21
1.4.2 Calcul automatique de la fréquence d'annonce de présence	21
1.4.3 Bilan	23
1.5 Base de données de proximité	23
1.5.1 Présentation des mécanismes proposés	25
1.5.2 Mise en œuvre de l'architecture PERSEND	27
1.5.3 Bilan	28
1.6 Glanage opportuniste au sein du Web de proximité	28
1.6.1 Gestion de la base d'information du terminal	29
1.6.2 Mécanismes de découvertes d'informations entre deux terminaux d'un S.I.S	30
1.6.3 Bilan	31
1.7 Conclusion	32

Partie 2 : Réseaux 4G	35
2 Mise en œuvre des réseaux à couverture discontinue	35
2.1 Gestion des flux descendants	36
2.1.1 Distribution de plusieurs niveaux de cache dans l'architecture	37
2.1.2 Analogie avec les systèmes multiprocesseurs	39
2.1.3 Une hiérarchie de trois niveaux de cache	40
2.1.4 Discrimination des débits dans les bulles radio	42
2.1.5 Analyse de l'impact de la hiérarchie de caches sur la délivrance des flux descendants	45
2.1.6 Bilan et principaux enseignements	46
2.2 Les débits augmentent ... et un niveau de cache disparaît	48
2.2.1 Traitement de la discontinuité au niveau de l'AC	48
2.2.2 Définition d'un protocole de transport de cache supportant la discontinuité	49
2.2.3 Bilan	51
2.3 Exploitation efficace des flux montants	52
2.3.1 Problématique terminal-AC	54
2.3.2 Problématique AC-serveur	55
2.3.3 Bilan	56
2.4 Déploiement à grande échelle d'un réseau à couverture discontinue	57
2.4.1 Mécanisme de découverte de l'AC candidat	58
2.4.2 Performances de la version distribuée du mécanisme de découverte	59
2.5 Conclusion	60
3 Couplage d'un réseau DVB avec un réseau cellulaire 3G	63
3.1 Analyse de scénarios de couplage	65
3.2 Basculement de contenus 3G vers un réseau DVB	68
3.2.1 Mise en œuvre d'un service programmé	69
3.2.2 Mise en œuvre d'un service non programmé	71
3.3 Délivrance de contenus personnalisés via le réseau 3G	73
3.4 Conclusion	75
4 Bilan et perspectives	77
4.1 Bilan des travaux	77
4.1.1 Mécanismes système pour les interactions sans fil courte portée	77
4.1.2 Mécanismes système pour les réseaux mobiles 4G	78
4.2 Perspectives	79
4.2.1 Réseaux 4G : vers une gestion du multi-attachement	79
4.2.2 Vers une exploitation d'un contexte local enrichi	80
Bibliographie	85
Table des figures	87

Introduction et fil conducteur de ma recherche

Ce document d'habilitation est consacré aux travaux de recherche que j'ai menés depuis 1998, date de mon arrivée à l'IRISA Rennes, au sein du projet SOLIDOR. Après une thèse soutenue au CNAM en 1996 dans le domaine des tests de protocoles réseaux, je souhaitais mettre à profit mon rattachement à un nouveau laboratoire pour initier des travaux de recherche en lien avec mes thématiques d'enseignements, principalement les technologies de l'Internet et les réseaux sans fil.

L'équipe SOLIDOR travaillait depuis plusieurs années sur des supports d'exécution pour les systèmes distribués. Ce domaine était en pleine évolution, principalement sur deux points. D'une part, l'avènement du réseau Internet amenait l'étude de systèmes distribués sur une grande échelle, tant du point de vue du nombre d'entités impliquées, que de leur étalement géographique. D'autre part, le développement conjoint des technologies mobiles et des réseaux sans fil introduisait le principe de nœuds mobiles connectés ou déconnectés du système.

Je me suis intéressé à ce deuxième aspect : l'avènement progressif d'une informatique distribuée et mobile. Au tout début des années 1990, les supports traditionnels de l'informatique distribuée, à savoir les ordinateurs personnels, les serveurs et les stations de travail, étaient largement cantonnés à un usage statique, sur un site donné. Ce panorama a commencé à évoluer de manière radicale, suite à deux tendances importantes :

- Les progrès technologiques constant dans le domaine de l'informatique embarquée (ordinateurs portables, puis téléphones portables) : miniaturisation des composants, accroissement des capacités de stockage, de calculs, d'affichage, intégration de nouvelles interfaces de communication.
- Le développement des communications sans fil.

Différentes familles de réseaux sans fil sont à même de servir de support aux services de l'informatique mobile. Les *réseaux cellulaires étendus*, utilisés initialement pour le support de la téléphonie mobile, sont constitués d'un ensemble de *stations de base*, définissant chacun une *cellule* de communication. Chaque station ne peut prendre en charge qu'un nombre limité de terminaux. En fonction de la densité des terminaux dans une zone, il est possible de jouer sur la taille et le nombre de cellules pour assurer efficacement une couverture étendue et continue. Cette approche implique une utilisation fine des ressources radio. En effet, les cellules proches se recouvrent partiellement, afin d'offrir en permanence à un utilisateur mobile un lien radio vers une station de base. En conséquence, ces cellules doivent exploiter des bandes de fréquence totalement disjointes. Les opérateurs des réseaux cellulaires parlent de planification des fréquences.

A l'opposé de cette complexité, on trouve les réseaux locaux sans fil (WLAN¹). Les plus répandus sont les réseaux WiFi, très simples à déployer et peu coûteux à mettre en œuvre. Ils reposent sur différentes normes (802.11b/g, 802.11a, 802.11n), offrant des débits de plusieurs dizaines de Mb/s. Ces réseaux peuvent être exploités sous forme de *HotSpots*. Ils sont alors centrés sur un point d'accès ou AP² (comparable à une station de base dans un réseau cellulaire), qui délimite une cellule radio de taille très limitée (de l'ordre d'une centaine de mètres). On parle alors de

1. WLAN : *Wireless Lan Area Network*
2. AP : *Access Point*

Pico Cell, ou de bulle radio. Un utilisateur se trouvant dans cette zone limitée de couverture peut alors s'attacher au point d'accès, et bénéficier d'une connectivité de plusieurs dizaines de Mb/s. On peut parler alors d'interactions sans fil courte portée. À noter que le nombre limité de bandes de fréquences disponibles sur lesquelles ces technologies WLAN sont autorisées à fonctionner, combiné à la faible taille des bulles radio, ne permet pas *a priori* d'envisager des déploiements continus et étendus comme c'est le cas pour des réseaux cellulaires exploités par des opérateurs. En plus de ce mode de fonctionnement centré sur un point d'accès, les communications WiFi peuvent également s'établir directement entre entités mobiles proches, sans passer par une infrastructure.

En 1998, les technologies WiFi venaient seulement d'être normalisées, et commençaient à équiper les ordinateurs portables. Leur intégration à des équipements plus réduits n'était encore qu'une hypothèse, alors que c'est maintenant une réalité. Mais l'utilisation des interactions sans fil courte portée qu'elles autorisaient, soient directes, soit via un AP, ouvraient de toute évidence de nouvelles voies de recherche dans le domaine de l'informatique mobile.

Exploitation des interactions sans fil courte portée « directes »

Le principe de l'informatique mobile est de mettre à disposition en permanence, indépendamment de la localisation de l'utilisateur, l'ensemble des services normalement accessibles au bureau ou au domicile. Les premiers travaux effectués dans ce domaine cherchent le plus souvent à masquer la variabilité des environnements mobiles, afin de rendre compatibles les services conçus initialement pour des supports traditionnels. La prise en compte de cette variabilité passe par des mécanismes d'adaptation distribués en différents points du système : adaptation au niveau du nœud mobile, via par exemple la prise en compte de la taille limitée de l'écran, adaptation aux fluctuations de performances du réseau etc. Cette variabilité est une conséquence directe de la mobilité. En la masquant, ce sont avant tout les effets de la mobilité qu'on cherche à masquer au système.

Ce constat m'a amené à poser la question suivante : *Est-il possible de construire un système distribué, en exploitant la mobilité des nœuds, plutôt qu'en la masquant ?* Des équipements comme des assistants numériques personnels ou PDA³, parce qu'ils sont portés par des utilisateurs mobiles, peuvent se trouver regroupés de manière temporaire dans un même voisinage physique. Et l'utilisation de technologies sans fil directes et courte portée, doit alors leur permettre d'échanger « spontanément » des informations pertinentes, tant que la distance qui les sépare le permet. Encore faut-il identifier des domaines d'application pouvant tirer parti d'un système d'information exploitant automatiquement la proximité des nœuds mobiles le composant.

Cette idée, très intuitive au départ, m'a amené à m'intéresser aux principes de l'*ubiquité numérique*. Ces principes reposent essentiellement sur la dualité entre monde réel et système d'information. Les entités mobiles étant engagées très souvent à la fois dans l'environnement réel et dans des environnements numériques, un couplage entre les deux est nécessaire pour faciliter les activités. Un tel couplage passe par une acquisition de l'état de l'environnement réel par les systèmes informatiques, afin de pouvoir automatiquement réagir aux évolutions de cet état. C'est ce qu'on appelle la sensibilité au contexte (*context-awareness*). Grâce à cette intégration, on passe d'un mode d'interactions explicites entre le système d'information et les utilisateurs mobiles (qui mobilise l'attention des utilisateurs), à un mode d'interactions implicites (automatiques). Dans beaucoup de systèmes exploitant les principes de l'ubiquité numérique, les communications entre les entités mobiles et les composants informatiques délivrant le service passe par une borne. Ici, l'idée est de définir des services automatiques entre utilisateurs mobiles, uniquement par le biais de communications directes et courte portée. En d'autres termes, le contexte repose sur la proximité physique des nœuds qui composent le système.

Nous parlons alors de *Systèmes d'Information Spontanés* (S.I.S). Mon premier axe de recherche concerne l'étude de tels systèmes. Avant d'envisager des solutions, il était indispensable d'analyser et de prendre en compte la nature très contraignante des communications sans fil directes et courte portée : la mobilité des nœuds n'est pas obligatoirement connue ou maîtrisée, et par voie

3. PDA : *Personal Digital Assistant*

de conséquence, les communications entre entités voisines peuvent être brèves et s'interrompre de manière brutale. Le système d'information, distribué entre les nœuds voisins, est donc par nature très dynamique. Sa composition dépend étroitement de l'entrée et de la sortie des entités mobiles d'un même voisinage physique. En premier lieu, j'ai donc cherché des outils et des mécanismes me permettant de prendre en compte toutes ces contraintes.

L'équipe de recherche SOLIDOR, que je venais d'intégrer, possédait une longue expertise dans la conception des systèmes d'exploitation, et plus précisément dans la construction de supports d'exécution pour des systèmes à grande échelle. Une des difficultés majeures pour ces systèmes est d'assurer l'*extensibilité*, c'est à dire la capacité à délivrer un service collectif pour un grand nombre d'utilisateurs, sans dégrader la qualité du service offerte individuellement à chaque usager. Différentes techniques système peuvent être utilisées pour traiter le problème de l'extensibilité, comme la distribution et la réplique des données, l'utilisation de caches, la répartition de charge etc. Les travaux décrits dans [22] avaient particulièrement retenu mon attention. Ils proposaient une organisation des informations de manière à favoriser le rapprochement des groupes de données et de groupes d'utilisateurs *a priori* intéressés par ces données. Et surtout, la constitution de ces groupes s'appuyait sur des techniques de *data mining*. Ces dernières permettent de découvrir des relations structurelles cachées entre les enregistrements d'une base de données. Et appliquées à un système distribué à grande échelle, ces techniques permettaient d'organiser sous forme arborescente les groupes de données, de manière à permettre aux groupes d'utilisateurs associés un accès efficace aux informations, autrement dit un accès *réactif*.

C'est sur ce dernier point que j'ai établi une analogie avec le principe des S.I.S. En effet, la volatilité des liens sans fil entre les nœuds voisins et mobiles impose que les informations puissent être échangées rapidement. La réactivité dans l'accès aux données est donc une propriété essentielle. Dans ces conditions, un nœud client doit savoir ce qu'il cherche dans le voisinage, et de l'autre côté, un nœud serveur doit pouvoir exprimer clairement quels types d'informations il peut offrir à ce nœud client. L'utilisation de techniques de *data mining* nous a conduit à proposer une gestion de la base d'information d'un terminal mobile reposant sur une indexation thématique et automatique de tous les documents stockés par l'utilisateur. Cette indexation est ensuite utilisée par les terminaux pour échanger spontanément des informations, dès lors que la proximité des utilisateurs l'autorise.

À partir de ces premiers travaux, j'ai acquis la conviction que des techniques utilisées dans le domaine des systèmes d'exploitation et des systèmes distribués pouvaient être exploitées pour mettre en œuvre efficacement les S.I.S. Dans un premier temps, j'ai donc cherché, en m'appuyant sur des *approches système*, à exploiter des communications sans fil directes et courte portée. Cette démarche, appliquée au cadre des S.I.S, m'a permis de traiter différents problèmes relevant de l'adressage et du contrôle de l'ensemble des informations distribuées entre des nœuds mobiles voisins. Je l'ai ensuite poursuivi dans le cadre d'autres architectures exploitant les communications sans fil. C'est l'objet de la suite de cette introduction.

Exploitation des interactions sans fil courte portée dans un réseau étendu

Mes travaux sur les S.I.S m'ont rapidement amené à envisager la possibilité de connecter ponctuellement une ou plusieurs bornes fixes de communication sans fil courte portée à un S.I.S existant, autrement dit de relier un AP à un ensemble de nœuds voisins et mobiles exploitant temporairement un système d'information. Nous avons fait le constat suivant : les terminaux voisins échangeant automatiquement de l'information disposent de capacités de stockage limitées. Phénomène aggravant, ils peuvent être équipés de mécanismes de capture, comme l'enregistrement de la voix, ou d'un appareil photo numérique. Ces fonctions ne font que renforcer la nécessité de dépasser les limites de stockage local d'un terminal.

De là découle l'idée de disposer d'îlots de connectivité haut débit dans l'espace physique, pour permettre à des terminaux mobiles de « décharger » automatiquement leurs données chaque fois qu'ils rencontrent une bulle radio. Les interactions sans fil courte portée offertes par un AP WLAN

autorisent ce type de déploiement. Si l'idée de relier ce principe à nos travaux sur les S.I.S n'a finalement pas été poussée plus loin, nous avons estimé que le déploiement d'un ensemble de bulles radio (donc d'APs) offrait de véritables perspectives dans le domaine de l'informatique mobile. On peut envisager une infrastructure mobile étendue, très simple à déployer, sans en passer par la complexité d'une planification radio inhérente aux réseaux cellulaires : les bulles, appelées également *Pico Cells*, sont interconnectées, mais ne se recouvrent pas, afin d'éviter les interférences au niveau radio.

La principale difficulté d'une telle architecture réside dans la mise en place d'une continuité de service entre les APs. Pris isolément, les points d'accès offrent simplement des îlots de connectivité, qu'on désigne souvent sous le terme de *HotSpot*. Dans le cadre d'une véritable infrastructure étendue, l'objectif est de fournir un service de transfert de données au terminal via la cellule à laquelle il se trouve temporairement rattaché. Autrement dit, toute la difficulté consiste à offrir un service continu, en dépit de la discontinuité de la couverture.

Cette problématique de recherche renvoie à l'exploitation d'architectures appelées *réseaux d'infostations*. De nombreux travaux ont proposé des solutions pour des topologies de déploiement très spécifiques, par exemple lorsque les utilisateurs se déplacent sur une route le long de laquelle un ensemble de points d'accès est disposé. Nous avons estimé qu'il devait être possible d'aller plus loin, sans imposer de topologie particulière pour disposer les APs dans l'espace, et sans connaître la mobilité des utilisateurs.

L'objectif de nos travaux est de proposer, dans le cadre de *réseaux à couverture discontinue* à faibles coûts de déploiement, des solutions systèmes permettant de masquer la discontinuité de la couverture radio. Pour ce faire, nous avons considéré la possibilité d'introduire des caches dans l'infrastructure réseau. Un cache est une mémoire de stockage de réplicas de données. Son exploitation a été étudiée dans différents domaines : les systèmes de fichiers, les processeurs, les bases de données etc. Disposé dans un système d'information distribué, il intercepte et conserve toutes les informations qui transitent par lui. Dans le cadre d'un processeur, l'introduction de différents niveaux de caches est motivée par des temps d'accès aux données très différents au niveau du cœur du processeur ou de la mémoire centrale. Partant d'une analogie entre les systèmes multiprocesseurs et un réseau d'infostations, nous avons montré qu'une hiérarchie de caches sur plusieurs niveaux permet de s'affranchir de la discontinuité de la couverture radio.

Mais cette hiérarchie n'est pas suffisante pour offrir un service continu à un nombre important d'utilisateurs mobiles. Il faut pour cela que les caches distribués dans le réseau soient correctement alimentés en données. Idéalement, cela revient à prédire la localisation de l'utilisateur dans l'infrastructure, et à positionner « à l'avance » ces données au plus près de l'utilisateur. On peut ici faire une analogie avec une politique d'anticipation des accès utilisée dans un système de fichiers. Une telle politique a pour but de diminuer les temps d'accès aux informations dont l'utilisateur a besoin, en profitant au mieux de la localité des données dans le système. Le plus souvent, elle repose sur trois phases : l'*analyse des accès*, la *prédiction des accès ultérieurs* et le *préchargement des données prédites*. Dans un réseau à couverture discontinue, l'anticipation peut se traduire par trois phases également : l'*analyse des déplacements d'un utilisateur*, la *prédiction de la prochaine bulle radio traversée par cet utilisateur*, et le *préchargement vers un cache proche de cette prochaine bulle radio*. Nous avons fait le constat que ce type de prédiction est très difficile à obtenir dans un réseau large échelle de *Pico Cells* dans lequel la mobilité des utilisateurs est peu contrainte. Nous avons du rechercher d'autres critères permettant l'anticipation des accès. Notre approche, présentée dans le chapitre consacré à nos travaux sur les réseaux à couverture discontinue, s'appuie sur une discrimination des débits radio offerts à l'intérieur de chaque *Pico Cell*.

Couplage système de deux réseaux mobiles

En démarrant des travaux sur les réseaux à couverture discontinue, nous avons identifié un aspect important : l'exploitation sur une grande échelle de bulles radio WLAN pouvait s'inscrire dans un thème de recherche plus large de l'informatique mobile en général, et des réseaux sans fil en particulier : les réseaux cellulaires de quatrième génération ou *réseaux 4G*. L'objectif de ces

recherches, initiées à la fin des années 90 dans différents domaines (problèmes radio, nouvelles infrastructures réseaux, gestion de la mobilité, définition de nouveaux usages et services), était de définir une infrastructure cellulaire large échelle, offrant des débits très élevés (plus de 100 Mb/s), et en mesure de servir une densité importante d'utilisateurs mobiles. Une approche possible consiste à considérer la future architecture 4G comme un nouveau réseau radio très haut débit, en s'appuyant sur des évolutions des infrastructures 3G (troisième génération) existantes. Mais les recherches sur les réseaux 4G dépassent le cadre de ces évolutions normatives. Une voie possible pour atteindre les objectifs affichés en termes de débits et de densité d'utilisateurs est de s'appuyer sur un ensemble de technologies d'accès radio, utilisées en fonction de leur complémentarité. On parle alors de *couplage de réseaux*.

Ainsi, les technologies sans fil courte portée comme WiFi permettent des échanges au delà de 100 Mb/s à proximité d'un point d'accès, et répondent aux critères de performances des futurs réseaux 4G. C'est ce constat qui a motivé mes travaux initiaux sur les réseaux à couverture discontinue. Mais au delà des WLANs, d'autres réseaux mobiles sont susceptibles de s'intégrer à la problématique 4G. Ainsi, DVB⁴ est la famille de normes qui a été adoptée pour la diffusion de services de télévision numérique et multimédia. Initialement dédiés à la diffusion de programmes « traditionnels » sur les réseaux terrestres, câblés et satellitaires, les travaux de normalisation se sont poursuivis afin de permettre la réception de programmes TV sur des terminaux mobiles.

Un réseau DVB offre un mécanisme de diffusion « de masse » : il permet, par le biais d'un seul canal, de servir l'ensemble des utilisateurs se trouvant dans la zone couverte par le réseau. Dans sa déclinaison mobile, une telle infrastructure présente un réel intérêt pour mettre en œuvre les services multimédia visés par les réseaux 4G. Le problème est que le canal de diffusion est unidirectionnel, l'interface DVB d'un terminal ne pouvant que recevoir des données. Il n'est donc pas possible d'offrir des services interactifs aux utilisateurs, comme dans le cadre des réseaux à couverture discontinue. La complémentarité des réseaux, évoquée plus haut, prend tout son sens. En effet, une voie montante existe au niveau des réseaux cellulaires. On parle alors d'*utilisation couplée d'un réseau mobile unidirectionnel de diffusion et d'un réseau cellulaire bidirectionnel*. Il s'agit d'une problématique relativement neuve dans le domaine des réseaux 4G. Afin de couvrir un spectre de recherche sur les réseaux 4G qui dépassent l'utilisation des communications sans fil courte portée, j'ai souhaité aborder la problématique du couplage d'un réseau DVB et d'un réseau cellulaire bidirectionnel.

Tout comme dans les réseaux à couverture discontinue, nous cherchons à intégrer des mécanismes systèmes au sein l'infrastructure pour traiter les aspects recherche que nous avons identifiés. Notre approche, présentée dans ce document, consiste à distribuer des mécanismes de gestion des données, leur but étant d'aiguiller de manière opportune et efficace les flux de données vers le réseau DVB ou vers le réseau cellulaire.

Structure du document

Ce document présente une synthèse de mes principaux résultats de recherche. Il est divisé en deux parties. La partie 1 présente mes travaux dans le domaine des Systèmes d'Information Spontanés. La partie 2, centrée sur la problématique des réseaux 4G, comporte deux chapitres. Le premier chapitre s'intéresse à la mise en œuvre des réseaux à couverture discontinue. Enfin, le second chapitre présente notre démarche pour le couplage système entre un réseau DVB et un réseau cellulaire bidirectionnel.

4. DVB : *Digital Video Broadcasting*

Partie 1 : Système d'Information Spontané

ÉTUDE ET MISE EN ŒUVRE DES SYSTÈMES D'INFORMATION SPONTANÉS (S.I.S)

Les technologies de communication sans fil à courte portée rendent possibles la mise en œuvre de communications directes entre tout type d'objets physiques se trouvant dans la même « bulle radio ». Dans un tel cadre, aucune borne fixe (*i.e.* point d'accès) n'est utilisée. On peut alors parler d'*interactions de proximité*. Cette notion de proximité reste bien entendu dépendante de la technologie utilisée, allant de quelques mètres avec une technologie comme Bluetooth, jusqu'à plusieurs dizaines de mètres pour les produits de la famille WiFi (802.11a/b/g/n). L'utilisation de telles interactions directes est longtemps restée marginale, tant dans le cadre de l'informatique mobile, que pour l'ubiquité numérique. Le plus souvent, les communications sans fil de proximité passent par une borne, qui joue le rôle de passerelle vers un réseau fixe (dans le cadre d'un *HotSpot* par exemple), ou qui fournit des services de nature contextuelle, comme c'est le cas de l'étude ParcTab [64], issue des travaux pionniers de Marc Weiser.

Pariant sur un nombre croissant de terminaux mobiles dotés (1) de capacités de calcul et de stockage croissantes, et (2) d'une interface de communication à courte portée, nous avons envisagé dès 1998 la mise en œuvre de nouvelles applications tirant pleinement parti de la proximité physique des calculateurs. On peut considérer l'exemple d'une conférence. Chaque participant est porteur d'un assistant numérique personnel (PDA¹), dans lequel il a stocké sa présentation, ses publications récentes, un lien vers sa page Web personnel etc. On peut supposer que toutes ces personnes présentes à une même conférence, partagent des centres d'intérêts communs. Leur proximité physique est donc une opportunité de « glaner » spontanément des informations pertinentes. Dans ce cadre, c'est la présence dans un même voisinage physique de plusieurs terminaux qui justifie l'échange d'informations. En d'autres termes, le contexte repose sur la proximité physique des terminaux mobiles.

Convaincus que les interactions de proximité sans fil trouvent toute leur place dans le cadre de travaux exploitant des mécanismes système, nous avons proposé « la génération spontanée de systèmes d'information ». Un *Système d'Information Spontané* (S.I.S) exploite des entités mobiles, et s'appuie sur des propriétés existant dans le monde réel, comme la proximité physique de ces entités, pour en déduire automatiquement un système d'information. La présentation de nos premières idées se trouve dans un premier rapport de recherche [13].

On peut présenter intuitivement ce principe de la manière suivante : il y a génération spontanée d'un système d'information lorsqu'au moins deux terminaux mobiles physiquement « voisins » peuvent établir une communication, et échanger de façon implicite ou explicite des informations. Ces échanges s'effectuent directement entre ces terminaux, sans s'appuyer sur une infrastructure externe comme des bornes fixes ou un réseau d'antennes. Un tel système d'information va ensuite évoluer, et éventuellement disparaître, en fonction de la mobilité des entités et de leurs interactions.

1. PDA : *Personal Digital Assistant*

L'objectif de ce premier chapitre est de donner une vue d'ensemble de nos travaux concernant les Systèmes d'Information Spontanés. Dans les sections 1.1 et 1.2, nous donnons nos motivations pour l'étude de tels systèmes. Dans la section 1.3, nous présentons les différents problèmes systèmes que nous avons identifiés dans le cadre de la mise en œuvre d'un S.I.S, en insistant tout particulièrement sur les contraintes à assurer pour chacun d'entre eux. Enfin, la synthèse des solutions proposées est rassemblée dans les sections 1.4, 1.5 et 1.6.

1.1 Principes généraux d'un S.I.S

L'idée d'une architecture pour des Systèmes d'Information Spontanés est partie de l'observation de la façon dont s'établit un dialogue entre des personnes et des mécanismes sous-jacents qui sont alors implicitement mis en place. L'échange d'informations, entre deux personnes et mobiles, ne peut avoir lieu que lorsqu'elles sont à « portée de voix », ceci caractérise dans notre approche « le voisinage ». Différentes situations sont donc possibles. Elles se rapprochent l'une de l'autre, s'arrêtent pour dialoguer, se découvrent progressivement des centres d'intérêts communs, s'échangent des informations, puis se séparent une fois la conversation terminée. La fin de l'échange est alors négociée par les deux personnes. Mais il se peut également qu'elles se croisent rapidement, sans possibilité de s'arrêter. Dans ce cas, le dialogue est bref, sans aucune garantie que la totalité des informations puisse être correctement échangée.

1.1.1 Analyse des mécanismes sous-jacents

Cette description de la communication humaine, bien que schématique, nous a permis de dégager les principaux mécanismes mis en œuvre pour échanger des informations. Ainsi, une analyse plus fine, mais toujours macroscopique, de ce scénario, fait ressortir trois éléments dans la mise en place et le déroulement de la communication.

Le premier mécanisme est *la phase d'identification*. Il concerne la détection par une personne de la présence d'un autre humain se trouvant « à portée de voix ». Généralement, ce problème est traité de manière inconsciente. Le second élément est *la phase de dialogue et d'interactions*. Cette étape est consciente et concerne l'échange effectif des informations. Elle traite plusieurs problèmes qui peuvent être résumés au travers des questions suivantes :

- De quelle manière une personne découvre-t-elle les informations que son voisin peut lui communiquer ?
- Comment cette personne détermine-t-elle que ces informations sont pertinentes en ce qui le concerne ?
- Comment enfin exploite-t-elle ces informations pertinentes ?

Le troisième mécanisme concerne la *gestion des interactions*, étroitement lié à la phase de dialogue qui vient d'être décrite. Les échanges sont effectués en estimant éventuellement de façon implicite la durée maximale de la communication. Du point de vue d'une personne, cette estimation s'appuie sur la trajectoire et la vitesse perçues du voisin.

La prise en compte de cette durée permet de déterminer la quantité d'information pouvant être transmise, tant que les humains restent à « portée de voix ». Dans le cas où ces derniers estiment qu'ils resteront peu de temps dans le même voisinage, l'un ou l'autre ne transmettra que des informations prioritaires pendant la phase de dialogue. L'approche est bien entendu différente si le temps de communication n'est pas borné (dans le cas par exemple où les personnes s'arrêtent pour dialoguer). Cet échange peut se terminer de manière explicite, par accord tacite. Mais il peut également s'interrompre de manière imprévue. Dans ce cas, l'information récupérée peut être partielle, utilisable ou non, ce qui peut donner lieu à une interprétation erronée de son contenu. L'ensemble de ces mécanismes est illustré par la figure 1.1.

Au final, la mise en place de la communication et l'échange d'information au sein d'un voisinage physique reposent sur trois étapes élémentaires : la *détection de présence* des entités se trouvant

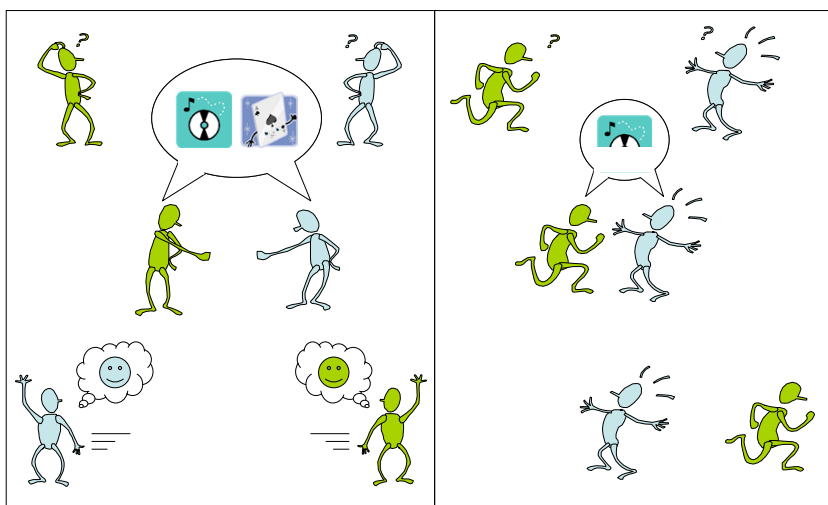


FIGURE 1.1 – Communication entre humains

à portée de communication, la *gestion de la communication* à « portée de voix », en prenant en compte ce voisinage physique, et enfin la *découverte et l'exploitation des informations pertinentes* au sein de ce voisinage physique. Cette analyse simple de la communication entre personnes proches nous a permis de caractériser *la communication par voisinage entre des systèmes informatiques*. C'est l'objet de la section suivante.

1.1.2 Définitions

Je définis ici les principales notions à partir desquelles j'ai conduit mes travaux sur les S.I.S. Pour plus de détails, le lecteur pourra se reporter à [13]. L'environnement considéré est composé d'entités mobiles E_i qui ne peuvent communiquer que lorsqu'elles sont physiquement proches (*i.e.* à portée de communication). Suivant l'application considérée, une entité E_i peut être par exemple un robot, un véhicule, un humain etc. ou tout objet du monde réel susceptible d'être couplé avec un ordinateur. Elle est définie par le couple (EE_i, EM_i) dans lequel EM_i (Entité Mobile) caractérise la partie « mécanique » d'une E_i , et EE_i (Entité autonome Embarquée) est la partie qui réalise la communication, le traitement et le stockage de l'information. Cette partie est mise en œuvre à l'aide d'un ordinateur, auquel on a adjoint les éléments de communication adéquats.

Une EE_i peut échanger des informations avec son EM_i par le biais d'une interface locale. La nature et le fonctionnement de cette interface dépendent bien entendu de l' EM_i associée (un robot, un bus etc.) et de la classe d'application considérée.

L'interface de communication implantée au sein de l' EE_i définit une zone de communication notée $ZC(EE_i)$. La taille de cette zone dépend de la technologie utilisée, allant de quelques mètres avec Bluetooth, jusqu'à plusieurs dizaines de mètres avec une des normes de la famille WiFi. Il est difficile de représenter cette zone de communication, tant sa forme dépend des impacts difficilement prévisibles de l'environnement de EE_i , tels que les problèmes liés au *shadowing* ou aux chemins multiples. Les normes de communication sans fil tiennent compte de ces phénomènes et tentent de les compenser. Malgré tout, la frontière de $ZC(EE_i)$ reste difficilement prévisible. Nos travaux s'appuient sur une approximation communément utilisée, qui consiste à modéliser la zone de communication comme une sphère centrée sur le terminal. La figure 1.2 résume le modèle retenu.

Communication par voisinage. Lorsqu'une entité autonome EE_j se trouve à l'intérieur de la zone de communication $ZC(EE_i)$ avec $i \neq j$, les mécanismes implantés au sein de EE_i lui permettent de détecter la présence de EE_j , puis de communiquer avec EE_j (et réciproquement) :

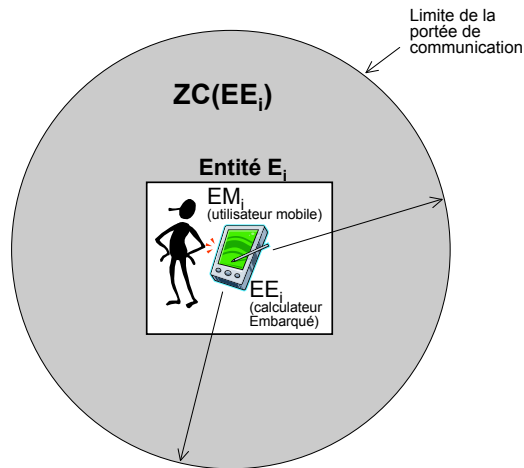


FIGURE 1.2 – Principe d'une entité mobile (exemple avec un être humain)

c'est ce que nous définissons comme une possibilité de « communication par voisinage ». Cette communication ne peut exister que si EE_i et EE_j se trouvent dans l'intersection des zones de communication $ZC(EE_i)$ et $ZC(EE_j)$. Ce mécanisme est illustré par la figure 1.3. Cette communication est dynamique, elle peut être maintenue ou se rompre, en fonction des déplacements respectifs de EE_i et EE_j . Elle est fonction du temps. $CV(EE_i, EE_j, t_i)$ est le prédicat qui caractérise le fait que EE_i et EE_j « communiquent par voisinage » à l'instant t_i (temps local de EE_i).

La notation $EE_i \in ZC(EE_j)$ signifie que EE_i appartient à la zone de communication de EE_j . On a donc la relation suivante :

$$CV(EE_i, EE_j, t_i) \iff \text{l'instant } t_i, EE_i \in ZC(EE_j) \wedge EE_j \in ZC(EE_i)$$

Système d'Information Spontané. Une entité autonome embarquée EE_i peut entretenir, à un instant donné, plusieurs communications par voisinage avec des calculateurs se trouvant à portée de communication. Dans ces conditions, cette même entité peut être à l'origine de la création dynamique d'un système d'information à l'instant t_i lorsqu'il existe au moins une autre entité autonome embarquée EE_j qui puisse communiquer par voisinage avec EE_i à l'instant t_i . De manière plus formelle, ce système d'information spontané est défini de la manière suivante :

$$SIS(EE_i, t_i) = \{EE_j \mid \forall j \text{ avec } i \neq j, CV(EE_i, EE_j, t_i)\}$$

D'après cette définition, un S.I.S donné est toujours défini par rapport à une entité autonome embarquée donnée. On notera également qu'une entité EE_i peut appartenir simultanément à plusieurs S.I.S. Le fait que EE_i soit mobile lui permet d'entrer ou de sortir spontanément de la zone de communication d'autres entités mobiles. La composition d'un S.I.S évolue dynamiquement dans le temps, en fonction des déplacements respectifs des entités la composant.

