SORMA-Project Meeting
University of Karlsruhe
11,12 March 2008
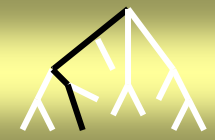
# Recent Research Activities

Viktor Yarmolenko

**People involved:**

Rizos Sakellariou        Viktor Yarmolenko

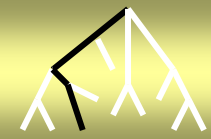rizos@cs.man.ac.uk    viktor@cs.man.ac.uk

The Univercity of Manchester
School of Computer Science
Kilburn Building, Oxford Road
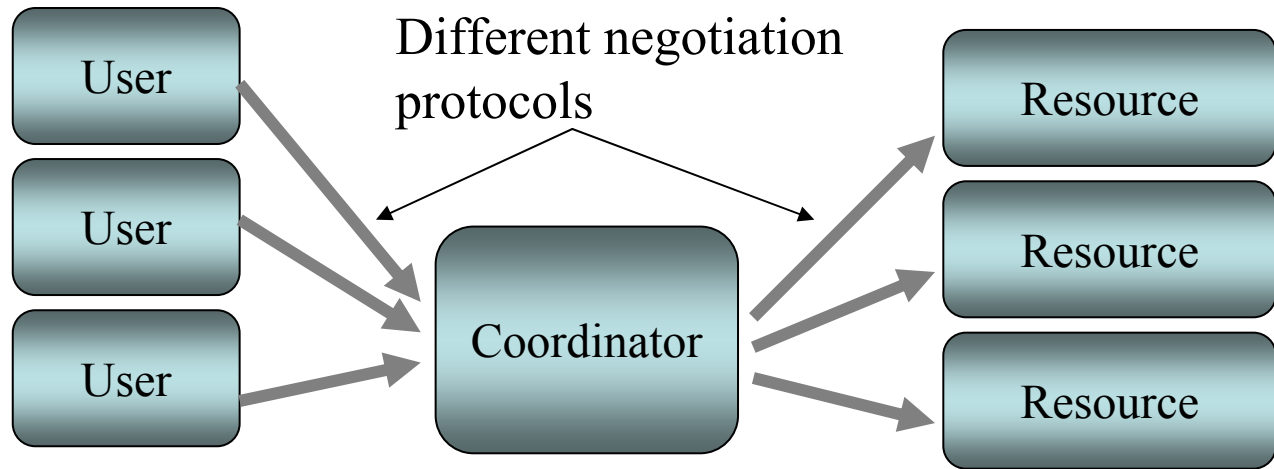Manchester M13 9PL
United Kingdom

# Main directions

- Dynamic Service Level Agreements
- SLA Aware Scheduling
- Synthetic SLA Workloads
- Other work in Progress

## Simulation Example

Job Generator – producing different user behaviour, job workloads, *etc.*

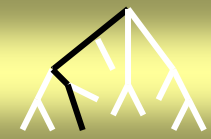Different number of Resources, each of different capacity, availability and other properties

Different negotiation protocols

Different scheduling algorithms, profit optimisations, coping with uncertainties, *etc.*

Different Coordinator Strategies



Viktor Yarmolenko, Rizos Sakellariou, Djamila Ouelhadj, Jonathan M Garibaldi, "SLA Based Job Scheduling: A Case Study on Policies for Negotiation With Resources", *Proceedings of the All Hands Meeting AHM'05*, Nottingham, UK (September 2005)

All this is in the context of Service Level Agreements.

SLA Example: Traditional Approach

CPU

$T_S$ – the earliest time the Job is allowed to start
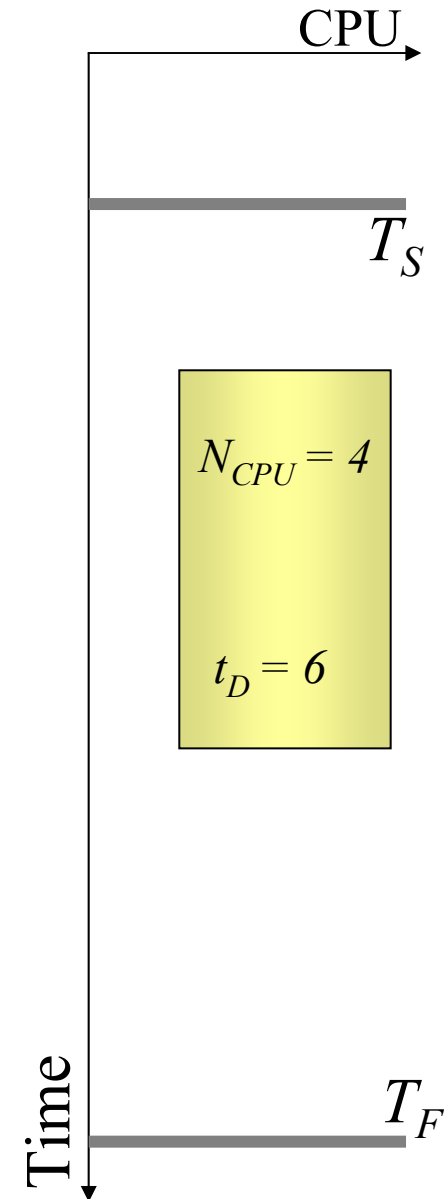
$T_F$ – the latest time the Job is allowed to finish

