# DIET: New Developments and Recent Results

A. Amar[1], R. Bolze[1], A. Bouteiller[1], A. Chis[1],
Y. Caniou[1], E. Caron[1], P.K. Chouhan[1], G.L. Mahec[2],
H. Dail[1], B. Depardon[1], F. Desprez[1], J. S. Gay[1], A. Su[1]

LIP Laboratory (UMR CNRS, ENS Lyon, INRIA, UCBL 5668) / GRAAL Project

LPC / PCSV (CNRS / IN2P3 UBP Clermont-Ferrand)

29 August 2006
CoreGrid Workshop

*Core* **GRID**

## Outline

1. Introduction

2. Distributed Interactive Engineering Toolbox

3. DIET Applications

4. Discussion

Presenter- Pushpinder Kaur Chouhan     DIET: New Developments and Recent Results

## Outline

1 Introduction

2 Distributed Interactive Engineering Toolbox

3 DIET Applications

4 Discussion

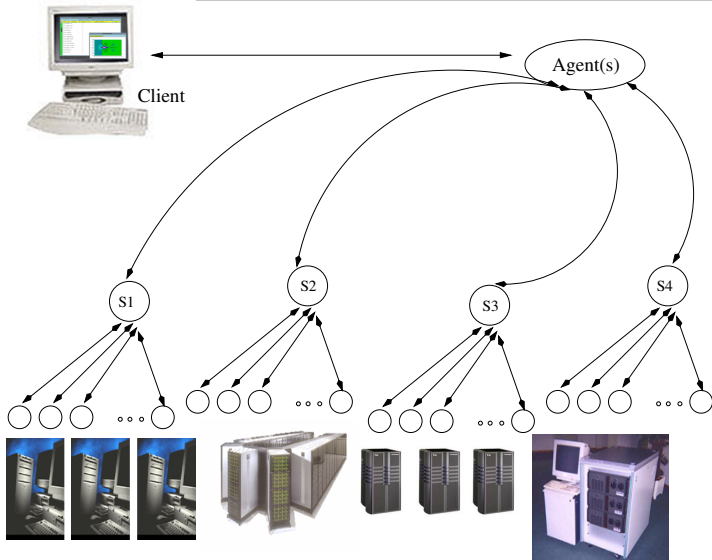Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results
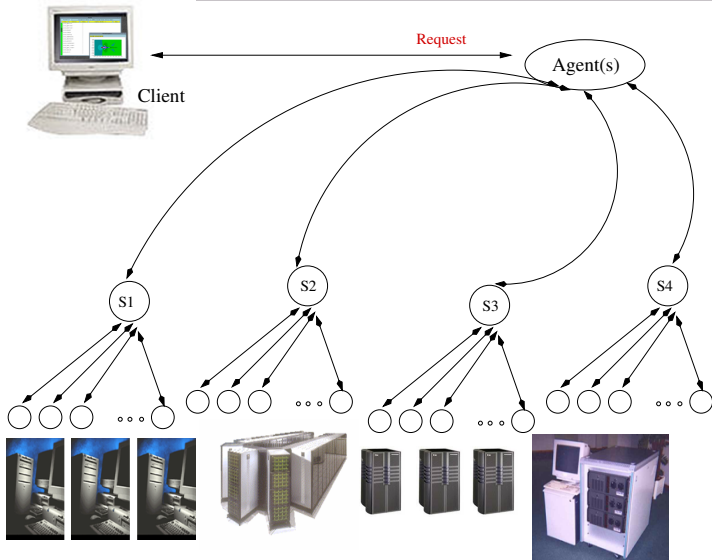
## Why we need tools like DIET?