Espace visible d'un S.I.S. Les terminaux embarqués sont capables de stocker et de traiter ces informations. Nous supposons que ces informations peuvent être modélisées sous forme d'objets. L'utilisateur EM_i du calculateur qui lui est associé EE_i peut souhaiter autoriser ses voisins à consulter une portion de son système d'information. Nous distinguons donc deux classes d'objets : les objets partagés, accessibles par des entités voisines dans le cadre d'un S.I.S, et les objets

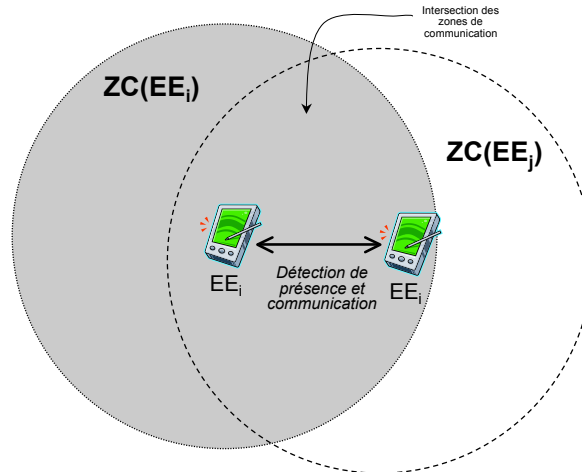
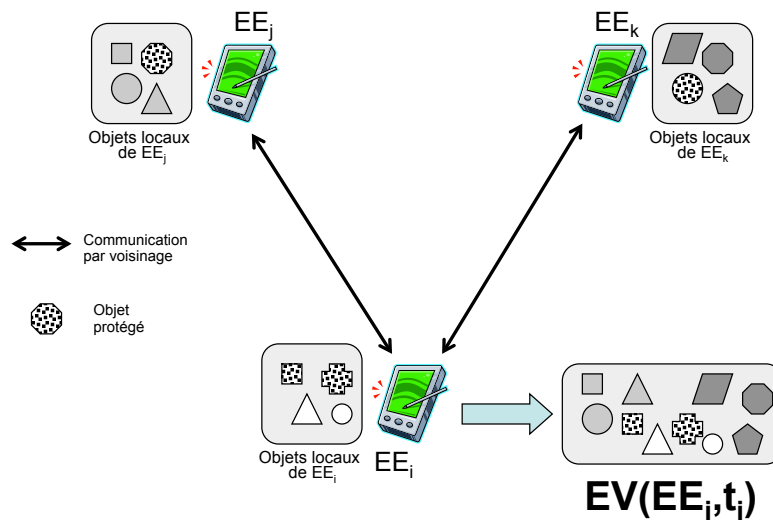


FIGURE 1.3 – Communication par voisinage entre deux entités mobiles

protégés, qui ne peuvent être utilisés que par l'entité qui les possède. Dans ce cadre, nous notons $\{ShO_i\}$ (resp. $\{PrO_i\}$) l'ensemble des objets partagés (resp. protégés) de EE_i . Nous désignons alors par *l'espace visible d'un S.I.S* comme l'ensemble des objets auquel l'entité référence de ce S.I.S peut accéder. Ce principe est illustré par la figure 1.4. De manière plus formelle :

$$EV(S.I.S(EE_i, t_i)) = \{PrO_i\} \cup \{ShO_j \mid EE_j \in S.I.S(EE_i, t_i)\}$$

FIGURE 1.4 – Espace visible de $SIS(EE_i, t_i)$

Adressage et contrôle au sein de l'espace visible d'un S.I.S. Le S.I.S géré par une entité se voit associer un espace visible. Je me suis alors posé la question de la mise en œuvre de ce dernier, afin de construire des applications tirant pleinement parti du système d'information. Ce problème recouvre deux aspects :

1. *L'adressage* des informations au sein de l'espace visible : de quelle manière les données partagées sont-elles publiées, rendues visibles en direction des entités participant au S.I.S ?

2. Le *contrôle* au sein de l'espace visible : suivant quel(s) mécanisme(s) les différentes entités composant un S.I.S vont-elles interagir, via les données de leur espace visible ?

Nous avons fait le choix d'aborder ces deux problèmes au travers de deux exemples : la vente de proximité et la conférence. Ces deux cas sont présentés dans la section suivante. Notre démarche consiste à dégager les principaux mécanismes à mettre en œuvre entre les calculateurs, afin de permettre une exploitation efficace du principe de l'espace visible. Il devient alors possible de mettre évidence les outils systèmes nécessaires à la réalisation de telles applications.

1.2 Analyse de deux scénarios

La vente de proximité. Considérons le cas d'un marché qui se déroule dans une zone géographique limitée, et dans laquelle se retrouve un ensemble d'acheteurs et de vendeurs potentiels. Un acheteur est mobile, à la recherche d'articles répondant à un ensemble de contraintes (type de l'article, fourchette de prix etc.). Et au gré de ses déplacements, il va rencontrer un ou plusieurs vendeurs susceptibles de répondre à sa demande. C'est la proximité de l'acheteur avec un vendeur qui va provoquer la réalisation de la transaction. En s'appuyant sur les définitions énoncées précédemment, on peut considérer qu'un vendeur gère son propre S.I.S, dans lequel il intègre en permanence les terminaux les plus proches. Ce même S.I.S doit lui permettre d'exprimer des requêtes du type « Je cherche une douzaine d'œufs à moins de trois euros » au sein de son espace visible, et doit réagir lorsque un ou plusieurs vendeurs rencontrés possèdent des articles qui répondent au(x) condition(s) exprimée(s). Nous qualifions ce type d'applications de *base de données de proximité*.

La conférence. Cet exemple a déjà été présenté brièvement dans l'introduction de ce chapitre. Nous considérons une conférence scientifique, dans laquelle chaque participant est porteur d'un terminal mobile équipé de capacités de communication sans fil à courte portée. Dans cet équipement, il a stocké un ensemble de documents lié à son domaine d'activité et en rapport avec la conférence. Cette dernière peut être assimilée à un espace clairement délimité, dans lequel se trouve une forte densité de personnes partageant des domaines d'intérêts communs. Leur proximité est donc une excellente opportunité de récupérer (de glaner) automatiquement des informations pertinentes (qui correspondent aux préoccupations de l'utilisateur). Là encore, chaque terminal gère son propre S.I.S, dans lequel il intègre dynamiquement les entités à proximité. Il va chercher, au sein de l'espace visible, à récupérer spontanément les informations et les documents susceptibles de l'intéresser sur les terminaux voisins. Nous parlons par la suite pour ce type d'applications de *Web de proximité*.

1.3 Synthèse des problèmes de recherche identifiés

Ainsi que l'illustre ces deux exemples, la gestion d'un S.I.S nous pose principalement deux problèmes : (1) la détection de la composition du S.I.S, et (2) l'accès aux informations partagées (*i.e.* l'espace visible) au sein de ce S.I.S. Nous parlons alors de *détection de présence des entités se trouvant à portée de communication*, et de *découverte des informations et interactions au sein du voisinage physique*.

Pour ce qui du premier problème, il est analogue pour les deux exemples d'application. Pour un S.I.S donné (*i.e.* géré par une entité donnée), il est nécessaire de détecter l'ensemble des nœuds mobiles avec lesquels le gestionnaire du S.I.S peut communiquer par voisinage. Il s'agit donc pour un terminal de construire une représentation informatique « fidèle » de son voisinage physique. Bien entendu, cette représentation est d'autant plus complexe à construire et à maintenir que les entités sont mobiles. La figure 1.5 illustre ce problème. Nos résultats sur ce sujet sont présentés dans la section 1.4 de ce chapitre.

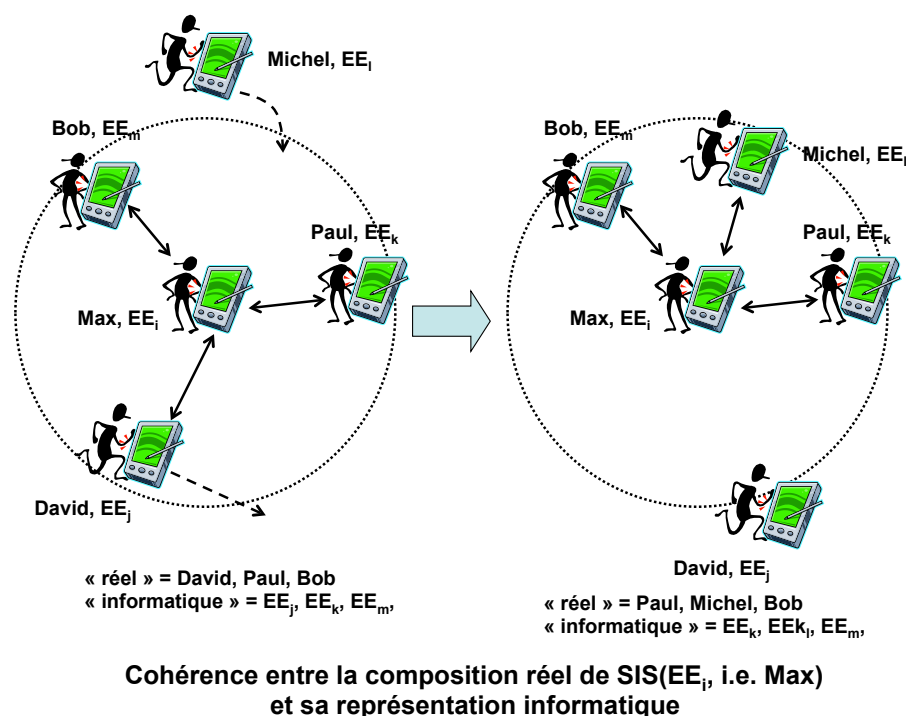


FIGURE 1.5 – Détection de la composition d'un S.I.S

Le deuxième problème recouvre les échanges d'informations au sein de l'espace visible d'un S.I.S. Sur ce plan, les deux exemples précédents diffèrent : les mécanismes requis pour offrir des mécanismes d'adressage et de contrôle de l'espace visible ne sont pas les mêmes dans les deux cas.

Concernant la vente de proximité, il s'agit pour un terminal de propager une interrogation, au sein de l'espace visible de son S.I.S, du type « Quelle entité saura répondre à ma demande ? ». En retour, il lui faut obtenir la liste des informations partagées au sein de son S.I.S, correspondant aux conditions qu'il a exprimées, par exemple un prix maximum sur un article donné. Du point de vue de l'adressage, nous avons fait le choix de traiter l'espace visible comme une base de données distribuées. Nous parlons alors de *base de données de proximité*. Et pour ce qui est de l'aspect contrôle, nous avons proposé une approche s'appuyant sur une interrogation continue de l'espace visible du S.I.S. La solution proposée est présentée dans la section 1.5 du chapitre suivant.

Dans le cadre de la conférence, un terminal doit identifier, au sein de l'espace visible de son S.I.S, les informations partagées les plus pertinentes (par rapport par exemple aux centres d'intérêts de l'utilisateur), et les récupérer spontanément dans son espace de stockage local. Toute la difficulté ici est d'être en mesure :

- de définir des profils utilisateurs associés aux informations que ces mêmes utilisateurs partagent au sein de l'espace visible (c'est à dire les sujets sur lesquels ils désirent trouver de l'information au sein d'un S.I.S),
- et d'exploiter ces profils pour découvrir et récupérer rapidement les documents pertinents partagés par les entités se trouvant (souvent de manière transitoire) dans son S.I.S.

Ceci passe par des mécanismes d'indexation des informations partagées, nous proposons, du point de vue de l'adressage, une technique d'analyse de données issue du *data mining*, avec l'objectif de dégager automatiquement le profil d'un utilisateur à partir des données qu'il stocke localement. Et sur cette base, nous proposons, pour mettre en œuvre les aspects relevant du contrôle de l'espace visible, un protocole de glanage au sein de l'espace visible d'un S.I.S. Notre solution est présentée dans la section 1.6 du chapitre suivant.

La figure 1.6 illustre les différents problèmes liés à la gestion de l'espace visible.

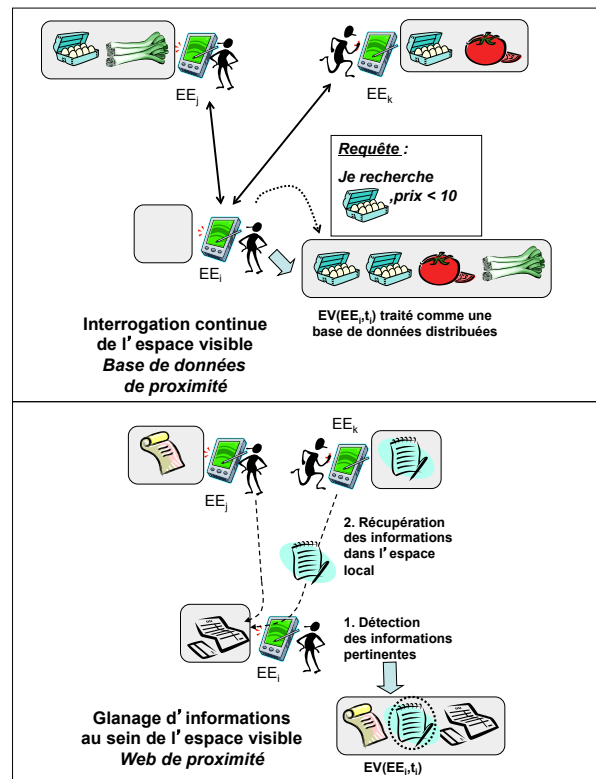


FIGURE 1.6 – Adressage et contrôle de l'espace visible d'un S.I.S

1.4 Découverte et construction du voisinage physique

Il est nécessaire, pour mettre en œuvre un S.I.S donné, de détecter l'ensemble des entités mobiles directement accessibles. Il s'agit donc pour un terminal de construire une représentation informatique de son voisinage physique, dont la composition varie en fonction des mouvements des nœuds voisins. Le principal problème à traiter est de maintenir une cohérence entre le monde informatique et le monde réel. Une approche classique consiste à ce qu'un terminal annonce sa présence en diffusant périodiquement un message contenant son identifiant, par exemple son adresse physique. Nous parlons indifféremment de *message de présence* ou de *message d'annonce*. Lorsqu'un autre terminal reçoit un tel message, il en conclut qu'il est à même d'initier une communication par voisinage avec le terminal distant. Cette communication est maintenue tant que les messages de présence continuent d'être reçus. Dès lors que ces messages d'annonces ne sont plus captés à l'expiration d'un délai fixé au niveau du terminal, on considère que la relation de voisinage est rompue.

Le problème est que les interactions sans fil de proximité sur lesquelles nous nous appuyons sont très sensibles à l'environnement (présence d'obstacles, mobilité des utilisateurs etc.). Dans ces conditions, il est possible qu'un terminal EE_i reçoive l'annonce de EE_j , mais que l'inverse ne soit pas vrai. Comment alors représenter de manière fiable le voisinage physique d'un S.I.S? Nous présentons une solution à ce problème dans la section 1.4.1.

De plus, la périodicité d'envoi des messages de présence est commune à tous les terminaux, et fixée de manière empirique, en fonction de la réactivité attendue de l'application. Plus ces messages sont envoyés fréquemment, plus le mécanisme de détection devient réactif. Bien entendu, il faut tenir compte du fait que des annonces trop fréquentes aboutissent à des consommations énergétique et de bande passante excessives. À l'inverse, une fréquence insuffisante dans l'envoi des annonces peut amener une forte dégradation au niveau de la cohérence de la représentation informatique du

voisinage physique du S.I.S. Comment alors cadencer automatiquement la fréquence des annonces au sein d'un SIS? Notre réponse à ce problème est présentée dans la section 1.4.2.

1.4.1 Définition d'une nouvelle relation de voisinage

Une communication par voisinage entre deux terminaux n'est pas obligatoirement symétrique : une entité peut tout à fait recevoir des données d'une autre (par exemple des messages de présence), sans pour autant être en mesure de lui répondre. On voit donc que la traditionnelle définition de voisinage reste trop floue pour caractériser tous les types de situations susceptibles d'être rencontrés au sein d'un S.I.S. Nous proposons donc de considérer trois types de voisinage : (1) **le voisinage unilatéral**, (2) **le voisinage réciproque**, et (3) **le voisinage unilatéral strict**.

Considérons EE_i et EE_j deux terminaux. EE_i est dit *voisin unilatéral* de EE_j si et seulement si EE_j est à portée de EE_i . Autrement dit, EE_j reçoit les messages envoyés par EE_i . Et par extension, deux terminaux EE_i et EE_j sont dits *voisins* ou encore *voisins réciproques* lorsqu'ils sont mutuellement voisins unilatéraux. Nous retrouvons alors la notion de communication par voisinage, telle que nous l'avons définie initialement.

EE_i est *voisin unilatéral strict* de EE_j lorsque EE_i est voisin unilatéral de EE_j sans en être un voisin réciproque. Ce cas de figure est illustré par la figure 1.7. Nous nous appuyons sur ces trois types de voisinage, pour définir des mécanismes de représentation de la composition d'un S.I.S.

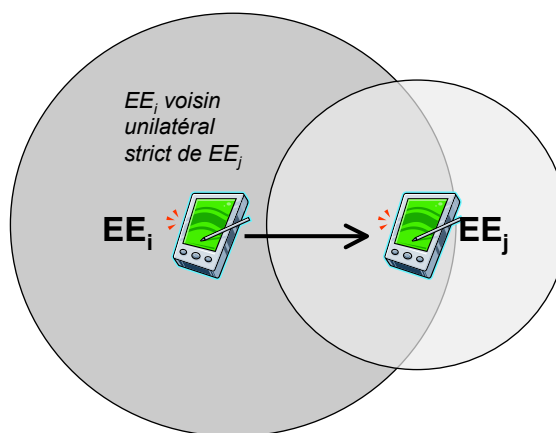


FIGURE 1.7 – Un exemple de voisinage unilatéral strict

1.4.2 Calcul automatique de la fréquence d'annonce de présence

La fréquence d'annonce « idéale » d'un terminal dépend directement de sa portée de communication, et de son profil de mobilité. Nous cherchons à lever cette limite, en calculant une fréquence d'émission permettant de construire une représentation fiable du voisinage physique pour un groupe de terminaux participant à un S.I.S. Pour cela, nous mettons tout d'abord en évidence la *condition de détection* d'un terminal par rapport à un autre. Cette condition est ensuite étendue à l'ensemble du voisinage de ce même terminal. A partir de là, nous dégageons un paramètre appelé *ratio de détection*, qui peut être fixé pour chaque terminal, indépendamment de sa portée de communication. Enfin, nous sommes en mesure de fixer pour chaque terminal la fréquence d'émission de ses messages d'annonce en utilisant ces différents paramètres. Les calculs justifiant la pertinence de ce *ratio de détection* ne sont pas décrits en totalité dans ce document, ils sont présentés de manière détaillée dans [85, 88].

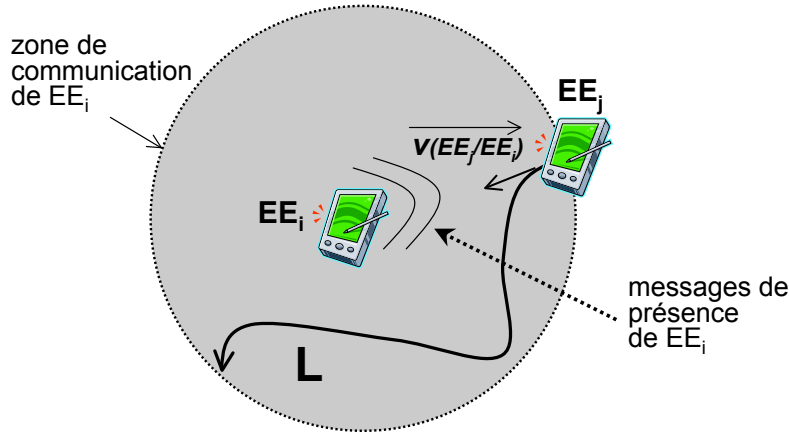


FIGURE 1.8 – Interprétation cinématique

Condition de détection des terminaux au sein d'un S.I.S. Considérons deux entités EE_i et EE_j . EE_i émet des messages de présence à une fréquence f_{EE_i} . Nous ne sommes pas en mesure de caractériser le temps de présence de EE_j dans $ZC(EE_i)$, la zone de communication de EE_i . Par contre, il est possible de définir la distance L que parcourt EE_j dans $ZC(EE_i)$. Une telle distance dépend de la vitesse relative de EE_j par rapport à EE_i . Cette dernière est notée $\overrightarrow{v_{EE_j/EE_i}}$. Nous notons également v_{EE_j/EE_i} l'intensité de cette vitesse relative.

Ces principes sont résumés dans la figure 1.8. Partant de cette représentation cinématique du voisinage de EE_i , nous avons dégagé une condition de détection de EE_i par EE_j , EE_i étant alors voisin unilatéral de EE_j . Une telle condition se traduit par la distance minimale que doit parcourir EE_j dans la zone de communication $ZC(EE_i)$, afin de recevoir au moins un message de présence de EE_i , envoyé avec une fréquence f_{EE_i} . Autrement dit, EE_j détecte EE_i s'il reste dans sa zone de communication pendant une durée supérieure à $1/f_{EE_i}$. Nous avons montré que EE_j détecte EE_i si $L > \frac{v_{EE_j/EE_i}}{f_{EE_i}}$.

Cette condition ne s'exprime que pour le terminal EE_j . Il nous est nécessaire de l'étendre à tous les terminaux se trouvant à portée de EE_i , *i.e.* formant le voisinage du S.I.S de EE_i . Nous utilisons pour cela la **vitesse relative du voisinage** :

On appelle *vitesse relative du voisinage* le maximum des vitesses relatives des nœuds dont EE_i est voisin unilatéral. On la note $V_{EE_i} = \max(v_{EE_k/EE_i}, \forall EE_k, EE_i \text{ voisin unilatéral de } EE_k)$.

Ratio de détection et fréquence d'annonce. Nous considérons ensuite un ratio noté R , dont la valeur est fixée quel que soit le type de terminal. L'interprétation de cette valeur est : *tout terminal parcourant R % de la portée d'une entité quelconque détecte nécessairement cette entité*. Avec cette valeur R fixée, un terminal est en mesure de calculer dynamiquement la fréquence d'émission de ses messages de présence [88] :

$$R = \frac{V_{EE_i}}{f_{EE_i} \cdot d_{EE_i}}$$

d'où $f_{EE_i} = \frac{V_{EE_i}}{R \cdot d_{EE_i}}$

Calcul effectif d'une fréquence d'annonce. Dans l'équation précédente, la valeur de R est connue, ainsi que la portée d_{EE_i} du terminal. Enfin, la vitesse relative maximale de EE_i peut être obtenue, à condition que chaque entité du voisinage calcule et transmette sa vitesse relative.

Un tel mécanisme peut être réalisé si une entité dispose en permanence d'une mesure de sa vitesse, et est capable de la communiquer à ses différents voisins. Dans le cas d'un bus par exemple, cette hypothèse est réaliste : le véhicule est équipé d'un tachymètre, qui peut être couplé au terminal. Dans le cas moins favorable où l'entité mobile n'a pas accès à cette information, on peut envisager d'associer une valeur fixe au terminal. Il est possible de dégager des *profils de mobilité*, pour lesquels on laisse la possibilité à l'utilisateur de spécifier son mode de déplacement : transport en commun urbain, véhicule routier, piéton etc. Chacun de ces modes se voit associer automatiquement une vitesse moyenne. Une telle approche présente l'avantage de pouvoir être mise en place dans tous les cas de figure.

Une autre possibilité consiste à s'appuyer, en plus de la vitesse, sur une mesure de la distance entre deux terminaux. On peut pour cela partir d'un mécanisme de géolocalisation tel que le GPS. Il est également possible d'utiliser des propriétés propres aux ondes radioélectriques (comme la puissance du signal reçu), ou bien encore la différence des vitesses de propagation entre des ondes infrarouges et des ondes radioélectriques [92]. Dans ce cas, les terminaux participant à un S.I.S disposent d'informations cinématiques plus étendues que précédemment.

1.4.3 Bilan

Cette partie de nos travaux concerne l'étude d'un mécanisme système permettant de construire et maintenir au niveau de chaque entité mobile une vue cohérente du S.I.S auquel elle participe. La principale difficulté est de déterminer est de maintenir un taux de détection acceptable des voisins au sein du S.I.S, tout en minimisant la bande passante consommée par les échanges de messages d'annonce. Après avoir mis en évidence les problèmes engendrés par une vision symétrique des communications par voisinage, nous avons proposé une approche reposant sur un ratio de détection, R , fixé par le concepteur du système, et qui peut être interprété de la manière suivante : un terminal parcourant $R\%$ de la portée d'un nœud voisin détecte nécessairement ce nœud. A partir du moment où chaque terminal connaît sa portée et sa vitesse, et qu'il diffuse ces deux informations à l'ensemble de ses voisins, il en mesure de calculer une fréquence d'annonce qui lui est propre.

Ce principe nous a permis de définir et évaluer des protocoles de découverte, en cherchant le meilleur compromis entre la qualité de la détection et son coût sur le réseau, pour différentes valeurs de ratio R . Les résultats détaillés sont présentés dans [85, 11]. A titre d'exemple, on a pu constater qu'une valeur de ratio R comprise entre 0,3 et 0,5 amène un compromis acceptable, la qualité de détection atteignant 80 %, avec un volume de messages échangés qui n'explose pas lorsque le nombre de terminaux voisins augmente.

Il me semble intéressant de soulever un aspect concernant l'évaluation de l'efficacité des protocoles de découverte exploitant le principe d'un ratio de détection. Nos résultats ont été obtenus uniquement par le biais de simulations. L'étude présentée ici date de 2004. Nous avons à ce moment là anticipé la capacité d'un ordinateur mobile à communiquer ses caractéristiques cinématique (comme sa vitesse) à ces voisins, condition nécessaire pour que notre approche fonctionne. Les limites technologiques des PDAs ne nous permettaient pas d'envisager une expérimentation réelle de notre solution. Force est de constater que les capteurs embarqués dans les *Smartphones* les plus récents (notamment l'accéléromètre et le gyroscope) doivent maintenant permettre une implémentation efficace de la gestion du voisinage physique d'un S.I.S.

1.5 Base de données de proximité

Dans le cadre d'un S.I.S donné, un ensemble de terminaux est susceptible de partager et d'échanger des données. La composition de cette collection de données, appelée l'espace visible du S.I.S, évolue dans le temps. La mobilité modifie les relations de voisinage entre les équipements, et par voie de conséquence, l'ensemble des données accessibles. De plus, rien n'empêche une entité de modifier ces données locales qu'elle partage au sein d'un S.I.S. Ceci a pour effet de faire évoluer la composition de l'espace visible correspondant.

Nous proposons dans un premier temps de traiter l'espace visible d'un S.I.S comme une base de données, et d'offrir un ensemble de mécanismes permettant à une application telle que la vente de proximité d'interroger cette base. On parle alors de **base de données de proximité**.

Considérant les capacités de stockage et de traitement sans cesse croissantes des terminaux mobiles, la gestion des informations d'un S.I.S peut être envisagée à différents niveaux, du système de fichiers à la base de données. C'est ce dernier domaine que nous avons retenu dans le cadre de notre approche. Les bases de données offrent en effet des mécanismes extrêmement efficaces, dès qu'il s'agit de stockage et d'interrogation de données. Plus précisément, notre choix s'est porté sur les bases de données relationnelles (SGBDR), et sur SQL, le langage d'interrogation qui leur est associé.

Afin de s'adapter au caractère très dynamique de l'espace visible, une application souhaitant disposer à tout moment d'informations à jour doit augmenter la fréquence de ses accès au système d'information. Une telle approche présente un inconvénient majeur : la base de données peut être interrogée, alors que l'ensemble des informations disponibles au sein de l'espace visible du S.I.S n'a pas évolué depuis la précédente requête. Ceci conduit à gaspiller les ressources locales des terminaux, et les capacités offertes par les canaux sans fil. Nous devons proposer aux applications des mécanismes d'interrogation efficace de l'espace visible, prenant en compte l'évolution des données accessibles au sein du S.I.S associé. Nous envisageons pour cela un système d'**interrogation continue de l'espace visible**.

Objectifs des systèmes à requêtes continues. Les requêtes continues ont été introduites pour palier aux insuffisances des systèmes d'information classiques, dès lors qu'il s'agit d'interroger des données dynamiques. Définie à l'origine dans le cadre du système *Tapestry* [69], une requête continue peut être vue comme une interrogation qui, une fois soumise au système d'information, est exécutée de manière continue. L'objectif ici est de permettre le filtrage d'un flux de messages électroniques. Ce système se limite à un modèle de base de données *append-only*, c'est à dire n'autorisant que l'insertion de nouvelles données (*i.e.* messages) dans la base. Cette approche fortement centralisée a largement évolué dans le cadre des réseaux de capteurs [17]. En effet, ces derniers constituent l'un des champs d'applications les plus évidents des requêtes continues. Il s'agit par exemple de collecter un ensemble de données d'environnement (températures, pressions etc.) soumis à des variations extrêmement fréquentes.

Les requêtes soumises à une base de données de capteurs portent très souvent sur un ensemble de valeurs mesurées au niveau de quelques capteurs. Il s'agit par exemple de récupérer périodiquement des températures dans un bâtiment. Et de plus en plus, les capteurs utilisés sont dotés de capacités de stockage, de traitement et de communication. Il devient donc envisageable de distribuer le processus d'interrogation. Le système Cougar [17] propose une approche pour laquelle seuls les capteurs concernés par une requête transmettent leurs données au serveur.

S.I.S et requêtes continues. Un S.I.S est un système totalement distribué, dans lequel chaque équipement mobile doit remplir une double fonction : (1) d'une part fournir des données aux autres nœuds du S.I.S, et (2) d'autre part participer à l'évaluation de la requête continue. A l'opposé, le moteur de requêtes continues de *Tapestry* est totalement centralisé. Dans le cadre des bases de données de capteurs, l'architecture est distribuée : les capteurs fournissent les données, et l'évaluation finale de la requête est localisée sur un serveur. La plupart des systèmes impose une connectivité globale entre tous les nœuds. Ceci passe par l'utilisation d'une infrastructure sans fil. Une telle approche ne peut pas être appliquée au contexte particulier des communications par voisinage.

De plus, l'interrogation continue de l'espace visible nécessite des mécanismes plus souples que ceux présentés précédemment. Un utilisateur doit avoir la possibilité d'insérer, de supprimer ou de mettre à jour des données qu'il souhaite partager dans le cadre d'un S.I.S. De tels mécanismes ne sont pas traités par le système *Tapestry*, le gestionnaire de requêtes continues se limitant à des bases *append-only*. De même, les relations virtuelles introduites par le système Cougar ont la même limitation, supportant uniquement l'insertion de nouvelles données.

On le voit, l'interrogation continue de l'espace visible d'un S.I.S nécessite de prendre en compte de nouvelles contraintes. Ainsi, une requête continue initiée par un terminal doit fournir à chaque instant une vue de l'ensemble des données de l'espace visible du S.I.S, qui vérifie les conditions associées à cette requête. Nous appelons cet ensemble de données le *Continuous Result Set (CRS)*. Nous présentons dans la suite de cette section les principaux mécanismes système étudiés.

1.5.1 Présentation des mécanismes proposés

Évolution de la composition du S.I.S. La composition du S.I.S d'une entité donnée exécutant une requête continue de l'espace visible associé a une influence directe sur le résultat de cette interrogation. La sortie d'un nœud du S.I.S provoque le retrait des informations qu'il partage au sein du S.I.S. Et si de plus, certaines de ces informations vérifient les conditions liées à la requête continue en cours, elles doivent donc être retirées du *CRS* associé. Considérons deux entités EE_i et EE_j . Dès que EE_j n'est plus en mesure de communiquer par voisinage avec EE_i , la requête continue initiée par EE_i ne doit plus prendre en compte les données partagées par EE_j . De manière analogue, si EE_k s'insère dans le S.I.S géré par EE_i , la requête continue initiée par EE_i doit intégrer les données partagées par EE_k . Ce problème est traité en s'appuyant sur la représentation informatique du S.I.S, dont la mise en œuvre est abordée dans la section 1.4. Dans le cas d'une sortie, les informations provenant de l'équipement en train de s'éloigner doivent être retirées du *CRS* par l'entité à l'origine de la requête, ainsi que l'illustre la figure 1.9. Et c'est cette même entité qui doit informer un nouvel entrant dans le S.I.S des conditions de la requête en cours.

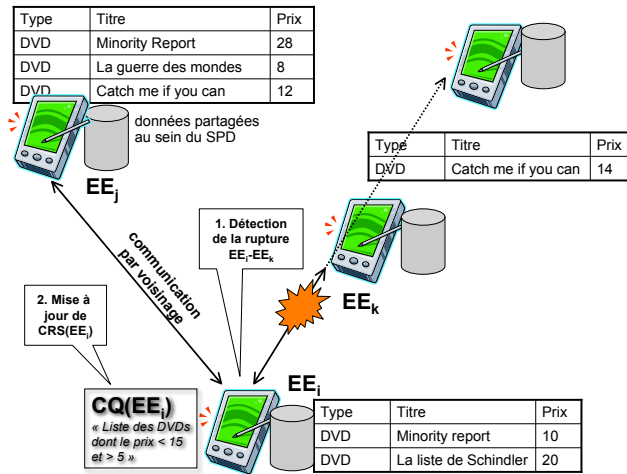


FIGURE 1.9 – Évolution physique du S.I.S

Mise à jour au sein du S.I.S. La valeur d'un *CRS* est amenée à évoluer en fonction de l'entrée et la sortie des nœuds au sein du S.I.S associé. Mais un autre facteur doit être pris en compte dans l'évaluation d'une requête continue : les informations partagées au sein d'un S.I.S peuvent être modifiées par l'entité qui les possède. Considérons l'exemple d'une vente aux enchères : une entité EE_j publie une liste d'objets qu'un utilisateur souhaite vendre. Un équipement voisin EE_i déclenche la requête continue CQ_{EE_i} suivante : « Je recherche un DVD dont le prix est compris entre 5 et 15 euros ». Nous notons $CRS(CQ_{EE_i})$ l'ensemble des données de l'espace visible correspondant à cette requête. Afin de maintenir ce *CRS* à jour, il est nécessaire de détecter trois types d'événements :

- La suppression d'une information. Si cette information correspond aux conditions associées à CQ_{EE_i} , elle doit être supprimée de $CRS(CQ_{EE_i})$.

- L'ajout d'une nouvelle information. Si cette dernière est valide du point de vue de la requête, elle doit être insérée dans le *CRS*.
- La modification d'une information. Imaginons que l'utilisateur de EE_j décide de diviser l'ensemble des prix qu'il publie par 2. Selon les cas, cela peut entraîner l'insertion ou le retrait de l'information de $CRS(CQ_{EE_i})$. Dans notre exemple, le DVD *La guerre des mondes* doit sortir du *CRS*, alors que *Minority report* doit y être inséré.

Le problème est de savoir comment ces modifications doivent être traitées pour mettre à jour un *CRS*. Dans le cas de l'évolution physique d'un S.I.S évoqué précédemment, les mises à jour du *CRS* sont effectuées par l'équipement à l'origine de la requête. Ici, c'est aux calculateurs stockant les informations modifiées, et non pas à l'entité interrogatrice, de réagir à ces trois types d'événements. En effet, dès qu'un nœud impliqué dans une interrogation continue modifie les informations qu'il partage au sein de l'espace visible, il doit s'assurer que les opérations effectuées n'ont pas de répercussion sur le *CRS* associé à la requête en cours d'exécution. Pour cela, il dispose des conditions caractérisant l'interrogation, qui lui ont été fournies par l'équipement initiateur de la requête. Quand cela est nécessaire, lors d'une mise à jour locale des informations qu'il partage au sein du S.I.S, un nœud est en mesure de reporter la nature de la mise à jour à apporter au *CRS* : *suppression* d'une information, *insertion* d'une nouvelle information, ou bien *modification* d'une information. Ce mécanisme est présenté dans la figure 1.10.

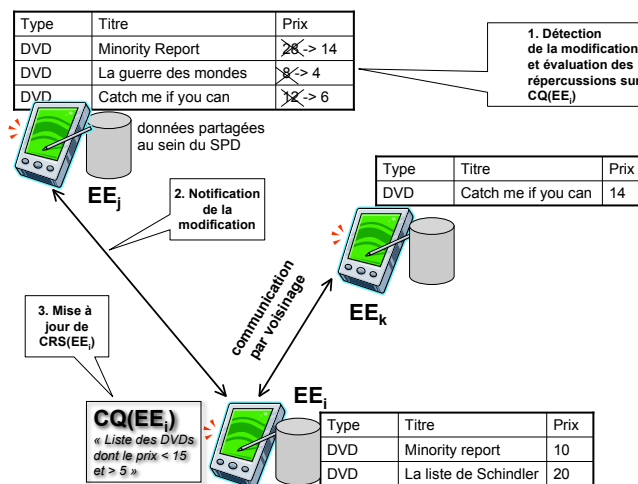


FIGURE 1.10 – Mise à jour au sein de l'espace visible du S.I.S

Mécanismes liés à la durée. Les requêtes continues amènent des problèmes liés à la durée de l'exécution. En effet, la construction d'un *CRS* évolue en fonction des données accessibles au sein de l'espace visible du S.I.S. Un langage de requête tel que SQL a été conçu pour définir et adresser des ensembles de données statiques. Les fonctions fournies par ce langage ne sont pas adaptées à la gestion d'ensembles de données soumis à des variations. Ainsi, SQL permet à un utilisateur d'invoquer des fonctions d'agrégations (comme `max`, `sum`, `count` etc.) dans ses requêtes. Ces fonctions calculent une valeur unique à partir d'un ensemble de données statiques.

Afin de pouvoir gérer des ensembles de données changeantes, nous avons proposé d'associer des sémantiques dynamiques à ces fonctions d'agrégations. La valeur retournée par ces fonctions doit refléter l'état courant du *Continuous Result Set*. Pour cela, cette valeur doit être réévaluée chaque fois que le *CRS* est mis à jour. Notons que dans la plupart des cas, la valeur à retourner peut être recalculée, sans que la totalité du *CRS* soit nécessairement balayée. Par exemple, la valeur donnée par l'appel `count(*)` (qui correspond au nombre d'enregistrements retourné par la requête `select`) peut être facilement calculée : elle doit être incrémentée chaque fois qu'un champ est inséré dans le *CRS*, et décrétementée pour chaque champ supprimé.

1.5.2 Mise en œuvre de l'architecture PERSEND

Les principes d'une base de données de proximité ont été mis en œuvre et expérimentés au travers de l'architecture PERSEND². Nous en présentons ici les principaux composants.

