



# Projet IGGI

## Infrastructure pour **G**rappe, **G**rille et **I**ntranet

CASCIMODOT - Novembre 2005

Fabrice Dupros



Géosciences pour une Terre durable

**brgm**

# CONTEXTE

## > **Etablissement Public à caractère Industriel et Commercial (EPIC)**

- Sous la tutelle des Ministères en charge de la Recherche, de l'Environnement et de l'Industrie

> **850 personnes**

## > **Domaines thématiques**

- Ressources minérales, Eau
- Aménagement et Risques naturels
- Environnement et pollutions
- Métrologie de l'environnement
- Cartographie et systèmes d'information

## > Réseau National des Technologies Logicielles

- Labellisation en Avril 2003
- Réunion de lancement au BRGM Novembre 2004

## > Les partenaires

- INRIA : laboratoire ID-IMAG (Grenoble)
  - Leader du projet : Jean-François Méhaut
  - Sous traitance ICATIS
- MandrakeSoft



## > Projet sur 2 ans



## Objectifs

- > **Accroître le taux d'utilisation des ressources informatiques d'un réseau Intranet d'entreprise**
  - Heures ouvrables
    - Mode Interactif (~ 8h/jour), principalement sous Windows  
Tableur, saisie de documents , Mail, Internet  
Bases de données
  - Inactivité des ordinateurs (**4 jours et demi par semaine**)
    - Nuits (**12h**), week-ends (**48h**)
    - Congés (**5 semaines**) + **23** jours RTT
- > **Périodes d'inactivité : ordinateurs de bureau deviennent les nœuds d'une grappe virtuelle**
  - Applications de **calcul scientifique**
  - Basculer d'un **mode interactif** vers un **mode calcul**
  - **Cloisonnement total** entre les modes interactif et calcul
    - **mode diskless**

# Intranet : Objet complexe

## > Réseau de l'Intranet

- Plusieurs bâtiments, étages, services
- Hétérogène (Ethernet 10, 100, Gb/s, FDDI)
- Protocole unique TCP/IP: **latence importante**
- Hiérarchisé
  - Physiquement (segments, routeurs,...)
  - Logiquement (v-lan,...)

## > Machines

- Plusieurs milliers de machines...
- Configurations hétérogènes (CPU, mémoire, disque)
- Performances hétérogènes
- Logiciels (OS, BD, applications,...)

## > Utilisateurs

- Profils utilisateurs
  - Horaires
  - Données sensibles
- Retrouver sa machine dans le même état qu'il l'avait laissée la veille!

## Parc Brgm

<b>Bureautique</b>	<b>700</b>
<b>Ingénieurs / développeurs</b>	<b>350</b>
<b>Portables</b>	<b>250</b>

<b><u>Sur le site d'Orléans</u></b>	
<b>Bureautique</b>	<b>350</b>
<b>Ingénieurs / développeurs</b>	<b>300</b>
<b>Portables</b>	<b>250</b>

→ Soit environ 500 postes de travail utilisables pour la grille

## Au sein du brgm :

- Optimiser l'utilisation des ressources de calcul existantes
- Extension des possibilités de modélisation (*adaptation outils*)
- Mutualisation de moyens hétérogènes
  - Grappes dédiées, PC de bureau

- **Renforcer les moyens de calcul**
  - **Si possible à faibles couts ....**



# Simulation Numérique

## *Environnement matériel et logiciel*

### → Grappe de calcul de 6 bi-pro Athlon - 1.8Ghz - interconnexion myrinet

→ Acquisition en décembre 2002.

→ Premier dimensionnement des besoins en calcul.

→ Système dédié au pilotage du système de Réalité Virtuelle.

### → Grappe de calcul de 12 bi-pro Xeon - 3.06 Ghz - interconnexion Gigabit

→ Acquisition en Janvier 2004.

→ Serveur dédié aux activités de calcul.

