# RADIC II

A Fault Tolerant Architecture with
Flexible Dynamic Redundancy

By Guna Santos

# Key Points

- The increasing number of applications demanding high performance and availability (*performability*)

- Fault tolerance plays a major role in order to keep high availability

- A fault tolerant solution must generates minimal overhead in its activities, in order to not affect the application performance

- After a fault ocurrence, some fault tolerant system generate system degradation, affecting the performance.

caos

# Challenges

- To provide a fault tolerant solution able to mitigate or to avoid the system degradation
  - Restoring the system configuration
  - Avoiding changes in number of active nodes
  - Allowing perform maintenance tasks, preventing faults
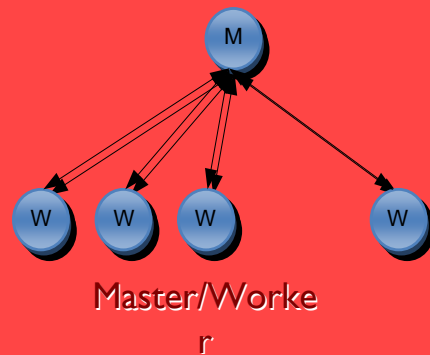
caos

# Contents

- Goals

- Fault Tolerance and the RADIC Architecture

- Recovery Side-Effect

- Protecting the System

- Implementation and Experimental Results

- Conclusions
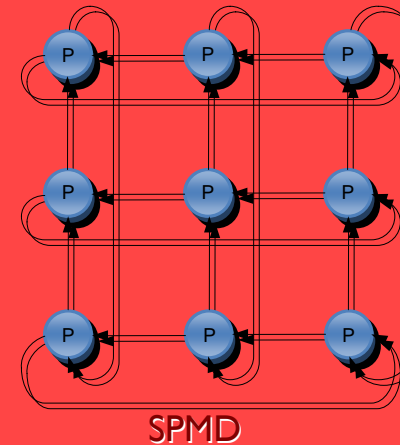
- Future Work and Open Lines

caos

# Goals

- To incorporate the ability to control the system degradation to a fault tolerant architecture with following characteristics:
  - Flexible, Decentralized, Transparent, Scalable
- To extend the functionality of this architecture implementation

caos

# A Fault Tolerance Approach

- Fault tolerant architecture for message-passi...

- It based on Rollba... protocol

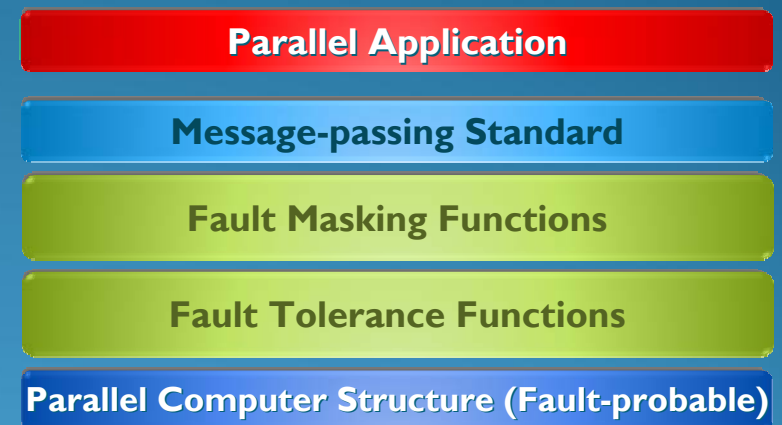- It creates a layer application from t... failures

Different Paradigms

M

W   W   W          W

Master/Worker

P   P   P

P   P   P

P   P   P

SPMD

P

caos

# A Fault Tolerance Approach

- Fault tolerant architecture for message-passing systems

- It based on Rollback-recovery protocol

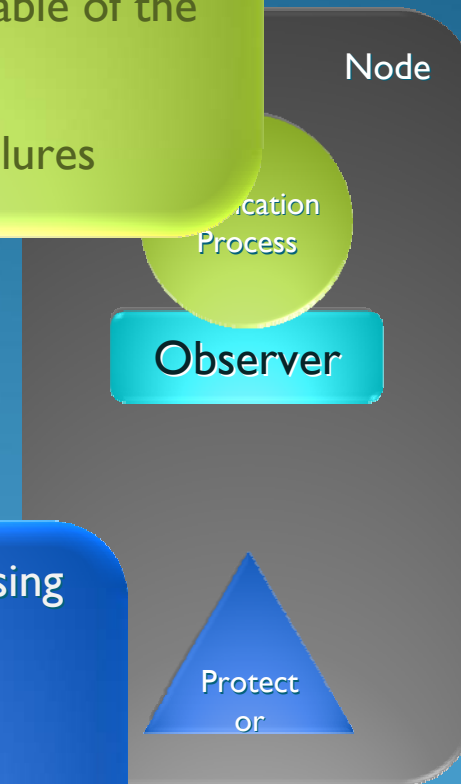- It creates a layer isolating the application from the cluster failures

| Parallel Application |
| --- |
| Message-passing Standard |
| Fault Masking Functions |
| Fault Tolerance Functions |
| Parallel Computer Structure (Fault-probable) |

caos

# The RADIC Architecture

- RADIC is a ... architecture ... passing systems
  - Two kind of process perform the fault tolerance activities
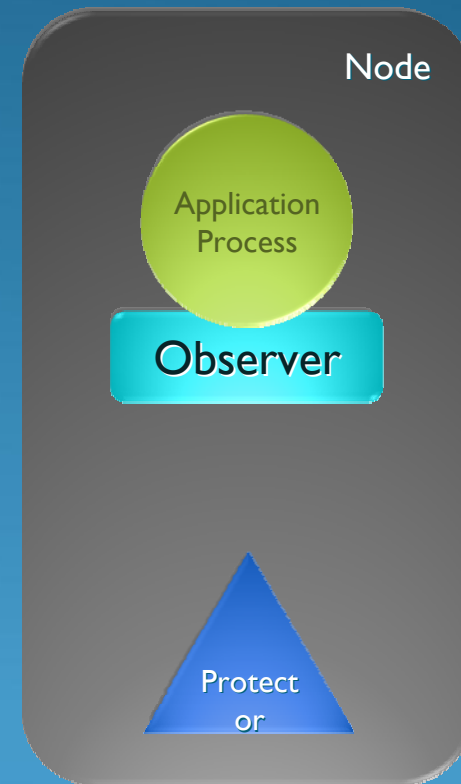
**Green callout:**
- It manages the communications between processes
- It takes checkpoints and event logs and sends to storage
- It masks the faults using a table of the processes.
- It detect communication failures

**Blue callout:**
- It performs fault detection tasks using a heartbeat/watchdog scheme
- It stores the redundancy data (checkpoints and logs)
- It starts the recovery process

