

**Observatoire de Strasbourg**

Mardi 6 mai 2003

# Etude de la « Clusterisation » d'un service de catalogues ( VizieR )

André Schaaff



# Introduction

- **Ce travail est actuellement réalisé dans le cadre d'un stage (5 mois) de dernière année d'école d'ingénieur**
- **2 phases :**
  - **Réflexion et prototypage (actuellement en cours)**
    - Répartition des calculs
    - Répartition également des données ?
      - Il est encore actuellement possible de dupliquer les données sur toutes les machines (par rapport à la configuration matérielle choisie)
      - Mais le volume de données risque de dépasser rapidement la taille unitaire des disques
  - **Implantation d'une première version publique (en août)**

# VizieR

- **VizieR fournit un accès à la plus complète librairie de catalogues organisés, documentés et disponibles en ligne.**
- **Des outils d'interrogation permettent à l'utilisateur de sélectionner des tables pertinentes, d'extraire et de formater des enregistrements correspondants aux critères de recherche.**
- **L'accès aux très grands catalogues a été optimisé : Guide Star Catalogs, USNO-B1, 2MASS...**
  - **Ordre de grandeur : plusieurs centaines de millions d'objets**
  - **Pour chacun de ces catalogues : stockage sous forme binaire + programme associé**



- **Le service VizieR existe depuis 1996.**

# Problématique

- La base de données ainsi que les traitements sont actuellement centralisés sur un serveur Sun.
- Nous souhaitons délocaliser les plus gros catalogues (format binaire, Sybase pour les autres catalogues) ainsi que les traitements les plus coûteux en temps machine sur plusieurs serveurs
- Notre attente : disposer d'une solution de grappes de calcul efficace, simple à mettre en œuvre et surtout peu coûteuse

# CLIC - Cluster Linux pour le Calcul

- **Partenariat MandrakeSoft, Bull et l'INPG/INRIA**
- **Financement RNTL**
- **Objectifs :**
  - simple et facile à installer
  - unifier l'ensemble des phases d'installation, de configuration de la couche d'interconnexion et de déploiement des applications parallélisées
- **3 phases :**
  - développement et publication d'une distribution Linux contenant tout les outils nécessaires au déploiement rapide d'une grappe de calcul prête à l'utilisation.
  - publication d'outils spécifiques d'administration, de contrôle et d'évaluation pour la solution de clustering
  - publication d'outils et d'applications spécialisés pour le développement en environnement parallèle.

# Travail en cours

- Une petite configuration de 6 PC (au total 6Go de RAM et 2,4 To de disque) pour la phase d'étude, de prototypage et la première version publique
- A partir d'un échantillon représentatif de grands catalogues :
  - Travail préparatoire : analyse fine des programmes associés aux fichiers binaires des grands catalogues
  - Parallélisation des traitements
    - NB : les données sont dans un premier temps dupliquées sur chaque noeud
  - Tests de performances
  - Etude de l'incidence de la répartition des données sur les performances obtenues précédemment

# Conclusion

- **Notre carnet de route :**
  - **La partie Analyse – parallélisation – performances pour la fin mai**
  - **Une étude de l'incidence de la répartition des données et la définition d'une démarche (généralisation à l'ensemble des grands catalogues) pour la fin juin**
  - **Une première version publique pour la fin août**
- **Une inconnue pour l'avenir : la pérennité de CLIC après la fin de la 3ème phase**