

# Contribution à l'algorithmique parallèle

## *Des bibliothèques de calcul au metacomputing*

---

Frédéric Desprez

Projet CNRS - ENS Lyon - INRIA ReMaP

LIP ENS Lyon

INRIA Rhône-Alpes

- **CV**
- **Introduction**
- **Recouvrements et pipelines**
  - Projets LOCCS/OPIUM
- **Parallélisation de programmes Fortran**
  - Projet TransTool
- **Problem Solving Environments**
  - Projets Scilab//, DIET
- **Conclusions et perspectives**

# Curriculum Vitae



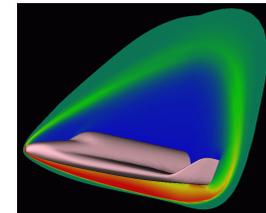
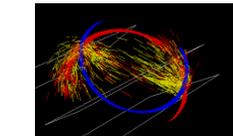
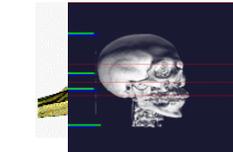
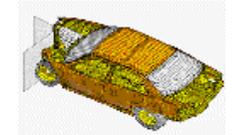
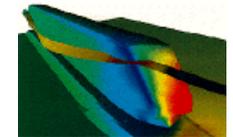
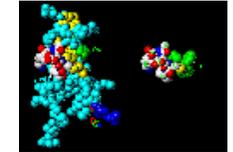
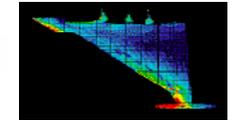
- [1994]            **Postdoctorat**  
à l'Université du Tennessee à l'ICL
- [1994-1995]    **Maître de Conférence** à l'ENSEIRB  
Recherche dans l'équipe ALiEnor
- [1995-]         **Chargé de Recherche** à l'INRIA, LIP ENS Lyon  
projet CNRS-INRIA-ENS Lyon ReMaP
- [Sep.2000-]    Responsable du projet ReMaP



- Encadrement de 6 thèses, 9 DEA, 3 DESS, ...
- Projets européens et internationaux EuroTOPS (93-98), TTN ProHPC (97-99) et Paralin (96-99)
- Projets nationaux ARC OURAGAN (98-00) et RNRT VTHD (99-01)
- Enseignement à l'ENS, à l'ENSEIRB, à SUPELEC,
- Comité d'édition de Parallel and Distributed Computing Practice
- Organisation des JIP'94-95, MSA'2000-01 dans ICPP

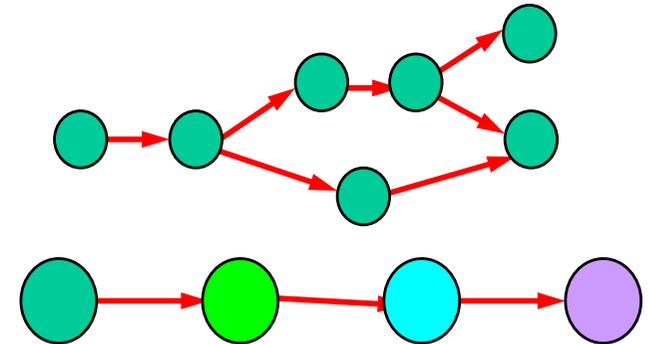
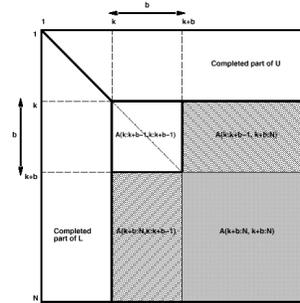
# Introduction

- Calcul numérique utilisé dans de nombreuses applications
- Simulations toujours plus précises  $\Rightarrow$  besoins importants en puissance de calcul et en capacité de stockage
- **Evolution rapide des architectures**
  - Processeurs superscalaires (du commerce) + augmentation des fréquences d'horloge, mémoires de plus en plus hiérarchiques
  - Machines parallèles (vectorielles, SM, DM, SIMD, MIMD, réseaux de stations, SMP, grappes de SMP, plates-formes de metacomputing)  $\Rightarrow$  de + en + d'hétérogénéité
- **Evolution plus lente du logiciel**
  - Pas de langage vraiment adapté au parallélisme
  - Pas de standard
  - La plupart des applications programmées en C/F77 + MPI
  - Problèmes de parallélisation automatique, debug, profiling, etc



- **Plusieurs sources de parallélisme**

- Parallélisme de données
- Parallélisme de tâches
- Parallélisme de flux



- **Plusieurs modèles de programmation en mémoire distribuée**

- Passage de messages
- Parallélisation automatique
- Langages data-parallèles
- Bibliothèques parallèles
- Java ? OpenMP + threads?
- Mélanges de langages, de bibliothèques et de supports d'exécution à hautes performances

# Bibliothèques numériques

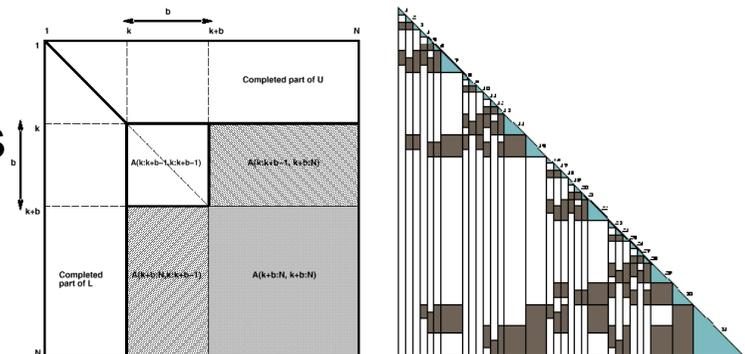
- **Avantages**

- Le parallélisme est masqué à l'intérieur de la bibliothèque
- Utilisation de noyaux numériques optimisés sur les processeurs (BLAS de niveau 3, ATLAS, PhiPAC)
- Eviter le choix de l'algorithme de résolution
- Maintenance et mise à jour

Problem type Ax = b	SDrv	EDrv	Factor	Solve	Inv	Cond Est	Iter Refin
<b>Triangular</b>				X	X	X	X
SPD	X	X	X	X	X	X	X
SPD Banded	X		X	X			
SPD Tridiagonal	X		X	X			
<b>General</b>	X	X	X	X	X	X	X
General Banded	X		X	X			
General Tridiagonal	X		X	X			
<b>Least squares</b>	X		X	X			
GQR			X				
GRQ			X				
<b>Ax = λx or Ax = λBx</b>							
Symmetric (2 types)	X	X	X	X			
General (2 types)			X	X			
Generalized BSPD		X	X	X			
SVD			X	X			

- **Inconvénients**

- Choix de la distribution laissé à l'utilisateur
- Insertion de redistributions entre les étapes de calcul
- Pas de choix de l'algorithme de résolution
- Transformation des structures de données des applications
- Encore limitées (creux)



# Problem Solving Environments

- **Besoin des utilisateurs**

- Applications simples à exécuter (parallélisme obligatoirement caché) dans un environnement intégré
- Exécuter n'importe quand des simulations sans limitation de taille
- Partager des codes entre les utilisateurs
- Pouvoir collaborer avec d'autres disciplines
- Utiliser efficacement et simplement des ressources disponibles à travers le réseau
- Mettre en commun des puissances de calcul

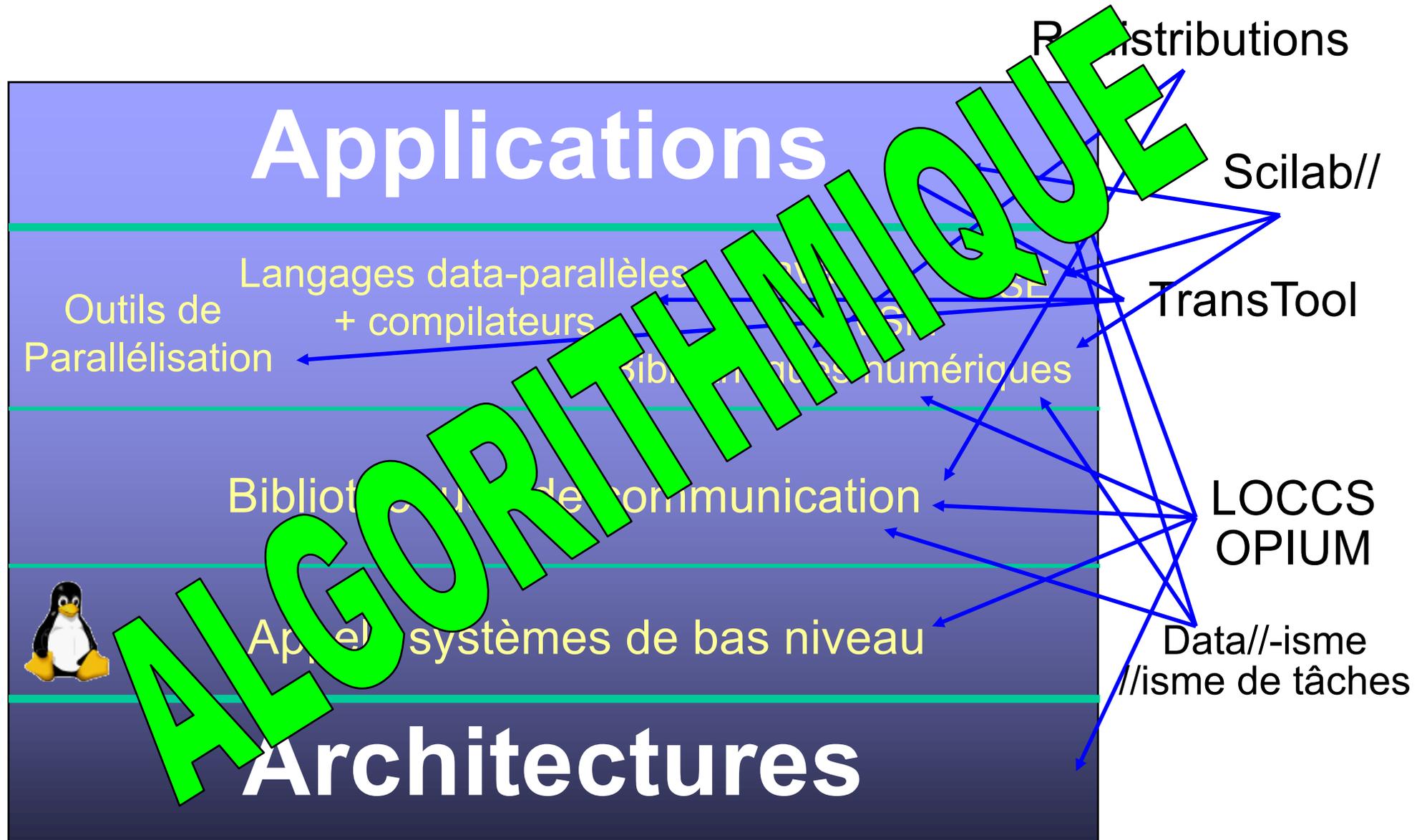
- Avoir des **environnements intégrés à hautes performances**

