

Découverte de services dans un environnement pair-à-pair pour les grilles de calcul

Cédric Tedeschi

ACI GDS - 06 juillet 2005

Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

Contexte

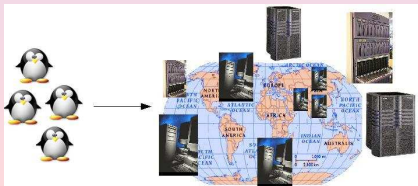
- **Constats**
 - explosion des besoins en puissance de calcul
 - coût important des supercalculateurs
- Apparition de nouvelles technologies
 - réseaux (très) haut débit
 - démocratisation des grappes de machines
 - ⇒ émergence des grilles de calcul
- Problématique des grilles
 - agréger les ressources
 - les rendre accessibles
 - ⇒ intergiciel de grille

Contexte

- **Constats**
 - explosion des besoins en puissance de calcul
 - coût important des supercalculateurs
- **Apparition de nouvelles technologies**
 - réseaux (très) haut débit
 - démocratisation des grappes de machines
 - ⇒ **émergence des grilles de calcul**
- **Problématique des grilles**
 - agréger les ressources
 - les rendre accessibles
 - ⇒ **intergiciel de grille**

Contexte

- Constats
 - explosion des besoins en puissance de calcul
 - coût important des supercalculateurs
- Apparition de nouvelles technologies
 - réseaux (très) haut débit
 - démocratisation des grappes de machines
 - ⇒ émergence des grilles de calcul
- Problématique des grilles
 - agréger les ressources
 - les rendre accessibles
 - ⇒ intergiciel de grille



Problématique

- Services \equiv ensemble des applicatifs disponibles
 - bibliothèques de calcul, exécutables
 - préinstallés sur la grille (serveurs)
- Découverte (par les clients)
 - multicritères (nom, système, architecture, adresse)
 - complétion automatique
 - rapide (minimisation des temps de communication)
- Environnement
 - statique (UDDI, Globus, GridRPC)
 - dynamique à large échelle (P2P)
 - tolérance aux pannes
 - répartition de la charge

Problématique

- Services \equiv ensemble des applicatifs disponibles
 - bibliothèques de calcul, exécutables
 - préinstallés sur la grille (serveurs)
- Découverte (par les clients)
 - multicritères (nom, système, architecture, adresse)
 - complétion automatique
 - rapide (minimisation des temps de communication)
- Environnement
 - statique (UDDI, Globus, GridRPC)
 - dynamique à large échelle (P2P)
 - tolérance aux pannes
 - répartition de la charge

Problématique

- Services \equiv ensemble des applicatifs disponibles
 - bibliothèques de calcul, exécutables
 - préinstallés sur la grille (serveurs)
- Découverte (par les clients)
 - multicritères (nom, système, architecture, adresse)
 - complétion automatique
 - rapide (minimisation des temps de communication)
- Environnement
 - statique (UDDI, Globus, GridRPC)
 - **dynamique à large échelle(P2P)**
 - tolérance aux pannes
 - répartition de la charge

Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

Réponses des technologies pair-à-pair (1)

● Définitions

- stricte (Gnutella, DHT)
- élargie (Napster)

● Propriétés

- répartition de la charge
- tolérance aux pannes (environnement dynamique)

● Pair-à-pair non structuré (Gnutella, KaZaA)

- inondation
- non exhaustivité
- surcharge du réseau

● Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)

- routage en fonction de la clef (toto.mp3, 140.77.23.58)
- exhaustivité
- routage logarithmique
- ignorance de la topologie du réseau physique
- interrogation rigide

Réponses des technologies pair-à-pair (1)

- Définitions
 - stricte (Gnutella, DHT)
 - élargie (Napster)
- Propriétés
 - répartition de la charge
 - tolérance aux pannes (environnement dynamique)
- Pair-à-pair non structuré (Gnutella, KaZaA)
 - inondation
 - non exhaustivité
 - surcharge du réseau
- Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)
 - routage en fonction de la clef (toto.mp3, 140.77.23.58)
 - exhaustivité
 - routage logarithmique
 - ignorance de la topologie du réseau physique
 - interrogation rigide

Réponses des technologies pair-à-pair (1)

