

Traçage événementiel pour l'analyse d'applications réparties

francoisgael.ottogalli@rd.francetelecom.fr

France-Télécom R&D DLT/ASR

Jean-Marc.Vincent@imag.fr

Projet Mescal

Guest stars : Jacques Chassin de Kergommeaux (LSR-IMAG)
Benhur Stein (Univ. Santa Maria)
Eric Maillet (Astra)



Laboratoire
Informatique et
Distribution



CENTRE NATIONAL
DE LA RECHERCHE
SCIENTIFIQUE



Institut National
Polytechnique
de Grenoble



INSTITUT NATIONAL
DE RECHERCHE EN
INFORMATIQUE ET
EN AUTOMATIQUE

INRIA



GRENOBLE 1
UNIVERSITÉ
JOSEPH FOURIER
SCIENCES, TECHNOLOGIE, MÉDECINE

Expérimentation et mesure






Je dis souvent que lorsqu'on peut mesurer ce dont on parle et l'exprimer en chiffres, on en sait quelque chose; en revanche, si on ne peut l'exprimer en chiffres, on en a une bien piètre connaissance.

Lord William Thomson Kelvin (1883)

On pourrait déterminer les différents âges d'une science par la technique de ses instruments de mesure.

Gaston Bachelard

Organisation

-  Motivations
-  Modèles et mesure
-  Mise en œuvre : application à objets répartis
-  Exemple d'analyse : serveur de données
client / courtier / serveur
-  Perspectives

Motivations

Contexte des applications parallèles/distribuées

Spécifications qualitatives :

→ Le résultat est correct (preuve, tests, ...)

Spécifications quantitatives :

→ Le résultat est-il obtenu dans de bonnes conditions

Identification du problème

→ Déboguage

→ Analyse de performances

Modification du code source / librairies / OS / architecture ...

Objectif scientifique

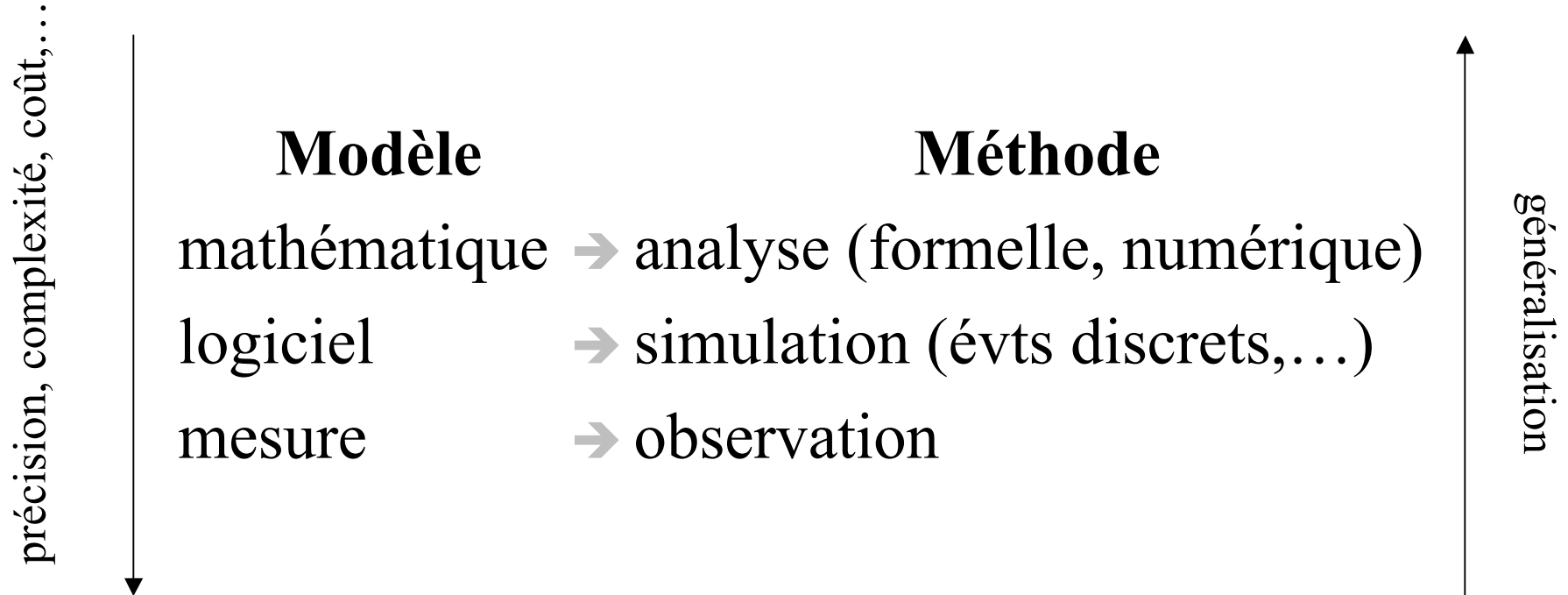
Comprendre le **comportement** dynamique d'une application **distribuée**

- identification de comportements (patterns répartis)
- vérification (aide au débogage)
- évaluation des durées
- analyse globale de l'exécution
- évaluation des performances
- politique de contrôle (monitoring)
- évaluation du coût

Analyser l'**utilisation** des **ressources**

- applicatif
- support exécutif
- système d'exploitation
- architecture matérielle

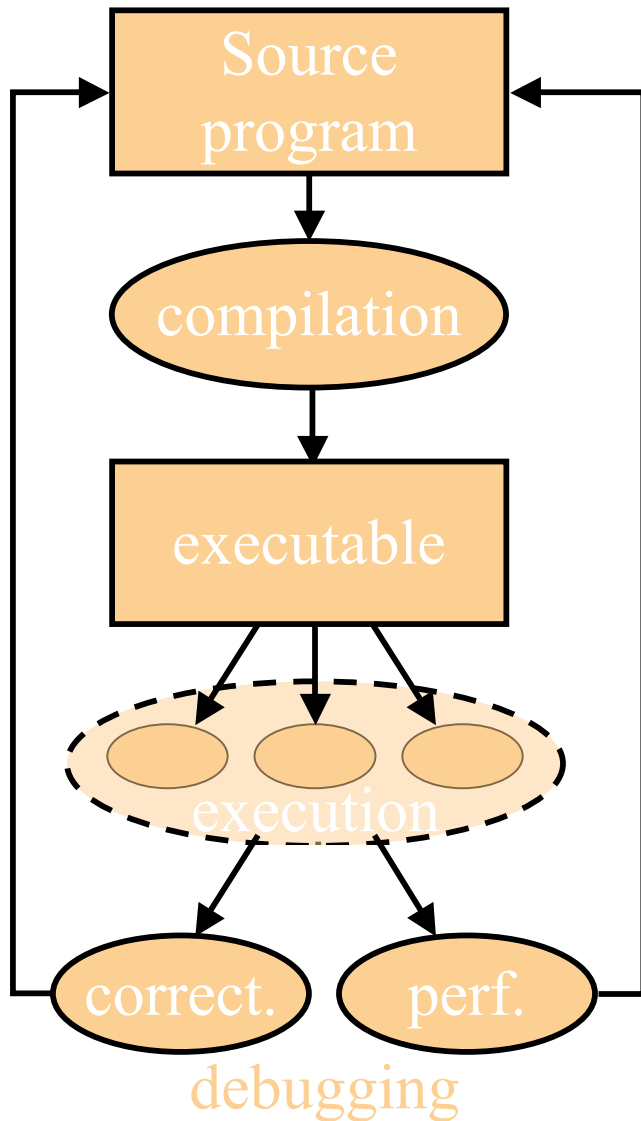
Approches



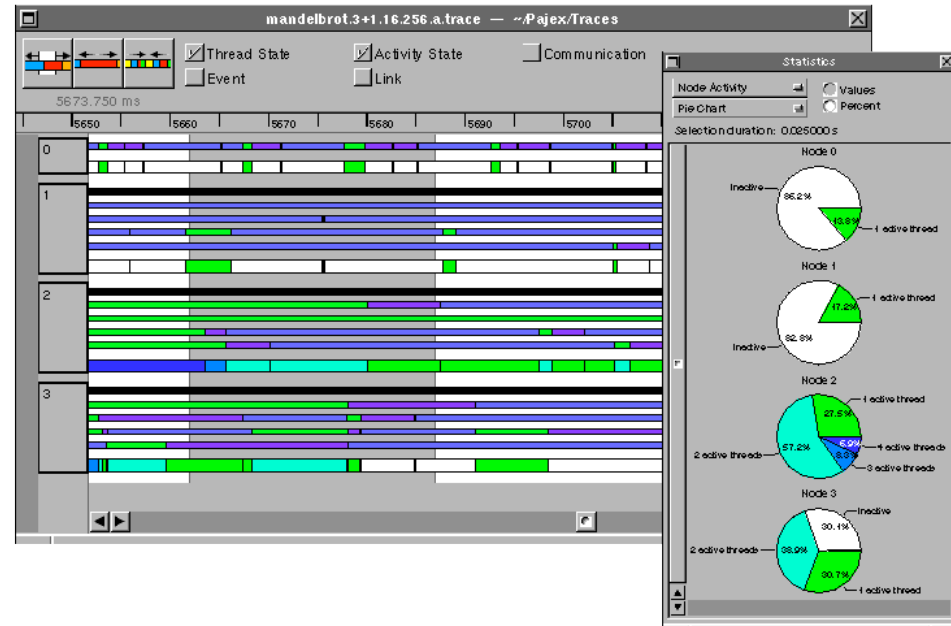
Remarque : méthodes hybrides (émulation, générateurs de trafic,...)

Expérimentation → Plan d'expérience
choix des paramètres, facteurs,...

Distributed program development cycle



Visualization



Context

(Potentially large size) parallel/distributed applications.

Executing on (potentially large size) distributed systems:

- Clusters of PCs.
- Grids of processors and clusters.

Keypoints

Distributed **heterogeneous** resources

Dynamicity of the architecture

Scalability of the middleware

Objective

Help users find performance errors:

Lack of parallelism, bottlenecks, overheads.

Behavior analysis methodology

Execution model

Measurement environment

Visualization model

Performance debugging tools

Monitoring to gather performance data

Data analysis:

- Collection of raw data

- Reconstruction (simulation) of the behavior of the application and computation of performance indexes.

- Presentation (visualization).

Existing visualization tools

Paragraph:

Large number of visualizations.

Not easily extensible/ Not interactive.

Pablo:

Easily extensible.

No behavioral visualization / Not extensible.

Paradyn:

Scalable. Automatic research of errors.

No behavioral visualization.

ParaGraph

The first widespread public domain visualization tool.

Execution traces recorded by the parallel programming library **PICL**.

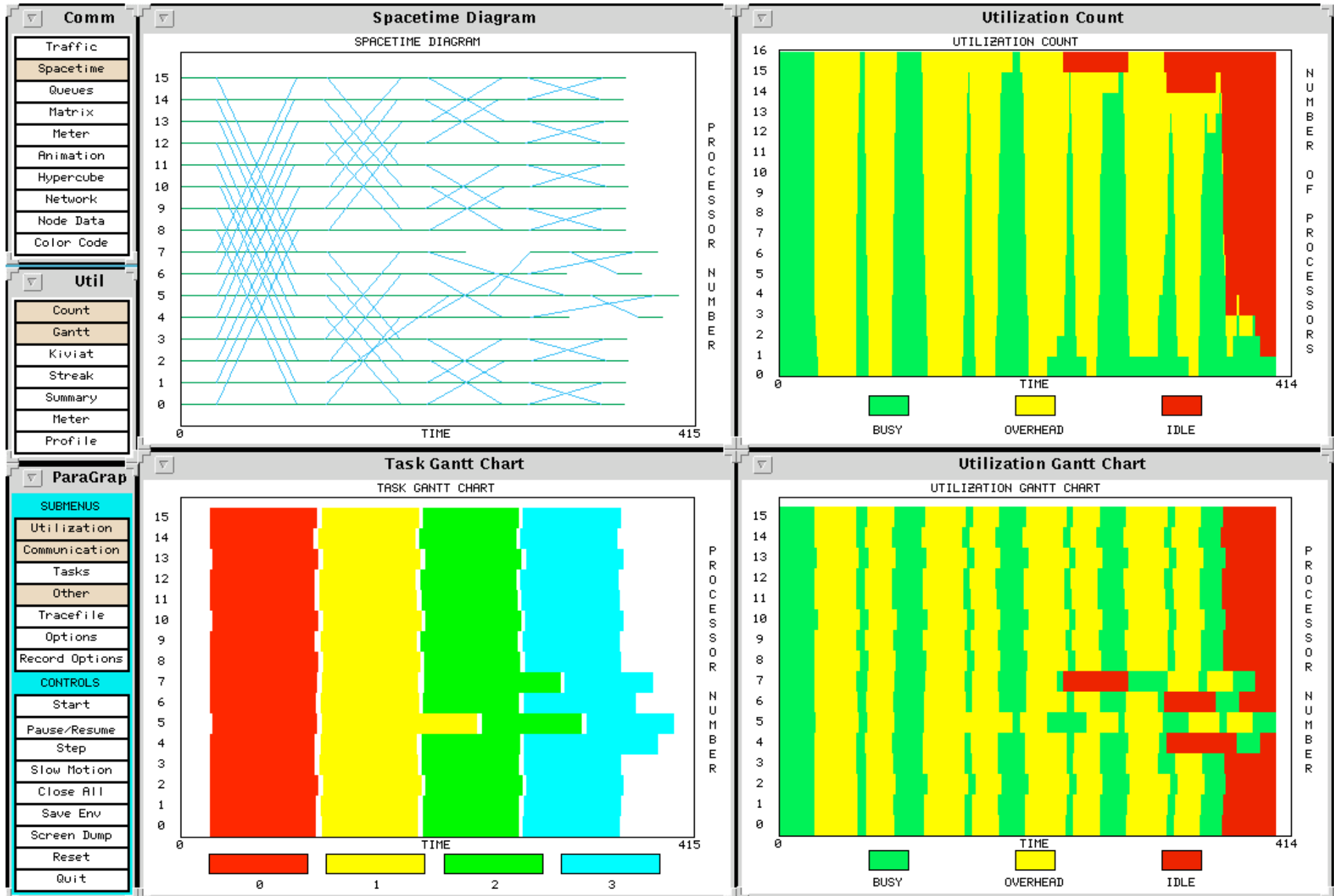
Large number of visualizations:

Processor utilization.

Communications: fréquence, volume, distribution.

Various: critical path, statistics, etc.

Some Visualizations of ParaGraph



Limitations of ParaGraph

Most visualizations are "instantaneous".

Lack of scalability:

Visualizations become cryptic for a couple of dozens of processors.

No filtering: impossible to visualize a subset of the system.

Limited extensibility: monolithic implementation.

PABLO

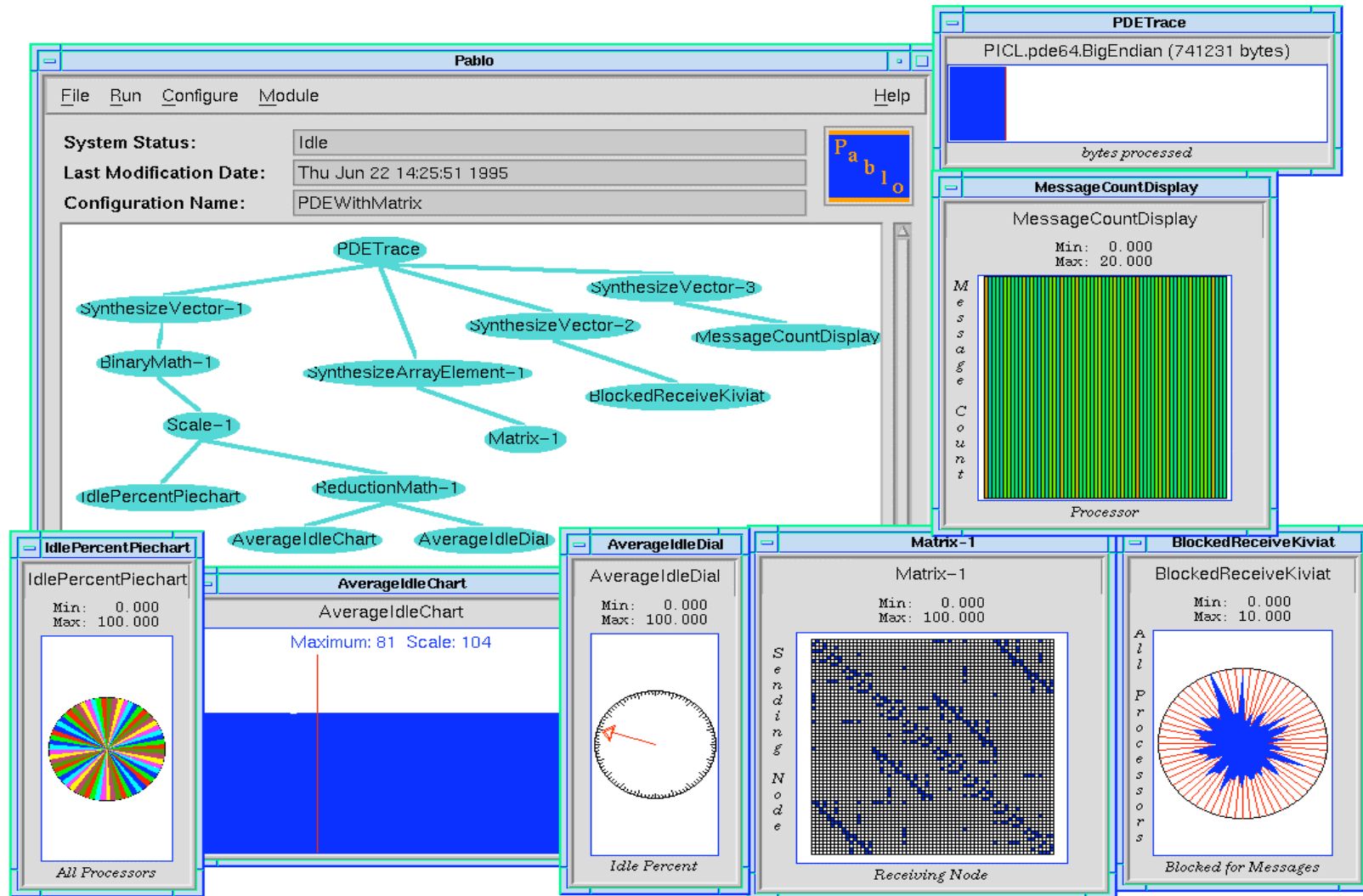
Architecture : **graph of components** connected by the user to build some visualizations.

Large number of **predefined components**:

Transformations : arithmetic opérations, statistics, etc.

Visualizations: pie chart, Kiviat, clouds of points, etc.

Pablo visualization



Assessment of Pablo

Extensibility, flexibility:

Relatively simple to add a new component.

Indépendence of programming model.

Scalable visualizations.

Non interactive, no relation between source code and visualizations.

No behavioral visualization of program executions.

Difficult to use? (meta tool).

SvPablo

One single component: monitoring, data analysis, etc.

Data monitoring: counting, timing.

Independence of the implementation language of the monitored application.

Correlation between input and application source program.

No trace \Rightarrow no behavioral visualization.

Scalability: counting + statistical visualizations.

SvPablo Visualization

The screenshot displays the SvPablo GUI interface. At the top, the window title is "svPablo". The menu bar includes "Project", "Instrument", "View", and "Help".

Project Description: Red Black SOR in C using MPI

Source Files: prbsor.c, prefax.c, p_io.c

Performance Contexts: Origin 2000 - 16 R10K processors - 800x800, Power Challenge - 8 R10K Processors - 125x125, NoW - 8 Sun UltraSparcs - 800x800, NoW - 4 Sun UltraSparcs - 125x125, Example: no instrumentation

Routines in Source File: main, MPI_Comm_size, MPI_Comm_rank, MPI_Get_processor_name, fprintf

Routines in Performance Data: MPI_Sendrecv, MPI_Send, MPI_Reduce

Source File: /home/reed/derose/D...
A grid visualization shows performance data for various lines of code, with colors corresponding to the routines in the performance data.

Specific Metric Dialog Box: R10K Statistics by Line
Data Cache Misses:
1300932.000000000 -- relax
370.000000000 -- updateOmega
246.000000000 -- LOOP
Buttons: Dismiss, Help

Code Snippet:

```
/*  
 * The fo  
 * to test  
 */  
for (i=1;  
    MPI_Send  
    u, n, MPI_DOUBLE, top, myid * blocksize,  
    MPI_COMM_WORLD, &status );  
MPI_Sendrecv( &u[n], n, MPI_DOUBLE, top, myid * blocksize  
              &u[(myend+1)*n], n, MPI_DOUBLE, bottom,  
              (myid+1) * blocksize + 1, MPI_COMM_WORLD, &  
              if (myid == 0) oldNorm = norm;  
              MPI_Reduce(&mynorm, &norm, 1, MPI_DOUBLE, MPI_SUM, 0, MPI  
              }  
saveOutput( u, n, i, norm, oldNorm );  
MPI_Finalize();  
}
```

Bottom navigation: Instrument/Clear Line, View Line Data

Vampirtrace and Vampir

Toolkit for performance visualization of **MPI** programs.

Scalability:

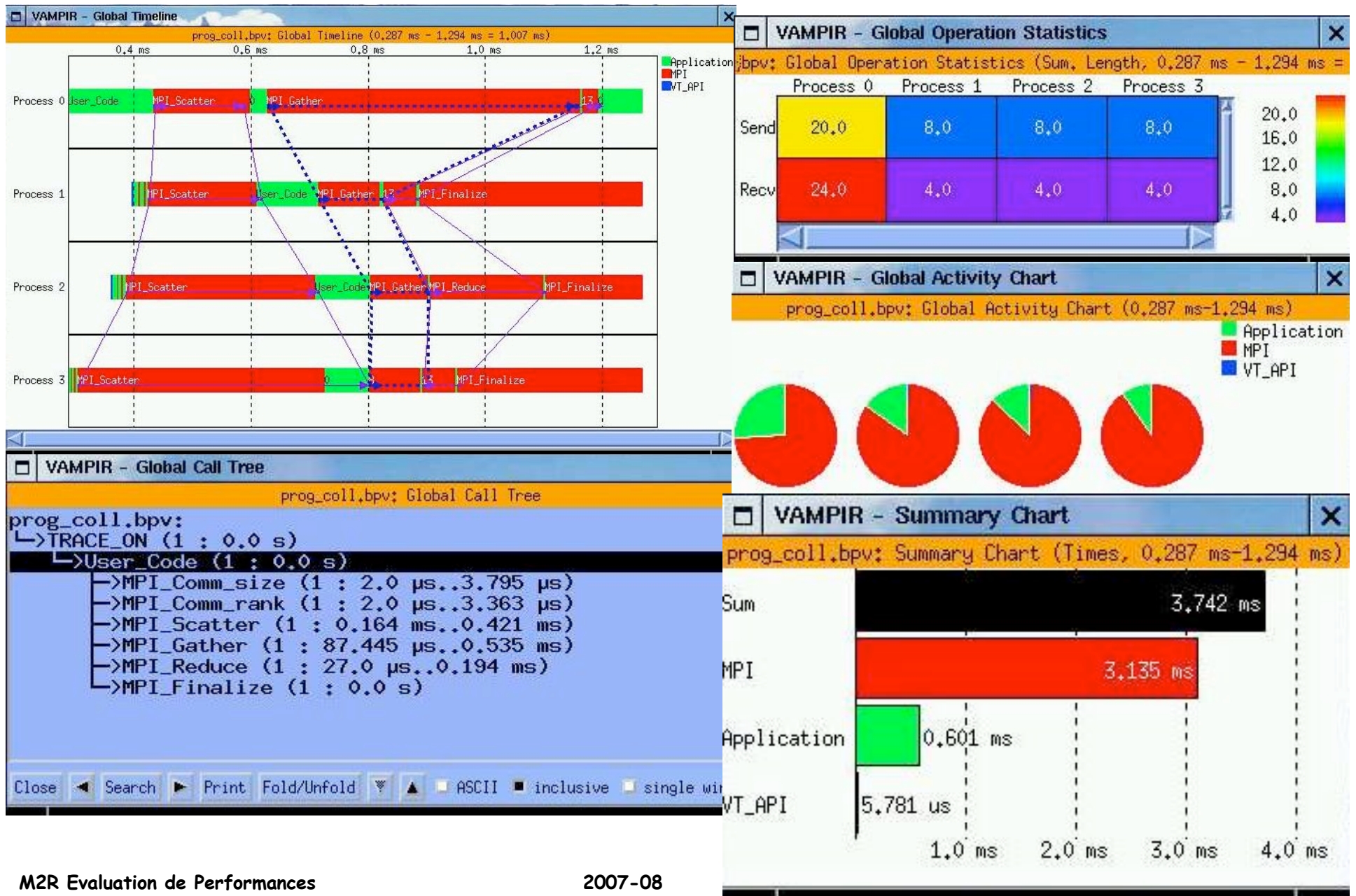
Numerous **zooming** possibilities.

Filtering of visualized objects.

Filtering of traces during monitoring (when using Vampirtrace).

Possibility to load a subset of a trace file: the set of events of a time period.

Vampir Visualizations



Vampir

Interactivity:

- Inspection of visual objects.

- Relation with source code (click back).

- Zooming along time axis.

- Not possible to move back and forth in time.

Extensibility? (commercial tool).

Standard?

- Translators to the Vampir format (ex. Tau \rightarrow Vampir).

- Widespread commercial tool.

Pajé*: interactive and extensible and scalable visualization tool



*doctor in Tupi language

Originally: visualizations of **communicating thread** program executions.