Node

...cation Process

Observer

Protect or

caos

# The RADIC Architecture

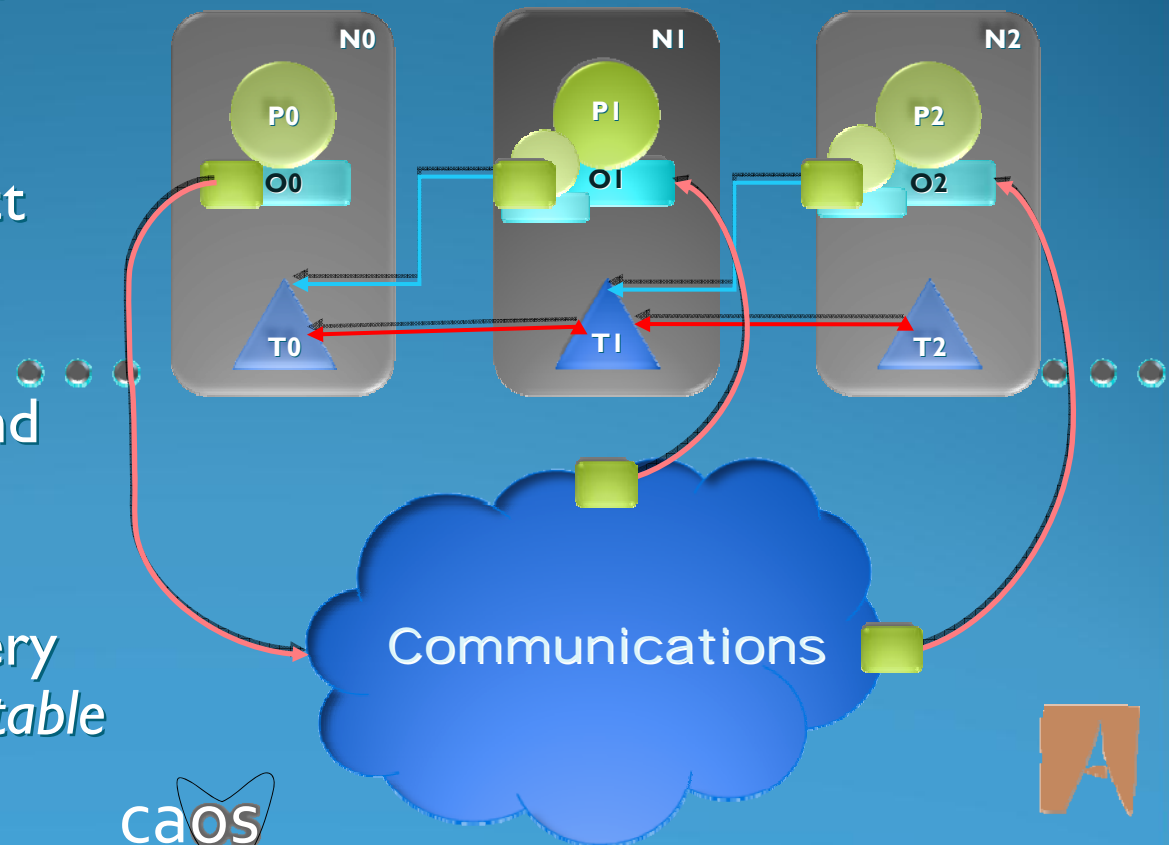- RADIC in Practice



Node

Application Process

Observer

Protect or
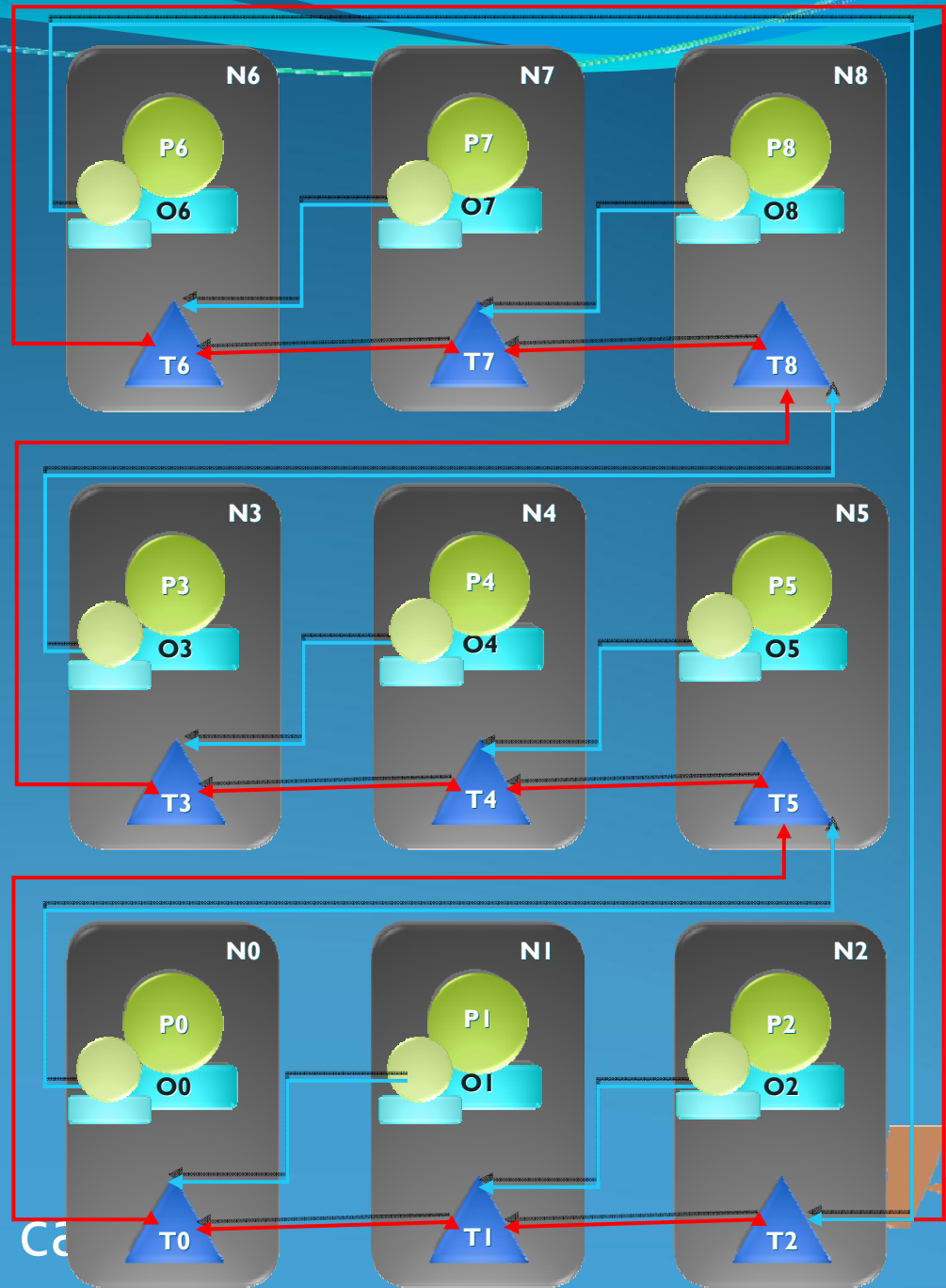
caos

# The RADIC Architecture

- RADIC functioning
  - Protectors establish the heartbeat/ watchdog
  - Observers connect with Protectors
    - It sends checkpoints and event logs
    - It manages the message delivery using the *radictable*

| RADICTABLE | | |
|---|---|---|
| PID | Address | Protector |
| 0 | Node 0 | Node 8 |
| 1 | Node 1 | Node 0 |
| … | … | … |

