

***A RISC approach  
to the GRID***



**M. Danelutto**  
*University of Pisa*

IFIP 10.3 Nov 1<sup>o</sup>, 2003



# Contents

---

- GRID: current status
- A RISC grid core
- Current experiments
- Conclusions



# Contents

---

- *GRID: current status*
- A RISC grid core
- Current experiments
- Conclusions

# ***GRID: current status***

## **Contents**

---



- Characterization of GRID
- Tools
- Abstract machine view
- Current, “GRID aware” applications



# Characterization of GRIDs

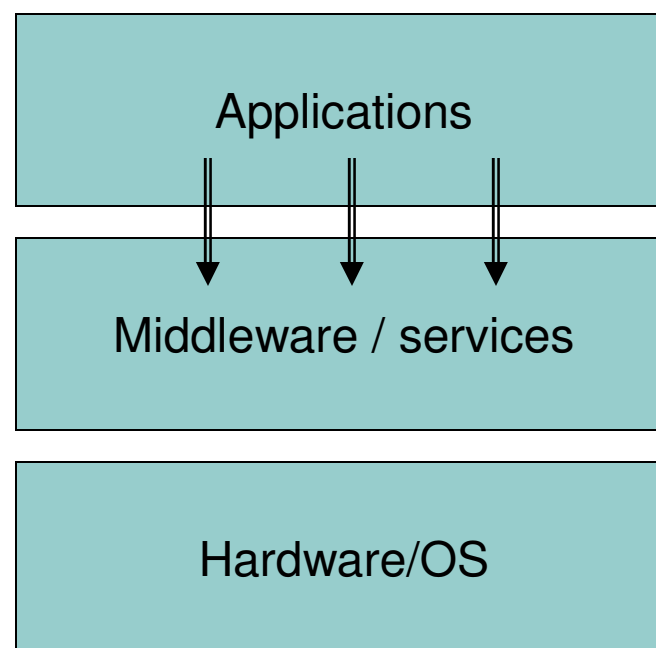
---

- Target architectures for complex computations
  - Complex, multidisciplinary, multilanguage ...
- Heterogeneous
  - HW & SW heterogeneous
- Dinamic
  - Instant (latency!) and long range (node up&down)
- Distributed
  - Geographic scale networks



# Tools

- Middleware
  - Between hw & sw !
- Many flavours
  - One winner
- Complex functionalities
  - Scheduling, resource manag., data/code staging, monitoring, etc.
- All is in charge to the appl. programmer !





# Abstract machine view

---

- Layer 0
  - OS services, common core (TCP, POSIX)
- Layer 1
  - Middleware (resource discovery/management, code/data staging, remote execution, security, monitoring)
- Layer 2
  - Programming environment (PSE, in some cases)
- Layer 3
  - Applications



# GRID aware applications

---

- Case A
  - Embarassingly parallel computations (CONDOR like), no heterogeneity, dinamicity
- Case B
  - Esplicitly needed resources (compiler, CPU power, ...), hand made placement, no dinamicity (but “instant” one)
- Case C
  - Data intensive/data driven applications (no dinamicity, heterogeneity)





# GRID aware applications

---

- Case D

- High level parallel code, high level requirement specs, automatic parallelism discovery, automatic adaptation to hardware, restructuring and automatic hardware lowering

**where PPP**



# Contents

---

- GRID: current status
- *A RISC grid core*
- Current experiments
- Conclusions



## **RISC GRID core (RGC)**

---

- Perspective
- Basic functionalities
- Implementation



## RGC: perspective

---

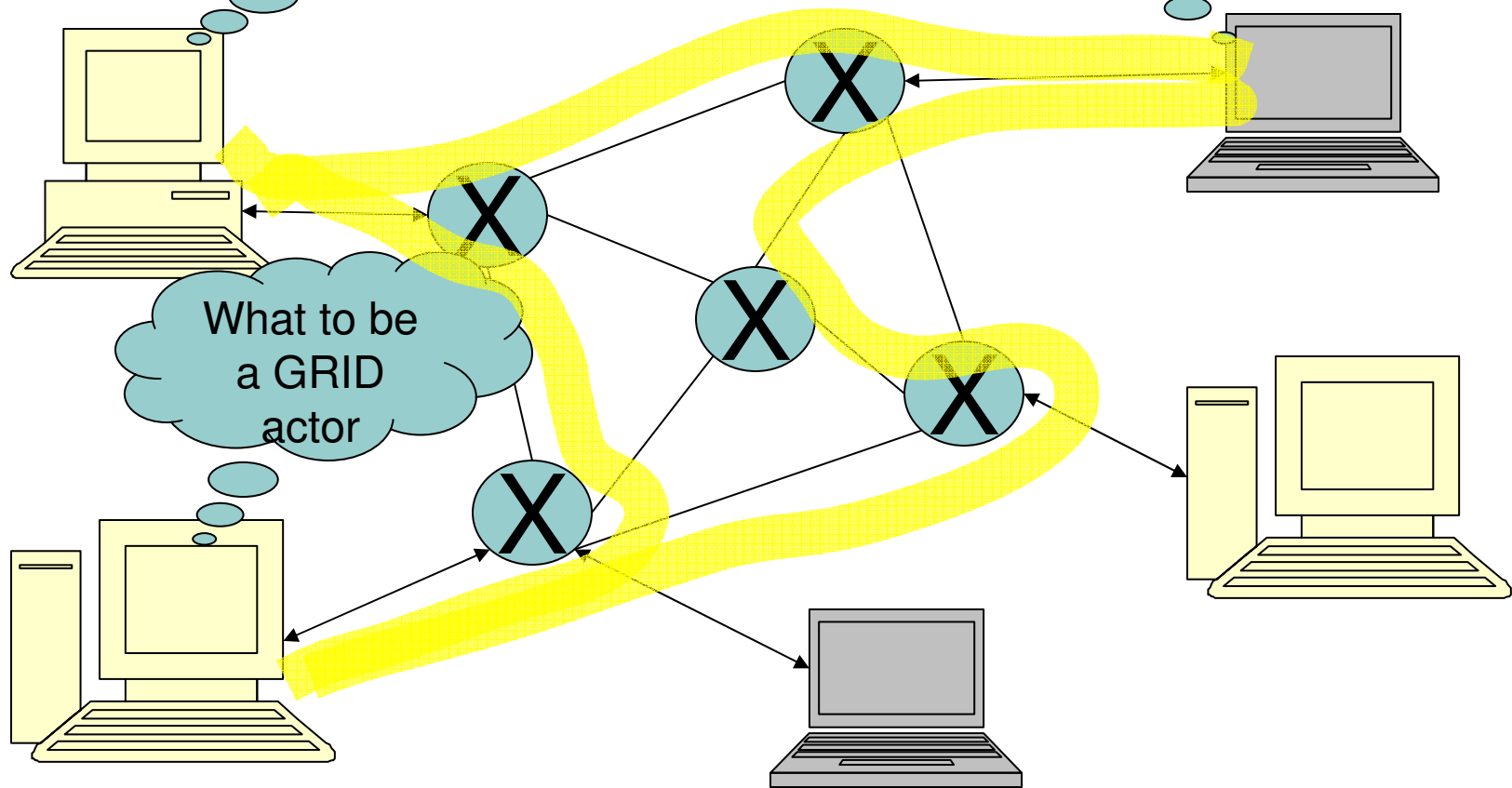
- GRID components actively participate to offer services
- Subject
  - Entity needing services → Entity providing services
- RGC basic service:
  - computing services