Expression d'une requête continue sur l'espace visible d'un S.I.S. La syntaxe du langage SQL n'a pas été conçue pour permettre l'expression de requêtes continues. Nous avons donc ajouté le mot clé `continuous`, de façon à distinguer les interrogations continues des demandes classiques. Dans une requête SQL, les données sources sont normalement identifiées en désignant explicitement les tables et les bases visées. Ceci n'est pas possible dans notre modèle. Un S.I.S donné évolue dans le temps, et les bases et les tables voisines ne sont pas connues *a priori*. Un terminal doit savoir quelle(s) table(s) de sa base locale il doit prendre en considération, et être capable de distribuer une requête vers les entités appartenant à son S.I.S. Dans ce but, nous avons introduit le mot clé `vicinity`. Il est placé au début de la requête, juste après le mot clé `continuous`. Il permet d'indiquer à notre système que la requête, continue ou non, doit être distribuée à tous les nœuds du voisinage physique. C'est ce que nous appelons une *requête continue de proximité*. La requête suivante permet par exemple l'*affichage continu des DVD en vente au sein de l'espace visible, dont le prix est compris entre 5 et 15 euros* :

```
continuous vicinity select titre, prix
from dvd_a_vendre
where prix is between 5 and 15 ;
```

Une fois l'interrogation lancée, le nœud origine doit propager cette requête vers l'ensemble des terminaux participant à son S.I.S.

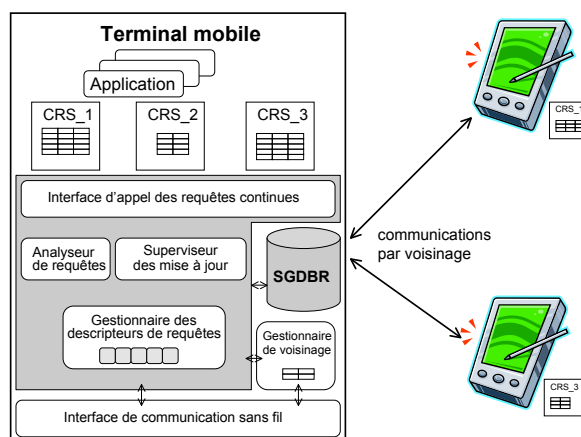


FIGURE 1.11 – Architecture du système d'interrogation PERSEND

Vue d'ensemble de l'architecture PERSEND. La figure 1.11 présente l'architecture globale de PERSEND. Cette dernière s'appuie sur une base de données relationnelles (SGBDR), ainsi que sur le gestionnaire de voisinage qui fournit la représentation informatique du S.I.S. PERSEND est organisé autour de quatre composants principaux : l'interface d'appel des requêtes continues, l'analyseur de requêtes, le superviseur des mises à jour, et le gestionnaire des descripteurs de requêtes. Le fonctionnement de ces composants est détaillé dans [83].

2. PERSEND : *PERsistent SEnsing for Neighbouring Data*

1.5.3 Bilan

Nous avons donc proposé dans un premier temps d'exploiter l'espace visible d'un S.I.S comme une base de données, appelée alors *base de données de proximité*. Ce type d'approche est tout à fait adapté pour « balayer » en permanence les informations partagées au sein d'un S.I.S répondant à des conditions précises. Dans le cas d'un système d'enchères, cela revient à chercher un ou des objets à vendre (par exemple un DVD) répondant à des critères plus ou moins larges (une fourchette de prix, le titre etc.). L'application embarquée sur le terminal et exploitant la base de données de proximité réagit en fonction du résultat de la requête soumise.

Notre approche s'appuie sur une interrogation continue de l'espace visible de l'espace visible, prenant en compte (1) l'évolution de la composition du système d'information (entrées et sorties des nœuds d'un S.I.S), (2) ainsi que les mises à jour effectuées localement par chaque entité sur les informations qu'elle partage au sein d'un ou plusieurs S.I.S. Partant d'une extension du langage SQL permettant d'exprimer des requêtes continues, nous avons proposé différents mécanismes système permettant de tirer parti des propriétés spécifiques des S.I.S (gestion du voisinage, gestion locale des informations par chaque nœuds etc.). Toutes nos propositions sont regroupées au sein d'une architecture appelée PERSEND. Le lecteur souhaitant plus de détails concernant son fonctionnement pourra se reporter à [80, 83, 77].

1.6 Glanage opportuniste au sein du Web de proximité

Les bases de données de proximité ne sont pas appropriées dès qu'il s'agit de profiter de la proximité entre des personnes pour **échanger spontanément** des documents pertinents. Dans l'exemple de la conférence évoqué précédemment, chaque participant, porteur d'un assistant numérique personnel, peut avoir stocké un ensemble de documents plus ou moins volumineux : sa présentation, des publications etc., qu'il aura par exemple extrait de son site Web personnel. Considérant la densité de personnes partageant de nombreux centres d'intérêts communs, le S.I.S d'un utilisateur constitue une excellente opportunité de « glaner » spontanément des informations intéressantes, au fur et à mesure des entrées et des sorties au sein de ce même S.I.S. Nous parlons alors de *Web de proximité*. Pour construire une telle application, il n'est pas envisageable de récupérer à chaque rencontre physique entre deux nœuds la totalité des informations partagées, et ce pour trois raisons :

1. Les temps de communication sont souvent limités entre des nœuds mobiles communiquant par voisinage. Le risque est donc de transmettre des informations peu pertinentes au moment de la création de la communication par voisinage, et d'empêcher des échanges plus fructueux par la suite, alors que la connexion est rompue.
2. La bande passante offerte par les communications sans fil est restreinte et soumise à de fortes fluctuations. Elle doit donc être utilisée efficacement pour échanger avant tout des informations utiles pour les utilisateurs.
3. Les capacités de stockage local des terminaux sont limitées. Il faut donc éviter à l'utilisateur, dans le cadre d'une application de « glanage », de stocker automatiquement des informations ne cadrant pas avec ces centres d'intérêts.

Les questions soulevées sont les suivantes :

- *Comment définir les centres d'intérêts d'un utilisateur ?* Ces derniers doivent être découverts avant de déclencher des échanges pertinents au sein de l'espace visible. Dans ce but, nous proposons de construire et de maintenir automatiquement *un profil de l'utilisateur*. Ce dernier doit être établi automatiquement, sans intervention explicite de l'utilisateur, de manière à ne pas freiner la spontanéité des échanges. Pour cela, nous nous appuyons sur l'ensemble des informations que l'utilisateur conserve dans son terminal. Il semble raisonnable de considérer que les documents stockés reflètent les centres d'intérêts d'une personne.

- *Comment organiser les informations partagées au niveau du terminal ?* Toujours dans le but de maintenir la spontanéité et l'efficacité des échanges, chaque terminal doit organiser automatiquement les informations qu'il souhaite partager au sein d'un S.I.S. Afin d'accélérer les communications entre nœuds voisins, nous proposons un système d'indexation thématique des documents locaux. Nous traitons là le problème de l'adressage au sein de l'espace visible.
- *Comment découvrir les informations pertinentes à échanger ?* Dans le cadre d'une communication par voisinage entre deux terminaux mobiles, un protocole doit permettre de sélectionner les documents à récupérer, sur la base des profils respectifs des deux utilisateurs distants.

Tout comme dans le modèle du Web, le principe du *Web de proximité* permet l'échange de documents, dont le contenu est caractérisé par les mots clés. Il existe dans le domaine de l'*information retrieval* des algorithmes permettant d'évaluer la proximité entre des mots clés, de façon à évaluer de manière souple la pertinence d'une requête par rapport au contenu d'un document. De tels mécanismes ne peuvent malheureusement pas être utilisés dans le contexte des S.I.S. Gourmands en ressource, ils ne sont pas adaptés aux capacités limitées des systèmes mobiles et des communications sans fil. Afin de contourner ce problème, les mots clés sont choisis dans une *ontologie* commune à tous les utilisateurs. Par exemple, dans le cadre restreint d'une conférence ou d'un salon de l'emploi, il m'a semblé envisageable pour tous les participants d'adopter une classification normalisée de sujets.

Le *Web de proximité* va permettre à un terminal de récupérer automatiquement des documents au sein de l'espace visible de son S.I.S. Et ce même terminal peut également participer au(x) S.I.S(s) des autres nœuds avec lesquels ils communiquent par voisinage. Donc si on observe deux terminaux dans le cadre d'une communication par voisinage, on constate qu'il joue un rôle symétrique : ils sont à la fois fournisseur et récupérateur de données. La figure 1.12 illustre les interactions symétriques instaurées dans le cadre du *Web de proximité*.

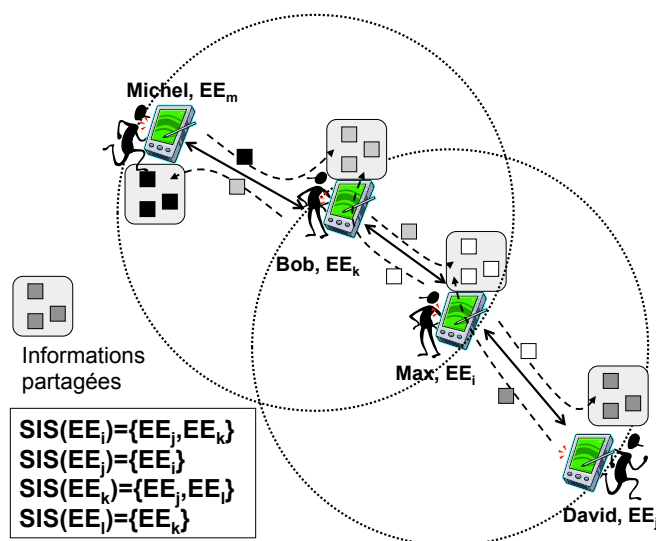


FIGURE 1.12 – Interactions entre terminaux dans le cadre du *Web de proximité*

Nos travaux sur le Web de proximité portent sur deux aspects, l'organisation des informations au sein d'un terminal et les mécanismes de découverte d'informations entre les terminaux voisins, présentés de manière synthétique dans la suite de cette section.

1.6.1 Gestion de la base d'information du terminal

Deux problèmes doivent être considérés : (1) la définition du profil d'un utilisateur, sachant que les informations qu'il partage au sein de son S.I.S reflètent ses domaines d'intérêts, et (2)

l'indexation thématique de ses documents. Pour aborder ces deux problèmes, nous avons adapté des techniques d'analyse de données issues du *data mining*. Ces techniques sont généralement utilisées pour découvrir des dépendances cachées entre des objets appartenant à des bases de données de grande taille. Reprenons l'exemple de la conférence traité dans le cadre du *Web de proximité*. Nous considérons une base de données regroupant des publications, caractérisées par le(s) auteur(s) référencé(s) dans les bibliographies. Dans ce cadre, le *data mining* peut être utilisé pour découvrir quelles sont les auteurs cités ensemble, et quelles sont ceux qui n'apparaissent jamais de manière groupée au sein d'une même bibliographie.

Le *data mining* permet principalement de découvrir des relations structurelles cachées entre *itemset* d'une base de données. Nous avons donc proposé d'appliquer cette technique pour extraire un ensemble de mots clés fréquents caractérisant les documents chargés. Chaque terminal susceptible de participer à un S.I.S contient donc un ensemble de pages Web partagées. Ces documents stockés correspondent aux enregistrements d'une base de données embarquée par chaque nœud mobile, et les mots clés les caractérisant sont vus comme des *itemset* associés à ces enregistrements.

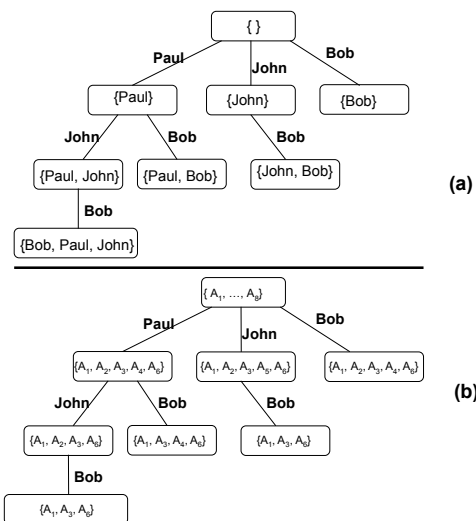


FIGURE 1.13 – Analyse des données via le *data mining*

Afin d'extraire un profil utilisateur de cette masse de données, l'algorithme de *data mining* est adapté pour fournir une structure arborescente des mots clés fréquents (*cf.* figure 1.13(a)). Seuls les *itemset* (*i.e.* les groupes de mots clés) fréquents sont indexés. Ainsi, la structure construite par l'algorithme fournit une représentation cohérente des principaux domaines d'intérêts de l'utilisateur. Cet arbre peut également être utilisé pour classer les documents caractérisés par les mots clés fréquents, ainsi que l'illustre la figure 1.13(b) : un document est indexé au niveau d'un nœud de l'arbre si l'ensemble des mots clés obtenu par le parcours de l'arbre jusqu'à ce nœud correspond à l'ensemble des mots clés caractérisant ce document. Ainsi, on dispose d'une représentation locale des informations qui peut être exploitée pour la construction du profil de l'utilisateur (*i.e.* ses centres d'intérêts), et également dans un stockage thématique des informations qu'il partage au sein de son S.I.S.

1.6.2 Mécanismes de découvertes d'informations entre deux terminaux d'un S.I.S

Dans l'exemple de la conférence, nous avons fait l'hypothèse que chaque utilisateur stocke sous forme de pages Web un ensemble de documents caractérisés par les auteurs référencés dans les bibliographies. A partir de ces informations, nous sommes en mesure d'extraire un profil de l'utilisateur (les auteurs référencés permettent de dégager finement ses domaines de recherche) et un classement thématique de ces informations s'appuyant sur ce profil.

Au sein d'un S.I.S, une communication par voisinage entre deux terminaux doit permettre d'échanger rapidement des documents intéressants. Décrit de manière simple, il s'agit de glaner en priorité des documents distants dans lesquels sont référencés des auteurs auxquels l'utilisateur s'intéresse fréquemment. Ce caractère fréquent va être évalué à l'aide des documents auxquels il possède localement dans son terminal. Nous avons proposé la mise en œuvre d'un mécanisme de découverte itératif, calculant les **intersections** successives entre les deux profils arborescents distants, autrement dit les documents pertinents au regard des profils des utilisateurs. Les documents associés sont alors automatiquement téléchargés. Le protocole recherche progressivement les intersections entre les profils utilisateurs, afin de prendre en compte le caractère volatile des communications par voisinage au sein du S.I.S.

Considérons une rencontre entre deux terminaux A et B . Le protocole démarre en sélectionnant un initiateur pour les échanges futurs. Supposons que A est élu dans le cas présent, il envoie alors à B l'ensemble des mots clés décrivant le premier niveau de son profil. B répond à A en envoyant les nœuds suivants de son arborescence, à partir des branches (*i.e.* mots clés) correspondants à l'intersection qu'il vient de calculer. A son tour, A calcule une intersection sur la base de la réponse envoyée par B , et ce processus est répété jusqu'à ce que l'intersection calculée soit vide. Les documents indexés par les intersections calculées à chaque itération sont transmis d'une entité vers l'autre (en prenant soin de ne transmettre qu'une seule fois ceux qui sont multi-indexés). Ces mécanismes sont illustrés par la figure 1.14. Plus de détails concernant ce protocole de découverte peuvent être trouvés dans [76].

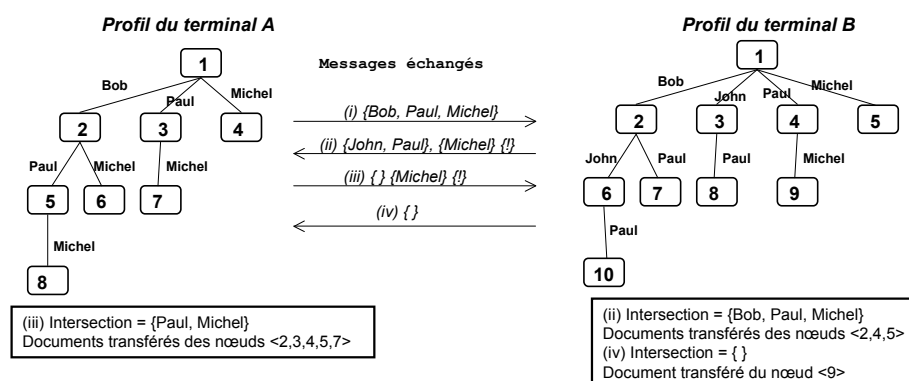


FIGURE 1.14 – Mécanisme de découverte des intérêts communs

1.6.3 Bilan

Tout comme dans le modèle du Web, le principe du *Web de proximité* permet d'échanger des documents, dont le contenu est caractérisé par des mots clés. Appliqué dans le contexte des S.I.S, cette approche permet à un terminal mobile de récupérer (*i.e.* charger) automatiquement des documents au sein de son espace visible (*i.e.* se trouvant sur les nœuds distants participant au S.I.S du terminal). Cette *découverte* et ce *glanage* automatiques induisent une gestion et une organisation spécifique de l'espace d'information locale du terminal [79].

Nous avons proposé un ensemble de mécanismes s'appuyant sur les mots clés associés aux pages Web stockées localement par l'utilisateur,

- qui extrait automatiquement les centres d'intérêts de l'utilisateur, on suppose pour cela que l'utilisateur stocke ses documents en fonction des domaines qui l'intéressent,
- et qui organise thématiquement ces documents, afin d'en accélérer le transfert au sein d'un S.I.S.

Enfin, nous avons défini un protocole permettant à des entités voisines d'échanger les centres d'intérêt respectifs de chaque utilisateur (appelés les *profils*), de les analyser, pour finalement échanger rapidement des informations pertinentes extraites des espaces d'information distants.

1.7 Conclusion

Ce chapitre présente mes premiers travaux de recherche portant sur les Systèmes d'Information Spontané (S.I.S). Initiés en 1998 à mon arrivée au sein de l'équipe SOLIDOR à l'IRISA, ils reposaient sur une idée très simple à résumer : tirer parti d'interactions sans fil courte portée entre calculateurs voisins, pour construire spontanément un système d'information. Bien entendu, les technologies embarquées, nécessaires à la réalisation des S.I.S, étaient nettement moins répandues et développées que ce que l'on connaît maintenant. Les premiers PDAs commençaient seulement à apparaître. Ils disposaient d'une mémoire et d'une puissance de calculs limitées. Et côté communications courte portée, les normes WiFi et Bluetooth n'étaient pas encore finalisées. Mais au sein d'une équipe de recherche à dominante « système » et enseignant dans le domaine des réseaux, j'ai souhaité explorer plus avant cette idée simple qui me semblait constituer une voie de recherche très prometteuse.

Cela m'a permis de découvrir et de m'intéresser au domaine de l'ubiquité numérique, inventé par Marc Weiser [96]. L'idée de l'ubiquité numérique, on parle également d'*informatique diffuse*, est de coupler le monde réel avec un ensemble de traitements automatiques. L'objectif est double : il s'agit de faciliter d'une part des activités réelles de la vie quotidienne, et d'autre part de profiter des activités du monde réel pour améliorer l'efficacité des systèmes informatiques. Cela suppose une connaissance par le système de la situation physique des entités qui participent à son fonctionnement. On parle alors de *sensibilité au contexte*, notion tout à fait essentielle en ubiquité numérique. Un aspect important des S.I.S était que le contexte reposait uniquement sur la proximité physique des calculateurs mobiles. C'était là un aspect original que j'ai souhaité aborder.

Cette première phase de mes travaux nous a permis de formaliser les premiers principes des S.I.S dans un rapport de recherche [13], de déposer un brevet [15], et de publier nos premières idées sur les principes systèmes des S.I.S [14, 94, 39, 12].

Nous avons rapidement souhaité associer un industriel du domaine des réseaux mobiles à nos travaux. Cette volonté s'est concrétisée à partir de 2001 par le démarrage d'une collaboration avec une équipe d'ALCATEL R&D. Au delà des aspects système dont nous avons la maîtrise, nous recherchions des compétences sur les architectures embarquées des terminaux mobiles et sur les communications sans fil, que les travaux communs avec ALCATEL R&D nous ont finalement apportés. Cette collaboration s'est avérée particulièrement fructueuse, puisqu'elle s'est poursuivie ensuite dans le cadre de mes travaux sur les réseaux 4G. Nous y revenons dans les chapitres suivants de ce document.

Sur les aspects recherche, comme cela est présenté dans ce chapitre, nous avons identifié deux problèmes principaux. Deux thèses, dont j'ai assurés le co-encadrement, ont démarré dans ce cadre, l'une portant sur la découverte de voisinage au sein du S.I.S et de la prise en compte du mouvement dans les interactions de proximité [85], l'autre s'intéressant aux architectures en ubiquité numérique et aux bases de données de proximité [77]. Tous nos résultats sur les SIS ne sont pas décrits dans ce chapitre, seuls les aspects les plus importants sont présentés. Le lecteur pourra se reporter à [84, 86, 87] pour des résultats concernant la gestion de la volatilité des interactions sans fil courte portée, et à [10] pour l'utilisation de bases de données événementielles pour l'exploitation de l'espace visible d'un S.I.S.

De manière formelle, cette phase de mes travaux sur les S.I.S s'est conclue par un ouvrage collectif [11], rédigé avec trois autres membres de l'équipe de recherche à laquelle j'appartiens. L'objectif principal de ce livre était de présenter nos recherches dans le domaine de l'ubiquité numérique. Nos principaux résultats concernant la mise en œuvre des S.I.S y sont intégrés.

Partie 2. Réseaux 4G : de la couverture discontinue au couplage de réseaux mobiles hétérogènes

MISE EN ŒUVRE DES RÉSEAUX À COUVERTURE DISCONTINUE

Un réseau à couverture discontinue, appelé également *réseau d'infostations*, peut être présenté de la manière suivante : on déploie dans l'espace physique un ensemble de points d'accès radio ou AP¹, chacun d'entre-eux définissant une cellule radio de taille réduite ou *Pico Cell* (de l'ordre de quelques dizaines de mètres) offrant des débits élevés (plusieurs dizaines de Mb/s). Ces points d'accès sont interconnectés via une infrastructure filaire, et sont déployés de telle sorte qu'ils offrent une connectivité intermittente aux usagers mobiles à certains endroits stratégiques (zones à forte densité de population par exemple). Dans un réseau cellulaire classique, la continuité de la couverture radio impose que les cellules adjacentes se recouvrent partiellement. Dans un réseau à couverture discontinue, cette obligation disparaît, et le déploiement est beaucoup plus simple : les bulles radio sont disjointes, la distribution des fréquences au niveau des points d'accès est donc nettement moins contraignante que dans une architecture où les cellules voisines doivent se recouvrir. Les technologies WiFi peuvent constituer un support pour ce type de déploiement.

Considérant les thèmes abordés dans le cadre des S.I.S, démarrer des travaux portant sur les réseaux à couverture discontinue n'a rien d'immédiat. A l'origine, mon intérêt pour cette nouvelle thématique est parti de réflexions autour des extensions possibles de nos travaux sur les S.I.S. Nous avons étudié la possibilité de connecter ponctuellement une ou plusieurs bornes fixes de communication sans fil courte portée à un S.I.S existant, autrement dit de relier un AP à un ensemble de nœuds voisins et mobiles. Nous avons fait le constat suivant : les terminaux participant à des S.I.S disposaient de capacités de stockage limitées, ce qui pouvait s'avérer très contraignant quand ces terminaux s'appuyaient uniquement sur des interactions de proximité. Par exemple, dans l'application du *Web de proximité* abordée dans chapitre précédent (*cf.* section 1.6), un nœud mobile « glane » en permanence via son S.I.S des informations pertinentes. Nos expérimentations ont montré que cela peut conduire très rapidement à une saturation de sa mémoire de stockage. Phénomène aggravant, les terminaux commençaient à être équipés de capacités de capture, comme l'enregistrement de la voix, ou un appareil photo numérique. Ces fonctions ne faisaient que renforcer la nécessité de dépasser les limites de stockage local d'un terminal. D'où notre idée de disposer d'îlots de connectivité haut débit dans l'espace physique, pour permettre à des terminaux mobiles de « décharger » automatiquement leurs données chaque fois qu'ils rencontrent une bulle radio. Si l'idée de relier ce principe à nos travaux sur les S.I.S n'a pas été poussée très loin, nous avons estimé que l'utilisation ponctuelle d'interactions sans fil courte portée offrait des perspectives intéressantes.

La principale difficulté d'une telle architecture réside dans la mise en place d'un lien logique entre les différentes bulles radio. Pris isolément, les points d'accès offrent simplement des îlots de connectivité IP, qu'on désigne souvent sous le terme de *HotSpot*. Dans le cadre d'une véritable infrastructure étendue, l'objectif est de fournir un service de transfert de données via la cellule à laquelle il se trouve temporairement rattachée. Autrement dit, toute la difficulté consiste à offrir un service continu, en dépit de la discontinuité de la couverture.

1. AP : *Access Point*

L'exploitation de réseaux d'infostations n'est pas neuve. Ainsi, les travaux pionniers menées par le WINLAB [29, 32, 30, 38] ont proposé des solutions pour des topologies très spécifiques, par exemple lorsque les utilisateurs se déplacent sur une route le long de laquelle un ensemble de points d'accès est déployé. Dans ce cas, il devient possible de précharger les données à l'avance au niveau de l'AP, de manière à ce qu'elle soit rendue disponible dès l'instant où un utilisateur entre dans la bulle radio. Cette approche fonctionne dans des cas très restrictifs où la trajectoire de l'utilisateur est connue, et où la bande passante d'une cellule n'est utilisée que par un utilisateur à la fois.

Nous avons estimé qu'il devait être possible d'aller plus loin dans le cadre d'une infrastructure discontinue et étendue exploitant des interactions courte portée, sans imposer de topologie particulière pour disposer les APs dans l'espace, et sans connaître *a priori* la mobilité des utilisateurs. L'objectif de nos travaux est de proposer des solutions système permettant le déploiement de réseaux à couverture discontinue à faibles coûts, tout en masquant la discontinuité de la couverture radio. Nous cherchons à traiter trois problématiques de recherche : mettre en place une gestion efficace des flux descendants en direction des terminaux mobiles, permettre à ces mêmes terminaux d'exploiter les flux montants en dépit de l'intermittence de la connectivité, et enfin définir les composants réseau nécessaires au déploiement à grande échelle d'une architecture discontinue. La solution proposée doit bien entendu offrir les débits élevés autorisés par les technologies sans fil courte portée à une densité importante d'utilisateurs mobiles.

Ces trois problèmes sont traités dans la suite de ce chapitre. Dans les sections 2.1 et 2.2, nous présentons notre approche pour traiter les flux descendants destinés aux terminaux. Puis nous abordons la gestion des données envoyées par le terminal en direction de l'infrastructure dans la section 2.3. Enfin, dans la section 2.4, nous proposons des mécanismes pour permettre le déploiement de cette architecture à une grande échelle.

2.1 Gestion des flux descendants dans un réseau à couverture discontinue : analyse et éléments de solutions

Pour démarrer ces travaux de recherche, j'ai souhaité m'appuyer sur une architecture réseau simple : les bulles haut débit sont déployées dans l'espace en connectant un ensemble de points d'accès à une infrastructure filaire IP grande échelle, via des liens xDSL offrant des débits de l'ordre de 6 Mb/s. Un tel choix peut sembler limitatif, considérant les performances offertes sur les infrastructures haut débit les plus récentes. Mais cette première architecture a été définie en 2004. A ce moment là, les liaisons xDSL étaient le moyen le plus simple pour un opérateur d'amener des débits élevés en n'importe quel point d'un territoire. De plus, l'utilisation d'une technologie largement déployée cadrerait avec notre objectif d'offrir un accès haut débit aux utilisateurs mobiles, pour un coût de déploiement et d'accès aux données relativement faible. Comme nous le verrons dans la suite ce chapitre, ce choix n'est pas sans conséquence pour nos travaux. Les bulles radio délivrent potentiellement un débit de plusieurs dizaines de Mb/s, alors que les performances du lien entre l'AP et l'infrastructure soient bien moindres.

Comme tout réseau cellulaire, les réseaux à couverture discontinue doivent répondre à deux objectifs : offrir une gestion efficace des flux descendants en direction des terminaux mobiles, et permettre à ces mêmes terminaux d'exploiter les flux montants vers l'infrastructure. Depuis le début des années 2000 se pose la question de « l'après 3G » dans le monde des réseaux cellulaires. Le déploiement des futurs réseaux mobiles doit permettre à une forte densité d'utilisateurs itinérants d'accéder efficacement à des flux audio et vidéo de bonne qualité, pour de la VoD², de la TMP³ etc. L'étude d'un service de *streaming* m'a semblé pertinente pour aborder les mécanismes système permettant de s'affranchir de la discontinuité de la couverture. En premier lieu, un tel service peut supporter un délai de démarrage dans la délivrance des flux descendants. Très souvent, ce délai vient de ce que le terminal doit acquérir une quantité de données afin de tolérer d'éventuelles fluctuations du débit sur le réseau. Dans le cas des réseaux à couverture continue, ce délai existe

2. VoD : *Video on demand*

3. TMP : *Télévision Mobile Personnelle*

par le fait que l'utilisateur ne rencontre pas toujours immédiatement une bulle de connectivité. Une fois le service démarré, tout le problème est de mettre à disposition des données de manière efficace quand un terminal se trouve dans une bulle radio, de façon à ce qu'il puisse les utiliser sans interruption par la suite, même en l'absence de connectivité radio. Ce problème est illustré par la figure 2.1 : un serveur de contenu transmet de manière continue un flux de données à destination de terminaux, connectés de manière intermittente à des bulles radio. Nous présentons les solutions étudiées pour traiter ce service dans la suite de ce chapitre.

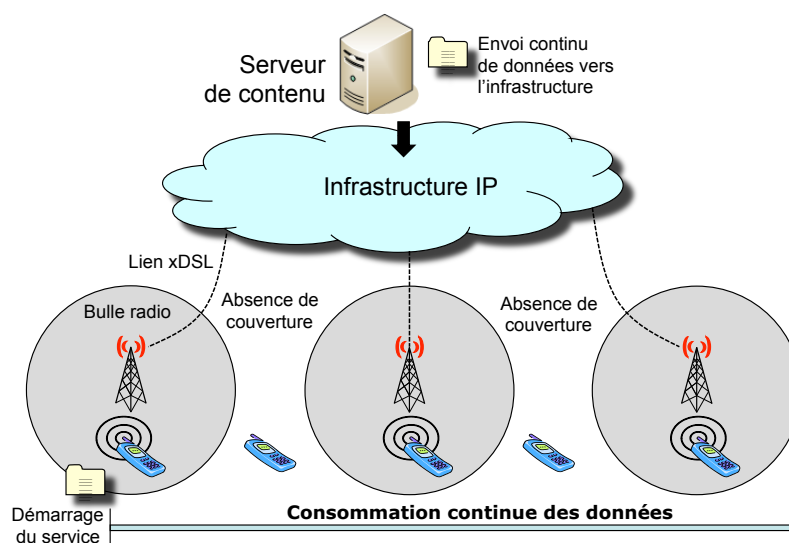


FIGURE 2.1 – Fonctionnement d'un service de *streaming* sur un réseau à couverture discontinue

2.1.1 Distribution de plusieurs niveaux de cache dans l'architecture

La délivrance d'un flux en *streaming* nécessite que le terminal client soit connecté au réseau, afin de recevoir et de consommer « au fil de l'eau » les données qui lui sont transmises par le serveur. Cette contrainte est *a priori* difficile à respecter dans le cadre d'un réseau à couverture discontinue. Les zones de connectivité offrent un débit très supérieur (plusieurs dizaines de Mb/s) à celui nécessaire pour un service de *streaming* sur mobile (quelques centaines de Kb/s). Il s'agit donc pour un terminal de « profiter » efficacement de la présence d'un point d'accès pour *précharger* une quantité *suffisante* de données dans sa mémoire, et les consommer ensuite de manière continue jusqu'à la prochaine zone de couverture. Cette notion de préchargement nous amène naturellement à considérer la possibilité d'introduire un cache dans le terminal. La quantité de données chargée dans le cache doit permettre à l'application de fonctionner sans interruption jusqu'à ce que l'utilisateur entre à nouveau dans une bulle radio. Ce principe est illustré par la figure 2.2.

Le volume de données stocké dans le cache du terminal est dépendant du débit offert par le point d'accès radio, et de la distance qui sépare les points d'accès. Si, dans le cadre du déploiement d'un réseau à couverture discontinue, ces deux paramètres sont connus et maîtrisés, un service de *streaming* peut fonctionner avec un cache uniquement dans le terminal mobile. Mais nous avons identifié deux problèmes qui font que cet unique cache n'est pas suffisant pour assurer un service continu.

Discontinuité des débits offerts dans la bulle radio. La bulle radio définie par un point d'accès offre aux usagers mobiles un débit théorique très élevé. Si on considère le cas des technologies WiFi, les deux normes les plus répandues au moment du démarrage de ces travaux, à savoir 802.11a/g, autorisent des taux de transfert de 54 Mb/s. La déclinaison la plus récente de ces normes, 802.11n, permet d'atteindre 450 Mb/s. Mais ces valeurs sont théoriques, et diminuent en

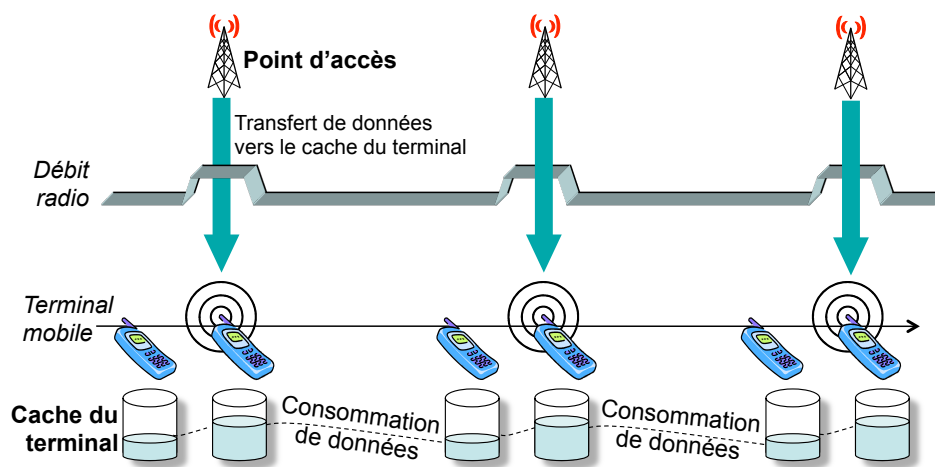


FIGURE 2.2 – Impact d'un cache dans le terminal

fonction de l'éloignement du terminal avec le point d'accès. Toujours dans le cas des réseaux WiFi 802.11, cela se traduit par un modèle des débits en couronne (*cf.* figure 2.3). Le terminal voit son débit automatiquement fixé par le point d'accès, suivant des valeurs de paliers fixées par les normes 802.11. Chaque palier correspond à une qualité du signal, et cette dernière diminue au fur et à mesure que le terminal s'éloigne du point d'accès.

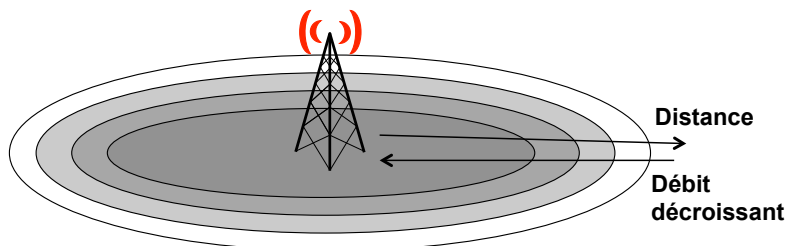


FIGURE 2.3 – Modèle des débits offert par un point d'accès

Le débit offert à un terminal dans une bulle radio utilisant une technologie WLAN n'est donc pas homogène, les performances étant très dégradées lorsqu'on s'approche de la limite de couverture. Par conséquent, la zone permettant de charger le cache terminal en utilisant un débit radio maximum est très réduite par rapport au modèle théorique envisagé. Ce problème est rendu plus critique si l'utilisateur se déplace rapidement, le temps de présence dans la zone radio « efficace » étant encore plus réduit.

Perte des données envoyées par le serveur de contenu. Le serveur de contenu envoie le flux de données de manière continue, y compris lorsque le terminal se trouve hors couverture radio. Dans ce cas, une partie des paquets transmis peut être perdue, faute de destinataire. De plus, l'infrastructure IP utilisée pour interconnecter les points d'accès peut être soumise à des congestions temporaires, provoquant des retards dans l'acheminement des flux. Dans le cas d'un service de *streaming* classique, les performances fluctuantes du réseau sont amorties en utilisant une mémoire tampon au niveau du terminal. Pour un réseau à couverture discontinue, ces fluctuations prennent une autre dimension : elles ont un impact direct sur l'arrivée du flux de données au niveau d'un point d'accès, et peuvent empêcher un terminal de profiter pleinement du débit radio quand celui-ci est disponible.