**N0** P0 O0 T0

**N1** P1 O1 T1

**N2** P2 O2 T2

**Communications**
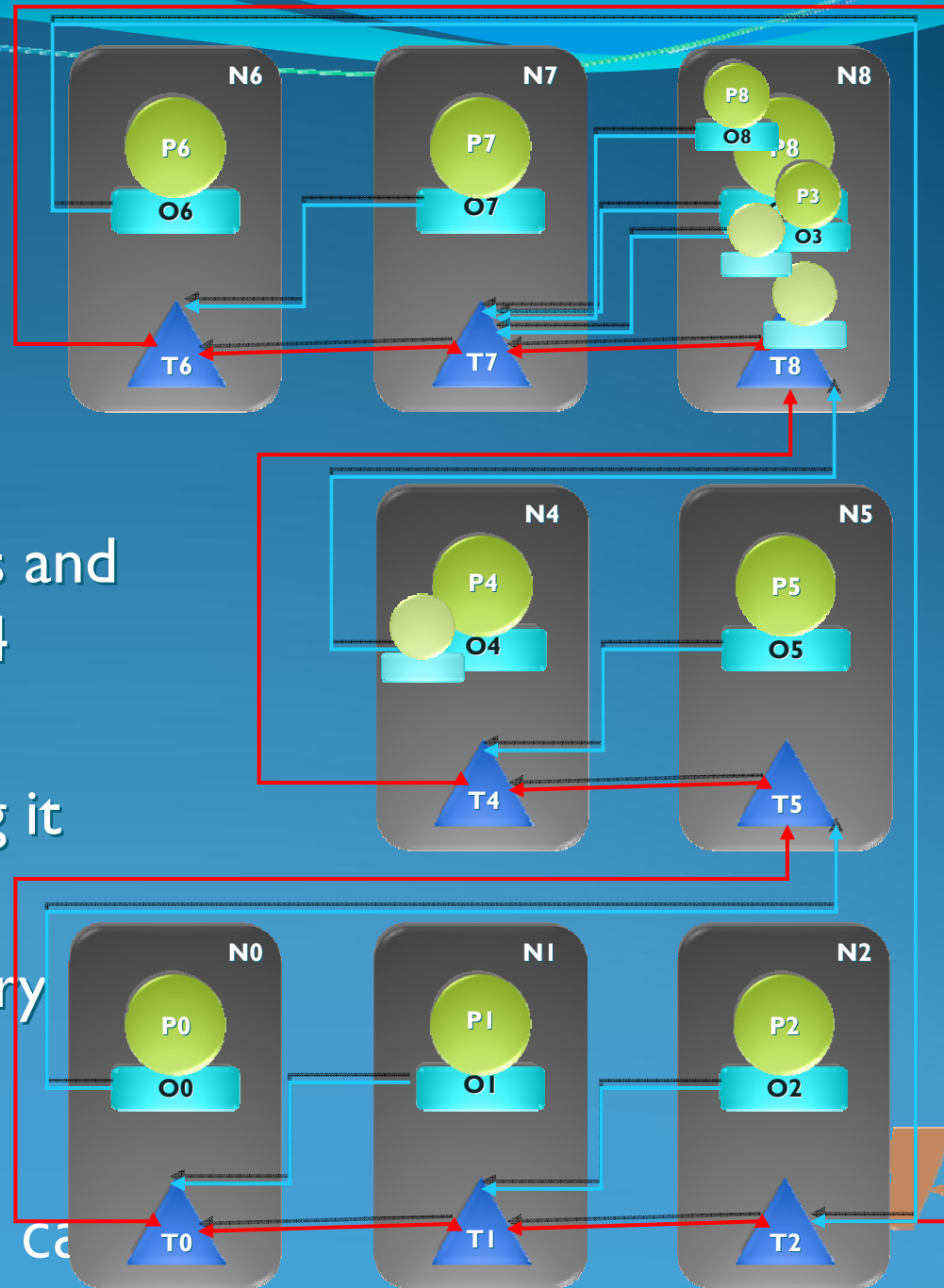
caos

# The RADIC Architecture
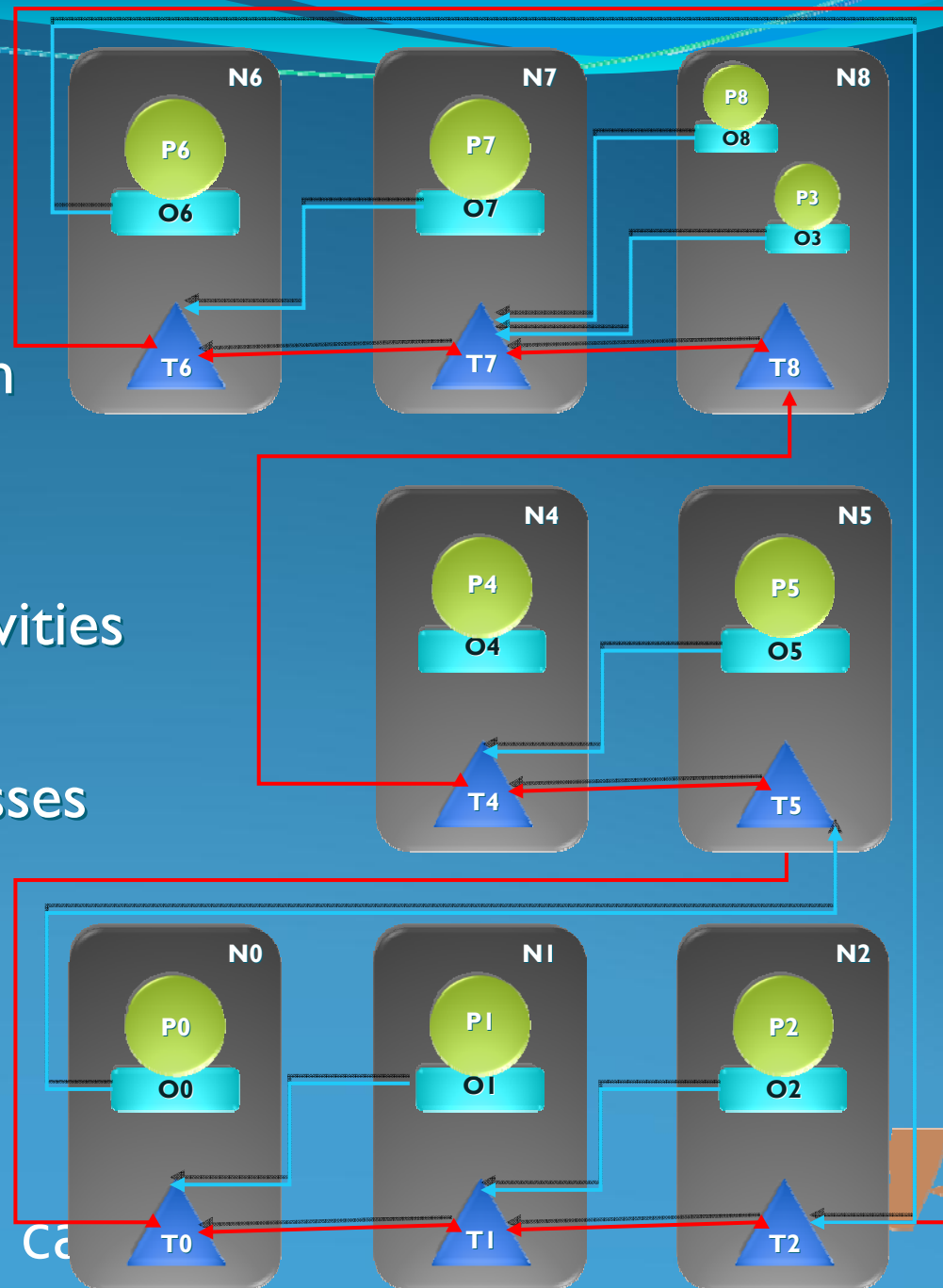
- The recovery process
  - N3 fails



11

# The RADIC Architecture

- The recovery process
  - N3 fails
  - T8 detects the faults and waits the T4 and O4 connection
  - T8 launches P3 using it checkpoint
  - O3 start the recovery process of P3

N6 P6 O6 T6

N7 P7 O7 T7

N8 P8 O8 P8 P3 O3 T8

N4 P4 O4 T4

N5 P5 O5 T5

N0 P0 O0 T0

N1 P1 O1 T1

N2 P2 O2 T2

# Recovery Side-Effect

- The sytem configuration was changed
  - Processes P3 and P8 slowdown their activities
  - Protector T7 now protects two processes
  - The memory usage in N8 rises