calcul, visualisation + couplage de codes

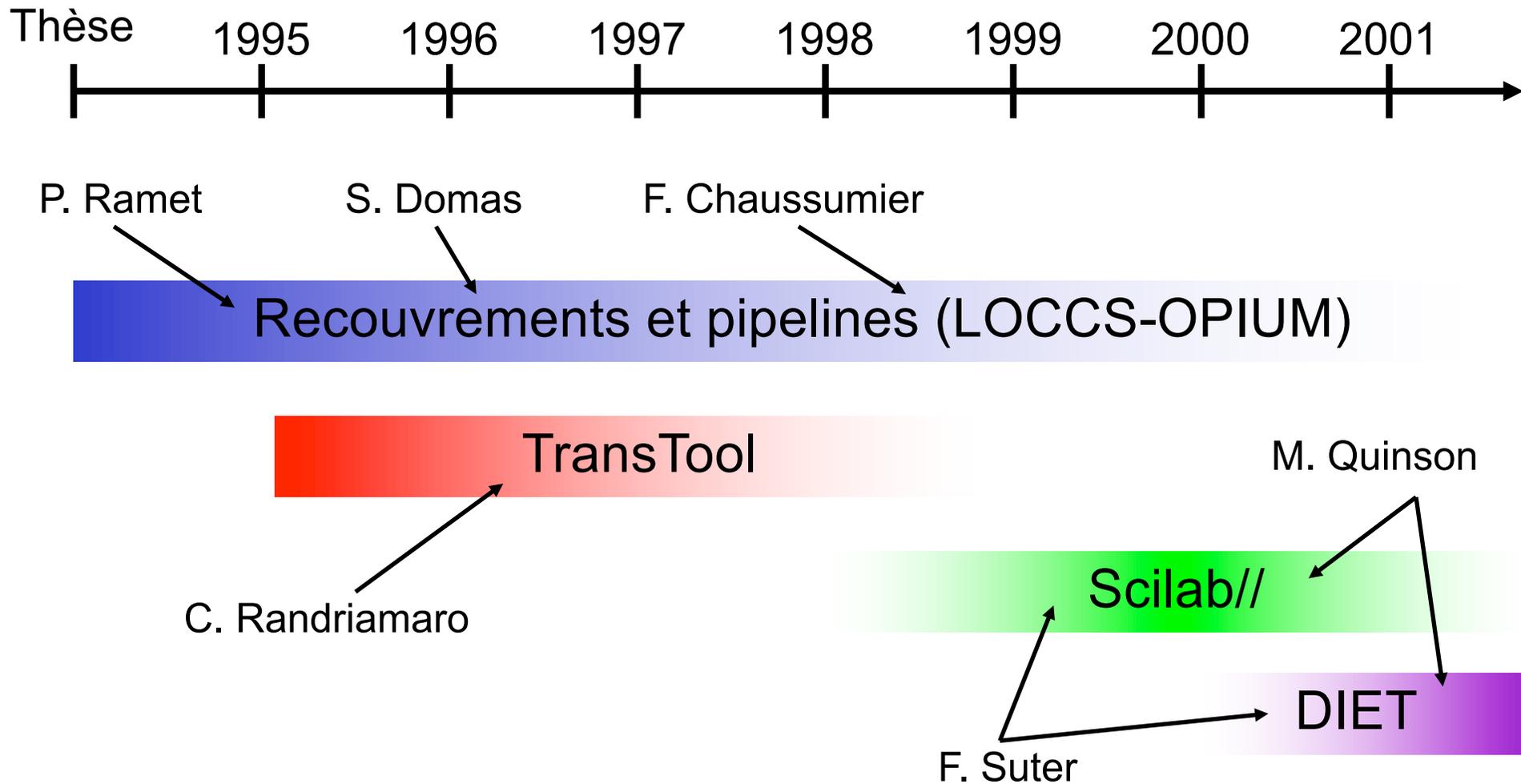


# ASP (*Application Service Provider*)

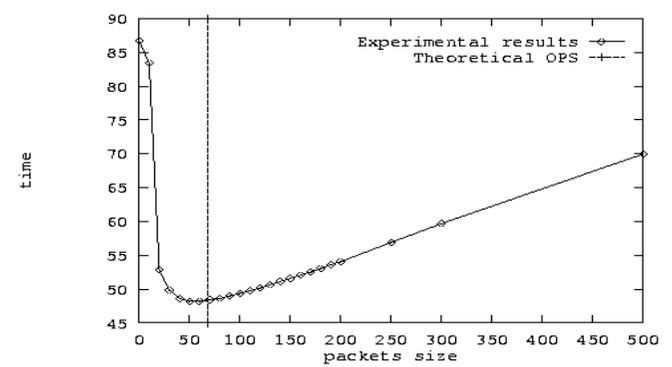
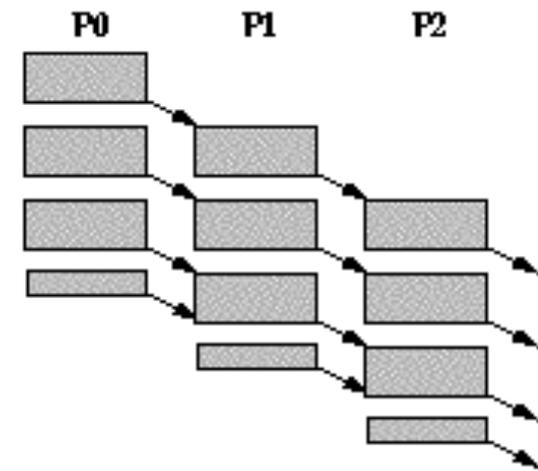
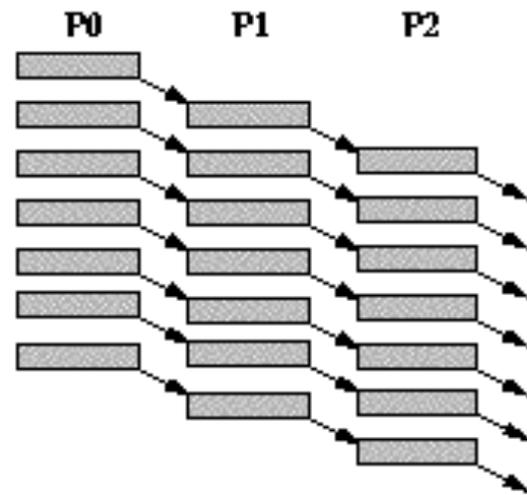
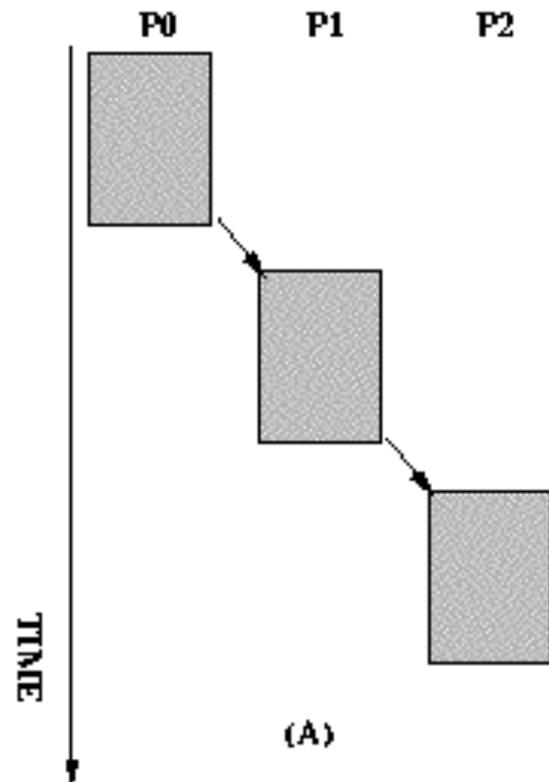
- ✓ **Une idée à long terme** : location de services de calcul via le réseau
  - Accès des performances en calcul et en capacité mémoire
  - Logiciels ou données protégés
  - Utilisation de serveurs de calcul accessibles via le(s) réseau(x)
  - Mais
    - ☹ Généralement autres domaines d'applications (BD, finances, ...)
  
- ✓ **GridRPC** (Netsolve, NINF, RCS, ...)
  - ☹ **Toujours difficiles à utiliser pour les non-spécialistes**
    - ☹ Pratiquement pas de transparence
    - ☹ Problèmes de sécurité généralement pas traités
  - ☹ **Souvent dépendants des applications**
  - ☹ **Manque de standards** (CORBA, JAVA/JINI, sockets, ...) pour construire les serveurs de calcul