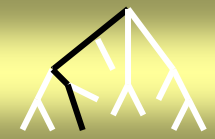$N_{CPU}$ – number of nodes required for the Job

$t_D$ – projected Job duration time for $N_{CPU}$ nodes

..........................................................................................................

$B_{job}$ – projected traffic that Job creates

$V_{pr}$ – the price for executing the Job

$V_{pn}$ – the penalty for failing the Job

$T_S$

$N_{CPU} = 4$
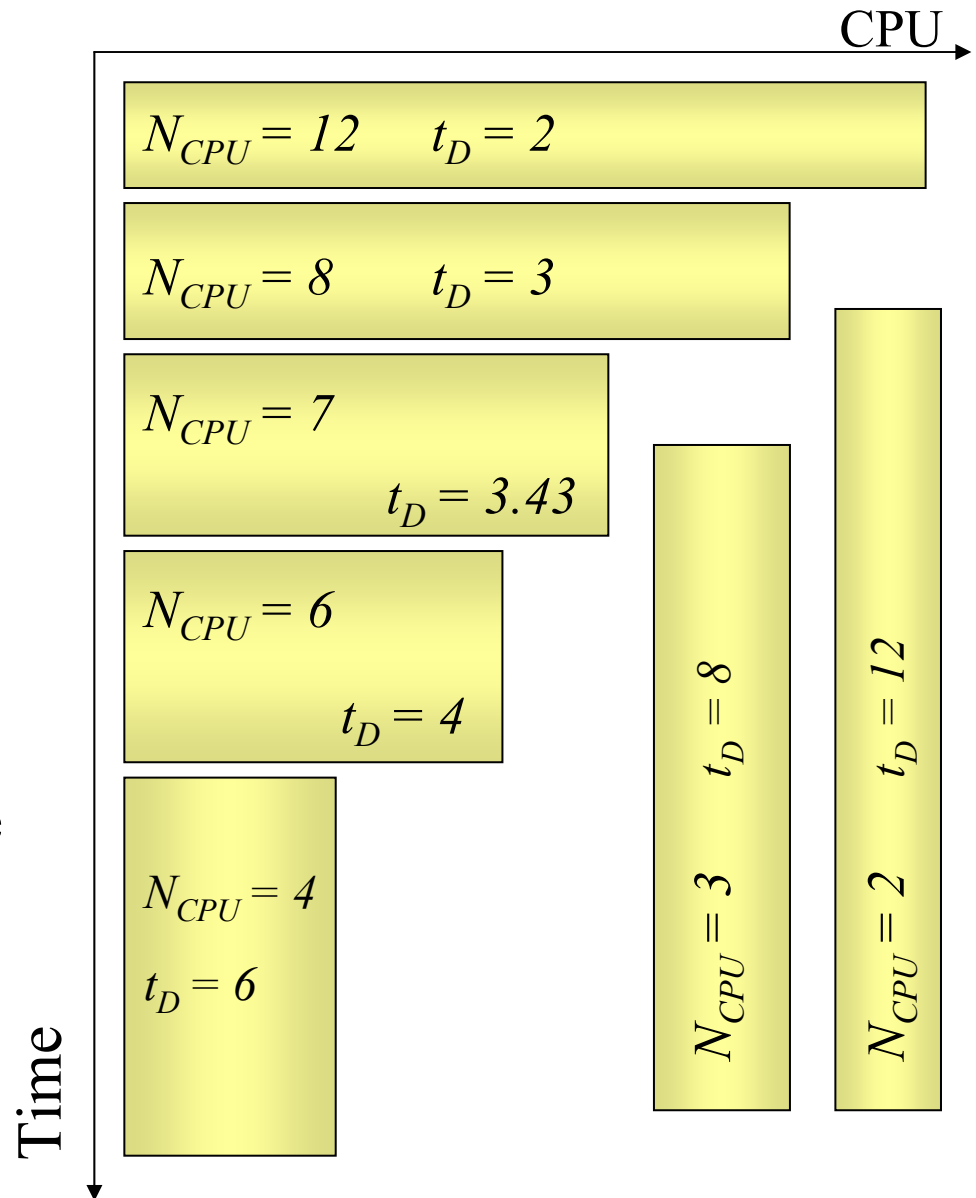
$t_D = 6$

Time

$T_F$

SLA Example: Our Approach

$T_S$, $T_F$, ... as before but ...