- To access large computational power and/or storage capacity through the Internet.
- Large scale computing over heterogeneous platforms was coined by Ian Foster in the mid-90's - Grid.
- Implementation of RPC programming model over the Grid - GridRPC.
- Existing implementations: NetSolve, Ninf, DIET, XtremWeb, OmniRPC, · · ·

# GridRPC Paradigm

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

# GridRPC Paradigm

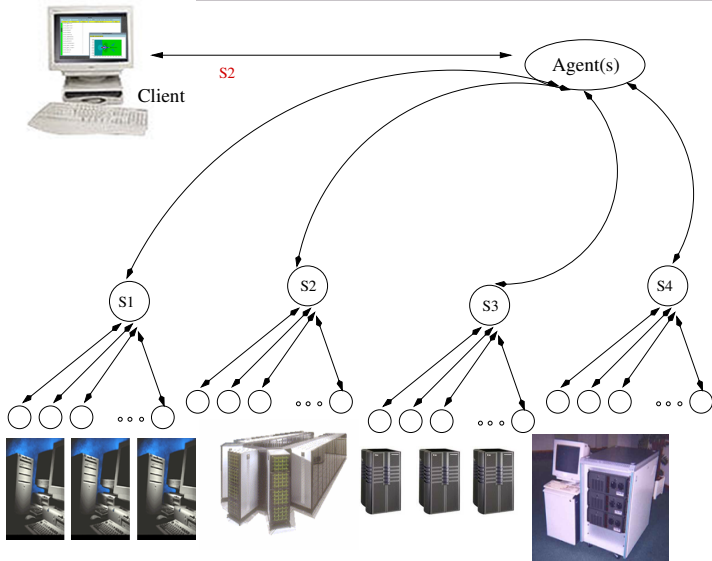Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

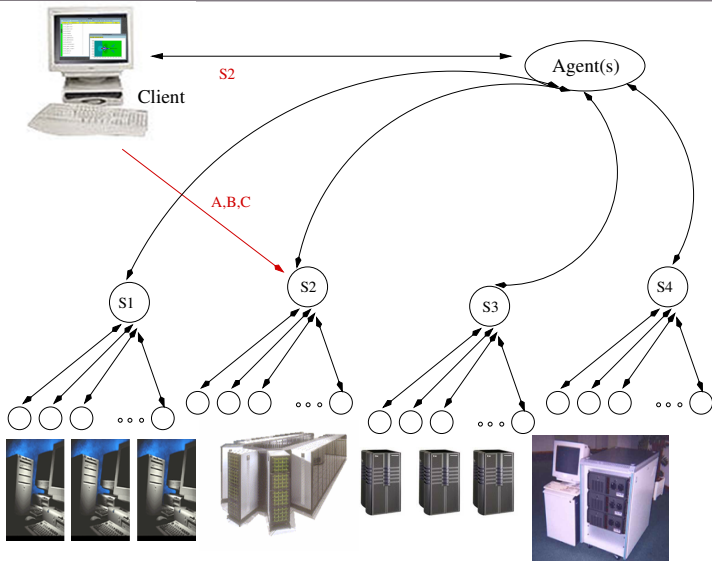# GridRPC Paradigm

# GridRPC Paradigm

# GridRPC Paradigm

# GridRPC Paradigm

## Why DIET?

Hierarchical architecture for an improved scalability.

## Why DIET?

Hierarchical architecture for an improved scalability.

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# Outline

**1** Introduction

**2** Distributed Interactive Engineering Toolbox
- DIET Scheduling
- DIET Deployment
- Fault Tolerance in DIET
- DIET Visualization

**3** DIET Applications

**4** Discussion

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET Components

- Client : An application that uses
  DIET to solve problems

  $\boxed{\text{Client}}$

- Master Agent (MA) :

  - Receives requests from clients

  - Collects computational abilities
    from servers and selects

  - Chosen server to the client

- Local Agents (LA) :

  - Act as transmitter

  - Share the workload of scheduling

- Servers (SeD) : Perform actual
  computation for client

*Core* GRID

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Store the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

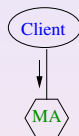## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Store the workload at scheduling
- Servers (SeD) : Perform actual computation for client

Client

MA

CoreGRID

8/ 30

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

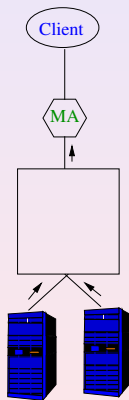## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Presenter- Pushpinder Kaur Chouhan      DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Presenter- Pushpinder Kaur Chouhan | DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
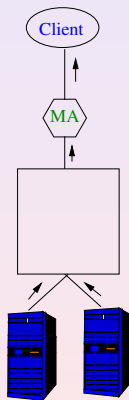DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Client