Main issue: **large number** of visual objects:

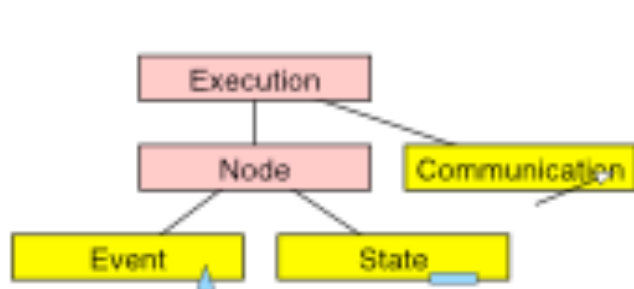
- Large number of nodes.

- One each node, dynamically created and terminated threads.

- Synchronizations and communications.**

Later: variety of programming models.

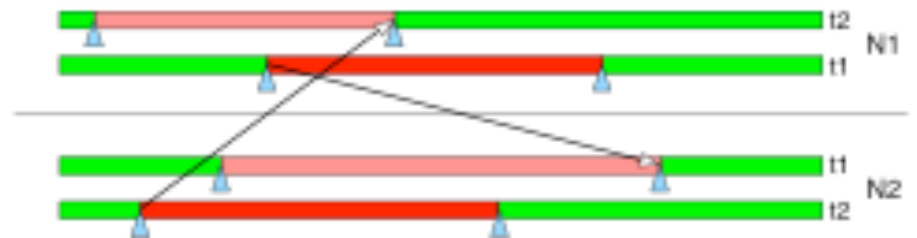
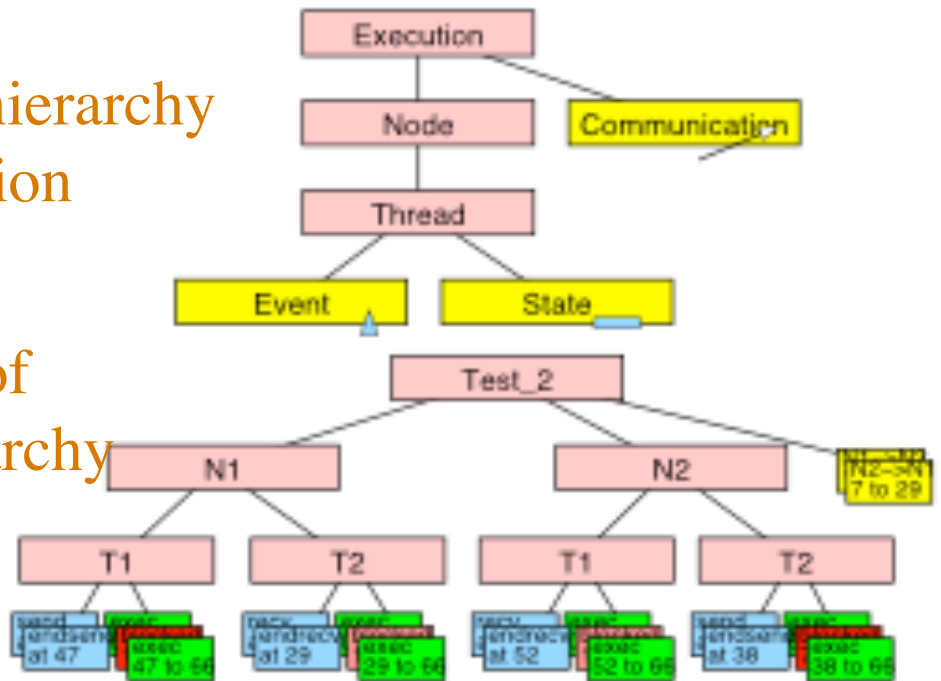
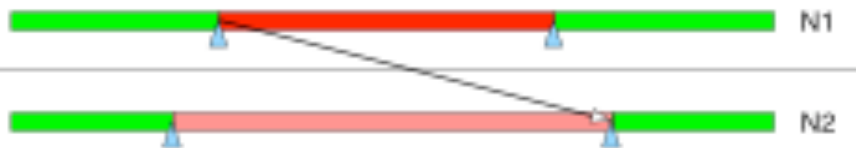
Example of type hierarchy definition and instance



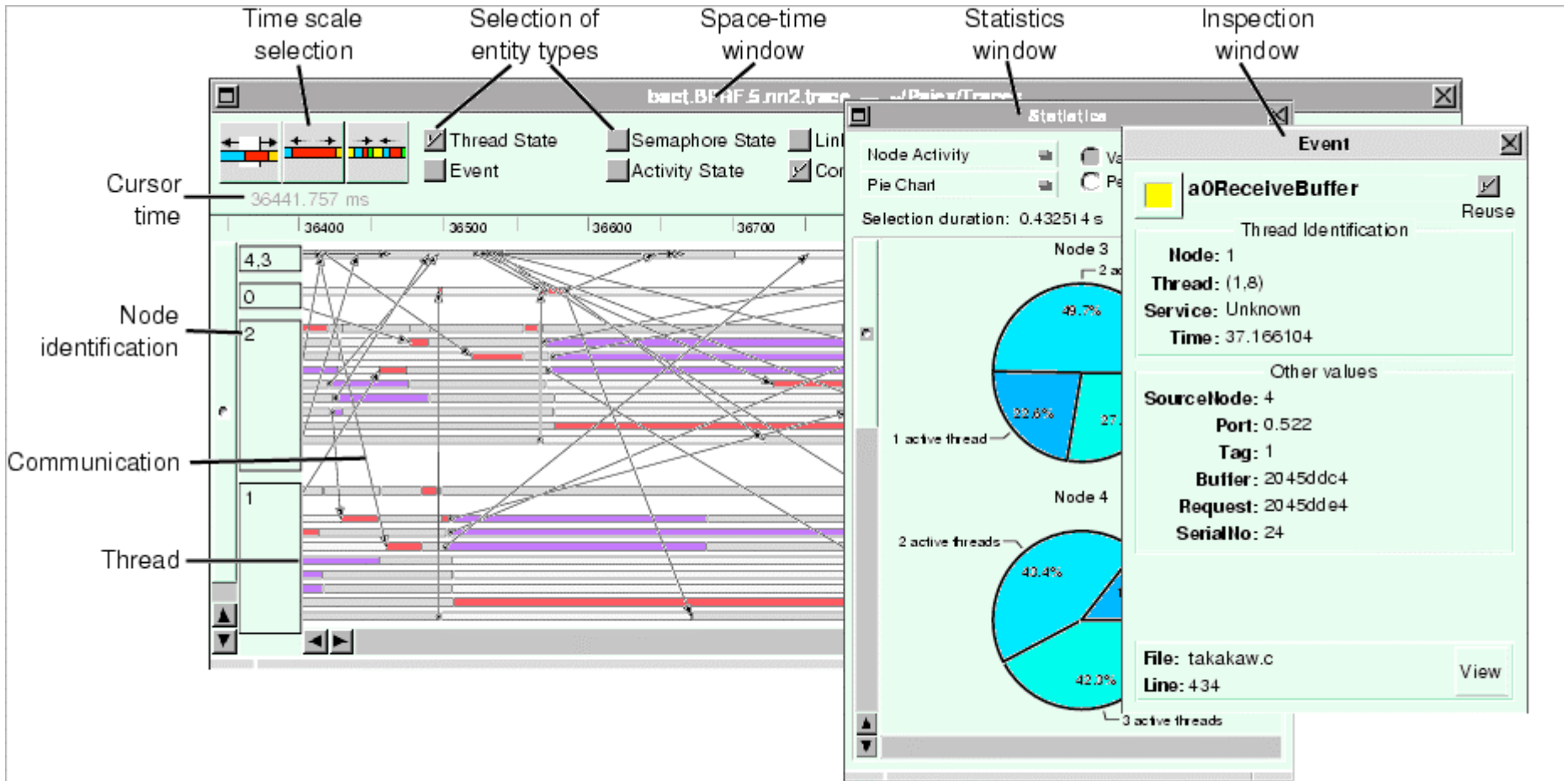
Type hierarchy definition



Instance of type hierarchy

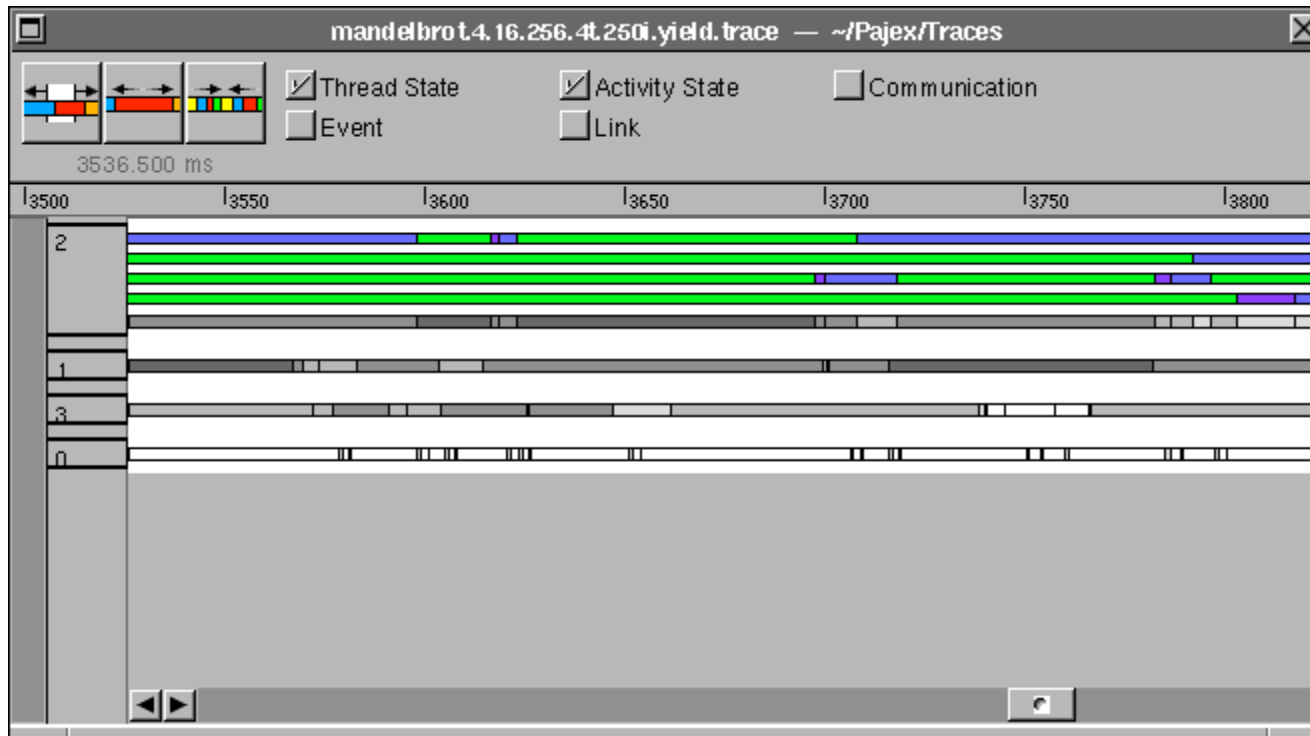


Pajé visualization example



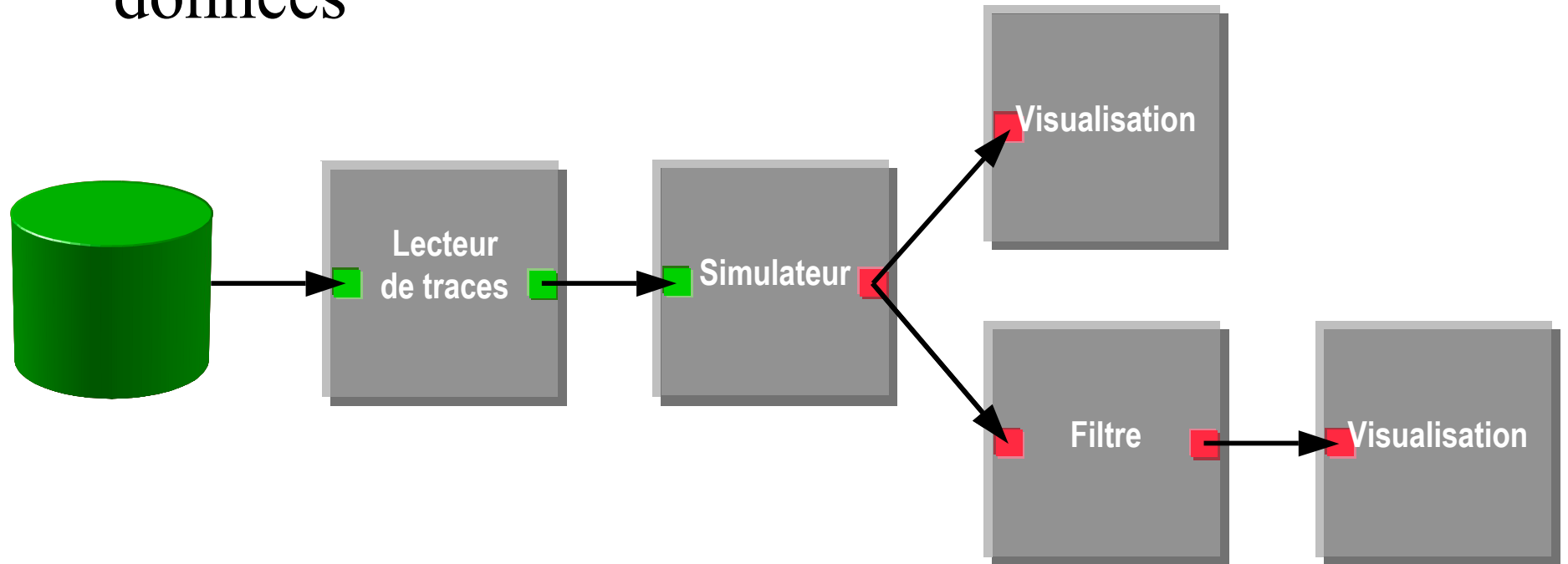
« Scalabilité »

■ Réduction des informations d'un nœud



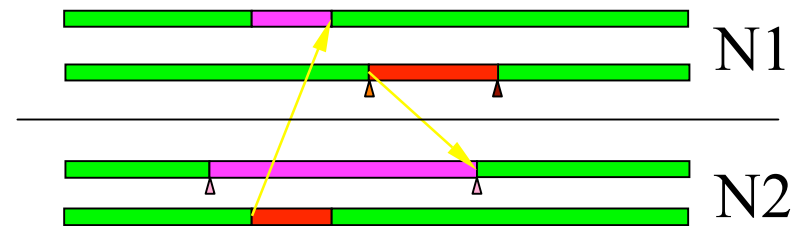
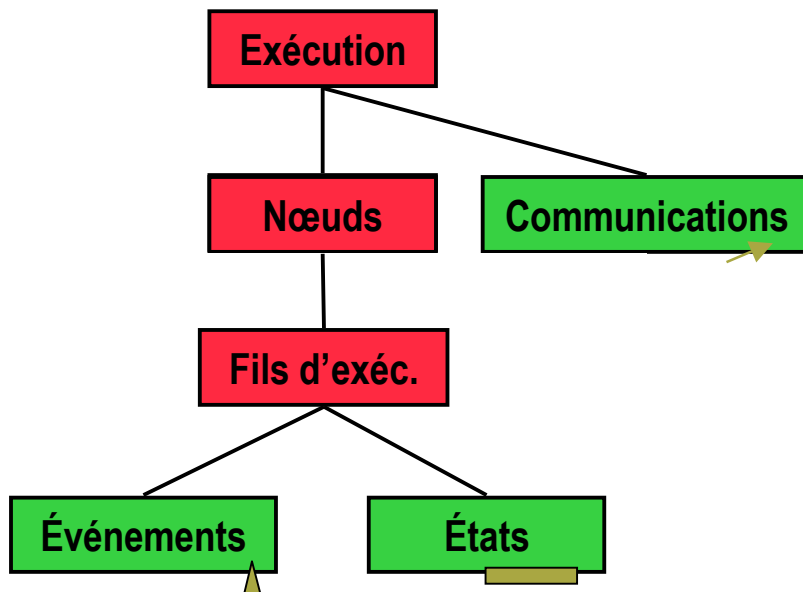
Extensibilité

- Architecture de l'outil : graphe flot de données



Extensibilité

- Pajé :
 - description de la structure des objets à visualiser
 - simulateur générique
 - description de la hiérarchie des types d'objets visualisés dans la trace.



Extensibilité

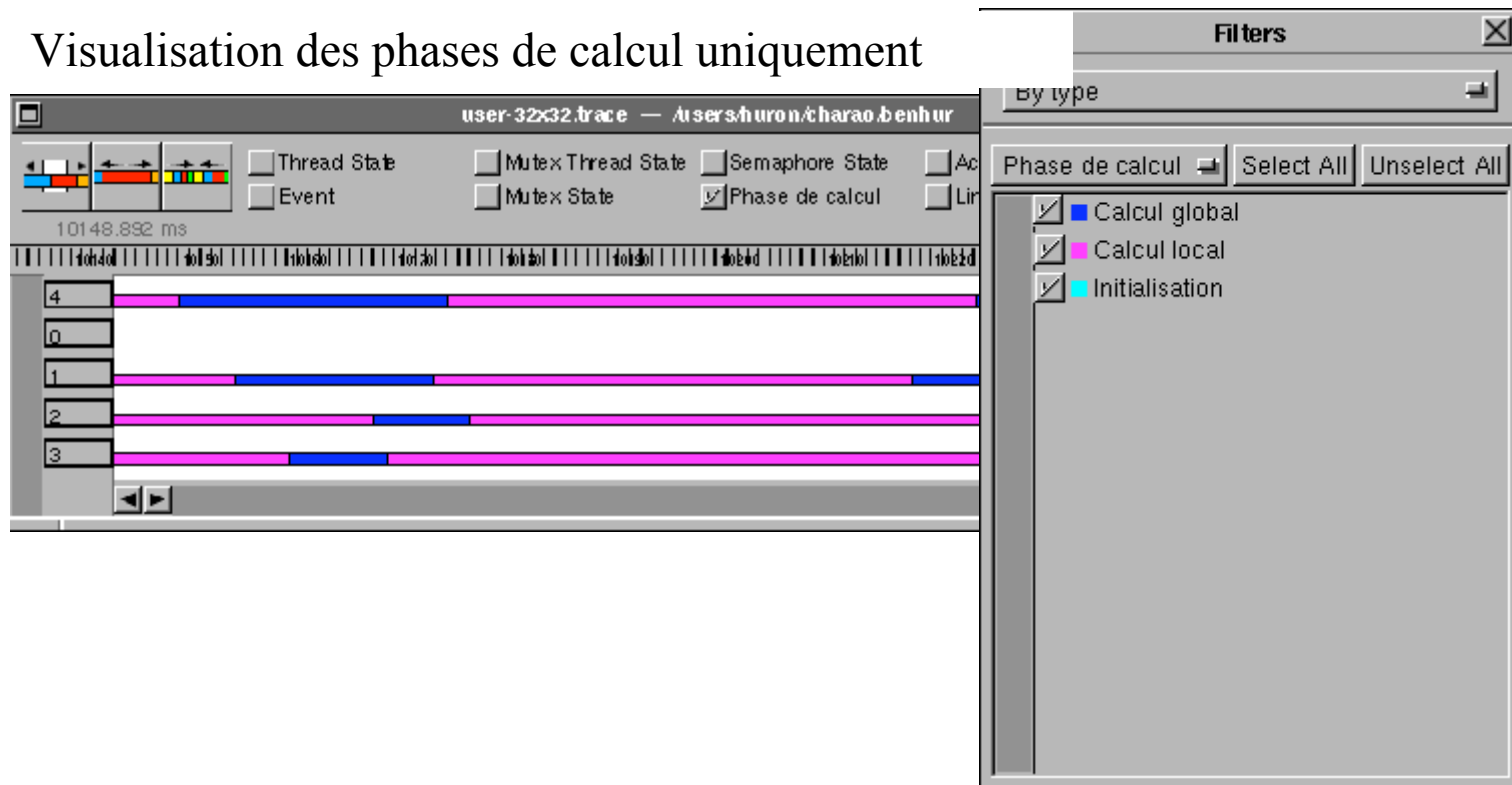
Création d'objets du nouveau type

```
unsigned phase_state, init_phase, local_phase, global_phase;
phase_state = pajeDefineUserStateType ( A0_NODE, "Phase de calcul" );
init_phase = pajeNewUserEntityValue ( phase_state, "Initialisation" );
local_phase = pajeNewUserEntityValue ( phase_state, "Calcul local" );
global_phase = pajeNewUserEntityValue ( phase_state, "Calcul global" );

pajeSetUserState ( phase_state, 0, init_phase, "" );
initialisation();
while (!converge) {
    iter++;
    pajeSetUserState ( phase_state, 0, local_phase, 0 );
    calcul_local();
    send (local_data);
    receive (remote_data);
    pajeSetUserState ( phase_state, 0, global_phase, 0 );
    calcul_global();
}
```

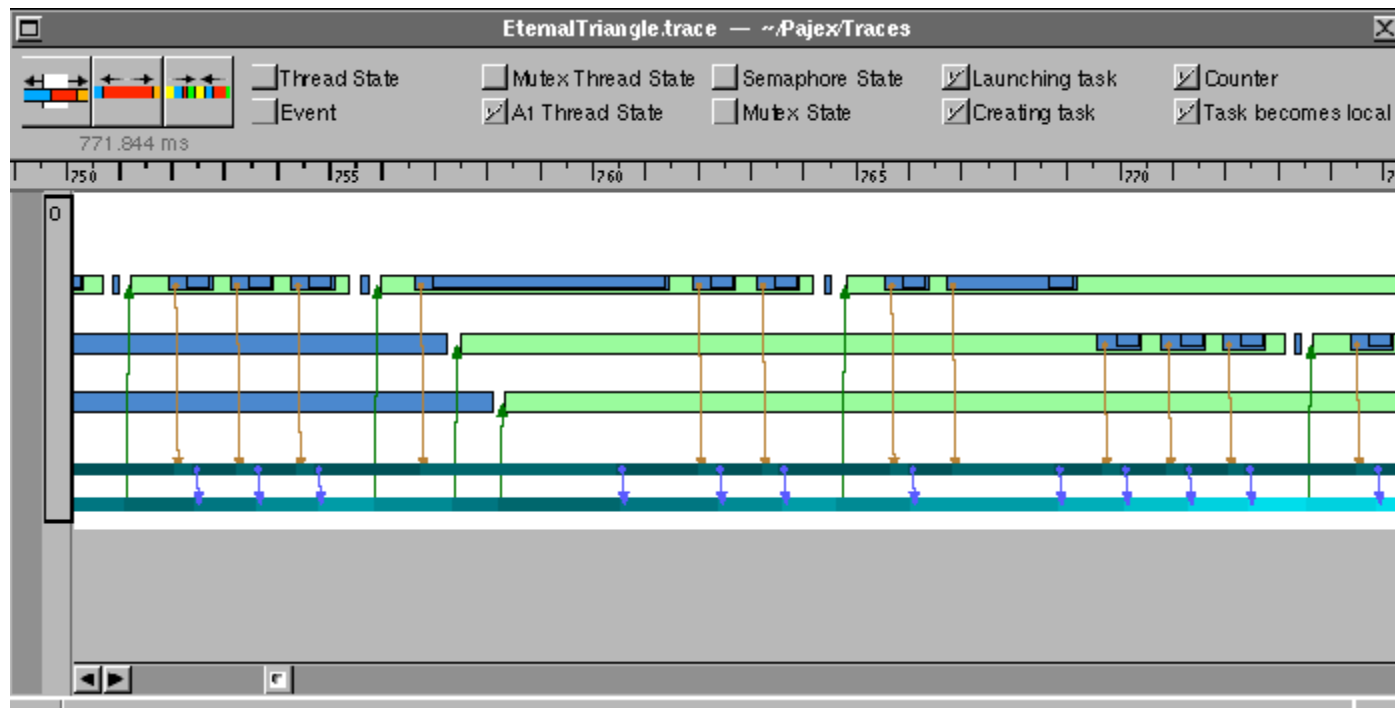
Extensibilité

Visualisation des phases de calcul uniquement



Extensibilité

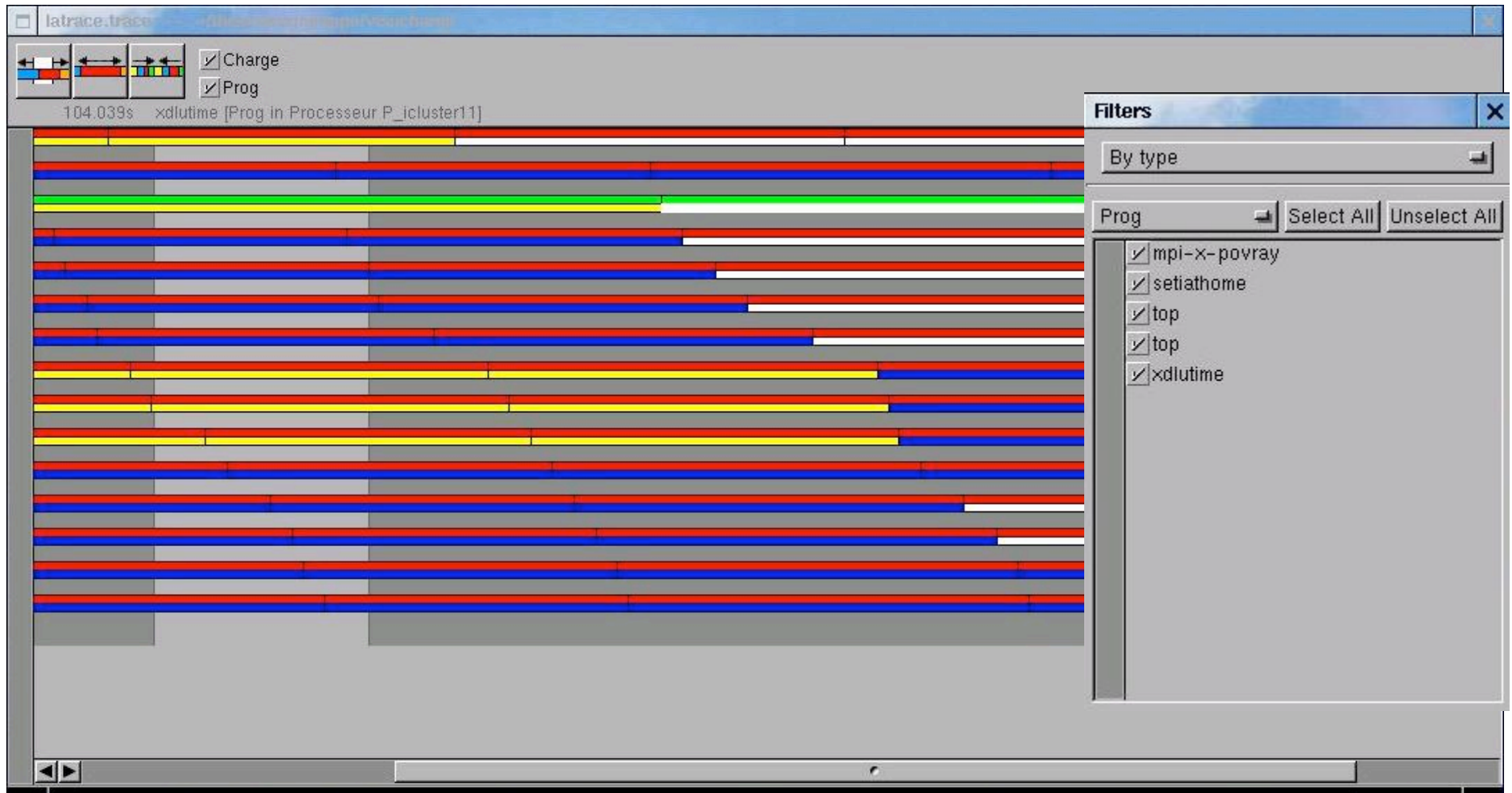
- Visualisation de l'ordonnanceur d'Athapascan-1
- Addition de quelques lignes dans Athapascan-1
- Aucun changement dans Pajé



Exploitation de la généricité de Pajé

- Utilisation de Pajé pour l'administration système de grappes (Cyril Guilloud).
- Mesure de performances de programmes Java distribués (François Ottogali):
 - Corrélation traces du niveau application et indices de performance du niveau système.
 - Utilisation de Pajé pour visualisation du comportement des applications.