→ Gestionnaire de batch OAR

## Risques Naturels

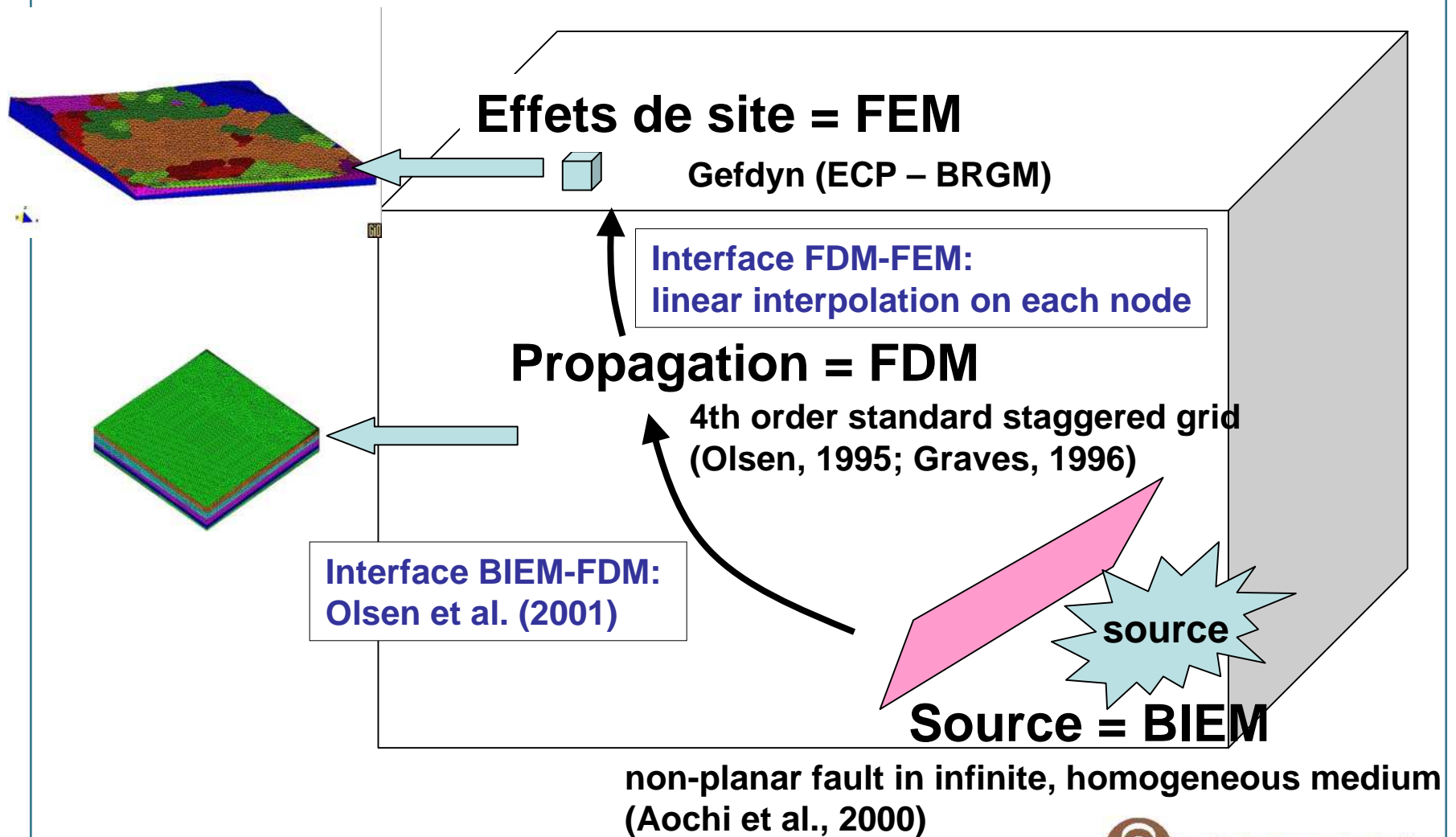
### → Propagation ondes

- Application aux risques sismiques
- Application à la géothermie ( Bouillante Guadeloupe et Soultz en Alsace )
- Outils brgm en version parallèle ou séquentielle