MA

LA

LA    LA

CoreGRID

8/ 30

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
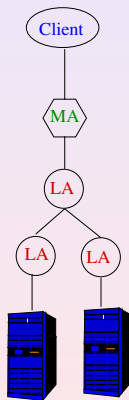DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
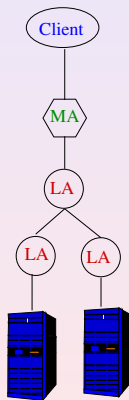DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

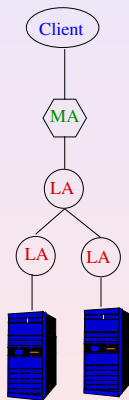## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
- Servers (SeD) : Perform actual computation for client

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

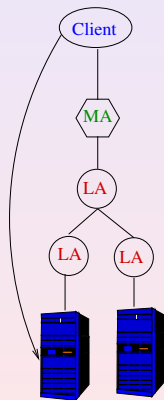## DIET Components

- Client : An application that uses DIET to solve problems
- Master Agent (MA) :
  - Receives requests from clients
  - Collects computational abilities from servers and selects
  - Returns the reference of the chosen server to the client
- Local Agents (LA) :
  - Act as transmitter
  - Share the workload of scheduling
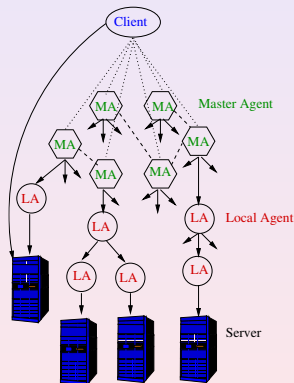- Servers (SeD) : Perform actual computation for client

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# **Co**llector of **R**esource **I**nformation (**CoRI**)

**CoRI-Easy**: provides basic measurements of the environment
**CoRI Manager**: manages the use of different collectors

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## Plug-in schedulers

- Applications vary in terms of performance demands
- Performance Estimation Vector (PEV) is used for scheduling
- SeDs sends PEV to MA as a response to a request

| Information tag starts with EST_ | Explanation |
|---|---|
| *TCOMP* | the predicted time to solve a problem |
| *LOADAVG* | CPU load average |
| *FREEMEM* | amount of free memory (Mb) |
| *NBCPU* | number of available processors |
| ⋮ | ⋮ |

Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# Verification of plug-in schedulers

- Six heterogeneous servers
- Sequential and independent requests
- Inter-arrival time for the request is 1 minute
- Compare two scheduler
  - Round Robin scheduler
  - CPU scheduler

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## Round Robin Scheduler

Presenter- Pushpinder Kaur Chouhan       DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## CPU Scheduler

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# Batch Scheduler Management

- System dependent
  - NFS: copy the code?
  - MPI: LAM, MPICH?
- Batch system dependent
  - No homogeneity
  - Scheduler behaviour
  - Information about the internal scheduling process
- Monitoring and performance prediction

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# Parallel and batch submissions
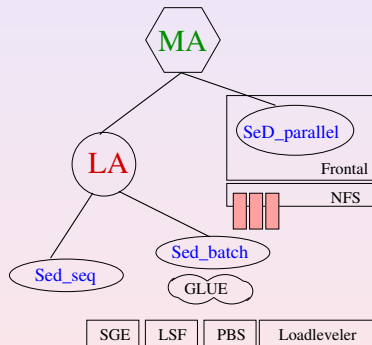
- Performance prediction: SIMBATCH
  - Batch, cluster, parallel tasks simulator plugin for SimGrid
  - Goal: test new grid schedulers more realistically and performance prediction
  - Implement FIFO, Conservative BackFilling, $\cdots$
- More problems
  - Asynchronous, long term production jobs
  - Performance prediction
    - How to decide number of processors for application?
    - If reservation available, how to compute deadline?
  - Co-scheduling?
  - Data and job migration?

*Core* **GRID**

15/

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## Workflow Management

- Workflow is represented as DAG
- Use different heuristic methods to solve scheduling problems
- Extensibility to address multi-workflows submission and large grid platform
- Manage heterogeneity and variability of environment

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET workflow engine

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET workflow engine

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET workflow engine

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET workflow engine

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

**DIET Scheduling**
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# DIET workflow engine

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## DIET Deployment

Mapping of DIET's components on available resources.