Cluster administration





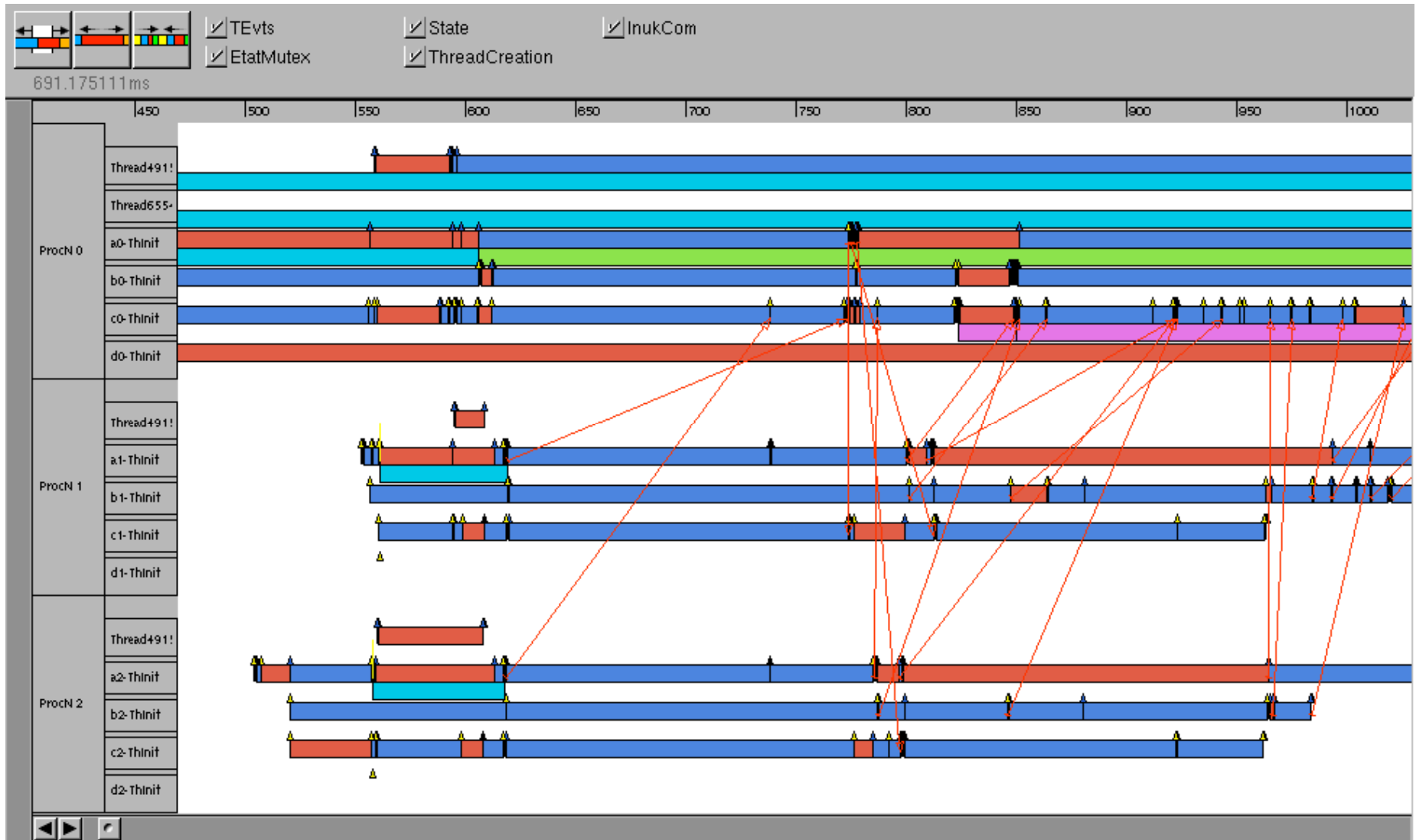
pbs-task

2.136081 d

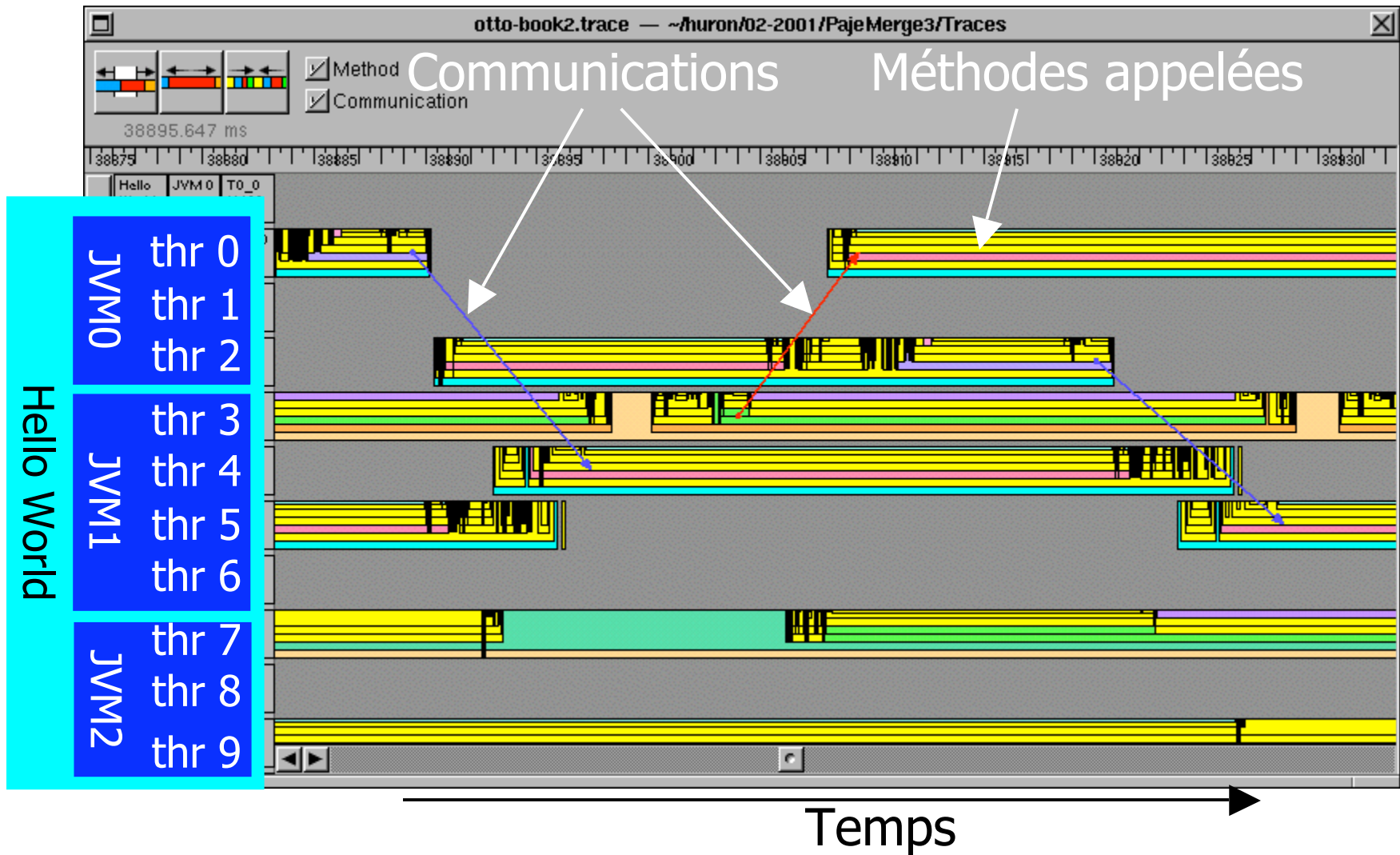
0 5 10 15 20



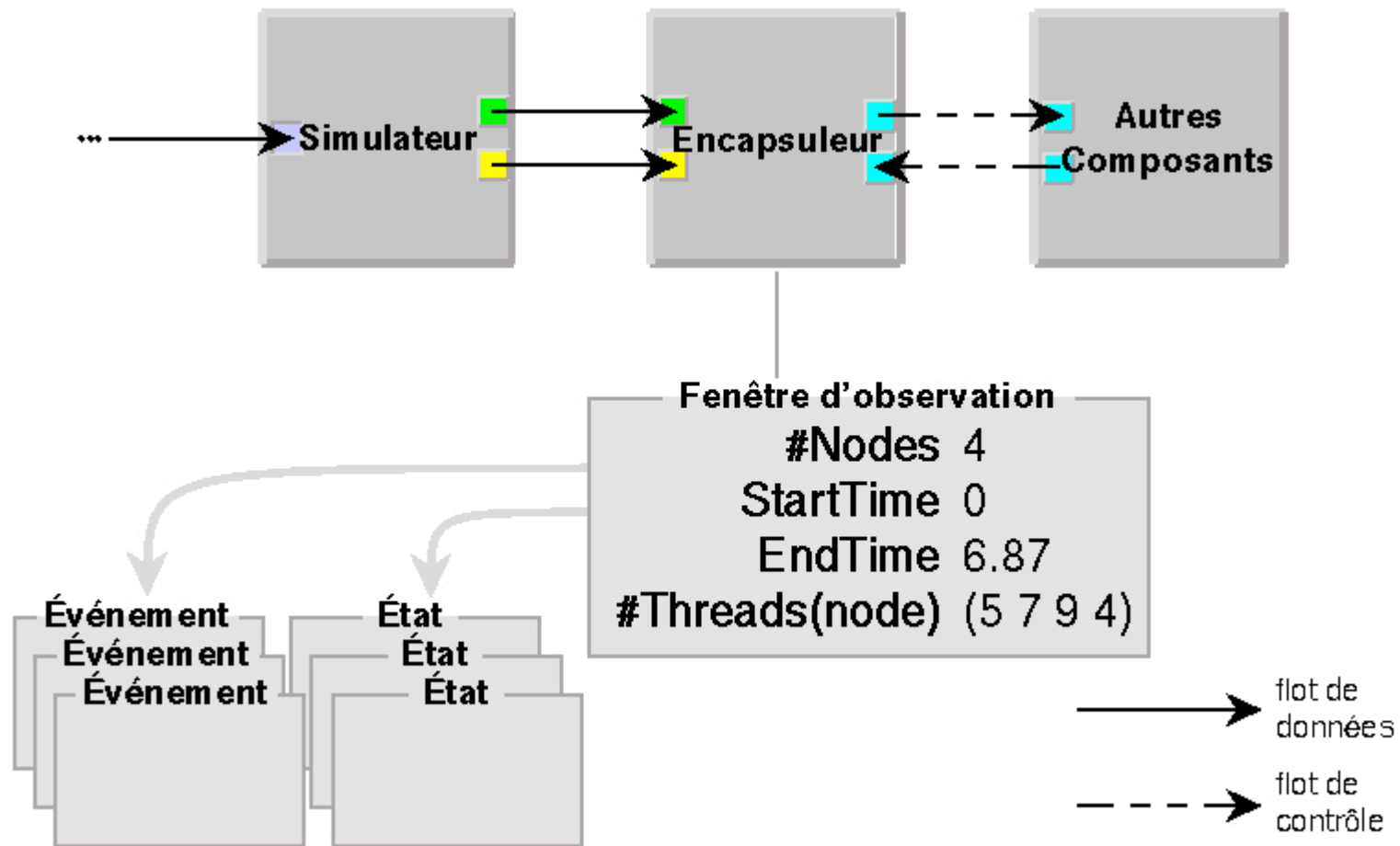
Molecular biology code



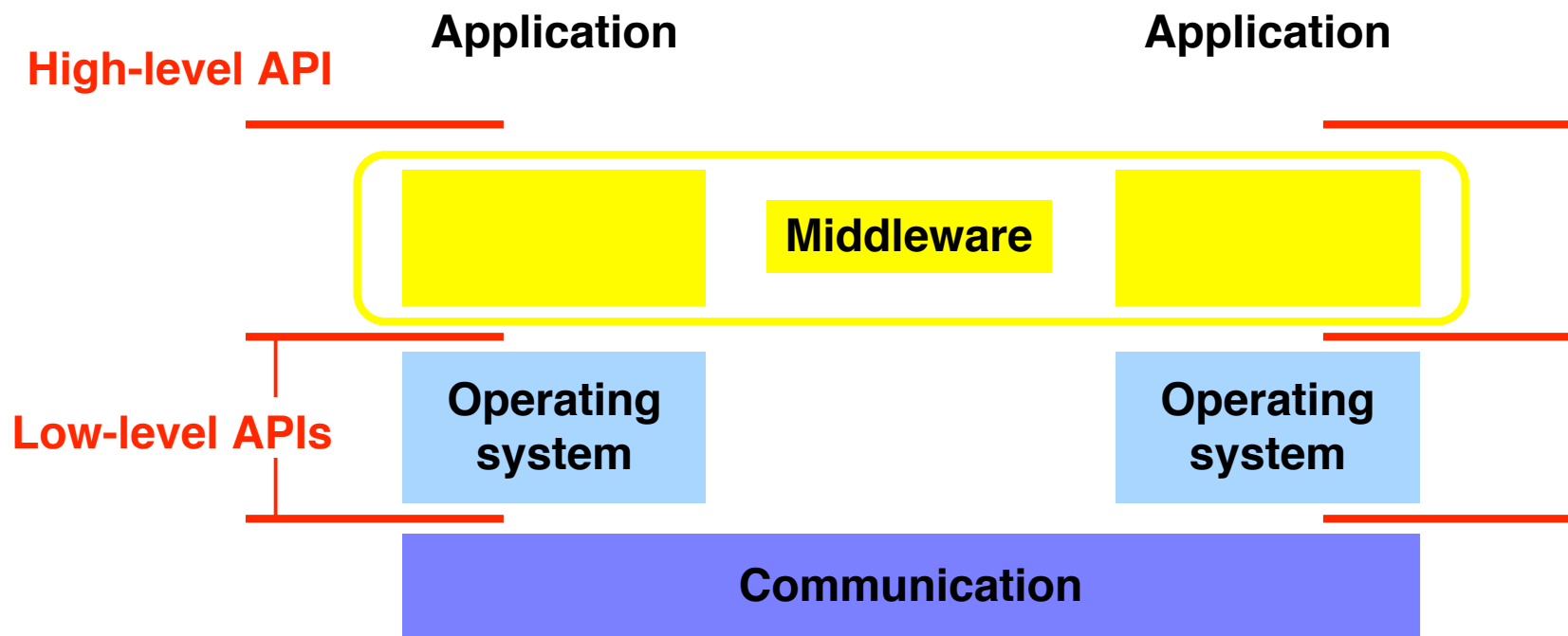
Programme Java distribué



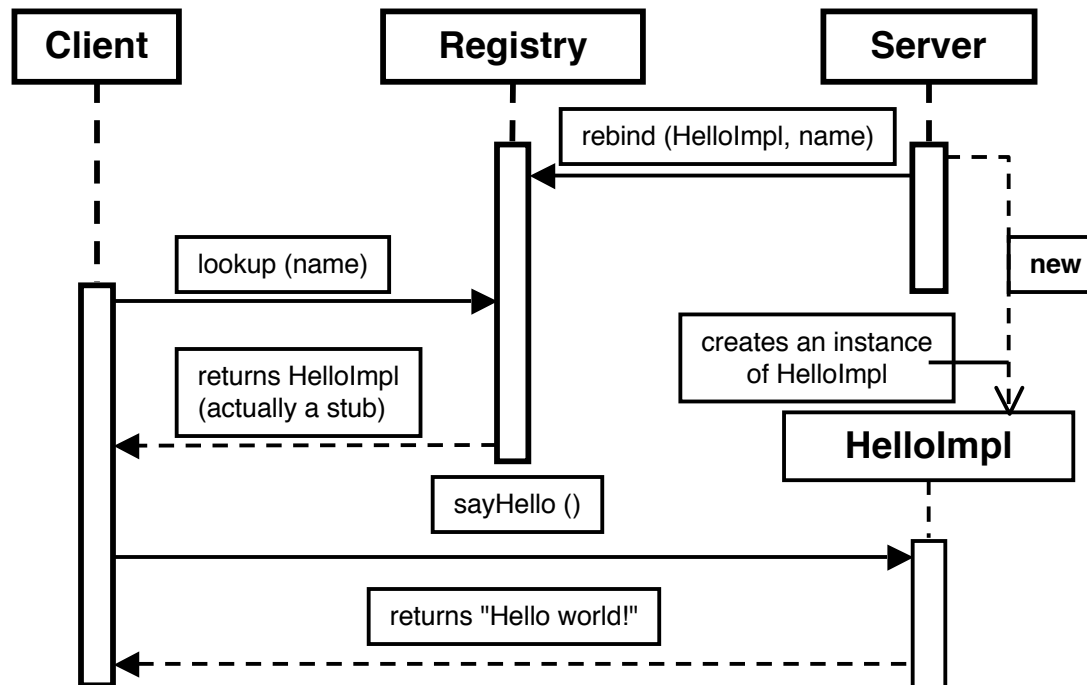
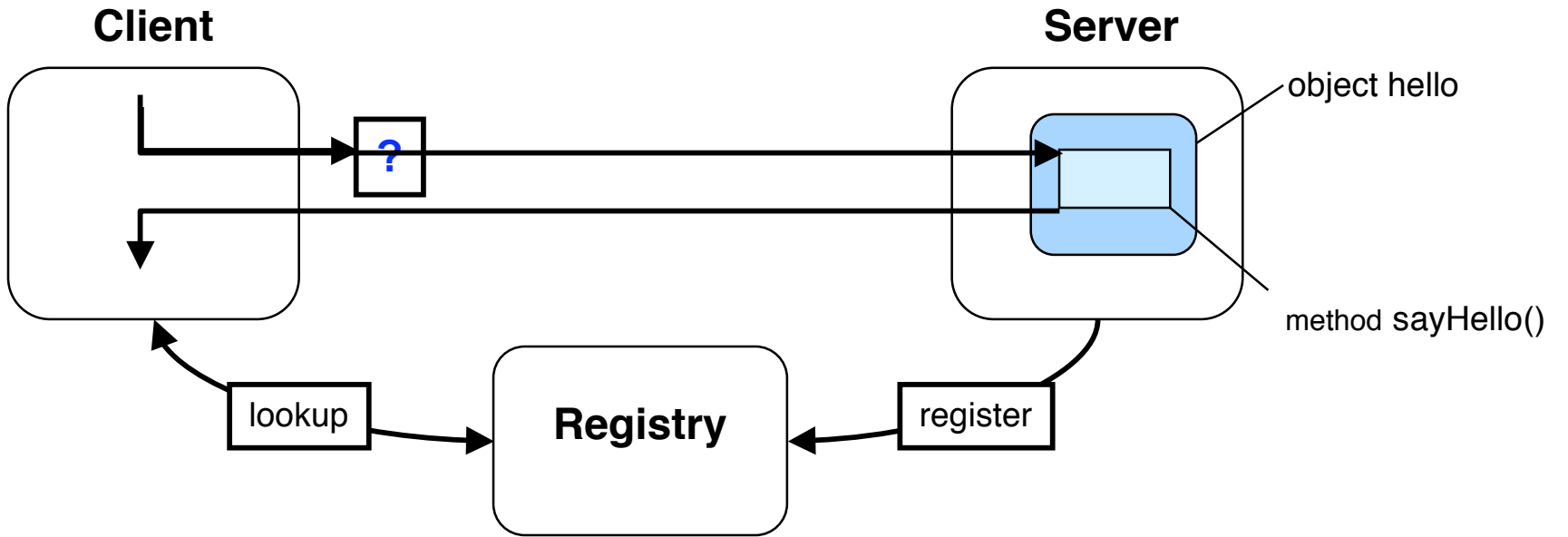
Extensibilité et Interactivité et «scalabilité»