$N_{CPU}$ = {2,3,4,..}    is a range

$t_D = \dfrac{t_{UP}}{N_{CPU}}$    is a function

$t_{UP}$ = 24, (CPU-hours) duration

$V_{tot} = X\, t_D\, V_{pr}$    is a final value
of the agreement (for example)

CPU →

$N_{CPU}$ = 12    $t_D$ = 2

$N_{CPU}$ = 8    $t_D$ = 3

$N_{CPU}$ = 7

$t_D$ = 3.43

$N_{CPU}$ = 6

$t_D$ = 4

$N_{CPU}$ = 4

$t_D$ = 6

$N_{CPU}$ = 3    $t_D$ = 8

$N_{CPU}$ = 2    $t_D$ = 12

Time ↓

SLA Example: Same as before, but …

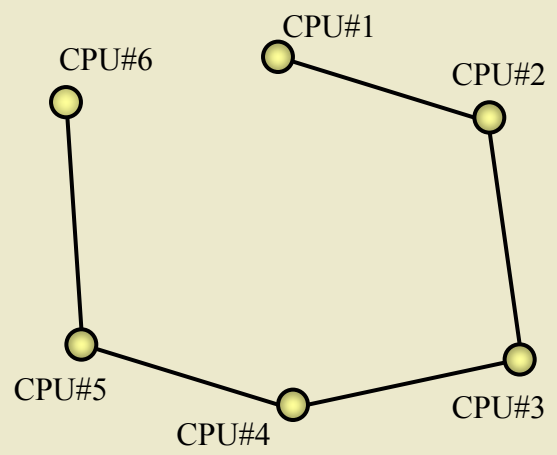$B_{RES}(t_{curr})$ , bandwidth provided by the Resource
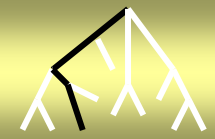
$d(n) = n + (n-1) + … + 2 + 1$ , *a known expression*

$B_{job} = B_0 \, d(N_{CPU} - 1)$ , traffic generated by the Job

$$t_D = \frac{B_{job} \, t_{UP}}{B_{RES} \, N_{CPU}} = \frac{B_0 \, t_{UP} \, (N_{CPU} - 1)}{2 B_{RES}}$$