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## DIET Deployment

Mapping of DIET's components on available resources.

- Homogeneous resources

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
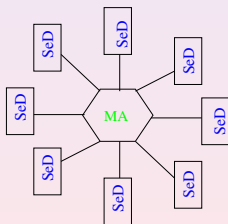DIET Visualization

## DIET Deployment

Mapping of DIET's components on available resources.

- Homogeneous resources

## DIET Deployment

Mapping of DIET's components on available resources.

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## DIET Deployment

Mapping of DIET's components on available resources.

- Homogeneous resources

### Complete d-ary Spanning tree

E.Caron, P.K.Chouhan,H.Dail and F.Vivien. Automatic Middleware Deployment Planning for Clusters. RR-2005-50, LIP,October 2005.

## DIET Deployment

Mapping of DIET's components on available resources.

- Heterogeneous resources

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## DIET Deployment

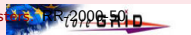Mapping of DIET's components on available resources.

- Heterogeneous resources

Find the best broadcast tree on a general graph, which is known to be NP-complete.

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

# DIET Deployment

Mapping of DIET's components on available resources.

- Heterogeneous resources



Thesis Automatic Deployment for Application Service Provider Environments by P. K. Chouhan

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

# GoDIET: DIET deployment tool

- XML file as Input

- Generate configuration files

- Launches services (name service, logging services)

- Launches DIET elements

- Destroy deployed platform

- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

## GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
**DIET Deployment**
Fault Tolerance in DIET
DIET Visualization

# GoDIET: DIET deployment tool

- XML file as Input
- Generate configuration files
- Launches services (name service, logging services)
- Launches DIET elements
- Destroy deployed platform
- Remote cleanup of launched processes

E.Caron, P.K.Chouhan and H.Dail GoDIET: A Deployment Tool for Distributed Middleware on Grid'5000. In IEEE,

editor, EXPGRID workshop, in conjunction with HPDC-15. pages 1-8.

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization

## Failure Detection in DIET

Failure detection

- Detection time: time between failure and definitive suspicion
- Accuracy: probability of the observer to be true at random time

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
**Fault Tolerance in DIET**
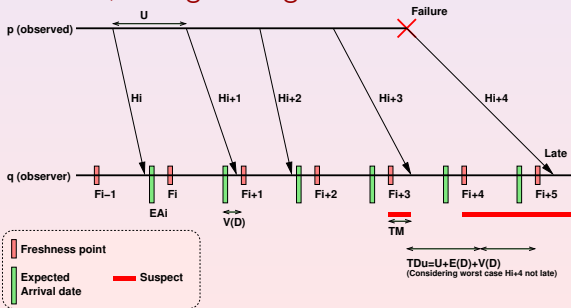DIET Visualization

## Failure Detection in DIET

Failure detection

- Detection time: time between failure and definitive suspicion
- Accuracy: probability of the observer to be true at random time

    Chandra, Toueg and Aguilera Failure Detector

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
DIET Visualization
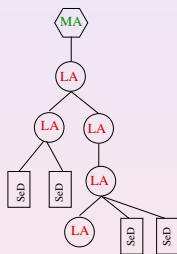
## Failure Detection in DIET

Failure detection

- Detection time: time between failure and definitive suspicion
- Accuracy: probability of the observer to be true at random time

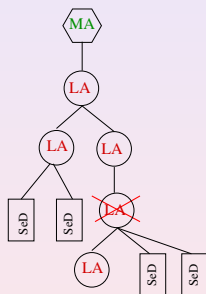### Chandra, Toueg and Aguilera Failure Detector

Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
**Fault Tolerance in DIET**
DIET Visualization

## Architecture Recovery

- Problem: Keep architecture connected
- Solution: Keep ancestors list
- tolerates up to f-1 simultaneous failures

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
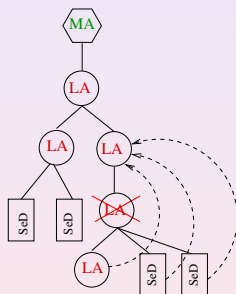DIET Deployment
**Fault Tolerance in DIET**
DIET Visualization