N6 | N7 | N8
P6 | P7 | P8
O6 | O7 | O8
| | P3
| | O3
T6 | T7 | T8

N4 | N5
P4 | P5
O4 | O5
T4 | T5

N0 | N1 | N2
P0 | P1 | P2
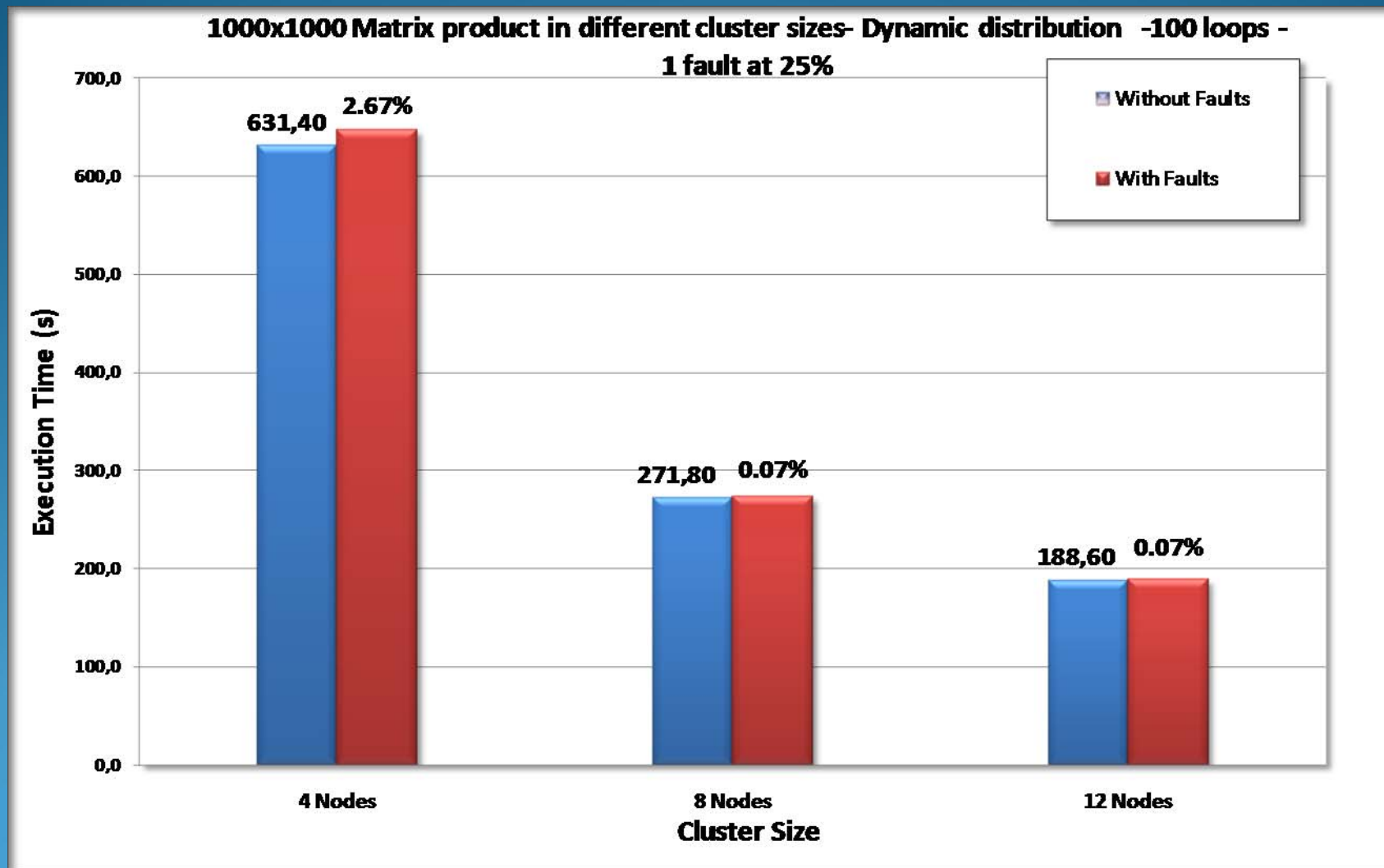O0 | O1 | O2
T0 | T1 | T2

13

# Recovery Side-Effect

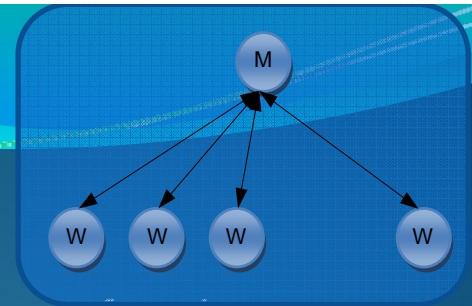- In order to quantify these side-effects facing different scenarios, we performed some experiments:
  - We applied three variations of a matrix product program:
    - Dynamic M/W – Loosely coupled
    - SPMD – High coupling in all processes
    - Static M/W – High dependency

caos

# Recovery Side-Effect

1000x1000 Matrix product in different cluster sizes- Dynamic distribution  -100 loops -
1 fault at 25%

# Recovery Side-Effect

SPMD – one fault in different moments

## 1500x1500 Matrix product using 9 nodes - SPMD cannon algorithm- 1 fault
### 160 loops -ckpt each 60s

Legend:
- Without failures
- Failures without spare

| Fault moment | Without failures | Failures without spare |
|---|---|---|
| 25% | 434,57 | 72,61% |
| 50% | 434,57 | 49,37% |
| 75% | 434,57 | 27,05% |

Time (s) axis: 0,0 to 800,0

# Recovery Side-Effect

## Static M/W



1000x1000 Matrix product in different cluster sizes– Static distribution – 1 fault at 25% — 100 loops - ckpt each 45s

# Recovery Side-Effect

## Static M/W

**1000x1000 Matrix product with 8 nodes- Static distribution - 1 fault at 25%**
**100 loops - ckpt each 45s**

Legend:
- Without faults
- With Faults

Y-axis: Execution Time (s) — 0, 100, 200, 300, 400, 500, 600

X-axis: Processes — P0, P1, P2, P3, P4, P5, P6, P7

**1000x1000 Matrix product in different cluster siz**
**100 loops – ckpt ea**

Legend:
- Failures without spare

Y-axis: Time (s) — 0,0, 200,0, 400,0, 600,0, 800,0, 1000,0, 1200,0, 1400,0, 1600,0

- 4 Nodes: 707,30 / 93.65%
- 8 Nodes: 324,30 / 73.38%
- 11 Nodes: 227,60 / 73.07%

X-axis: Cluster Size — 4 Nodes, 8 Nodes, 11 Nodes

# Protecting the System

- We incorporate a new protection level in RADIC using a dynamic redundancy functionality, allowing to mitigate or to avoid the recovery side-effects
  - Restoring the system configuration
  - Avoiding system configuration changes
  - Allowing to prevent faults
    - It incorporates a transparent management of **spare nodes**
    - It does not mantain any central information about the spare nodes
    - The spare nodes use does not affect the scalability

caos

# Protecting the System

- Restoring the configuration

- Insert a replacement node running a protector in spare mode
- The Protector announces itself by message forwarding
- T8 requests the new node because it is overloaded
- The new node incorporates the protecting scheme
- T8 tells to P3 to take its checkpoint in N9
- T8 commands T9 to recover P3
- The old P3 suicides
- From now, the message forwarding informs the new spare state

# Protecting the System

- Avoiding changes in number of active nodes

- Insert a new node running a Protector in spare mode (waiting a request)
- The spare Protector announces itself by message forwarding
- Each Protector adds the new spare in its *spare table* and forward it
- This procedure when a protect receives a repe information