Au final, un cache unique dans le terminal ne permet de masquer que très partiellement la

discontinuité de la couverture. Afin de tirer pleinement parti du débit radio, nous avons finalement identifié trois aspects à prendre en compte dans la conception de notre architecture :

- Le débit offert à un terminal à l’entrée d’une bulle radio est faible, mais augmente significativement si ce terminal s’approche du point d’accès.
- Le terminal doit pouvoir « retrouver » les données envoyées par le serveur alors qu’il se trouvait hors couverture.
- L’exploitation du débit radio ne doit pas être limitée par d’éventuelles fluctuations de performances au niveau de l’infrastructure IP.

Une solution est de stocker temporairement les données envoyées par le serveur, autrement dit d’introduire un ou plusieurs niveaux de cache au sein de l’infrastructure. Cela permet de rapprocher les données à transmettre de la prochaine bulle traversée par le terminal, et donc de bénéficier du débit radio, sans être bridé par le débit d’émission constant du serveur de *streaming* et/ou par d’éventuelles congestions au sein du réseau. Le ou les caches peuvent également conserver les données alors que le terminal se trouve hors couverture radio.

La problématique d’utilisation et de distribution de caches au sein d’une architecture n’existe pas uniquement dans le domaine des réseaux. En architecture des ordinateurs, l’utilisation de caches multi-niveaux a été largement étudiée. L’efficacité des solutions repose avant tout sur des considérations liées aux débits de transferts des données entre les différents niveaux de cache. Pour conduire nos travaux, nous avons donc cherché des analogies possibles entre l’utilisation de caches en architecture des ordinateurs et les réseaux à couverture discontinue. Nous présentons de manière synthétique nos conclusions dans la section suivante.

2.1.2 Analogie avec les systèmes multiprocesseurs

Dans le domaine de l’architecture des ordinateurs, l’introduction de caches est motivée par des temps d’accès aux données très différents au niveau d’un processeur ou d’une mémoire principale. Le plus souvent, une hiérarchie mémoire s’appuie sur deux niveaux de cache : un cache primaire et un cache secondaire. Le cache primaire est intégré au processeur. De taille réduite, il offre au processeur un accès très rapide aux données. Le cache secondaire, de taille plus importante, est placé entre le processeur et la mémoire centrale, et permet un accès moins rapide aux données. Cette architecture a été étendue pour les systèmes multiprocesseurs [35, 66], comme l’illustre la figure 2.4.

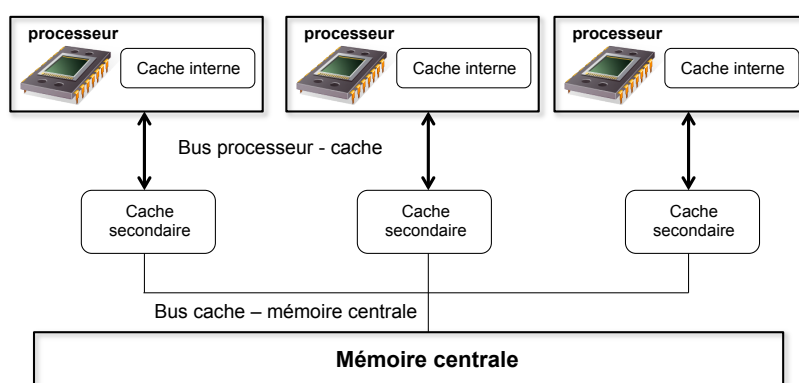


FIGURE 2.4 – Utilisation de caches dans les architectures multiprocesseurs

Cette hiérarchie mémoire n’est efficace que si elle permet de minimiser les défauts de cache : les données réclamées par un processeur doivent se trouver dans son cache interne, dans la mesure où les temps d’accès au cache secondaire ou à la mémoire centrale sont pénalisants. De plus, les processeurs échangent des données et communiquent avec une mémoire centrale par le biais

d'un bus commun, ce qui peut entraîner des problèmes de congestion. Les caches permettent de minimiser les accès directs à ce bus par les processeurs, et minimisent les goulots d'étranglement.

Nous avons cherché à identifier des analogies entre les systèmes multiprocesseurs et les réseaux à couverture discontinue. En premier lieu, il existe un parallèle évident entre processeurs et terminaux d'une part, et mémoire centrale et serveur de données d'autre part. Les terminaux (les processeurs) doivent accéder sans interruption aux données d'un serveur partagé (une mémoire commune). L'infrastructure (le bus commun) assurant la communication entre les terminaux (les processeurs) et le serveur (la mémoire) est partagée par les terminaux (les processeurs). Des problèmes de latence ou de congestion peuvent survenir lors de la délivrance des données du serveur (la mémoire centrale).

A la condition qu'il soit correctement alimenté, le cache interne du terminal (du processeur) doit permettre un accès rapide aux données du serveur (de la mémoire centrale), en s'affranchissant du débit d'émission du serveur de *streaming* (des temps d'accès à la mémoire centrale) et des latences de l'infrastructure (du bus commun). Ces analogies nous ont conduit à étudier la distribution d'un ou plusieurs caches intermédiaires distribués au sein de l'architecture, comme l'illustre la figure 2.5. Cette première proposition, formulée dans [16], est présentée dans les sections suivantes, 2.1.3, 2.1.4 et 2.1.5.

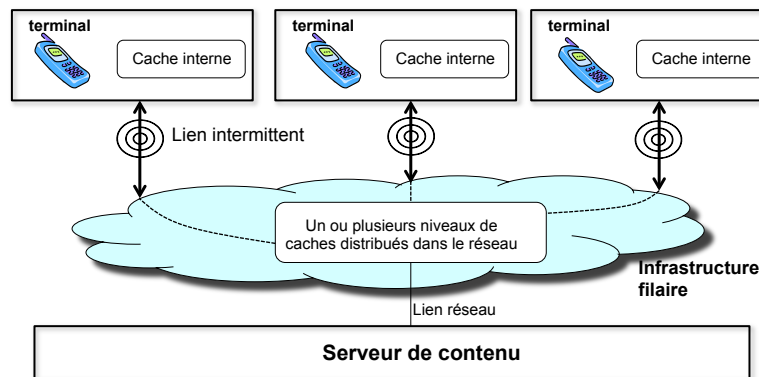


FIGURE 2.5 – Utilisation de caches dans un réseau à couverture discontinue

2.1.3 Une hiérarchie de trois niveaux de cache

Nous présentons ici la démarche qui nous conduit à introduire un ou plusieurs niveaux de cache au sein de l'architecture, en plus du cache interne intégré au terminal mobile.

Un cache dans les points d'accès radio ? Dans une architecture à couverture discontinue, c'est le lien radio qui offre le débit le plus élevé. Une approche consiste donc à intégrer un cache secondaire directement dans les points d'accès. Dans ce cas, l'interface entre le cache primaire du terminal et ce cache de niveau 2 est immédiatement disponible, dès que l'utilisateur entre dans une zone couverte. Mais cette approche n'est efficace que si les données sont présentes dans le cache de l'AP, au moment où l'utilisateur entre dans la bulle attachée à cet AP. Dans le cas contraire, il peut se produire un défaut du cache secondaire, et le remplissage du cache terminal est limité par les capacités de l'infrastructure et le débit du serveur de *streaming*.

Autrement dit, les données doivent être présentes dans le cache de l'AP, avant que le terminal n'entre dans la bulle radio associée. Une telle hypothèse est valable quand la trajectoire de l'utilisateur est connue ou prédictible. C'est le cas de l'étude du WINLAB, citée dans l'introduction de ce chapitre [29, 30, 32]. Dans le cadre de nos travaux, la mobilité des utilisateurs est non contrainte. La prédiction du « prochain point d'accès » traversé n'est pas possible. Une solution consistant à diffuser le flux en permanence vers la totalité des APs n'est pas acceptable, en regard de la charge de trafic engendrée dans l'infrastructure. Nous cherchons donc d'autres solutions de distribution.

Un cache dans l'infrastructure filaire ? Une approche possible est de placer un cache au dessus des points d'accès, « plus haut » dans l'infrastructure. Un tel cache secondaire reçoit en permanence les données de *streaming*, et peut potentiellement les communiquer aux APs connectés à l'infrastructure IP, sans savoir *a priori* quelles seront les zones radio traversées par le terminal. Bien entendu, cette approche n'est efficace que si le lien entre le cache secondaire et le cache primaire n'est pas limitant.

L'introduction d'un cache secondaire dans l'infrastructure implique l'existence d'un équipement réseau en mesure de « l'héberger », cet équipement devant pouvoir communiquer avec l'ensemble des points d'accès offrant la connectivité radio. Afin de valider cette approche, nous avons étudié les évolutions des technologies sans fil susceptibles de supporter la mise en œuvre d'un réseau à couverture discontinue, principalement WiFi et WiMax. Pour ce qui est de WiFi, l'IETF a proposé le standard CAPWAP [97, 57, 3, 33], dont l'objectif est de simplifier le déploiement de réseaux sans fil à grande échelle. Un équipement spécifique, un *contrôleur*, réalise l'interconnexion d'un groupe de points d'accès, et centralise différentes fonctions, comme la planification des ressources radio, la gestion de l'authentification des utilisateurs etc. Pour le déploiement de réseau WiMax 802.16e [6], un équipement appelé MAP⁴ [21] est connecté à plusieurs stations de base WiMax, et prend en charge la mobilité des utilisateurs au niveau IP. Ces deux composants réseau, en charge du pilotage de groupes de points d'accès ou de stations de base, peuvent potentiellement héberger les mécanismes d'un cache secondaire, à savoir l'interception des flux du serveur et leur transmission vers un cache primaire quand le terminal cible entre dans une zone radio. Dans la suite de ce document, nous appelons l'équipement réalisant ces mécanismes un contrôleur d'accès (ou AC⁵). L'introduction de ce niveau de cache est représentée dans la figure 2.6.

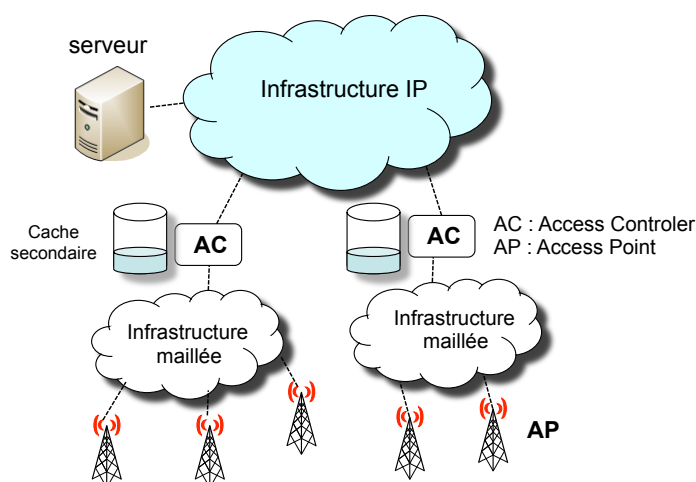


FIGURE 2.6 – Utilisation d'un cache intermédiaire dans l'infrastructure

Au final, trois niveaux de cache. Un cache de second niveau dans l'AC n'est efficace que si le lien entre l'AC et l'AP ne limite pas l'usage du lien radio censé alimenter le cache interne du terminal. Or, comme cela a été montré par la figure 2.1, le déploiement des APs est assuré via des liaisons offrant un débit de l'ordre de 6 Mb/s. C'est donc ce dernier qui conditionne le remplissage du cache du terminal, lorsqu'il entre dans une zone de couverture. Considérant cette limite, nous avons fait le choix d'utiliser deux niveaux de cache dans l'infrastructure : un cache de niveau 2 au sein de l'AP, et un cache de niveau 3 dans un AC. Cette hiérarchie sur trois niveaux est théoriquement efficace, considérant que :

- Le cache de l'AC permet d'intercepter le flux continu envoyé par le serveur, sans nécessiter que le terminal cible se trouve dans une bulle radio.

4. MAP : *Mobility Anchor Point*

5. AC : *Access Controller*

- Le cache de l'AP permet d'alimenter le cache du terminal à hauteur du débit radio offert.
- Enfin, le cache interne permet de masquer la discontinuité de la couverture.

Ces trois principes sont illustrés dans la figure 2.7.

Un telle hiérarchie n'est que « théoriquement efficace » dans la mesure où le problème attaché à l'utilisation d'un cache uniquement dans l'AP demeure : la mobilité des usagers étant non contrainte, il n'est pas possible d'identifier l'AP de rattachement d'un utilisateur avant que celui-ci ne rentre dans la zone de couverture. L'alimentation du cache de l'AP ne peut pas démarrer avant que l'entrée de l'utilisateur dans une bulle radio soit détectée. Dans ces conditions, le transfert des données entre le cache du terminal et le point d'accès est limité par la capacité du lien AC-AP. Nous présentons notre solution à ce problème dans la section suivante.

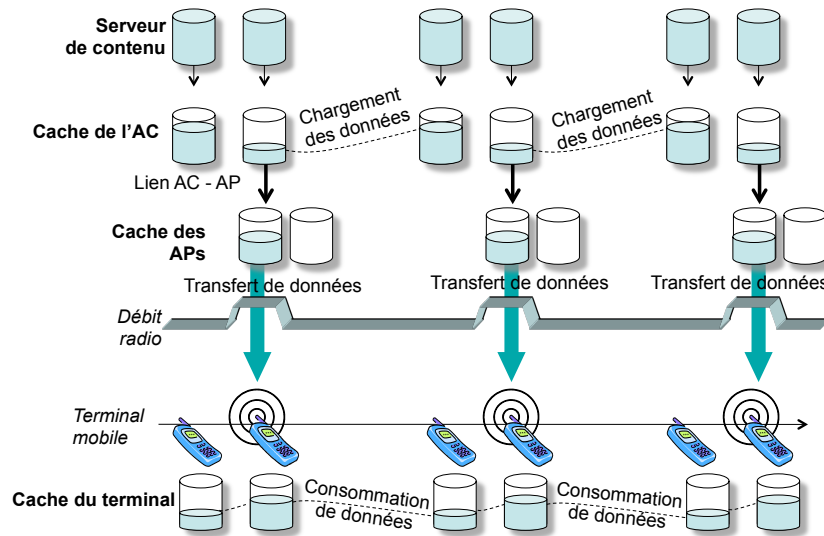


FIGURE 2.7 – Utilisation de trois niveaux de caches distribués dans l'infrastructure

2.1.4 Discrimination des débits dans les bulles radio

La mobilité d'un utilisateur dans le réseau n'étant pas contrainte, le cache d'un point d'accès ne peut être alimenté qu'à partir du moment où l'utilisateur entre dans la zone radio attachée à ce point d'accès. Partant de cette contrainte, nous avons reconsidéré l'exploitation qui peut être faite d'une zone radio de taille réduite. L'approche la plus évidente consiste à utiliser la totalité de la zone de couverture pour envoyer des données au cache du terminal. Mais comme cela a déjà été évoqué dans la section 2.1.1 de ce chapitre, les débits offerts dans la zone de couverture ne sont pas homogènes : les performances sont d'autant plus élevées que la distance entre l'utilisateur mobile et le point d'accès est faible. Pour les réseaux WiFi, cela se traduit par le modèle en couronne présenté au début de ce chapitre (*cf.* figure 2.3). Les écarts de performances peuvent être conséquents. A titre d'exemple, [5] montre qu'un usager se trouvant à 80 mètres d'un AP WiFi dispose d'un débit de 2 Mb/s, alors qu'un utilisateur distant de 13 mètres de ce même point d'accès bénéficie d'un débit de 54 Mb/s.

En conséquence, les débits offerts à un ensemble de terminaux se trouvant dans la même bulle radio dépendent de la position de ces terminaux par rapport à l'AP. De plus, la ressource radio est limitée, elle doit être partagée entre les utilisateurs. Dans le cas d'un réseau WiFi, un terminal mobile monopolise la totalité du canal radio quand il obtient le droit de transmettre et recevoir des données, les autres restant en attente. Ce droit est géré au niveau de la couche MAC, via le protocole CSMA/CA. Ainsi, un terminal se trouvant loin de l'AP remplira son cache avec un débit faible, au détriment d'un autre se trouvant « en attente » au centre de la bulle, et ne pouvant disposer d'un taux de transmission nettement plus élevé.

Nous avons donc fait l'analyse, dans un premier temps de manière intuitive, que traiter les terminaux de manière identique sous couverture d'un même d'AP n'est pas efficace pour alimenter leurs caches internes. Ceci nous a amené à considérer la possibilité de différencier les terminaux en fonction de leur position dans la zone radio. L'approche est la suivante : il n'est pas nécessaire d'exploiter la totalité de la bulle radio pour alimenter efficacement les terminaux. Seule la zone centrale (*i.e.* les couronnes offrant les débits les plus élevés) est exploitée pour transférer des données du cache de l'AP vers le cache du terminal. Le reste de la zone de couverture (*i.e.* les couronnes délivrant les débits les moins élevés) peut être utilisée pour détecter l'entrée d'un terminal dans la zone radio, et déclencher alors le préchargement des données du cache de l'AC vers le cache de l'AP.

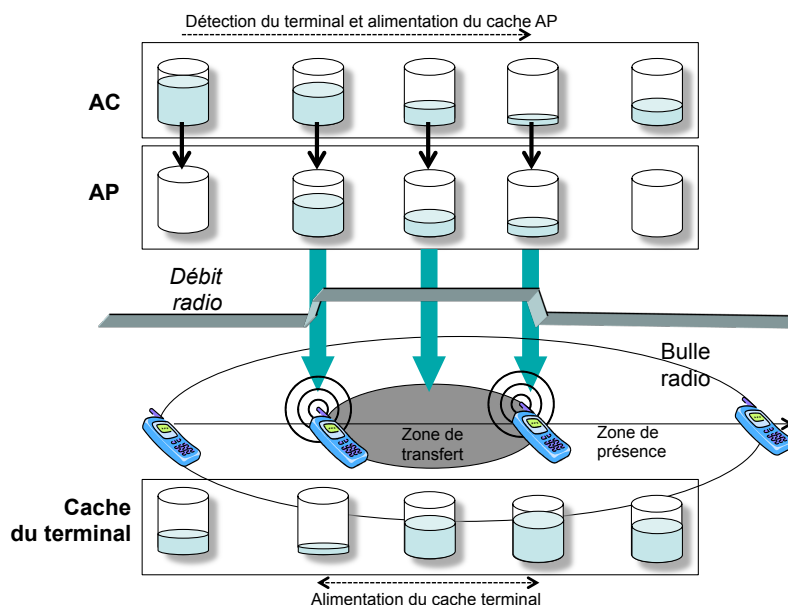


FIGURE 2.8 – Chargement des caches AP via une discrimination des débits radio

Au final, nous proposons donc de diviser la bulle radio associée à un AP en deux zones distinctes :

1. La *zone de présence (ZP)* correspond à la zone dans laquelle les débits offerts sont les plus faibles, comparable au débit du lien AC-AP. Cette zone permet à un terminal de signaler sa présence, et de déclencher le transfert de données de l'AC vers l'AP.
2. La *zone de transfert (ZT)* recouvre les couronnes offrant les débits les plus élevés, permettant ainsi au terminal d'alimenter efficacement son cache interne. Cette efficacité doit être renforcée par le fait que le cache de l'AP se remplit depuis l'entrée du terminal dans la bulle radio.

Ces principes sont présentés dans la figure 2.8.

Pertinence d'un modèle radio à deux zones. Avant d'aller plus loin, nous avons cherché à valider plus formellement l'efficacité du modèle proposé. Restreindre les transmissions vers les terminaux aux zones de présence des APs ne doit pas être pénalisant. Autrement dit, la quantité de données chargée par un terminal en exploitant deux zones distinctes dans les bulles radio ne doit pas être inférieure à la quantité de données reçues par ce même terminal quand la totalité des zones de couverture sont utilisées pour alimenter le cache interne.

Pour ce faire, reprenons le modèle des débits en couronne, déjà présenté dans la figure 2.3. La couverture radio est représentée par un ensemble de zones concentriques, chacune d'entre-elles étant associée à un taux de transfert des données. Notons B_i le débit offert par la zone i , et r_i son rayon, B_1 étant la zone offrant les meilleures performances. En supposant que les utilisateurs

B : capacité moyenne de la bulle radio B_i : débit de la zone i r_i : rayon de la zone i $S_i = \pi \cdot (r_i^2 - r_{i-1}^2)$: surface de la zone i $B = \frac{\sum_{i=0}^n B_i \cdot S_i}{\sum_{i=0}^n S_i}$ <p>(a) Distribution spatiale</p> $d_i = r_i - r_{i-1}$: largeur de la zone i $B = \frac{\sum_{i=0}^n B_i \cdot d_i}{\sum_{i=0}^n d_i}$ <p>(b) Distribution linéaire</p>

FIGURE 2.9 – Capacité moyenne d’un point d’accès

sont uniformément distribués dans la bulle radio, la capacité B du point d’accès (*i.e* la quantité de données que la cellule va traiter par unité de temps) peut être calculée par la formule 2.9(a). Considérons un modèle de mobilité très simple, illustré par la figure 2.10, dans lequel un utilisateur se déplace en ligne droite avec une vitesse constante v , du centre d’un AP vers un autre AP distant de d . La densité linéaire des utilisateurs se trouvant dans le réseau (nombre d’utilisateurs / m²) est notée λ . La capacité peut alors être évaluée via la formule 2.9(b).

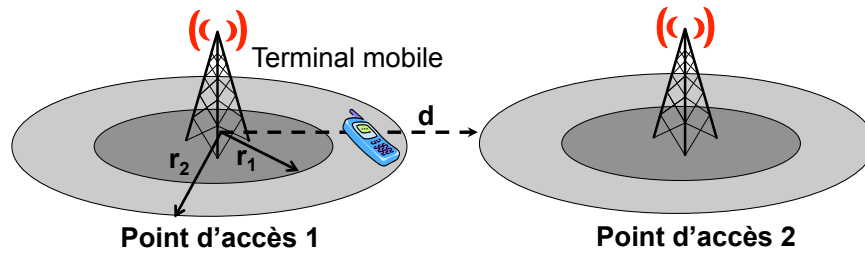


FIGURE 2.10 – Modèle de mobilité entre deux APs

Nous cherchons à estimer la bande passante moyenne utilisable pour un utilisateur durant son trajet d . Si on exploite la totalité de la zone, l’utilisateur dispose durant son déplacement d’un débit égal à la capacité B de la bulle radio divisée par le nombre d’utilisateurs se trouvant sur sa trajectoire, soit $bw = \frac{B}{d \cdot \lambda}$. En appliquant le modèle à deux zones, le terminal se déplace durant $\frac{d}{v}$, et son cache est alimenté durant son passage en zone de transfert, soit $\frac{r_1}{v}$. La bande passante « perçue » par l’utilisateur sur son trajet devient $vbw = \frac{B_1}{r_1 \cdot \lambda} \frac{r_1}{v} = \frac{B_1}{d \cdot \lambda}$. Plus de détails concernant ces calculs peuvent être trouvés dans [16, 93].

Au final, $\frac{vbw}{bw} = \frac{B_1}{B} > 1$. Sur un plan théorique, un utilisateur dispose donc d’une bande passante moyenne supérieure s’il n’accède aux données qu’en zone de transfert.

Mécanisme de remplissage des caches. Dans le cadre de la hiérarchie que nous proposons, le cache de l’AC stocke de manière continue le flux transmis par le serveur de *streaming*. Le cache de l’AP est alimenté dès qu’un terminal entre dans la zone de présence de la bulle radio associée. Reste à traiter le problème plus spécifique du cache terminal.

En effet, plusieurs utilisateurs peuvent se trouver dans la zone de transfert d’un AP, et sont donc tous susceptibles de recevoir des données du cache de cet AP. Dans un réseau WiFi, la ressource radio est partagée de manière équitable : les terminaux sont servis par la couche MAC (la couche gérant l’accès au medium radio) à tour de rôle, sans privilège d’accès particulier. Dans ce cas, le remplissage d’un terminal sera d’autant retardé que la densité d’utilisateurs présents en zone de transfert est élevée. De plus, l’état du cache interne d’un terminal à son entrée dans une bulle radio

dépend de plusieurs facteurs : le temps passé par ce terminal hors couverture radio, et le volume de données transmis dans la bulle radio précédemment traversée. Dans ces conditions, un terminal peut entrer en zone de couverture avec un cache pratiquement en famine, et devoir malgré tout attendre l'accès à la ressource radio à cause d'une densité élevée de terminaux déjà présents dans la zone de transfert.

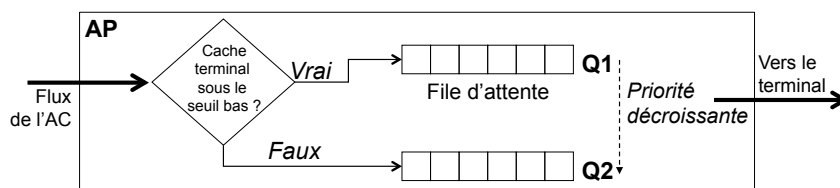


FIGURE 2.11 – Politique de distribution des données au niveau du cache AC

Pour traiter ce problème, nous avons proposé une politique d'ordonnancement simple couplée à la couche MAC au sein de l'AP, dont le but est de servir en priorité les terminaux dont l'état du cache risque d'entraîner une interruption du service. Cette politique est illustrée par la figure 2.11. Elle s'appuie sur deux files d'émission, Q1 et Q2. Dès qu'un terminal entre dans la zone de transfert d'un point d'accès, il indique le niveau de son cache. Si ce dernier se situe en dessous d'un *seuil bas*, les données nécessaires pour amener son cache au dessus de ce seuil sont placées dans la file Q1, la plus prioritaire. Le reste des données est mis en attente dans la file Q2. Cette dernière est vidée uniquement si la file Q1 est vide.

2.1.5 Analyse de l'impact de la hiérarchie de caches sur la délivrance des flux descendants

Notre première contribution dans le cadre des réseaux à couverture discontinue repose donc sur une hiérarchie de caches à trois niveaux. Nous avons cherché à évaluer cette proposition par le biais de simulations. Nous avons fait le choix de développer notre propre environnement de simulation, afin de traiter les mécanismes spécifiques à l'architecture étudiée : gestion de caches, politique de distribution de données en fonction de l'état de ces caches, découpage en zones des couvertures radio etc. Notre outil s'appuie sur DESMO-J [1], un moteur à événements discrets conçu par l'université de Hambourg. L'architecture du simulateur de réseaux à couverture discontinue a été conçue afin de correspondre à l'infrastructure présentée au début de ce chapitre. Chaque nœud (serveur, AC, AP, terminal) est représenté par un composant DESMO-J spécifique, auquel est attaché un ensemble de paramètres qui permet d'en décrire les caractéristiques : présence d'un cache dans un nœud du réseau, débit des liens entrants et sortants, profil des mobilité des terminaux. Cet environnement de simulation a constamment évolué avec l'avancement de nos travaux, et a également été utilisé par ALCATEL R&D avec qui nous avons collaboré sur une partie de nos recherches. Pour plus détails concernant l'architecture interne de l'outil de simulation, le lecteur peut se reporter à [16, 70].

Définition des critères d'évaluation. La hiérarchie de caches proposée a été évaluée suivant trois critères, les deux premiers étant rattachés à une mesure objective de la qualité du service délivré, le dernier permettant de mesurer l'impact des caches dans l'infrastructure.

En premier lieu, la présence des caches doit permettre de minimiser le *nombre d'interruptions de service par terminal (critère 1)*, autrement dit le nombre de fois où le cache d'un terminal se trouve en état de famine. Mais une interruption n'est pas forcément préjudiciable si elle est de courte durée. Ce principe renvoie à un critère important pour les opérateurs, celui du « ressenti » de l'utilisateur par rapport au déroulement d'un service. On parle également de *user experience*. Nous avons donc cherché à évaluer la *durée moyenne d'un service sans interruption (critère 2)*.

En ce qui concerne la gestion des caches, l'utilisation de deux zones dans la bulle radio (zone de présence et zone de transfert) doit amener une utilisation efficace du lien radio offert par un

AP : un terminal entrant en zone de transfert doit pouvoir charger ses données depuis le cache de l'AP, donc avec un débit égal à celui du lien radio. Ce qui nous a amené à mesurer *l'utilisation de la bande passante du lien radio (critère 3)*.

Analyse des principaux résultats. Pour un service de streaming à 256 kb/s (ce qui correspond à un flux de bonne qualité sur un terminal mobile), les simulations ont permis de mesurer un peu moins de 4 interruptions par heure et par utilisateur, pour une durée moyenne de service d'un peu plus de 8 minutes. Ces chiffres sont à comparer aux mesures effectuées avec un cache dans l'AP désactivé : plus de 50 interruptions, pour une durée de service moyenne de moins de 60 secondes. Pour ce qui est du troisième critère, le débit radio utilisé est équivalent au débit théorique offert par l'AP dans la zone de transfert lorsque la hiérarchie de caches sur trois niveaux est utilisée.

Ces simulations ont mis en évidence l'impact attendu des caches dans l'AC, les APs et les terminaux, combiné à une bulle radio divisée en deux parties, une zone de présence et une zone de transfert. Néanmoins, le niveau des interruptions reste trop élevé du point de vue des critères d'un opérateur.

Introduction d'un seuil de déclenchement et d'une rafale de données. Partant de ce constat, nous avons analysé la manière dont ces interruptions sont distribuées dans le temps. De manière très nette, il est apparu que les famines au niveau des caches des terminaux se produisaient principalement en début de simulation, ce qui peut être interprété de la manière suivante :

- les terminaux effectuent leur demande de service avec un cache vide.
- Par conséquent, le flux reçu dès la première bulle radio traversée commence immédiatement à être consommé par l'application de *streaming* sur le terminal, ce qui provoque des interruptions de service.

En d'autres termes, le cache d'un terminal ne peut pas toujours atteindre un niveau « suffisant » pour lequel la quantité de données reçues dans les bulles traversées lui permet de compenser le flux consommé entre deux zones de couverture. Ce constat nous amène à introduire le principe d'un *seuil de déclenchement* dans le fonctionnement du service. Ce seuil correspond à la quantité de données minimale que doit contenir le cache du terminal, en dessous de laquelle l'application de *streaming* ne peut pas démarrer la lecture du flux. Une fois ce niveau atteint par le cache du terminal, le *seuil de déclenchement* ne joue plus aucun rôle, et c'est le seuil bas, introduit dans la section 2.1.4 qui régule la gestion des données vers les terminaux. Ce principe a été validé par le biais de simulations. Pour mémoire, sans seuil de déclenchement, la durée moyenne d'un service sans interruption est d'un peu plus de 8 minutes, avec un nombre moyen d'interruptions de service par terminal légèrement inférieur à 4. En fixant un seuil de déclenchement équivalent à 210 secondes à partir de la demande de service, la durée moyenne d'un service continu va au delà de 39 minutes, avec une seule interruption observée par terminal.

La continuité dans la délivrance du flux est améliorée de manière significative par l'introduction de ce second seuil. Mais sa prise en compte par le cache terminal retarde le démarrage du service pour l'utilisateur. Si le seuil est à 0, la durée moyenne d'accès au service est de 38 secondes. Cette valeur monte à 235 secondes pour un seuil fixé à 210 secondes, ce qui n'est pas acceptable en terme de « ressenti » pour l'utilisateur. Aussi, nous avons étudié un mécanisme permettant au serveur de transmettre une rafale de données (un *burst*) à un débit supérieur à celui imposé par le service de *streaming*. Ce *burst* est maintenu de manière temporaire, de la demande de service par le terminal, jusqu'à ce que le cache de ce même terminal atteigne le seuil de déclenchement. Si ce dernier est fixé à 210 secondes, la durée moyenne d'accès au service tombe de 235 secondes (sans rafale) à 165 secondes (avec rafale). L'exploitation combinée d'un seuil de déclenchement et d'une rafale de données est illustrée par la figure 2.12.

2.1.6 Bilan et principaux enseignements

En démarrant ses travaux de recherche, mon objectif était de montrer qu'il était possible de déployer une infrastructure étendue

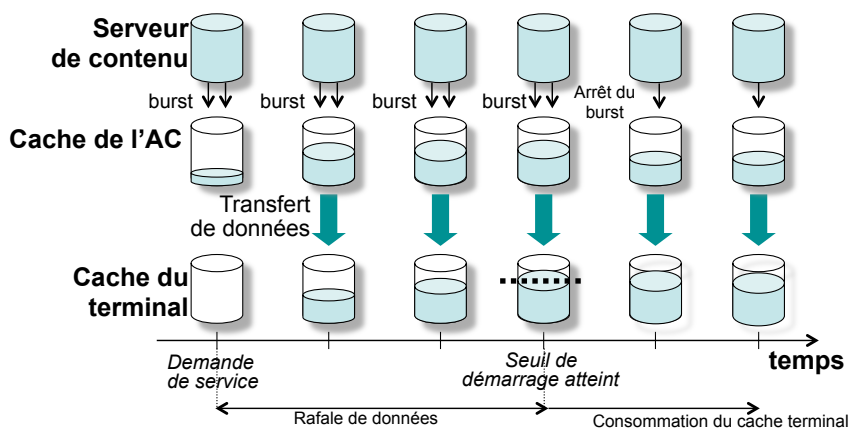


FIGURE 2.12 – Utilisation d'un seuil de déclenchement et d'une rafale de données

- en utilisant un ensemble de petites cellules radio fournies par des technologies sans fil courte portée,
- sans rechercher une continuité de la couverture.

L'intérêt de cette démarche est double : la discontinuité (*i.e.* le non-recouvrement) des cellules simplifie la planification radio du réseau (*i.e.* la distribution des fréquences), et l'utilisation des technologies sans fil courte portée permet d'offrir des débits très élevés aux utilisateurs se trouvant sous couverture radio.

Cette démarche repose sur une conviction de départ finalement assez simple : en dépit de la couverture discontinue, un flux descendant peut être traité de manière continue par un terminal, en utilisant des mécanismes système de gestion de caches au sein du réseau. Partant d'une analogie avec les architectures multiprocesseurs, nous avons proposé une hiérarchie de caches sur trois niveaux, dans le terminal, dans le point d'accès et dans le contrôleur d'accès, un équipement réseau qui joue le rôle de point de rattachement entre les APs et le serveur. La gestion de ces caches repose sur trois principes :

1. La bulle radio peut être divisée en deux zones, une zone de présence, et une zone de transfert.
2. Les terminaux sont alimentés uniquement en zone de transfert.
3. La distribution des données, de l'AP vers le terminal, privilégie les terminaux ayant le niveau de cache le plus bas.

Ce début de chapitre présente les bases de nos résultats concernant les réseaux à couverture discontinue. Une partie de ces résultats ont donné lieu à une première thèse [16]. En parallèle, le démarrage d'une nouvelle collaboration avec Alcatel R&D nous a permis de discuter et poser les caractéristiques d'une première architecture réseau : nature des liens dans le réseau filaire, performances attendues des bulles sans fil en fonction de différentes technologies courte portée, validation d'un modèle radio à deux zones, analyse des équipements réseaux susceptibles d'accueillir des caches.

Un aspect important de cette architecture est que la présence d'un cache dans un point d'accès se justifie par la faiblesse du lien AC-AP. Cette hypothèse, pertinente au début de nos recherches sur les réseaux à couverture discontinue, doit être reconsidérée : l'utilisation de plus en plus massive des technologies Ethernet dans les réseaux d'opérateurs, permet d'accroître les débits offerts dans les infrastructures étendues (plusieurs centaines de Mb/s, voir plusieurs Gb/s). De plus, notre approche repose sur une vision très simple des liens entre les différents composants (serveur, AC, AP, terminal). Nous avons considéré ces liens comme de simples canaux d'échanges de données, caractérisés simplement par un débit. Or, leur exploitation par le réseau repose sur l'utilisation d'un ensemble de protocoles. Par exemple, un flux de *streaming* utilise une pile de protocoles RTP

/ UDP / IP. Il nous a semblé donc essentiel d'étudier les interactions possibles entre ces protocoles et les mécanismes de distributions de données entre les caches.

Tous ces constats doivent nous amener à reconsidérer la gestion des caches au sein de notre architecture. C'est cette deuxième phase de nos travaux sur les liens descendants que nous présentons dans la suite de ce chapitre.

2.2 Les débits augmentent ... et un niveau de cache disparaît

Un cache dans l'AC n'est efficace que si le lien entre l'AC et l'AP ne limite pas l'usage du lien radio alimentant le cache du terminal. Avec l'augmentation des débits offerts au sein de l'infrastructure filaire, il devient envisageable d'alimenter directement le cache du terminal avec le cache de l'AC. Dans ces conditions, le cache intermédiaire au sein de l'AP ne se justifie plus. Le fonctionnement d'une hiérarchie de caches sur trois niveaux est illustré par la figure. Cette dernière peut être comparée à la figure 2.7 dans laquelle le rôle d'un cache dans l'AP est mis en évidence.