## Architecture Recovery

- Problem: Keep architecture connected
- Solution: Keep ancestors list
- tolerates up to f-1 simultaneous failures

Presenter- Pushpinder Kaur Chouhan     DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
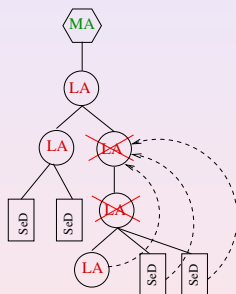DIET Deployment
**Fault Tolerance in DIET**
DIET Visualization

## Architecture Recovery

- Problem: Keep architecture connected
- Solution: Keep ancestors list
- tolerates up to f-1 simultaneous failures

Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
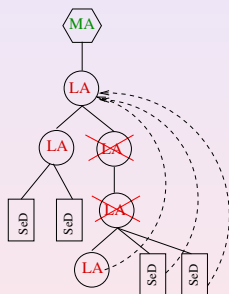DIET Deployment
**Fault Tolerance in DIET**
DIET Visualization

## Architecture Recovery

- **Problem**: Keep architecture connected
- **Solution**: Keep ancestors list
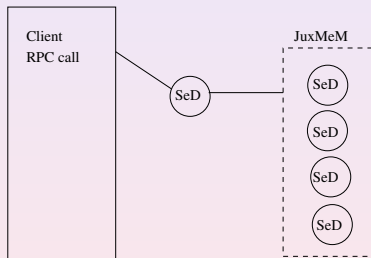- tolerates up to f-1 simultaneous failures

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
**Fault Tolerance in DIET**
DIET Visualization

## Architecture Recovery

- Problem: Keep architecture connected
- Solution: Keep ancestors list
- tolerates up to f-1 simultaneous failures

Presenter- Pushpinder Kaur Chouhan        DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
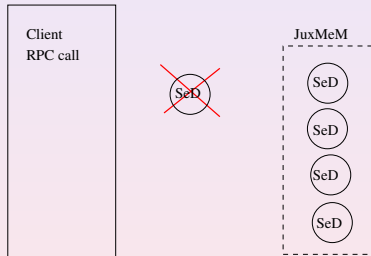DIET Visualization

## Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
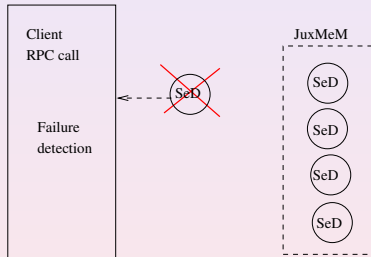Fault Tolerance in DIET
DIET Visualization

## Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan | DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
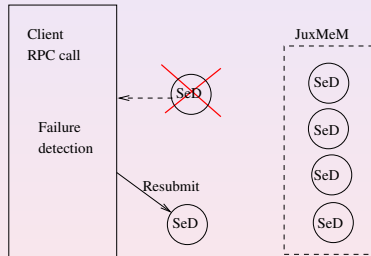DIET Visualization

## Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
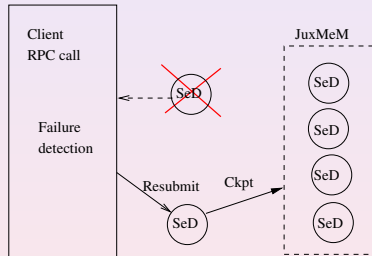Fault Tolerance in DIET
DIET Visualization

## Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan    DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
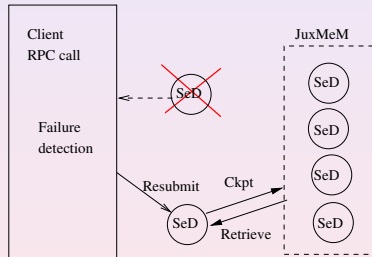**Fault Tolerance in DIET**
DIET Visualization

## Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
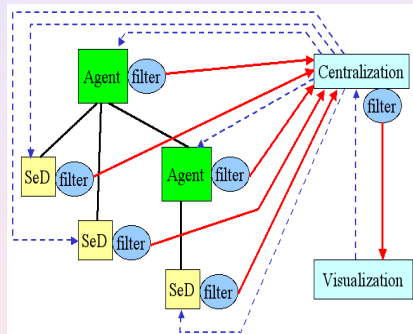**Fault Tolerance in DIET**
DIET Visualization