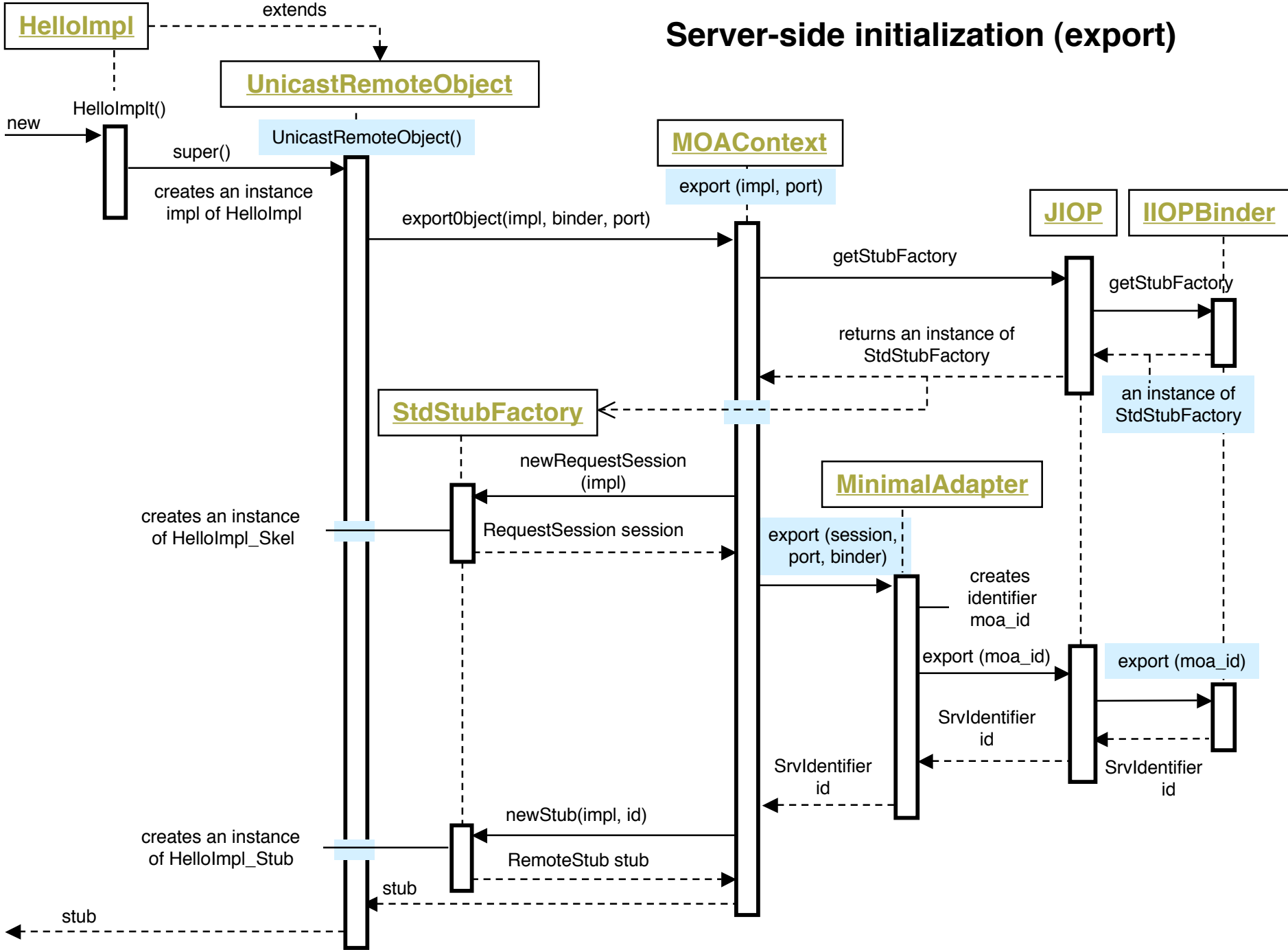
Architecture logicielle

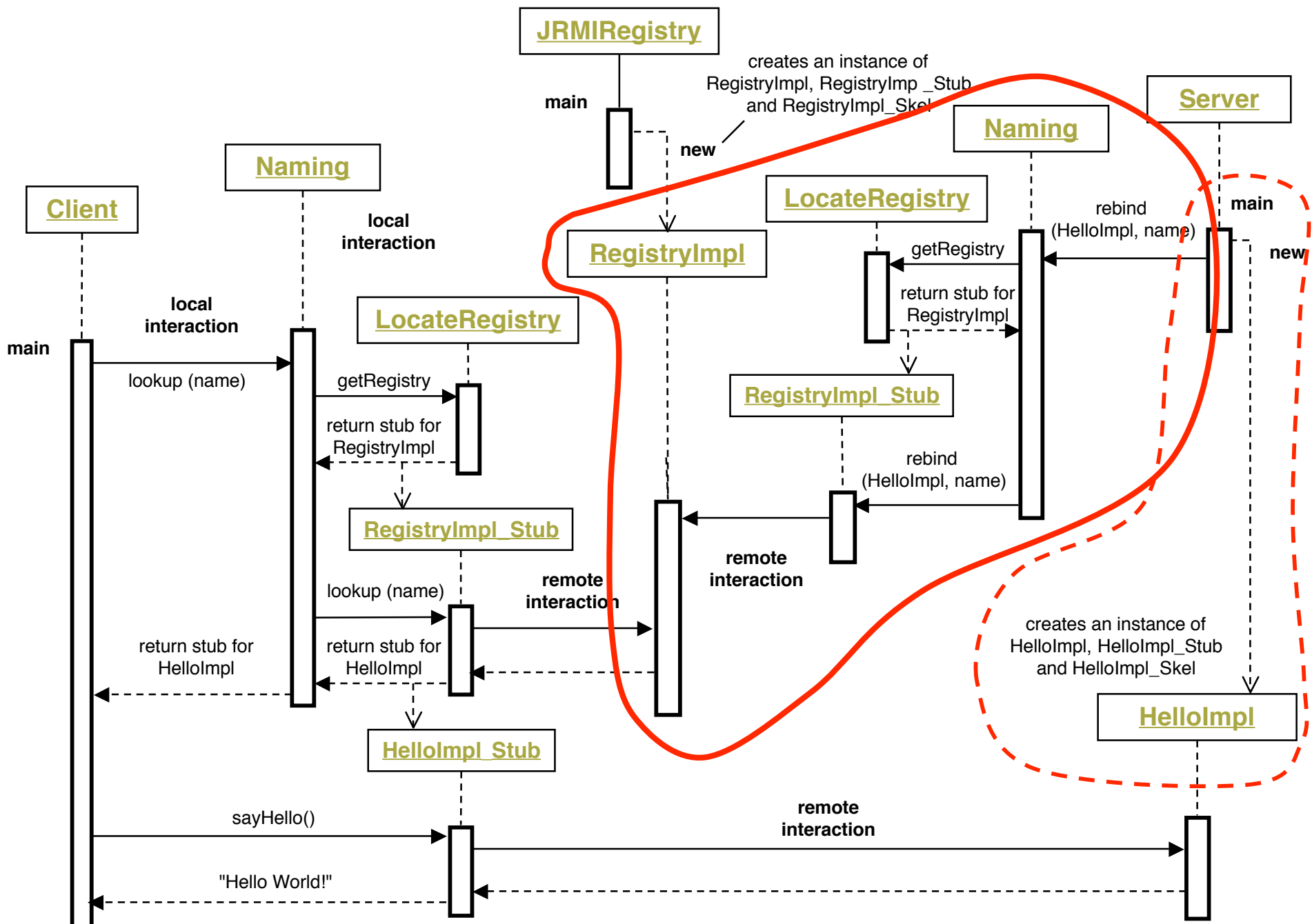


© 2002, S. Krakowiak

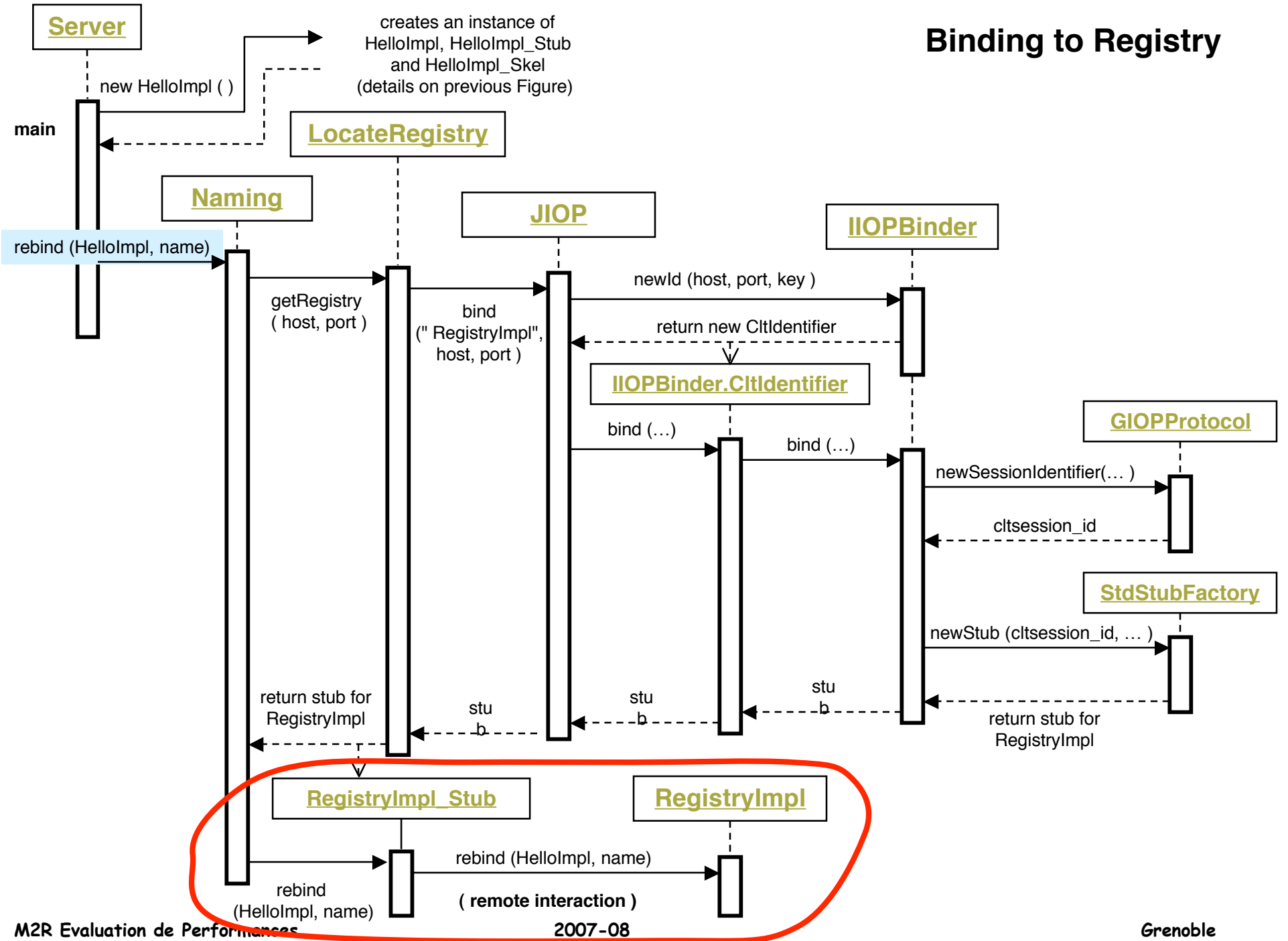


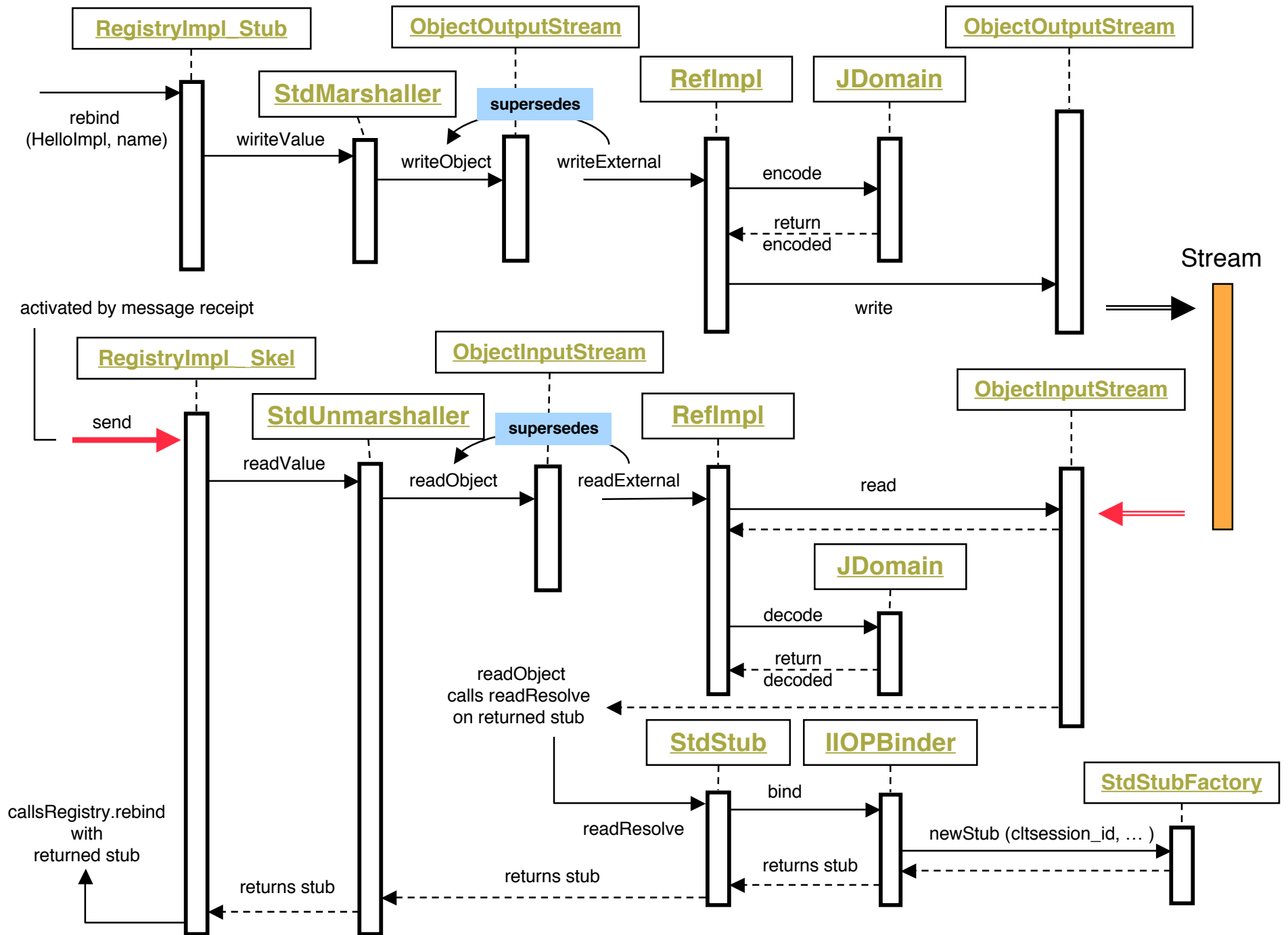
Server-side initialization (export)





Binding to Registry





Modèles et mesure

Définition états – transitions

→ niveau d'abstraction

→ abstraction de l'exécution

État intéressant

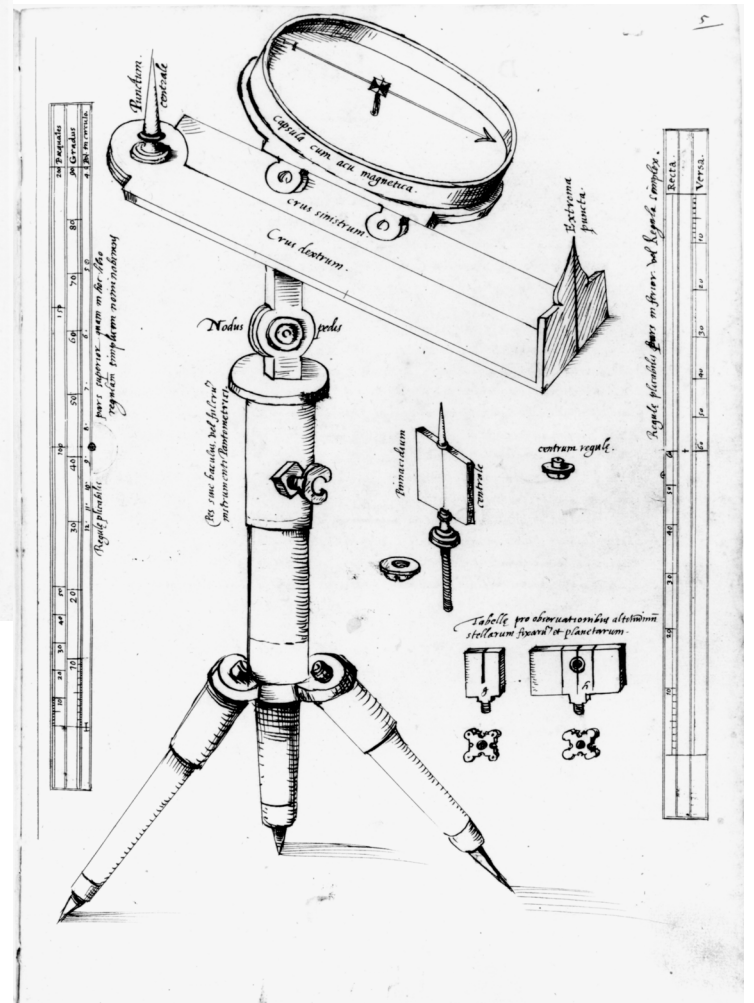
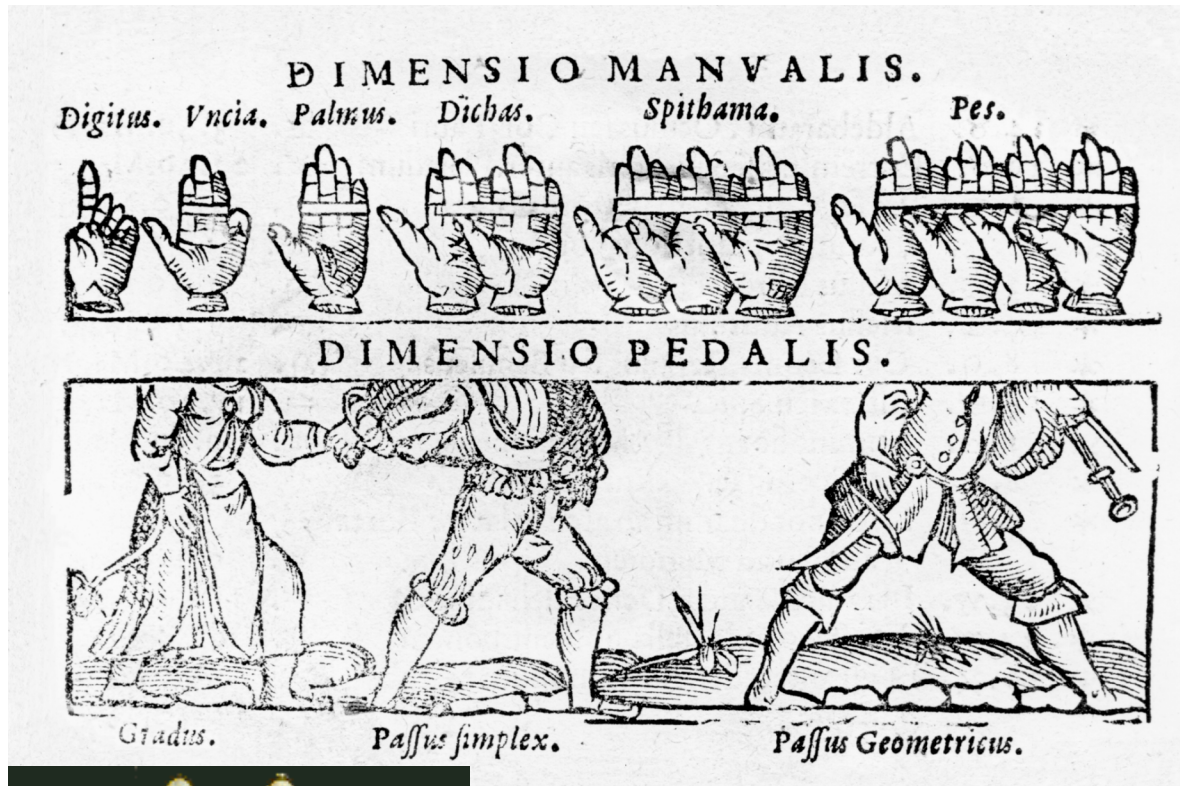
→ action de l'utilisateur (code, ordonnancement,...)

État observable

→ instrument de mesure associé

Conception et instrumentation ↔ Intégration d'instruments

Exécution séquentielle

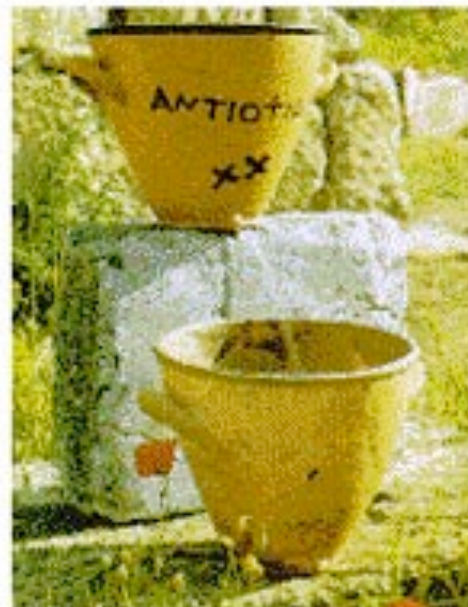


Techniques : time references

› Clock driven measurements



(a)



(b)

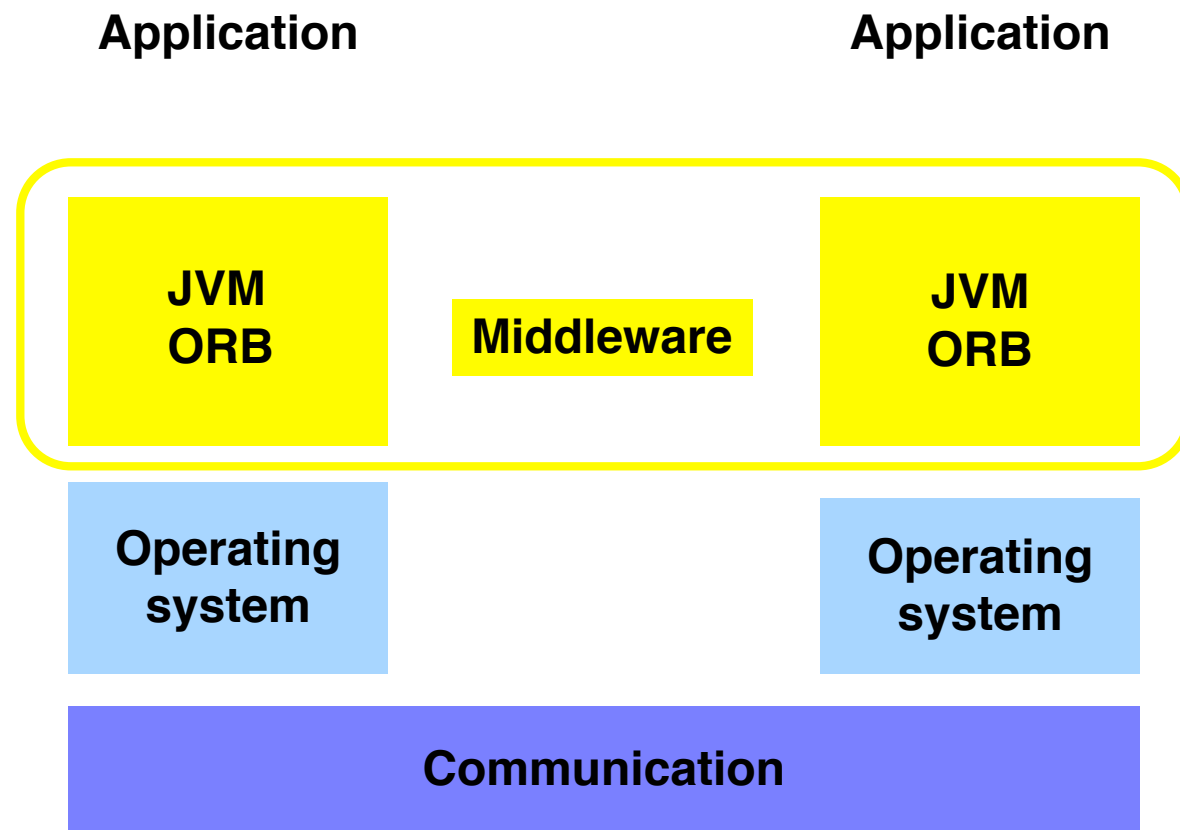


(c)



(d)

Architecture logicielle



événements **utilisateur**
- état de variables,...

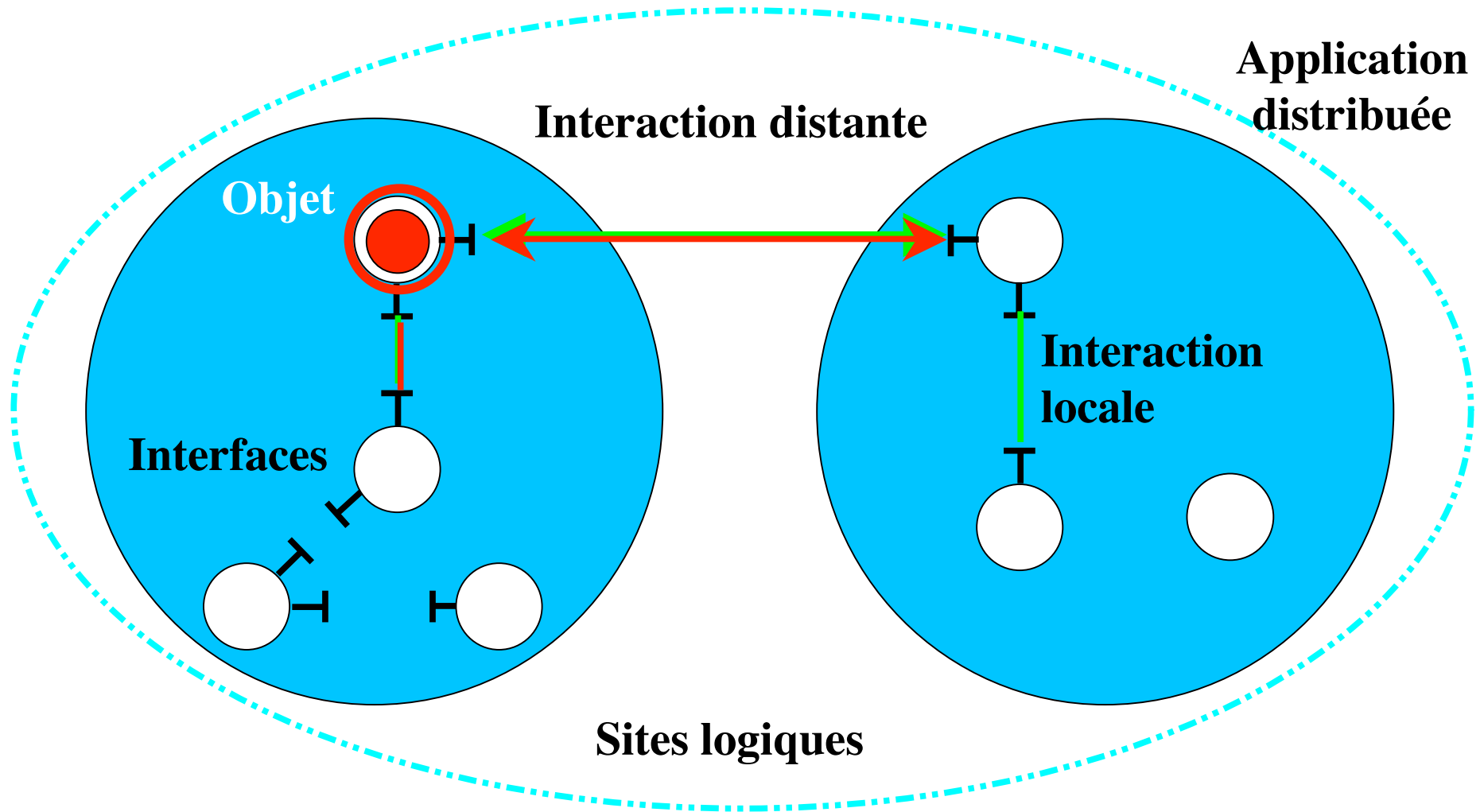
événements **utilisateur JVM**
-appels de méthode,
-chargement de classe,...

événements **gestion JVM**
-threads,
-moniteurs,...