**RGC**

Need to compute A

What to be a GRID actor



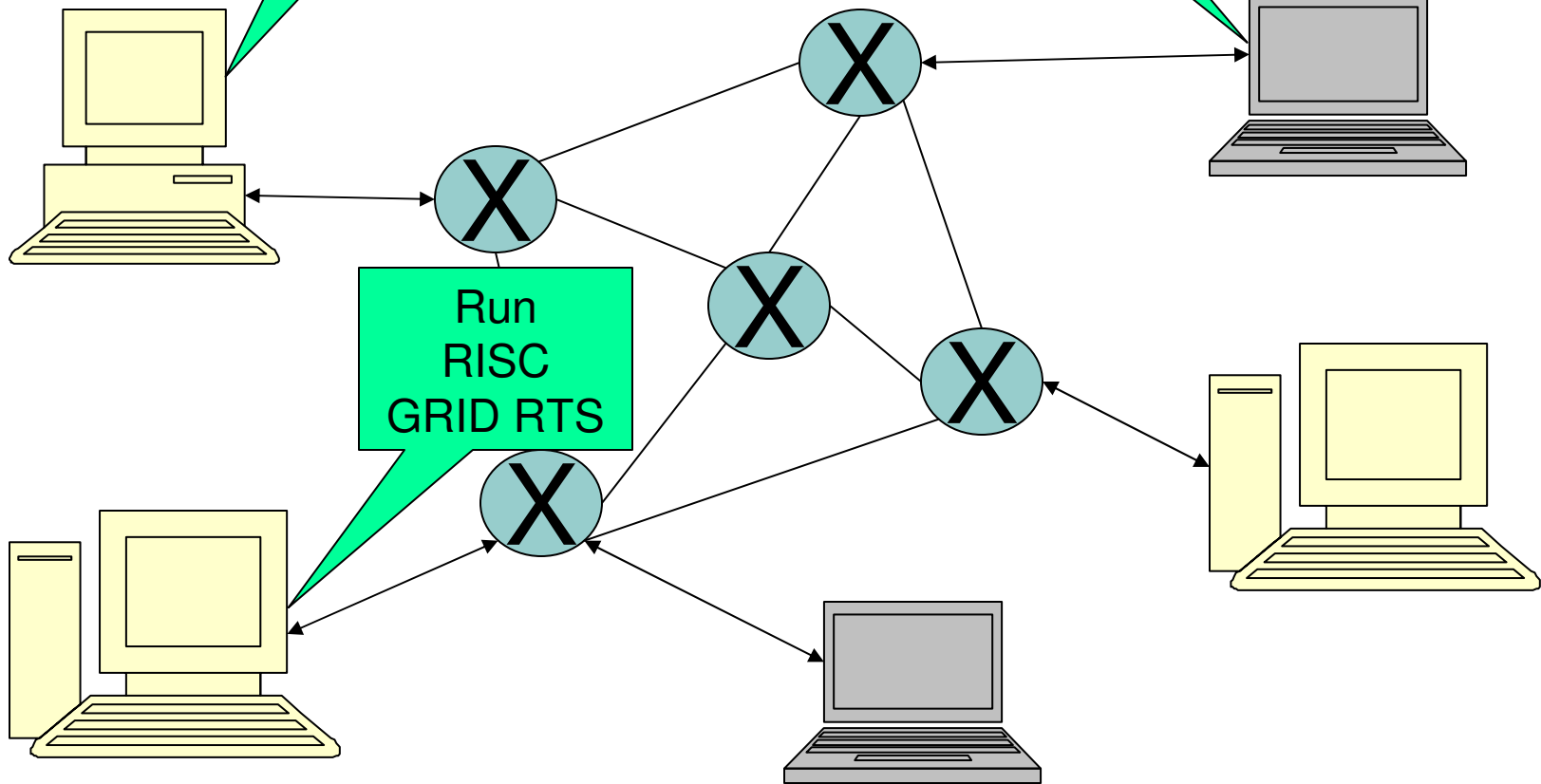


**RGC**

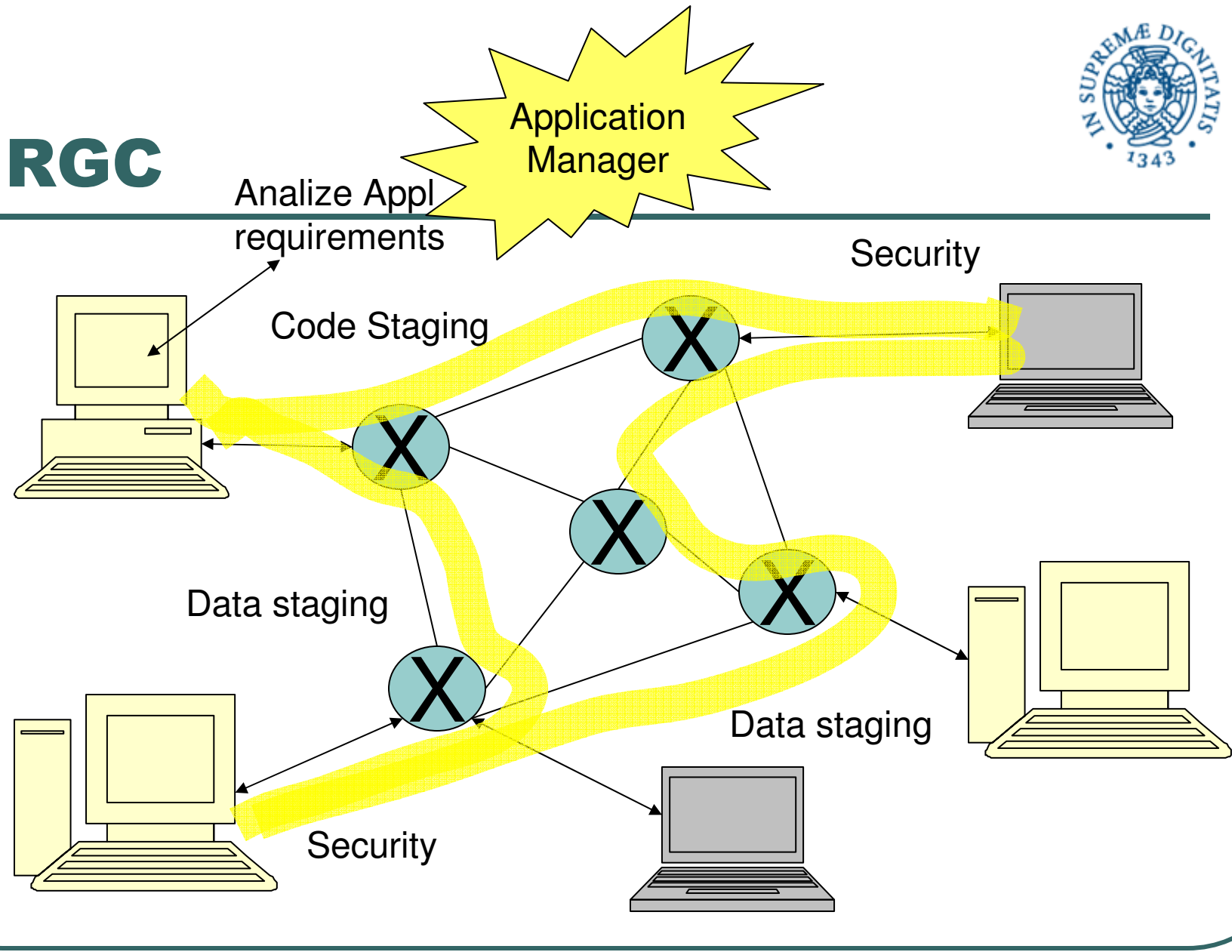
Run RISC  
GRID RTS

Run RISC  
GRID RTS

Run  
RISC  
GRID RTS



# RGC





## **RGC: *basic functionalities***

---

- *Secure access*  
(user classes, certificates, etc.)
- *Code & Data staging* (on demand!)
- *Secure data/code staging*
- *Introspection/reflection/meta info*
- *Remote control*  
(monitor, checkpoint, interrupt&resume)
- *Accounting*
- All services managed from application !