- Définitions
 - stricte (Gnutella, DHT)
 - élargie (Napster)
- Propriétés
 - répartition de la charge
 - tolérance aux pannes (environnement dynamique)
- Pair-à-pair non structuré (Gnutella, KaZaA)
 - inondation
 - non exhaustivité
 - surcharge du réseau
- Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)
 - routage en fonction de la clef (toto.mp3, 140.77.23.58)
 - exhaustivité
 - routage logarithmique
 - ignorance de la topologie du réseau physique
 - interrogation rigide

Réponses des technologies pair-à-pair (1)

- Définitions
 - stricte (Gnutella, DHT)
 - élargie (Napster)
- Propriétés
 - répartition de la charge
 - tolérance aux pannes (environnement dynamique)
- Pair-à-pair non structuré (Gnutella, KaZaA)
 - inondation
 - non exhaustivité
 - surcharge du réseau
- Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)
 - routage en fonction de la clef (toto.mp3, 140.77.23.58)
 - exhaustivité
 - routage logarithmique
 - ignorance de la topologie du réseau physique
 - interrogation rigide

Réponses des technologies pair-à-pair (1)

- Définitions
 - stricte (Gnutella, DHT)
 - élargie (Napster)
- Propriétés
 - répartition de la charge
 - tolérance aux pannes (environnement dynamique)
- Pair-à-pair non structuré (Gnutella, KaZaA)
 - inondation
 - non exhaustivité
 - surcharge du réseau
- Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)
 - routage en fonction de la clef (toto.mp3, 140.77.23.58)
 - exhaustivité
 - routage logarithmique
 - ignorance de la topologie du réseau physique
 - interrogation rigide

Réponses des technologies pair-à-pair (1)

- Définitions
 - stricte (Gnutella, DHT)
 - élargie (Napster)
- Propriétés
 - répartition de la charge
 - tolérance aux pannes (environnement dynamique)
- Pair-à-pair non structuré (Gnutella, KaZaA)
 - inondation
 - non exhaustivité
 - surcharge du réseau
- Pair-à-pair structuré : DHT (CAN, Chord, Pastry, Tapestry)
 - routage en fonction de la clef (toto.mp3, 140.77.23.58)
 - exhaustivité
 - routage logarithmique
 - ignorance de la topologie du réseau physique
 - interrogation rigide

Réponses des technologies pair-à-pair (2)

- **Localité physique dans les DHT (DHT hiérarchiques)**
 - approche des **supernœuds** (Brocade, Jelly, etc.)
 - approche des **noeuds repères** (CAN, ECAN, Hieras, etc.)
- **Système d'interrogation complexe**
 - description semi-structurée (INS/Twine, etc.)
 - requête sur un ensemble de clefs numériques
 - Squid
 - auto-complétion
 - recherche multi-critères
 - perte de la répartition uniforme
 - surcouche des DHT → pire cas $O(M \log(N))$
 - nombre fixe de critères
 - pas de localité physique

Réponses des technologies pair-à-pair (2)

- Localité physique dans les DHT (DHT hiérarchiques)
 - approche des **supernœuds** (Brocade, Jelly, etc.)
 - approche des **noeuds repères** (CAN, ECAN, Hieras, etc.)
- Système d'interrogation complexe
 - description semi-structurée (INS/Twine, etc.)
 - requête sur un ensemble de clefs numériques
 - Squid
 - autocomplétion
 - recherche multicritères
 - perte de la répartition uniforme
 - surcouche des DHT → pire cas $O(N \log(N))$
 - nombre fixe de critères
 - pas de localité physique

Réponses des technologies pair-à-pair (2)

- Localité physique dans les DHT (DHT hiérarchiques)
 - approche des **supernœuds** (Brocade, Jelly, etc.)
 - approche des **noeuds repères** (CAN, ECAN, Hieras, etc.)
- Système d'interrogation complexe
 - description semi-structurée (INS/Twine, etc.)
 - requête sur un ensemble de clefs numériques
 - Squid
 - autocomplétion
 - recherche multicritères
 - perte de la répartition uniforme
 - surcouche des DHT → pire cas $O(N \log(N))$
 - nombre fixe de critères
 - pas de localité physique

Réponses des technologies pair-à-pair (2)

