

# Deadline Scheduling with Priority for Client-Server Systems on the Grid

Eddy CARON, Pushpinder Kaur CHOUHAN, Frédéric DESPREZ



8 November 2004  
GRID Computing

# Outline

- 1 Introduction
- 2 Client-Server Scheduler
- 3 Simulation
- 4 Inference

# Outline

- 1 Introduction
  - Grid Overview
  - Notations and constraints
- 2 Client-Server Scheduler
- 3 Simulation
- 4 Inference

# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- GridRPC : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely

# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- **GridRPC** : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely

# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- **GridRPC** : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely

# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- **GridRPC** : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely

# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- **GridRPC** : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely



# Grid and GridRPC

- **Grid** : Platform resulted from aggregating distributed computers and storage units
  - Renting computation power and memory capacity
  - Need of Problem Solving Environments
- **GridRPC** : Implement Remote Procedure Call model over the Grid
  - Good and simple paradigm to implement the Grid
  - Run computation remotely

# Introduction

- **Goal** : To find a scheduling algorithm that can consider both priority and deadline of a task.
- **Foundation** : A Study of Deadline Scheduling for Client-Server Systems on the Computational Grid by Takefusa et al.
- **Motivation** : Allocate the tasks on server where they can meet their deadline with in the priority constraint.
- **Requirement** : Performance prediction tool FAST and modelization

# Introduction

- **Goal** : To find a scheduling algorithm that can consider both priority and deadline of a task.
- **Foundation** : A Study of Deadline Scheduling for Client-Server Systems on the Computational Grid by Takefusa et al.
- **Motivation** : Allocate the tasks on server where they can meet their deadline with in the priority constraint.
- **Requirement** : Performance prediction tool FAST and modelization

# Introduction

- **Goal** : To find a scheduling algorithm that can consider both priority and deadline of a task.
- **Foundation** : A Study of Deadline Scheduling for Client-Server Systems on the Computational Grid by Takefusa et al.
- **Motivation** : Allocate the tasks on server where they can meet their deadline with in the priority constraint.
- **Requirement** : Performance prediction tool FAST and modelization

# Introduction

- **Goal** : To find a scheduling algorithm that can consider both priority and deadline of a task.
- **Foundation** : A Study of Deadline Scheduling for Client-Server Systems on the Computational Grid by Takefusa et al.
- **Motivation** : Allocate the tasks on server where they can meet their deadline with in the priority constraint.
- **Requirement** : Performance prediction tool FAST and modelization

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)



# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Notation

$$T_{aS_i} = \frac{W_{send}}{P_{send}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{P_S}$$

- $T_{aS_i}$ : the execution time for the task  $a$  on server  $i$ .
- $W_{send}$ : the size of the data transmitted from the client to the server
- $P_{send}$ : the predicted network throughputs from the client to the server
- $P_{recv}$ : the predicted network throughputs from the server to the client and  $P_S$  is the server performance
- $W_{recv}$ : the data size transmitted from the server to the client
- $W_{aS_i}$ : the number of floating point operations of the task
- $P_S$ : the server performance (floating point operations per second)

# Outline

## 1 Introduction

## 2 Client-Server Scheduler

- With load measurements
- With a Forecast Correction Mechanism
- With a Priority Mechanism

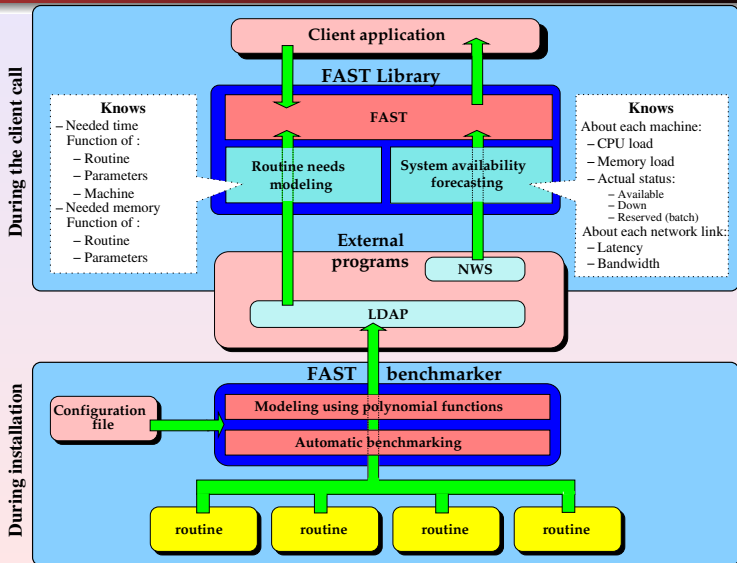
## 3 Simulation

## 4 Inference

# Straightforward algorithm

```
1: repeat  
2:   for all server  $S_i$  do  
3:     if can_do( $S_i, T_a$ ) then  
4:        $T_{aS_i} = \frac{W_{send}}{F_{Bd}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{F_{S_i}}$   
5:       List=sort_insert(List,  $T_{aS_i}, T_a, S_i$ )  
6:     end if  
7:   end for  
8:    $num\_submit = \text{task\_ack}(List, \frac{2 \times W_{send}}{F_{Bd}})$   
9:   task_submit( $List[num\_submit]$ )  
10: until the end
```