**N6** P6 O6 T6
**N7** P7 O7 T7
**N8** P8 O8 T8
**N9** T9
**N3** P3 O3 T3
**N4** P4 O4 T4
**N5** P5 O5 T5
**N0** P0 O0 T0
**N1** P1 O1 T1
**N2** P2 O2 T2

### Spare Table

| SID | Address | #Observers |
|-----|---------|------------|
| 0   | Node 0  | Node 8     |
| 1   | Node 1  | Node 0     |
| ... | ...     | ...        |

# Protecting the System

- Recovering with spare node

- T8 detects a fault and connects to spare, activating it
- It queries about its state (if still is a spare)
- It commands the spare to join to the protecting scheme
- O4 also connects the spare
- T8 sends the check-point and log to spare
- T8 commands to recover P3

# Protecting the System

- Restoring the configuration

# Protecting the System

- Preventing Faults

- N3 is a fault probable node
- Insert a spare node
- Oportuniscally inject a fault in N3

N6    P6   O6   T6
N7    P7   O7   T7
N8    P8   O8   T8

N9   T9

N3    P3   O3   T3
N4    P4   O4   T4
N5    P5   O5   T5

N0    P0   O0   T0
N1    P1   O1   T1
N2    P2   O2   T2

# Protecting the System

- Preventing Faults

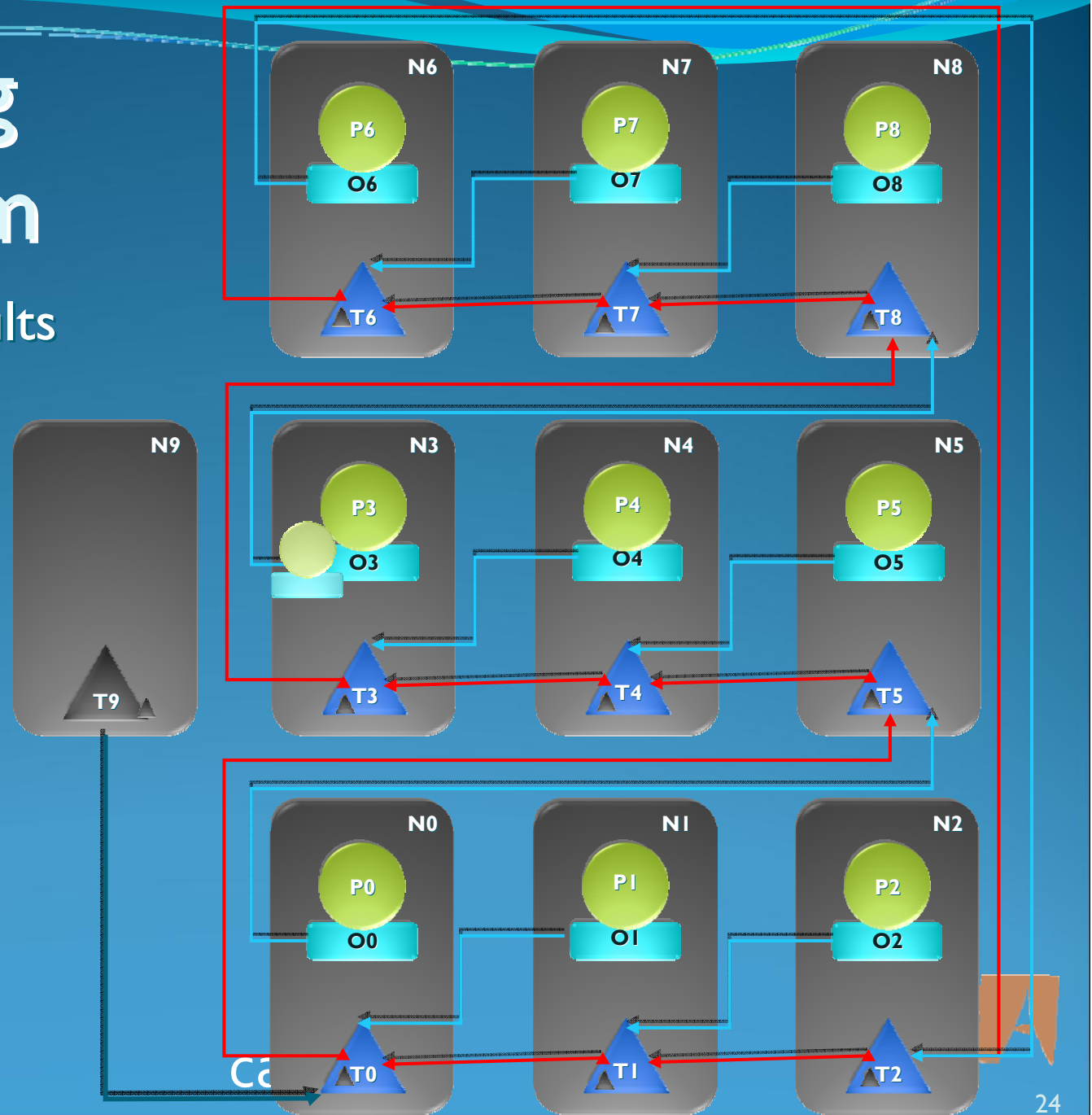- N3 is a fault probable node
- Insert a spare node
- Oportuniscally inject a fault in N3
- P3 will recover in N9

# Implementation

- We adapted the RADIC prototype (RADICMPI) in order to incorporate the Flexible Dynamic Redundancy
  - Creation of management functions
  - Definition of a new protocol to communicate with spares
  - Increment of communication between observers and local protector
  - Changes in the RADIC fault masking procedure including search in the spare nodes

caos

# Implementation

- We also implemented a set of MPI non-blocking functions, allowing to run more kind of application
  - We take care about the fault tolerance issues
  - We had to change the message management kernel of RADICMPI to deal with asynchronous communications
  - Other issues are solved too:
  - We changed the original message log approach of RADICMPI to an event log approach in order to assure the determinism of the recovery process.

caos

# Experiments Methodology

- We conduced two kind of experiments:
  - Validation
    - Spare adding task
    - Recovery task using spare

    { Using the RADICMPI Debug Log

    Running a Ping-Pong program

  - Evaluation
    - According with the fault moment
    - According with the number of nodes
    - Throughput behavior in continuous running applications

    { 
    - Matrix Product – Static
    - Matrix Product – Dynamic
    - Matrix Product – SPMD
    - N-Body simulation

caos

# Experiments Methodology

- Validation metodology

```
        ┌──────────┐
        │   Init   │
        └──────────┘
         │        │
         ▼        ▼
┌──────────────┐ ┌──────────────┐
│ Define test  │ │Define expected│
│  scenario    │ │ Log results  │
└──────────────┘ └──────────────┘
         │
         ▼
┌──────────────┐
│Define Events to│
│  be logged   │
└──────────────┘
         │
         ▼
┌──────────────┐
│ Execute tests│
└──────────────┘
         │
         ▼
  ┌─────────────────┐
  │Analyze Debug Log│
  └─────────────────┘
         │
         ▼
    ┌──────────┐
    │   End    │
    └──────────┘
```

| Col | Field Name | Description |
|-----|------------|-------------|
| 1 | Event Id | Identifies the Event Type |
| 2 | Event time | Elapsed time since startup |
| 3 | Process ID | Rank of Observer/Protector |
| 4 | Function name | Internal Function triggering this event |
| 5 | Event | Description of the event |