## **RGC bf: *secure access***

---

- Identify class of users
- Allow different sandbox levels
- Guarantee unique id & certificates
- Use session certificates for the single computation



## **RGC bf: *staging***

---

- On demand
- No long range consistency
- Completely application driven
- Caching allowed/enforced upon application directives
  - To enhance appl. performance as well as multiple runs performance



## **RGC bf: *secure staging***

---

- Session certificates
- Data & code
  - on demand, default → encode
- Performance issues
  - Latency/bandwidth vs. cypher/decypher time
  - Critical ? Very often latency much larger!



## **RGC bf: *introspection***

---

- Needed to gather compatibility info
- Needed to set up application deployment
- Open format
  - E.g. XML
- Service
  - To be asked from application manager
- Announce
  - When announcing node availability w.r.t. grid



## **RGC bf: *remote control***

---

- Available to the application manager
  - To monitor code execution
  - To preempt no more useful computations
  - To force different behaviour upon performance changes (loss w.r.t. theoretical model)

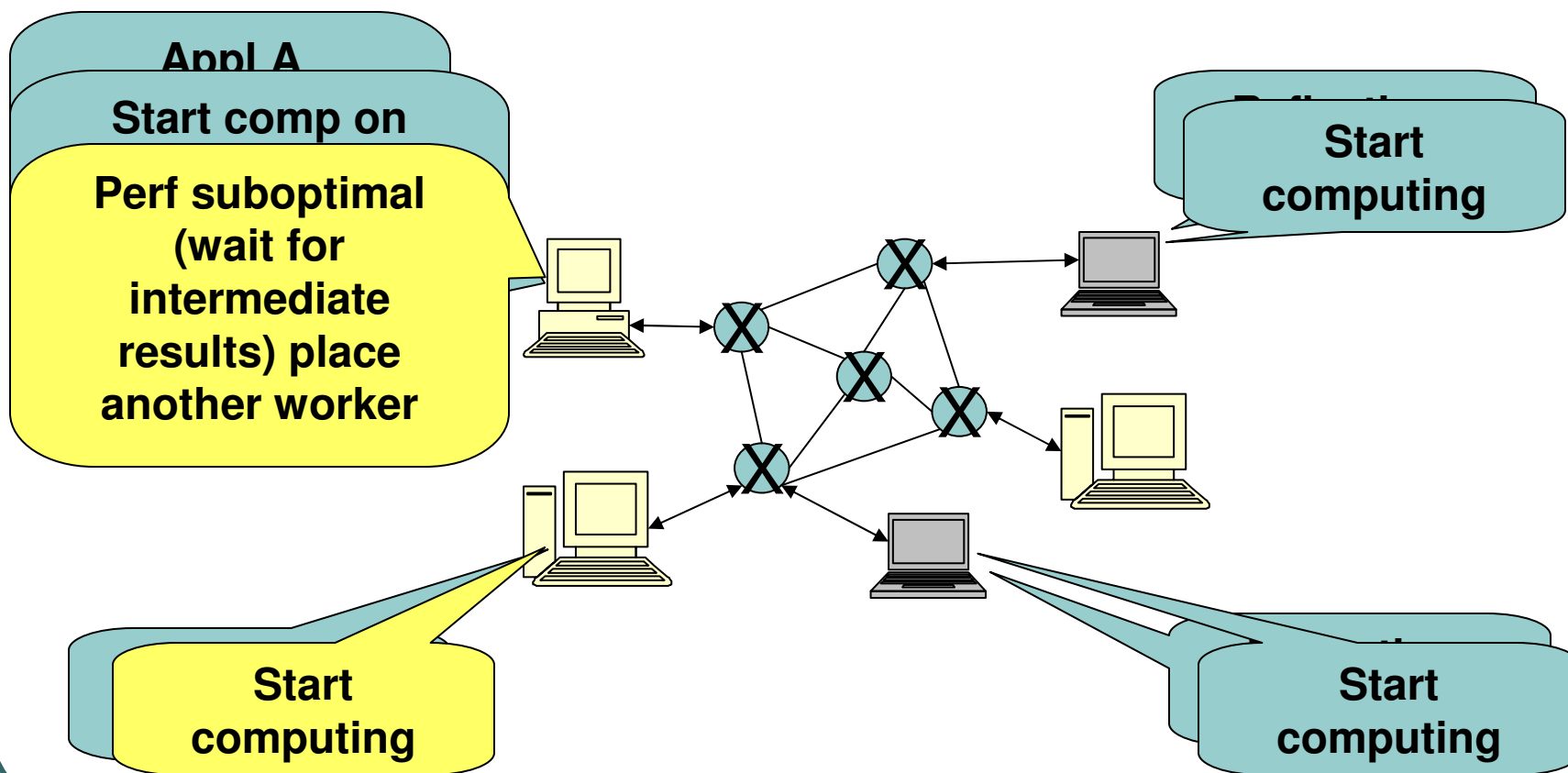


## **RGC bf: *accounting***

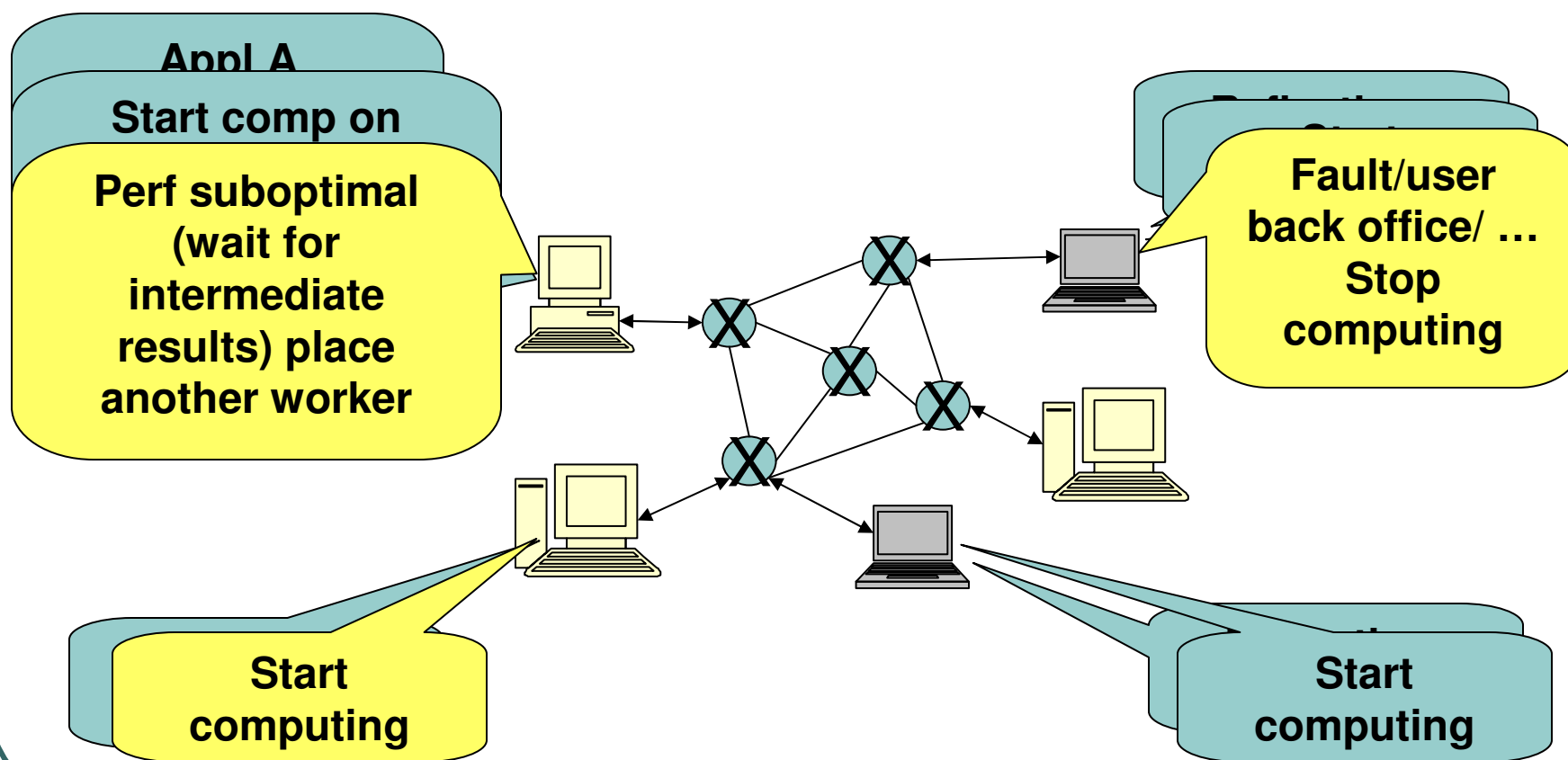
---

- Needed to bill users!
- Account CPU time
  - Different policies per user class
  - Different policies per type (idle time, ...)
- *Must be present to make the approach acceptable*
  - Differently from other points above !