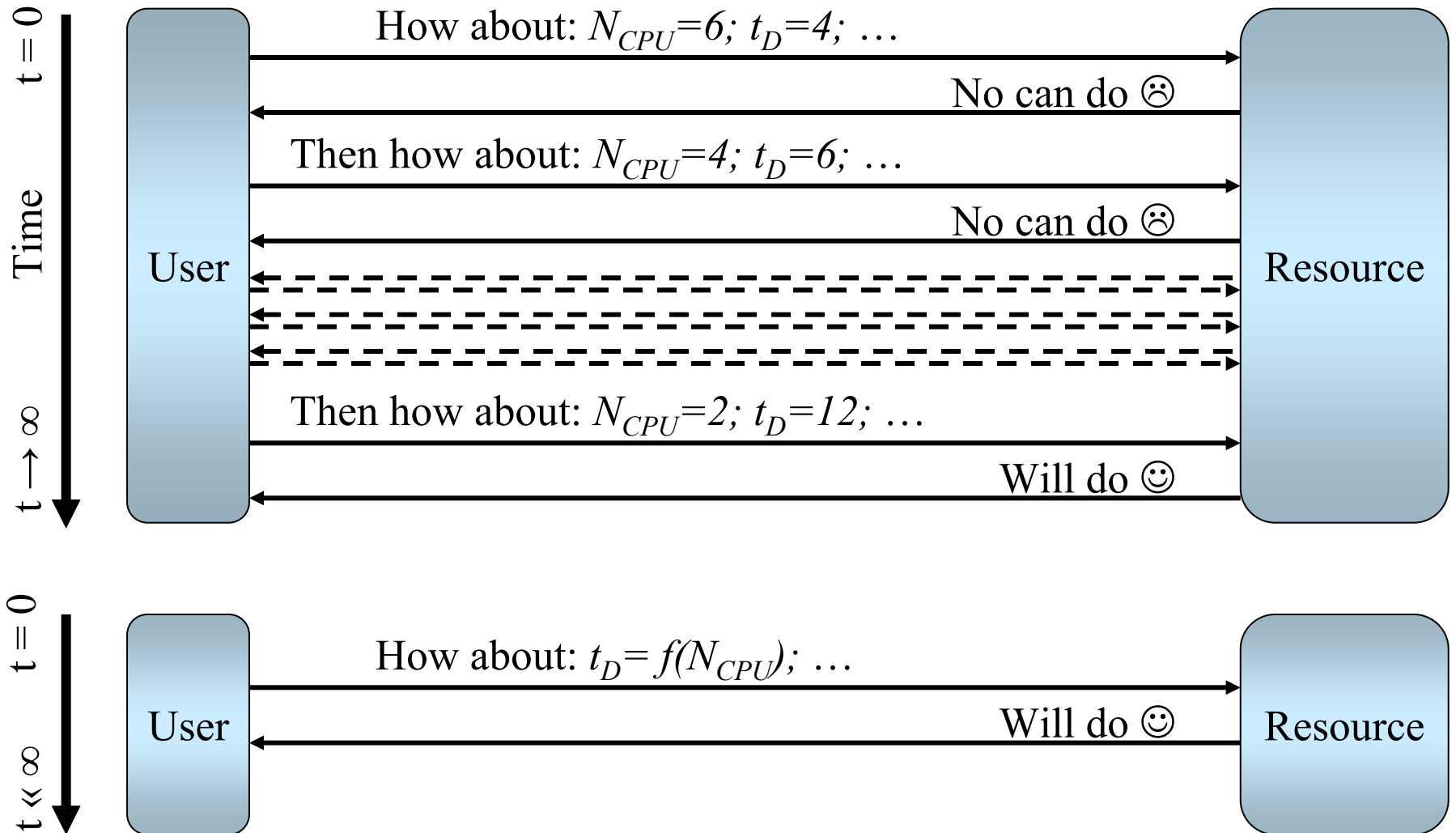$$t_D = \frac{B_{job} \, t_{UP}}{B_{RES} \, N_{CPU}} = \frac{B_0 \, t_{UP} \, (N_{CPU} - 1)}{N_{CPU} \, B_{RES}}$$
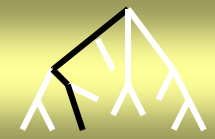
$$B_{job} = B_0 \, (N_{CPU} - 1)$$
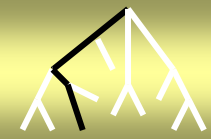
Variable CPU Scenario (Traditional vs. Expressive SLA)
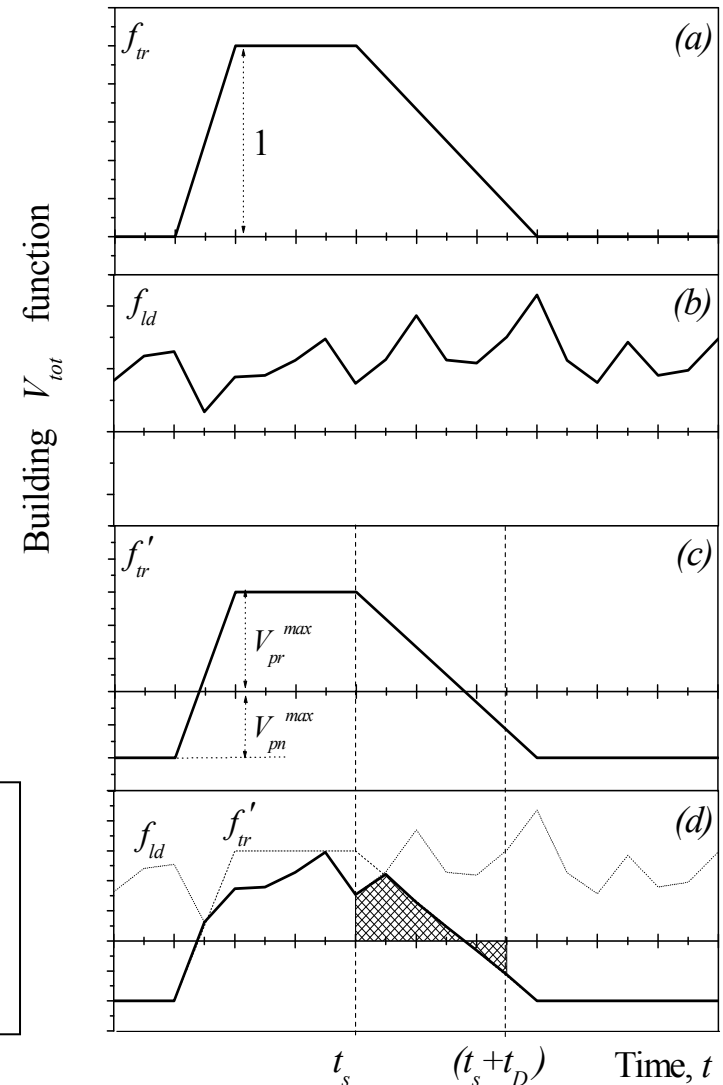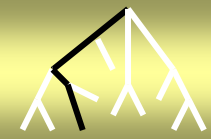
## Only Single Negotiation is Allowed