caos

# Experiment Design

- **Evaluating according the fault moment**
  - The experiments were conduced in a twelve node cluster
  - We executed a product of two 1500x1500 matrix using the SPMD paradigm over 9 nodes and a product of two 1000x1000 matrixes using a static distribution over 11 nodes
  - We injected one fault at 25%, 50% and 75%
  - We measured the execution times in three situations
    - Fault-free
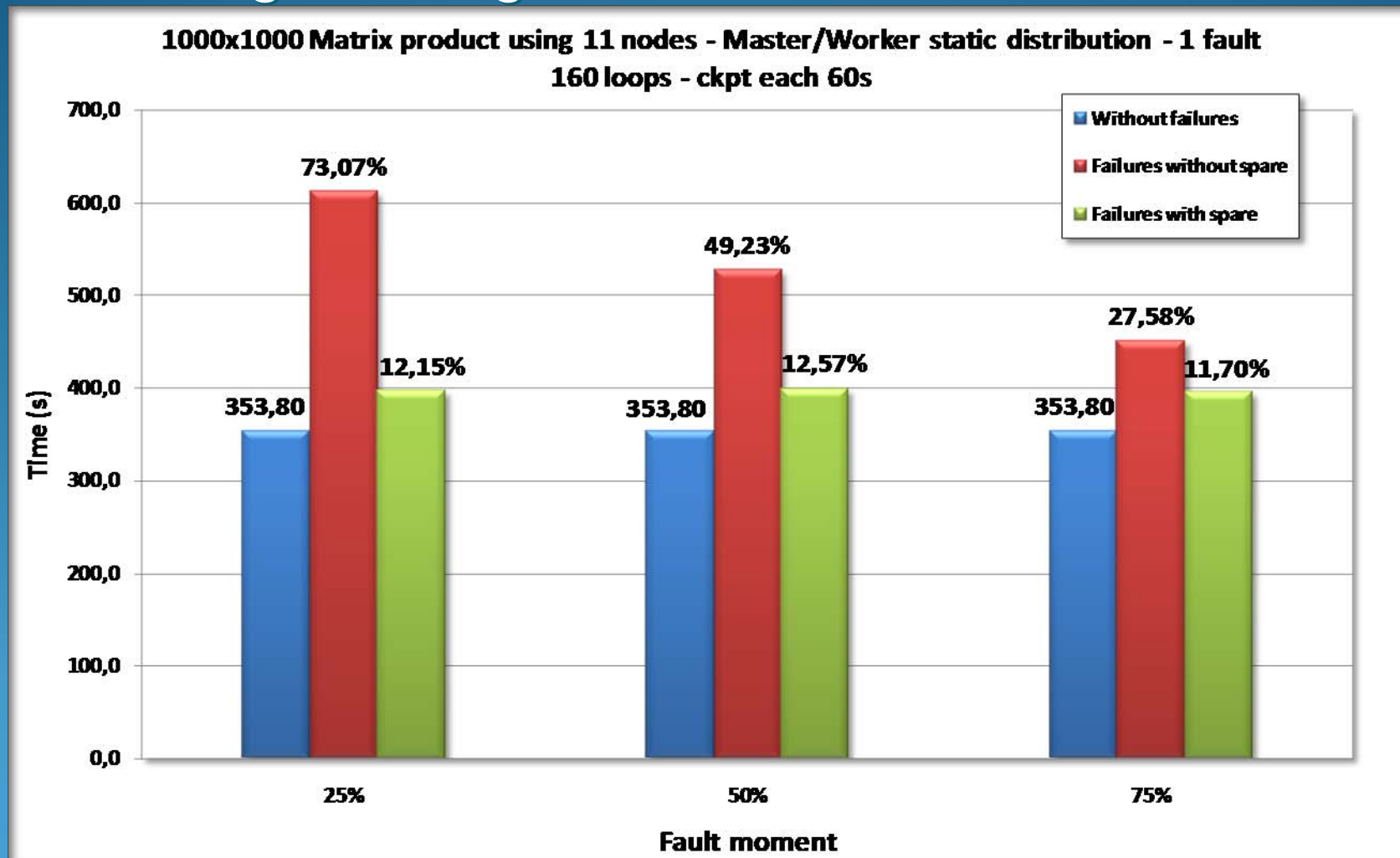    - Without spare
    - With spare

caos

# Experimental Results

- Evaluating according with the fault moment



1500x1500 Matrix product using 9 nodes - SPMD cannon algorithm- 1 fault
160 loops -ckpt each 60s

# Experimental Results

- Evaluating according with the fault moment

**1000x1000 Matrix product using 11 nodes - Master/Worker static distribution - 1 fault**
**160 loops - ckpt each 60s**

Legend:
- Without failures
- Failures without spare
- Failures with spare

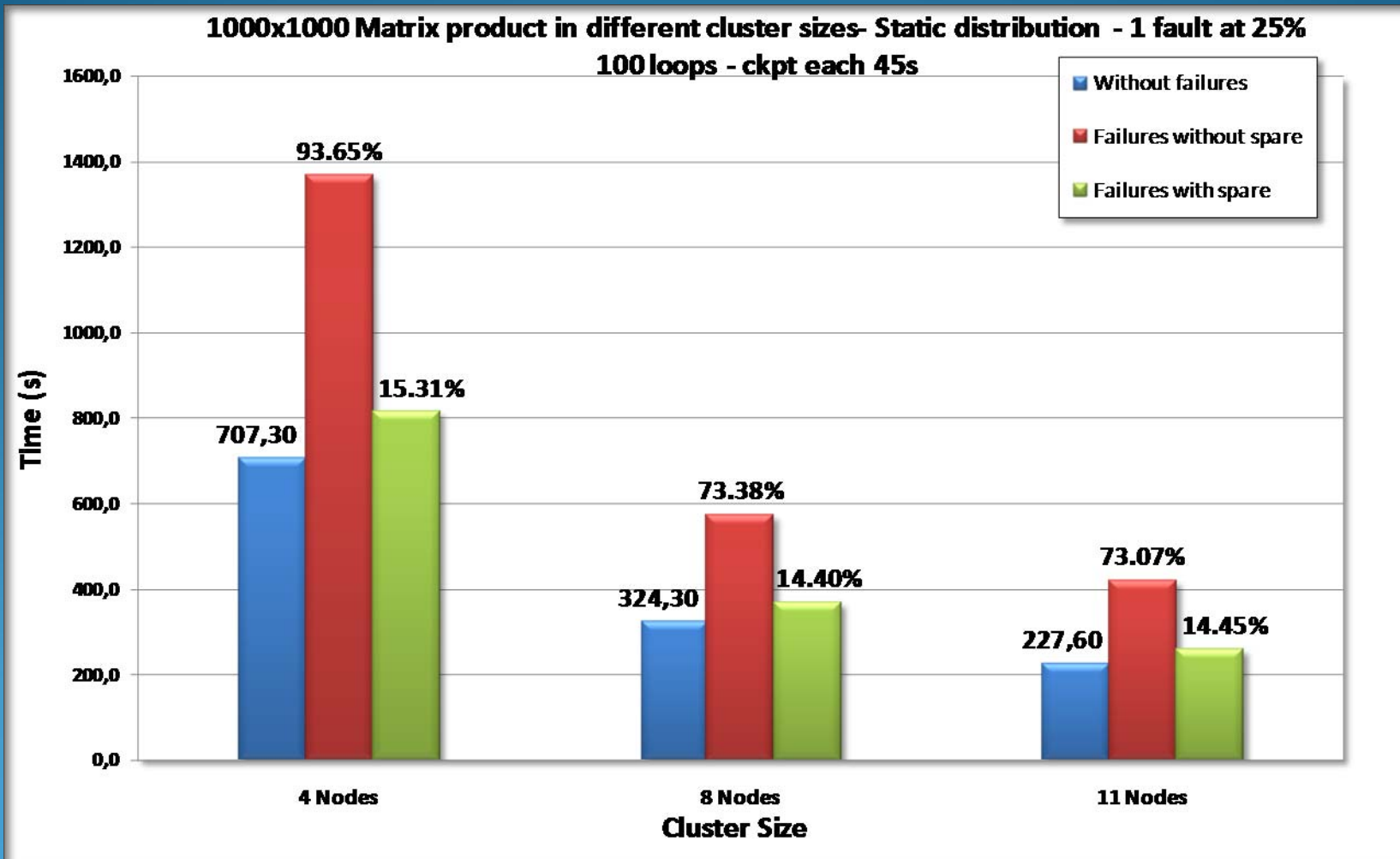| Fault moment | Without failures | Failures without spare | Failures with spare |
|---|---|---|---|
| 25% | 353,80 | 73,07% | 12,15% |
| 50% | 353,80 | 49,23% | 12,57% |
| 75% | 353,80 | 27,58% | 11,70% |