- Localité physique dans les DHT (DHT hiérarchiques)
 - approche des **supernœuds** (Brocade, Jelly, etc.)
 - approche des **noeuds repères** (CAN, ECAN, Hieras, etc.)
- Système d'interrogation complexe
 - description semi-structurée (INS/Twine, etc.)
 - requête sur un ensemble de clefs numériques
 - Squid
 - **autocomplétion**
 - **recherche multicritères**
 - perte de la répartition uniforme
 - surcouche des DHT → pire cas $O(N \log(N))$
 - nombre fixe de critères
 - pas de localité physique

Réponses des technologies pair-à-pair (2)

- Localité physique dans les DHT (DHT hiérarchiques)
 - approche des **supernœuds** (Brocade, Jelly, etc.)
 - approche des **noeuds repères** (CAN, ECAN, Hieras, etc.)
- Système d'interrogation complexe
 - description semi-structurée (INS/Twine, etc.)
 - requête sur un ensemble de clefs numériques
 - Squid
 - **autocomplétion**
 - **recherche multicritères**
 - **perte de la répartition uniforme**
 - **surcouche des DHT → pire cas $O(N \log(N))$**
 - **nombre fixe de critères**
 - **pas de localité physique**

Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD**
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

Description d'un service

Val = { nom_service,	(clef, valeur)
système,	(nom_service, val)
architecture,	(système, val)
adresse }	(architecture, val)
	(adresse, val)

Exemple

Val = { DGEMM,	(clef, valeur)
Linux Debian 3.0,	(DGEMM, val)
PowerPC G5,	(Linux Debian 3.0, val)
ble.ens-lyon.fr }	(PowerPC G5, val)
	(ble.ens-lyon.fr, val)

Table de Placement Lexicographique Distribuée

- Similaire à une DHT
- Arbre lexicographique des clefs des services
- Gestion de l'arrivée **dynamique** des services
- 1 nœud logique \equiv 1 clef
- Deux types de clefs
 - clef **réelle**
 - clef **virtuelle**
- Exemple
 - DGEMM
 - DTRSM
 - DTRMM
 - DGEMM

Table de Placement Lexicographique Distribuée

- Similaire à une DHT
- Arbre lexicographique des clefs des services
- Gestion de l'arrivée **dynamique** des services
- 1 nœud logique \equiv 1 clef
- Deux types de clefs
 - clef **réelle**
 - clef **virtuelle**
- Exemple
 - DGEMM
 - DTRSM
 - DTRMM
 - DGEMM



DGEMM

Table de Placement Lexicographique Distribuée

- Similaire à une DHT
- Arbre lexicographique des clefs des services
- Gestion de l'arrivée **dynamique** des services
- 1 nœud logique \equiv 1 clef
- Deux types de clefs
 - clef **réelle**
 - clef **virtuelle**
- Exemple
 - DGEMM
 - DTRSM
 - DTRMM
 - DGEMM

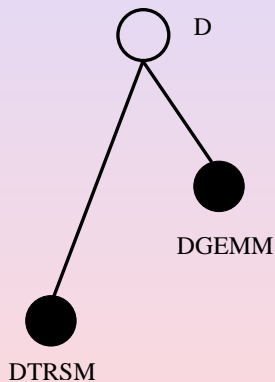


Table de Placement Lexicographique Distribuée

- Similaire à une DHT
- Arbre lexicographique des clefs des services
- Gestion de l'arrivée **dynamique** des services
- 1 nœud logique \equiv 1 clef
- Deux types de clefs
 - clef **réelle**
 - clef **virtuelle**
- Exemple
 - DGEMM
 - DTRSM
 - DTRMM
 - DGEMM

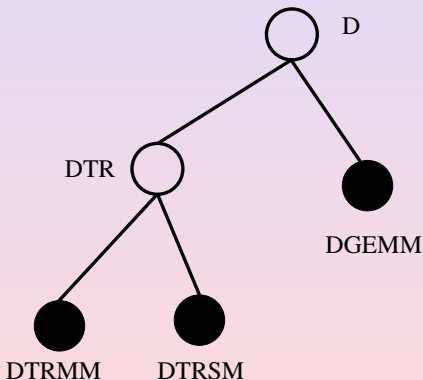
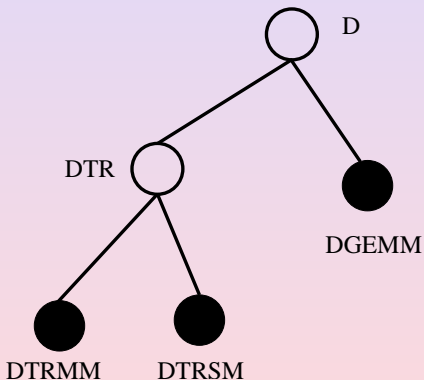