# SLA Example: Defining the Price of the Service as Function

$$t_{curr}$$

$$B_{RES}(t_{curr})$$

$$R_{ld}(t_{curr}) = f_{ld}$$

$$V_{tot} = f(R_{ld}, t_{D,} N_{CPU}, ...)$$

Don't stop here, add more functions!!!

Viktor Yarmolenko, Rizos Sakellariou, "Towards Increased Expressiveness in Service Level Agreements", *Concurrency and Computation: Practice and Experience*, vol.19, 1975-1990 (2007)
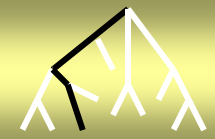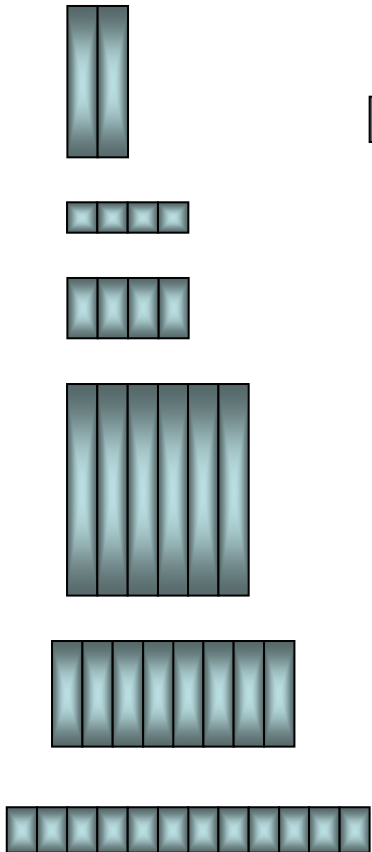
The Problem

• When client pays money, client wants guarantees, QoS, promises, etc

• These can be defined in Service Level Agreements, which preferably are legally binding – contracts.

• Once provider agreed to the terms described in SLA, provider better keep the agreement. How? If failures occur, what to do? Which SLAs to brake? How to schedule more efficiently? How to schedule to generate more income? …
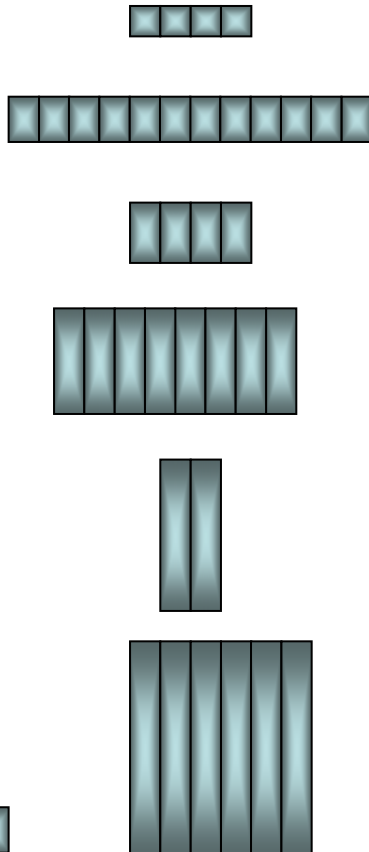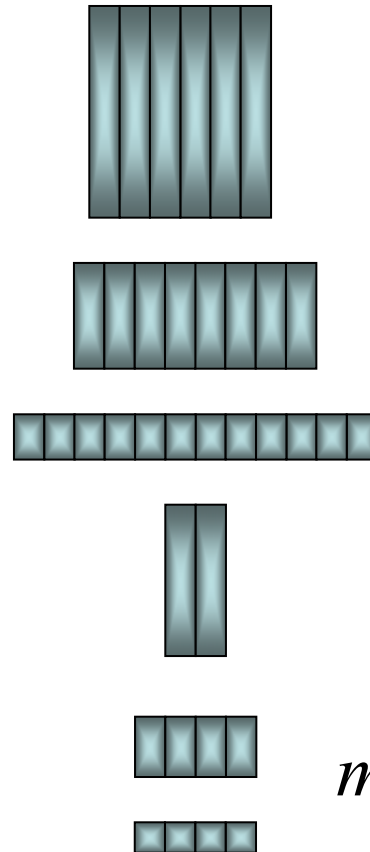
Simple and Fast Heuristics: Step 1 - Prioritising Jobs