# RGC scenario



# RGC scenario (2)







## RGC implementation

---

- Server process on well known port
  - Peer2peer discovery/announce
  - Implements all basic RGC services
  - Possibly built on top of existing middleware
    - Globus
    - JXTA
    - ...
  - Explicitly run by machine owner !



## Why RISC ?

---

- Compare Globus services with our ones!
- Compare the amount of code in the p2p server placed on the machines!
- Compare the usage of services (application driven, even concerning server code staging)!



# GT3 online man ...

Overview (GT3 OGSA 3.0.2 API) - Microsoft Internet Explorer

Address <http://www-unix.globus.org/toolkit/3.0/ogsa/impl/java/build/javadocs/>

Overview Package Class Use Tree Deprecated Index Help

PREV NEXT FRAMES NO FRAMES

## GT3 OGSA 3.0.2

### Packages

<a href="#">org.apache.axis.deployment.wsdd</a>	
<a href="#">org.globus.ogsa</a>	
<a href="#">org.globus.ogsa.client</a>	
<a href="#">org.globus.ogsa.client.managers</a>	
<a href="#">org.globus.ogsa.client.reflection</a>	
<a href="#">org.globus.ogsa.config</a>	
<a href="#">org.globus.ogsa.core.admin</a>	
<a href="#">org.globus.ogsa.core.admin.bindings</a>	
<a href="#">org.globus.ogsa.core.admin.service</a>	
<a href="#">org.globus.ogsa.core.logging</a>	
<a href="#">org.globus.ogsa.core.logging.bindings</a>	

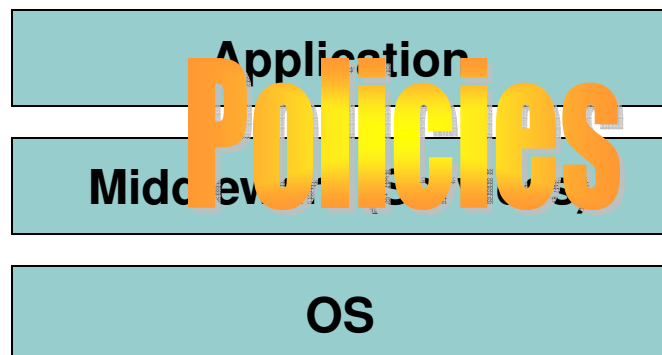


# ... vs. Lithium man !

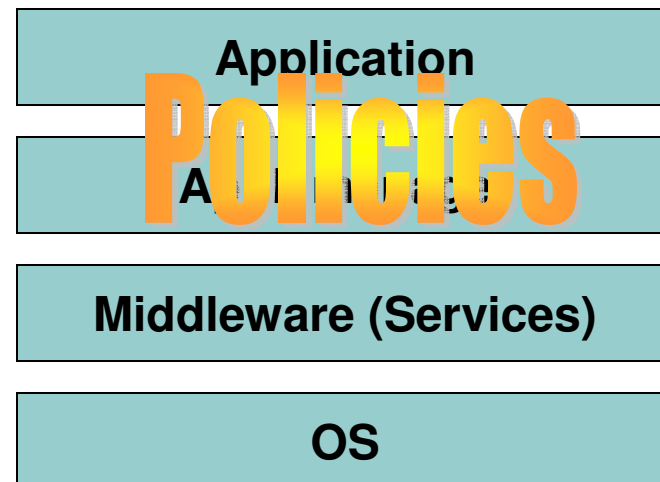


## Why RISC (2)

- Control/policies moved up to application manager



*currently*



*RGC*



# Contents

---

- GRID: current status
- A RISC grid core
- *Current experiments*
- Conclusions



## Framework



- FIRB GRID.it
  - Three year Italian National Project
  - Basic research on grids
  - WP8 “advanced, structured, component based parallel programming model for grids”
- In the meanwhile
  - Several prototypes already available
    - P3L, Lithium, ASSIST, ...