# Projets de recherche et thèses

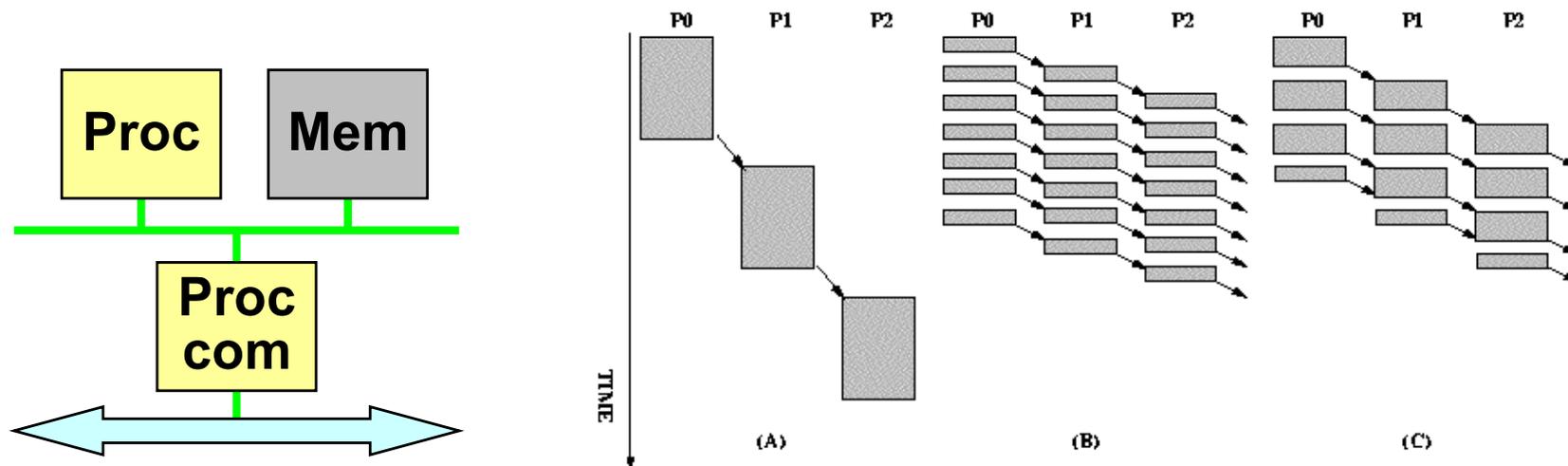


# RECOUVREMENTS ET PIPELINES

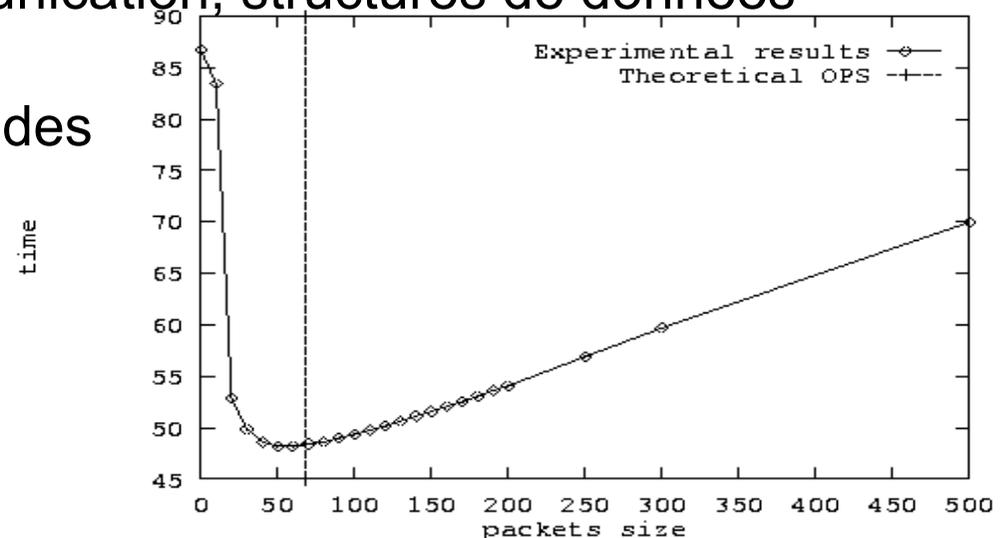


# Recouvrements et pipelines de calculs

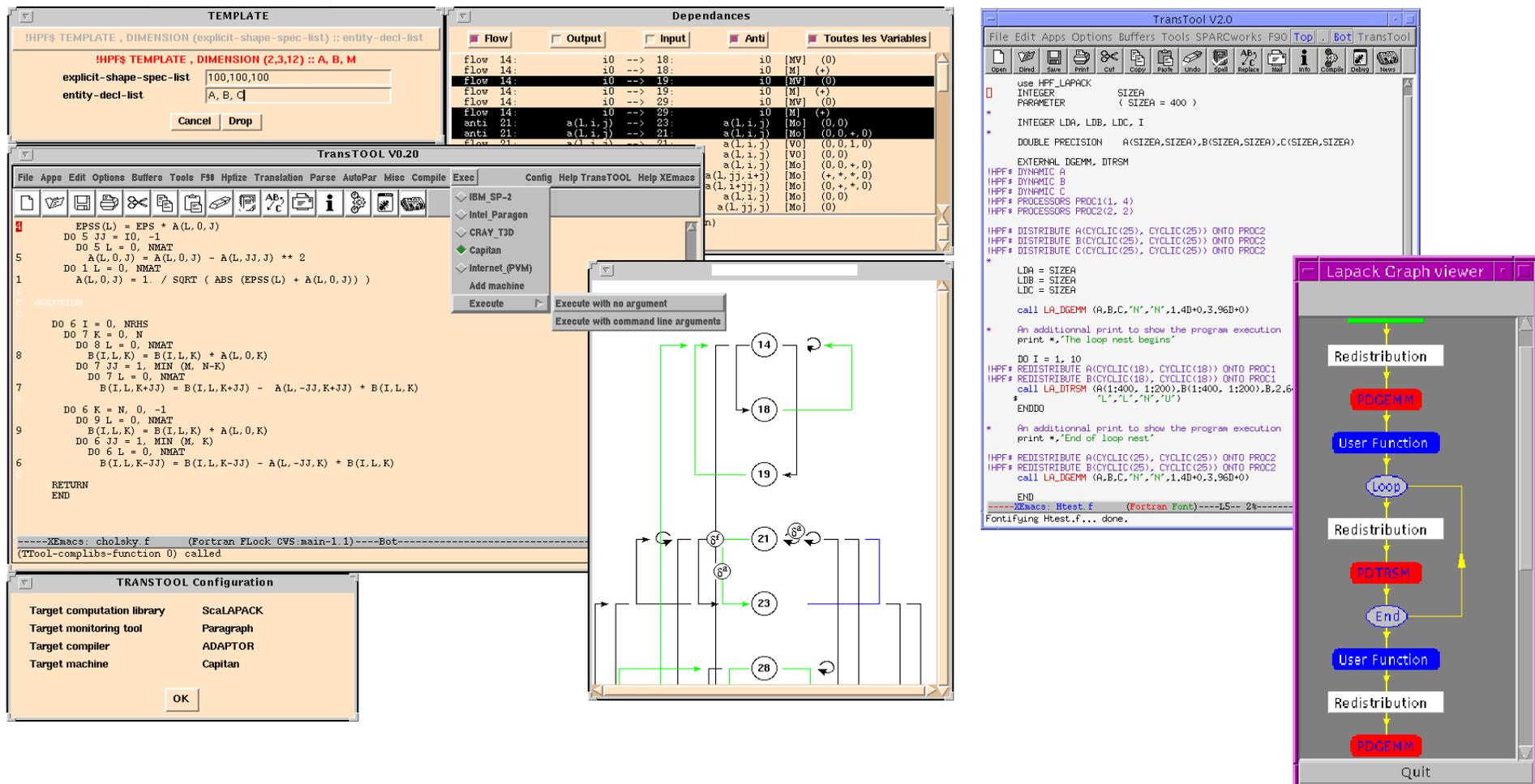
- **Buts**
  - Cacher les communications
  - Briser la séquentialité d'un noyau numérique
- Nombreuses applications dans le calcul numérique et l'imagerie
- **Deux approches combinées**
  - Recouvrements des communications (archi + support d'exécution)
  - Pipelines des calculs (algorithmique)



- **Recherche de boucles « pipelinables »**
  - Boucles *cross-processors* par analyse de dépendances
- **Analyse du gain**
  - Modélisation de l'architecture et de l'algorithme
- **Calcul du grain de calcul optimal**
  - Minimisation de la fonction de temps (OPIUM)
- **Génération de code** ou appel à une bibliothèque spécialisée
  - Différents schémas de communication, structures de données régulières ou irrégulières, supports d'exécution, gestion des buffers (LOCCS)
  - Calcul de la taille de paquets pour différentes fonctions
  - Génération de code HPF



# PARALLELISATION DE PROGRAMMES FORTRAN



The image displays the TransTool V2.0 interface, which is used for the parallelization of Fortran programs. The main window shows the source code of a program, with various annotations and comments. The code includes a main program and several subroutines, such as `LA_DGEMM` and `LA_DTRSM`, which are used for matrix operations. The code is annotated with `HPF` (High Performance Fortran) directives, such as `HPF DYNAMIC A`, `HPF DYNAMIC B`, and `HPF DYNAMIC C`, which are used to specify the distribution of data across processors. The code also includes `HPF PROCESSORS` and `HPF DISTRIBUTE` directives to specify the number of processors and the distribution of data across them.

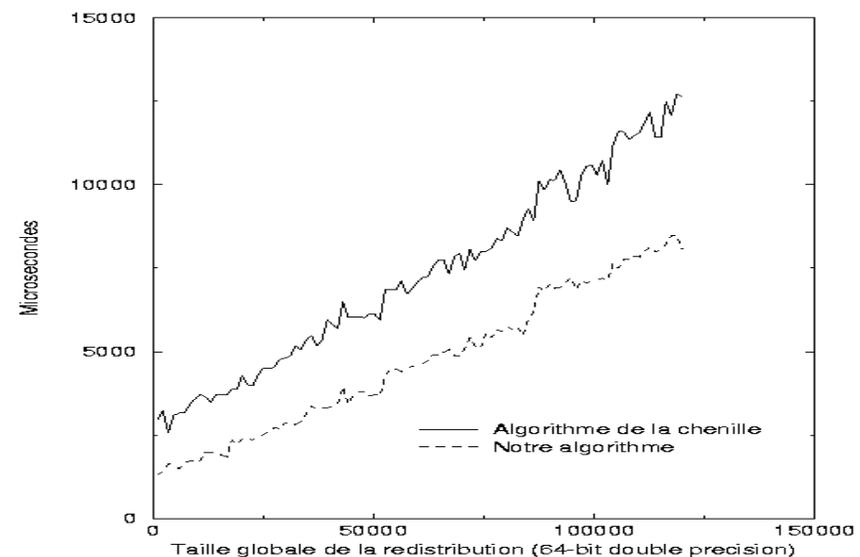
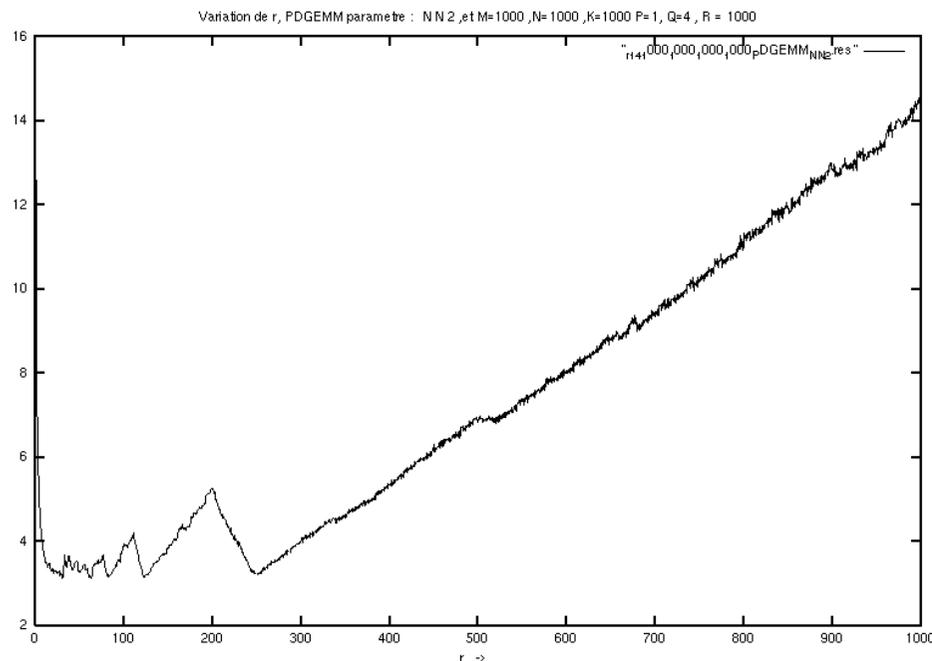
Below the code, there is a **TRANSTOOL Configuration** window showing the target computation library (ScaLAPACK), target monitoring tool (Paragraph), target compiler (ADAPTOR), and target machine (Capitan). The **Dependencies** window shows a table of dependencies between variables and data types. The **Lapack Graph viewer** window shows a flow graph of the program, with nodes representing different operations and data flows. The graph shows a sequence of operations: `Redistribution`, `PDGEMM`, `User Function`, `Loop`, `Redistribution`, `PDTRSM`, `End`, `User Function`, `Redistribution`, and `PDGEMM`. The graph is used to visualize the execution flow and to identify opportunities for parallelization.

# TransTool

- Environnement intégré de parallélisation interactive de programmes Fortran et transformations source à source
- Développé au dessus de XEmacs
- **Plate-forme de présentation de nos développements**
  - ▣▣▣▣➤ Parallélisation et distribution automatique
  - ▣▣▣▣➤ Traduction d'appels de bibliothèques numériques
  - ▣▣▣▣➤ Optimisation de macro-pipelines
- Génération de code HPF (lié au compilateur cible)
- Plusieurs API pour l'intégration de nouveaux outils (édition, analyse de dépendances, génération d'arbre syntaxique, visualisation de données distribuées)
- Projet EuroTOPS



- **Optimisation des redistributions** pour des distributions cycliques par blocs (de  $P \times Q$  avec taille de bloc  $r$  vers  $P' \times Q'$  avec  $r'$ )
  - résultats optimaux en nombre et en volume de messages
- **Traduction d'un code Fortran 77** contenant des appels BLAS et LAPACK vers
  - Un code F77 avec des appels à ScaLAPACK
  - Un code HPF avec des appels à l'interface HPF de ScaLAPACK
  - Trouver les meilleures (re)distributions