Time (s) — y-axis from 0,0 to 700,0

# Experiment Design

- **Evaluating according the number of nodes**
  - The experiments were conduced in a twelve node cluster
  - We executed a product of two 1000x1000 matrix using a dynamic distribution and a product of two 1000x1000 matrixes using a static distribution
  - We injected one fault at 25%
  - We ran the programs over three cluster sizes
    - 4 nodes + spare
    - 8 nodes + spare
    - 11 nodes + spare
  - We measured the execution times in three situations
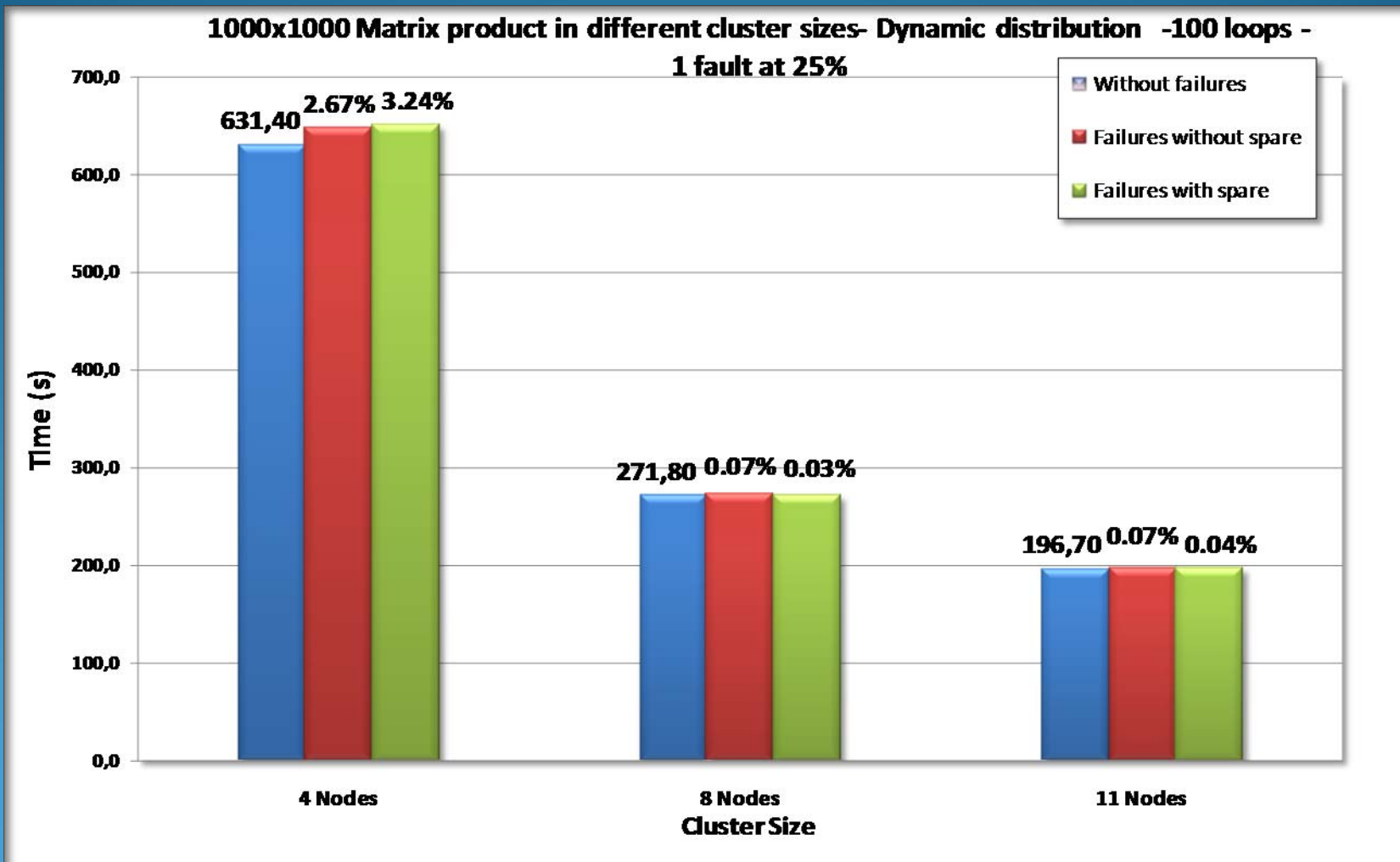    - Fault-free
    - Without spare
    - With spare

caos

# Experimental Results

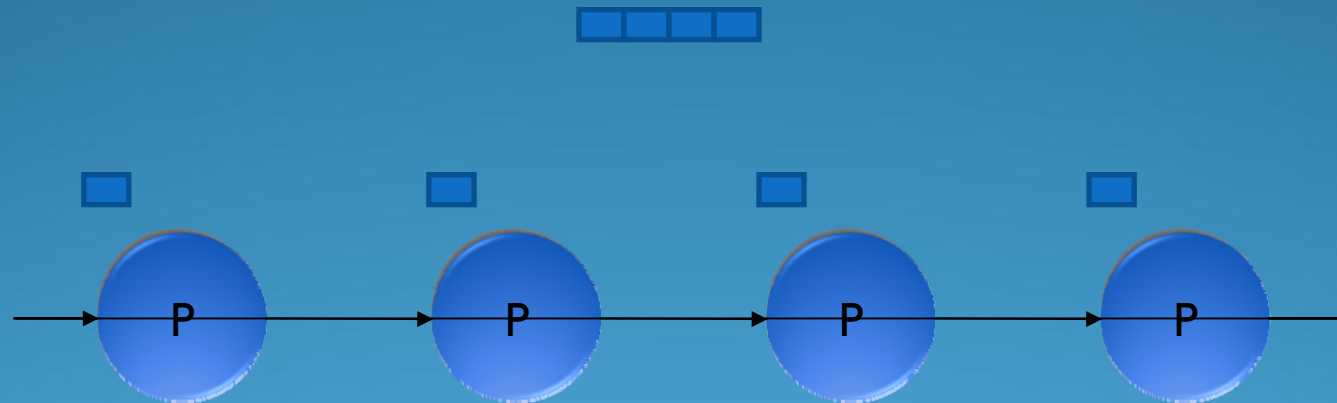- Evaluating according with the number of nodes



1000x1000 Matrix product in different cluster sizes- Static distribution - 1 fault at 25%
100 loops - ckpt each 45s