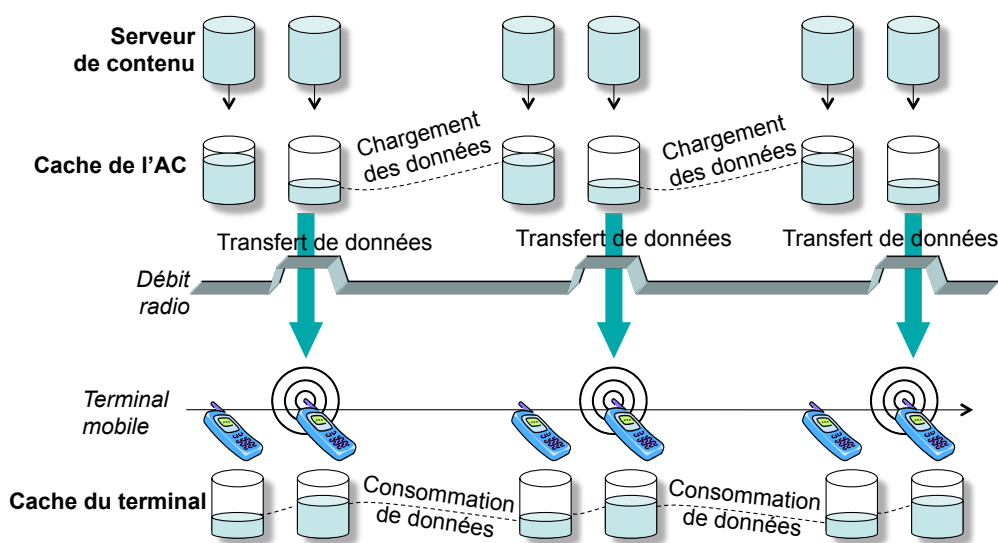


FIGURE 2.13 – Utilisation de deux niveaux de caches distribués dans l'infrastructure

Au delà de cette simplification de notre architecture, nous avons cherché à prendre en compte de manière plus fine les propriétés du canal radio de l'AP, ainsi que l'environnement protocolaire pouvant être embarqué dans l'AC. Nous présentons nos résultats sur ces sujets dans les sections suivantes.

2.2.1 Traitement de la discontinuité au niveau de l'AC

Le principe d'un protocole de *streaming* comme RTP est de transmettre des paquets de données de manière régulière, le rythme d'émission étant imposé par le débit du flux audio/vidéo. Cette régularité permet au récepteur de restituer le flux d'origine, un tampon de réception permettant de compenser des décalages éventuels induits par des fluctuations de performances du réseau. De plus, une voie de retour permet au terminal client d'informer le serveur de ces fluctuations et l'état de son tampon récepteur. Le serveur peut alors sur cette base adapter son rythme d'émission.

Dans l'architecture considérée, le lien serveur - AC est stable, alors que les échanges AC - terminal reposent sur une liaison radio disponible de manière intermittente. Le stockage dans le cache de l'AC peut amener des retards importants dans la délivrance des données au terminal. De plus, ces données sont transmises en rafale au terminal dès lors que ce dernier entre dans une zone

de transfert. Un serveur de *streaming* utilisant un protocole comme RTP risque alors d'adapter de manière erronée le rythme d'émission de ses paquets.

Nous proposons donc de masquer la présence du terminal au serveur : l'AC joue le rôle de *client virtuel* vis à vis du serveur, permettant ainsi au protocole de *streaming* de fonctionner correctement entre le serveur et l'AC. Et c'est un protocole spécifique, appelé protocole de transport de cache, qui traite l'envoi des données entre le cache de l'AC et le cache du terminal. Ce principe est résumé dans la figure 2.14 [70].

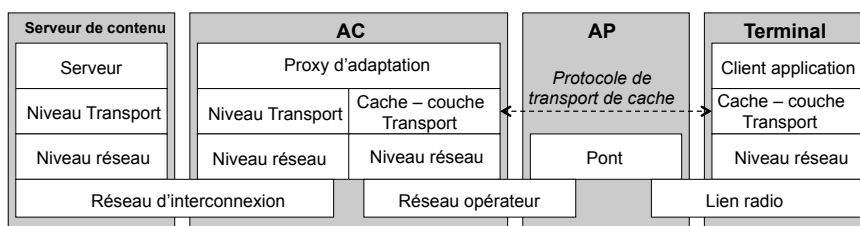


FIGURE 2.14 – Piles de protocoles au sein du serveur, de l'AC et du terminal

2.2.2 Définition d'un protocole de transport de cache supportant la discontinuité

La discontinuité de la couverture modifie la « forme » traditionnelle de la bande passante utilisée par le client d'un service de *streaming*, ainsi que l'illustre la figure 2.15. Cette bande est représentée par une valeur moyenne lorsque la couverture est continue. Par contre, dans le cas d'une couverture discontinue, elle est nulle tant que l'utilisateur se trouve hors couverture. Dès que le terminal arrive en zone de transfert, on peut la représenter sous forme d'une rafale (*burst*) qui permet au cache interne de se remplir. La bande passante reprend ensuite la valeur moyenne imposée par le serveur de *streaming*.

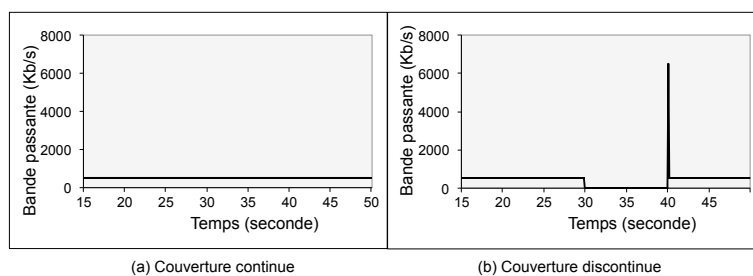


FIGURE 2.15 – Évolution de la bande passante durant le trajet d'un terminal

Le protocole de transport de cache doit être en mesure de traiter ce profil de trafic. La rafale de données se produisant à l'entrée du terminal en zone de transfert peut provoquer des pertes de données, notamment à cause de congestions provoquées au niveau d'un AP pendant une rafale de données.

Nous proposons l'approche suivante pour traiter ce problème : le protocole de transport de cache doit réguler le flux des données de l'AC vers le terminal, en fonction de la bande passante disponible entre ces deux nœuds. Ce protocole doit être *suffisamment fiable* pour assurer le transfert efficace des données en dépit des rafales qui peuvent se produire. L'utilisation d'une estimation de la bande passante de « bout en bout » est une approche classique pour prévenir les congestions dans un réseau. On parle de mécanismes de contrôle de flux au niveau Transport. Ce problème a été traité dans des protocoles comme TCP ou SCTP [31]. PR-SCTP [67], une extension du protocole

SCTP, propose un mécanisme de fiabilité partielle, qui permet de supprimer des paquets arrivant en retard, et qui donc adapté au transport de flux audio/vidéo dans notre cas.

Calcul de la bande passante entre l'AC et les terminaux. Cette évaluation de la bande passante, couplée à un protocole de transport, pose un problème spécifique dans le cadre des réseaux sans fil. La plupart des approches interprète la perte de paquets comme une preuve de congestion dans le réseau. Cette hypothèse est réaliste dans un réseau filaire où les liens présentent des taux d'erreurs très faibles. Un lien sans fil, comme WiFi par exemple, offre une qualité de liaison très variable dans le temps. Nous avons pu constater ce problème une première fois dans le cadre de l'expérimentation des S.I.S. Dans une architecture discontinue, cela signifie qu'un protocole de gestion de cache s'appuyant sur SCTP ne sait pas distinguer les pertes de données dues à des congestions, de celles provoquées par l'instabilité du lien sans fil.

Nous avons donc étudié différentes solutions alternatives d'estimation de la bande passante prenant en compte le caractère spécifique d'une transmission sans fil [19]. La formule *Westwood* [52, 20] est l'approche qui nous a semblé la plus adaptée à notre problématique, et nous l'avons retenu dans le cadre de notre architecture. Nous ne décrivons pas ici le mode de calcul retenu. Le lecteur peut se reporter à [70, 75].

Intégration du contrôle de flux dans l'AC et dans le terminal. Comme cela est illustré dans la figure 2.16, un AP est attaché à un ordonnanceur spécifique dans l'AC. Chaque terminal cible attaché à cet AP se voit associer au niveau de l'AC une fenêtre de transmission et un temporisateur de retransmission. Ces deux éléments permettent au protocole SCTP de réguler les flux transmis vers chaque terminal, tout en assurant la fiabilité de la connexion.

Les données à envoyer vers chaque terminal sont soumises à une file d'attente unique associée à l'ordonnanceur. C'est le débit mesuré entre le terminal et l'AC par l'estimateur de bande passante (utilisant la formule *Westwood*) qui autorise la fenêtre de transmission du terminal à se vider plus ou moins rapidement, l'objectif étant d'éviter les congestions au niveau du point d'accès.

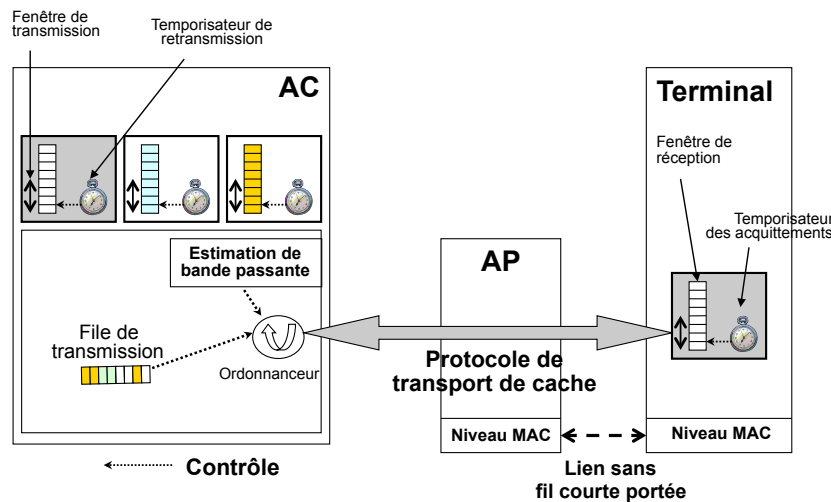


FIGURE 2.16 – Intégration du contrôle de flux au sein de l'AC

Analyse des principaux résultats. Notre environnement de simulation de réseaux à couverture discontinue a été étendu pour supporter le protocole de gestion de cache et les mécanismes d'évaluation de bande passante. A titre d'exemple, nous présentons les résultats obtenus pour un service de *streaming* de 512 Kb/s, plus exigeant que lors de nos premières campagnes de simulations où nous avons alors considéré un service de 256 Kb/s. La figure 2.17 illustre l'évolution du

nombre d'interruptions de service en fonction de la densité des utilisateurs présents dans le réseau, dans trois cas : avec un protocole de transport de cache non fiable (courbe 1), avec le protocole de transport fiable que nous proposons (courbe 2 - SCTP *like* combiné avec une évaluation de bande passante de l'AP « Westwood »), et enfin ce même protocole combiné avec la rafale de données (courbe 3) dont le principe a été discuté dans la section 2.1.5. On constate l'impact très positif du protocole de transport entre l'AC et les terminaux mobiles. Sans fiabilité, les premières interruptions interviennent pour une densité de 10 utilisateurs/ km^2 . Ce chiffre est multiplié par 4 avec un protocole de transport fiable, et par 5 si ce même protocole est combiné avec le mécanisme de rafale de données.

La totalité de nos résultats concernant la gestion des flux descendants n'est pas présentée dans ce document. Des améliorations de l'algorithme d'ordonnancement au sein de l'AP, ainsi qu'une utilisation optimisée du seuil de déclenchement (introduit dans la section 2.1.5) ont été proposés et publiés dans [75]. Nous avons également étudié la possibilité d'améliorer la transmission de flux audio/vidéo en exploitant des propriétés de la norme H.264 SVC⁶ [40, 42].

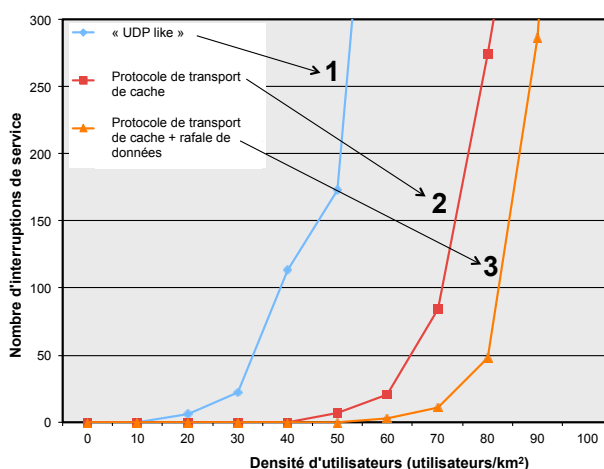


FIGURE 2.17 – Évolution du nombre d'interruptions de service

2.2.3 Bilan

En abordant la thématique des réseaux à couverture discontinue, notre premier objectif était de montrer que des flux descendants (*i.e.* en direction du terminal) pouvaient être délivrés efficacement à des terminaux mobiles, en dépit de la discontinuité de la couverture. Notre proposition repose sur une hiérarchie de caches, distribuée entre le serveur de données et les terminaux. Cette hiérarchie a évolué au cours de nos travaux, passant de trois à deux niveaux de caches, en fonction de l'évolution des débits offerts par l'infrastructure filaire.

L'architecture système proposée s'articule autour d'un équipement central, le contrôleur d'accès ou AC, qui peut être vu comme un *point d'attachement* entre le serveur et les terminaux mobiles en attente d'être servis. L'AC assure la politique de distribution des données transmises par le serveur. Cette politique privilégie les zones des bulles radio offrant les débits les plus élevés. On parle de zones de transfert, par opposition au « reste » des bulles qualifié de zones de présence. Ces dernières sont utilisées uniquement pour détecter l'entrée des terminaux sous couverture radio.

En couplant la politique de distribution de données à un protocole de transport fiable, dont le rythme d'émission est asservi à une évaluation continue de la bande passante disponible dans chaque bulle radio, nous avons montré qu'un service descendant est à même de servir une densité importante d'utilisateurs sans rupture de service. Le respect de ces deux conditions est essentiel pour intégrer les réseaux à couverture discontinue dans les futurs réseaux cellulaires.

6. SVC : Scalable Video Coding

Nous abordons dans la section suivante nos travaux concernant l'exploitation des flux montants d'un réseau à couverture discontinue.

2.3 Exploitation efficace des flux montants

Nous avons envisagé différentes approches pour tirer parti des flux montants dans une infrastructure discontinue, notamment la possibilité pour un utilisateur mobile de mettre à disposition en temps réel un flux audio/vidéo en cours de capture [53]. Je présente ici ce que je considère comme le cœur de notre contribution concernant la gestion des flux montants : nous proposons d'exploiter la *popularité* des données qu'un client souhaite « télécharger » dans l'infrastructure. Ce principe est présenté dans la suite de cette section.

Les terminaux mobiles récents embarquent de nombreuses fonctions d'acquisition de données : appareil photo numérique, caméscope, enregistreur audio, GPS etc. Partant de là, considérons le cas où des milliers de touristes prennent des photos et des séquences vidéo dans une grande ville. Il est facile d'identifier de nombreuses raisons pour lesquels on cherche à transmettre ces données vers le réseau : ces touristes doivent pouvoir envoyer ces données à leur famille, ou bien encore les terminaux mobiles peuvent ne pas disposer de suffisamment d'espace mémoire pour stocker toutes ces données. Ce problème, déjà évoqué dans l'introduction de ce chapitre, est une des raisons qui nous a amenée à étudier les réseaux à couverture discontinue. Ces utilisateurs mobiles souhaitent envoyer des photos et des séquences vidéo vers un site personnel ou le serveur de stockage d'un fournisseur de services. Le transfert de données peut prendre plus ou moins de temps selon la bande passante sans fil offerte par la liaison montante. Le temps moyen nécessaire pour transférer des données vers le serveur applicatif conditionne directement la disponibilité de ces données pour d'autres utilisateurs. Ce principe est illustré par la figure 2.18.

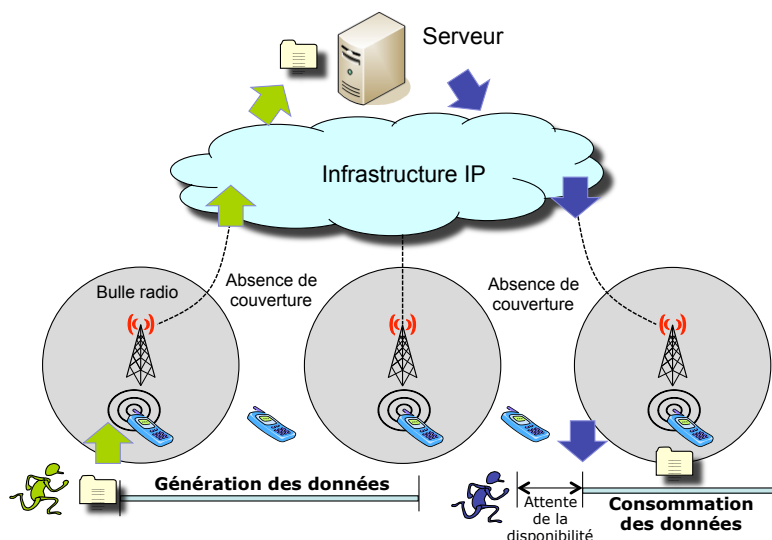


FIGURE 2.18 – Exploitation d'un flux montant dans un réseau à couverture discontinue

Les données captées par le terminal via un appareil photo doivent être « déchargées » vers l'infrastructure filaire. Dans le cadre de l'architecture discontinue proposée dans la section 2.2, on peut identifier deux entités susceptibles de recevoir les données produites par un terminal, le serveur de contenu et le cache de l'AC auquel est attaché l'utilisateur mobile. Si les données sont maintenues dans un AC, elles seront rapidement accessibles pour tous les utilisateurs connectés au réseau, car l'opérateur du service contrôle cette partie de l'infrastructure. Cet aspect est intéressant pour les applications où le partage est motivé par la proximité géographique des utilisateurs. Dans ce cadre, ce groupe d'utilisateurs proches est rattaché au même AC, et bénéficie de la localisation

des données dans le cache de cet AC. Cette approche n'est pas nouvelle, on peut trouver des problématiques similaires pour les services de *streaming vidéo*. En effet, tout comme un service exploitant les flux montants, le *streaming* repose sur un producteur unique de données (un écrivain), et un ensemble de consommateurs (des lecteurs) distribués dans l'infrastructure réseau. A titre d'exemple, dans [49], un flux vidéo produit est divisé en segments élémentaires, diffusés vers des caches distribués dans l'infrastructure réseau. Une telle segmentation présente des avantages : elle permet de distribuer la charge, et il est possible d'implémenter une gestion des caches qui dépasse le simple principe du « tout ou rien ». Si certains segments sont plus populaires que d'autres, seuls les segments populaires sont mis en cache. Autrement dit, seules les parties les plus demandées du flux vidéo doivent être présentes dans les caches aux extrémités du réseau.

Mais au delà de cette co-localisation des lecteurs, un terminal producteur peut également souhaiter voir remonter rapidement les données jusqu'au serveur, afin que ces dernières soient accessibles au plus grand nombre d'utilisateurs. On voit donc, en ce qui concerne la gestion de la liaison montante, qu'il existe diverses approches possibles pour la relation entre les différents éléments de l'architecture discontinue. Le terminal qui a capturé les données, s'attend à un comportement spécifique de l'AC et du serveur dont il dépend. Nous avons traduit ce comportement par la notion de classes de service (*CoSs, Class of Services*) [70]. Une CoS est une information attachée à une donnée créée au niveau d'un terminal. Elle conditionne la façon dont cette donnée doit être gérée tout au long du chemin depuis le terminal jusqu'au serveur applicatif. Dans le cadre de nos travaux, la CoS jugée la plus intéressante est celle liée à la *popularité* des données. Pour les liaisons montantes, nous caractérisons la *popularité* des données via le nombre des demandes reçues par le serveur applicatif pour lire ces données.

A partir de là, deux problèmes doivent être traités. Le premier concerne les données qui ne sont pas encore entièrement téléchargées vers le serveur de contenu. Certaines données sont demandées par les utilisateurs plus que d'autres. Or, dans la plupart des réseaux sans fil, la bande passante radio offerte par le lien montant est très limitée. Il est donc nécessaire de favoriser le transfert de ces données « demandées » en direction du serveur. Le second problème est d'être en mesure de vider les caches des terminaux mobiles quand les utilisateurs ont besoin de libérer leur mémoire de stockage locale. Face à ces deux problèmes, nous avons montré que les ACs distribués dans l'infrastructure, initialement utilisés pour masquer la discontinuité, peuvent jouer un rôle significatif pour les applications utilisant la liaison montante. L'emplacement d'un AC contrôlant un ensemble de points d'accès, lui permet de mettre en œuvre différents mécanismes en liaison avec les terminaux et les serveurs de contenu.

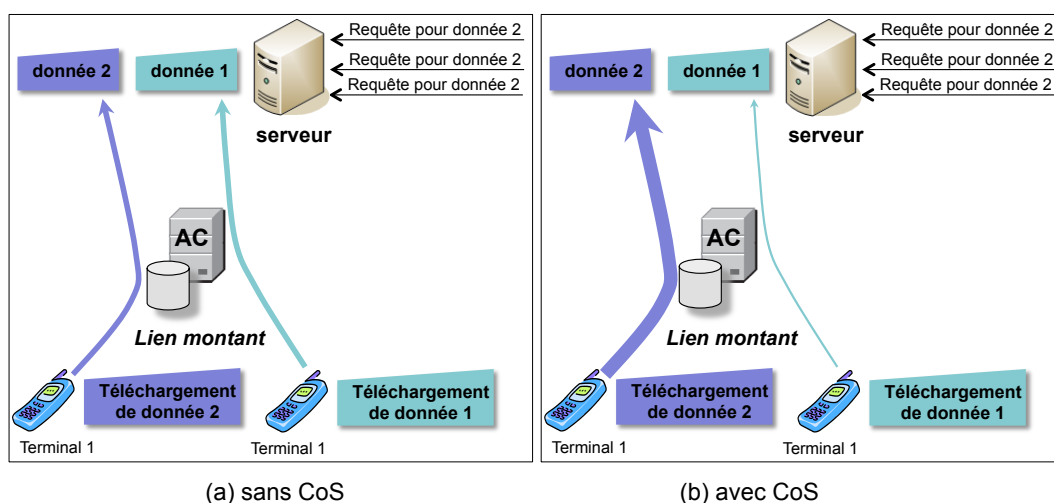


FIGURE 2.19 – Principes des CoSs

Une CoS impose donc un comportement spécifique à un AC. Ce dernier échange des messages avec les terminaux et des serveurs, afin que le service délivré aux utilisateurs soit amélioré. Ainsi, l'*indice de popularité* attaché à une donnée correspond au nombre d'utilisateurs demandant l'accès

à cette donnée pendant qu'elle est en cours d'envoi du terminal vers le réseau. L'AC doit récupérer cet *indice de popularité* au niveau du serveur de contenu. Dans un premier scénario, *indice de popularité* = 0 signifie que les données n'ont pas été demandées par des terminaux « lecteurs ». Elles ne nécessitent donc pas un surcroît de débit au niveau du lien montant. Par contre, *indice de popularité* > 0 signifie que les données sont *populaires*, il est donc pertinent de leur octroyer davantage de bande passante. La figure 2.20 illustre le principe des classes de service. L'information attachée à une CoS permet de lier l'exploitation du flux montant au sein de la bulle radio, et les mécanismes de caches déployés au sein d'un AC.

L'impact de l'exploitation des flux montants s'appuyant sur des CoSs doit se mesurer

- sur le temps moyen nécessaire pour transférer les fichiers *populaires* et *non populaires* (que nous qualifions d'*impopulaires* par la suite),
- sur le temps nécessaire pour les rendre accessibles sur le serveur applicatif,
- et sur la moyenne de bande passante allouée pour transférer les flux de données populaires et impopulaires.

L'architecture à couverture discontinue repose sur deux couples d'équipements : terminal-AC et AC-serveur de contenu. Dans la suite de cette section, nous présentons une synthèse de nos principaux résultats portant sur l'exploitation de ces deux couples via le principe des CoSs.

2.3.1 Problématique terminal-AC

Le lien intermittent entre un terminal et son AP de rattachement exploite des ressources radio limitées. De plus, la couche protocolaire (la couche MAC) qui régit l'accès à ce lien peut s'appuyer sur un partage de ces ressources qui privilégie le lien descendant, comme c'est le cas par exemple pour la technologie WiMax.

Nous proposons donc via les CoSs d'accorder la bande passante du lien montant en priorité aux terminaux

- générant les données présentant un *indice de popularité* non nul,
- ou dont la mémoire de stockage locale est proche de la saturation.

Ainsi, la bande passante radio allouée à chaque terminal doit être augmentée ou diminuée au cours de la communication, en fonction de ces paramètres. Pour cela, il est possible de s'appuyer sur des mécanismes de contrôles de flux analogues à ceux des protocoles TCP et SCTP [31]. Nous n'avons pas traité cette problématique dans le cadre des flux montants. Mais elle a été abordée « dans l'autre sens », dans le cadre du protocole de transport entre l'AC et les terminaux, pour la gestion des flux descendants (*cf.* section 2.2.2). Un mécanisme de calcul de bande passante comme *Westwood* [52, 20], combiné à un contrôle de flux via une fenêtre glissante, doivent permettre de réguler la quantité de donnée qu'un terminal est autorisé à transmettre sur le lien qui le rattache à l'AC, sans perturber le contrôle d'accès MAC à la couche radio utilisé par le terminal et l'AP.

Afin qu'une telle approche ne dégrade pas le ressenti du service pour l'utilisateur, les données *populaires* sont traitées de manière équitable. Ainsi, elles se voient offrir la même bande passante à partir du moment où leur *indice de popularité* est non nul. Cette approche a été validée via différentes simulations. A titre d'exemple, la figure 2.20 présente quelques résultats significatifs. Durant cette phase de simulation, nous avons évalué le transfert des fichiers de différentes tailles (1 MB, 10 MB) et d'un flux de données. Plusieurs taux de popularité ont été testés, 0%, 10% et 50%. Comme attendu, les résultats mettent en évidence que le temps moyen nécessaire pour télécharger un fichier populaire est nettement inférieur à celui d'un fichier impopulaire. Si la densité d'utilisateurs est faible, le temps moyen nécessaire pour télécharger un fichier populaire est assez similaire à celui d'un fichier impopulaire. Par contre, lorsque la densité des utilisateurs augmente, le temps moyen nécessaire pour transférer les fichiers populaires augmente lentement, alors que celui des fichiers non populaires croît rapidement. Plus il y a de fichiers populaires, moins le service fourni aux utilisateurs est efficace [73, 75].

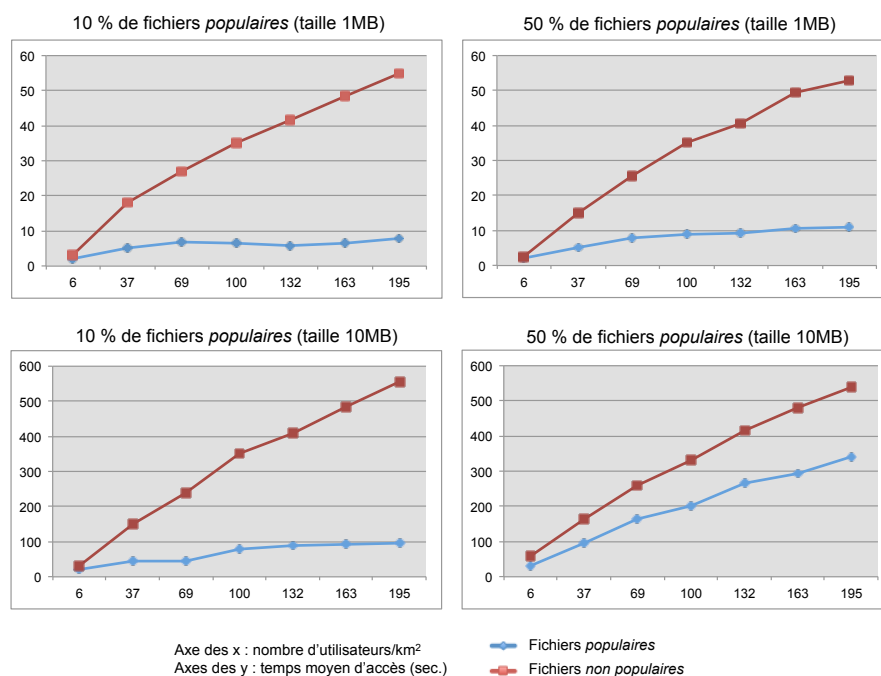


FIGURE 2.20 – Analyse des performances sur le lien terminal - AC

2.3.2 Problématique AC-serveur

Comme cela a été évoqué dans la section 2.1.1, l'infrastructure IP qui assure l'interconnexion entre les serveurs et les contrôleurs d'accès peut être soumise à des congestions temporaires, provoquant retards et pertes dans la remontée des données. L'opérateur qui déploie le réseau à couverture discontinue n'a pas obligatoirement la maîtrise du lien entre l'AC et le serveur de contenu. Et bien entendu, lorsque de nombreux ACs sont fortement chargés, cela peut entraîner des goulots d'étranglement vers le serveur. Le problème est donc est de garantir le transfert des données en dépit de ces limitations. C'est sur ce plan que le cache de l'AC peut jouer un rôle bénéfique : capable de stocker temporairement des données dans sa mémoire, il peut retarder la transmission des données vers le serveur de contenu. Pour ce faire, l'AC utilise l'*indice de popularité*, pour décider quelles données doivent être envoyées au serveur, et quelles données peuvent être temporairement maintenues dans son cache.

Gestion des interruptions au niveau des ACs. Un flux montant est intermittent en raison de la discontinuité de la couverture. Par voie de conséquence, le flux de données (ou *session* de données) envoyé par un terminal peut être divisé de manière plus ou moins aléatoire sur un ou plusieurs ACs, en fonction de la mobilité de l'utilisateur. Par exemple, lorsqu'un terminal sort d'une bulle radio attachée à AC_1 , l'envoi des données est interrompu. Et si ce même terminal entre ensuite dans une zone radio reliée toujours au même AC (AC_1), la session de données doit être conservée par AC_1 . Si par contre, le terminal change d'AC de rattachement (par exemple AC_2), une nouvelle session doit être établie dans AC_2 . Dans notre approche, les données transmises sont traitées comme des fichiers divisés en segments, et le contrôleur d'accès est responsable de la gestion et du recueil de ces segments.

Considérant une architecture discontinue, il est possible qu'un AC contienne des segments multiples, consécutifs ou non, d'un même fichier. De plus, ces segments peuvent être distribués sur plusieurs ACs. Sur la figure 2.21, pour une même session, le terminal quitte AC_1 pour AC_2 , puis s'attache de nouveau à AC_1 . Afin de traiter ce type d'interruptions, chaque AC de l'architecture maintient une table des fichiers en cours de chargement et des segments associés. Chaque élément de la table est associé avec une CoS, un serveur de contenu, la taille du fichier associé et l'index du

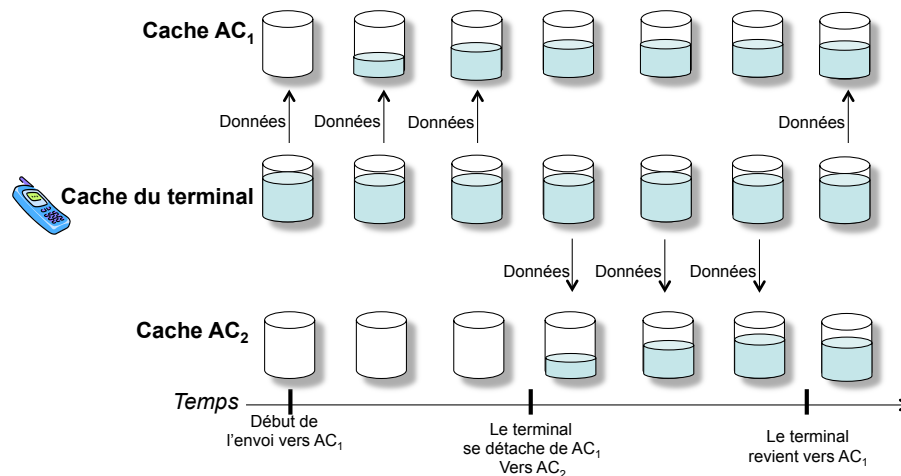


FIGURE 2.21 – Gestion de la discontinuité

premier et du dernier segment. Dans le cas où le terminal reste attaché au même AC après chaque interruption de couverture (on parle d'interruptions *intra-AC*), le fichier est reconstruit dans le cache de l'AC. Et face à des interruptions *inter-AC*, le fichier est divisé entre plusieurs ACs. Dans un schéma impliquant un écrivain unique (le terminal qui produit les données) et plusieurs lecteurs (les terminaux qui consomment les données au niveau des ACs ou des serveurs de contenu), ce type de distribution ne pose pas de problème de cohérence, la reconstruction du fichier étant assurée au niveau du serveur de contenu.

Mise à disposition des données au niveau du serveur. Retarder la délivrance de données vers le serveur de contenu les rend indisponible pour les autres utilisateurs, tout particulièrement si les segments attachés sont distribués sur plusieurs ACs. Dès qu'un utilisateur décharge ses données dans le réseau, le serveur de contenu référence le fichier correspondant, avec un *indice de popularité* initialisé à 0. Tant que cet indice conserve cette valeur initiale, les données peuvent être maintenues dans le ou les ACs sans compromettre le fonctionnement du service. Par la suite, quand ce fichier est réclamé par d'autres utilisateurs, l'*indice de popularité* est incrémenté, et le serveur doit récupérer l'ensemble des segments d'un ou plusieurs ACs, afin de mettre le fichier à disposition des autres utilisateurs. Pour cela, comme l'illustre la figure 2.22, le terminal écrivain transmet avant l'envoi des premiers segments un ensemble de paramètres décrivant le comportement attendu du service. Parmi ces paramètres figure un temps maximum au delà duquel le fichier est mis à disposition sur le serveur, même si l'*indice de popularité* est resté à 0.

Les performances de ces mécanismes ont été évaluées via des simulations. Nous avons mis en évidence que, lorsque les données sont de grandes tailles, les performances des services rendus aux utilisateurs ayant des données populaires et impopulaires sont proches. Et l'utilisation de CoSs uniquement entre les ACs et les serveurs applicatifs n'apporte pas de réelle amélioration au ressenti de l'utilisateur. Par contre, l'utilisation des CoSs « tout au long du chemin » entre les terminaux mobiles et les serveurs applicatifs permet une amélioration sensible de la qualité du service délivré. Pour plus de détails, le lecteur pourra se reporter à [73, 75].

2.3.3 Bilan

Après avoir abordé la thématique des réseaux à couverture discontinue au travers des flux descendants (*i.e.* en direction du terminal mobile), nous avons étudié des éléments de solution pour l'exploitation des flux montants (*i.e.* partant du terminal mobile). Cette démarche partait du constat suivant : un terminal mobile doté de mécanismes de capture est capable de produire une quantité importante de données, mais dispose d'une capacité de stockage limitée pour les conserver.

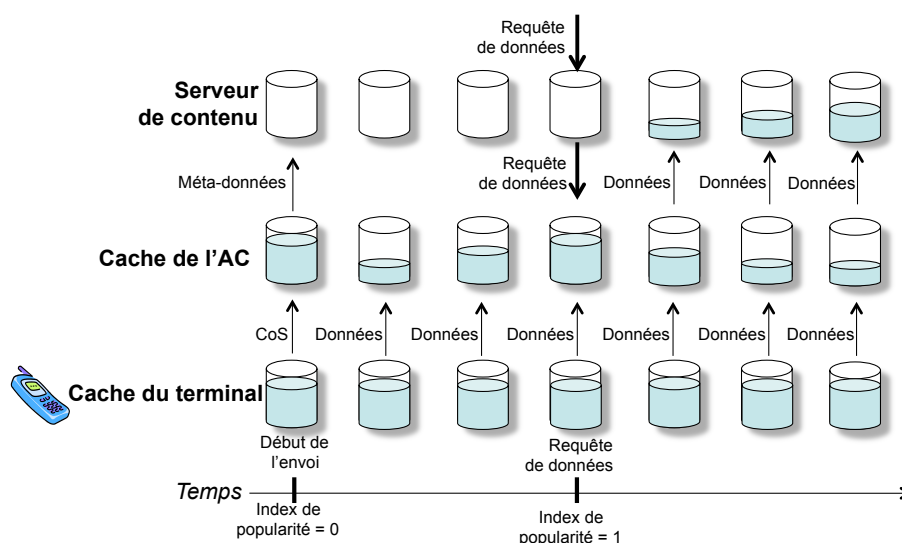


FIGURE 2.22 – Mise à disposition des données vers le serveur

Nous avons cherché à nous appuyer sur le principal élément système introduit au début de nos travaux : un contrôleur d'accès ou AC, vu comme un « point d'attachement » entre le serveur de contenu et les terminaux mobiles. La présence intermittente d'un lien montant doit permettre à un terminal mobile producteur de données de décharger ses informations dans le cache de l'AC auquel il est rattaché, afin de les mettre à la disposition des terminaux attachés au même AC. Mais si on souhaite dépasser cette co-localisation des utilisateurs matérialisée par l'attachement à un même AC, le terminal producteur doit pouvoir demander à ce que ces données soient remontées rapidement au delà de l'AC, jusqu'au serveur de contenu.