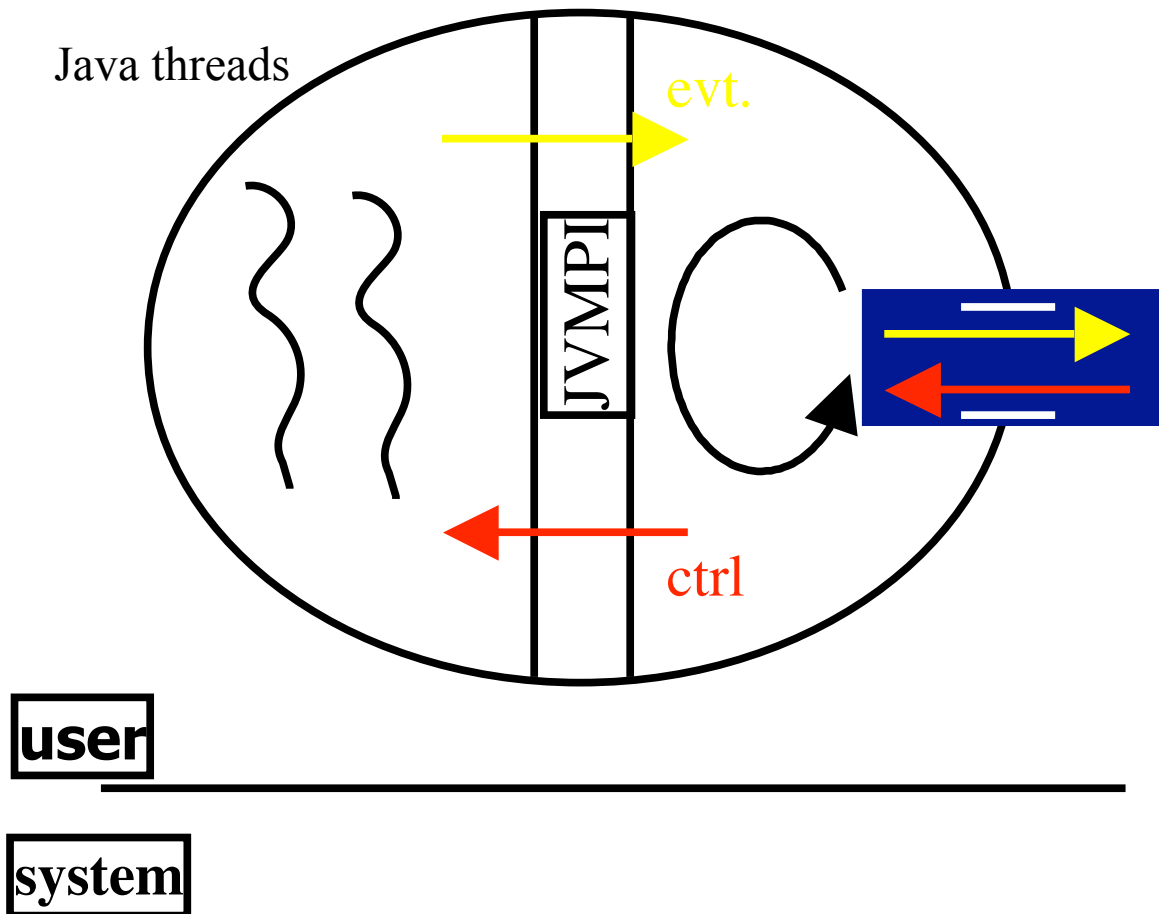
événements **système**
- état des threads système,
- sockets, mémoire,...

événements **réseau**
- état des liens,routeurs,...

Application distribuée à objets



JVM Profiling Interface



Timestamped events:

classes
objects
methods
monitors
GC
threads
JVM

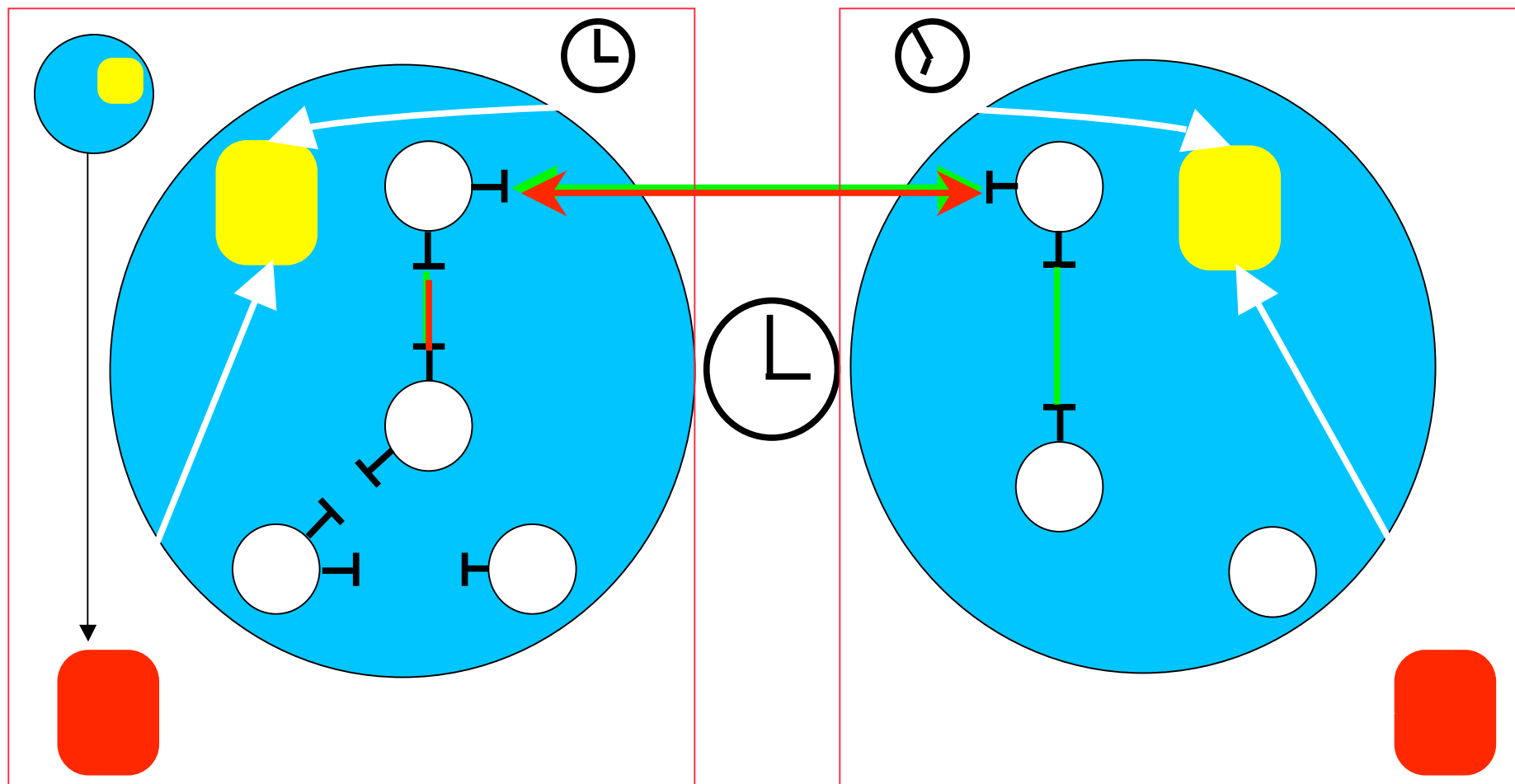
Control :

threads scheduling
GC
execution context

Trace locale filtrée et indentée

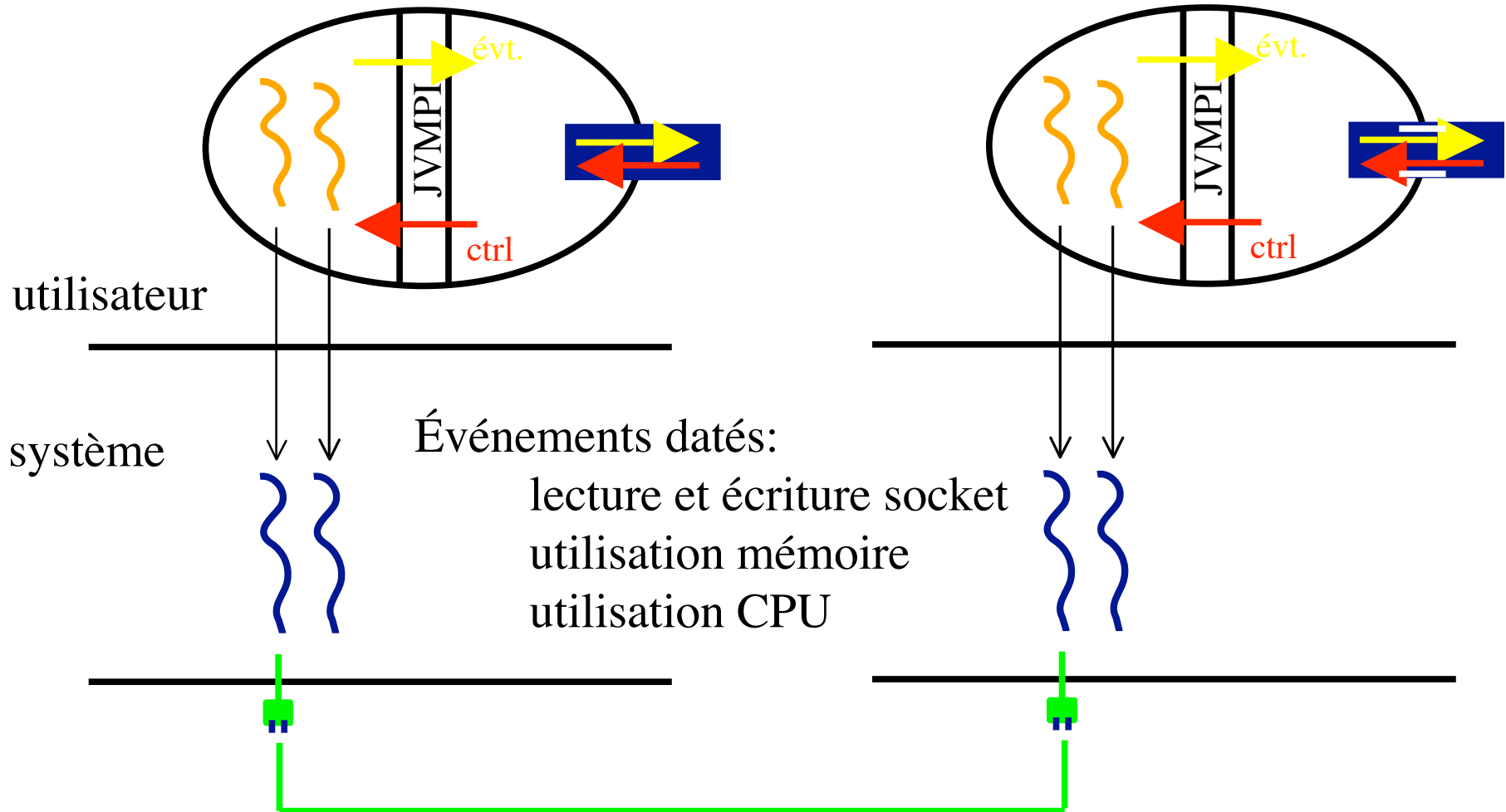
```
Client/<clinit> : D F
Client/main : D
| Client/<init> : D F
| org/objectweb/david/libs/contexts/orbs/ORBSingletonClass/<init> : D F
| org/objectweb/david/libs/contexts/orbs/ORBClass/<clinit> : D F
| org/objectweb/david/libs/contexts/orbs/iiop/IIOPORB/<clinit> : D F
| org/objectweb/david/libs/contexts/orbs/iiop/IIOPORB/<init> : D
| | org/objectweb/david/libs/contexts/orbs/ORBClass/<init> : D
| | | org/objectweb/david/libs/contexts/orbs/ORBSingletonClass/<init> : D F
| | | org/objectweb/david/libs/contexts/orbs/ORBSingletonClass/setDefault : D F
| | org/objectweb/david/libs/contexts/orbs/ORBClass/<init> : F
| org/objectweb/david/libs/contexts/orbs/iiop/IIOPORB/<init> : F
```

Cohérence d'observations distribuées

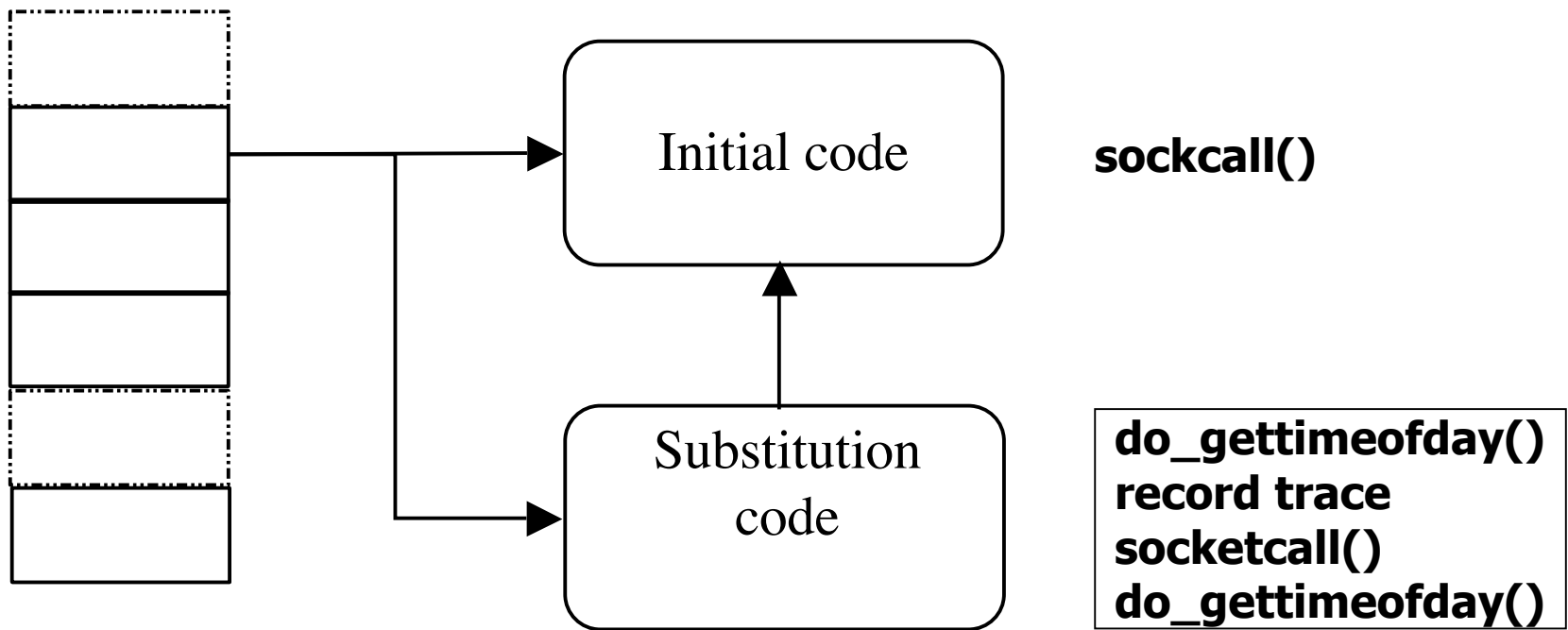


Architecture distribuée

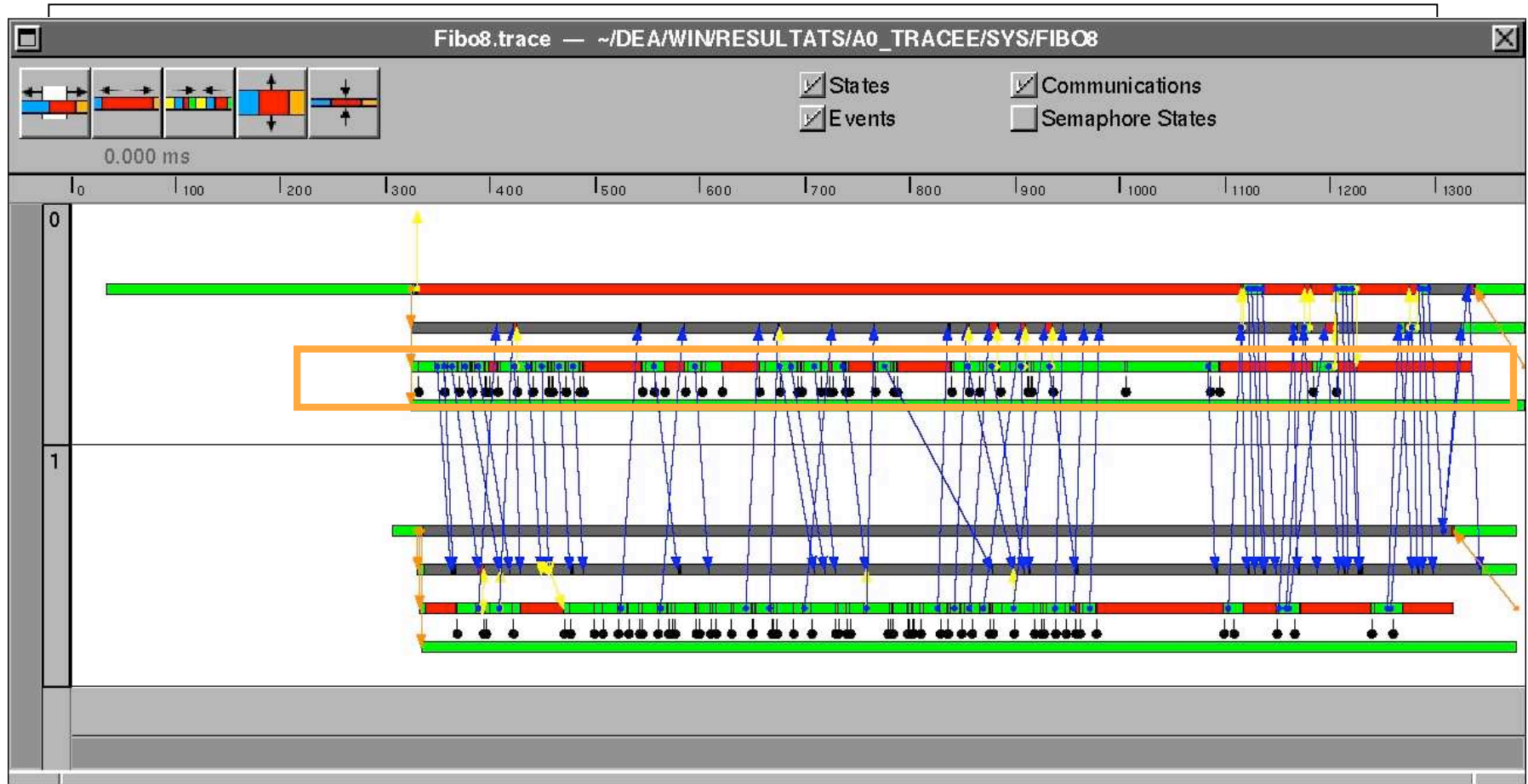
Observations systèmes



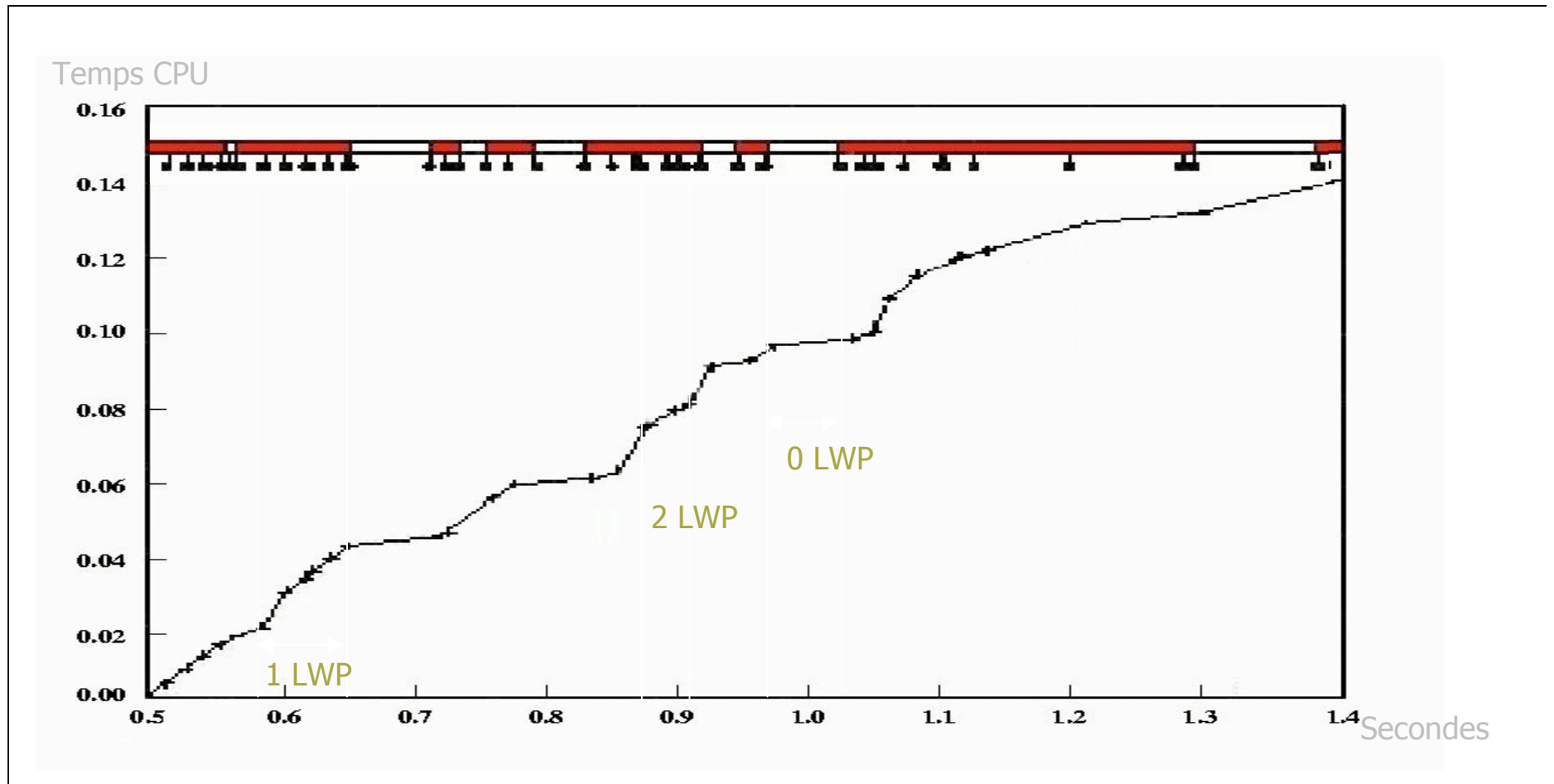
Derouting of system calls



Modèle Threads + communications



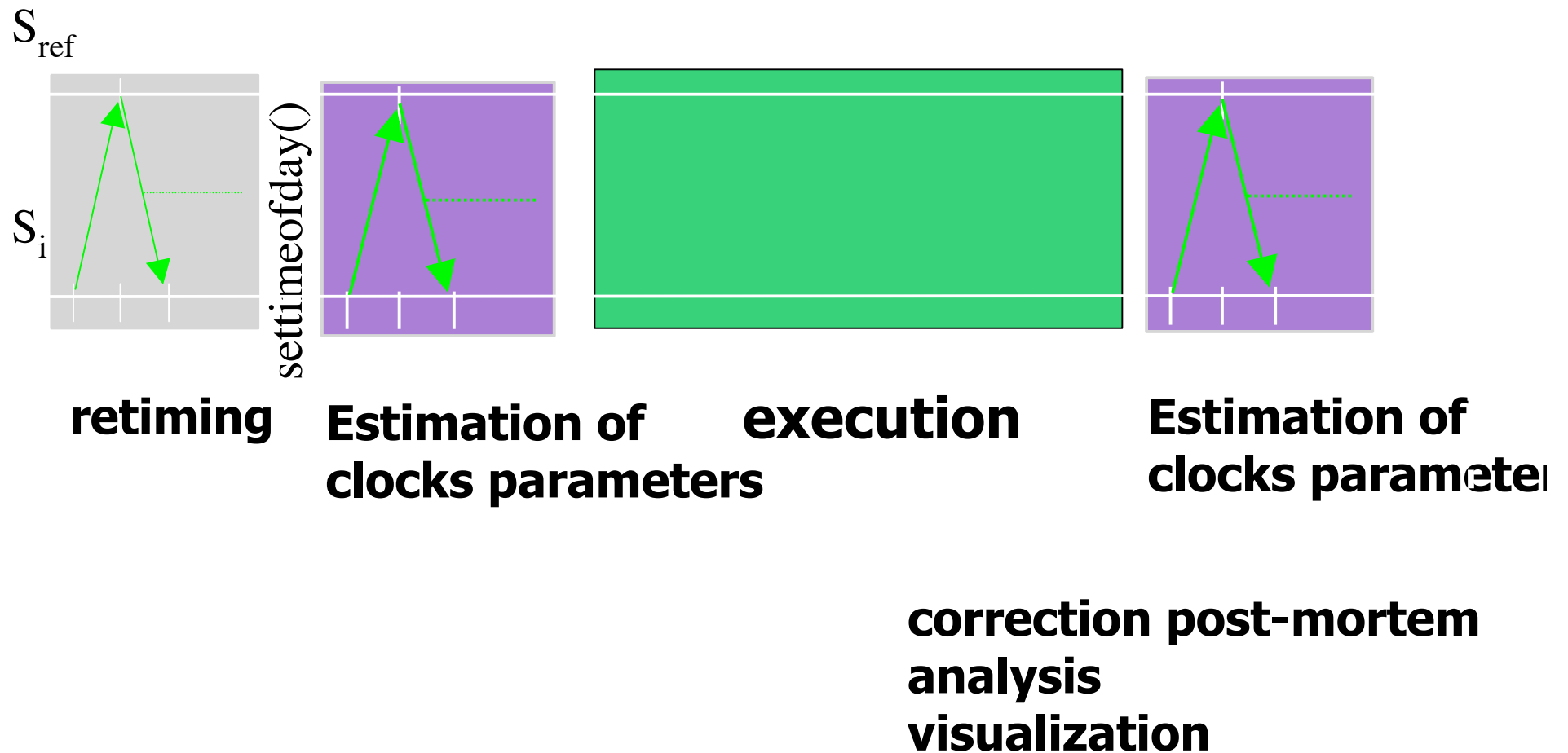
Traceurs systèmes / événements applicatifs