# Application Recovery

- Replication and checkpointing
- Replication divides the available computing by f
- Checkpointing takes periodic snapshot of process state, saving it to another place
- JuxMem manages data persistence across failures by replicating it on nodes

Presenter- Pushpinder Kaur Chouhan     DIET: New Developments and Recent Results

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
**DIET Visualization**

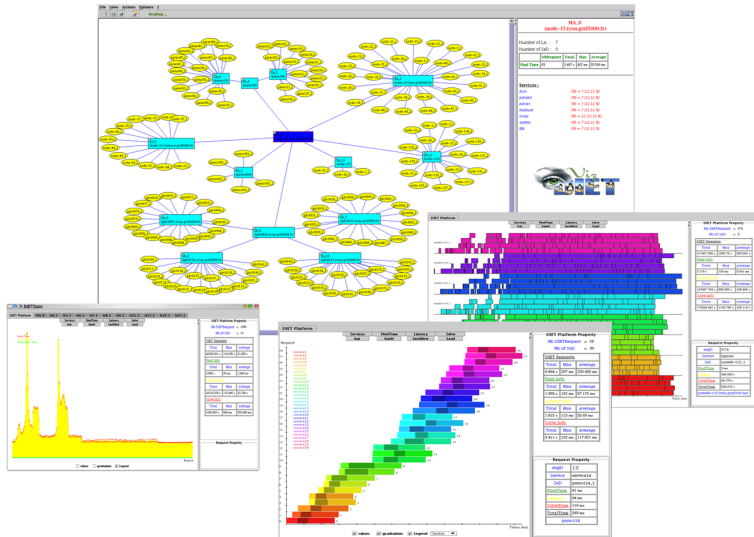# VizDIET: DIET Visualization Tool

- Current view of the DIET platform
- Communication between agents
- State of SeDs
- Scalability
- Available services
- Data persistency
- Name information
- CPU, memory and network load



R.Bolze, E.Caron, F.Desprez, G.Hoesch, and C.Pontvieux. A Monitoring and Visualization Tool and its
Applications. Computational Science and Its Applications-ICCSA 2006, volume 3984 of LNCS, pages 202-213.

Introduction
**Distributed Interactive Engineering Toolbox**
DIET Applications
Discussion

DIET Scheduling
DIET Deployment
Fault Tolerance in DIET
**DIET Visualization**

# Screen-shot of VizDIET

## Outline

**1** Introduction

**2** Distributed Interactive Engineering Toolbox

**3** DIET Applications

**4** Discussion

## Applications

- BLAST application using DIET
  - Basic Local Alignment Search Tool.
  - Compare biological sequences such as nucleotides sequences
  - N-sequences versus one database
  - Multi-request files are partitioned intro several smaller requests
- DIET to analyse Cosmological Results
  - Study large scale structure and galaxy formation
  - RAMSES collects raw data
  - Galices analysis the data

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
**Discussion**

Conclusion
Future Work

# Outline

Core GRID

Presenter- Pushpinder Kaur Chouhan     DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
**Discussion**

**Conclusion**
Future Work

## Conclusion

- Scalable, open-source, and multi-application platform
- Concentration on several issues
  - performance evaluation (CoRI, FAST )
  - scheduling (plugin schedulers, workflow management)
  - deployment (planning and tool GoDIET)
  - data management and replication (JuxMem)
  - fault toleration (architecture and application)
  - monitoring (LogService and VizDIET)
- Large scale validation on the Grid'5000 platform

http://graal.ens-lyon.fr/DIET          http://www.grid5000.org

Presenter- Pushpinder Kaur Chouhan          DIET: New Developments and Recent Results

Introduction
Distributed Interactive Engineering Toolbox
DIET Applications
**Discussion**

Conclusion
**Future Work**

# Future work

- Need of implementation and validation of algorithms from the scheduling literature in real-life middleware infrastructures?
- Still room for fundamental research on algorithms
- Some problems waiting for solution
    - Finding accurate models
    - Large scale validation of algorithms (simulators, real grids?)
- Take a look at other applications
- ...