GRID.IT



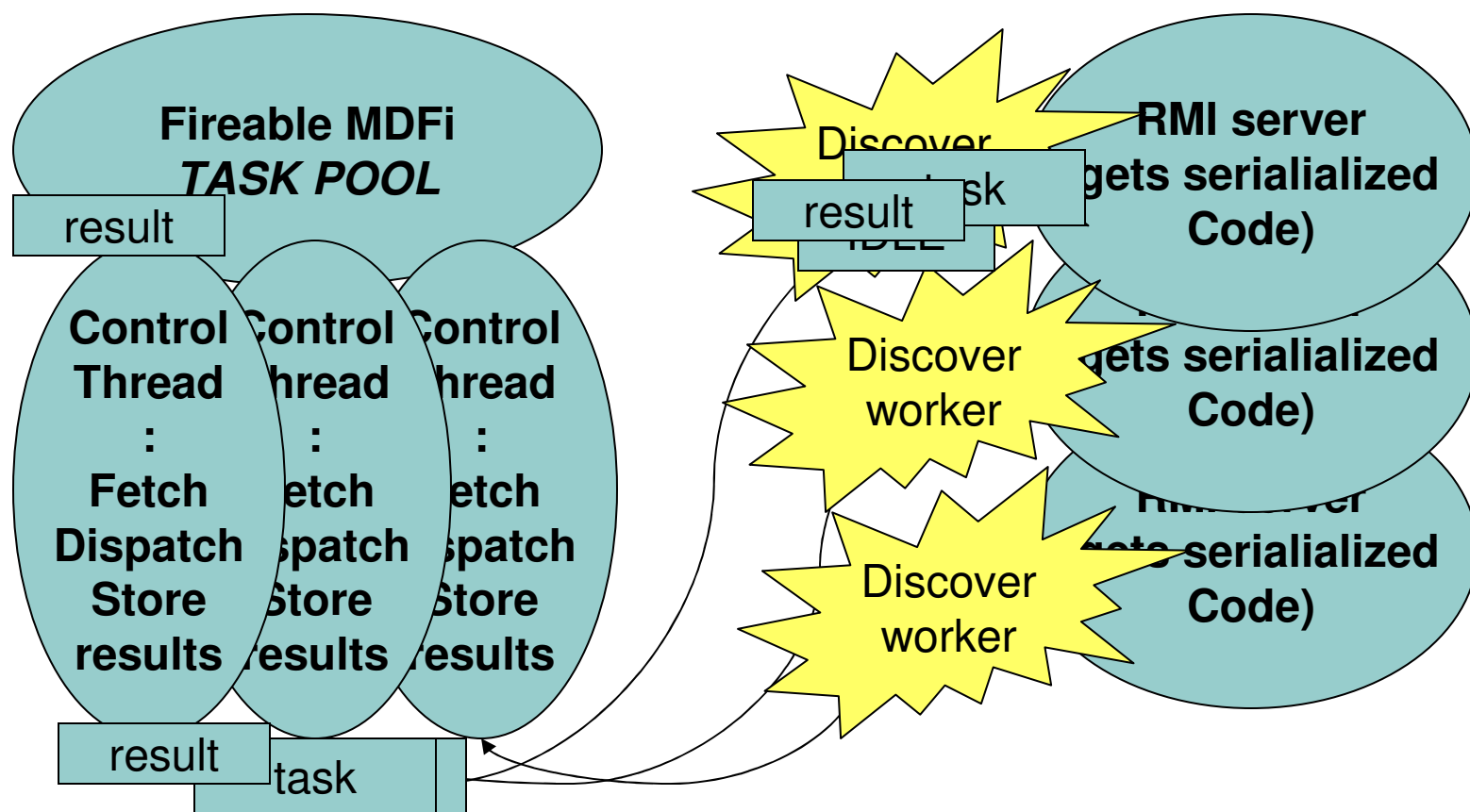
# Lithium

---

- Java based, parallel, structured programming environment
  - RMI
  - Algorithmical skeletons (including farms, pipelines & divide&conquer)
  - Full Java library
  - On since 2001 (experimental)



# Lithium structure





## Lithium GRID evolution

---

- Discovery → *JINI, JXTA, ...*
- Code/Data staging → *Java serialization*
- Security → *SSL*
- Reflection → *Beans*
- Remote Control → *RMI*
- Accounting → *???*



## Application manager sample

---

- Embarassingly parallel computations
- Data at the application site
- Results at the application site
- *TASK FARM dynamic* template
  - *Compute tasks ASAP*
  - *On the available resources*
  - This is a CONDOR-like, Case A ***application schema!***



## Usually ...

---

- Task list (pool)
- Code computing f
- List of possible workers (static, dynamic required to the GRID resource manager, available once and for all)
- ... go
- & wait for completion



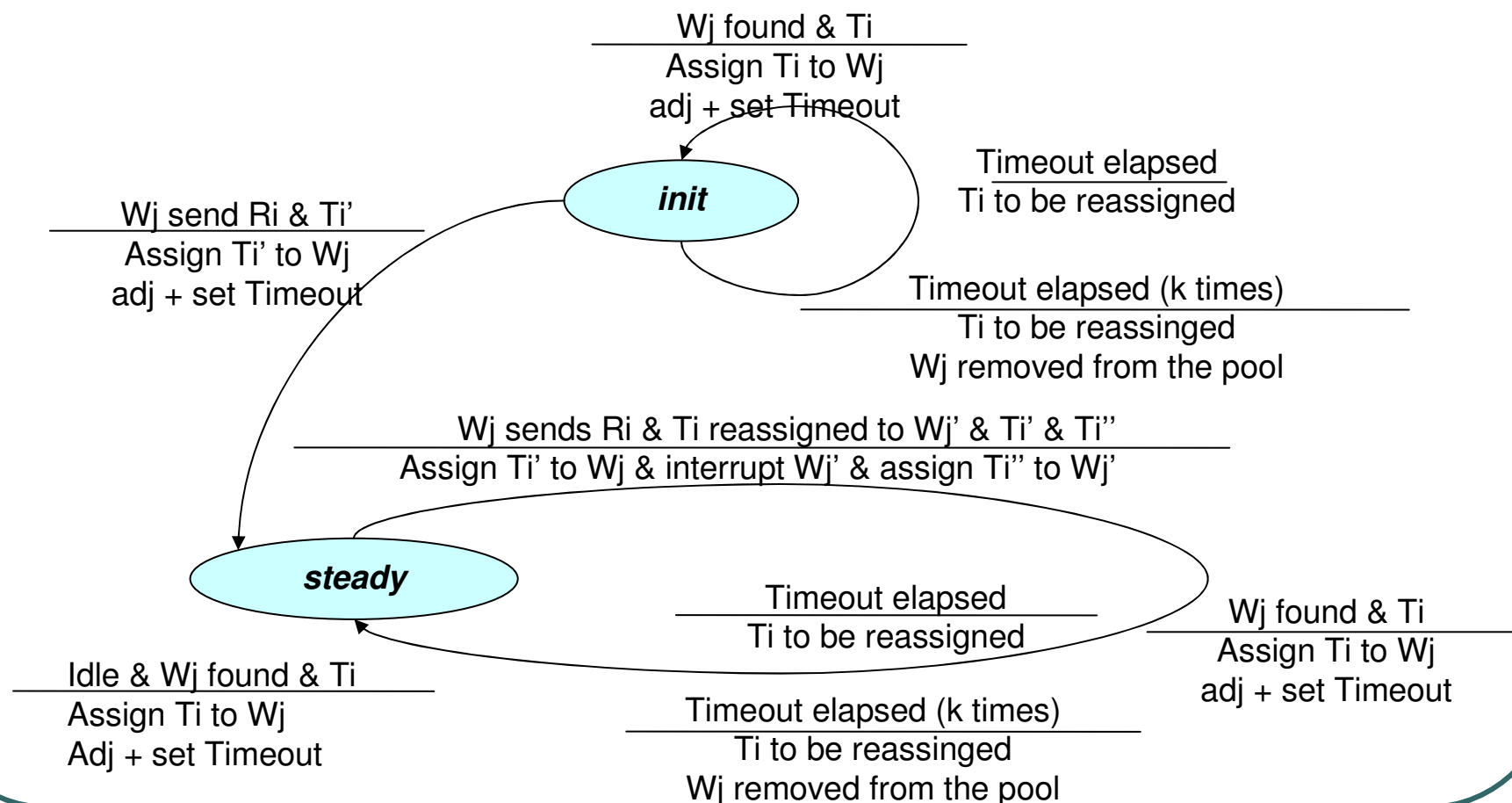
# RISC approach

---

- Task farm manager
  - Looks for workers (p2p)
  - Assigns tasks for execution
  - When results come back
    - Stops looking for new workers (steady state)
    - Assigns tasks for execution to idle workers
- Workers
  - Just receive work to be computed  
(RGC: computation services only at the GRID peers!)