Notre contribution repose sur la notion de classes de services ou *CoS*. Il s'agit d'une information attachée par le terminal producteur à une donnée qu'il a créée, et qui va conditionner l'exploitation des liens terminal - AC et terminal - serveur. Nous avons présenté dans cette section un cas de *CoS* spécifique, lié à la *popularité* des données. En complément des simulations qui nous ont permis de valider les performances de notre proposition, l'exploitation des *CosS* a été mise en œuvre dans le cadre d'une plate-forme de démonstration. Des compléments à ce sujet peuvent être trouvés dans [70, 26].

2.4 Déploiement à grande échelle d'un réseau à couverture discontinue

Plus le déploiement d'un réseau à couverture discontinue est étendu, plus le nombre de points d'accès nécessaire est élevé. Ces APs sont reliés logiquement les uns aux autres au travers de leur connexion à un contrôleur d'accès. Cet AC constitue un point de stockage intermédiaire pour les flux montants (venant des terminaux) et descendant (en direction des terminaux). Pour des raisons d'efficacité, de capacités de stockage, de réactivité, un AC ne peut piloter qu'un nombre limité d'APs. Une architecture discontinue et étendue doit donc s'appuyer sur plusieurs contrôleurs d'accès. Dans ce contexte, nous avons mis en évidence deux problématiques.

Découverte de l'AP « courant ». La première est liée à la façon dont les terminaux mobiles, qui rejoignent le réseau, peuvent découvrir l'AC « courant » auquel ils doivent se rattacher. Nous avons proposé un mécanisme de découverte s'appuyant sur le protocole SIP⁷ [61]. Ce dernier, standardisé par l'IETF, a été conçu pour établir, modifier et terminer des sessions multimédia.

7. SIP : *Session Initiation Protocol*

Utilisé notamment pour les services de VoIP⁸, un des intérêts de SIP est qu'il offre des mécanismes de signalisation très souple entre les entités souhaitant établir une session (pour échanger de la voix par exemple). Il est possible par exemple d'utiliser un *proxy* (appelé proxy SIP), qui peut servir d'intermédiaire entre deux terminaux qui ne connaissent pas leur localisation respective (représentée pour chacun par leur adresse IP).

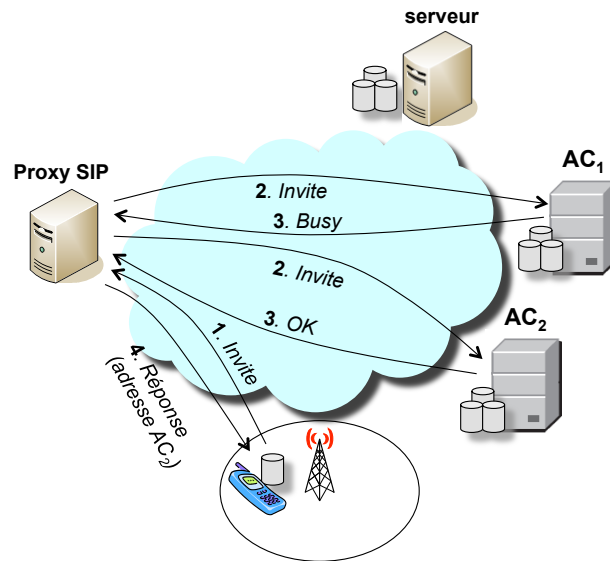


FIGURE 2.23 – Découverte de l'AC courant

Dans notre proposition, un proxy SIP est utilisé pour permettre au terminal qui entre dans le réseau de découvrir automatiquement son AC de rattachement. Le terminal envoie un message d'ouverture de session *Invite*. Et c'est le proxy qui détermine la localisation « la plus proche » du terminal, à savoir l'adresse IP de l'AC de rattachement. Ces échanges sont résumés de manière simple sur la figure 2.23. Notre approche est détaillée dans [70, 63, 26].

Passage entre deux ACs. Le second problème est lié à la mobilité entre deux ACs. En effet, les terminaux mobiles peuvent se déplacer d'un AC à l'autre, ce qui conduit à ajouter une latence importante dans la livraison de données. Ce délai peut être réduit via des mécanismes de prédiction du prochain AC de rattachement. C'est là le cœur de notre contribution pour assurer le déploiement à grande échelle d'un réseau à couverture discontinue. Nous y consacrons la suite de cette section.

2.4.1 Mécanisme de découverte de l'AC candidat

Dans cette partie, les travaux présentés reposent sur l'étude d'un protocole visant à déterminer le prochain AC de rattachement. Les mécanismes recherchés doivent permettre d'assurer le transfert proactif des données destinées à un terminal, d'un point d'attache du réseau mobile à un autre. Par exemple, le protocole IAPP⁹ [4, 55], défini dans le cadre des réseaux WiFi, a été conçu pour assurer le transfert des données entre deux APs. Cette approche n'est pas applicable dans notre architecture. En effet, dans l'infrastructure étudiée dans la section 2.2, un point d'accès n'est qu'un simple pont radio, sans cache, les mécanismes de distribution des données étant placés au niveau de l'AC.

Notre proposition, appelée *Neighbor Discovery Protocol (NPD)*, tient compte de différentes contraintes de l'architecture, principalement la couverture discontinue et la présence de contrôleurs d'accès équipés d'une mémoire cache. Contrairement à IAPP, quand un terminal se déplace entre

8. VoIP : *Voice over IP*

9. IAPP : *Inter-Access Point Protocol*

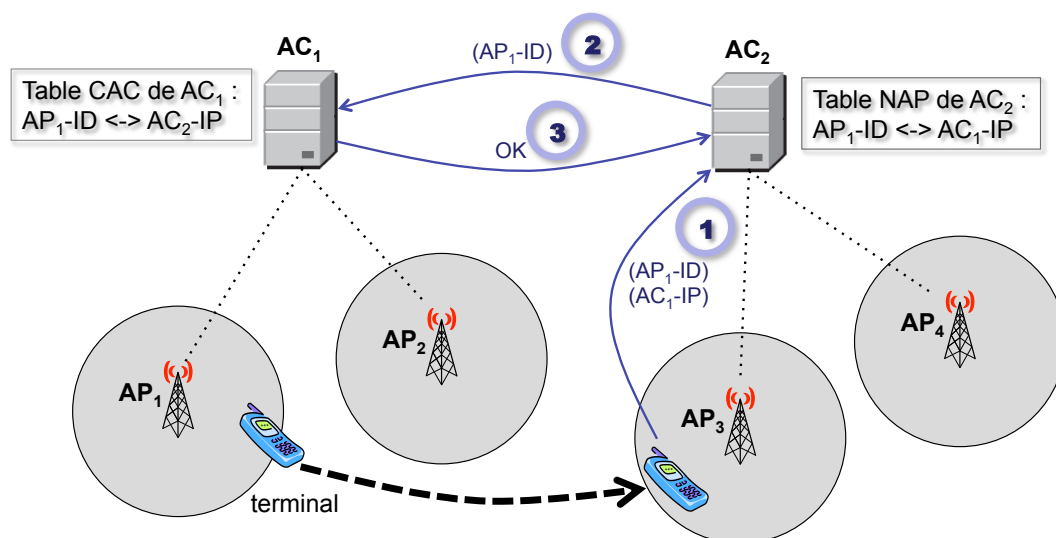


FIGURE 2.24 – Découverte distribuée de l'AC suivant

deux APs attachés au même AC, il n'est pas nécessaire d'effectuer un transfert proactif des données, dans la mesure où aucune information n'est conservée par un point d'accès.

Le principe de NPD est simple : il s'agit de construire et maintenir une table qui associe à chaque AP les ACs *candidats* au rattachement, c'est à dire l'ensemble des ACs auquel le terminal est susceptible de se lier après être sorti de la bulle radio d'un AP. Cette table, configurée de manière dynamique, est appelée *CAC*¹⁰. Elle peut être construite soit d'une manière distribuée [74, 72], soit via des mécanismes centralisés [71]. Dans la version distribuée, des messages de contrôle sont échangés entre le terminal et l'AC, et entre les ACs, afin d'assurer l'apprentissage entre un AC et les AC candidats « potentiels ». Dans l'approche centralisée, une entité spécifique appelée *Mobility Proxy* centralise les messages de contrôle provenant des terminaux.

2.4.2 Performances de la version distribuée du mécanisme de découverte

A titre d'exemple, nous présentons dans cette section les mécanismes principaux de l'approche distribuée pour construire la table *CAC* : chaque AC communique et « apprend » de ses ACs voisins. Dans ce but, chaque AC configure deux tables :

1. Une table *CAC* qui associe chaque AP avec les ACs candidats : $\{AP_{courant}, AC_{suivant}\}$. Via cette table, l'AC est en mesure de prendre une décision concernant la réalisation de transferts de données en direction des ACs candidats.
2. Une table appelée *NAP*¹¹. L'AC y enregistre les APs appartenant à d'autres ACs et s'associant avec lui comme prochain AC. Cette table permet à l'AC de déterminer s'il est encore un candidat pour ces APs ou non.

La figure 2.24 présente de quelle manière un AC apprend à connaître les ACs candidats de ses propres points d'accès. Le terminal est initialement rattaché à AP_1 et il se déplace vers AP_3 , il envoie un message d'enregistrement ❶ au nouvel AC (AC_2). Ce message est utilisé pour désigner auprès de AC_2 l'ancien AP (AP_1-ID) et l'adresse IP du précédent AC (AC_1-IP). AC_2 vérifie sa table *NAP*, si AP_1-ID n'est pas déjà enregistré en tant que voisin, AC_2 envoie à AC_1 l' AP_1-ID ❷. Dès réception de ce second message, AC_1 actualise sa table *CAC* et associe l' AP_1-ID avec l'adresse IP de AC_2 . Puis AC_1 répond positivement à AC_2 ❸. A la réception de la réponse, AC_2 ajoute AP_1-ID et l'adresse IP de AC_1 à la table *NAP*. La procédure, représentée par les messages ❶, ❷

10. *CAC* : *Candidate Access Controller*

11. *NAP* : *Neighbor Access Points*

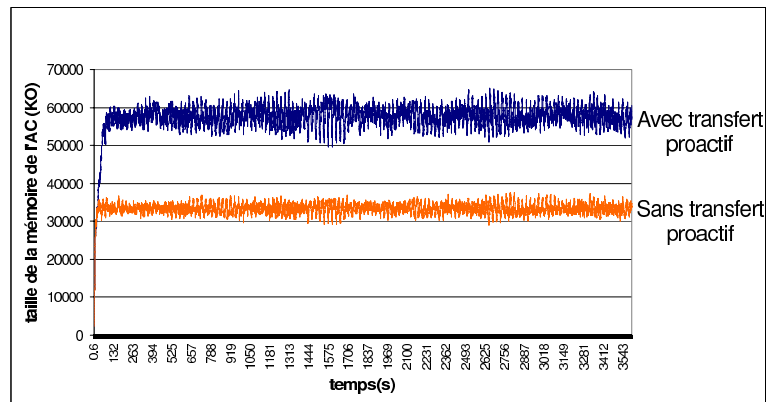


FIGURE 2.25 – Taille de la mémoire de stockage de l'AC

et ③, est exécutée une seule fois, au démarrage du réseau. Par conséquent, à mesure que le temps passe, chaque AC affine sa connaissance de la topologie du réseau et se configure dynamiquement pour construire ses deux tables.

Ce processus d'apprentissage a été évalué par le biais de simulations. Nous nous sommes intéressés à deux paramètres : le nombre d'interruptions de service et l'impact du processus d'apprentissage au niveau de mémoire de stockage de l'AC. La figure 2.25 illustre l'évolution de la taille de la mémoire de stockage de l'AC, avec et sans le processus d'apprentissage. L'exécution du protocole de découverte permet de maintenir la continuité du service lors du transfert inter-AC, et dans le même temps, la mémoire de stockage de l'AC conserve une taille tout à fait acceptable, considérant les capacités d'un terminal mobile.

2.5 Conclusion

Ce chapitre présente les principaux résultats obtenus dans le cadre des réseaux à couverture discontinue, appelés également *réseaux d'infostations*. Le principe est de déployer dans l'espace physique un ensemble de cellules de petites tailles offrant des débits élevés, et de les interconnecter via une infrastructure filaire. L'intérêt d'une telle approche réside dans la simplicité de déploiement. Mais cette simplicité implique une couverture radio discontinue. De nombreuses études ont apporté des solutions à l'intermittence de la connectivité pour des topologies et des services très spécifiques, par exemple pour un service de *streaming* en direction des terminaux le long d'une route. Dans ce cadre contraint, la mobilité des utilisateurs est connue, et il devient possible d'anticiper vers quelles bulles radio les données doivent être acheminées.

Notre principale contribution a été d'envisager les réseaux à couverture discontinue dans un cadre général : un déploiement à grande échelle, pour servir efficacement une densité importante d'utilisateurs dont les trajectoires dans le réseau mobile ne sont pas connues à l'avance. De plus, nous avons abordé le traitement des données dans les deux sens, montant et descendant.

Dans les réseaux d'infostations, l'approche système permettant de masquer la discontinuité radio peut être résumée de la manière suivante : introduire des mémoires de stockage temporaire (*i.e.* des caches) dans l'infrastructure réseau. Partant d'une analogie avec les systèmes multiprocesseurs, nous avons proposé une stratégie de distribution des caches entre le terminal et le réseau. Au niveau de ce dernier, le cache est hébergé par un équipement intermédiaire, le contrôleur d'accès ou AC, défini comme un *point d'attachement* entre les terminaux mobiles et le reste de l'infrastructure. Et comme dans une architecture multiprocesseurs, les caches distribués dans le réseau à couverture discontinue sont efficaces uniquement s'ils sont alimentés correctement. Il n'est pas possible de prévoir la trajectoire des utilisateurs dans le réseau, donc d'aiguiller les données *a priori*. Nous avons donc proposé une stratégie de distribution des données en direction des terminaux s'appuyant sur une discrimination des débits dans les bulles radio traversées. Nous avons également montré

que le cache de l'AC peut être utilisé pour exploiter efficacement les flux montants générés par les terminaux mobiles. Enfin, nous avons proposé des mécanismes permettant aux données en cache de migrer entre les ACs dès lors que la mobilité des utilisateurs l'impose.

Ces travaux ont servi de cadre à deux thèses, la première portant sur l'introduction de caches dans l'infrastructure [16], la seconde s'intéressant à la gestion des flux montants et au déploiement à grande échelle d'un réseau d'infostations [70].

Enfin, ces travaux sur les réseaux à couverture discontinue m'ont permis d'aborder la problématique du *couplage* de réseaux mobiles. Coupler des réseaux revient à offrir des mécanismes permettant de tirer parti de manière transparente des complémentarités fonctionnelles de ces réseaux. On parle alors de réseau « sans couture » ou *seamless network*. Ces mécanismes peuvent être envisagés à différents niveaux. Sans être exhaustif, on peut citer à titre d'exemple :

- Couplage au niveau radio, permettant ainsi le passage automatique et transparent d'une infrastructure sans fil à une autre, sans perte de connectivité. On parle alors de *handover vertical*, par opposition au *handover horizontal*, qui désigne le passage de l'utilisateur mobile entre deux cellules radio d'une même technologie.
- Couplage au niveau de la couche réseau, pour assurer une connectivité IP quelle que soit l'architecture de rattachement. On qualifie ce type de mécanisme de *soft handover*.

De manière connexe à nos travaux sur les réseaux discontinus, nous avons envisagé dans le cadre de nos travaux communs avec ALCATEL R&D un couplage système entre les bulles radio haut débit et une architecture cellulaire classique. Notre idée était qu'une demande d'accès à un contenu multimédia peut être effectuée à tout moment via le réseau cellulaire, ce dernier offrant une couverture continue. Et la délivrance du service est volontairement retardée, jusqu'à ce que le terminal destinataire rencontre une bulle radio haut débit. En repartant de l'analogie avec les architectures multiprocesseurs présentées dans la section 2.1.2, le réseau cellulaire est exploité comme un *bus de commandes* permettant d'effectuer des requêtes en direction de la mémoire centrale (le serveur de contenu). La délivrance des données (le service multimédia) est retardée jusqu'à ce que l'utilisateur rencontre une bulle radio haut débit. Ces travaux ont fait l'objet d'une publication commune [51] et d'un dépôt de brevet [65]. Ils n'ont pas été poussés plus avant. Mais cette problématique de couplage système, qui renvoie également à la notion de multi-attachement côté terminal, nous a semblé très pertinente. Elle constitue, dans un autre cadre que les réseaux d'infostations, l'axe directeur du dernier chapitre de ce mémoire.

COUPLAGE D'UN RÉSEAU DVB AVEC UN RÉSEAU CELLULAIRE 3G

Dans le domaine des réseaux cellulaires, de nombreuses recherches sont menées pour définir les futures architectures des réseaux de quatrième génération (4G). L'objectif est de déployer une infrastructure cellulaire large échelle, offrant des débits très élevés (plus de 100 Mb/s), et en mesure de servir une densité importante d'utilisateurs mobiles. Une approche possible consiste à considérer la future architecture 4G comme un réseau radio très haut débit, en s'appuyant sur des évolutions des infrastructures 3G (troisième génération) existantes. Deux normes sont actuellement en compétition pour les réseaux 4G mobiles : LTE¹ et WIMAX mobile.

Mais les recherches sur les réseaux 4G dépassent le cadre de ces évolutions normatives. Une voie possible pour atteindre les objectifs affichés en termes de débits et de densité d'utilisateurs est de s'appuyer sur un ensemble de technologies d'accès radio, utilisées en fonction de leur complémentarité [62]. Cette approche est illustrée par la figure 3.1. On retrouve les réseaux cellulaires 3G, LTE et WIMAX. Il est prévu que ces standards coexistent dans les années à venir, en fonction des zones de couverture. Mais d'autres technologies radio peuvent être envisagées pour offrir des zones de connectivité répondant aux exigences de la 4G.

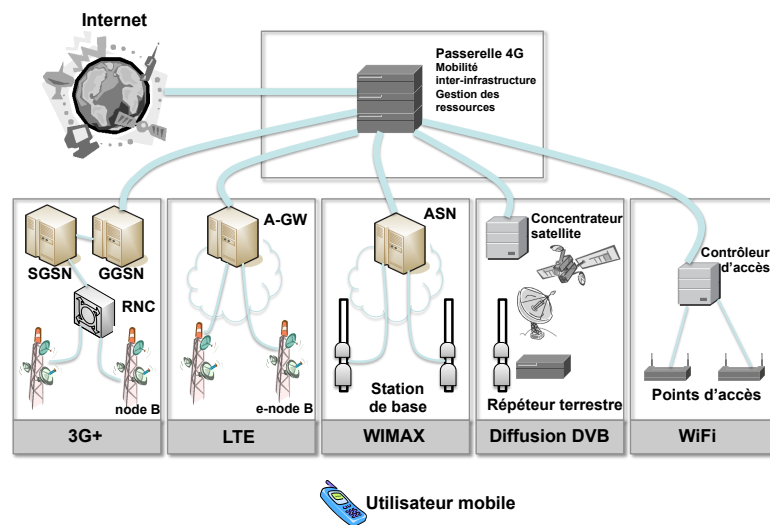


FIGURE 3.1 – Utilisation conjointe de différentes technologies dans les réseaux 4G

Utilisation des technologies WLANs. Les technologies sans fil courte portée comme WiFi permettent des échanges au delà de 100 Mb/s à proximité d'un point d'accès, et se rapprochent

1. LTE : *Long Term Evolution*

donc des critères de performances des futurs réseaux 4G. C'est ce constat qui a motivé nos travaux sur les réseaux à couverture discontinue, présentés dans le chapitre précédent (*cf.* chapitre 2).

Utilisation des technologies DVB. En Europe, DVB² est la famille de normes qui a été adoptée pour la diffusion de services de télévision numérique et multimédia. Initialement dédiés à la diffusion de programmes « traditionnels » sur les réseaux terrestres, câblés et satellitaires, les travaux de normalisation se sont poursuivis afin de permettre la réception de programmes TV sur des terminaux mobiles. Ainsi, la norme DVB-H [27] permet une diffusion via une infrastructure terrestre. Le standard DVB-SH [27], plus récent, s'appuie sur une infrastructure hybride, satellitaire et terrestre.

Un réseau DVB offre un mécanisme de diffusion « de masse » : il permet, par le biais d'un seul canal, de servir l'ensemble des utilisateurs se trouvant dans la zone couverte par le réseau. Dans sa déclinaison mobile, une telle infrastructure présente un réel intérêt pour mettre en œuvre les services multimédia visés par les réseaux 4G. Le problème est que le canal de diffusion est unidirectionnel, l'interface DVB d'un terminal ne pouvant que recevoir des données. Il n'est donc pas possible d'offrir de services interactifs aux utilisateurs, comme cela a été traité dans le cadre des réseaux à couverture discontinue. La complémentarité des réseaux, évoquée plus haut, prend tout son sens. En effet, une voie montante existe au niveau des réseaux cellulaires. On parle alors d'*utilisation couplée d'un réseau unidirectionnel de diffusion vers mobile et d'un réseau cellulaire bidirectionnel*. Ce principe est présenté dans la figure 3.2. Il s'agit d'une problématique relativement neuve dans le domaine des réseaux 4G.

J'avais pu aborder cette problématique du *couplage* entre deux réseaux offrant des fonctionnalités complémentaires, dans le cadre d'une étude menée en collaboration avec ALCATEL R&D. L'objectif était de proposer des mécanismes permettant une utilisation conjointe d'un réseau 2G et d'une bulle radio WiFi, pour la mise en œuvre d'un service de messages multimédia [51, 65]. Ces travaux ont été brièvement évoqués dans la conclusion du chapitre précédent (*cf.* section 2.5).

Afin de couvrir un spectre de recherches sur les réseaux 4G qui dépassent l'utilisation des communications sans fil courte portée, j'ai souhaité aborder la problématique du couplage d'un réseau DVB et d'un réseau cellulaire bidirectionnel. Tout comme dans les réseaux à couverture discontinue, notre démarche s'appuie sur une vision « système » de l'infrastructure pour traiter les aspects recherche.

Le reste du chapitre présente une synthèse de nos résultats. Après une analyse de deux scénarios de couplage, de « 3G vers DVB » et de « DVB vers 3G » présentée dans la section 3.1, nous en étudions la mise en œuvre dans les sections 3.2 et 3.3.

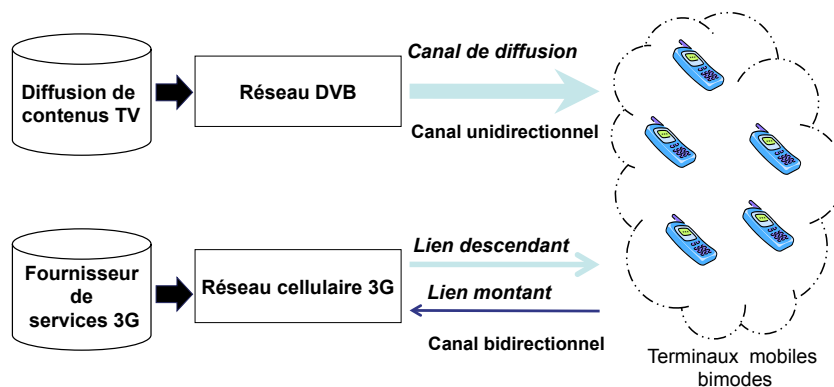


FIGURE 3.2 – Couplage d'un réseau DVB et d'un réseau cellulaire

2. DVB : *Digital Video Broadcasting*

3.1 Analyse de scénarios de couplage

Coupler deux réseaux revient souvent à envisager un scénario de service tirant parti des complémentarités fonctionnelles de ces réseaux. Ainsi, dans nos travaux antérieurs exploitant les interactions WiFi [51], la demande d'accès à un contenu multimédia peut être faite à tout moment par un utilisateur via un réseau cellulaire 2G. Par contre, la délivrance du service est volontairement retardée, jusqu'à ce que le terminal entre dans une bulle radio WLAN. L'idée est de combiner la continuité de couverture d'un réseau cellulaire avec le débit élevé offert ponctuellement par les technologies WiFi.

Nos travaux sur le couplage 4G ont été initiés en suivant une démarche analogue : nous avons commencé par analyser des scénarios d'usage réalistes pouvant justifier l'utilisation conjointe d'un réseau DVB et d'un réseau cellulaire. Pour plus de simplicité, nous désignerons le réseau cellulaire dans la suite de ce chapitre par *réseau 3G*. Un des enjeux identifiés dans le cadre d'une utilisation conjointe 3G et DVB est de tirer parti du couplage « dans les deux sens », ce qui peut se résumer simplement par deux questions :

- Un service sur un réseau 3G peut-il exploiter les capacités d'une infrastructure DVB mobile ?
- Un service sur un réseau DVB mobile peut-il exploiter les liens bidirectionnels d'un réseau cellulaire 3G ?

Du réseau 3G vers le réseau DVB. L'objectif est de faire basculer des flux 3G vers un canal de diffusion 3G. C'est le sens du couplage le plus évident : la voie montante du réseau cellulaire est utilisée par les usagers pour effectuer une demande de service, et le canal de diffusion DVB est exploité pour délivrer un flux de données. Il s'agit ici d'un principe de couplage proche de celui étudié précédemment dans un cadre 2G-WLAN : le réseau cellulaire fournit la voie montante pour acheminer les requêtes des utilisateurs, alors que les capacités de la voie descendante du second réseau sont exploitées pour améliorer la délivrance du service. Dans le cas d'un réseau couplé 3G-DVB, le canal descendant offre un mécanisme de diffusion : le flux n'est transmis qu'une fois, et est reçu par la totalité des utilisateurs se trouvant sous la couverture du réseau DVB.

Le couplage entre une infrastructure DVB et un réseau bidirectionnel n'est pas un problème nouveau. Il a été envisagé dès la normalisation des réseaux de diffusion terrestre comme DVB-T, pour l'envoi de logiciels par exemple. Ceci implique que les flux basculés depuis le réseau bidirectionnel (un réseau 3G dans notre cas) soient insérés de manière systématique dans le canal radio DVB [56]. Une solution à ce problème est d'allouer statiquement une partie du canal radio DVB à des nouveaux services, soit au détriment des programmes « classiques » de la télévision numérique, ou bien via l'obtention de nouvelles bandes de fréquences. Or, l'exploitation du spectre radio est soumise à une législation rigoureuse, et la distribution de nouvelles bandes de fréquences pour une technologie donnée est souvent problématique. On a pu le constater ces dernières années pour ce qui est de la planification des futurs réseaux 4G LTE.

Les bandes de fréquences dédiées au déploiement d'une infrastructure mobile sont donc des ressources « chères » qui doivent être exploitées de manière optimale. Dans ce contexte, le principe d'un scénario de couplage requérant l'allocation statique d'une large bande passante dans un réseau DVB mobile ne nous a pas semblé pertinent. Les opérateurs disposeront d'un canal radio de taille fixe, dédié en priorité à la diffusion de chaînes de télévision. Néanmoins, un des objectifs de ces opérateurs sera d'optimiser (*i.e.* minimiser) l'exploitation de la bande passante totale disponible, tout en offrant une qualité d'image satisfaisante pour la totalité des programmes. L'amélioration du rapport entre le volume du flux audio/vidéo et la qualité du flux audio / vidéo passe par l'utilisation de technique d'encodage à *débit variable*.

Considérons l'exemple de DVB-SH, un des standards les plus récents pour la TMP³. Cette norme prévoit l'utilisation de H.264 [58], pour l'encodage des programmes audio/vidéo transmis par les utilisateurs mobiles. Nous avons déjà abordé certains mécanismes propres à H.264 dans le cadre de nos travaux sur les réseaux à couverture discontinue [42].

3. TMP : Télévision Mobile Personnelle

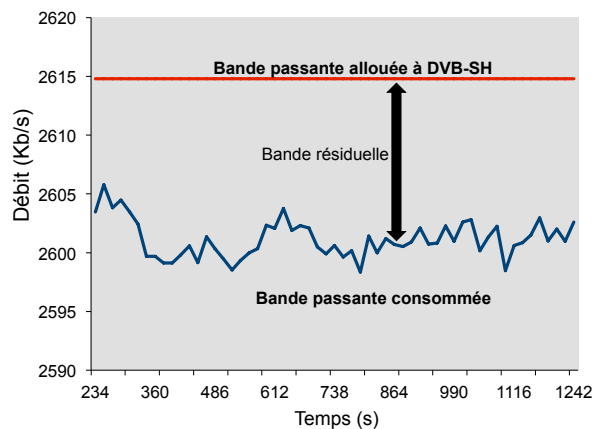


FIGURE 3.3 – Calcul de la bande résiduelle dans un réseau DVB-SH

En exploitant des captures réseau à la sortie d'un encodeur vidéo H.264 respectant les préconisations du standard DVB-SH, il nous a été possible de mesurer le trafic effectivement envoyé sur le canal radio DVB [41, 46]. Le scénario suivant a été considéré : le canal DVB dispose d'une bande passante totale de 2615 kb/s, et le flux est constitué de 10 canaux (programmes) distincts, diffusés chacun à un débit variable de 260 kb/s (débit comparable aux 256 kb/s envisagés pour le service de *streaming* dans les réseaux à couverture discontinue). Les mesures obtenues sont présentées dans la figure 3.3. Il existe donc une bande résiduelle dans une voie descendante DVB-SH d'une valeur moyenne de 13,5 kb/s, très faible en comparaison des débits offerts sur le réseau cellulaire 3G. Nous avons choisi de l'exploiter dans le cadre d'un premier scénario de couplage 3G-DVB. Nos travaux sont présentés dans la section 3.2 de ce chapitre.

Du réseau DVB vers le réseau 3G. De quelle manière les flux DVB peuvent-ils être enrichis via un réseau cellulaire 3G ? Ce couplage « dans l'autre sens » est peu abordé dans la littérature. Très souvent, le réseau 3G est considéré comme un réseau de « complément » capable de délivrer des contenus analogues à ceux diffusés via une infrastructure DVB. Un exemple classique consiste à basculer des programmes à faible audience sur le réseau 3G. Dans ce cas de figure, on peut parler d'utilisation *conjointe* des deux réseaux, mais pas d'exploitation *couplée*.

Les réseaux DVB mobiles sont conçus pour délivrer un flux (*i.e.* un ensemble de programmes) unique à l'ensemble des utilisateurs se trouvant dans la zone de couverture. Il s'agit là d'une contrainte forte. Afin d'amener davantage de souplesse dans la diffusion des contenus DVB, des travaux récents vise la personnalisation des contenus, afin de répondre aux besoins et centres d'intérêts des utilisateurs mobiles. Les approches utilisées consistent à multiplexer des contenus spécifiques avec des contenus classiques. Le flux transporté mêlent alors programmes généraux et programmes personnalisés. C'est pendant l'acheminement de ce flux que les programmes spécifiques sont filtrés et aiguillés vers différentes zones de couverture du réseau.

Par exemple, la société UDCAST propose une solution s'appuyant sur un équipement spécifique appelé *i-Splicer* dans l'infrastructure DVB [91], afin de transmettre des contenus régionaux sur DVB-H. Ce type d'approche est complexe à mettre en œuvre : elle nécessite de modifier l'infrastructure DVB en insérant des équipements spécifiques capable de filtrer les contenus spécialisés. Ce constat nous a amené à proposer une approche alternative : le réseau 3G peut être utilisé pour transmettre des contenus spécifiques en direction des terminaux, sans modifier l'architecture DVB existante. Et c'est au niveau des terminaux, doublement attachés aux réseaux DVB et 3G, que doit s'effectuer la substitution des flux. L'étude de ce second scénario de couplage est présentée dans la section 3.3 de ce chapitre.

Identification d'un point de connexion réaliste entre les deux réseaux. Comme nous le montrons dans la suite de ce chapitre, la mise en œuvre de services couplant les réseaux DVB et

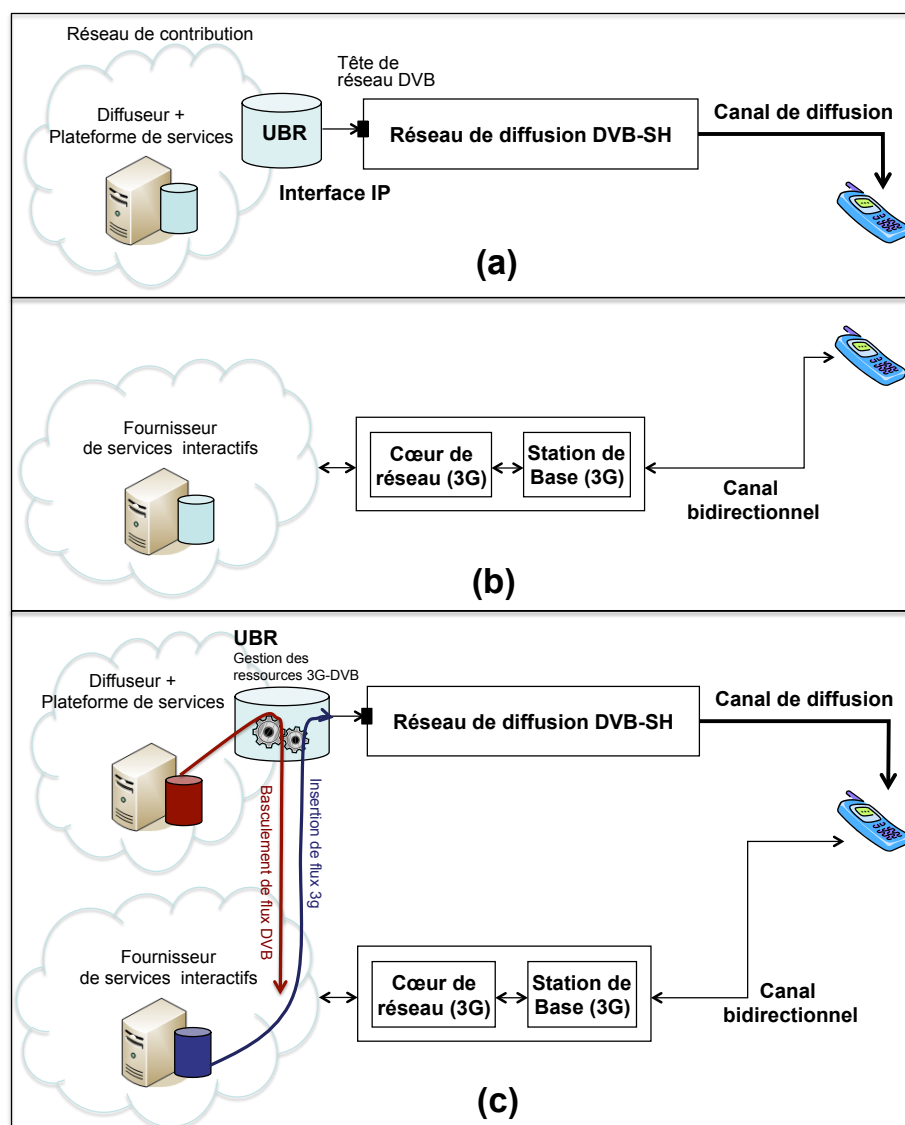


FIGURE 3.4 – Interconnexion entre les réseaux 3G et DVB

3G doit s'appuyer sur différents mécanismes système : stockage et distribution des données entre les deux réseaux, gestion de la signalisation entre les deux réseaux.

Mais préalablement à la définition de ces mécanismes se pose la question de leur implémentation dans une infrastructure 3G-DVB. Nous cherchons une réponse au problème suivant : dans quelle composant du réseau peut-on implanter des mécanismes de couplage, en prenant en compte les cadres normatifs stricts qui régissent le déploiement d'infrastructures 3G et DVB ? Nous avons rencontré un problème analogue dans le cadre de nos travaux sur les réseaux à couverture discontinue (*cf.* section 2.1.3). Le positionnement d'un cache intermédiaire, dans l'infrastructure filaire entre le serveur et le terminal renvoyait à l'identification d'un équipement réseau susceptible d'assurer la mise en cache des données envoyées par le serveur. Le contrôleur d'accès (AC), groupant un ensemble de point d'accès (AP) pouvait jouer ce rôle.

Nous avons abordé le problème de la même façon en ce qui concerne les mécanismes de couplage 3G-DVB. La figure 3.4 (a) donne une vue synthétique d'une infrastructure DVB mobile, qu'on désigne par *chaîne de diffusion DVB*. En amont, on trouve la plate-forme du diffuseur, en charge de la création des programmes destinés à être diffusés. Cette plate-forme est installée dans une

infrastructure IP, le réseau de contribution, qui assure le lien avec l'entrée (*i.e.* la tête) du réseau DVB.

L'équipement réalisant l'interface entre la plate-forme de service et la tête de réseau DVB est un routeur IP. Son rôle est d'assurer la diffusion des flux vers le reste de la chaîne DVB. Dans [68], ce routeur doit permettre également de diffuser des flux IP *unicast* en direction d'un réseau mobile, dans le cas par exemple de programmes temporaires ou à faible audience. Nous désignons cet équipement par la suite comme routeur de flux *unicast* et *broadcast* ou UBR⁴ [47].