Time in distributed systems



Experimental protocol



Intégration et mise en cohérence

- Événements utilisateurs → Threads Java
(filtre JVMPI sur appel de méthode)
- Événements JVM locaux → JVMPI
- Interactions distantes
Thread Java → thread système → socket id
- Horloge globale

Observation = perturbation

Traceur et application partagent les mêmes ressources

Calcul : modification du temps d'exécution

- › estimation, correction

Mémoire : stockage local des traces

- › gestion de cache, ralentissement,

Réseau : extraction des traces locales

- › latences, débits,

Analyse des perturbations

	Delays (μ s)	Size (octet)
gettimeofday()	1.049	8
do_gettimeofday()	0.195	8
event JVMPI	3	60
Writing « socket »	2.13	36 + 32
Reading « socket »	1.55	36 + 32

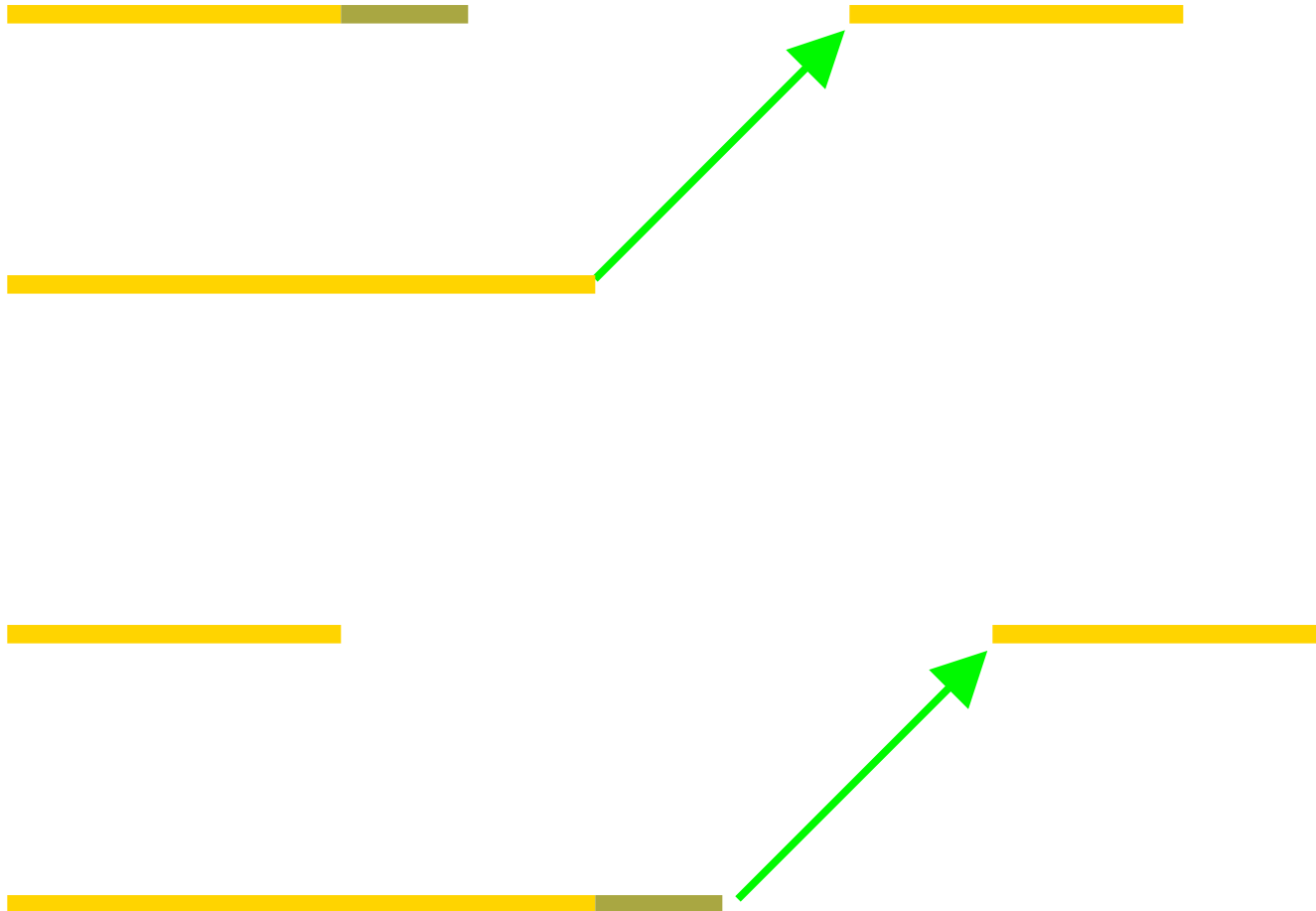
Perturbations application level : \approx **6 %**

Perturbations OS level : \approx **3 %**

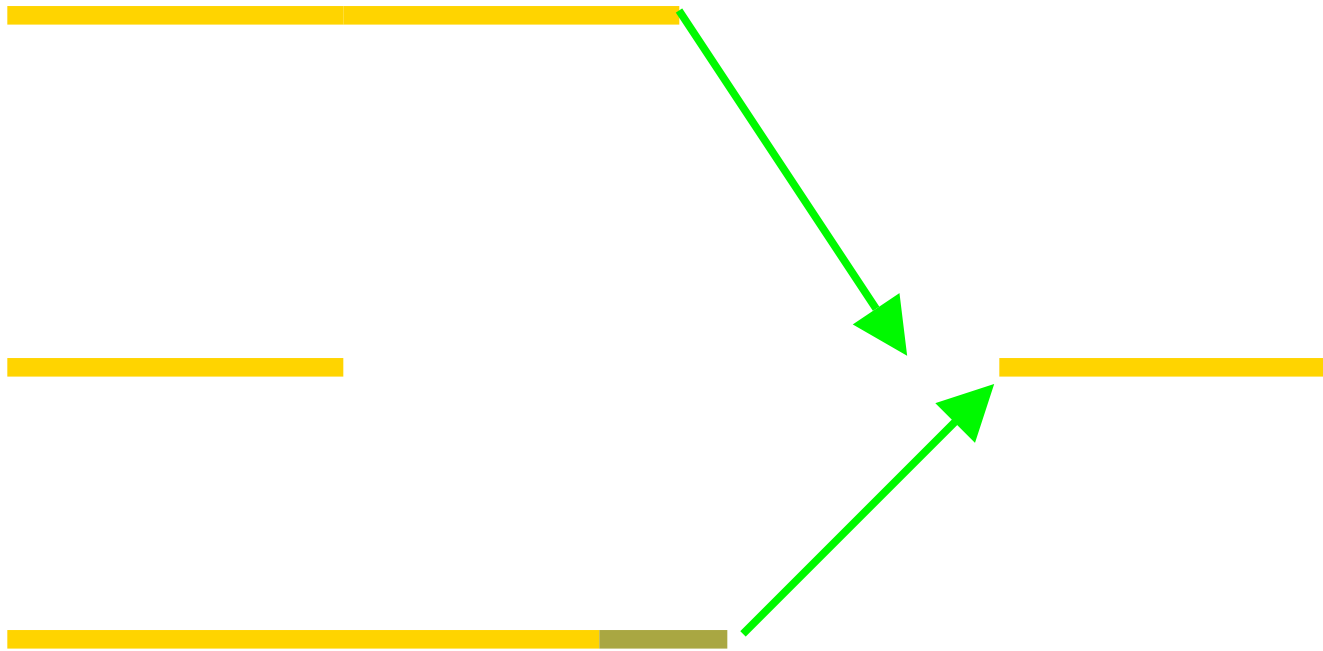
Correction post-mortem



Estimation of recording time
Perturbation accumulation
Possible transfert



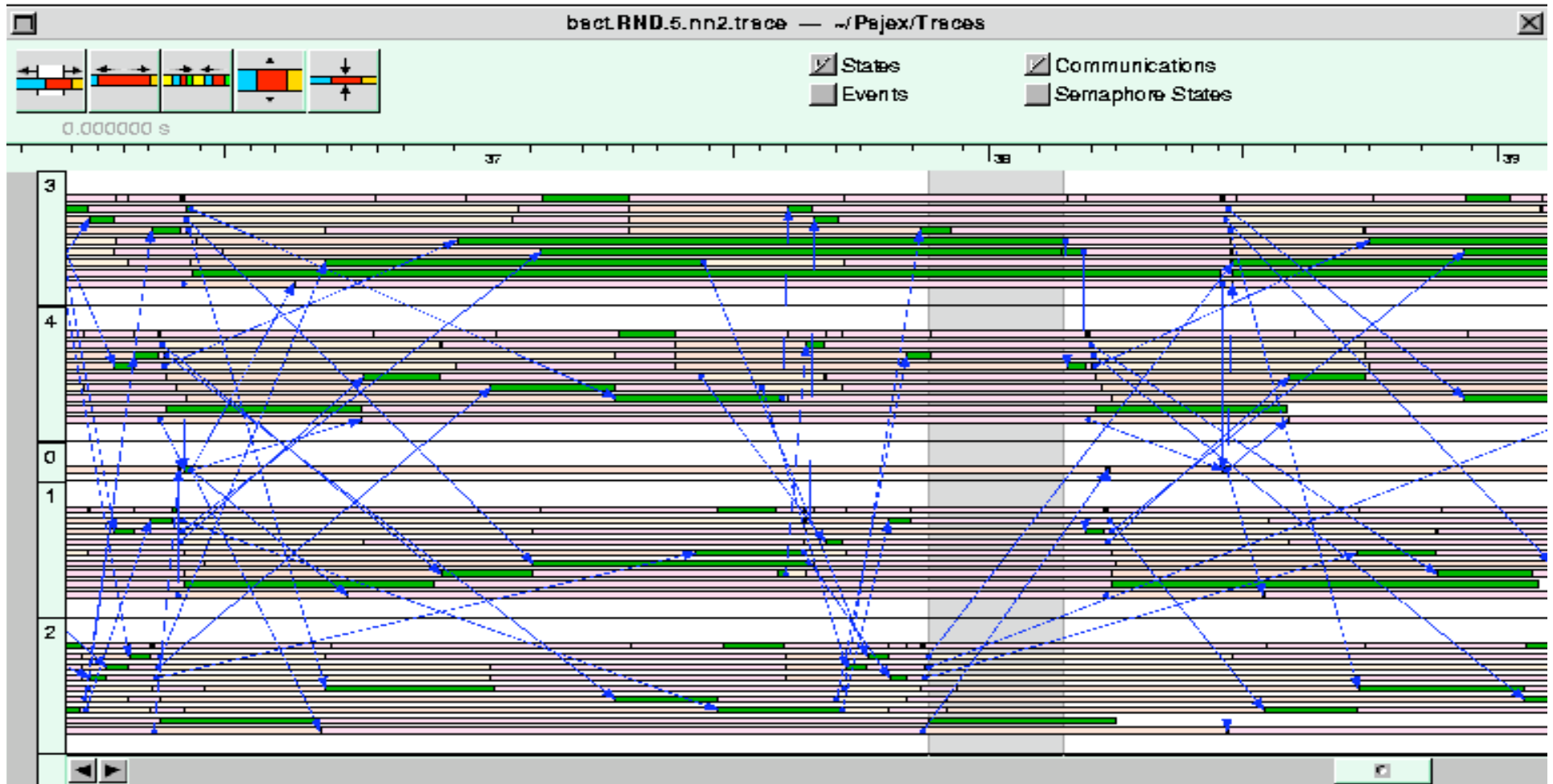
Race conditions



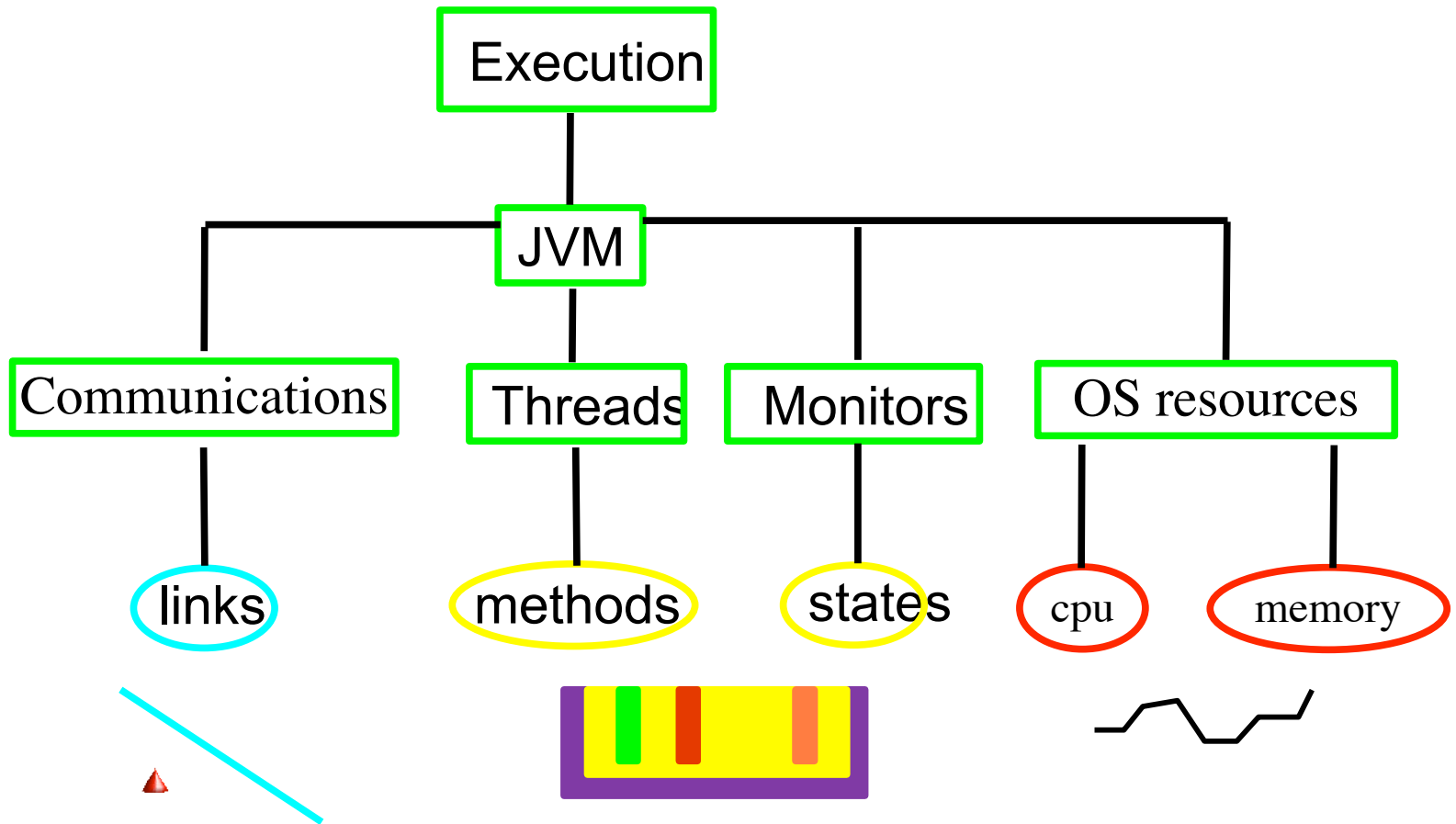
Conservative policy

Visualisation « Paje »

[J. Chassin, B. Stein]



Hierarchy of visual objects in « Pajé »

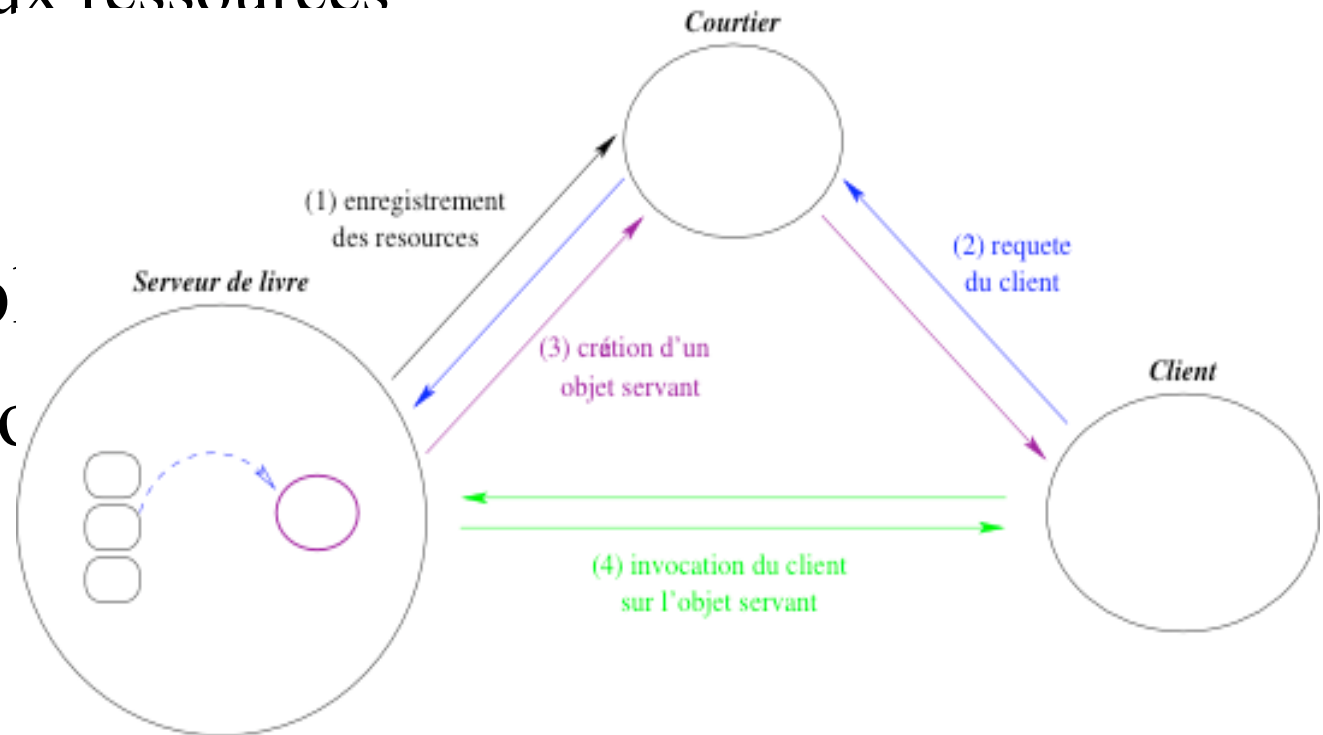


Traçage événementiel pour
l'analyse d'exécutions réparties :
de la théorie à la pratique

Application au serveur multimédia

- Serveur : livre, musique
- Courtier : enregistrement des ressources et requêtes
- Client : accès aux ressources

- Modèle d'app
- Modèle d'invoc



Problèmes à résoudre

■ Exécution locale :

- causalité des appels de méthode
- utilisation des objets de synchronisation
- gestion des ressources d'exécution : thread, GC...

■ Exécution répartie :

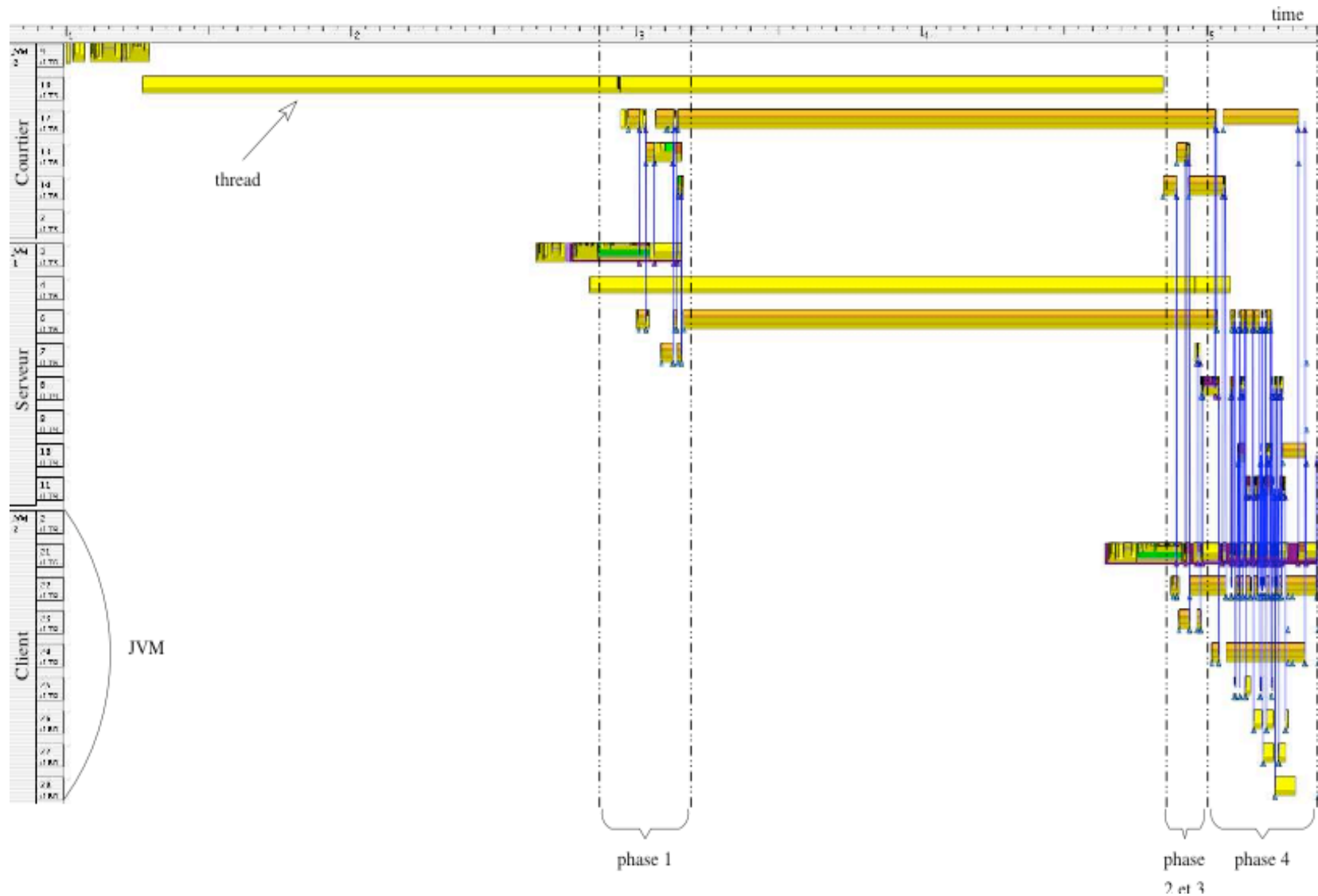
- interaction distante : communications
- blocage : du retour d'une invocation

■ Performances :

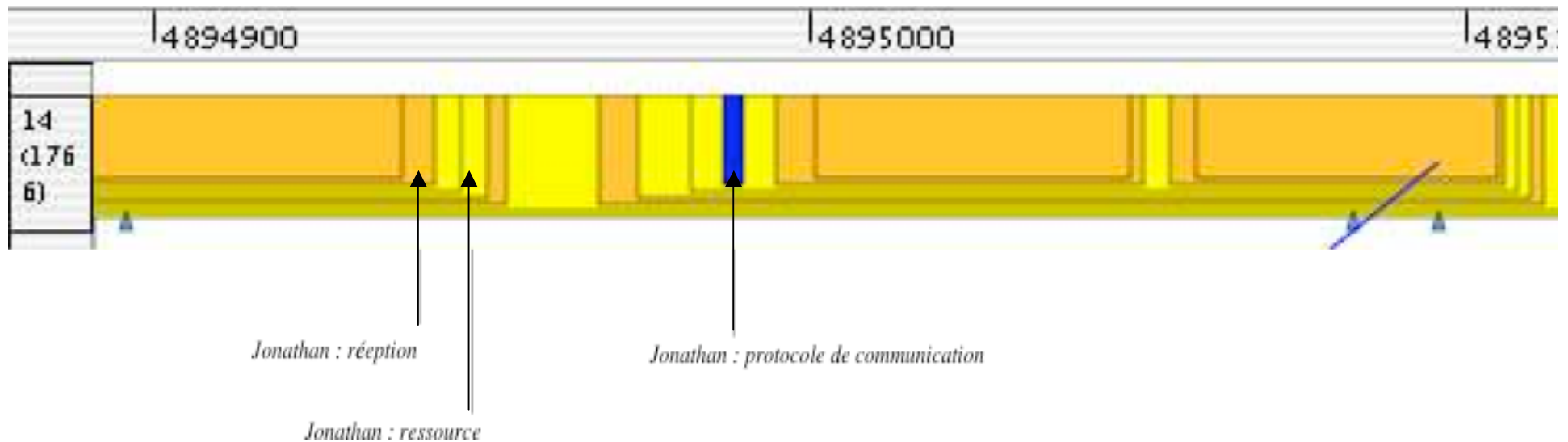
- utilisation des ressources systèmes : thread, CPU, mémoire
- utilisation des ressources d'infrastructure : ORB,