### → Modélisation géomécaniques

- Risques sismiques - Stockage de CO2 - Ouvrages
  - Approche Eléments Finis ou Meshfree
- **Besoin de discrétisations fines sur des domaines toujours plus grands**

# Procédure numérique : risques sismiques



## Stockage géologique du CO2

Expertise brgm reconnue ( projet Europeen CO2STORE - WEYBURN etc .. )

- Outils de modélisation Phreeqc ( USGS) ou TOUGHREACT (LBNL)
- Code purement séquentiel et runs longs ( plusieurs jours )

### • Environnement

Stockage de déchets radioactifs en grande profondeur

- Outils Phast ( USGS ) ou TOUGHREACT ( LBNL )

## Hydrogéologie

- Outil de modélisation Marthe - Étude des hydrosystèmes
- Développé au brgm depuis une vingtaine d années
- Large diffusion : Andra - EDF - Antéa etc ..

## Pesticides

- ⇒ **Modèles de transfert de pesticides dans les sols**
- ⇒ Modèles majoritairement sous Windows
- Besoins importants en termes de ressources de calcul ( campagnes de jobs )
  - Plusieurs semaines en temps CPU

# Briques de base : mise en oeuvre

# Compute mode



- > Bascule des PC en fonction des périodes d'inactivités ( reboot réseau )
- > Mode bureautique et mode nœud de calcul cloisonné ( OS différents)
  
- > **Nécessite un serveur matériel installé sur le reseau**
  - Système de calcul en grappe centralisée ( diskless ) : **intranet**
  - Serveur CM héberge système, applications et données
  - Linux standard sur les machines cibles ( distribution des images )
  - Serveur cache de données
    - Images système , données utilisateurs
  
- > **Connexion à développer avec batch-scheduler**
  - Reboot en fonction file d'attente





# Gestion des ressources

## > OAR

<http://oar.imag.fr/>

- Gestionnaire de batch
- Base de données : medium d'échanges entre composants
- Scalabilité

## > CIGRI

- Solution pour grilles légères
- Gestion de jobs multi-paramétriques ( 10k jobs )
  - Pas de communication entre travaux
- Expérience ACI-GRID et communauté Ciment ( Grenoble )

## > Connexion forte de ces outils avec GRID5000

# Checkpoint – Migration des applications

## > Application séquentielle

- Utilisation des possibilités condor
- Checkpoint utilisateur disponible pour certaines applications
  - Comparaison en cours

## > Application parallèle

- Solution LAM + BLCR ( Checkpoint/Restart Lib )
  - Solution bas niveau (Kernel)

## > Librairie SAMORY

- Checkpoint système version parallèle ou séquentielle



Laboratoire  
Informatique et  
Distribution



Géosciences pour une Terre durable

**brgm**

# En pratique pour l'utilisateur :

## > **Portail de soumission unique des travaux :**

- pour la grappe de calcul dédiée
- pour les grappes virtuelles

## > **Scheduling**

- Applications parallèles sur cluster dédié en priorité
- Prise en compte mémoire/cpu/disk pour choix plate-forme

## > **Checkpoint**

- Au niveau système ou au niveau applicatif
- Prise en compte au niveau du batch scheduler

# Premiers tests et perspectives

## Déploiement

14 Postes en test depuis mai 2005

- 7 PC dédiés + 7 PC salle de formation

Fin 2005

- *Passage en production à l'échelle d'un service ( 50 PC )*

### ✓ Gestion retour utilisateur

- Restauration contexte Windows ( hibernation )