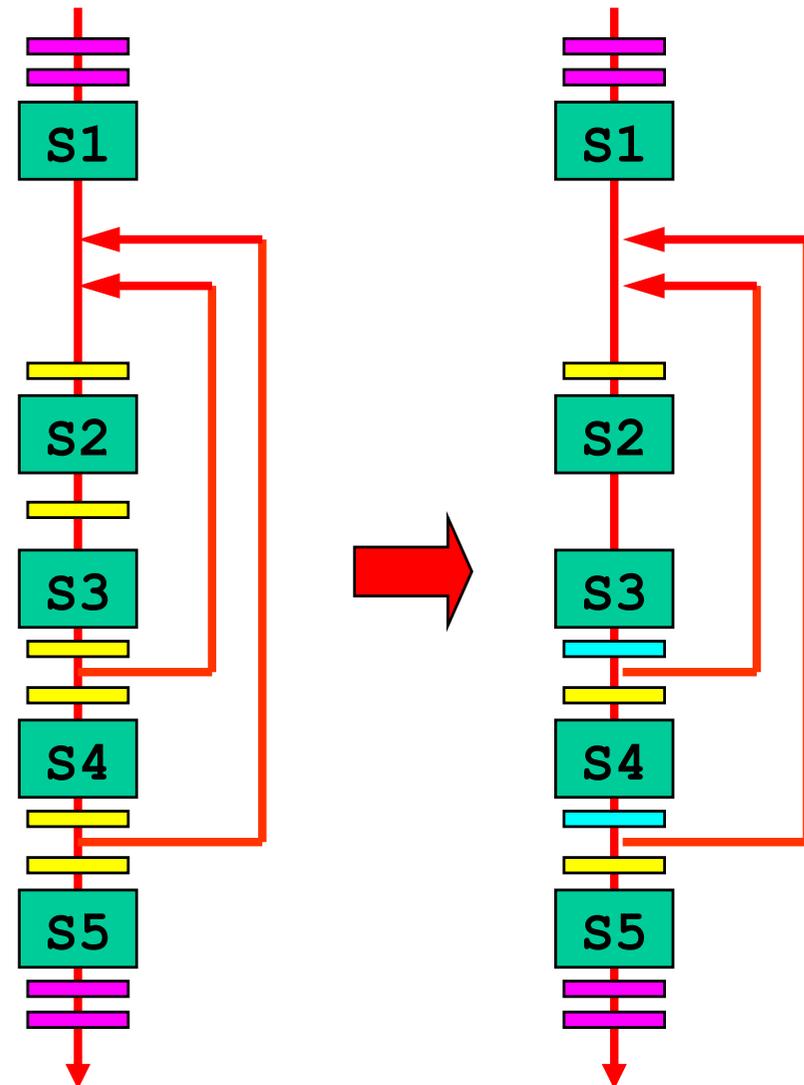


# Exemple

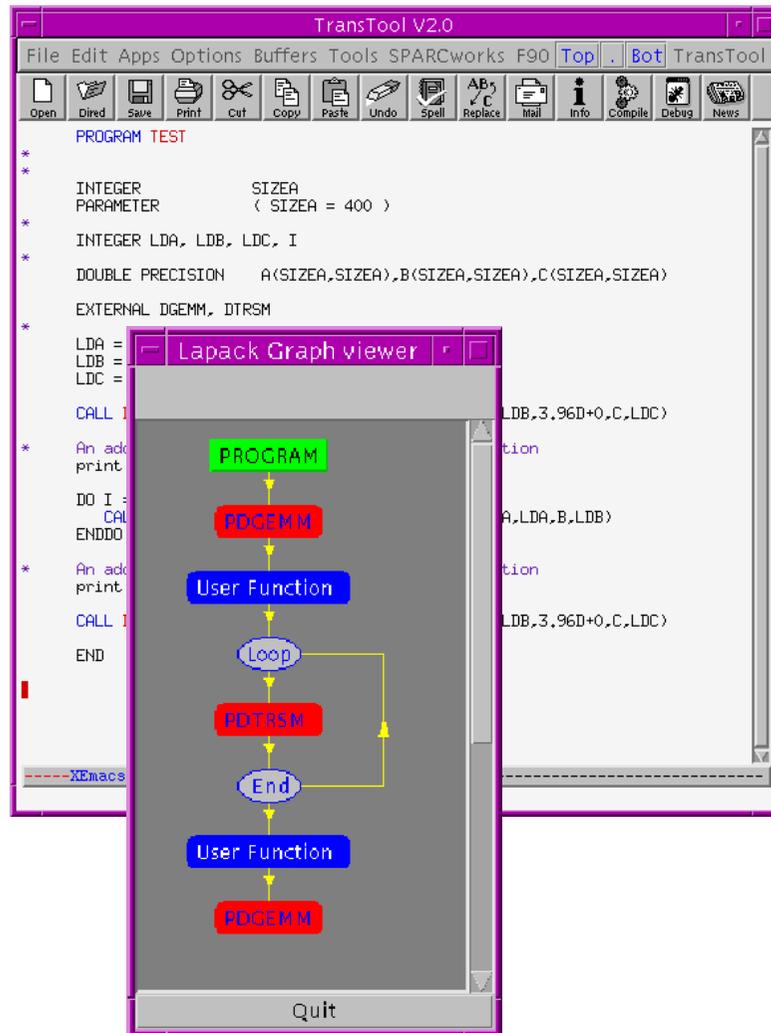
```

begin
  Initialize A
  Initialize B
  ...
  DGEMM(A, B, C)      S1
  do i=1, m
    do j=1, n
      DGEMM(C, C, D)  S2
      DGEMM(A, B, C)  S3
    enddo
    (E, F) = LU(D)    S4
    ...
  enddo
  DSYMM(F, A, C)      S5
  ...
  Gather matrices C, E, F
end

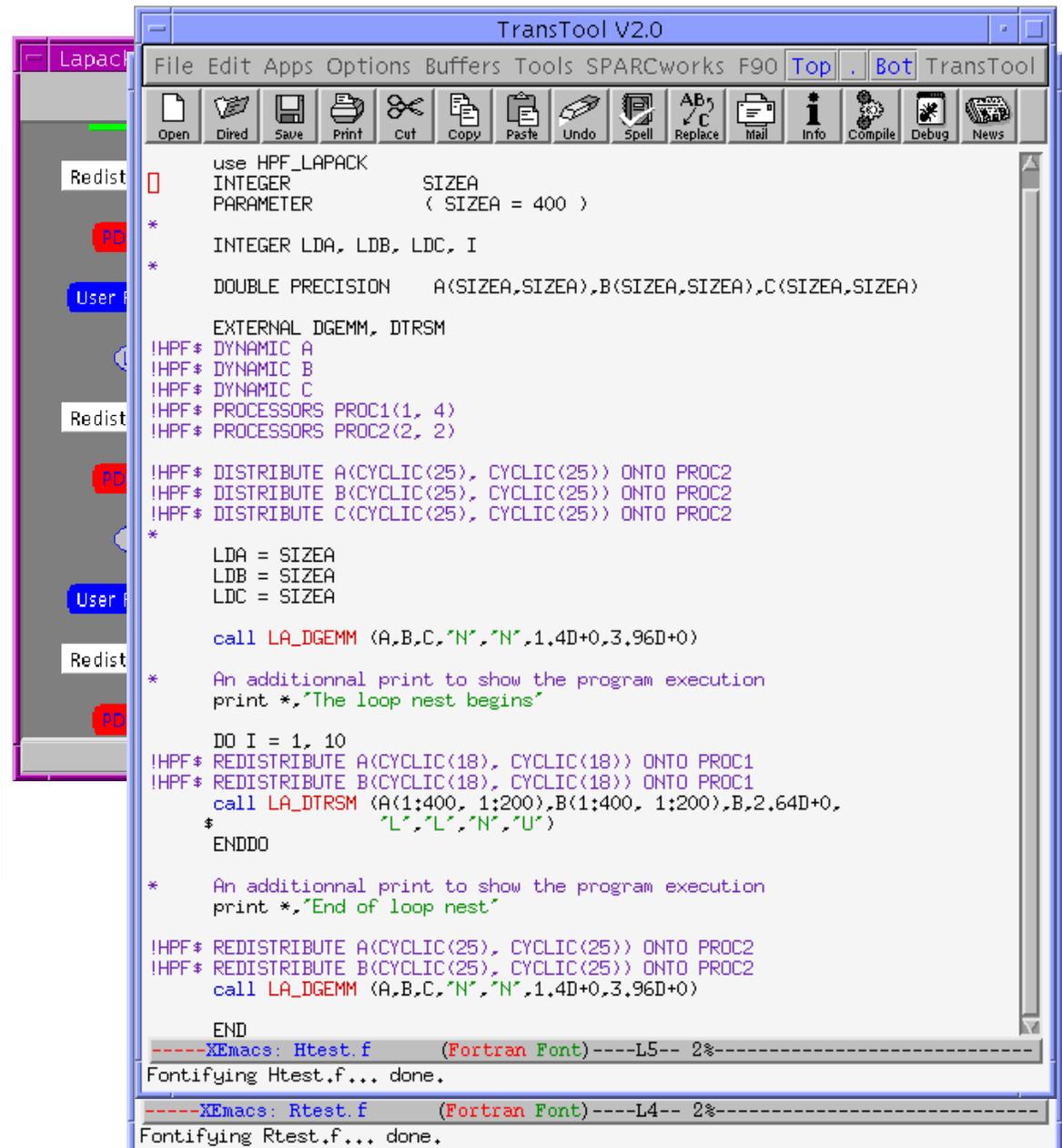
```



# Integration d'ALASca



The screenshot shows the TransTool V2.0 interface with a Fortran program named 'PROGRAM TEST'. The code includes parameter declarations, integer and double precision variables, and calls to LAPACK routines like DGEMM and DTRSM. A 'Lapack Graph viewer' window is overlaid on the code, showing a flow graph of the program's execution. The graph starts with a 'PROGRAM' node, followed by 'DGEMM', a 'User Function', a 'Loop' node, 'DTRSM', 'End', another 'User Function', and finally 'DGEMM'. The 'Quit' button is visible at the bottom of the graph viewer.



The screenshot shows the TransTool V2.0 interface with a Fortran program named 'Htest.f'. The code includes parameter declarations, integer and double precision variables, and calls to LAPACK routines like DGEMM and DTRSM. The code is annotated with comments and directives for parallel execution. The status bar at the bottom indicates that the font for 'Htest.f' and 'Rtest.f' has been fontified.

```
use HPF_LAPACK
INTEGER          SIZEA
PARAMETER       ( SIZEA = 400 )

*
INTEGER LDA, LDB, LDC, I
*
DOUBLE PRECISION A(SIZEA,SIZEA),B(SIZEA,SIZEA),C(SIZEA,SIZEA)

EXTERNAL DGEMM, DTRSM
!HPF$ DYNAMIC A
!HPF$ DYNAMIC B
!HPF$ DYNAMIC C
!HPF$ PROCESSORS PROC1(1, 4)
!HPF$ PROCESSORS PROC2(2, 2)

!HPF$ DISTRIBUTE A(CYCLIC(25), CYCLIC(25)) ONTO PROC2
!HPF$ DISTRIBUTE B(CYCLIC(25), CYCLIC(25)) ONTO PROC2
!HPF$ DISTRIBUTE C(CYCLIC(25), CYCLIC(25)) ONTO PROC2
*
LDA = SIZEA
LDB = SIZEA
LDC = SIZEA

call LA_DGEMM (A,B,C,"N","N",1,4D+0,3,96D+0)

*
An additional print to show the program execution
print *, "The loop nest begins"

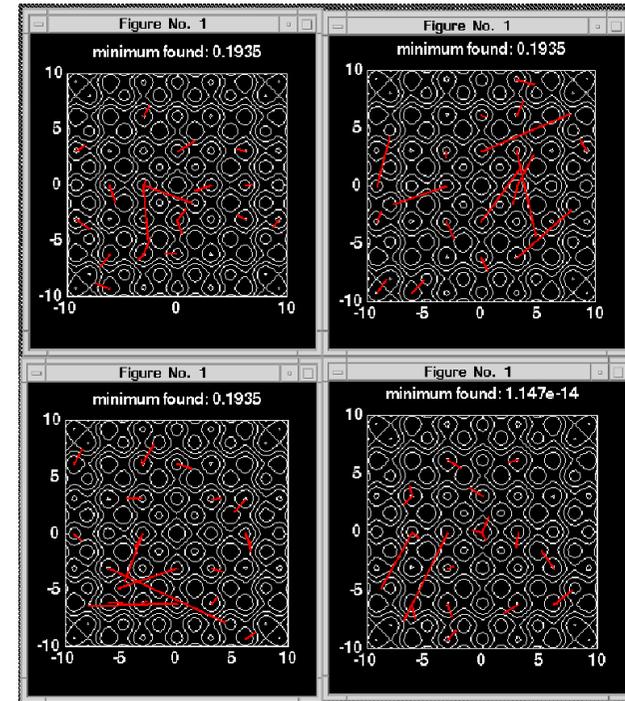
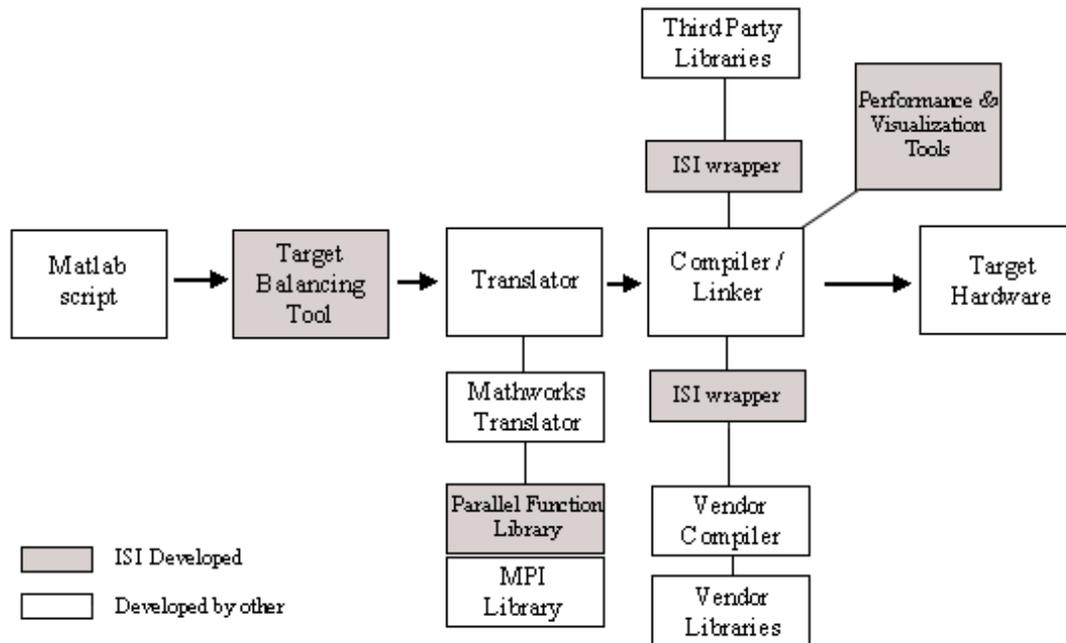
DO I = 1, 10
!HPF$ REDISTRIBUTE A(CYCLIC(18), CYCLIC(18)) ONTO PROC1
!HPF$ REDISTRIBUTE B(CYCLIC(18), CYCLIC(18)) ONTO PROC1
call LA_DTRSM (A(1:400, 1:200),B(1:400, 1:200),B,2,64D+0,
$ "L","L","N","U")
ENDDO

*
An additional print to show the program execution
print *, "End of loop nest"

!HPF$ REDISTRIBUTE A(CYCLIC(25), CYCLIC(25)) ONTO PROC2
!HPF$ REDISTRIBUTE B(CYCLIC(25), CYCLIC(25)) ONTO PROC2
call LA_DGEMM (A,B,C,"N","N",1,4D+0,3,96D+0)

END
```

# PARALLELISATION D'OUTILS de type MATLAB

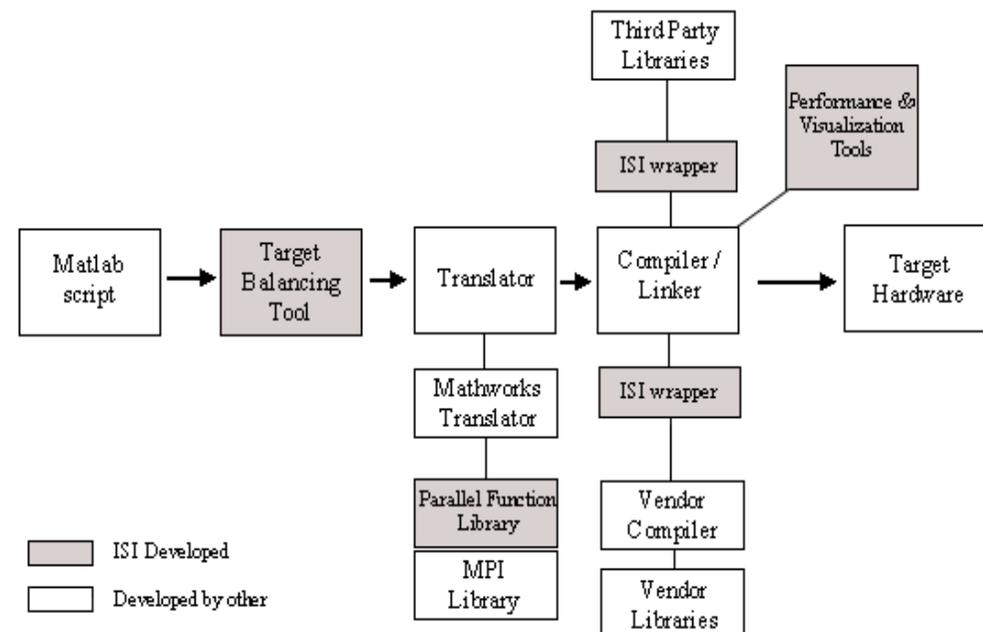


# Approche compilation

- **Entrée:** un script Matlab
- On ajoute des directives de compilation
- On le compile vers un langage de type Fortran ou C (en s'occupant des problèmes de type et de tailles de variables)
- On ajoute des appels à des bibliothèques séquentielles ou parallèles

- **Buts:**

- Eviter l'interprétation
- Plus de problèmes de types
- Appels à des bibliothèques à hautes performances
- Utiliser la technologie des compilateurs paralléliseurs



Falcon, Menhir, Match, Paradigm, Conlab, ...