JVM

Interprétation d'une trace



Causalité des appels de méthode

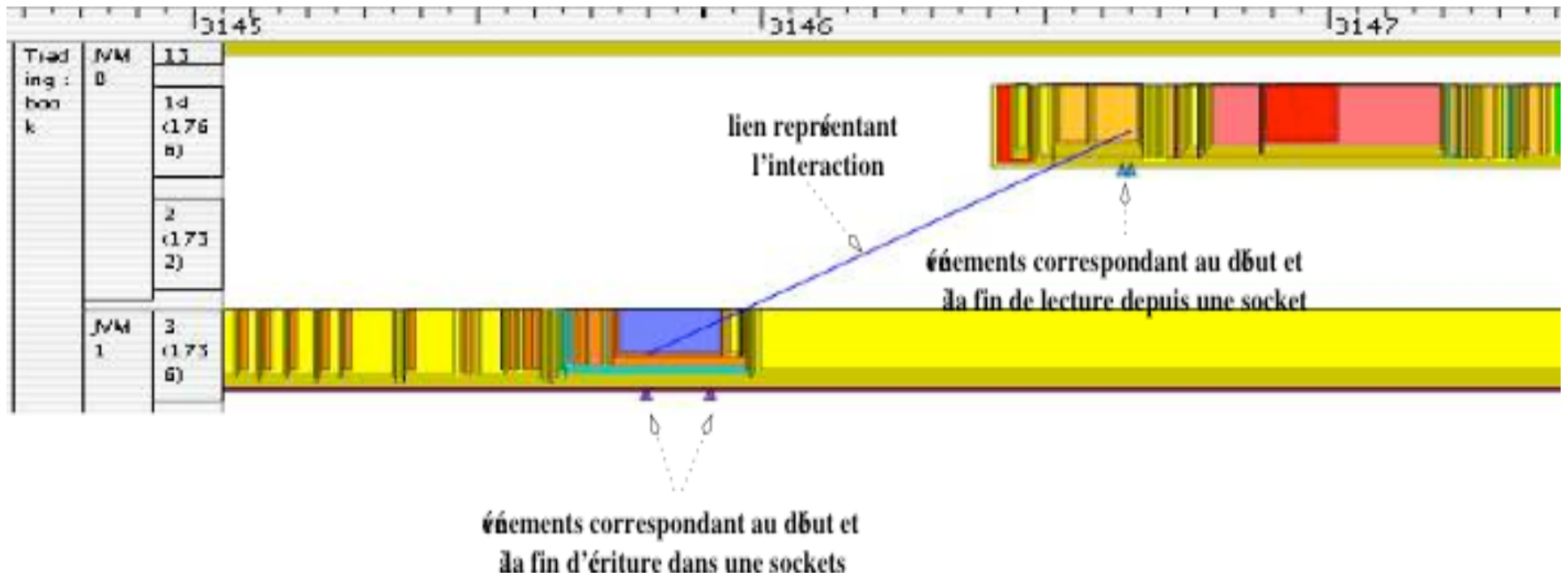


- Thread : exécution séquentielle
- Emboîtement = causalité
- Enchaînement = séquentialité

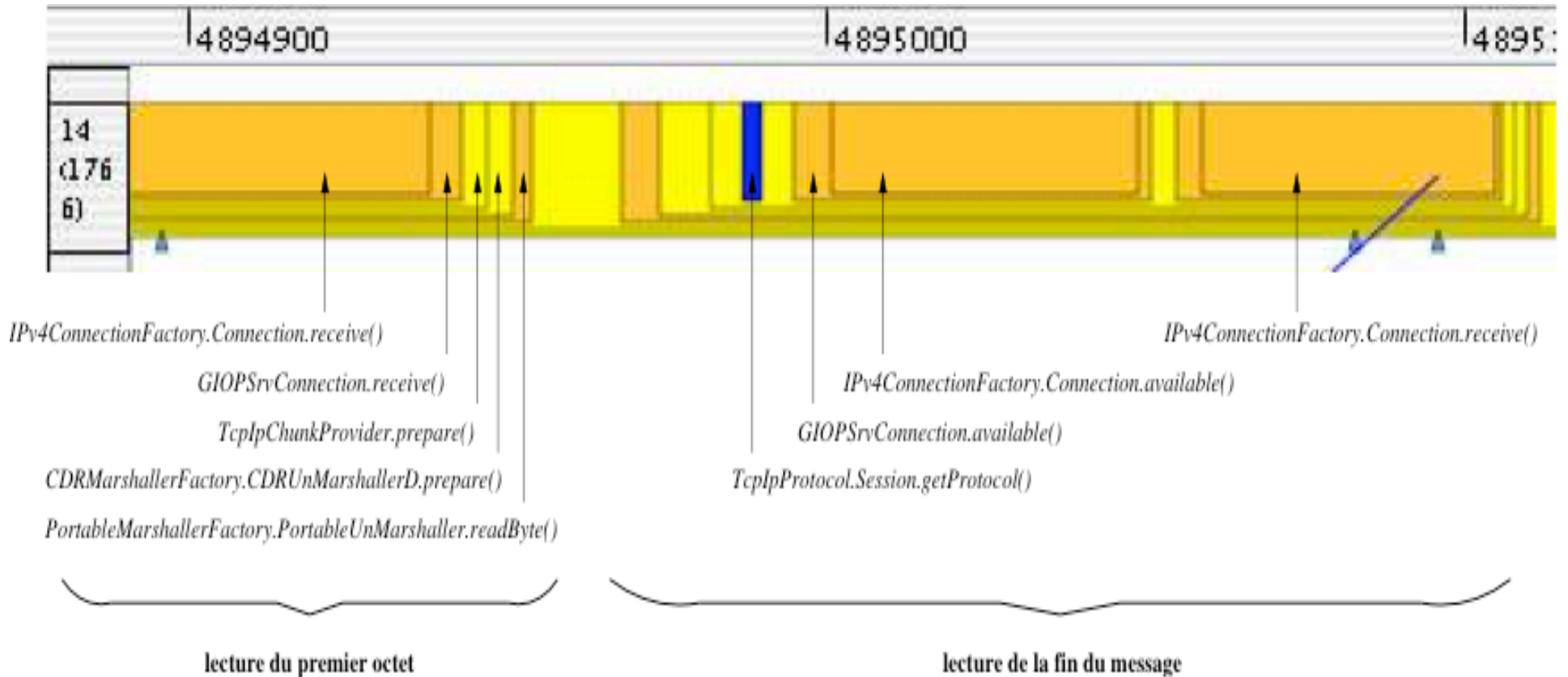
Objet de synchronisation



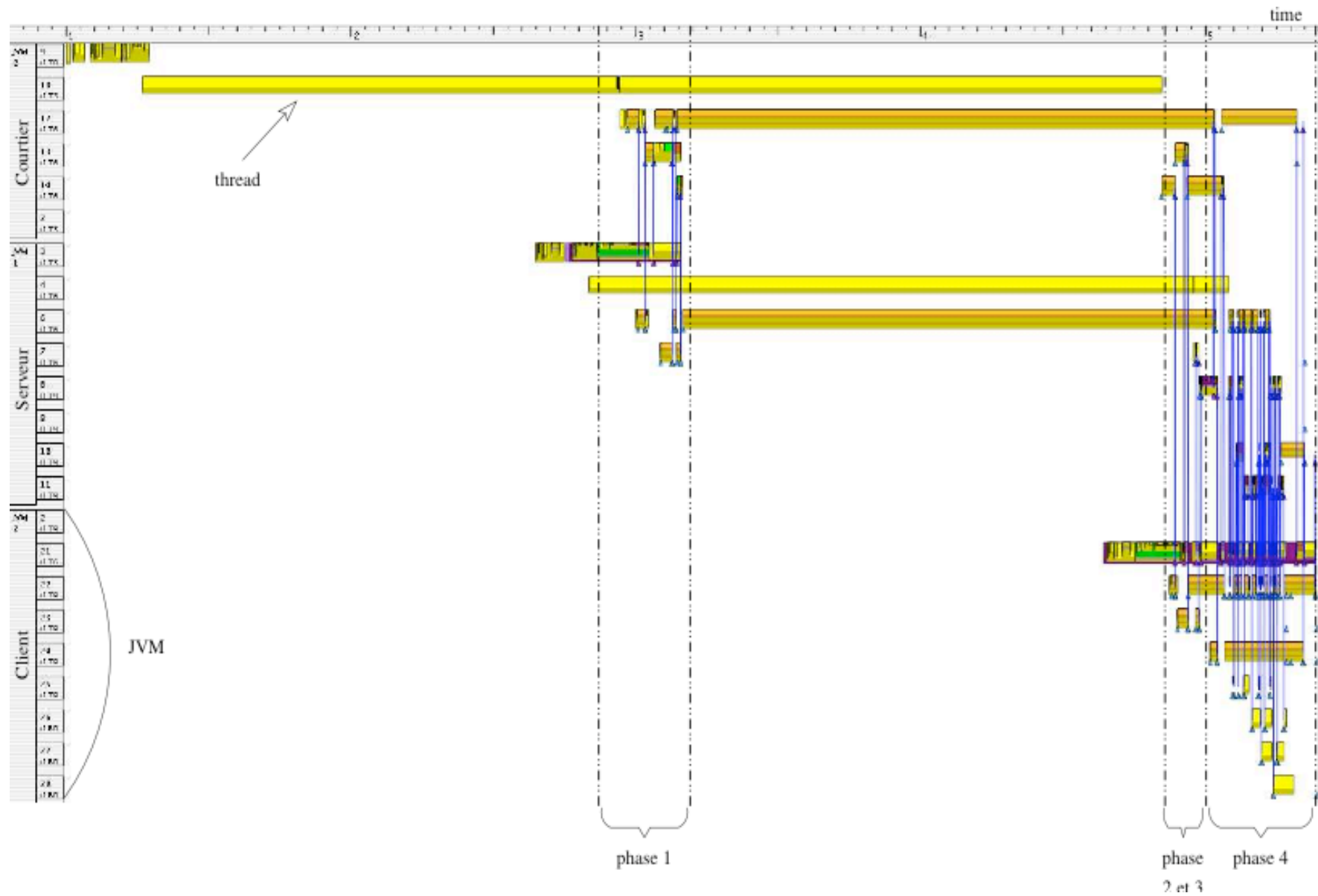
Interaction distante



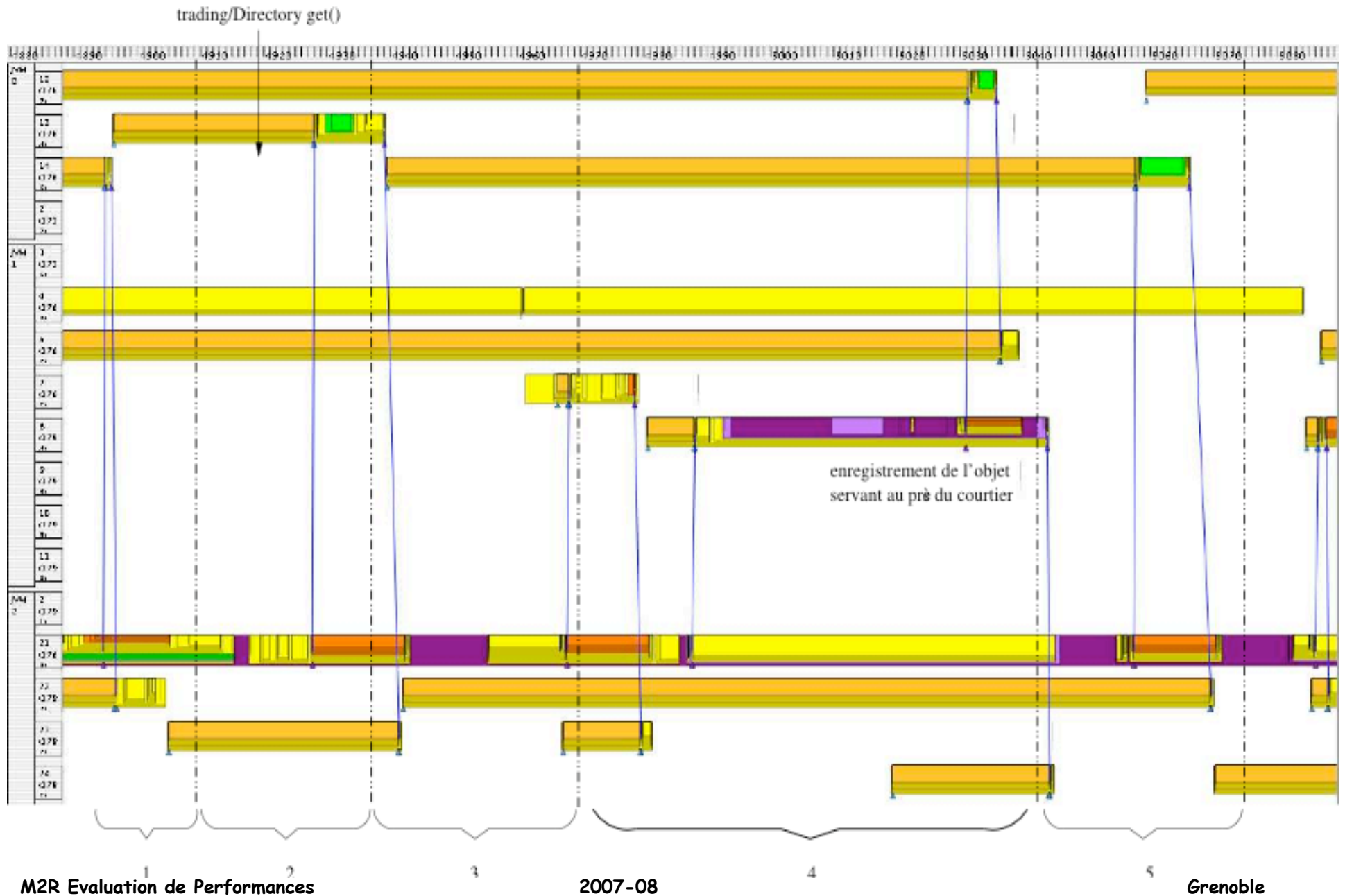
Pattern de réception



Vue globale



Interprétation de la phase 2



Utilisation des ressources systèmes



Approche d'analyse proposée

Phase d'étude

Identification des besoins
Identification du modèle de performances à étudier

Phase d'instrumentation

Instrumentation du programme à observer

Phase d'observation

Exécution du programme et observation
des évènements (traçage des évènements)

Phase de traitement

Fusions, correction, compression des traces

Phase de présentation et d'analyse

Visualisation et analyse des
traces collectées

État de l'art et choix technologiques

Phase d'étude

Phase d'instrumentation

Phase d'observation

Phase de traitement

Phase de présentation
et d'analyse

•Qu'est ce qu'on veut étudier?

->Les threads et les objets de synchronisation Java

•Quel est le modèle de performances à étudier?

->Le modèle de programmation concurrente

Portable/non portable

Code source/code binaire

Automatique/manuelle/semi-automatique

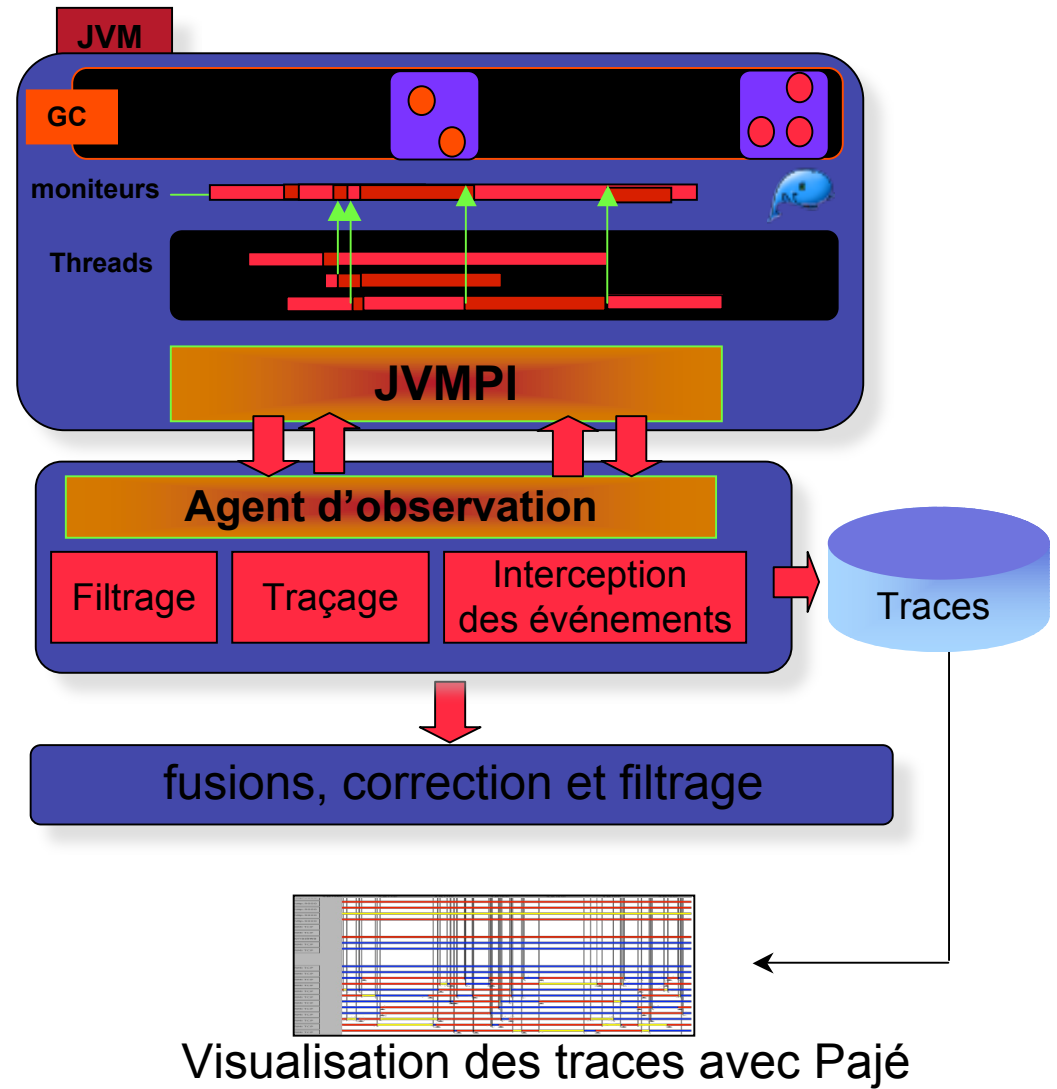
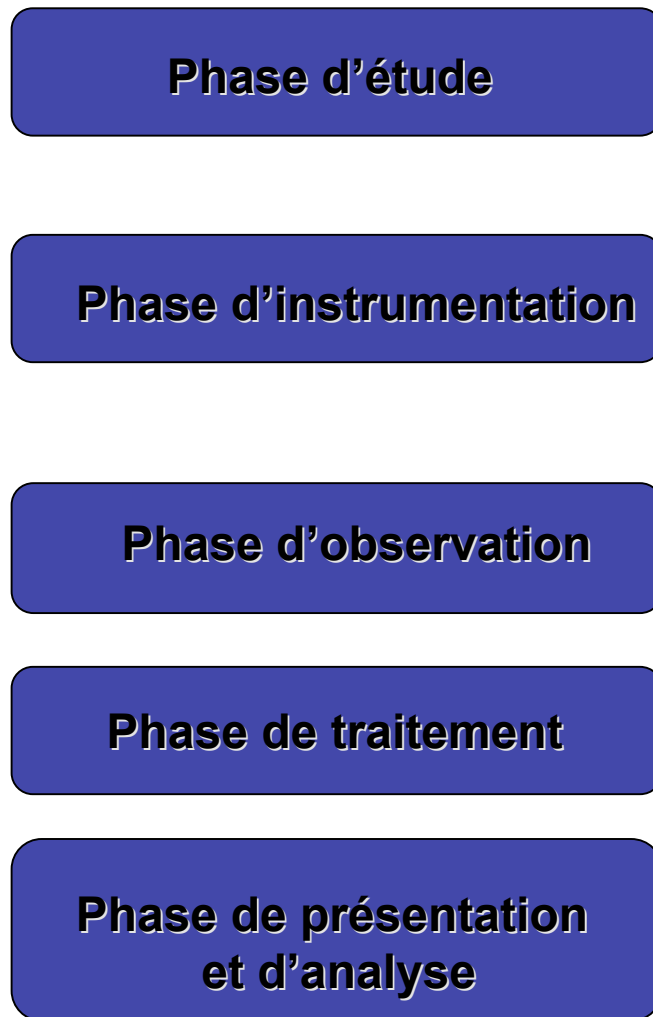
Traçage/échantillonnage/chronométrage/comptage

Capture événementielle/à la demande/périodique

Filtrage à la volé + filtrage post-mortem

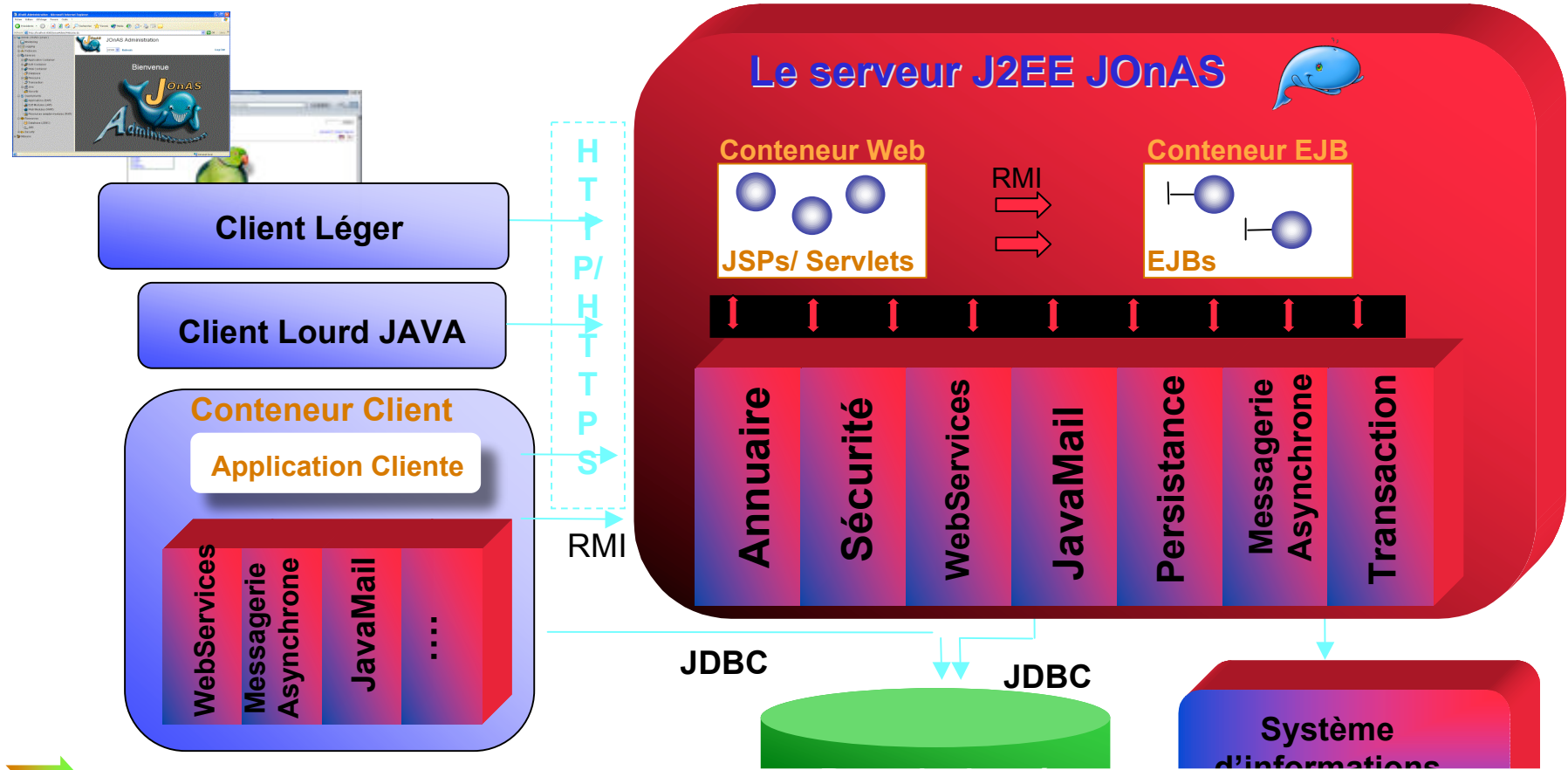
Pajé

Plateforme proposée



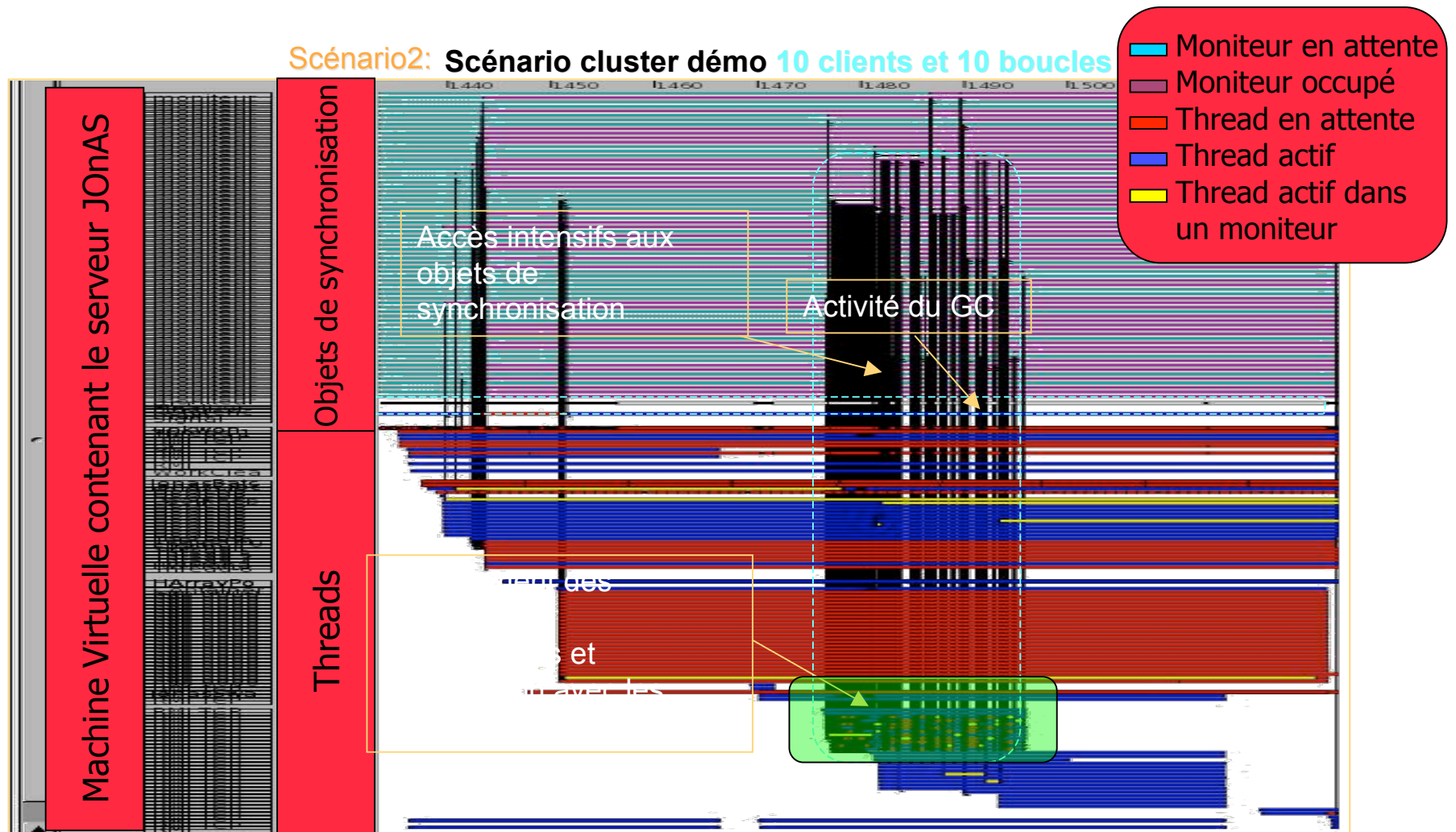
Le serveur J2EE JOnAS (2)

