

Application BASTRI

Fiches Equipes

KERDATA (SR0346XR)

Gestion de données distribuées à très grande échelle pour les grilles et les clouds

KERDATA KERDATA (SR0523GR)

Statut: Terminée

Responsable : Gabriel Antoniu

Mots-clés de "A - Thèmes de recherche en Sciences du numérique - 2024" : *Aucun mot-clé.*

Mots-clés de "B - Autres sciences et domaines d'application - 2024" : *Aucun mot-clé.*

Domaine : Réseaux, systèmes et services, calcul distribué
Thème : Calcul distribué et applications à très haute performance

Période : 01/07/2009 -> 30/06/2012

Dates d'évaluation :

Etablissement(s) de rattachement : CNRS, ENS RENNES

Laboratoire(s) partenaire(s) : IRISA (UMR6074)

CRI : Centre Inria de l'Université de Rennes

Localisation : Centre Inria de l'Université de Rennes

Code structure Inria : 031087-0

Numéro RNSR : 200920935W

N° de structure Inria: SR0346XR

Présentation

Les recherches de l'équipe KerData abordent le domaine de la gestion de données réparties à de très grandes échelles, en particulier sur les clouds et les infrastructures post-pétaflopiques. Nous visons des applications à haute performance à accès intensif aux données, qui nécessitent le traitement de données massives non structurées - BLOBs : binary large objects (de l'ordre du téraoctet) - stockées en grand nombre (des milliers voire des dizaines de milliers), accédées de manière hautement concurrente par un très grand nombre de processus (des milliers voire des dizaines de milliers d'accès quasi-simultanés), avec un grain fin d'accès (de l'ordre du mégaoctet). Voici quelques exemples de telles applications :

- Fouille de données massives (typiquement en utilisant le paradigme MapReduce) sur des clouds.
- Stockage réparti pour des applications numérique exécutées sur des architectures post-pétaflopiques.
- Services avancés de gestion de données sur des clouds (i.e. optimisés pour la concurrence, orientés versioning, etc.) pour le stockage des données applicatives et pour la gestion des machines virtuelles au niveau IaaS.
- Stockage pour des applications s'exécutant sur des "desktop grids" nécessitant une très haute bande passante en écriture.

Axes de recherche

- **Gestion de BLOBs multiversion**
Nous travaillons à la conception, à la mise en oeuvre et à la validation expérimentale d'une plate-forme de partage de données non structurées de grande taille (BLOBs) appelée **BlobSeer**. Cette plate-forme est destinée à nous permettre d'aborder les défis mentionnés ci-dessus : des données massives, des accès à grain fin hautement concurrents, tout en fournissant un support pour le multiversioning et pour la gestion décentralisée des méta-données.
- **Gestion de données sur les clouds**
Sur des infrastructures de type Infrastructure-as-a-Service (IaaS), les ressources du cloud sont exploitées sur la base d'un modèle "à la demande": au lieu d'acheter et de devoir assurer la maintenance du matériel, les utilisateurs louent des machines virtuelles et de l'espace de stockage. Il faut alors gérer le stockage et le traitement des données sur des ressources de stockage externes, virtualisées. Un des défis à relever dans ce contexte est de concilier performance, passage à

Contact

- **Responsable :** Gabriel Antoniu
- **Tél :** 02.99.84.72.44
- **Secrétariat Tél :** 02.99.84.71.86

En savoir plus

- Site de l'équipe
- Site sur inria.fr
- Site du [responsable](#)
- Derniers Rapports d'Activité : [2016](#), [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)

Documents sur la structure

- [Intranet](#)
- [Privés](#)

Décisions

- **6847** (01/07/2009) : création
- **7376** (05/07/2010) : prolongation
- **8014** (16/06/2011) : prolongation
- **8333** (20/12/2011) : prolongation

Localisation

- **Adresse postale :** Centre Inria de l'Université de Rennes 263, avenue du Général Leclerc Campus universitaire de Beaulieu 35042 Rennes Cedex France
- **Coordonnées GPS :** 48.116, - 1.64

l'échelle, sécurité et qualité de service. Il est important, par ailleurs, d'étudier l'impact du partage des ressources physiques sur ces différents plans. L'objectif est d'évaluer les bénéfices du stockage à base de BLOBs apportés par l'approche BlobSeer (gestion décentralisée des méta-données, multiversioning, efficacité dans de conditions de haute concurrence) à la conception de systèmes de fichiers répartis au service des applications de fouille de données massives, notamment via le paradigme MapReduce.

- **Gestion de données pour les systèmes post-pétaflopiques**
En parallèle avec l'émergence des infrastructures de type cloud, des efforts considérables portent aujourd'hui sur la construction de machines post-pétaflopiques, telles que **Blue Waters**. Ces systèmes visent à soutenir des performances de l'ordre du pétaflop pour un large spectre d'applications scientifiques. Sur de telles infrastructures, la gestion des données a un impact fort sur la performance des applications. Ces infrastructures ont des caractéristiques architecturales particulières (i.e. hiérarchie mémoire multi-niveaux supportant des centaines de milliers de coeurs) destinées à permettre un degré de parallélisme sans précédent. Le système de stockage des données doit alors être conçu de manière à ne pas limiter le passage à l'échelle. Nos recherches se focalisent sur les besoins des applications numériques qui doivent s'exécuter sur ces architectures: l'objectif est d'évaluer les bénéfices que l'approche BlobSeer peut apporter à différents niveaux (gestion des données, des métadonnées) pour assurer des entrées/sorties performantes dans des conditions de très haute concurrence des accès aux données.

Logiciels

- **BlobSeer**

Relations industrielles et internationales

- AzureBrain : projet collaboratif avec Microsoft sur l'analyse conjointe de données génétiques et de neuroimagerie sur des clouds Azure dans le cadre du Centre Commun de Recherche INRIA-Microsoft Research
- MapReduce : projet ANR sur la gestion avancée de données sur les clouds à l'aide du paradigme MapReduce, avec des partenaires internationaux et industriels : Argonne National Lab (USA), University of Illinois at Urbana-Champaign (UIUC, USA), IBM
- FP3C : un projet ANR-JST sur la programmation de machines post-pétaflopiques, regroupant les acteurs majeurs académiques français et japonais du domaine. Collaboration forte avec l'Université de Tsukuba.
- NCSA/UIUC : collaboration dans le cadre du JLPC (Urbana-Champaign) sur la gestion de données sur les architectures post-pétaflopiques, optimisée pour les accès concurrents
- SCALUS : Marie Curie Initial Training Network (FP7).
- DataCloud@work : équipe associée avec l'Université Polytechnique de Bucarest, Roumanie.