# FAST's overview



# Scheduling with Forecast Correction Mechanism

$$T_{aS_i} = \frac{T_{aS_i} \times \text{CorrecFAST}}{100}$$

$$\text{CorrecFAST} = \frac{\text{nb\_exec} \times \text{CorrecFAST} + \frac{100 \times T_r}{T_{aS_i}}}{\text{nb\_exec} + 1}$$

- *CorrecFAST*: an error average between the prediction time and the actual execution time.
- $T_r$  : the actual execution time
- $T_{aS_i}$ : the time predicted execution time



# Scheduling with Forecast Correction Mechanism

$$T_{aS_i} = \frac{T_{aS_i} \times \text{CorrecFAST}}{100}$$

$$\text{CorrecFAST} = \frac{\text{nb\_exec} \times \text{CorrecFAST} + \frac{100 \times T_r}{T_{aS_i}}}{\text{nb\_exec} + 1}$$

- *CorrecFAST*: an error average between the prediction time and the actual execution time.
- $T_r$ : the actual execution time
- $T_{aS_i}$ : the time predicted execution time

# Scheduling with Forecast Correction Mechanism

$$T_{aS_i} = \frac{T_{aS_i} \times \text{CorrecFAST}}{100}$$

$$\text{CorrecFAST} = \frac{\text{nb\_exec} \times \text{CorrecFAST} + \frac{100 \times T_r}{T_{aS_i}}}{\text{nb\_exec} + 1}$$

- *CorrecFAST*: an error average between the prediction time and the actual execution time.
- $T_r$ : the actual execution time
- $T_{aS_i}$ : the time predicted execution time

# Scheduling with Forecast Correction Mechanism

$$T_{aS_i} = \frac{T_{aS_i} \times \text{CorrecFAST}}{100}$$

$$\text{CorrecFAST} = \frac{\text{nb\_exec} \times \text{CorrecFAST} + \frac{100 \times T_r}{T_{aS_i}}}{\text{nb\_exec} + 1}$$

- *CorrecFAST*: an error average between the prediction time and the actual execution time.
- $T_r$ : the actual execution time
- $T_{aS_i}$ : the time predicted execution time

# Scheduling with Forecast Correction Mechanism

$$T_{aS_i} = \frac{T_{aS_i} \times \text{CorrecFAST}}{100}$$

$$\text{CorrecFAST} = \frac{\text{nb\_exec} \times \text{CorrecFAST} + \frac{100 \times T_r}{T_{aS_i}}}{\text{nb\_exec} + 1}$$

- *CorrecFAST*: an error average between the prediction time and the actual execution time.
- $T_r$ : the actual execution time
- $T_{aS_i}$ : the time predicted execution time

# Scheduling with Forecast Correction Mechanism

```
1:  $CorrecFAST = 100$ 
2:  $nb\_exec = 0$ 
3: for all server  $S_i$  do
4:   if  $can\_do(S_i, T_a)$  then
5:      $T_{aS_i} = \frac{W_{send}}{F_{Bd}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{F_{S_i}}$ 
6:      $T_{aS_i} = \frac{T_{aS_i} \times CorrecFAST}{100}$ 
7:      $List = \text{sort\_insert}(List, T_{aS_i}, T_a, S_i)$ 
8:   end if
9: end for
10:  $num\_submit = \text{task\_ack}(List, \frac{2 \times W_{send}}{F_{Bd}})$ 
11:  $T_r = \text{task\_submit}(List[num\_submit])$ 
12:  $CorrecFAST = \frac{nb\_exec \times CorrecFAST + \frac{100 \times T_r}{T_{aS_i}}}{nb\_exec + 1}$ 
13:  $nb\_exec++$ 
```

# Priority and Deadline Mechanism

```
1: repeat
2:   for all server  $S_i$  do
3:     if can_do( $S_i, T_a$ ) then
4:        $T_{aS_i} = \frac{W_{send}}{F_{Bd}} + \frac{W_{recv}}{P_{recv}} + \frac{W_{aS_i}}{F_{S_i}}$ 
5:     end if
6:     if  $T_{aS_i} < TD_a$  then
7:       count_fallback_tasks( $T_a, T_{aS_i}, TP_a, TD_a$ )
8:       if  $TF_{aS_i} < TD_a$  then
9:         best_server( $S_i, \text{best\_server\_name}$ )
10:      end if
11:    end if
12:  end for
13:  task_submit(best_server_name, task_name)
14:  Re-submission(task_name)
15: until the end
```

# Priority and Deadline Mechanism

- **can\_do** This function returns true if server  $S_i$  have the resource required to compute task  $T_a$ . This function takes into account the availability of memory and disk storage, the computational library etc.
- **Count\_fallback\_tasks** This function counts the fallbacked tasks. Tasks that cannot meet their deadline after the insertion of the new task are called fallbacked tasks. Task  $T_a$  is placed according to its priority  $TP_a$  on the server task queue, which may change the execution time of the tasks on the queue.