Table de Placement Lexicographique Distribuée

- Similaire à une DHT
- Arbre lexicographique des clefs des services
- Gestion de l'arrivée **dynamique** des services
- 1 nœud logique \equiv 1 clef
- Deux types de clefs
 - clef **réelle**
 - clef **virtuelle**
- Exemple
 - DGEMM
 - DTRSM
 - DTRMM
 - DGEMM

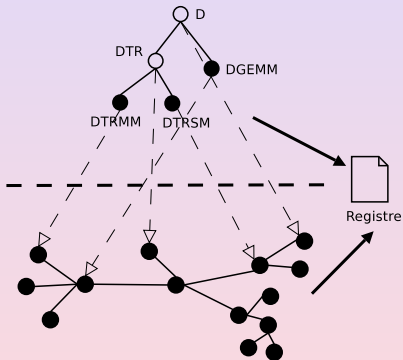


Placement des nœuds logiques

- Placement des nœuds logiques sur les **pairs**
- Deux approches
 - semi-centralisée
 - totalement distribuée

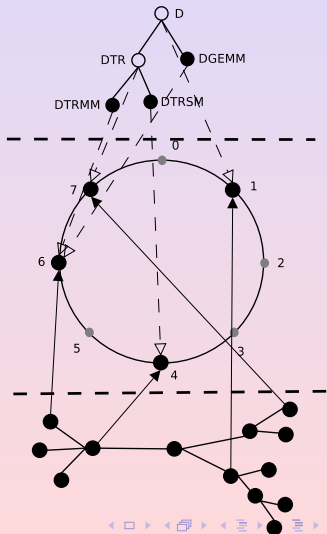
Placement des nœuds logiques

- Placement des nœuds logiques sur les **pairs**
- Deux approches
 - **semi-centralisée**
 - **totalemt distribuée**



Placement des nœuds logiques

- Placement des nœuds logiques sur les **pairs**
- Deux approches
 - semi-centralisée
 - **totale**ment distribuée



Tolérance aux pannes et arbre couvrant

- Prévention par réplication
 - k duplications des nœuds et des liens
 - répartition uniforme de la charge
 - registre
 - k fonctions de hachage aléatoires uniformes
- Minimisation des temps de communication (arbre couvrant)
 - heuristique gloutonne
 - minimiser localement pour minimiser globalement
 - intégration au processus de réplication
 - sans en modifier la complexité (linéaire)

Tolérance aux pannes et arbre couvrant

- Prévention par réplication
 - k duplications des nœuds et des liens
 - répartition uniforme de la charge
 - registre
 - k fonctions de hachage aléatoires uniformes
- Minimisation des temps de communication (arbre couvrant)
 - heuristique gloutonne
 - minimiser localement pour minimiser globalement
 - intégration au processus de réplication
 - sans en modifier la complexité (linéaire)

Algorithmes distribués

- Insertion (dans un arbre éventuellement répliqué)
 - contact
 - routage
 - insertion
- Réplication
 - lancement périodique depuis la racine
 - réplication de la racine
 - descente asynchrone dans l'arbre
- Mis en place par un ensemble de messages

Algorithmes distribués

- Insertion (dans un arbre éventuellement répliqué)
 - contact
 - routage
 - insertion
- Réplication
 - lancement périodique depuis la racine
 - réplication de la racine
 - descente asynchrone dans l'arbre
- Mis en place par un ensemble de messages

Algorithmes distribués

- Insertion (dans un arbre éventuellement répliqué)
 - contact
 - routage
 - insertion
- Réplication
 - lancement périodique depuis la racine
 - réplication de la racine
 - descente asynchrone dans l'arbre
- Mis en place par un ensemble de messages

Exemple : insertion

Arrivée consécutive des services

- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol

Exemple : insertion

Arrivée consécutive des services

- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol

Exemple : insertion

Arrivée consécutive des services

- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol

●
DGEMM

Exemple : insertion

Arrivée consécutive des services

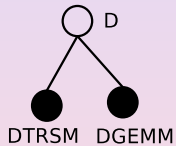
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol

●
DGEMM

Exemple : insertion

Arrivée consécutive des services

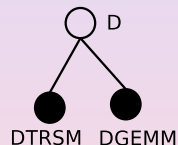
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



Exemple : insertion

Arrivée consécutive des services

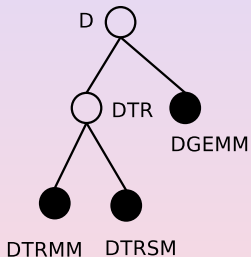
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult_chol



Exemple : insertion

Arrivée consécutive des services

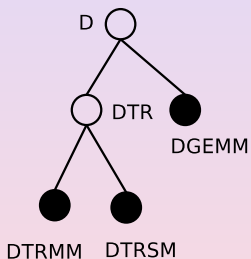
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



Exemple : insertion

Arrivée consécutive des services

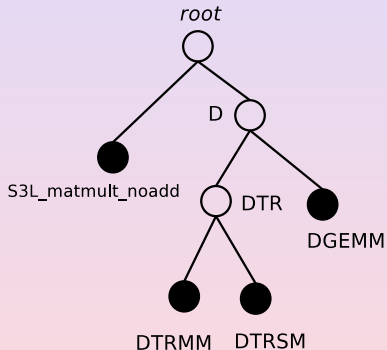
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



Exemple : insertion

Arrivée consécutive des services

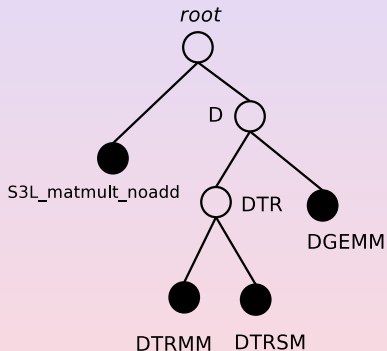
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



Exemple : insertion

Arrivée consécutive des services

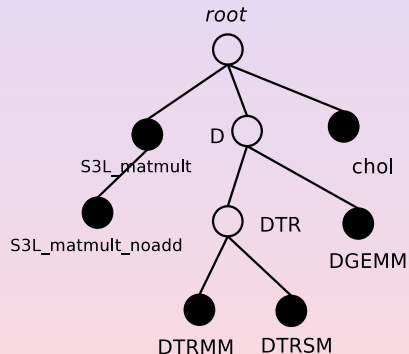
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



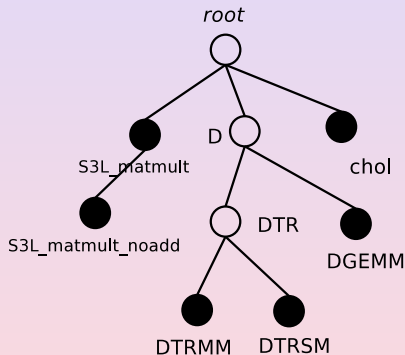
Exemple : insertion

Arrivée consécutive des services

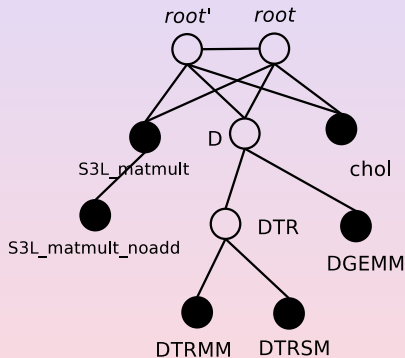
- DGEMM
- DTRSM
- DTRMM
- S3L_matmult_noadd
- S3L_matmult, chol