# Approche maître-esclave

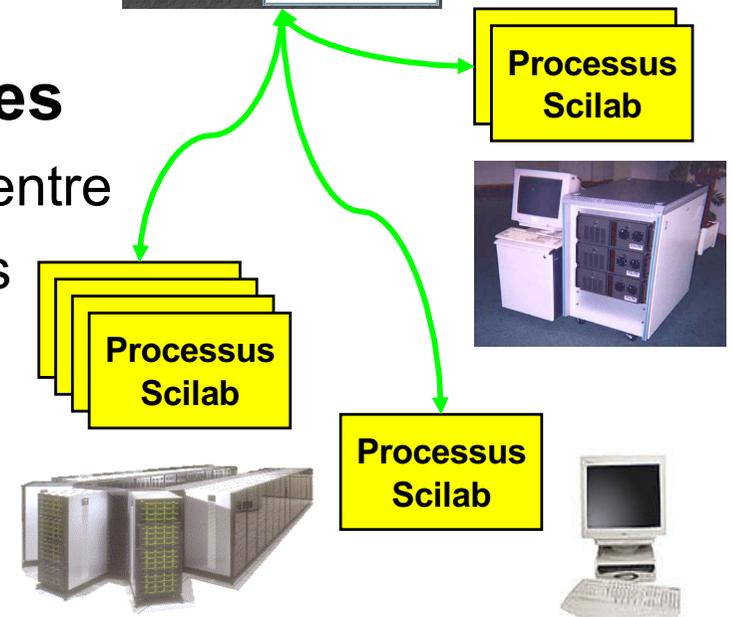
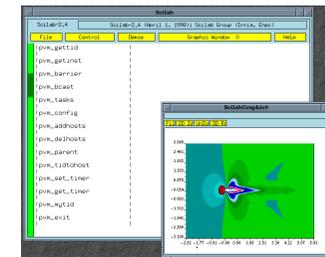
Conserver l'interactivité de Matlab

{ **Dupliquer des processus** Matlab sur chaque processeur et échanger des messages entre eux (en SPMD ou en maître-esclave)

- 😊 Facile à développer
- ☹️ Surcoût d'interprétation

} **Utiliser des serveurs de bibliothèques**

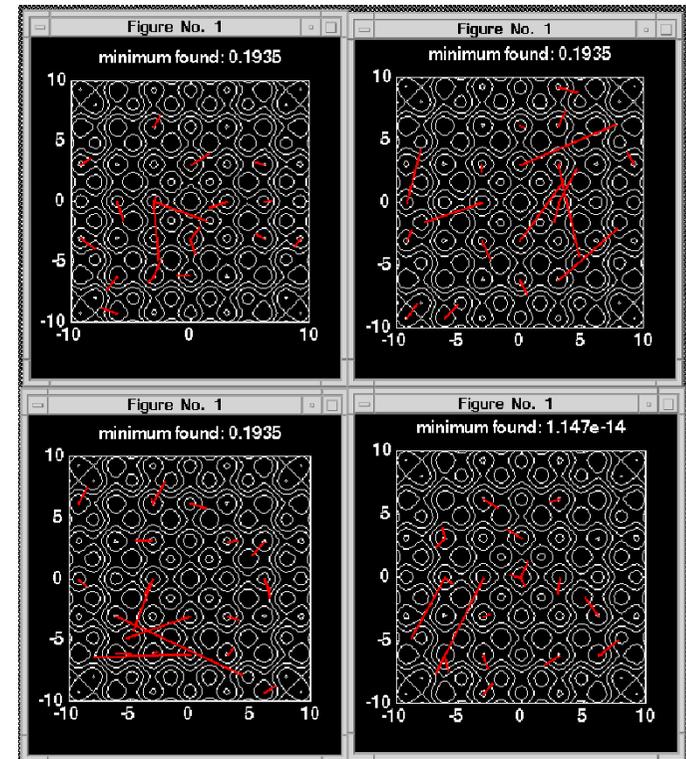
- ☹️ Besoin d'un protocole de communication entre le client (la fenêtre Matlab) et les serveurs
- 😊 Plus simple à modifier et à améliorer
- 😊 Indépendant de Matlab



MultiMatlab, PPServer, MatPar, PSI, Netsolve

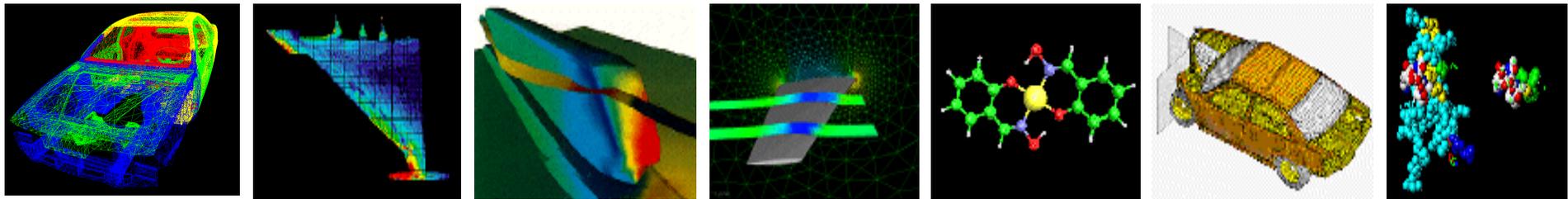
# Projets qui laissent Matlab interactif

- **MultiMatlab** (Cornell): Duplication + PVM, MPI, BLACS
- **PPServer** (MIT): Serveurs
  - ☺ Interfaces ScaLAPACK, PLAPACK et PETSc
  - ☹ Pas de transparence complète
  - ☹ Environnement homogène
  - ☹ Serveurs non-partagés
- **Matpar** (JPL)
  - Serveur ScaLAPACK
- **PSI** (U. Austin)
  - Interface PLAPACK avec serveur
- **NETSOLVE**





# Les différentes couches



SCILAB<sub>//</sub>

**PM-2  
High-Perf**

**Serveurs de calcul / Netsolve**

**Bibliothèques denses**

**Bibliothèques creuses**

**BIP-MPI / MADELEINE / PACX-MPI /  
PVM / coms Netsolve / CORBA**

**Outils d'évaluation de  
performances  
NWS, ...**



# Versions actuelles

## Passage de messages

■ PVM and MPI

☹ Assembleur du parallélisme !

## Interface Netsolve

😊 Facile à développer

☹ Netsolve limité

## Interface ScaLAPACK + duplication Scilab

■ In-core (ScaLAPACK) et out-of-core

😊 Relativement facile à développer

☹ Duplication de Scilab

☹ Types à rajouter (matrice distribuées)

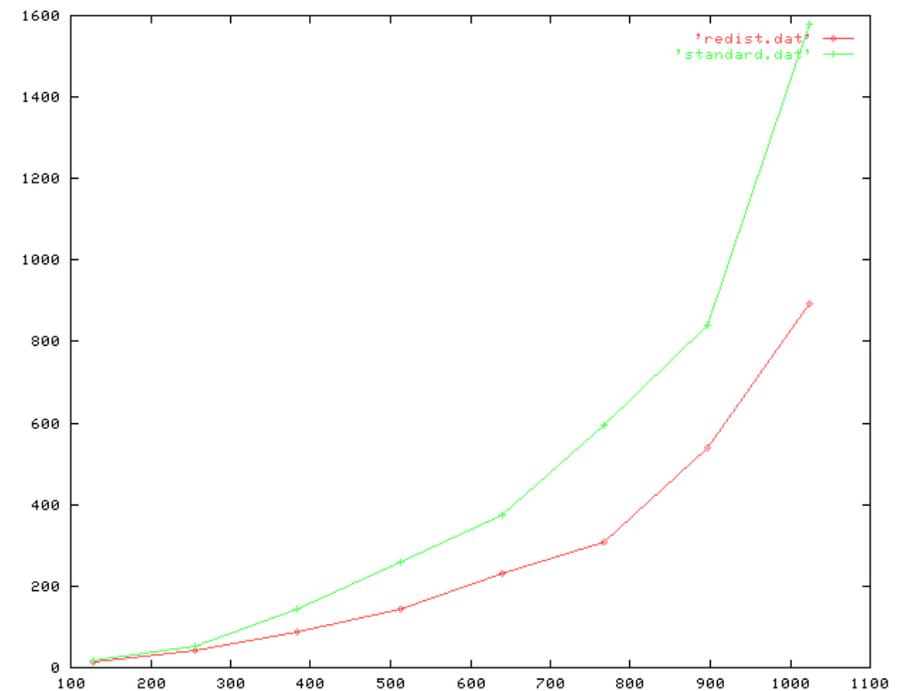
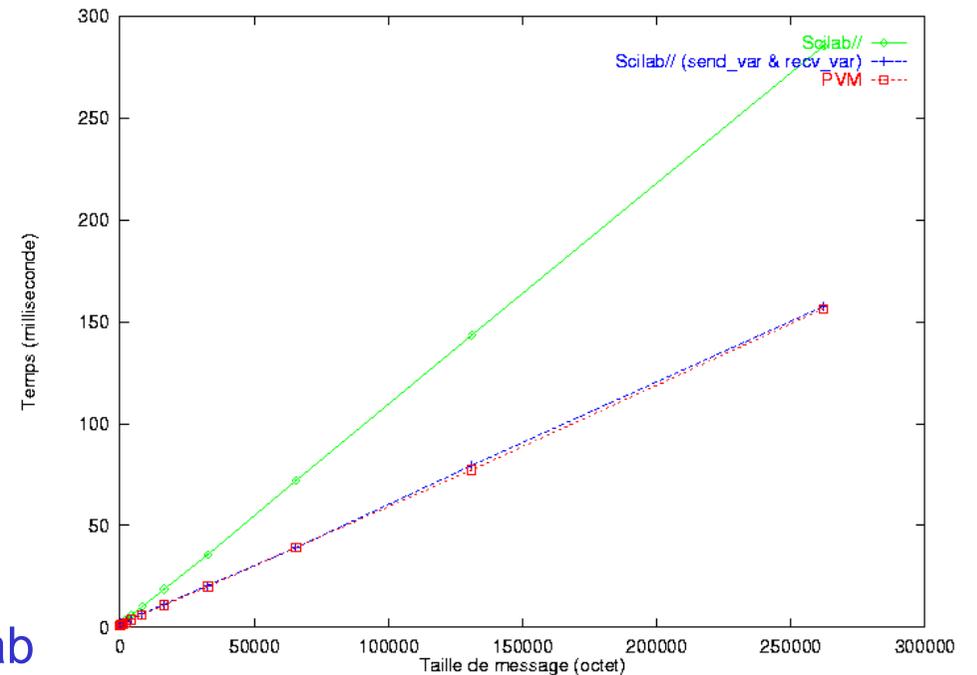
## Serveurs de calculs