$min(N_{CPU})$ $\quad$ $min(t_D)$ $\quad$ $max(A)$ $\qquad\qquad$ $min(T_S)$

$min(T_F)$
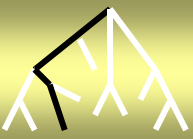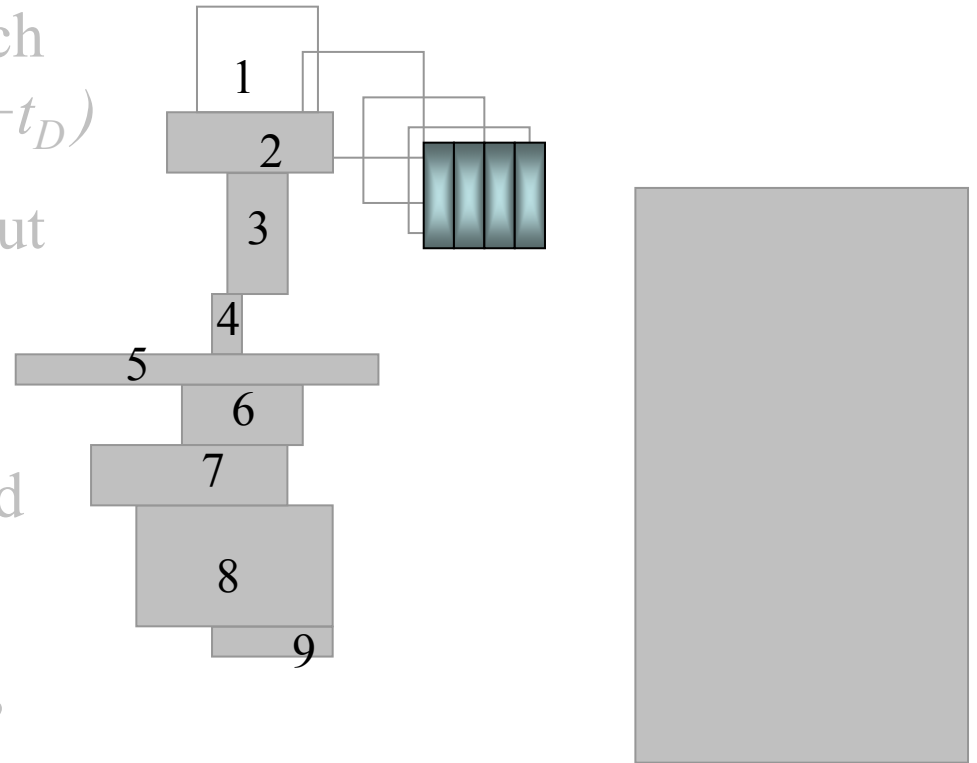
$min(t_L)$

$min(T_S + wt_D)$

$max(T_F + wA)$

$etc.$

$min(T_F + w_1A + w_2t_L)$
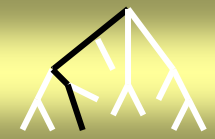
Simple and Fast Heuristics: Step 2 - Allocating Jobs