Architecture

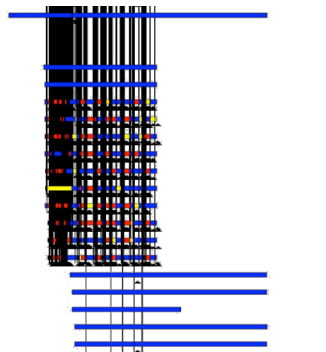


JOnAS est l'exemple type d'un programme complexe très intéressant à étudier dans le cadre du test logiciel et l'analyse des comportements et des performances

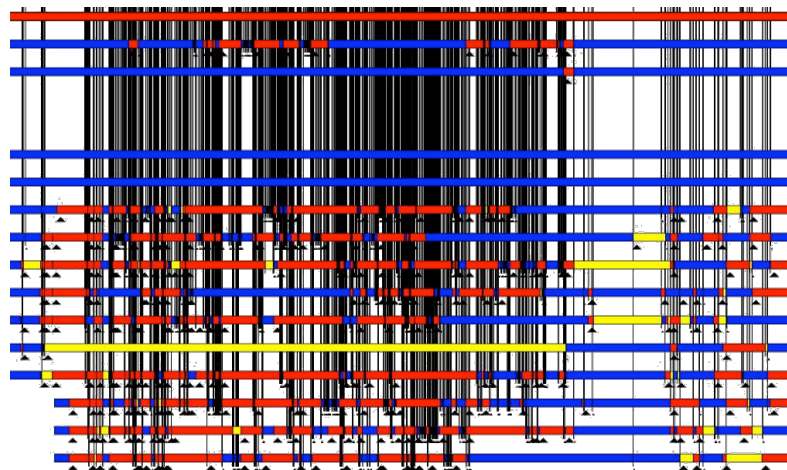
Exemple de visualisation des threads JOnAS



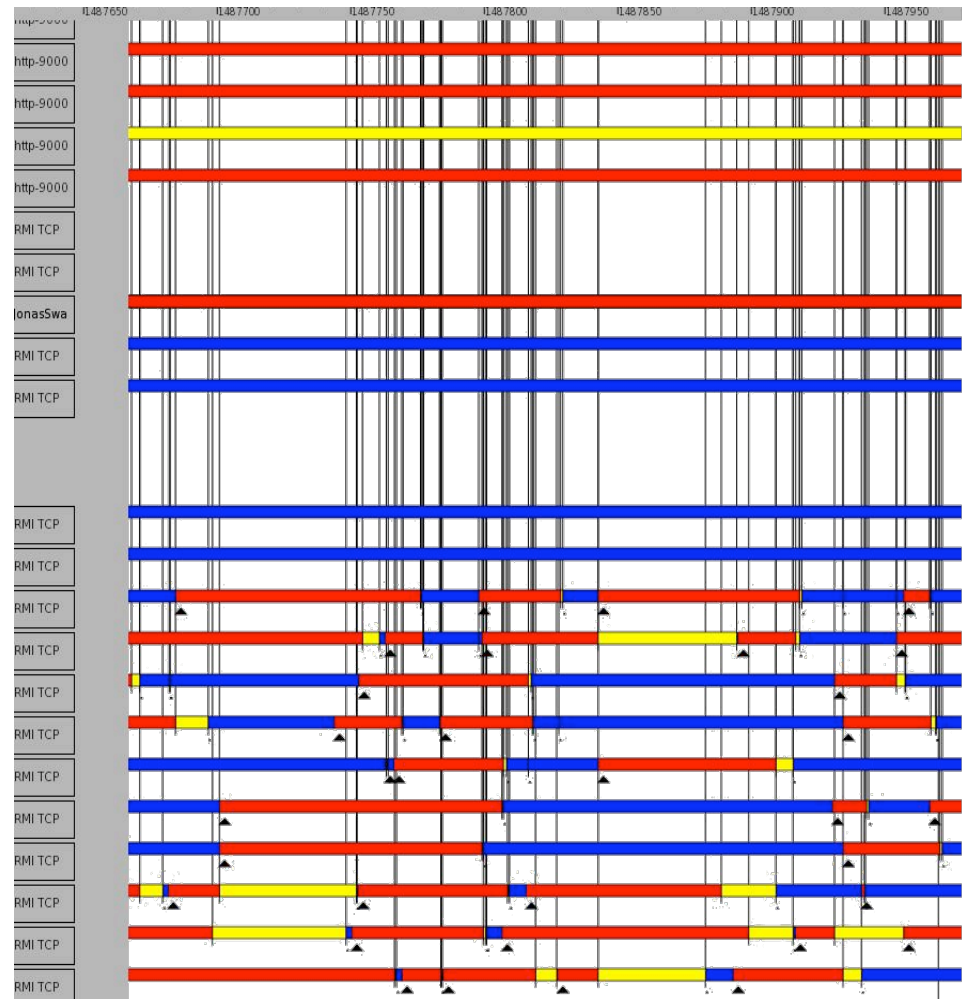
Agrandissement des visualisations



Visualisation initiale



Premier niveau d'agrandissement

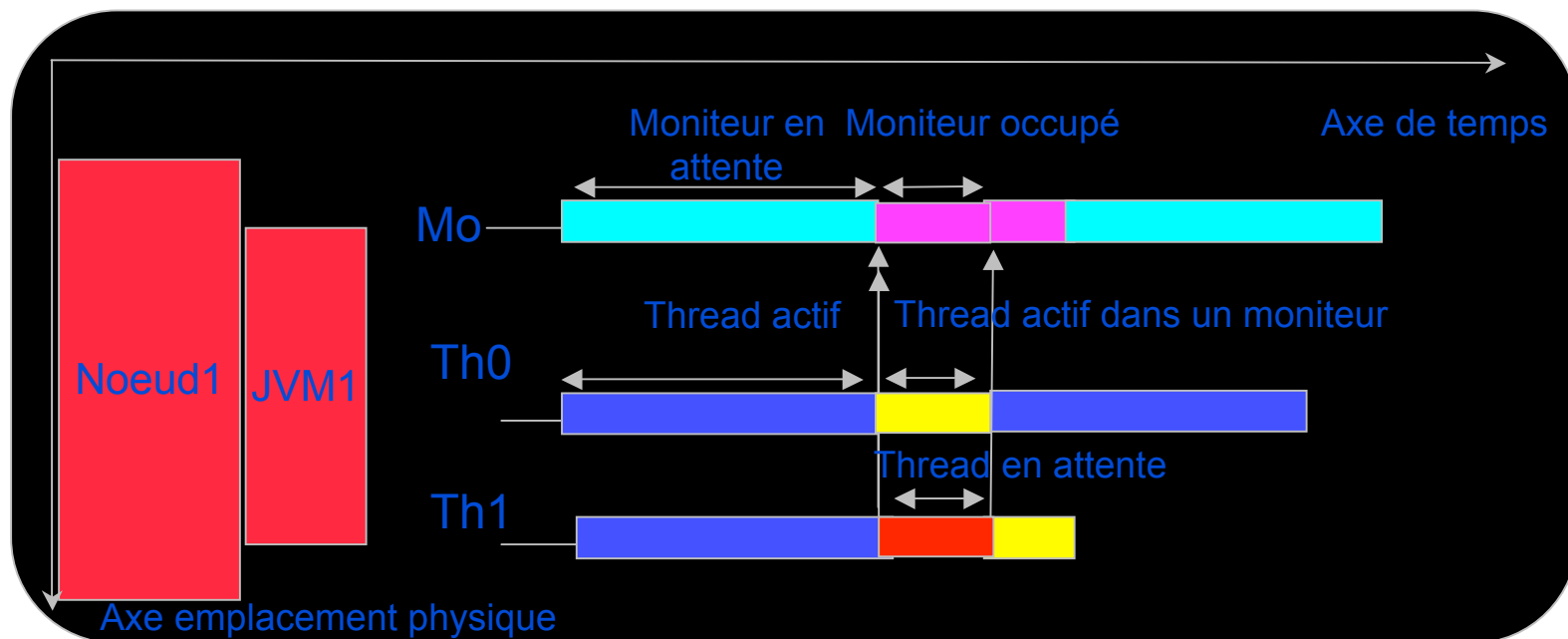


Deuxième niveau d'agrandissement

Interprétations des visualisations

Pajé permet la visualisation :

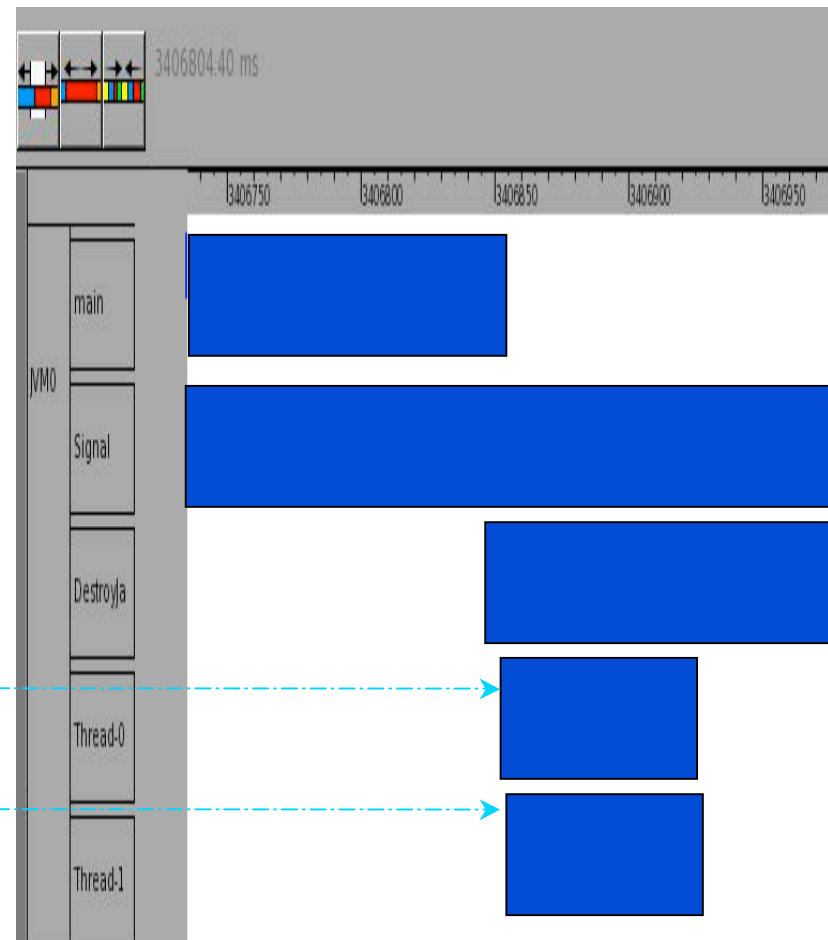
- des threads dans un diagramme spatio-temporel
- des objets de synchronisation
- de l'interaction entre threads et objets de synchronisation



Correspondance entre code Java et visualisation Pajé

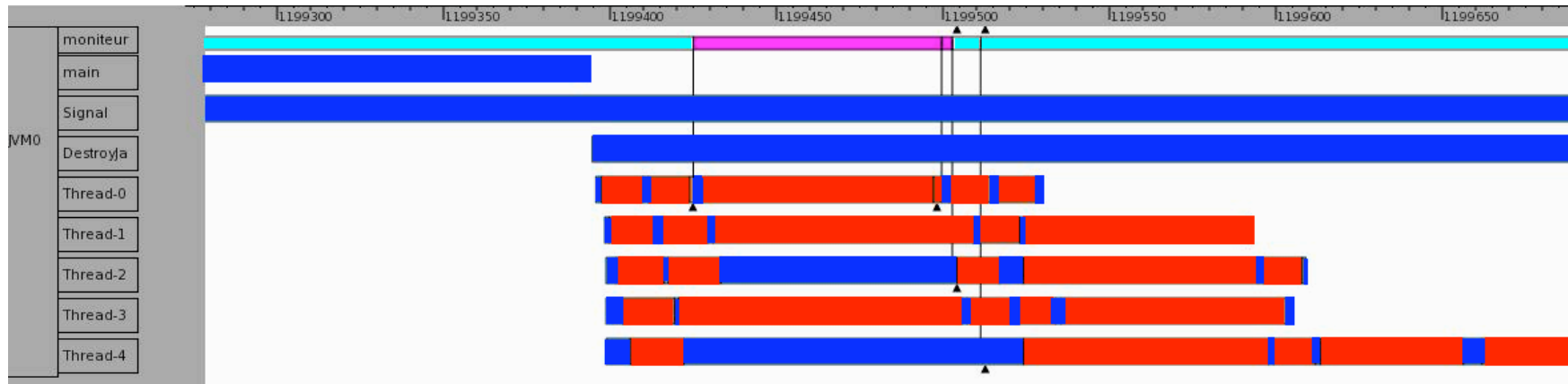
```
public class Demo1 extends Thread {  
    public int x;  
    public Demo1(int param){  
        x=param;  
    }  
    public synchronized void Task1(){...}  
    public void run(){  
        for (int i=0;i<5;i++){  
            Task1();  
        }  
    }  
  
    public static void main(String[] args) {  
        Demo1 Thread0=new Demo1(1);  
        Thread0.start();  
        Demo1 Thread1=new Demo1(2);  
        Thread1.start();  
    }  
}
```

Code Java

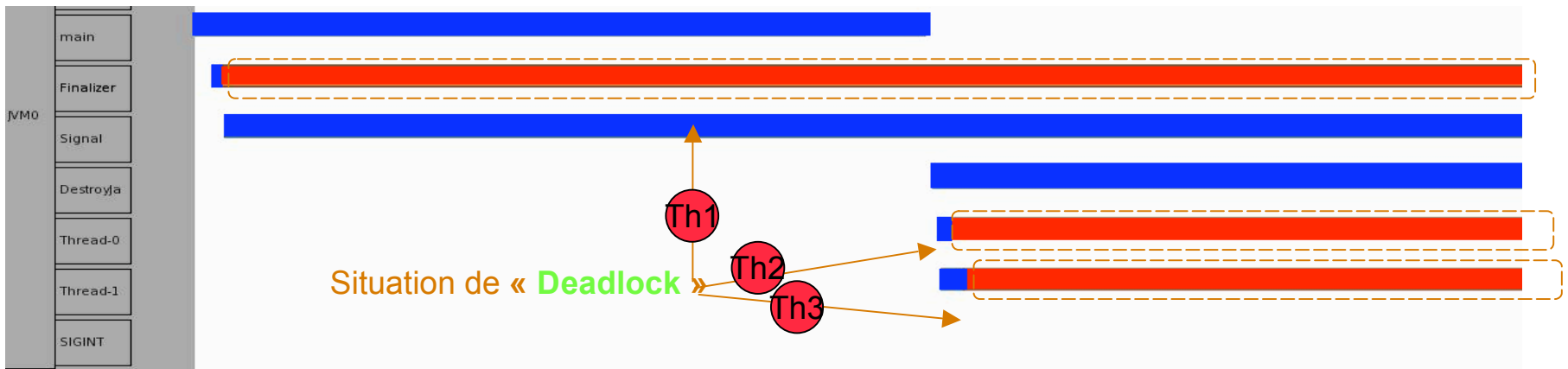


Environnement de visualisation Pajé

Exemples de situations possibles



Entrelacement des threads



Inter blocage

Évaluation de la plateforme

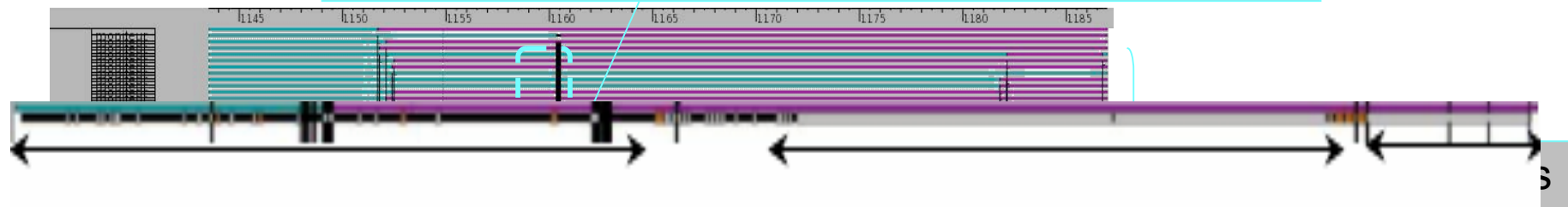
Taux de perturbation:

« Le surcoût du à l'instrumentation du programme »

	Nbr évènements	Temps d'exécution sans instrumentation (Milliseconde)	Temps d'exécution avec instrumentation (Milliseconde)	Taux de perturbation
Démarrage et arrêt de JOnAS	2081	16.341	29.763	82.112%
Scénario src	4008	27.57	32.778	18.89%
Scénario Cluster Démo (10 Clients/10 boucles)	26253	81.311	151	85,7%
Scénario Cluster Démo (10 Clients/100 boucles)	81077	462.677	562.92	17.80 %

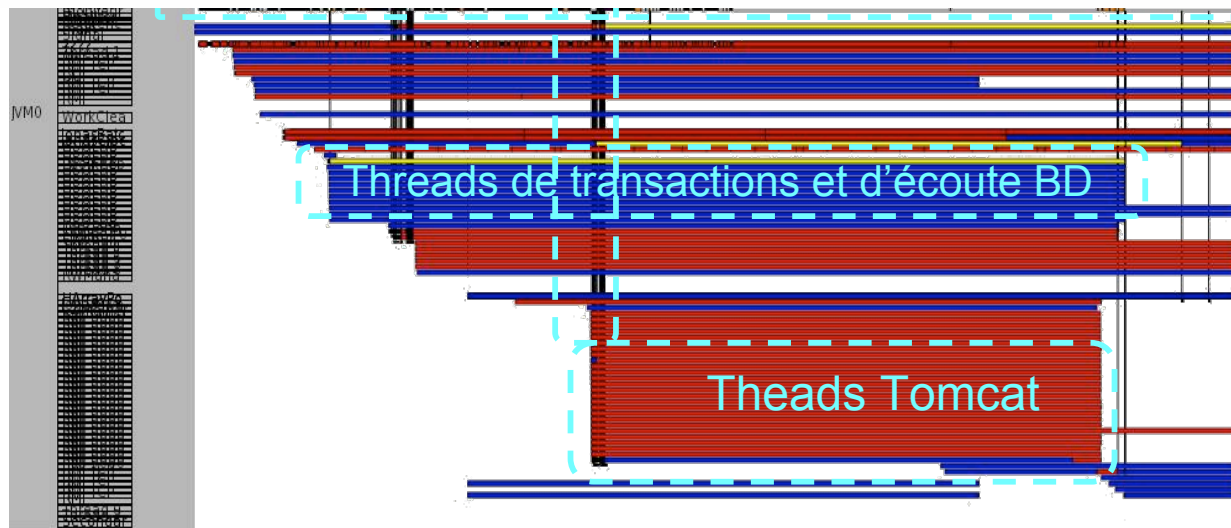
Analyse de l'activité du GC de JOnAS

Communication Threads-Objects de synchronisation



Activité intensive du GC

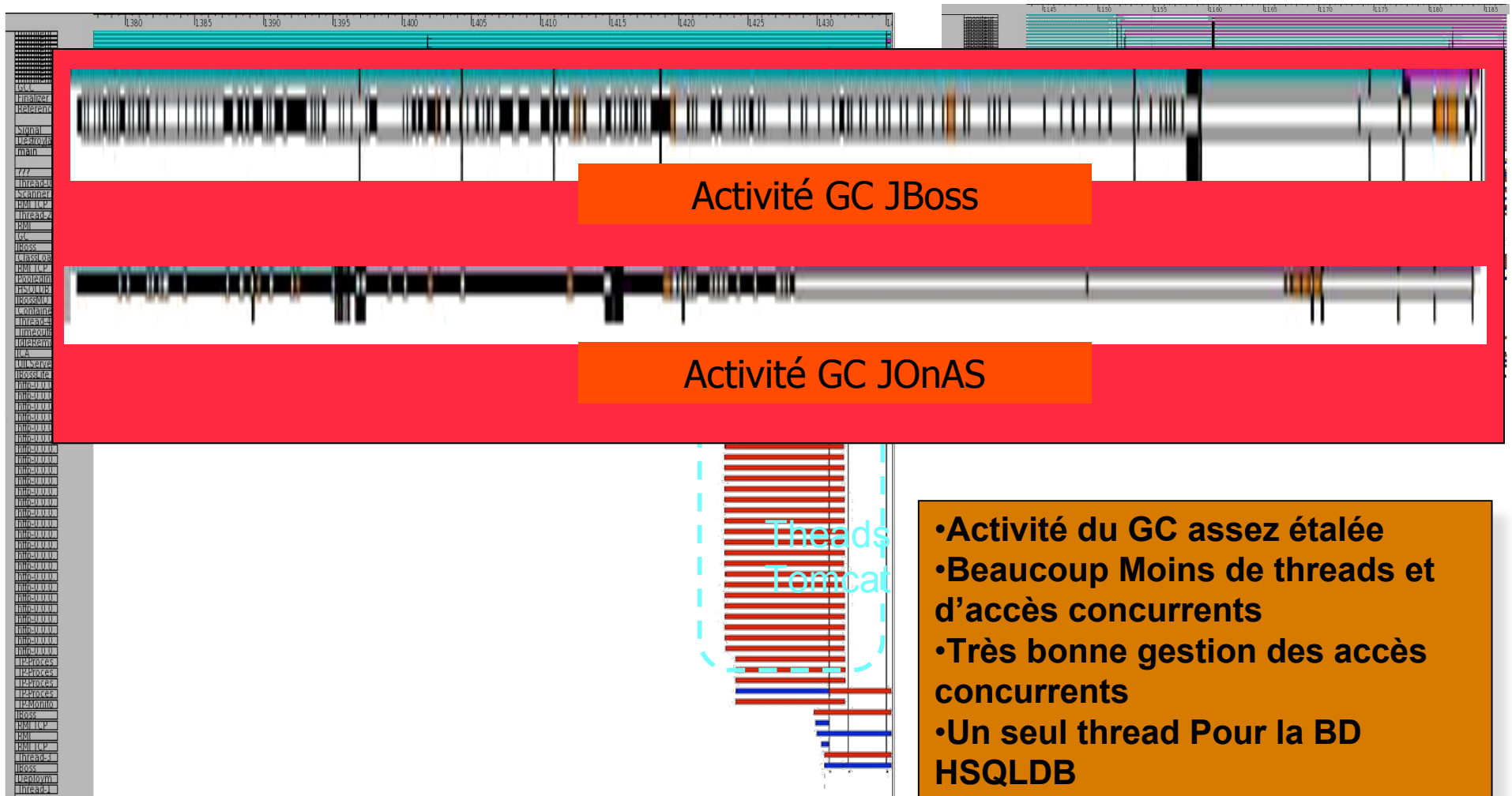
Stabilisation de l'activité du GC



visualisation de l'activité du GC

Liste des Threads JOnAS

Comparaison avec le serveur JBoss



Scénario démarrage du serveur JBoss

Perspectives

Observation : Définition d'un modèle
état/événements **MULTI-NIVEAUX**

Observation = perturbation \leftrightarrow compromis

Challenges :

facteurs d'échelle : clusters, grids, grid5000;

niveaux hétérogènes pour l'instrumentation

hétérogénéité des environnements

Analyse de données et visualisation

recherche de motifs automatique, abstraction

automatique, modèles de programmation...

Sonorisation

■ Possibilités du son :

- traitement passif : repérage d 'anomalie, schéma répétitif, etc.
- combinaison de sons (traitement parallèle par exemple).
- Plusieurs dimensions utilisables : timbre, hauteur, durée, intensité, positionnement spatial.

■ Représentation du comportement temporel d 'un programme : communications, charge des processeurs, etc.

■ Exemple