😊 Style client-serveur

😊 Plus souple

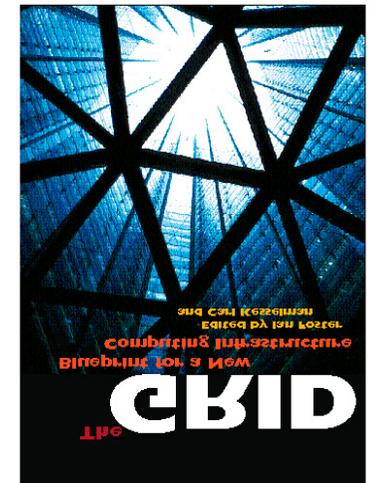
😊 Choix des couches de liaisons

😊 Utilisable avec d'autres interfaces



# Metacomputing

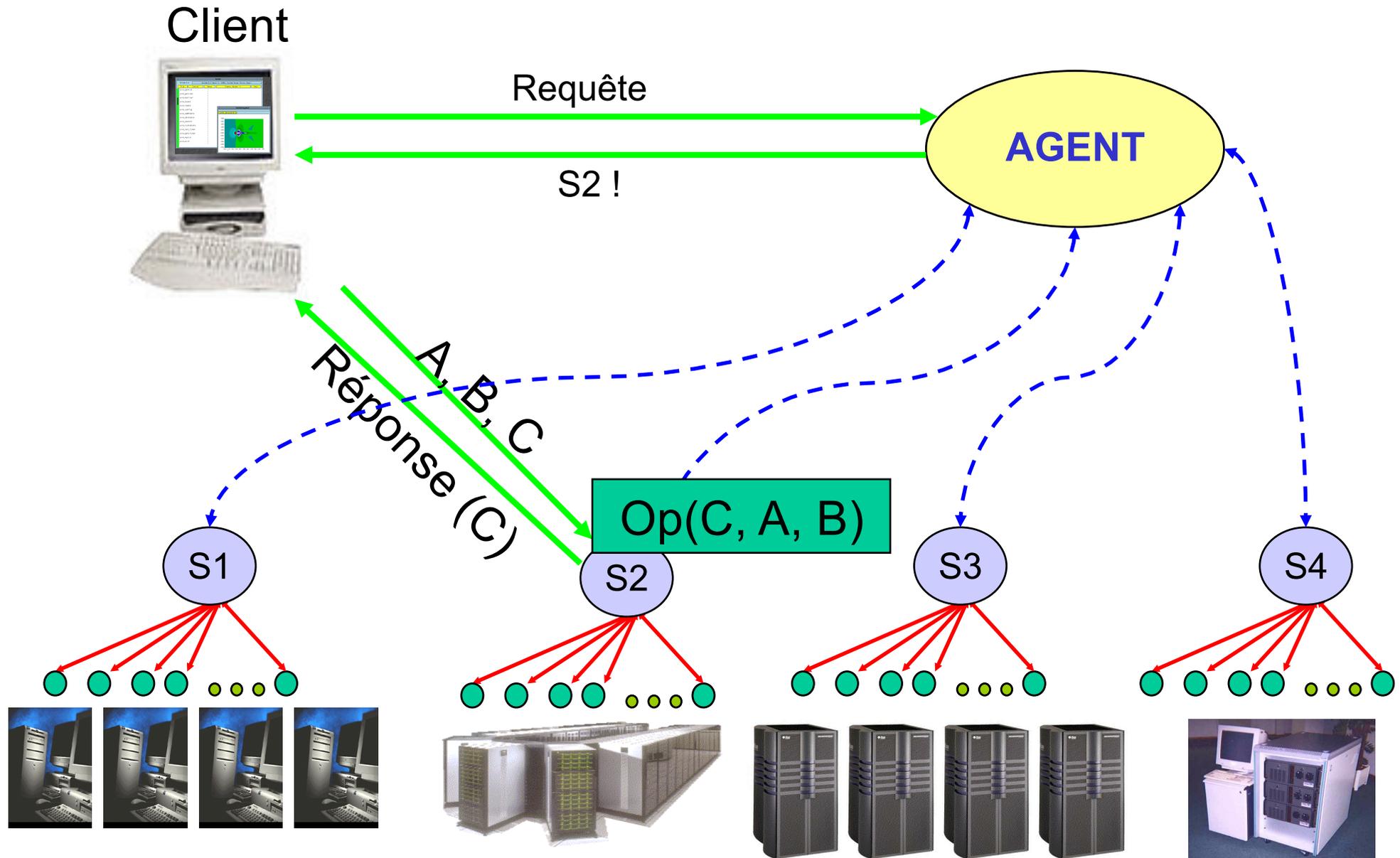
- De la recherche de la vie extraterrestre aux applications numériques les plus classiques
- De PACX-MPI à Globus
- Des grappes de grappes à l'Internet dans son ensemble
- **Idées:**
  - Utiliser les ressources inutilisées
  - Mettre en commun (et agréger) des ressources de calcul
  - Offrir à l'utilisateur de la puissance de calcul et de la capacité de stockage de manière (relativement) transparente
- **Nombreux problèmes** liés à l'hétérogénéité et à la mise en commun de ressources dispersées
- **Mais** focalisation des travaux sur les couches logicielles et « utilisation » (un peu trop) massive de Globus





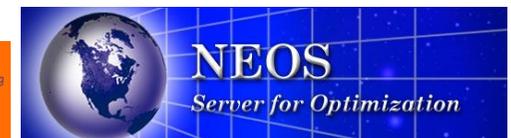
« Mon équipe a développé une solution très innovante, mais nous sommes encore en train de chercher un problème qui va avec. »

# RPC Big picture V1



# RPC et *grid-computing* : GridRPC

- **Une idée simple**
  - Implémenter le modèle de programmation RPC sur la grille
- **Fonctionnalités nécessaires**
  - Équilibrage des charges (localisation de ressources et évaluation de performances, ordonnancement), sécurité, tolérance aux pannes, IDL, distribution et migration de données, interopérabilité avec d'autres systèmes, ...
- **Outils existants**
  - **Ninf** (Electrotechnical Lab, Umezono, Japon)
    - Calcul numérique
  - **Netsolve** (University of Tennessee, USA)
    - Calcul numérique
  - **NEOS** (Argonne National Lab., USA)
    - Problèmes d'optimisation
  - **RCS** (ETH Zürich)
    - Serveurs ScaLAPACK
  - **NIMROD-G** (Monash University in Melbourne, Australie)
    - Au dessus de Globus
  - **OVM** (LRI)
    - Orienté grappes



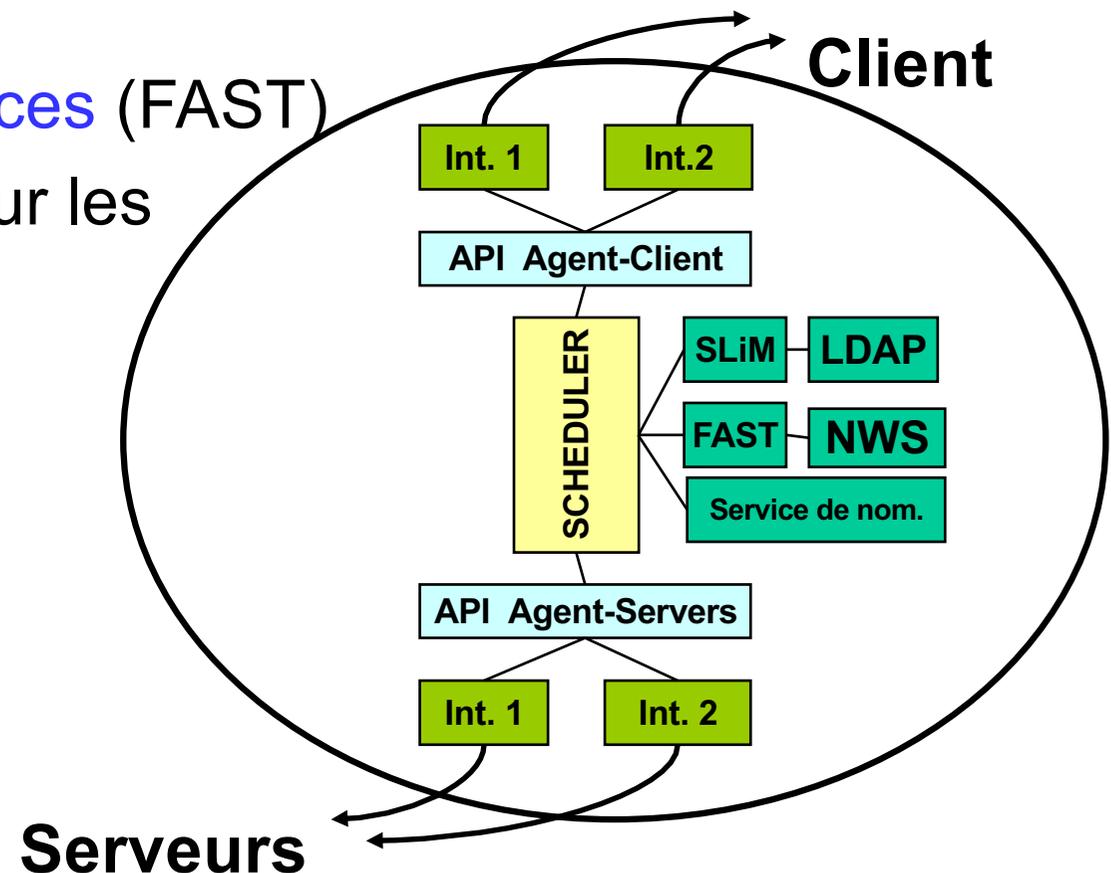
- *Distributed Interactive Engineering Toolbox*
- Bâtir une infrastructure logicielle générique pour un ensemble portails de supercalcul et de *Problem Solving Environments*
- En utilisant des couches logicielles portables (Java, JINI, CORBA, bibliothèques à hautes performances, ...)
- **Plusieurs composants**
  - Ordonnanceur à hautes performances (heuristiques dépendantes des applications)
  - Localisation de ressources
  - Evaluation et prédiction de performances (basé sur NWS)
- VTHD comme plate-forme cible
- RNTL GASP et ACI GRID GRID-ASP





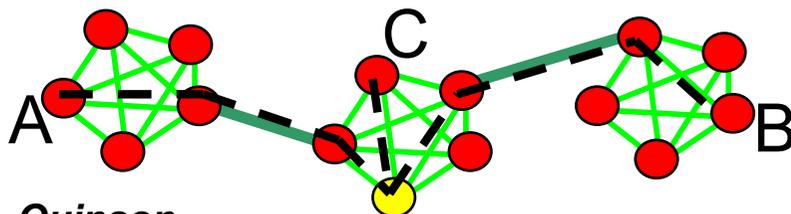
# Améliorer l'agent(s)

- Ordonnanceur
  - Heuristiques générales ou dépendantes de l'application
- Localisation de ressources logicielles (SLiM)
  - Basé sur LDAP
- Evaluation de performances (FAST)
- Service de nommage pour les données génériques
- Différentes API pour différentes interfaces
- Structure hiérarchique pour la scalabilité

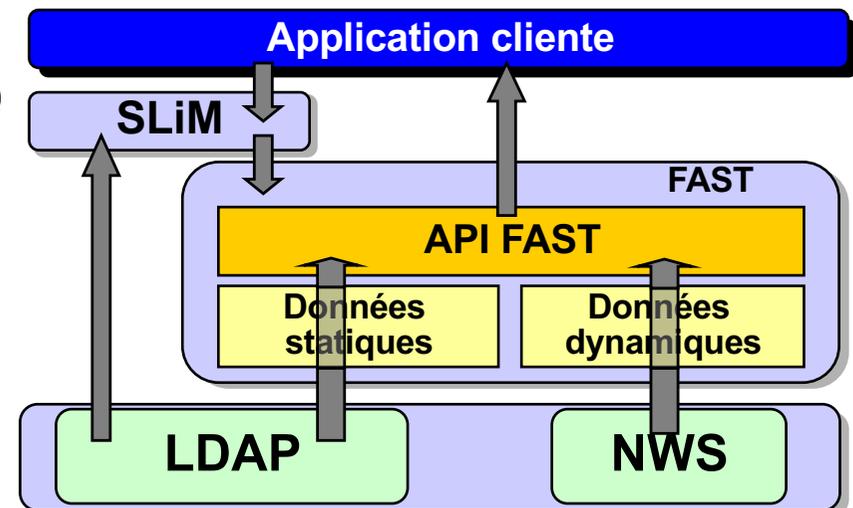


# FAST: Fast Agent's System Timer

- Evaluation des performances de la plate-forme pour être capable de trouver un serveur efficace (coûts de redistribution et de calcul) sans tester toutes les configurations → base de données de performances pour l'ordonnanceur
- Basé sur NWS (*Network Weather Service*)
- **Performances de calcul**
  - Charge de la machine, capacité mémoire et performances des queues de batch (données dynamique)
  - Tests des bibliothèques pour différentes données (statique)
- **Performances de communication**
  - Pour être capable de deviner le coût de la redistribution de données entre deux serveurs en fonction de l'architecture réseau et des informations dynamiques
  - Latence et bande-passante (hiérarchique)
- **Ensemble hiérarchique d'agents**
  - Problèmes de scalabilité