1.  Pick up the next job on the list

2.  Try to find $N_{CPU}$ nodes which are available from $T_S$ to $(T_S+t_D)$

3.  If unsuccessful, try step 2 but with $(T_S + \Delta t)$ to $(T_S+\Delta t+t_D)$

4.  Repeat steps 2 and 3 while $(T_S+\Delta t+t_D) < T_F$ or until find enough free nodes

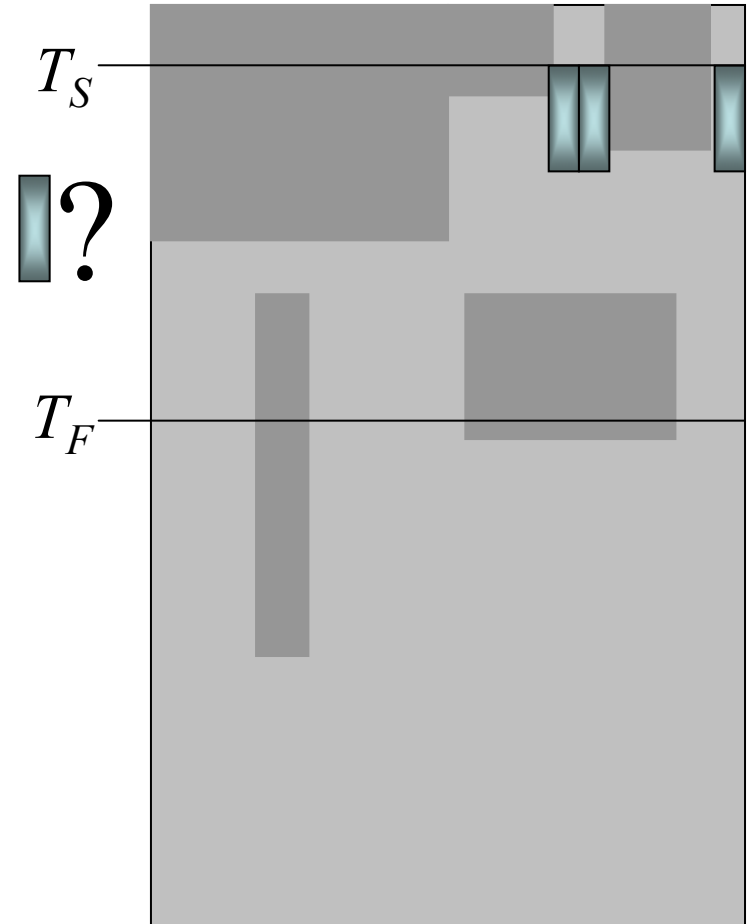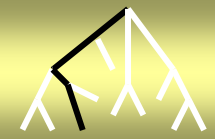5.  If failed to find $N_{CPU}$ nodes, reject the request.

## Simple and Fast Heuristics: Step 2 - Allocating Jobs

1. Pick up the next job on the list

2. Try to find $N_{CPU}$ nodes which are available from $T_S$ to $(T_S+t_D)$

3. If unsuccessful, try step 2 but with $(T_S + \Delta t)$ to $(T_S+\Delta t+t_D)$

4. Repeat steps 2 and 3 while $(T_S+\Delta t+t_D) < T_F$ or until find enough free nodes

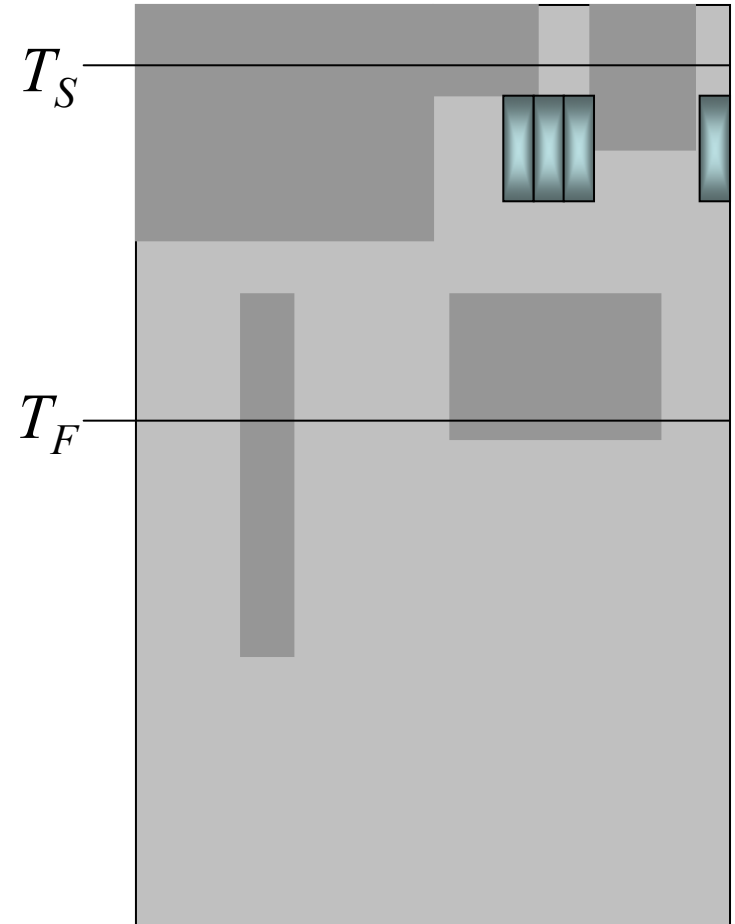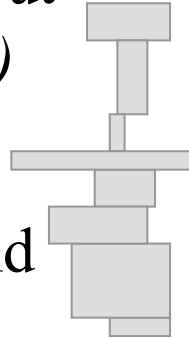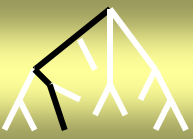5. If failed to find $N_{CPU}$ nodes, reject the request.