La figure 3.4 (b) donne une vue synthétique d'une architecture 3G offrant un accès à une plate-forme distante. Le réseau cellulaire s'appuie sur un ensemble de stations de base connectées à un cœur de réseau. Ce dernier est lié à une infrastructure IP assurant l'accès à la plate-forme de services de l'opérateur 3G.

Au delà la tête de réseau, la chaîne DVB est un système fermé : les flux, après avoir été traités (multiplexés) par la plate-forme de services, sont envoyés vers des relais terrestre ou satellite. Insérer des données provenant du réseau 3G est impossible sans une modification fonctionnelle de la chaîne de diffusion DVB. Par contre, l'infrastructure en amont de la tête de réseau s'appuie sur une infrastructure IP ouverte, tout comme l'architecture 3G connecté au cœur de réseau. Partant de cette analyse, nous avons identifié l'UBR comme équipement d'interconnexion entre la chaîne DVB et la chaîne 3G. Ce principe est illustré par la figure 3.4 (c). Au niveau du réseau de diffusion DVB, l'UBR intervient dans la gestion et l'acheminement des flux vers le réseau de diffusion.

Au final, nos travaux sur le couplage 3G-DVB nous conduisent à étendre ces mécanismes au sein de l'UBR pour permettre

- l'insertion de flux *initialement 3G* vers la chaîne DVB (*cf.* section 3.2),
- et le basculement de flux *initialement DVB* vers le cœur de réseau 3G (*cf.* section 3.3).

Nous présentons nos principaux résultats dans la section suivante.

3.2 Basculement de contenus 3G vers un réseau DVB

Comme cela a été évoqué dans la section précédente, un réseau DVB mobile dispose d'une bande passante résiduelle dont la valeur est très faible. Par conséquent, la transmission de fichiers volumineux est exclue. En prenant en compte cette contrainte, nous avons envisagé la possibilité d'un service « retardé » : les données transmises ne sont mises à disposition de l'utilisateur qu'après avoir été reçues en totalité par le terminal. Le problème est alors d'obtenir un « ressenti » du service acceptable. Tout comme dans les réseaux à couverture discontinue où la notion de *user experience* avait été abordée, il s'agit d'offrir aux utilisateurs un temps d'accès au service raisonnable. On ne peut donc envisager que la transmission de données de tailles limitées (quelques Mb), comme une courte séquence vidéo que le terminal charge avant de la lire.

Un autre aspect à prendre en considération concerne les capacités de diffusion offerte par le réseau DVB. Transmettre ponctuellement un contenu via la bande résiduelle permet d'atteindre l'ensemble des utilisateurs mobiles se trouvant sous la couverture du réseau. Cela n'a de sens que si un nombre important de ces utilisateurs a souscrit au service proposé. Cette *popularité* du service peut être évaluée simplement en utilisant le réseau 3G. Les utilisateurs effectuent leur requête de service en utilisant la voie montante du réseau cellulaire. Et au delà d'un seuil de demandes, le fournisseur de service 3G cherche à basculer la délivrance des données vers le réseau DVB. Ce service présente deux avantages pour les opérateurs. En premier lieu, le réseau 3G est déchargé d'une partie de son trafic, on s'appuie sur le seul lien DVB pour transmettre les données, au lieu d'exploiter un ensemble de liaisons *unicast* dans le réseau 3G. Enfin, le nombre de destinataires du service justifie pleinement l'utilisation des capacités de diffusion du réseau DVB. Bien entendu, le temps de délivrance augmente fortement si les données sont basculées vers DVB. Un tel service doit donc être présenté comme « bas coût » pour l'utilisateur, en compensation de son caractère retardé.

Dans nos travaux, nous distinguons deux types de service exploitant la bande résiduelle DVB : les *services programmés*, et les *services non programmés*. Ils font l'objet des deux sections suivantes.

4. UBR : *Unicast Broadcast Router*

3.2.1 Mise en œuvre d'un service programmé

Un *service programmé* est associé à un flux de données pour lequel l'opérateur 3G garantit la délivrance avant une échéance fixée. Cette dernière est annoncée aux utilisateurs au moment de la demande de service. Si un flux 3G « programmé » est souscrit au delà d'un seuil fixé par l'opérateur, le fournisseur de service 3G envoie une requête à l'UBR pour demander l'insertion de ce flux dans la bande résiduelle DVB.

Dans la réalisation d'un mécanisme de couplage 3G vers DVB, l'UBR assure deux fonctions [46] :

1. Accepter ou refuser le flux, cette décision reposant sur la capacité de l'UBR à transmettre un **volume** de données (dépendant de la taille du flux) tout en respectant le **retard toléré** dans la délivrance (fonction de l'échéance annoncée par l'opérateur).
2. Insérer le flux dans la bande résiduelle DVB.

Ces deux fonctions sont illustrées par la figure 3.5.

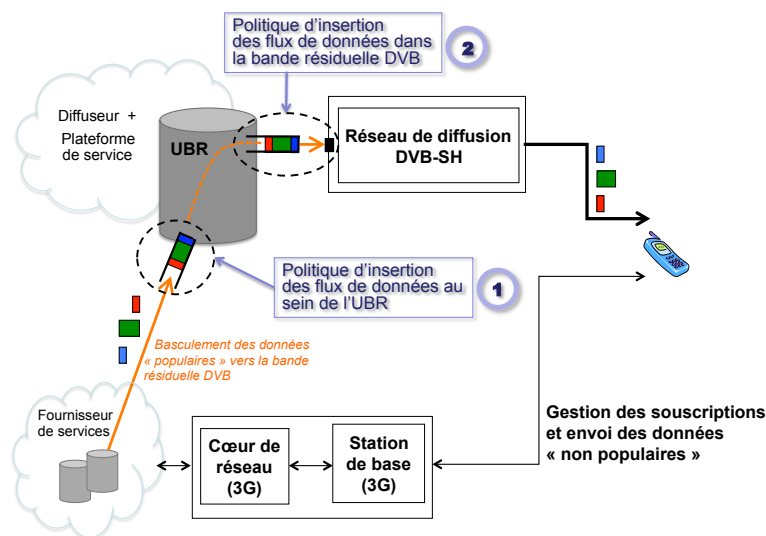


FIGURE 3.5 – Fonctions remplies par l'UBR dans un couplage 3G vers DVB

Politique d'insertion des flux dans la bande résiduelle DVB. Pour traiter ce problème, nous sommes partis de l'architecture d'une chaîne de diffusion DVB se trouvant en amont de la tête de réseau. La partie grisée de la figure 3.6 en donne une vue fonctionnelle. Trois types de flux peuvent être traités : les flux de signalisation (permettant de transmettre notamment le guide des programmes), les flux temps réels (les programmes TV), et les flux non temps réel ou contenus (permettant l'envoi de données comme des programmes de mise à jour, des jeux proposés au téléchargement).

Les programmes temps réel sont multiplexés dans un flux unique, tout comme les programmes non temps réel. Les trois flux résultants (avec la signalisation) sont ensuite soumis à un mécanisme d'ordonnancement qui traite en priorité les messages de signalisation. Sont ensuite transmis vers la tête de réseau le flot contenant les programmes TV multiplexés, la bande passante allouée devant permettre de vider la file d'attente associée à chaque cycle de fonctionnement de l'ordonnancement. Enfin, avec le plus faible niveau de priorité, la capacité restante est allouée à la transmission des données non temps-réel.

Cette analyse fonctionnelle des mécanismes d'insertion des flux DVB nous a permis de comprendre à quelle niveau le débit résiduel DVB est exploitable par l'UBR. C'est le second multiplexeur, en charge de constituer le flot des contenus non temps-réel, qui exploite la capacité restante

du canal DVB à chaque fin de cycle de l'ordonnanceur. L'UBR doit donc insérer les flux 3G programmés au niveau de ce second multiplexeur, avec le plus faible niveau de priorité. Ce principe est représenté sur le bas de la figure 3.6. Pour plus de détails concernant ce mécanisme, le lecteur peut se reporter à [41, 46].

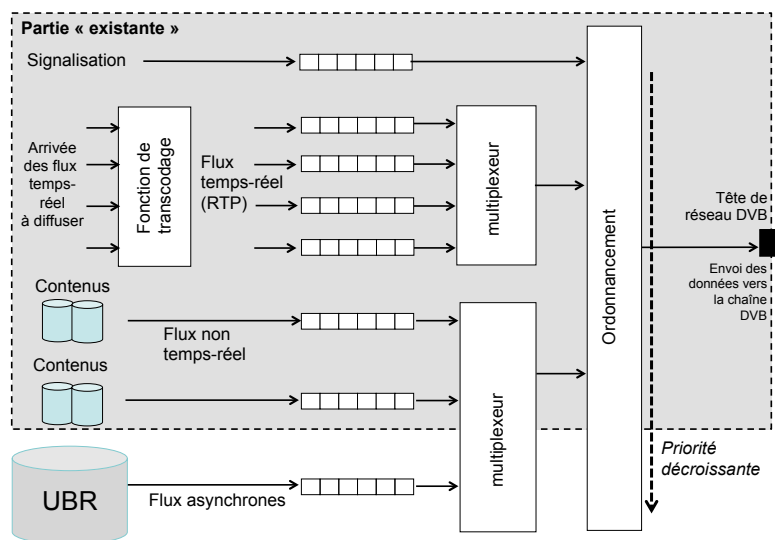


FIGURE 3.6 – Politique d'insertion des flux dans la tête de réseau DVB

Politique d'insertion des flux dans l'UBR. En prenant en compte la valeur moyenne de la bande résiduelle, l'UBR ne doit accepter un flux que si ce dernier peut être transmis avant l'échéance annoncée par l'opérateur. L'UBR maintient donc une liste ordonnée de demandes d'accès à la bande résiduelle DVB, ainsi que l'illustre la figure 3.7. Dès qu'une nouvelle demande est transmise par la plate-forme de services 3G (étape 1), un ordonnanceur cherche à insérer la requête dans la liste, en s'assurant que l'échéance associée au flux puisse être respectée (étape 2). Dès que la bande résiduelle est disponible, l'ordonnanceur extrait la requête en tête de liste (étape 3) et informe la plate-forme 3G que le flux associé à la requête peut être transmis vers la file asynchrone la moins prioritaire (*cf.* figure 3.6) connectée à l'UBR.

Le principe d'un classement des requêtes par l'UBR nous a amené à évaluer différentes politiques d'ordonnement. Nous ne les présentons pas en détails dans ce document, nos résultats à ce sujet sont publiés dans [45]. Deux paramètres permettent de caractériser la qualité du service délivrée via le réseau 3G-DVB : le nombre de requêtes de basculement refusées par l'UBR (qualité du service « perçue » par l'opérateur 3G) et le nombre de flux délivrés aux utilisateurs après expiration du délai (qualité du service « perçue » par l'utilisateur). Nous avons choisi de chercher avant tout à minimiser le second paramètre.

Pour chaque demande de basculement émanant de la plate-forme de services 3G, la politique d'ordonnement doit donc insérer la requête dans la file d'attente de l'UBR, en respectant les deux contraintes suivantes :

- Le placement dans la file permet à la requête d'être traitée à temps.
- Le placement dans la file de la nouvelle requête ne remet pas en cause le traitement du reste des demandes déjà acceptées.

Comme le montre nos simulations [45] effectuées via l'outil de simulation OPNET, cette politique, appelée *Nearest Deadline First*, offre le meilleur compromis entre les deux critères précédemment cités, (1) le nombre de requêtes refusées par l'UBR et (2) le nombre de délivrance de services hors délai.

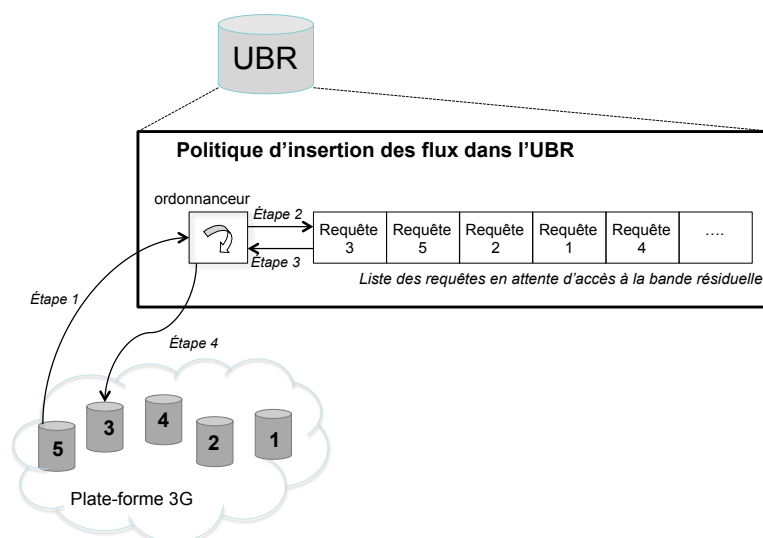


FIGURE 3.7 – Politique d'insertion des flux dans l'UBR

3.2.2 Mise en œuvre d'un service non programmé

Faute de pouvoir allouer statiquement une large bande de fréquences pour des services de données DVB mobile, il est possible d'envisager une « légère » extension de la bande résiduelle. Pour appuyer cette hypothèse, on peut se référer par exemple à un document publié par le Conseil Supérieur de l'audiovisuel. Dans [23], il est prévu qu'une bande passante de 120 Kb/s soit réservée pour un usage futur dans les réseaux DVB-SH, dans le cadre de services innovant pour la TMP. En complément des services programmés, nous avons souhaité étudier les perspectives offertes par une (faible) extension statique de la bande résiduelle. On parle dans la suite de cette section de *bande résiduelle étendue*.

Avec le débit offert, il devient possible de transmettre en permanence des flux temps réel, pour offrir par exemple un service de *streaming* de qualité moyenne. Tout comme les services programmés, cette proposition s'inscrit dans le cadre de services « bas coûts ». Par contre, la délivrance des contenus se fait de manière continue, sans échéance annoncée. Nous les qualifions donc de *services non programmés*.

Un service non programmé s'appuie donc sur un flux audio/vidéo diffusé en permanence par l'opérateur 3G, reçu soit via un ensemble de liens *unicast* 3G, soit par le biais de la diffusion DVB. Nous avons considéré qu'au démarrage du service, la communauté d'utilisateurs est restreinte, et le programme est transmis sur le réseau cellulaire. Mais tout comme dans le cas d'un service programmé, l'augmentation du nombre d'abonnés peut justifier le basculement vers le canal de diffusion DVB.

Un service non programmé soulève deux problèmes :

1. Comment insérer le flux de *streaming* dans la bande résiduelle étendue ?
2. Comment décider du basculement du flux du réseau 3G vers le réseau DVB, et inversement ?

Insertion d'un flux temps réel dans la bande résiduelle étendue. Nous avons traité ce problème en reprenant l'architecture fonctionnelle précédente (*cf.* figure 3.6). L'extension de la bande résiduelle s'inscrit en complément de la bande DVB. Elle devient donc disponible à partir du moment où la signalisation, les programmes TV et les contenus de données ont été transmis. Cette hiérarchie peut se traduire par l'insertion d'une nouvelle file d'attente, traitée « juste avant » la file des services programmés. La figure 3.8 résume ce principe.

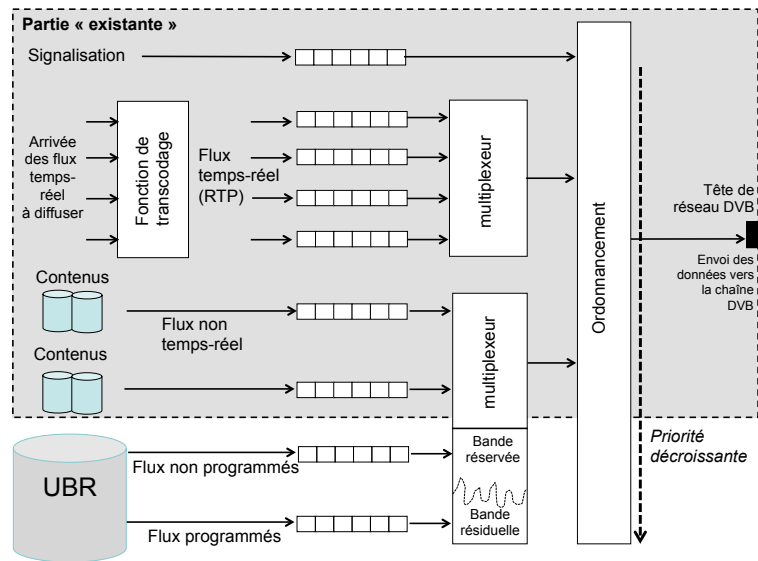


FIGURE 3.8 – Politique d'insertion des flux avec une bande résiduelle étendue

Stratégie de basculement du flux temps réel. Dans un service programmé, la décision de basculer vers le réseau DVB résulte du passage au dessus d'un seuil de souscription. Dans le contexte d'un service non programmé, le nombre de souscripteurs doit être observé en permanence. En effet, le passage au dessus d'une valeur fixée doit provoquer le passage du flux via le réseau DVB (le mécanisme de diffusion DVB est utile). Par contre, la diminution de ce même nombre doit amener le retour du flux vers l'infrastructure 3G (l'utilisation du mécanisme de diffusion ne se justifie plus).

L'utilisation d'un seuil fixe n'est pas efficace. Si le nombre d'utilisateurs oscille autour de la valeur de seuil, on peut se retrouver avec une série de basculements d'un réseau à l'autre. Afin de proposer un mécanisme de basculement stable et efficace d'un réseau à l'autre, nous sommes partis du principe qu'un contenu devient progressivement populaire. L'intérêt ou le désintérêt pour un service peut s'évaluer en prenant en compte les évolutions du nombre d'utilisateurs.

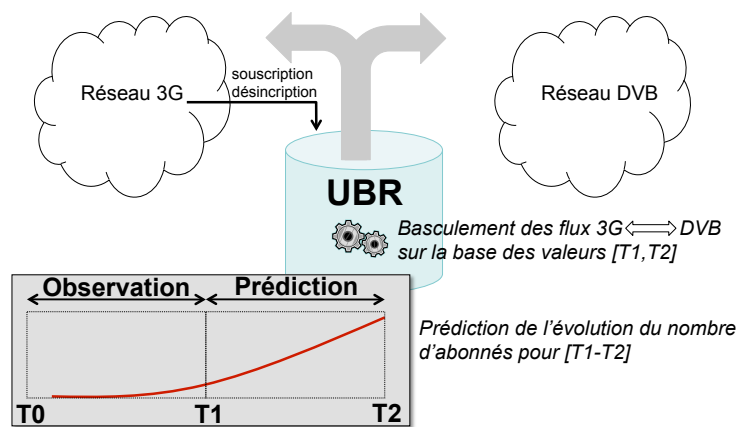


FIGURE 3.9 – Prédiction du nombre de souscripteurs à un service non programmé

L'approche proposée, implémentée au niveau de l'UBR, mesure le nombre de souscriptions sur un intervalle de temps $[T_0, T_1]$, et évalue sur cette base le nombre potentiel d'utilisateurs pour l'intervalle de temps suivant $[T_1, T_2]$. Cette évaluation s'appuie sur une régression par la méthode des moindres carrés. Cette méthode, très simple à mettre en œuvre, permet de calculer une courbe

approchant au mieux un nuage de points (le nombre d'abonnés en fonction du temps dans notre cas). Le principe de cette approche est illustré par la figure 3.9.

Le principe des services non programmés exploitant une bande résiduelle étendu a été validé par le biais de simulations, effectuées à l'aide de l'outil OPNET. L'ensemble des résultats obtenus est détaillé dans [47].

3.3 Délivrance de contenus personnalisés via le réseau 3G

En complément d'un couplage « 3G vers DVB », nous avons étudié le principe d'un couplage « DVB vers 3G ». Comme cela a été présenté en introduction de ce chapitre, ce type de couplage doit permettre de délivrer ponctuellement des contenus personnalisés. Un terminal bi-mode est en mesure de recevoir deux versions d'un même programme, une *version standard* diffusée en permanence via la chaîne DVB vers tous les utilisateurs, et une *version personnalisée* transmise ponctuellement sur le réseau 3G. Plus précisément, pendant une séquence temporelle fixée par la plate-forme DVB, les terminaux ayant souscrits au service de contenus personnalisés reçoivent des contenus spécifiques par le biais de connexions 3G *unicast*, contenus qu'ils substituent automatiquement au flux standard transmis via DVB.

Analyse du fonctionnement du service. La principale difficulté est de déterminer quels sont les mécanismes système nécessaires pour effectuer ces substitutions dans le cadre d'une architecture centrée sur l'UBR. L'étude de ce service nous a conduit à distinguer cinq étapes élémentaires dans son fonctionnement, représentées dans la figure 3.10.

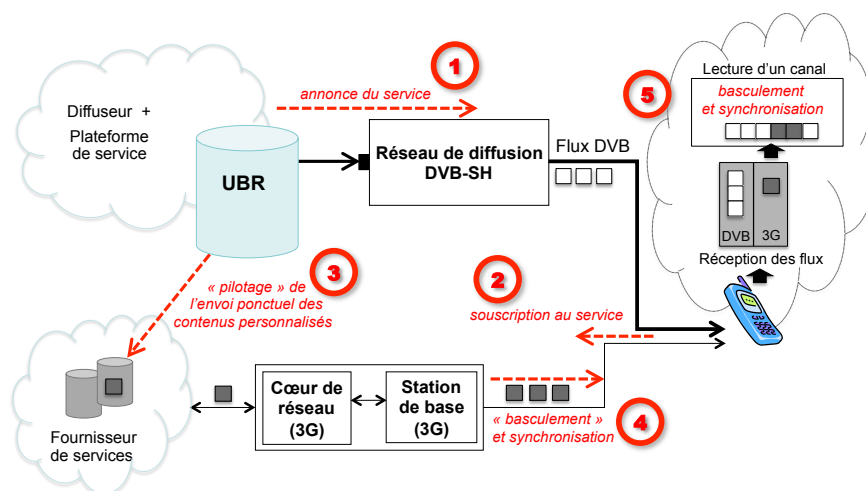


FIGURE 3.10 – Substitution des flux au niveau d'un terminal 3G-DVB

- 1. Annonce du service.** L'UBR annonce la disponibilité du service de contenu personnalisé aux utilisateurs du réseau DVB.
- 2. Souscription.** Les utilisateurs peuvent s'inscrire auprès de l'UBR.
- 3. Pilotage de l'envoi des contenus personnalisés.** Ces contenus sont stockés au niveau de la plate-forme de services 3G. C'est l'UBR, en charge de la gestion des mécanismes de couplage, qui pilote l'envoi des flux spécifiques vers les utilisateurs ayant souscrits au service. Pour cela, l'UBR a besoin de transmettre deux paramètres : une *date* indiquant le début de la lecture du contenu personnalisé par le terminal (*i.e.* le début de la substitution), et une *date* précisant la fin de la lecture (*i.e.* la fin de la substitution).

4. **Transmission des contenus personnalisés.** La plate-forme de service 3G envoie les flux, en prenant en compte les *dates* transmises par l'UBR.
5. **Substitution des flux.** Chaque terminal gère deux zones mémoire distinctes, l'une pour la réception du flux DVB, l'autre pour stocker les contenus personnalisés reçus via le lien 3G. Aux *dates* indiquées, le terminal doit être en mesure de basculer la lecture du programme d'une zone mémoire à l'autre.

Gestion de la signalisation et de la synchronisation. L'analyse du service montre qu'il est nécessaire de véhiculer différentes informations de signalisation entre le diffuseur DVB, l'UBR, la plate-forme de service 3G et le terminal mobile : annonce de disponibilité du service, envoi des informations de *date* de début et de *date* de fin.

Une plate-forme de diffusion DVB véhicule un ensemble de messages appelé ESG⁵. Leur objectif initial est de décrire les services et programmes offerts aux utilisateurs. Les informations fournies concernent à la fois la description des contenus, et le moyen de les acquérir. L'ESG offre beaucoup de possibilités : un message repose sur un ensemble de fragments XML [18], qui peuvent être combinés pour fournir des informations plus riches. De plus, les messages ESG ne sont pas uniquement destinés aux terminaux. Ils peuvent être pris en compte par des composants intermédiaires comme l'UBR, et rien n'empêche leur transmission à la fois sur la chaîne DVB et sur le réseau 3G.

Nous avons étudié la transmission des informations de signalisation du service de contenu personnalisé via l'ESG. Notre analyse, avec l'ensemble des messages ESG spécifiés se trouve dans [44]. C'est donc un message ESG qui permet à l'UBR de transmettre les informations de synchronisation (les *dates*) à la plate-forme de services (étape 3 sur la figure 3.10). L'efficacité du service repose sur la précision du basculement du flux DVB vers le flux 3G par le terminal. Le problème est donc d'identifier le type d'information dont le terminal a besoin pour déclencher la substitution.

Une synchronisation s'appuyant sur une référence temporelle nécessite l'existence d'une horloge commune entre les deux plates-formes de services, DVB et 3G. Nous avons retenu une approche plus simple, et sans doute plus réaliste. La structure du flux DVB est utilisée comme référence, et nous avons identifié deux marqueurs (correspondants à la *date* de début et à la *date* de fin) renvoyant à deux positions dans le flux DVB. Plus précisément, la structure d'un flux transporté par une chaîne de diffusion DVB est définie par la norme MPEG TS⁶. De manière assez simple, un flux de transport TS peut être vue comme un ensemble de paquets pouvant correspondre à des programmes différents. Dans l'entête de chacun de ses paquets, un champ spécifique, le PTS⁷, est un marqueur temporel indiquant à quel moment une image doit être présentée (lue). Autrement dit, l'ensemble des PTS d'un flux permet de retrouver l'axe des temps pendant la lecture d'un programme. C'est donc cette information que nous avons sélectionnée pour exprimer les *dates* de basculement au niveau du terminal [44, 43].

Analyse de l'impact sur l'infrastructure 3G. Un aspect essentiel des futures architectures 4G réside dans leur capacité à servir efficacement des densités élevées d'utilisateurs. Dans le cas du service étudié dans cette section, un nombre important de souscripteurs dans une cellule implique d'utiliser un nombre équivalent de liens *unicast* 3G dans cette cellule. Nous avons pu observer ce mécanisme par le biais de simulations sous OPNET : les ressources mobilisées par le service de contenu personnalisé augmentent de manière quasi linéaire avec le nombre de souscripteurs de la cellule.

Or, les ressources radio déployées au niveau d'une station de base 3G sont limitées. La difficulté pour un opérateur de réseau cellulaire est d'allouer une partie de ces ressources pour de nouveaux services, tout en continuant à servir efficacement les flux *unicast* existants. Très souvent, les contenus personnalisés correspondent à des décrochages régionaux. Et la majorité des utilisateurs intéressés par ce type de décrochage se trouve concentrée dans une même zone géographique, donc

5. ESG : *Electronic Service Guide*

6. TS : *Transport Stream*

7. PTS : *Presentation Time Stamp*

rattachée au même groupe de cellules du réseau. La transmission des programmes personnalisés risque d'entraîner une saturation au niveau des stations de base associées.

De cette analyse est née notre idée de coupler le service étudié avec un réseau MBMS⁸. MBMS est une évolution possible d'un réseau 3G, afin d'offrir des capacités de *broadcast* et de *multicast*. Le principe est de permettre aux abonnés mobiles de s'inscrire à un groupe *multicast* dans leur cellule de rattachement. Les contenus sont alors diffusés une seule fois en direction du groupe d'abonnés. Ce type d'approche, courante dans les réseaux filaires, permet d'économiser la consommation de bande passante dans l'infrastructure, à la fois dans le cœur de réseau 3G, et également dans la cellule attachée à la station de base. Ce principe est résumé dans la figure 3.11. La partie gauche présente l'envoi d'un contenu en utilisant un ensemble de liens *unicast*. A droite, on peut observer l'envoi du même contenu exploitant les services offerts par les composants MBMS.

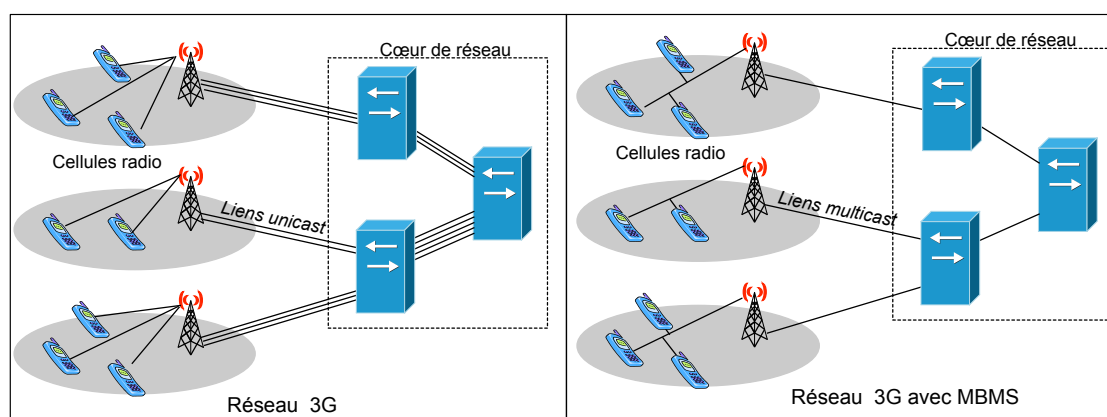


FIGURE 3.11 – Apport d'un réseau MBMS

Dans une architecture cellulaire « 3G + MBMS », les programmes personnalisés ne sont plus transmis via un ensemble de liens *unicast*, mais par le biais de flux MBMS. L'impact attendu sur la charge des cellules concentrant un nombre important d'abonnés au service a été mesuré via un ensemble de simulations, dont les résultats sont présentés dans [41, 44]. L'environnement de simulation, développé sous OPNET, a été étendu pour supporter les services MBMS. A titre d'exemple, la figure 3.12 présente l'évolution de la charge dans une cellule avec 10 % de la population ayant souscrit au service de contenu personnalisé, puis avec 90 %. On peut observer l'impact positif de MBMS : les deux courbes sont pratiquement identiques, la charge induite par la délivrance du contenu personnalisé reste stable quel que soit le nombre d'abonnés au service se trouvant dans la cellule.

3.4 Conclusion

Ce chapitre présente une synthèse de nos derniers travaux sur les réseaux 4G. Ces derniers doivent respecter de nombreuses contraintes : couverture étendue, des débits élevés, support d'une densité importante d'utilisateurs mobiles, personnalisation des services. Un aspect essentiel de ces nouveaux réseaux est qu'ils ne marquent pas une rupture technologique franche avec les générations précédentes des réseaux cellulaires étendus. Ils cherchent à exploiter la convergence et la coopération entre différents réseaux sans fil existants. Après avoir abordé cette problématique au travers des réseaux à couverture discontinue, nous avons inscrit nos travaux dans le cadre d'un couplage entre un réseau de diffusion de télévision mobile personnel DVB et un réseau cellulaire 3G. Le premier permet la diffusion unidirectionnelle vers un nombre important d'utilisateurs. Le second offre la possibilité de cibler plus finement les utilisateurs, et autorise l'interactivité grâce à la bidirectionnalité des liens radio.

8. MBMS : *Multimedia Broadcast Multicast Service*

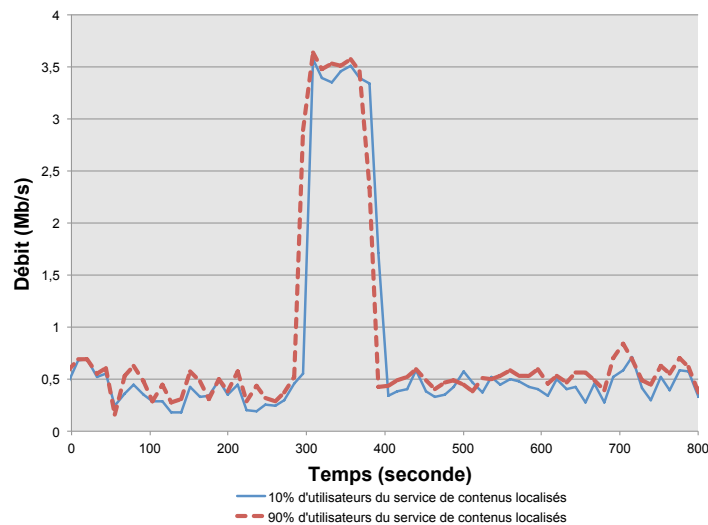


FIGURE 3.12 – Impact du service sur le réseau via l'utilisation de MBMS

Un aspect important est d'avoir recherché un couplage **mutuel** des deux infrastructures dans le cadre de nouveaux services, et donc d'avoir étudié des solutions système pour les « deux sens » du couplage. Nous nous sommes attachés à montrer sur quelle architecture, sur quels mécanismes de gestion de données et sur quels types d'échanges il était possible d'établir une collaboration entre les deux réseaux.

Une des difficultés était d'inscrire notre approche dans un cadre normatif très restrictif. Une chaîne DVB est un système fermé, dans lequel il est complexe d'ajouter de nouveaux équipements. Cette contrainte était beaucoup moins forte dans le cadre de nos travaux sur les réseaux à couverture discontinue. A titre d'exemple, j'avais initialement envisagé d'introduire des mécanismes de cache au sein de la chaîne DVB, dans le but d'exploiter plus efficacement le codage hiérarchique offert par la norme H.264. Cette approche s'est avérée finalement irréaliste, la chaîne DVB ne permettant pas d'intégrer les composants nécessaires.

Notre implication dans le projet TVMSL⁹ nous a fortement aidé à appréhender les contraintes fortes des architectures existantes des réseaux DVB mobiles. Projet de trois ans financé par l'OSEO, son objectif était de développer un standard DVB mobile hybride, en s'appuyant sur une transmission terrestre et satellitaire. Il regroupait des partenaires industriels et académiques. Une partie de nos travaux sur le couplage DVB-3G a pu être discutée et présentée dans le cadre de ce projet, notamment la gestion de la signalisation et les fonctions assurées par le routeur UBR.

Nos travaux se sont concrétisés dans le cadre d'une thèse [41], et l'ensemble des résultats présentés dans ce chapitre a été publié [46, 47, 44, 45].

9. TVMSL : TV Mobile Sans Limite

BILAN ET PERSPECTIVES

4.1 Bilan des travaux

Nous avons présenté dans ce document nos travaux concernant l'exploitation d'approches système dans différents réseaux sans fil. Cherchant à inscrire notre thématique de recherche dans le cadre des systèmes mobiles et distribués, nous nous sommes intéressés aux perspectives offertes par les réseaux locaux sans fil. A l'opposé de la complexité de déploiement d'une infrastructure cellulaire, les interactions sans fil courte portée peuvent être utilisées de manière très simple, sans infrastructure. Ainsi, elles permettent à des calculateurs proches d'échanger automatiquement des informations. Nous avons proposé des supports système prenant en compte la volatilité des communications sans fil, et permettant de développer des applications tirant spontanément parti de la proximité physique des nœuds mobiles.

Nous avons poursuivi ces travaux dans le cadre d'autres familles de réseaux sans fil. Ainsi, nous sommes intéressés aux réseaux à *couverture discontinue*. La technologie support est la même que celle de notre première étude. Simplement, les communications entre les nœuds mobiles ne sont plus directes, mais passent par une borne fixe, un point d'accès ou AP. Cet AP définit une bulle radio de taille limitée. C'est l'interconnexion de ces bulles, sans souci de continuité de la couverture radio, qui permet d'envisager un réseau étendu et simple à déployer. Dans ce cadre, les mécanismes système étudiés permettent de masquer l'intermittence de la connectivité, et autorisent le support d'applications exploitant les flux montants et descendants dans le réseau.

Enfin, nous avons abordé le problème du *couplage système* de deux architectures sans fil hétérogènes. Un tel couplage présente des objectifs comparables à ceux des réseaux à couverture discontinue : offrir des nouveaux services sur une couverture large, à des densités importantes d'utilisateurs mobiles. Cette dernière phase de nos travaux a débouché sur la définition de mécanismes permettant de coupler au sein d'un même service des propriétés fonctionnelles de chacune de ces deux infrastructures.

4.1.1 Mécanismes système pour les interactions sans fil courte portée

Les mécanismes étudiés ici doivent permettre la mise en œuvre de Systèmes d'Information Spontanés ou S.I.S. Le principe est très simple : il y a génération spontanée d'un système d'information dès lors que deux terminaux mobiles au moins se retrouvent suffisamment proches pour établir une communication, et échanger de manière implicite ou explicite des informations. La mise en place d'un S.I.S pose deux problèmes : la détection de la composition du S.I.S, et l'accès aux informations distribuées (appelées *l'espace visible* du S.I.S) sur les nœuds mobiles composant le S.I.S.

Le premier problème concerne la détection de présence des entités mobiles se trouvant à portée de communication. Il s'agit pour un terminal de reproduire une représentation informatique « fidèle » de son voisinage radio. Les mécanismes de découverte dans le domaine des réseaux reposent le plus souvent sur des envois périodiques de messages d'annonce. Le problème est de maintenir un taux de détection acceptable des voisins au sein du S.I.S, tout en minimisant la bande passante

consommée par les messages d'annonce. En prenant en compte la portée radio et la vitesse de déplacement des terminaux, l'approche proposée permet de calculer automatiquement une fréquence d'annonce pour chaque terminal.