# Experimental Results

- Evaluating according with the number of nodes



**1000x1000 Matrix product in different cluster sizes- Dynamic distribution -100 loops - 1 fault at 25%**

Legend:
- Without failures
- Failures without spare
- Failures with spare

4 Nodes: 631,40 — 2.67% 3.24%
8 Nodes: 271,80 — 0.07% 0.03%
11 Nodes: 196,70 — 0.07% 0.04%

Y-axis: Time (s)
X-axis: Cluster Size

# Experimental Results

- We implemented a N-Body simulation in order to study the behavior of our solution over continuous running applications

P → P → P → P

caos
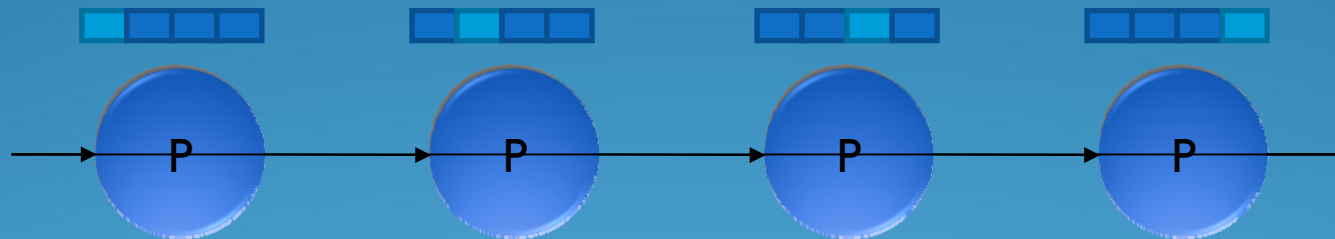
# Experimental Results

- We implemented a N-Body simulation in order to study the behavior of our solution over continuous running applications
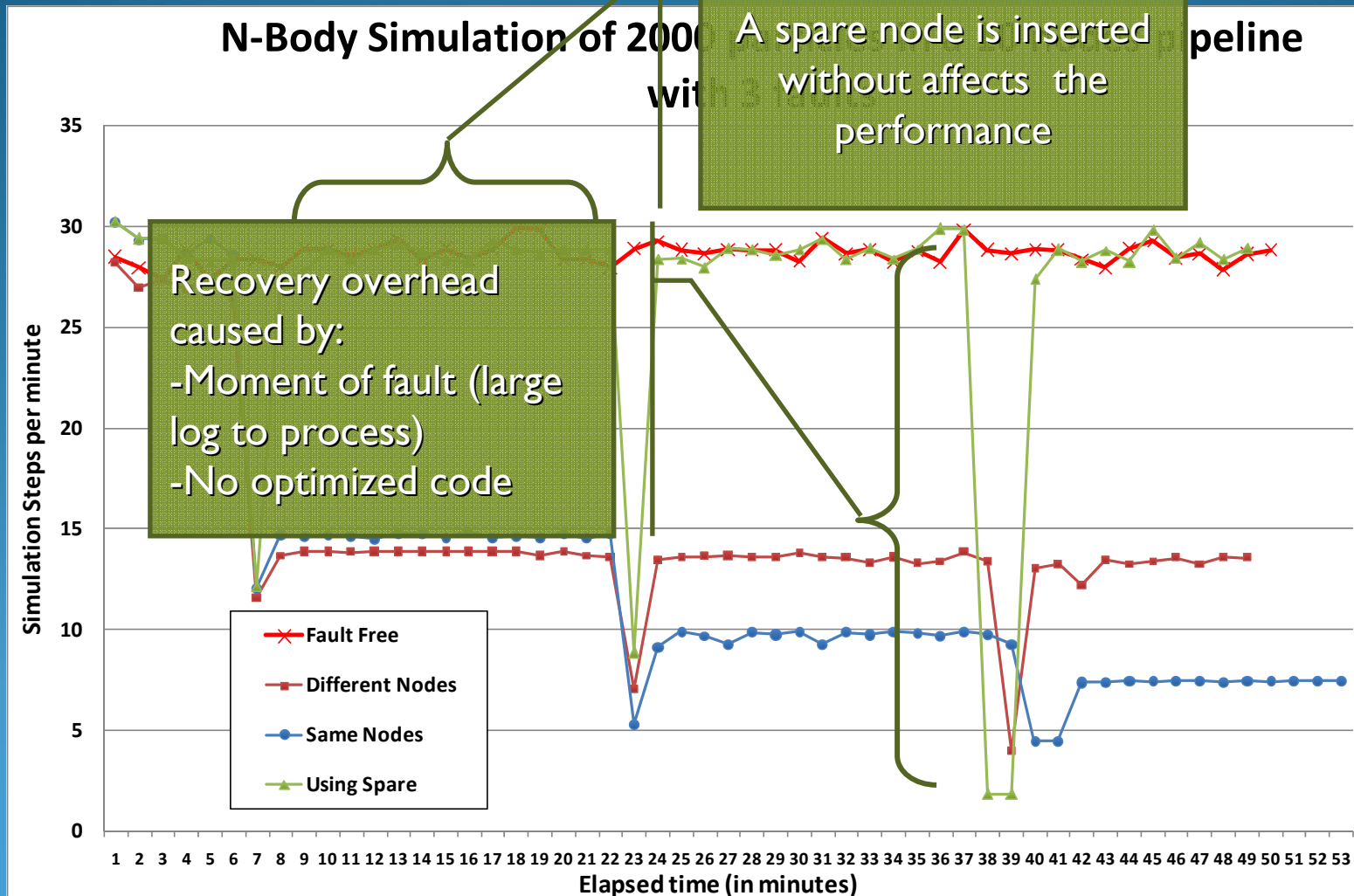
caos

# Experiment Design

- **Throughput of 24x7 applications**
  - We executed this application simulating 2000 particles in a ten node pipeline
  - We measured the throughput (in simulation steps per minute) in four scenarios
    - Fault-free
    - Injecting three faults without spare:
      - Processes recovering in different nodes
      - Process recovering in same node without spare
    - Injecting three faults using two initial spares and re-insert a "repaired" one after the first fault.

caos

# Experimental Results

- Throughput of 24x7 applications with three faults

**N-Body Simulation of 2000 particles with instruments pipeline with 3 faults**

A spare node is inserted without affects the performance

Recovery overhead caused by:
-Moment of fault (large log to process)
-No optimized code

Y-axis: **Simulation Steps per minute** (0, 5, 10, 15, 20, 25, 30, 35)

X-axis: **Elapsed time (in minutes)** (1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53)

Legend:
- Fault Free
- Different Nodes
- Same Nodes
- Using Spare

# Conclusions

- We implemented a dynamic redundancy functionality that avoids or mitigates the recovery side-effects
- This functionality is flexible because
  - It allows system configuration restablishment by dynamic insertion of spare nodes - **RESTORING**
  - It allows a transparent management of the spare nodes use - **AVOIDING**
  - It incorporates a **maintenance** feature that also allows to prevent failures by replacing fault-probable nodes – **PREVENTING**
  - Configurable, allowing from 0 to N number of spares
- We extend the RADICMPI functionality, implementing a set of MPI non-blocking functions

caos

# Conclusions

- Our results show the benefits of the dynamic redundancy solution in different scenarios.

- The results also show a strong dependency between the recovery side-effects and the application characteristics and how we can adapt to each one.

caos

# Future Work

- To study the spare nodes allocating facing factors like degradation level acceptable or memory limits of a node
- To integrate a fault prediction mechanism in the maintenance feature
- To continue expanding the RADICMPI functionality

caos

# Open Lines

- To investigate how adapt and use RADIC II with new HPC trends like the clusters of multicore computers.

- To achieve a RADIC II analytical model, allowing to determine better parameter values.

- To develop a RADIC II simulator allowing to assess its behavior in large clusters

- To incorporate other features towards an autonomic fault tolerant system.

caos