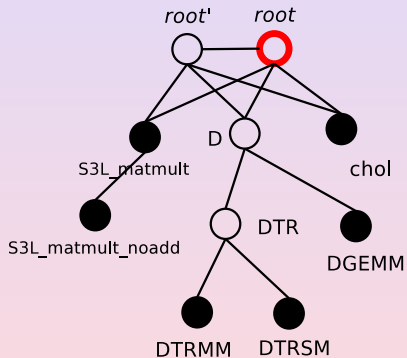
Exemple : réplication et arbre couvrant



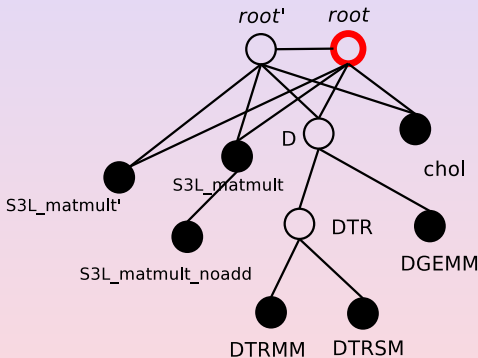
Exemple : réplication et arbre couvrant



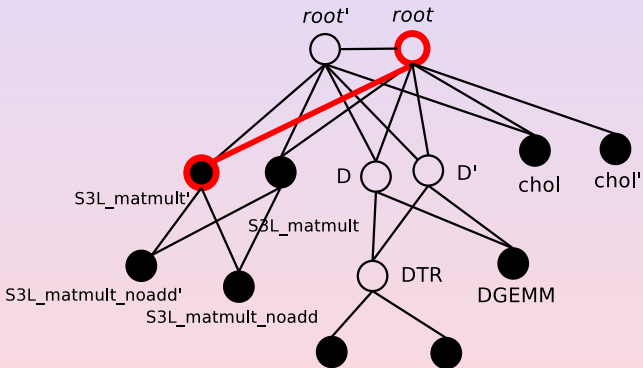
Exemple : réplication et arbre couvrant



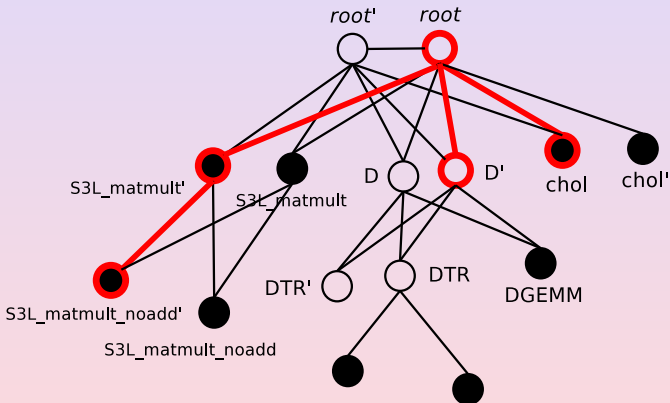
Exemple : réplication et arbre couvrant



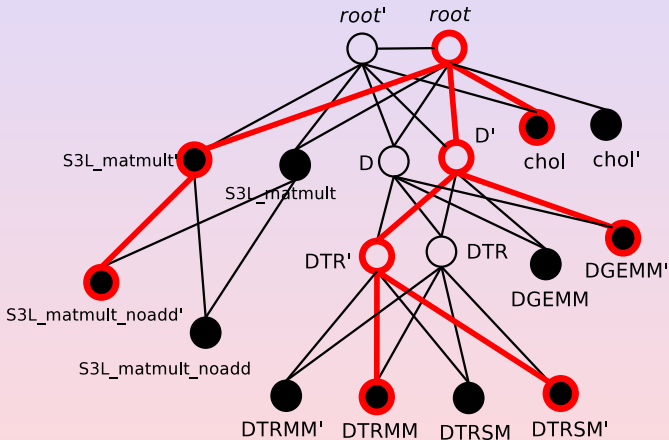
Exemple : réplication et arbre couvrant



Exemple : réplication et arbre couvrant

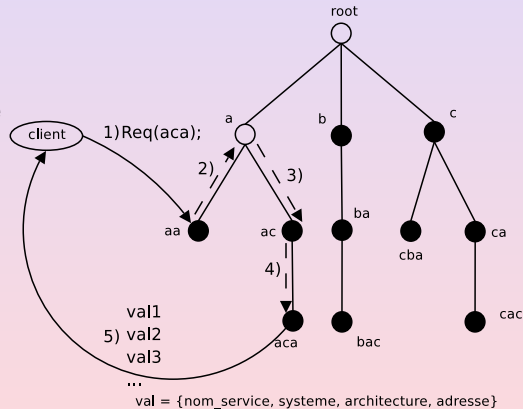


Exemple : réplication et arbre couvrant



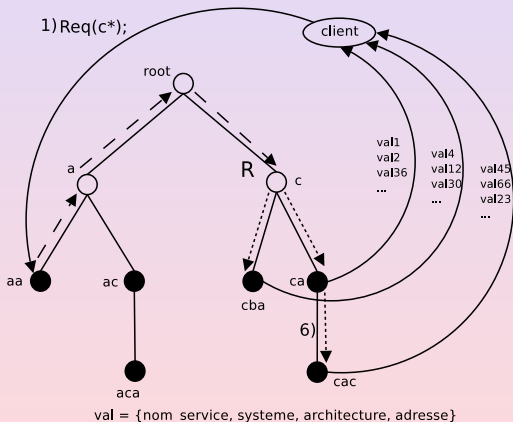
Mécanismes d'interrogation de la TPLD

- Recherche sur chaîne complète
- Recherche sur chaîne partielle
- Recherche multicritères



Mécanismes d'interrogation de la TPLD

- Recherche sur chaîne complète
- Recherche sur chaîne partielle
- Recherche multicritères



Mécanismes d'interrogation de la TPLD

- Recherche sur chaîne complète
- Recherche sur chaîne partielle
- Recherche multicritères

Mécanismes de cache