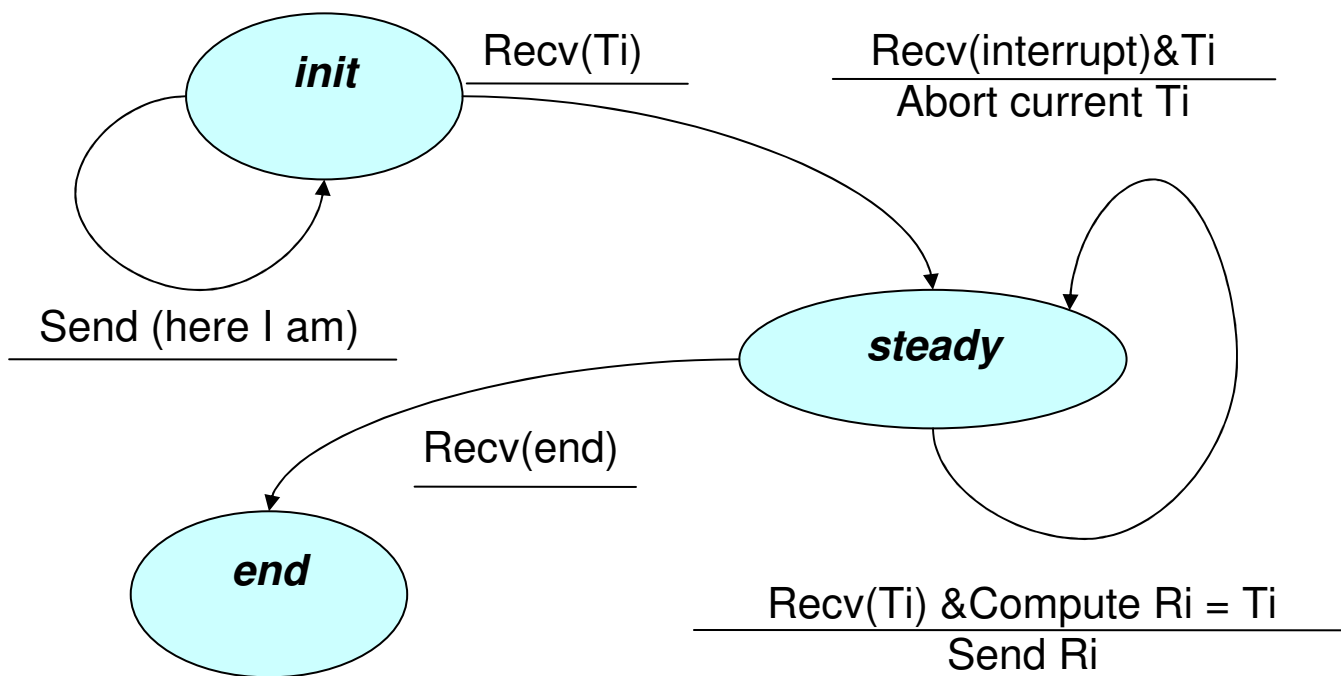


# Task manager (FSA)





# Worker (FSA)





## therefore ...

---

- **Adaptivity**
  - To faulty workers
  - To dynamic changes in network performance
- **Fault tolerance**
  - Faulty workers
- **Heterogeneity**
  - Java ...





# Experiments

---

- Modified version of the Lithium prototype
  - RMI, Serialization, ...
- Simulator
  - Java program
  - Exact knowledge
  - ... compared to measured

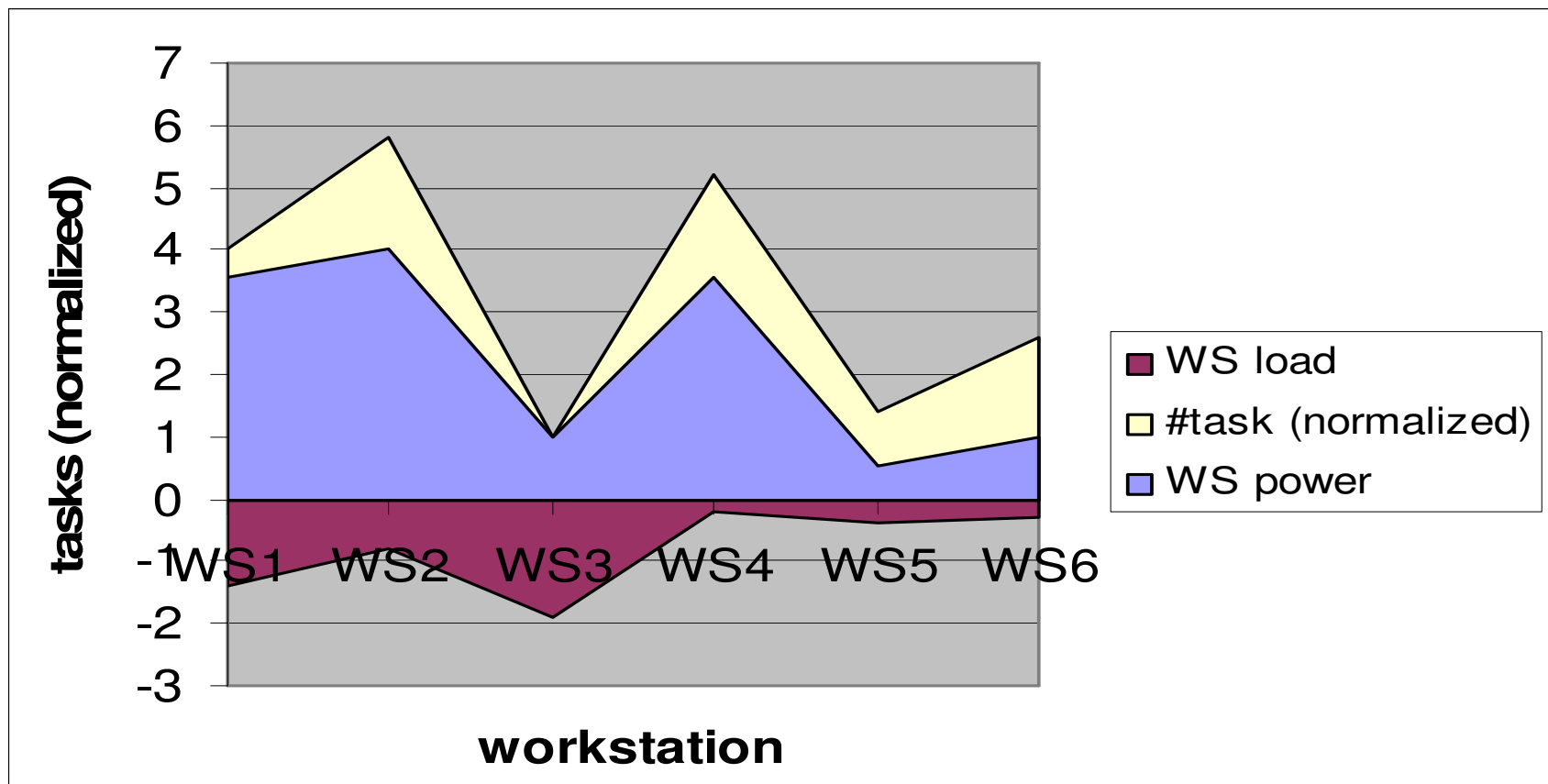
# Results: *load balancing*

---



- Execute
  - a stream of independent task
    - Average execution time not known
    - Distribution not known
  - On a set of *production* workstations
    - Different Hw
    - Different load