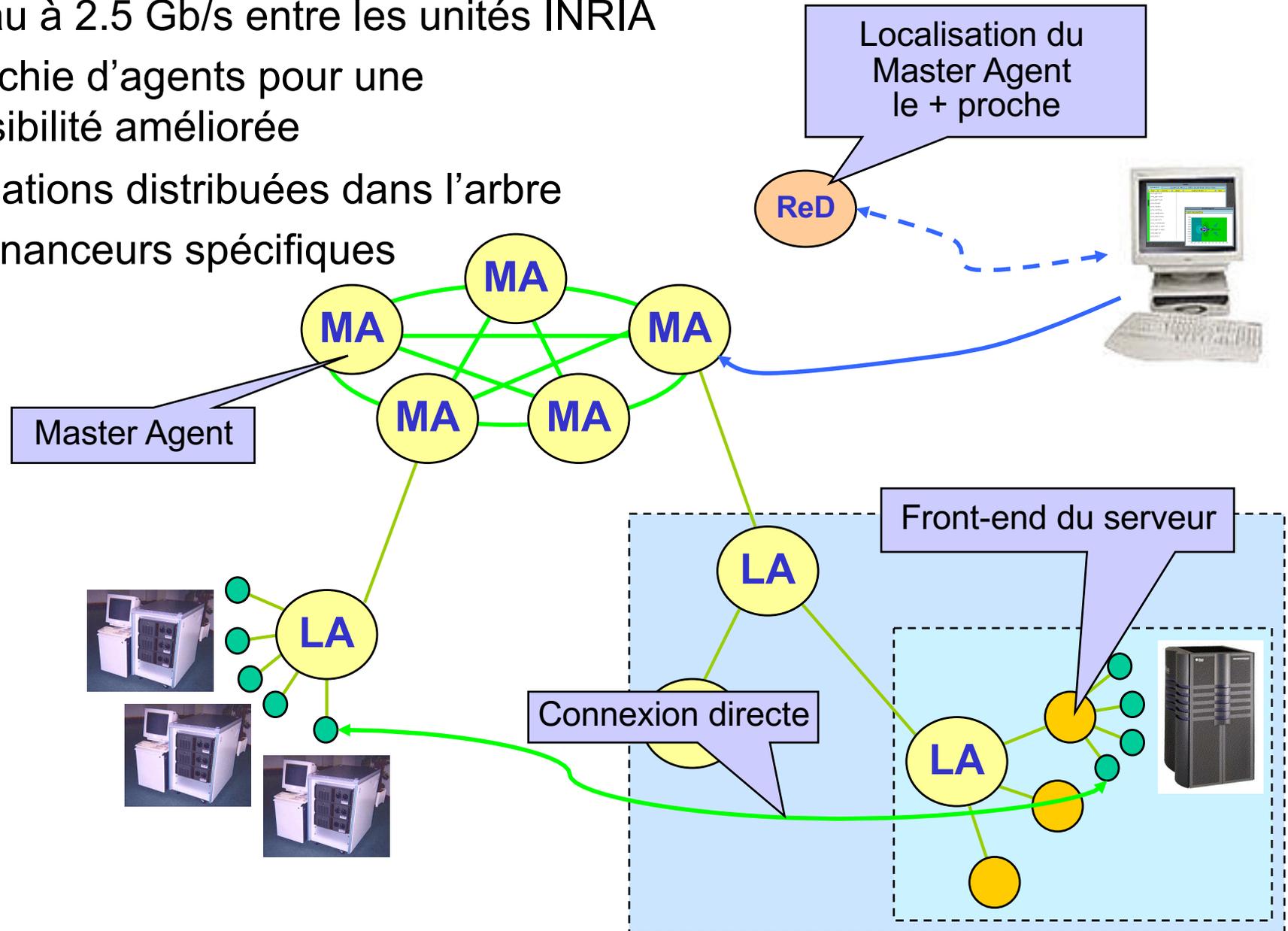


Avec M. Quinson



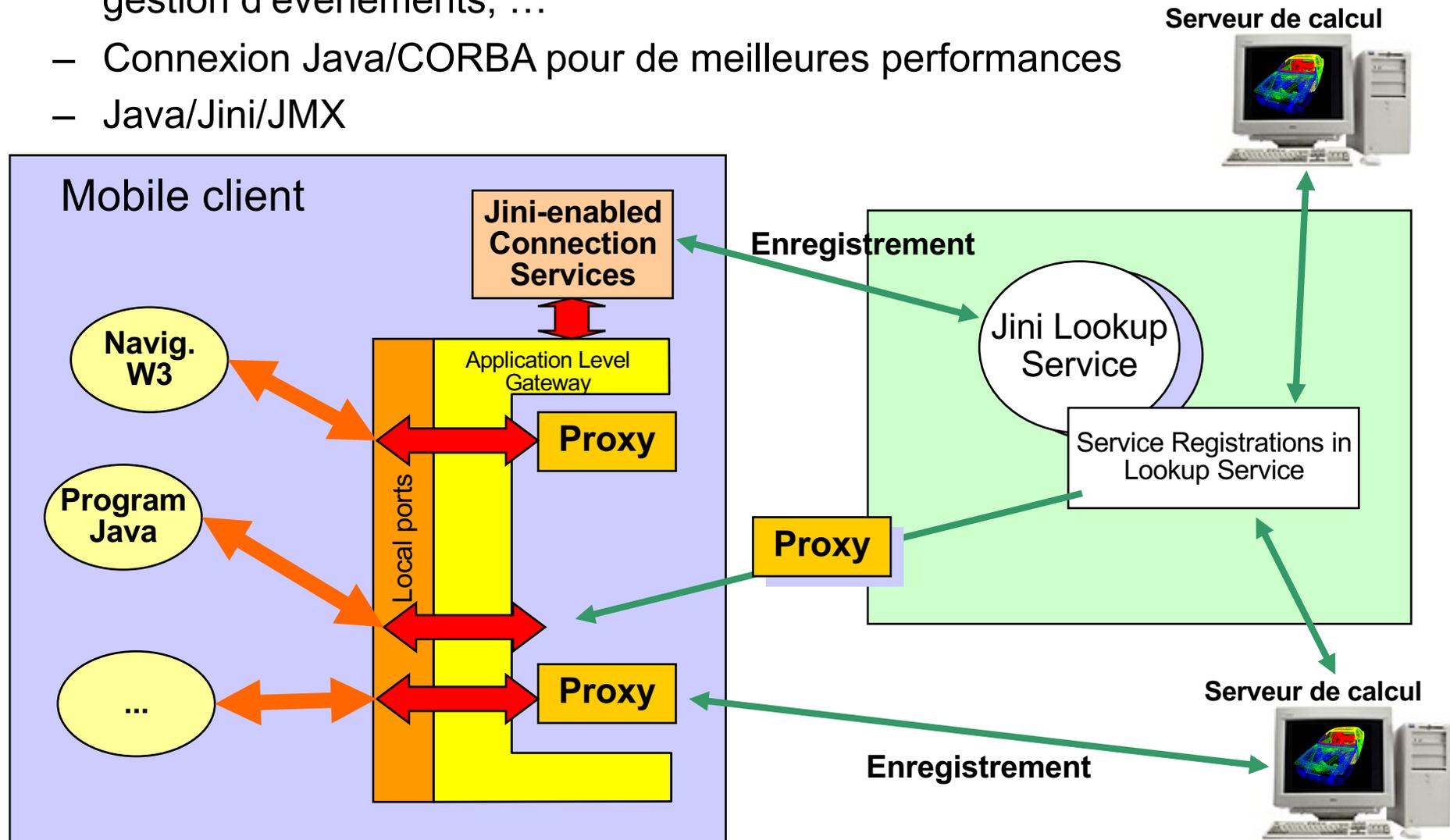
# Architecture hiérarchique pour VTHD

- Réseau à 2.5 Gb/s entre les unités INRIA
- Hiérarchie d'agents pour une extensibilité améliorée
- Informations distribuées dans l'arbre
- Ordonnanceurs spécifiques



# Et après ?

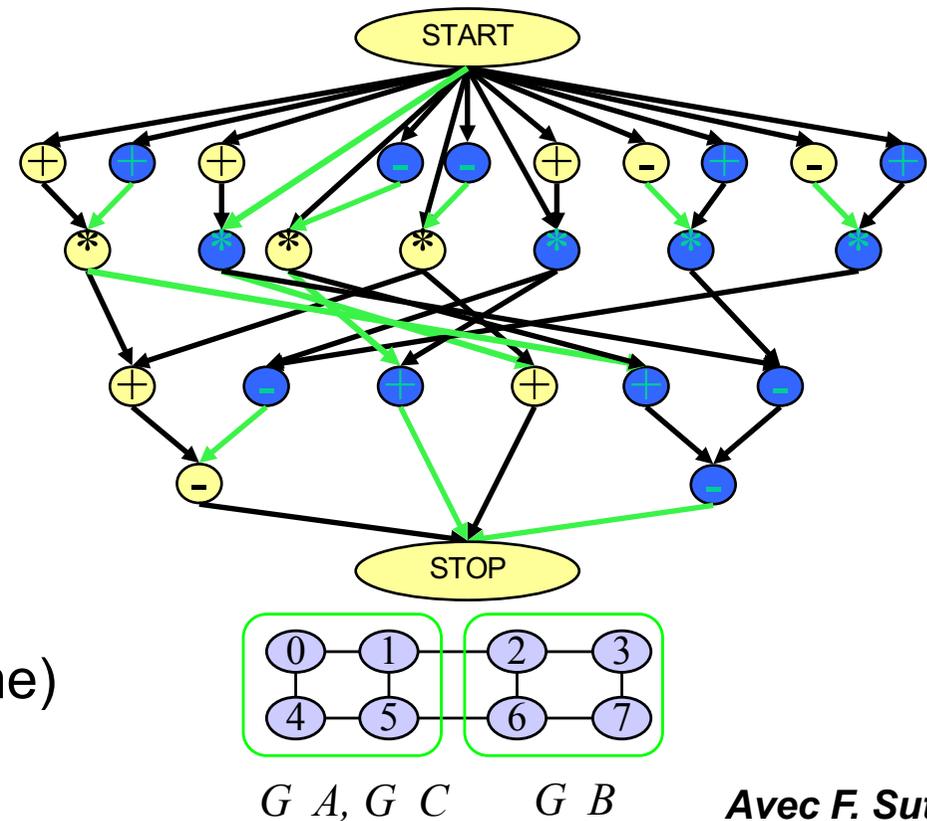
- **Technologie Java pour le *portal supercomputing***
  - Utilisation de Jini pour l'enregistrement de serveurs, lookup service, gestion d'événements, ...
  - Connexion Java/CORBA pour de meilleures performances
  - Java/Jini/JMX



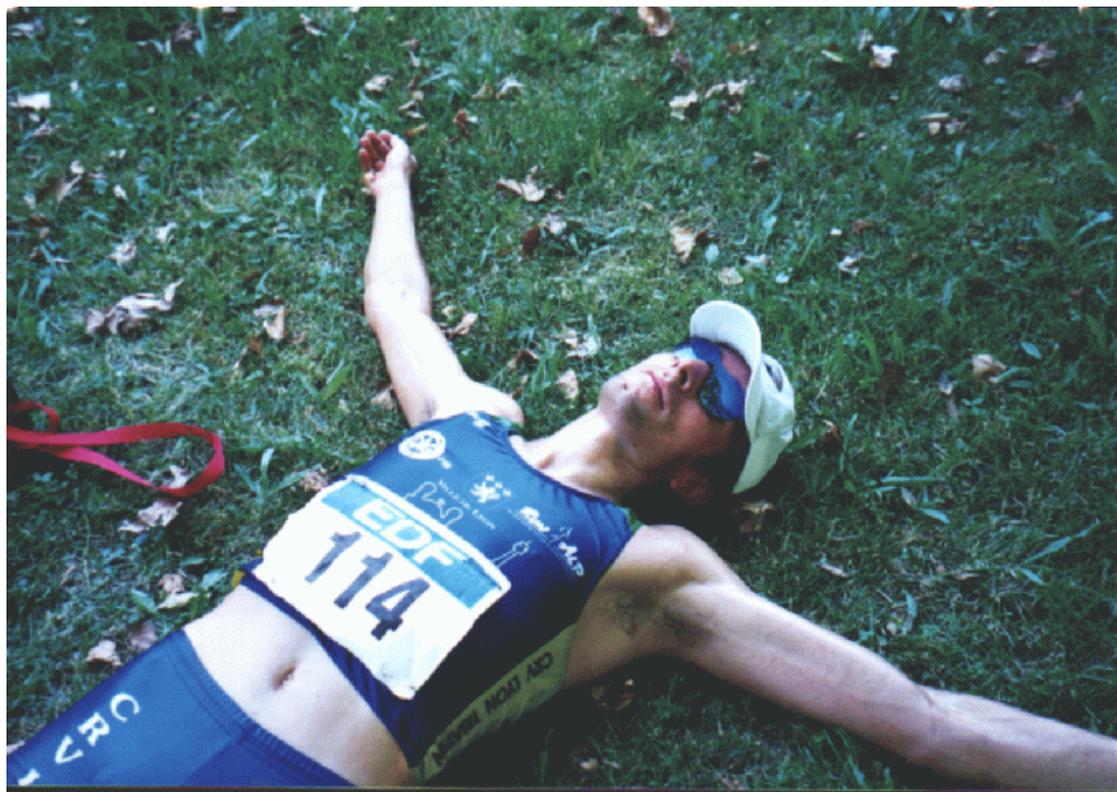
- **Sécurité**
  - Comptes utilisateurs et accès
  - Transferts de données
- **Tolérance aux pannes**
  - Serveurs ou agents
- **Intéropérabilité**
  - Description de problèmes
  - **Redistributions de données entre les serveurs**
- **Check-pointing**
  - Entrées/sorties parallèles rapides
- **Scalabilité**
  - **Hiérarchie de serveurs/agents**
- **Recherche de ressources**
  - Matériel et logiciel
- **Ordonnancement**
  - **Ordonnancement on-line d'ordonnements off-line**
- **Partage de serveurs entre utilisateurs**
  - Problèmes de sécurité
  - Lock/unlock, consistance des données, race conditions
- **Evaluation de performances**
  - **Hétérogénéité**
  - **Systemes batch**
- **Visualisation de données**
  - Problèmes d'extensibilité
- ...