- A la demande
- Cache lors d'une requête sur chaîne partielle
 - réplication sur la racine du sous-arbre concerné
 - ⇒ réponse au client dès l'arrivée sur la racine
- Cache lors d'une requête sur chaîne complète
 - problème des clefs populaires
 - sur la route de la requête
 - routage en arrière
- Politique LRU par pair

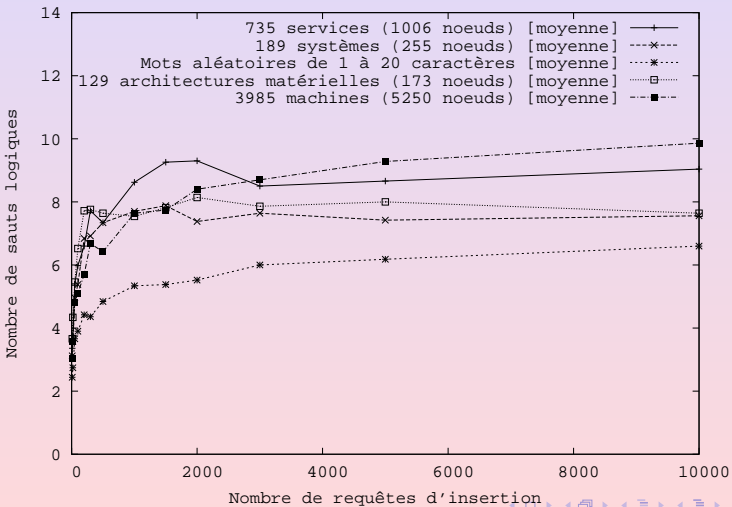
Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations**
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

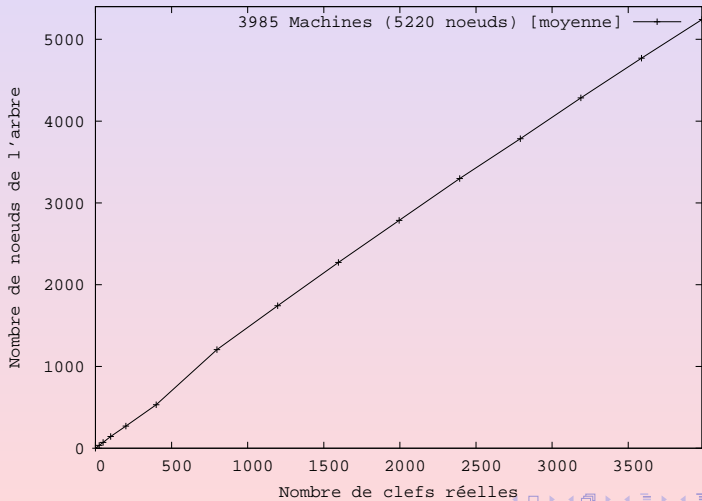
Complexités de la TPLD

- Arbre lexicographique
 - T : taille max des mots
 - A : alphabet utilisé
- Complexité du routage
 - **bornée par $2T$, indépendant de la taille du réseau**
 - en $O(\log(N))$ pour le réseau le plus grand possible
- Taille de la table de routage
 - **bornée par $\text{Card}(A)$, indépendant de la taille du réseau**
 - en $O(\sqrt{N})$ pour le réseau le plus grand possible
- Choix local du routage
 - une case de tableau par caractère possible de A
 - **en $O(1)$**