Le second problème se concentre sur la découverte des informations et des interactions au sein du voisinage d'un terminal. Cela revient à proposer des mécanismes de *contrôle* et d'*adressage* au sein de l'espace visible d'un S.I.S.

Dans un premier temps, l'espace visible a été traité comme une base de données, appelée *base de données de proximité*. Les mécanismes permettant à une application d'exploiter sa base reposent sur une interrogation continue de l'espace visible. Également exploitées dans les réseaux de capteurs, les requêtes continues permettent d'interroger des données dynamiques. Elles offrent les mécanismes de contrôle et d'adressage nécessaires à la gestion de l'espace visible. Dans nos propositions, les requêtes continues ont été étendues, afin de prendre en compte la composition dynamique d'un S.I.S.

Nous avons ensuite considéré une seconde classe d'applications, en complément des bases de données de proximité, appelée *Web de proximité*. Cette approche permet à un terminal mobile de récupérer automatiquement des informations « pertinentes » au sein de son espace visible. Pour juger de cette « pertinence », nous avons proposé un ensemble de mécanismes exploitant des mots clés associés aux documents stockés localement par un utilisateur dans son terminal. À partir de ces mots clés, il devient possible d'extraire automatiquement les centres d'intérêt d'un utilisateur, et d'organiser thématiquement les documents, l'objectif final étant d'accélérer les échanges des informations pertinentes extraites de l'espace visible.

L'ensemble de ces travaux a été validé par la construction de logiciels prototypes. La découverte du voisinage physique a été implémentée sous forme d'un service système au sein d'un terminal mobile, suivant un schéma client/serveur, le rendant disponible pour toutes les applications exploitant le principe des S.I.S. Les bases de données de proximité ont été expérimentées au travers de l'architecture PERSEND (*PERsistent SEnsing for Neighbouring data*), implémentée au dessus d'une base de données relationnelles embarquée sur chaque nœud mobile.

4.1.2 Mécanismes système pour les réseaux mobiles 4G

L'exploitation d'un réseau à couverture discontinue implique souvent des déploiements contrôlés des bulles radio dans l'espace, par exemple le long d'une route. Dans ce cadre, il devient possible de précharger des données en fonction de l'évolution de la trajectoire d'un utilisateur. Nous avons cherché à aller plus loin, en exploitant des interactions sans fil courte portée, sans imposer de topologie particulière pour disposer les points d'accès dans l'espace, et sans contraindre la mobilité des utilisateurs. Trois problèmes ont été traités : mettre en place une gestion efficace des flux descendants en direction des terminaux mobiles, permettre à ces mêmes terminaux d'exploiter les flux montants, et enfin identifier les composants réseau permettant le déploiement à grande échelle d'une architecture discontinue.

En ce qui concerne le traitement des flux descendants, nous sommes partis d'une analogie avec les systèmes multiprocesseurs. Nous avons montré qu'une hiérarchie de caches distribuée sur plusieurs niveaux dans l'infrastructure permet de s'affranchir de la discontinuité de la couverture radio. L'architecture système proposée s'articule autour d'un équipement central appelé le contrôleur d'accès ou AC. Ce dernier joue le rôle de *point d'attachement* entre le serveur et les terminaux mobiles en attente de données. C'est l'AC qui assure la politique d'accès anticipés aux données. Dans un système de fichiers, l'anticipation des accès repose sur l'analyse des accès précédents, la prédiction des accès ultérieurs, et le préchargement des données prédites. Dans un réseau à couverture discontinue, une approche analogue est difficile à envisager dès lors que la mobilité des utilisateurs est non contrainte. Nous avons donc proposé une politique d'accès privilégiant les zones des bulles radio offrant les débits les plus élevés, le « reste » des bulles étant utilisé uniquement pour détecter l'entrée des terminaux sous couverture radio. Pour ce qui est du problème des flux montants, nous avons montré que le cache de l'AC peut être utilisé pour exploiter efficacement les données générées à partir des terminaux mobiles. Et au final, des mécanismes ont été proposés pour

permettre la migration des données entre ACs, dès lors que la mobilité des utilisateurs l'impose. Une telle migration permet au réseau de s'appuyer sur plusieurs ACs sans perte de données pour les utilisateurs. Ce qui autorise un déploiement à grande échelle d'un réseau à couverture discontinue.

Il nous est apparu que ces travaux pouvaient s'inscrire dans un thème de recherche plus large dans le domaine des réseaux sans fil, celui des réseaux cellulaires de quatrième génération ou *réseau 4G*. L'objectif général est d'offrir une infrastructure cellulaire grande échelle, en s'appuyant sur un ensemble de technologies d'accès radio, utilisées en fonction de leurs complémentarités. On parle de *couplage de réseaux*. Parmi les différentes formes de couplage envisagées, nous avons travaillé sur la combinaison d'un réseau de diffusion *unidirectionnel* de télévision mobile, et d'un réseau cellulaire *bidirectionnel*. Le premier propose une diffusion massive vers les utilisateurs se trouvant sous couverture, alors que le second offre la possibilité de cibler plus finement des groupes d'utilisateurs, notamment via l'utilisation d'une voie montante.

Tout comme dans les réseaux à couverture discontinue, nous avons proposé l'intégration de mécanismes de distribution de données au sein de l'infrastructure réseau. Un aspect important de notre contribution est d'avoir recherché un couplage *mutuel* des deux infrastructures dans le cadre de nouveaux services, c'est à dire un couplage dans « les deux sens ». La principale difficulté a été d'inscrire nos propositions dans un cadre normatif beaucoup plus restrictif que celui des S.I.S ou des réseaux à couverture discontinue. Notre proposition s'appuie sur un *point* de connexion « réaliste » entre les deux réseaux couplés. Ce point, un routeur en l'occurrence, réalise l'interface entre la plate-forme de diffusion et le cœur du réseau cellulaire. Tout comme un AC dans un réseau à couverture discontinue, il héberge des politiques de distribution de données. Ces politiques ont pour but d'aiguiller les flux automatiquement entre les deux infrastructures que l'on cherche à coupler.

4.2 Perspectives

Nous dressons ici quelques perspectives des travaux présentés dans ce mémoire.

4.2.1 Réseaux 4G : vers une gestion du multi-attachement

La problématique du couplage n'a été abordée que de manière partielle dans nos travaux. Dans le cadre des réseaux à couverture discontinue, l'infrastructure cellulaire a été envisagée comme un simple réseau de repli, exploité uniquement si aucune bulle radio n'est accessible. Et pour ce qui est de l'utilisation conjointe d'un réseau de diffusion et d'un réseau cellulaire, l'attachement aux deux architectures est utilisé seulement en alternance : on bascule d'un réseau DVB vers un réseau 3G, et inversement, en fonction du service visé.

Une perspective intéressante pourrait être de permettre à un terminal d'être attaché **simultanément** à plusieurs réseaux, afin d'exploiter des mécanismes de tolérance aux fautes et/ou de répartition de charge. On parle alors de *multihoming*.

A titre d'exemple, nous avons envisagé la distribution d'un flux multimédia via le réseau DVB, accompagné d'un flux *secondaire* transitant par le réseau 3G. L'intérêt de ce second flux est de renforcer ponctuellement la qualité du flux traduit au niveau du terminal. Des mécanismes de codage hiérarchique comme H.264 SVC autorisent ce type d'approche. Mais bien entendu, cette utilisation combinée des deux réseaux nécessite des extensions des politiques de distribution de données que nous avons proposées. Dans cet exemple, le renforcement d'un flux DVB par le biais d'un flux 3G doit s'appuyer sur des mécanismes de synchronisation que nos solutions ne prennent pas en compte.

De plus, cette utilisation conjointe et simultanée de deux réseaux mobiles doit s'accompagner d'une réflexion autour des protocoles à utiliser pour aiguiller les flux de données. Nous avons étudié dans le cadre des réseaux à couverture discontinue l'utilisation du protocole SCTP pour l'envoi des lignes de caches à partir de l'AC. Si on envisage des connexions simultanées de l'AC avec un réseau

discontinu et un réseau cellulaire, le protocole transportant des lignes de cache doit supporter le multi-attachement. Des extensions de SCTP proposent le support d'*associations* multiples. L'utilisation de ces extensions en conjonction avec une politique de distribution de données de l'AC vers deux réseaux différents pourrait être envisagée.

4.2.2 Vers une exploitation d'un contexte local enrichi

Dans le domaine de l'ubiquité numérique, de nombreuses solutions reposent sur un modèle logique de l'environnement. Par exemple, un système de navigation routière dispose d'une base de données cartographiques, afin de pouvoir guider l'utilisateur en fonction de sa position dans l'espace. Un système d'information est chargé d'interpréter en permanence ce modèle, et d'effectuer les traitements associés, comme par exemple d'avertir le conducteur de la présence d'une zone dangereuse. On peut parler d'approche *logique*.

Une autre approche consiste à définir un système dans lequel les informations sont construites automatiquement à partir d'interactions entre entités du monde réel. Le principe consiste à exploiter des propriétés existant dans l'espace physique, comme la proximité physique, pour en déduire automatiquement un système d'information. Le couplage entre le monde réel et le système d'information s'effectue de manière plus immédiate qu'avec une approche logique. On peut alors parler d'approche *physique*.

Les applications construites via cette seconde approche reposent sur des propriétés très simples comme la disposition physique d'objets et de calculateurs dans un même espace, ou la proximité physique des personnes etc. C'est ce principe que nous avons exploité dans le cadre des S.I.S. Les champs applicatifs ouverts, l'extensibilité, la simplicité de déploiement font que les architectures reposant sur une approche physique sont très prometteuses. Elles l'étaient déjà en 2004, au moment où nos travaux sur S.I.S se terminaient. Nous n'avons pas poussé plus loin nos propositions pour l'exploitation d'un contexte local fondé sur la proximité physique de nœuds mobiles, principalement pour nous concentrer sur les réseaux à couverture discontinue. Mais c'est aussi des raisons d'ordre technologique qui nous avaient poussées à nous tourner vers des architectures plus étendues que les S.I.S. Pour exploiter nos mécanismes de découverte et de construction du voisinage physique, nous avons besoin de connaître précisément les caractéristiques cinématiques d'une entité mobile. Les limites des terminaux ne nous permettaient pas d'expérimenter complètement nos solutions. Et plus généralement, nous souhaitions traiter le terminal comme un nœud de capteurs, capables de sonder en permanence de manière fine son environnement. Notre motivation était d'exploiter un contexte local plus riche que celui découlant uniquement d'interactions de proximité entre des entités voisines. Mais le support technologique d'une telle approche ne nous paraissait pas réaliste à court terme.

Ce cadre technologique a largement évolué, et de nombreux verrous sont maintenant levés. Les terminaux les plus récents embarquent de nombreuses fonctions de capture (GPS, accéléromètre, gyroscope) et des interfaces de communications courte portée (WiFi, Bluetooth, NFC¹). Il est maintenant possible d'embarquer des capteurs sur un microcontrôleur basse consommation, et équipé d'une interface sans fil. Considérant ce panel technologique, nos travaux s'inscrivent à nouveau depuis 2010 dans le domaine des supports système pour l'ubiquité numérique. Plus précisément, nous étudions de nouvelles applications exploitant un contexte local. A titre d'exemple, nous travaillons sur des mécanismes de localisation *indoor* [36, 37]. Nous explorons également de nouvelles classes d'application dans le cadre de la *maison intelligente*. L'objectif est de proposer une organisation spatiale des informations au sein de l'habitat, dans le but de piloter des éléments électriques, en prenant en compte de manière fine le contexte utilisateur dans la mise en œuvre de mécanismes de régulation énergétique [25, 24, 59].

1. NFC : *Near Field Communication*

Bibliographie

- [1] A framework for discrete-event modelling and simulation. <http://sourceforge.net/projects/desmoj/>.
- [2] The network simulator - ns-2. http://nsnam.isi.edu/nsnam/index.php/Main_Page.
- [3] Control and provisioning of wireless access points (capwap). <http://www.ietf.org/html.charters/capwap-charter.html>, [06 July 2007].
- [4] Ieee trial-use recommended practice for multi-vendor access point interoperability via an inter-access point protocol across distribution systems supporting ieee 802.11(tm) operation. Technical report, IEEE Standard 802.11f, 2003.
- [5] Capacity coverage and deployment considerations for ieee 802.11g. White paper, Cisco Systems, 2005.
- [6] Ieee standard for local and metropolitan area networks part 16 : Air interface for fixed and mobile broadband wireless access systems amendment 2 : Physical and medium access control layers for combined fixed and mobile operation in licensed bands and corrigendum 1. 2006.
- [7] Ieee standard for local and metropolitan area networks part 16 : Air interface for fixed and mobile broadband wireless access systems amendment 2 : Physical and medium access control layers for combined fixed and mobile operation in licensed bands and corrigendum 1. 2006.
- [8] A-S. Adde. Réalisation d'une application de démonstration pour les systèmes d'information spontanés. Master's thesis, Ecole d'ingénieur de l'IFSIC, Université de Rennes 1, Juin 2003.
- [9] A-S. Adde. Délivrance continue de données sous une architecture discontinue. Master's thesis, DEA d'informatique de l'Université de Rennes 1, Juin 2004.
- [10] M. Avila. Mise en œuvre d'une base de données ad-hoc en s'appuyant sur les systèmes d'information spontanés. Master's thesis, Mémoire de l'Institut d'Informatique d'Entreprise (IIE), Juin 2002.
- [11] M. Banâtre, C. Bryce, P. Couderc, and F. Weis. *Informatique diffuse : des concepts à la réalité*. Hermès, 2007.
- [12] M. Banâtre, P. Couderc, and F. Weis. Spontaneous communications. In *Proceeding of IST Mobile Communications Summit*, Galway, Ireland, October 2000.
- [13] M. Banâtre and F. Weis. Système d'information spontané (sis) : Problématique et premiers éléments de solutions. Technical Report 1222, IRISA, December 1998.
- [14] M. Banâtre and F. Weis. Sis : a new paradigm for mobile computer systems. In *Proceedings of the Information Society Technologies Conference (IST'99)*, Helsinki, Finland, November 1999.
- [15] M. Banâtre and F. Weis. Poste mobile de traitement de données, à module de communication local, brevet déposé au titre de l'inria. Extension pour la protection à l'étranger du brevet Français PCT/FR99/02315, October 2000.
- [16] C. Bonan. *Délivrance continue de données sur une architecture de réseaux sans fil à couverture discontinue*. PhD thesis, Université de Rennes 1, 2006.
- [17] P. Bonnet, J. Gherke, and P. Seshadri. Querying the physical world. 7(5) :10–15, October 2000.
- [18] Digital Video Broadcasting. Ip datacast over dvb-h : Electronic service guide (esg) implementation guidelines. *ETSI TS 102 592-1*, V1.1.2, July 2009.

- [19] A. Capone, L. Fratta, and F. Martignon. Bandwidth estimation schemes for tcp over wireless networks. *IEEE Transaction on Mobile Computing*, 3(2) :129–143, April 2004.
- [20] C. Casetti, M. Gerla, S. Mascolo, M. Sanadidi, and R. Wang. Tcp westwood : Bandwidth estimation for enhanced transport over wireless links. *ACM Wireless Networks*, (8) :467–479, September 2002.
- [21] C. Chang. A mobile-ip based mobility system for wireless metropolitan area network. In *Proceedings of the IEEE Int. Conference on Parallel Processing Workshops (ICPPW'05)*, June 2005.
- [22] B. Charpiot. *L'extensibilité par la répartition thématique des accès à un système d'informations distribuées*. PhD thesis, Université de Rennes 1, December 1998.
- [23] CSA. Synthèse de la consultation sur le développement des services interactifs en télévision mobile personnelle. *Conseil supérieur de l'audiovisuel, rapport annuel*, 2008.
- [24] M. Dominici, M. Fréjus, J. Guibourdenche, B. Pietropaoli, and F. Weis. Towards a system architecture for recognizing domestic activity by leveraging a naturalistic human activity model. In *In Proceedings of the Third IEEE International Conference on Smart Spaces, ruSMART 2010*, Freiburg, Germany, May 2011.
- [25] M. Dominici, G. Zecca, F. Weis, and M. Banâtre. Physical approach in smart homes : a proposition and a prototype. In *In Workshop on Goal, Activity and Plan Recognition at the International Conference on Automated Planning and Scheduling (ICAPS)*, St Petersburg, Russia, August 2010.
- [26] W. Drizet. Simulation d'un flux montant dans un réseau à couverture discontinue. Master's thesis, Juillet 2007.
- [27] ETSI. Digital video broadcasting (dvb) ; transmission system for handheld terminals (dvb-h). *European Standard EN 302 304*, November 2004.
- [28] ETSI. Digital video broadcasting (dvb) ; system specifications for satellite services to handheld devices (sh) below 3 ghz. *European Standard ETSI TS 102 585, V1.1.2*, 2008.
- [29] R.H. Frenkel and T. Imielinski. Infostations : the joy of the "many-time, many-where" communications. Technical Report TR-119, WINLAB, Avril 1996.
- [30] R. H. frenkiel, B. R. Badrinath, J. Borràs, and R. D. Yates. The infostations challenge : Balancing cost and ubiquity in delivering wireless data. *IEEE Personal Communications*, 297(2), April 2000.
- [31] S. Fu and M. Atiquzzaman. Sctp : State of the art in research, products, and technical challenges. *IEEE Communication Magazine*, April 2004.
- [32] D. Goodman, J. Borràs, N. Mandayam, and R. Yates. Infostations : A new system model for data and messaging services. In *Proceedings of the IEEE Vehicular Technology Conference*, pages 969–973, Mai 1997.
- [33] S. Govindan, H. Cheng, ZH. Yao, WH. Zhou, and L. Yang. Objectives for control and provisioning of wireless access points (capwap). *RFC 4564, IETF informational*, 2006.
- [34] C. Guillemetto. Définition et évaluation d'une architecture pour les systèmes d'information spontanés. Master's thesis, Mémoire de l'Institut d'Informatique d'Entreprise (IIE), Juin 2000.
- [35] J. L. Hennessy and D. A. Patterson. *Computer architecture : a quantitative approach*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2002.
- [36] M. T. Ho, S. Roche, and F. Weis. Design of an indoor location service in museums using rfids and wlan connectivity. In *4th Symposium of Ubiquitous Computing and Ambient Intelligence (UcamI 2010)*, Valencia, Spain, September 2010.
- [37] M. T. Ho, F. Weis, and P. Couderc. Design of a smart information diffusion service for museums using rfid-based location system. In *6th International Conference on Wireless Communications, Networking and Mobile Computing (Wicom 2010)*, Chengdu, China, September 2010.
- [38] AL. Iacono and C. Rose. Bounds on file delivery delay in an infostations system. *IEEE VTC Proceedings*, 2000.

- [39] V. Issarny, M. Banâtre, F. Weis, G. Cabilic, P. Couderc, T. Huigerra, and F. Parain. Providing an embedded software environment for wireless pdas. In *Proc. of the 9th ACM SIGOPS European Workshop, Beyond the PC : New Challenges for the Operating System*, pages 49–54, Kolding, Denmark, September 2000.
- [40] A. Jedidi. Gestion de caches pour la distribution de flux h.264 dans des réseaux à couverture discontinue. Master's thesis, Master recherche d'informatique de l'Université de Bordeaux I, Juin 2006.
- [41] A. Jedidi. *Mise en œuvre de nouveaux services dans le cadre du couplage d'un réseau de diffusion de télévision mobile personnelle et d'un réseau cellulaire 3G*. PhD thesis, Université européenne de Bretagne, 2010.
- [42] A. Jedidi and F. Weis. Caching and scheduling mechanisms for h.264 video flow delivery over discontinuous coverage wireless networks. In *IEEE International Conference on Wireless Information Networks and Systems (WINSYS)*, Porto, Portugal, July 2008.
- [43] A. Jedidi and F. Weis. Efficient switched services over a dvb-sh/3g network. Deliverable L2.11a, TVMSL project, September 2009.
- [44] A. Jedidi and F. Weis. Customized contents service over a dvb-sh/3g network. In *2nd ICST International Conference on Mobile Lightweight Wireless Systems (mobilight 2010)*, Barcelona, Spain, May 2010.
- [45] A. Jedidi and F. Weis. Efficient scheduling of low cost popular services over a dvb-sh/3g network. In *Second AIRCC International Conference on Wireless & Mobile Networks (WiMo 2010)*, Ankara, Turkey, June 2010.
- [46] A. Jedidi, F. Weis, and M. Tlais. Coupling 3g with dvb networks for low cost services. In *3rd International Conference on Engineering Management and Service Sciences (EMS 2009), Special Session on Media Interactivity and Network Convergence*, Beijing, China, September 2009.
- [47] A. Jedidi, F. Weis, M. Tlais, and S. Kerboeuf. Efficient switched services over a dvb-sh/3g network. In *5th ACM International Mobile Multimedia Communications Conference (mobimedia 2009)*, London, UK, September 2009.
- [48] R. Knopp and P. A. Humblet. Information capacity and power control in single-cell multiuser communications. In *Proceedings of the IEEE ICC'95*, June 1995.
- [49] S. Lee, W. Ma, and B. Shen. An interactive video delivery and caching system using video summarization. *Computer Communications*, 25(4) :424–435, March 2002.
- [50] F. Leleu. *Techniques d'anticipation des accès à un service d'informations distribuées. Application à un service de presse écrite électronique*. PhD thesis, Université de Rennes 1, July 1997.
- [51] H. Maillard, C. Bazin, G. Jaupitre, R. Skraba, and F. Weis. Enhanced multimedia messaging service (mms) delivery over heterogeneous networks. In *Proc. of the 9th International Conference on Intelligence in service delivery Networks (ICIN'2004)*, Bordeaux, France, October 2004.
- [52] S. Mascolo, M.Y. Sanadidi, C. Casetti, M. Gerla, and R. Wang. Tcp westwood : End-to-end congestion control for wired wireless networks. *Wireless Networks*, 8 :467–479, 2002.
- [53] R. Ménard. Lecture / écriture d'un flux rtp sur des réseaux à couverture discontinue. Master's thesis, Master 2 STS de Recherche en Informatique, Juin 2005.
- [54] M. Minard. Conception et réalisation d'un schéma d'exécution de méthodes distantes pour les systèmes d'information spontanés. Master's thesis, DEA d'informatique de l'Université de Rennes 1, Juin 2002.
- [55] A. Mishra, M. Shin, and WA. Arbaugh. Context caching using neighbor graphs for fast handoffs in a wireless network. *Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE INFOCOM)*, 2004. Hong Kong.
- [56] D. Négru, A. Mehaoua, Y. Hadjadj-Aoul, and C. Berthelot. Dynamic bandwidth allocation for efficient support of concurrent digital tv and ip multicast services in dvb-t networks. *Elsevier*, July 2005.

- [57] B. O'Hara, P. Calhoun, and J. Kempf. Configuration and provisioning for wireless access points (capwap) problem statement. *RFC 3990, IETF informational*, February 2005.
- [58] J. Ostermann, J. Bormans, P. List, D. Marpe, M. Narroschke, F. Pereira, T. Stockhammer, and T. Wedi. Video coding with h.264/avc : tools, performance, and complexity. *Circuits and Systems Magazine, IEEE*, 4(1) :7 – 28, quarter 2004.
- [59] B. Pietropaoli, M. Dominici, and F. Weis. Multi-sensor data fusion within the belief functions framework - application to smart home services. In *In Proceedings of the fourth IEEE International Conference on Smart Spaces, ruSMART 2011*, St Petersburg, Russia, August 2011.
- [60] X. Qin and R. Berry. Exploiting multiuser diversity for medium access control in wireless networks. In *Proceedings of the 22nd IEEE Infocom conference*, March 2003.
- [61] J. Rosenberg and al. Sip : Session initiation protocol. RFC 3261, June 2002.
- [62] D. Rouffet, S. Kerboeuf, L. Cay, and V. Capdevielle. 4g mobile. *Alcatel Telecommunications Review*, 2005.
- [63] S. Ben Sassi. Découverte des contrôleurs d'accès dans un réseau à couverture discontinue. Master's thesis, Master recherche d'informatique de l'Université de Bordeaux I, Juillet 2007.
- [64] B. N. Schilit, N. Adams, R. Gold, M. M. Tso, and R. Want. The ParcTab Mobile Computing System. In *Proceedings of the Fourth Workshop on Workstation Operating Systems*, pages 34–39, October 1993.
- [65] R. Skabra, G. Watts, F. Weis, and M. Banâtre. Dynamic network interface selection in heterogeneous networks, european patent. number : 04 291 679.1., July 2004.
- [66] P. Stenstrom. A survey of cache coherence schemes for multiprocessors. 23(6) :12–24, June 1990.
- [67] R. Stewart and al. Streamcontrol transport protocol (sctp) - partial reliability extension. RFC 3578, May 2004.
- [68] ALCATEL LUCENT BELL LABS Research Project Technical annex. Unlimited mobile tv with dvb-sh : Enabling high-quality mobile tv anywhere and on any device. 2008.
- [69] D. Terry, D. Goldberg, D. Nichols, and B. Oki. Continuous queries over append-only databases. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 321–330, June 1992.
- [70] M. Tlais. *Discontinuous coverage architecture : challenges, design and evaluation*. PhD thesis, Université de Rennes 1, 2008.
- [71] M. Tlais and F. Weis. Centralized mobility prediction support in a hierarchical architecture. In *Proc. of the International Conference on Wireless and Mobile Communications (ICWMC'06)*, Bucharest, Romania, August 2006.
- [72] M. Tlais and F. Weis. Distributed communication model in hierarchical infostations systems. *Concurrency and Computation : Practice and Experience*, 9(19) :1183–1192, June 2007.
- [73] M. Tlais and F. Weis. Exploitation of up-link flows in hierarchical infostation systems. Technical Report 1851, IRISA, May 2007.
- [74] M. Tlais, F. Weis, and C. Bonan. Mobility prediction in 4g d-cov networks. In *Proc. of the 1st IEEE International Workshop on Performance Analysis and Enhancement of Wireless Networks (PAEMN'06)*, Vienna, Austria, April 2006.
- [75] M. Tlais, F. Weis, and S. Kerboeuf. Enhancing the users' experience in a discontinuous coverage architecture. In *IEEE International Wireless Communications and Mobile Computing Conference (IWCMC)*, Crete Island, Greece, August 2008.
- [76] D. Touzet. Découverte d'informations dans les systèmes d'information spontanés. Master's thesis, DEA d'informatique de l'Université de Rennes 1, Juin 2000.
- [77] D. Touzet. *Interrogation continue des systèmes d'information de proximité*. PhD thesis, Université de Rennes 1, March 2004.

- [78] D. Touzet, J.M. Menaud, F. Weis, P. Couderc, and M. Banâtre. Side surfer : a spontaneous information discovery and exchanges system. In *Proceedings of International Workshop on Smart Appliances and Wearable Computing (IWUCC'01)*, Barcelona, Spain, September 2001.
- [79] D. Touzet, JM. Menaud, Frédéric Weis, and Michel Banâtre. Side surfer : Enriching casual meetings with spontaneous information gathering. *ACM SIGARCH CAN*, 29(5) :76–83, December 2002.
- [80] D. Touzet, F. Weis, and M. Banâtre. Persend : Enabling continuous queries in proximate environments. In *Proceedings of the Workshop on Mobile and Ubiquitous Data Access (WMUIA)*, Udine, Italy, March 2003.
- [81] D. Touzet, F. Weis, and M. Banâtre. Persend : Enabling continuous queries in proximate environments. In *Proceedings of the Workshop on Mobile and Ubiquitous Data Access (WMUIA)*, Udine, Italy, March 2003.
- [82] D. Touzet, F. Weis, and M. Banâtre. Architectures pour l'ubiquité numériques. *TSI*, 23(4) :439–478, 2004.
- [83] D. Touzet, F. Weis, and M. Banâtre. *Mobile and Ubiquitous Information Access, Fabio Crestani, Mark Dunlop, Stefano Mizzaro (Eds)*, volume 2954, chapter Sensing and Filtering Surrounding Data : The PERSEND Approach, pages 283–297. Springer Verlag, 2004.
- [84] A. Troël. Etude et mise en œuvre d'un traitement prédictif dans les systèmes d'information spontanés. Master's thesis, Mémoire de l'Institut d'Informatique d'Entreprise (IIE), Juin 2000.
- [85] A. Troël. *Prise en compte de la mobilité dans les interactions de proximité entre terminaux à profils hétérogènes*. PhD thesis, Université de Rennes 1, March 2004.
- [86] A. Troël, M. Banâtre, P. Couderc, and F. Weis. Predictive scheme for proximate interactions. In *Proceedings of the International Workshop on Smart Appliances and Wearable Computing (IWSAWC'01)*, pages 235–239, Mesa, AZ, United States, April 2001.
- [87] A. Troël, M. Banâtre, P. Couderc, and F. Weis. Progressive html for proximate and automatic interactions. In *Proceedings of the International Workshop on Smart Appliances and Wearable Computing (IWSAWC'02)*, pages 723–727, Vienna, Austria, July 2002.
- [88] A. Troël, F. Weis, and M. Banâtre. Prise en compte du mouvement dans les systèmes de communication sans fil. Technical Report 1508, IRISA, January 2003.
- [89] A. Troël, F. Weis, and M. Banâtre. Représentation du voisinage physique dans les interactions de proximité. Technical Report 1551, IRISA, September 2003.
- [90] A. Troël, F. Weis, and M. Banâtre. Prise en compte du mouvement dans les systèmes de communication sans fil. *TSI*, 24(1) :65–94, 2005.
- [91] UDCAST. Dvb-h mobile tv flexible satellite distribution. *European Standard ETR 154*, January 2007.
- [92] A. Ward, A. Jones, and A. Hopper. A new location technique for the active office. In *Proceedings of Mobicom'97*, September 1997.
- [93] F. Weis, F. Allard, A. Luu, M. Tlais, ML. Alberi-Morel, and S. Kerboeuf. Improving wireless network capacity by logical discontinuous coverage. In *Proceedings of the IEEE International Symposium on Wireless Communication Systems (ISWCS)*, Island, November 2008.
- [94] F. Weis, M. Banâtre, P. Couderc, and J.M. Menaud. Proximate interactions for wireless appliances. In *Proc. of the 9th ACM SIGOPS European Workshop, Beyond the PC : New Challenges for the Operating System*, pages 127–132, Kolding, Denmark, September 2000.
- [95] F. Weis and S. Roche. Specification of user device software architecture. Deliverable D3.1a, FP7 EU SmartMuseum project, January 2009.
- [96] M. Weiser. Some Computer Issues in Ubiquitous Computing. *Communications of the ACM*, 36(7) :74–84, July 1993.
- [97] L. Yang, P. Zerfos, and E. Sadot. Architecture taxonomy for control and provisioning of wireless access points. *RFC 4118, IETF informational*, 2005.

Table des figures

1.1	Communication entre humains	15
1.2	Principe d'une entité mobile (exemple avec un être humain)	16
1.3	Communication par voisinage entre deux entités mobiles	17
1.4	Espace visible de $SIS(E_i, t_i)$	17
1.5	Détection de la composition d'un S.I.S	19
1.6	Adressage et contrôle de l'espace visible d'un S.I.S	20
1.7	Un exemple de voisinage unilatéral strict	21
1.8	Interprétation cinématique	22
1.9	Évolution physique du S.I.S	25
1.10	Mise à jour au sein de l'espace visible du S.I.S	26
1.11	Architecture du système d'interrogation PERSEND	27
1.12	Interactions entre terminaux dans le cadre du <i>Web de proximité</i>	29
1.13	Analyse des données via le <i>data mining</i>	30
1.14	Mécanisme de découverte des intérêts communs	31
2.1	Fonctionnement d'un service de <i>streaming</i> sur un réseau à couverture discontinue	37
2.2	Impact d'un cache dans le terminal	38
2.3	Modèle des débits offert par un point d'accès	38
2.4	Utilisation de caches dans les architectures multiprocesseurs	39
2.5	Utilisation de caches dans un réseau à couverture discontinue	40
2.6	Utilisation d'un cache intermédiaire dans l'infrastructure	41
2.7	Utilisation de trois niveaux de caches distribués dans l'infrastructure	42
2.8	Chargement des caches AP via une discrimination des débits radio	43
2.9	Capacité moyenne d'un point d'accès	44
2.10	Modèle de mobilité entre deux APs	44
2.11	Politique de distribution des données au niveau du cache AC	45
2.12	Utilisation d'un seuil de déclenchement et d'une rafale de données	47
2.13	Utilisation de deux niveaux de caches distribués dans l'infrastructure	48
2.14	Piles de protocoles au sein du serveur, de l'AC et du terminal	49
2.15	Évolution de la bande passante durant le trajet d'un terminal	49
2.16	Intégration du contrôle de flux au sein de l'AC	50

2.17	Évolution du nombre d'interruptions de service	51
2.18	Exploitation d'un flux montant dans un réseau à couverture discontinue	52
2.19	Principes des CoSs	53
2.20	Analyse des performances sur le lien terminal - AC	55
2.21	Gestion de la discontinuité	56
2.22	Mise à disposition des données vers le serveur	57
2.23	Découverte de l'AC courant	58
2.24	Découverte distribuée de l'AC suivant	59
2.25	Taille de la mémoire de stockage de l'AC	60
3.1	Utilisation conjointe de différentes technologies dans les réseaux 4G	63
3.2	Couplage d'un réseau DVB et d'un réseau cellulaire	64
3.3	Calcul de la bande résiduelle dans un réseau DVB-SH	66
3.4	Interconnexion entre les réseaux 3G et DVB	67
3.5	Fonctions remplies par l'UBR dans un couplage 3G vers DVB	69
3.6	Politique d'insertion des flux dans la tête de réseau DVB	70
3.7	Politique d'insertion des flux dans l'UBR	71
3.8	Politique d'insertion des flux avec une bande résiduelle étendue	72
3.9	Prédiction du nombre de souscripteurs à un service non programmé	72
3.10	Substitution des flux au niveau d'un terminal 3G-DVB	73
3.11	Apport d'un réseau MBMS	75
3.12	Impact du service sur le réseau via l'utilisation de MBMS	76

Résumé

Les travaux présentés s'inscrivent dans le cadre des systèmes mobiles et distribués, et s'intéressent tout particulièrement aux perspectives offertes par les réseaux locaux sans fil.

A l'opposé de la complexité de déploiement d'une infrastructure cellulaire étendue, les interactions sans fil courte portée peuvent être utilisées de manière très simple, sans infrastructure. Ainsi, elles permettent à des calculateurs proches d'échanger automatiquement des informations. Nous proposons des supports système prenant en compte la volatilité des communications sans fil, et permettant de développer des applications tirant spontanément parti de la proximité physique des noeuds mobiles.

Ces travaux sont ensuite étendus dans le cadre d'autres familles de réseaux sans fil. Ainsi, nous nous intéressons aux réseaux à couverture discontinue. La technologie support est la même que celle de notre première étude. Simplement, les communications entre les noeuds mobiles ne sont plus directes, mais passent par une borne fixe. Cette borne définit une bulle radio de taille limitée. C'est l'interconnexion de ces bulles, sans souci de continuité de la couverture radio, qui permet d'envisager un réseau étendu et simple à déployer. Dans ce cadre, les mécanismes système étudiés permettent de masquer l'intermittence de la connectivité, et autorisent le support d'applications exploitant les flux montants et descendants dans le réseau.

Enfin, dans une dernière partie, nous abordons le problème du couplage système de deux architectures sans fil hétérogènes. Un tel couplage présente des objectifs comparables à ceux des réseaux à couverture discontinue : offrir des nouveaux services sur une couverture large, à des densités importantes d'utilisateurs mobiles. Ces travaux débouchent sur la définition de mécanismes permettant de coupler au sein d'un même service des propriétés fonctionnelles des deux infrastructures sans fil.

Abstract

The research presented in this document centers on the management of mobile distributed systems, and in particular, on emerging wireless local area networks and the possibilities they offer for new applications.

In contrast to cellular wireless networks, with their technical complexity and expensive deployment overhead, local area wireless networks can be installed very simply, without the need for existing infrastructure. Wireless local-area networks permit laptops and mobile phones that are in close proximity to spontaneously exchange information. In this document, systems programming techniques are presented to facilitate the development of applications that exploit the proximity of devices, and that notably handle the inherent network volatility.

The themes presented in the document are then extended to other types of mobile network, notably to discontinued coverage networks. These are wireless networks with the particularity that user devices do not exchange information directly, but rather via base stations that integrate radio receivers and transmitters. Base stations can be easily deployed, thus facilitating a scalable and highly dynamic network. The systems programming techniques presented in this part of the document are aimed at rendering communication breakages transparent using caching techniques that manage the up-flow (device to base station) and down-flow (base station to device) of application data.

The final part of this document looks at the merging two types of heterogeneous and wireless networks. The merging of these type of networks, on a technical level, presents immense possibilities for incorporating a large number of users in a wireless infrastructure at low cost. The work presented in this final part concludes with the presentation of a series of systems programming mechanisms for integrating the two kinds of network.