# Parallélisme mixte

- Exploitation simultanée du parallélisme de tâches et du parallélisme de données
- **Entrée:**
  - Un graphe de tâches
  - Un ensemble fini de processeurs connectés par un réseau
  - Des données (potentiellement) déjà placées sur des sous-grilles
- **Sortie:**
  - Un ordonnancement des calculs
  - Un placement des données
- Première application sur Strassen et Winograd
- **Perspectives:**
  - Automatisation (heuristiques)
  - Intégration dans Scilab//
  - Approche mixte (compil., run-time)



# CONCLUSION ET PERSPECTIVES



# « Quelques » défis logiciels et autres problèmes difficiles

- **Les applications « classiques » du parallélisme supposaient**
  - Une architecture homogène de processeurs (processeurs, réseaux)
  - Des ressources statiques
  - Encore de nombreux problèmes ouverts !
- **Sur la grille**
  - Architecture totalement hétérogène (de collections d'éléments homogènes à l'Internet), ressources connues à l'exécution
  - Comportement dynamique (ajout et disparition de ressources)
  - Encore plus de problèmes ouverts
- **Trop de projets actuels oublient**
  - L'algorithmique
  - La modélisation des algorithmes

# « Quelques » défis logiciels et autres problèmes difficiles, suite

- **Hétérogénéité**

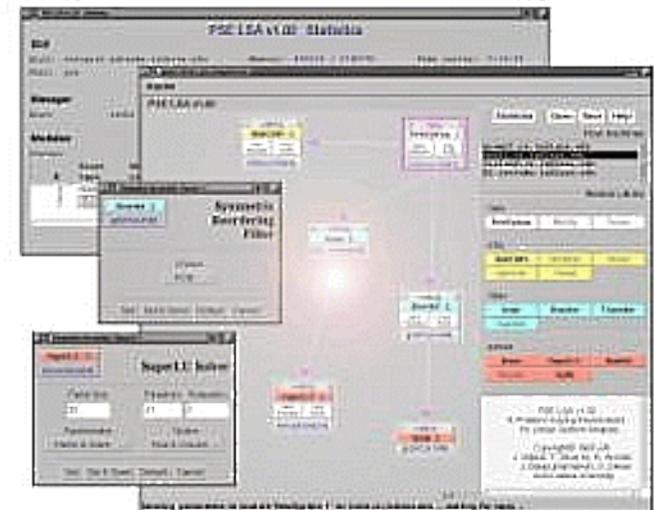
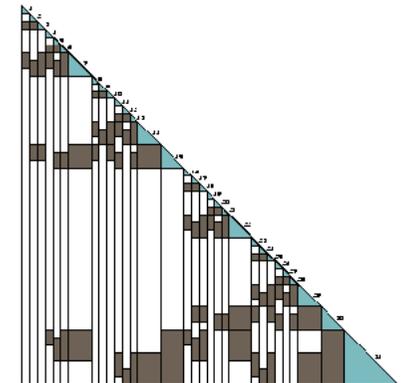
- Réseaux (couches de communications multi-protocoles), processeurs, hiérarchies mémoire profondes
- Equilibrage des charges (ordonnancement)
- Distribution des données (statique et dynamique)
- Informations sur la charge en temps réel
- Evaluation des performances
- Modélisation des architectures (déjà des modèles pour les grappes de SMP, cf P. Fraigniaud (HiHCoHP))
- Couplage de codes
- Interopérabilité de logiciels
- Bases de données de codes (GAMS, Netlib)

- **Administration**

- Sécurité (transferts, accès, firewalls)
- Équité, disponibilité, équilibrage des charges

# Quelques défis logiciels et autres problèmes difficiles, suite

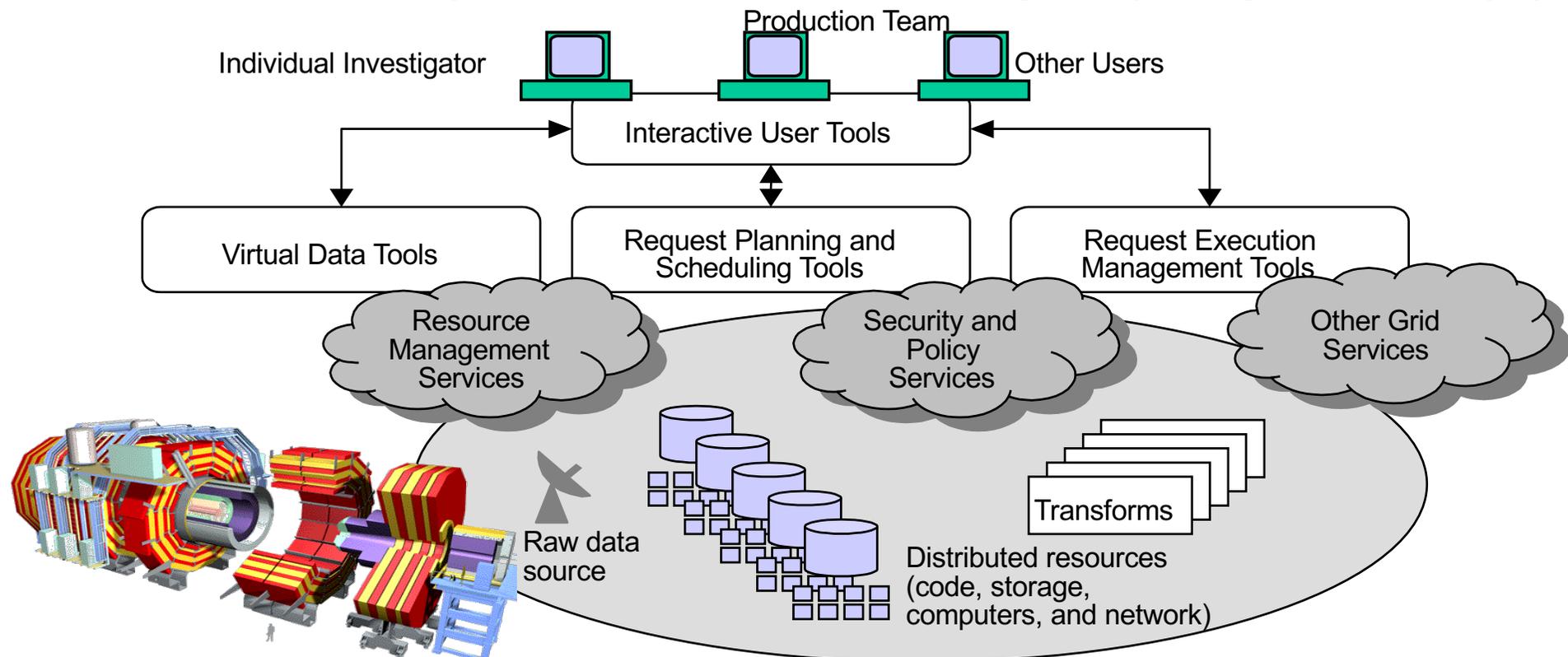
- **Algorithmique numérique et bibliothèques**
  - Algorithmique pour le creux
  - Interopérabilité entre les logiciels (redistributions de données, structures de données)
  - Méthodes hybrides (directes, itératives)
  - Liaisons plus importantes avec les supports d'exécution à hautes performances
  - Auto-adaptabilité des logiciels à tous les niveaux  
(voir ATLAS)
  - Analyse de problèmes et recherche de solveurs
  - Abstraction (programmation par composants) !



# Quelques défis logiciels et autres problèmes difficiles, suite

- **Problem Solving Environments**

- Plus de transparence, généricité, sécurité
- Unification des couches logicielles (Grid Forum)
- Couplage algorithmique et modèles avec les supports d'exécution et les architectures
- Nombreux projets applicatifs autour de la grille (Datagrid, GridPhyn)



# Conclusions et perspectives

- De l'algorithmique parallèle homogène et des aspects bas niveaux des couches de communications ....
- ... aux environnements de parallélisation et aux *Problem Solving Environments*
- Orientation vers les aspects middleware pour les PSE ...
- ... tout en regardant l'algorithmique liée au *Grid Computing* (notamment autour du creux)
- Etudier d'autres applications (MNT) et d'autres domaines (travail collaboratif, télé-enseignement, ...)
- RNRT VTHD++, RNTL GASP et ACI GRID-ASP

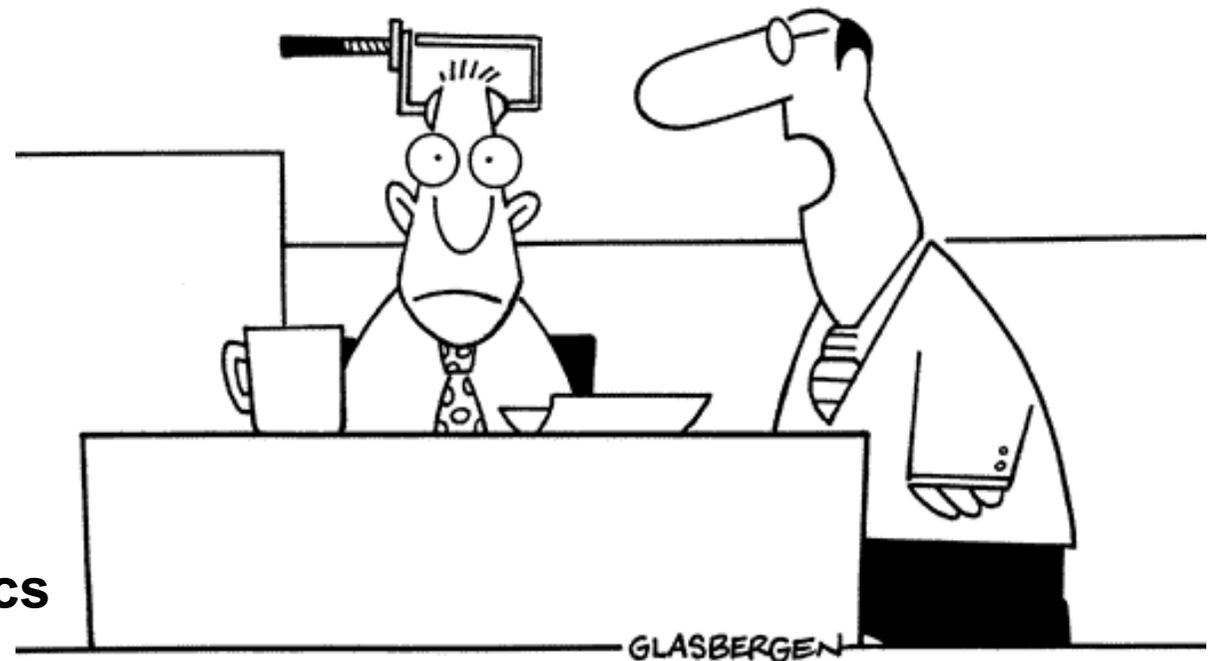


# Remerciements

- **Les étudiants en DEA et en thèse**

- Pierre Ramet
- Stéphane Domas
- Julien Zory
- Frédérique Chaussumier
- Fabrice Rastello
- Cyrille Randriamaro
- Frédéric Suter
- Martin Quinson, ...

© 2000 Randy Glasbergen. www.glasbergen.com



- **Les ingénieurs et postdocs**

- Jean-Christophe Mignot
- Eddy Caron

**"Stop whining — most people  
do their best work under pressure!"**

- **Les stagiaires de tous horizons**

- Nicolas Bert, Lionel Tricon, ...