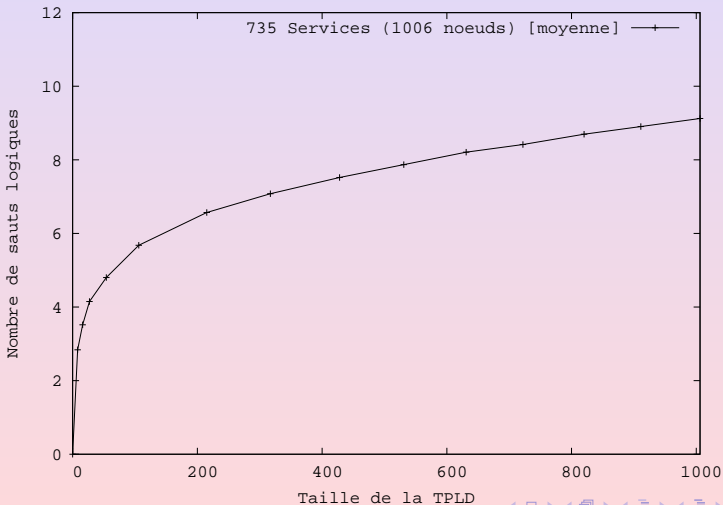
Taille de l'arbre / nombre d'insertions



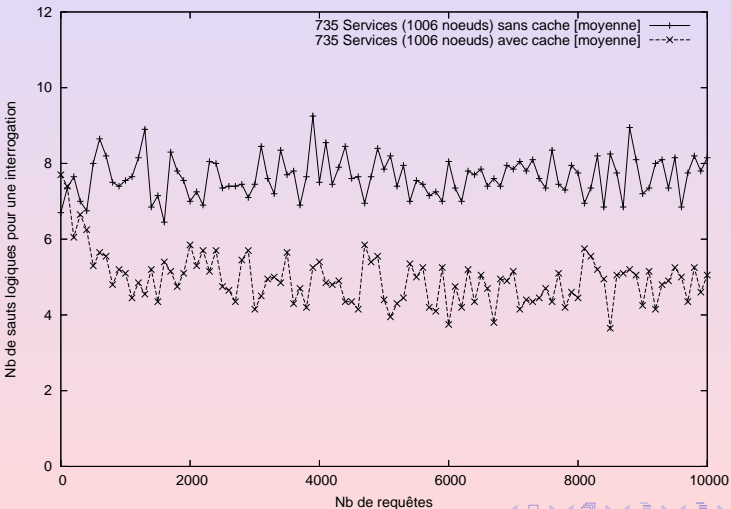
Proportionnalité clefs réelles / taille de la TPLD



Routage d'une requête d'interrogation



Cache (requête sur chaîne complète)



Gains sur l'existant

Complexités (TPLD/DHT)

DHT	Taille table routage	Taille routage	Décision routage
CAN	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$
Chord	$O(\log(N))$	$O(\log(N))$	$O(\log(N))$
Pastry/Tapestry	$O(\log(N))$	$O(\log(N))$	$O(1)$
TPLD	$O(\text{Card}(A))$	$O(T)$	$O(1)$

Gains sur l'existant

Fonctionnalités (multi-TPLD/Squid)

-	SQUID	Multi-TPLD
Requêtes multicritères	X	X
Chaînes partielles	X	X
Nombre variable de dimensions	-	X
Adaptation au réseau physique	-	X
Optimisation par cache	-	X

Plan

- 1 Introduction
- 2 État de l'art
- 3 Contribution : TPLD
 - Création et maintenance dynamiques
 - Interrogation
- 4 Validations
 - Analyse
 - Expériences de simulation
 - Gains sur l'existant
- 5 Conclusion

Conclusion

- Contribution : TPLD
 - architecture distribuée pour la découverte de services dans un environnement dynamique à large échelle
 - dynamique
 - tolérante aux pannes (par k-réplication)
 - s'adaptant aux performances du réseau (arbre couvrant)
 - optimisations (cache)
- Perspectives
 - ordonnancement
 - gestion de données
 - implantation