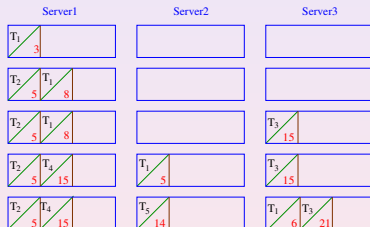
# Priority and Deadline Mechanism

- **best\_server** This function selects the best server among the servers that can execute the task within the deadline time. The best server is selected by comparing the number of fallbacked tasks. Server with less fallbacked tasks is selected for the task execution. If the servers have same number of fallbacked tasks, then the time to compute the task is compared and the server that takes less time is selected.
- **task\_submit** It performs the remote execution on the server given by the `best_server` and the argument of the function is one server and not a list of servers.
- **Re-submission** This function submits the fallbacked task to the servers, for recomputing the execution time. If any server can meet the task's deadline then the task is allocated to that server.



# Priority and Deadline Mechanism

Task	Priority	Deadline	Exec. time on server		
			$S_1$	$S_2$	$S_3$
1	3	15	3	5	6
2	5	10	5	12	9
3	2	30	11	20	15
4	4	20	10	np	17
5	5	15	12	14	np

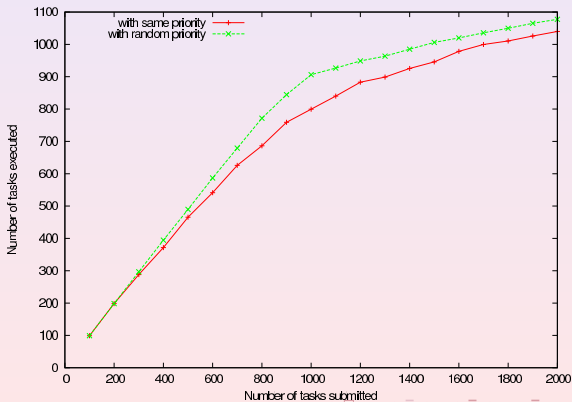


# Outline

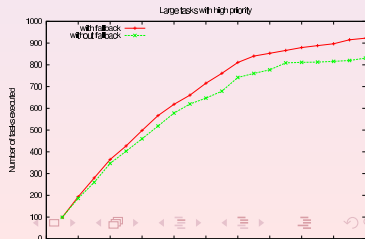
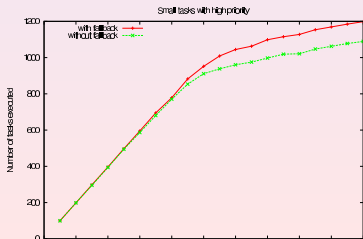
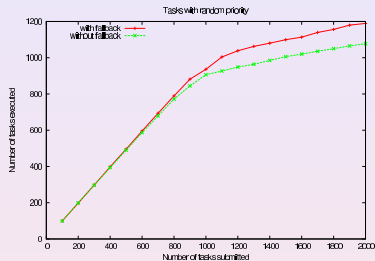
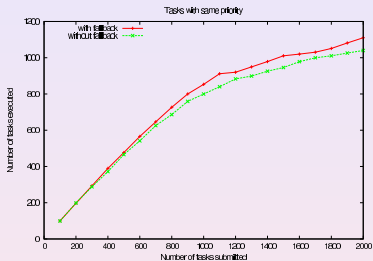
- 1 Introduction
- 2 Client-Server Scheduler
- 3 Simulation**
- 4 Inference

# Priority based tasks are executed without fallback mechanism

Testbed 100 servers, priority range 1-10, deadline =  $5 \times T_{aS_i}$



# Tasks executed with and without fallback



# Outline

- 1 Introduction
- 2 Client-Server Scheduler
- 3 Simulation
- 4 Inference

# Inference

- Conclusion
  - Introduction of the performance prediction tool FAST
    - <http://graal.ens-lyon.fr/FAST>
  - Focused on the management of tasks with respect to their priorities and deadlines
  - Increase the number of tasks that can meet their deadlines by
    - Fallback mechanism to reschedule tasks
    - Load correction mechanism (using FAST)
    - Using task's priority
- Future Work
  - Implement these scheduling algorithms in DIET platform
    - ASP middleware
    - Grid-RPC middleware
    - <http://graal.ens-lyon.fr/DIET>