# Results: load balancing

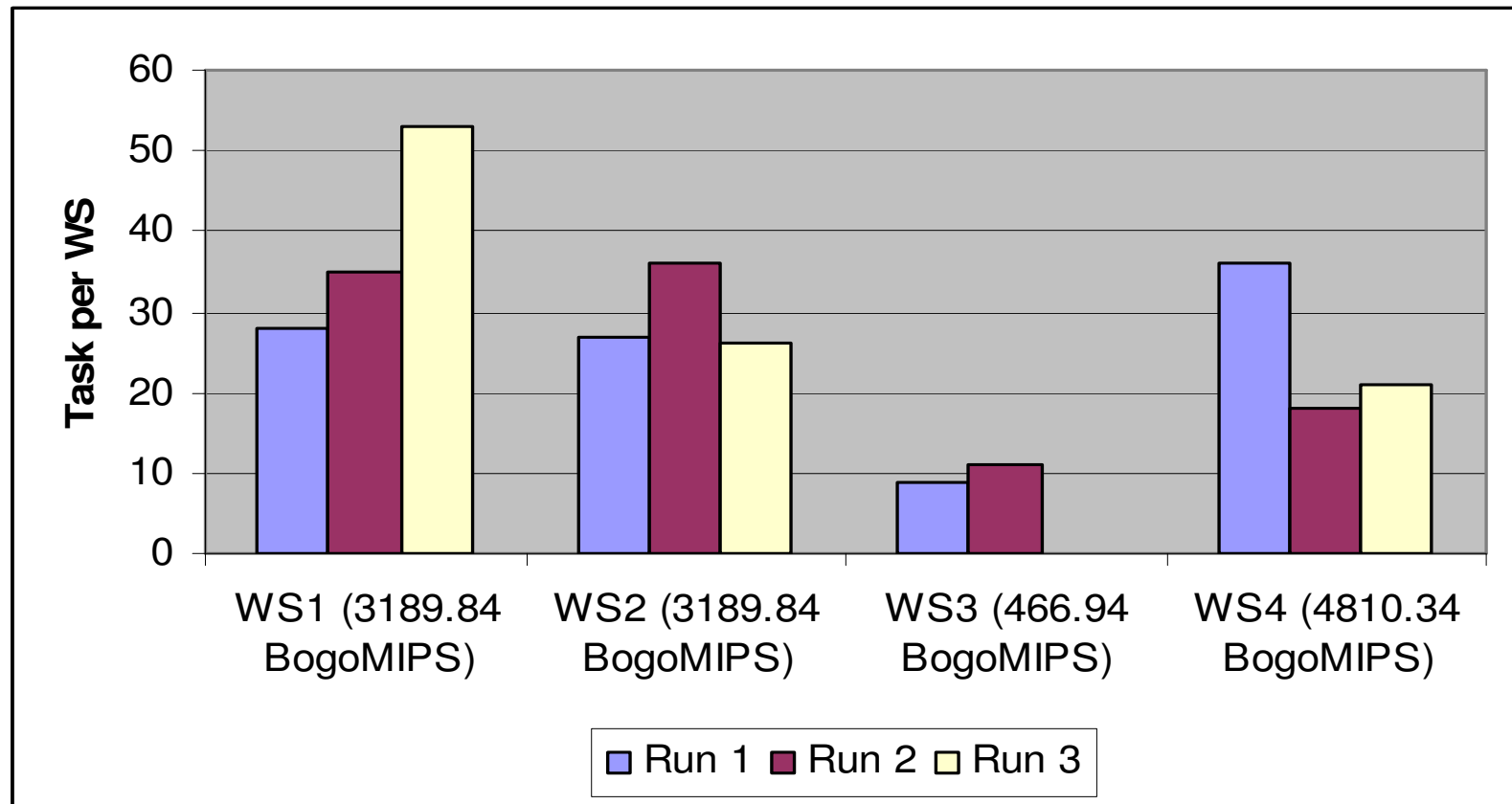


## Results: *discovery*



- Run a stream of tasks
- On a set of production workstations
- First run
  - All discovered since T0
- Second run
  - WS4 discovered after  $\frac{1}{2}$  tasks
- Third run
  - WS3 never discovered
  - WS2 discovered after  $\frac{1}{3}$  tasks
  - WS4 discovered after  $\frac{1}{2}$  tasks

# Results: *discovery*



## Results: *fault tolerance*

---

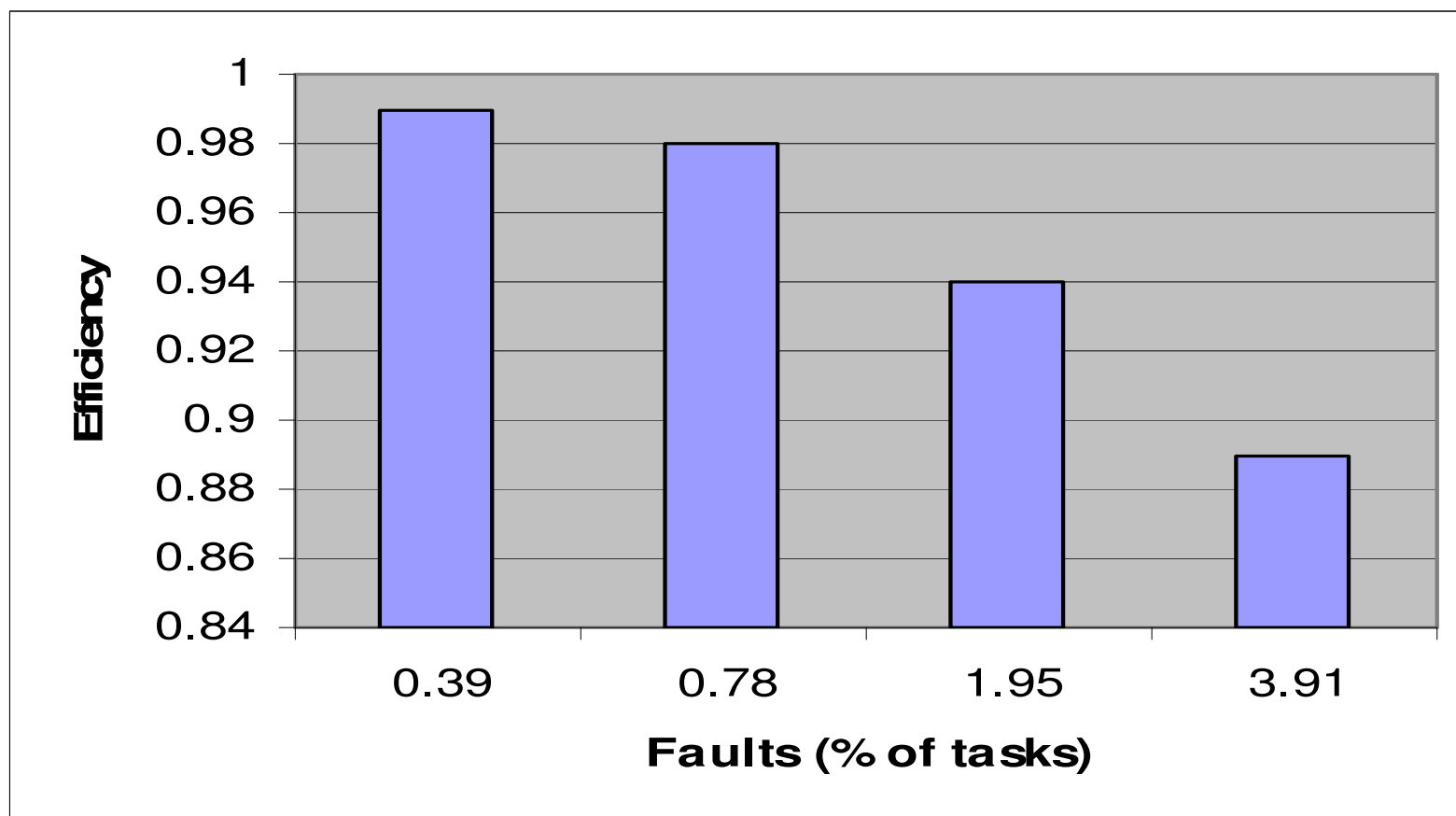


- Workers assumed to fail
  - Due to internal problems
  - Due to network problems/delays
- Faults % to the total number of tasks executed
- Efficiency measured