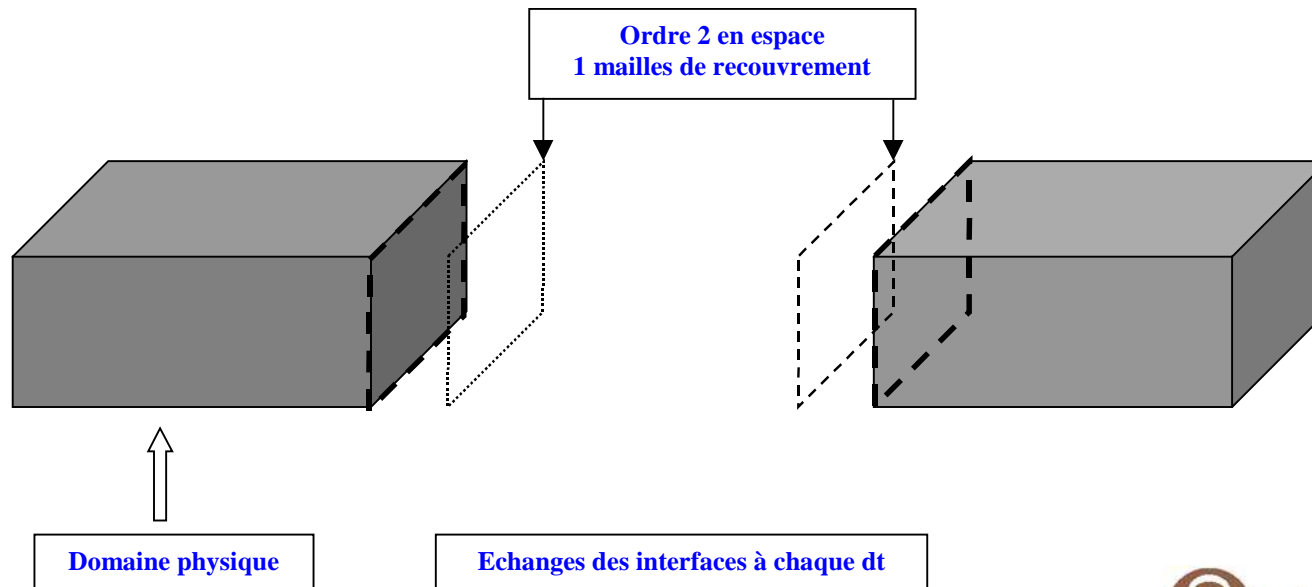
### ✓ Gestion modes de mise à disposition des postes

- Base volontariat ( type inscription à la grille )
- Base automatique ( type salle de TP )

# Performances

## • *Equation de la chaleur 2D en version parallèle*

- Différences Finis = données régulières : bonnes performances parallèles ( mémoire - CPU )
- Equation et méthodologie classique
- code simple à disséquer pour évaluer performances



## Équation de la chaleur 2D

Performances comparatives : 4096\*4096 - 100 dt.

	2-1	1 – 2	2-2	4-1	1-4	8-1	1-8
<b>UFR-Grenoble ComputeMode</b>	<b>23.4 0.57</b>	<b>24.08 0.33</b>	<b>11.85 0.56</b>	<b>11.65 1.55</b>	<b>12.03 0.91</b>	<b>5.88 1.43</b>	<b>6.04 0.94</b>
<b>ComputeMode brgm</b>	<b>35.1 0.43</b>	<b>46.41 0.31</b>	<b>17.9 0.31</b>	<b>17.61 1.40</b>	<b>23.12 0.89</b>	X	X
Icluster2 - grenoble	<b>46.43 0.08</b>	<b>45.6 0.02</b>	<b>22.86 0.04</b>	<b>23.26 0.10</b>	<b>22.83 0.04</b>	<b>11.42 0.08</b>	<b>11.47 0.05</b>
Cluster dédié brgm	<b>14.49 0.32</b>	<b>15.13 0.25</b>	<b>6.96 0.14</b>	<b>6.95 0.57</b>	<b>7.09 0.31</b>	<b>3.55 0.46</b>	<b>4.22 0.13</b>

## Intégration / appropriation de l'architecture

### ⇒ **Intégration informatique**

⇒ Point unique de soumission ( grappe dédiée - grappe virtuelle )

⇒ Important pour nos modélisateurs souvent néophytes

### ⇒ **Portail de suivi et d'administration**

⇒ Intégrer les différentes interfaces graphiques ( CM – OAR – Cigri )

### ➤ **Communication interne et adoption**

➤ Retour d'expérience et prise en compte spécificités