## Results: *fault tolerance*

**Simulation**



## Results: *heterogeneity*



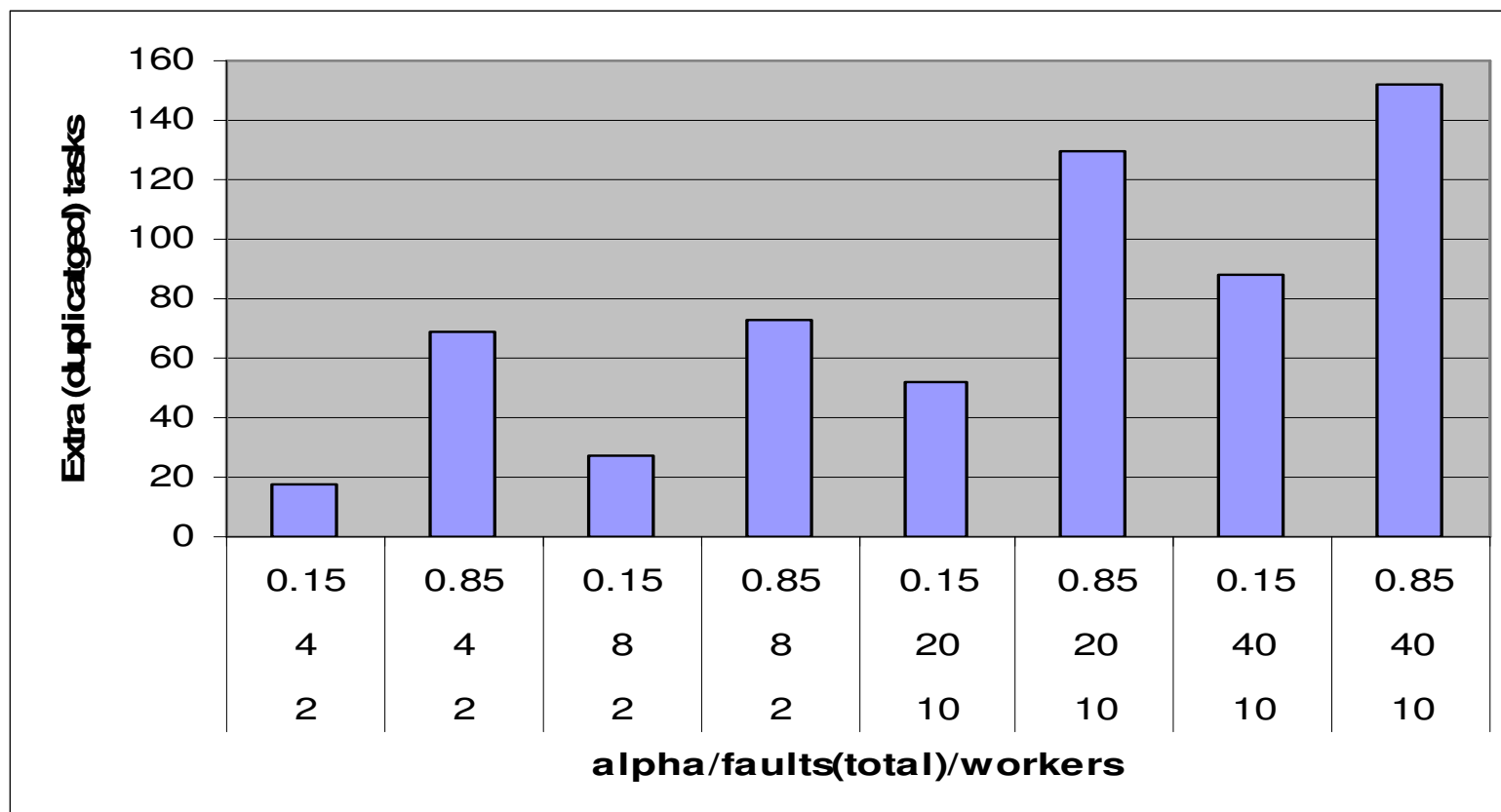
- A stream of 1024 tasks executed
- With 2 or 10 workers
- 2 or 4 faults per worker (at random time)
- $\alpha = 0.15$  or  $\alpha = 0.85$   
RTT =  $\alpha$  RTT + (1- $\alpha$ ) RTT<sub>current</sub>





# Results: *heterogeneity*

**Simulation**





# Contents

---

- GRID: current status
- A RISC grid core
- Current experiments
- **Conclusions**



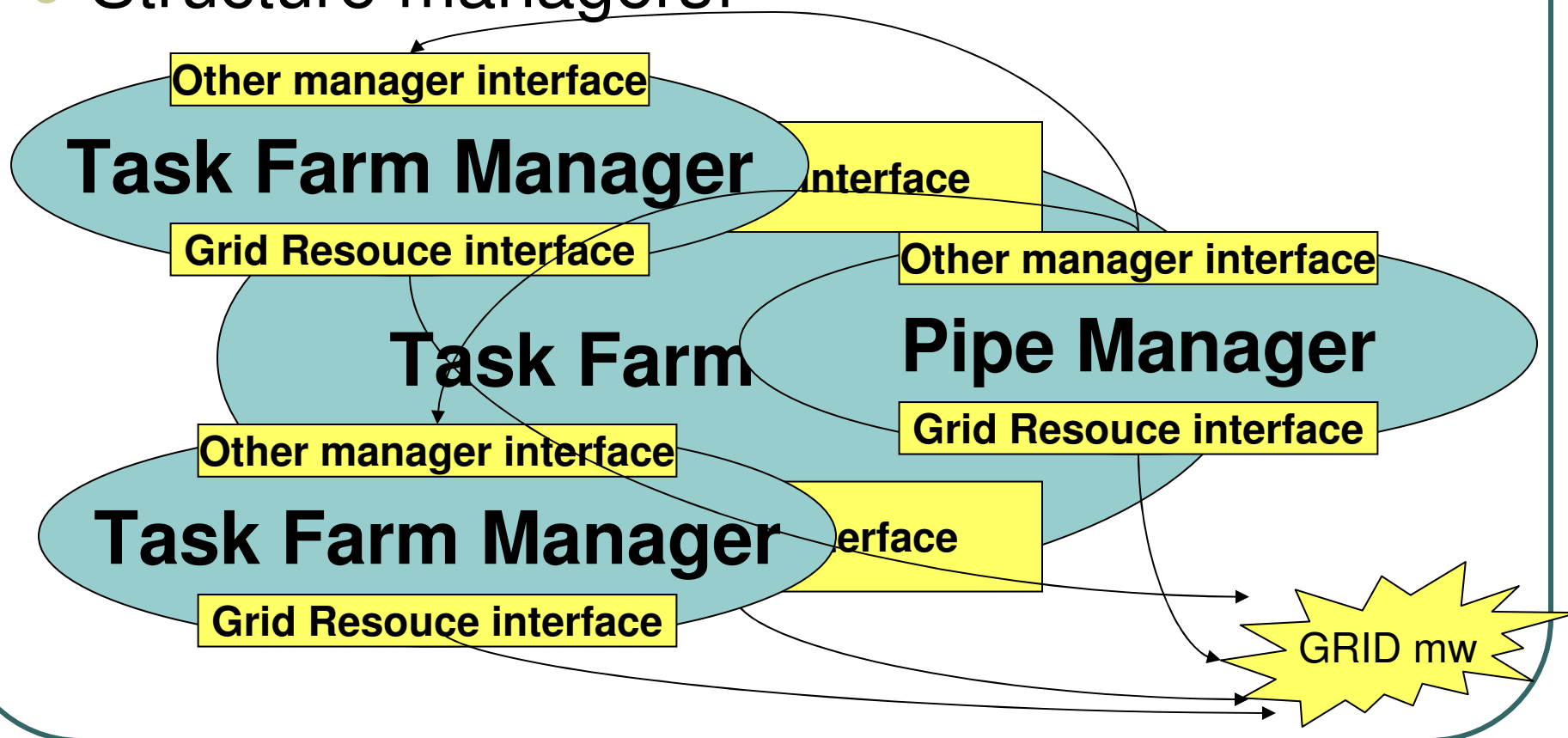
# Evolution

---

- Task farm manager  
→ application schema manager
- Other managers (e.g. pipeline)  
→ more complex application schemas
- Become an *environment*
  - Suitable to handle common grid aware applications

## Evolution (2)

- Structure managers!



## Evolution (3)

---

GRID.IT

- Programming environment
  - Target GRIDs
  - GRID aware
  - Structured
  - Component based
  - With application managers
  - Sitting on top of *common sense* middleware (say GT3, Java/Jini/JXTA, ...)



## Related work

---

- So much activity on GRID ... !
- Want to cite CONDOR
  - Limited application schema
  - Very efficient RTS
  - Included in recent toolkits
- **Whereas RGC:**
  - Unlimited application schemas
  - (hopefully) efficient RTS
  - ???



## Conclusions

---

- New methodology for GRID tool development proposed
- Join skeleton/design pattern results with the GRID world
- Preliminary results
- Currently evolving
- ...

[marcod@di.unipi.it](mailto:marcod@di.unipi.it)

[www.di.unipi.it/~marcod](http://www.di.unipi.it/~marcod)

---



***A RISC approach to the GRID***

**any questions ?**