$T_S$

$T_F$

$?$

## Simple and Fast Heuristics: Step 2 - Allocating Jobs

1. Pick up the next job on the list

2. Try to find $N_{CPU}$ nodes which are available from $T_S$ to $(T_S+t_D)$

3. If unsuccessful, try step 2 but with $(T_S + \Delta t)$ to $(T_S+\Delta t+t_D)$

4. Repeat steps 2 and 3 while $(T_S+\Delta t+t_D) < T_F$ or until find enough free nodes

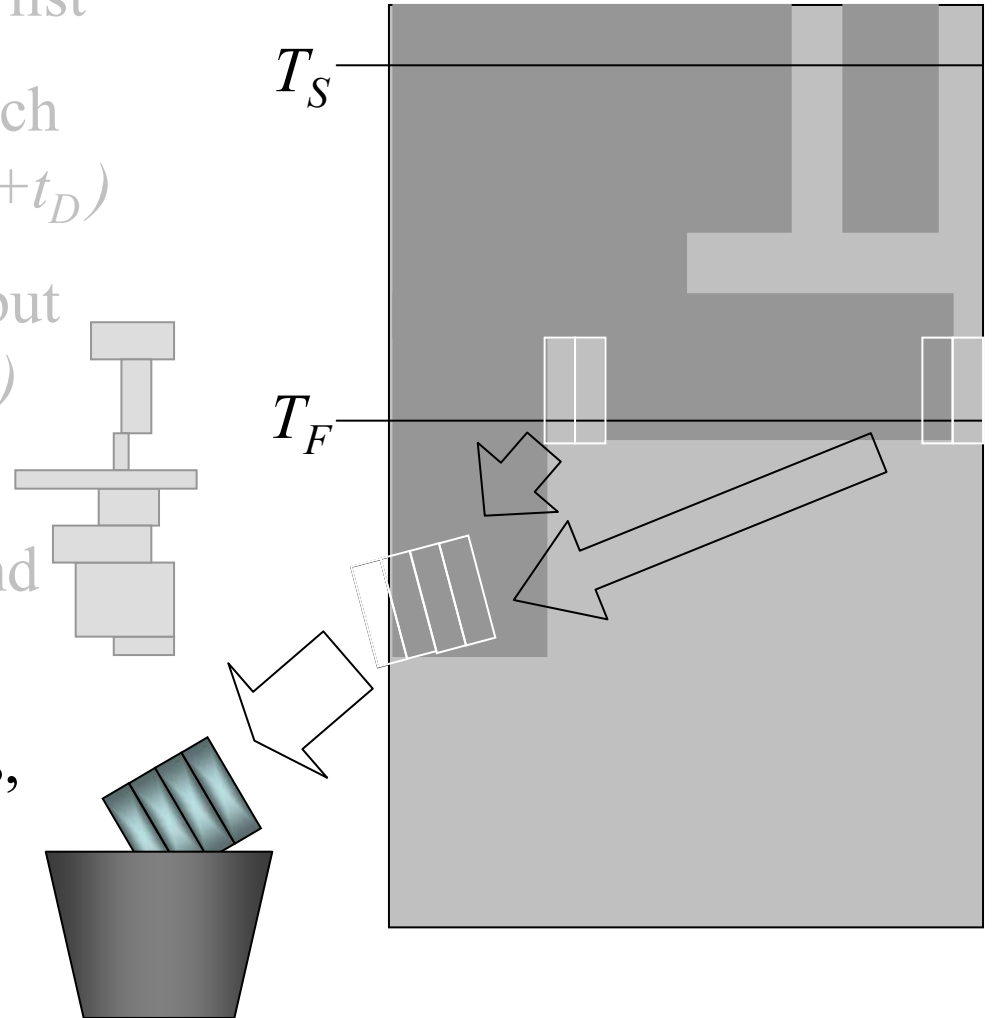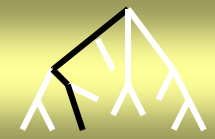5. If failed to find $N_{CPU}$ nodes, reject the request.

$T_S$

$T_F$

Simple and Fast Heuristics: Step 2 - Allocating Jobs

1. Pick up the next job on the list

2. Try to find $N_{CPU}$ nodes which are available from $T_S$ to $(T_S+t_D)$

3. If unsuccessful, try step 2 but with $(T_S + \Delta t)$ to $(T_S+\Delta t+t_D)$

4. Repeat steps 2 and 3 while $(T_S+\Delta t+t_D) < T_F$ or until find enough free nodes

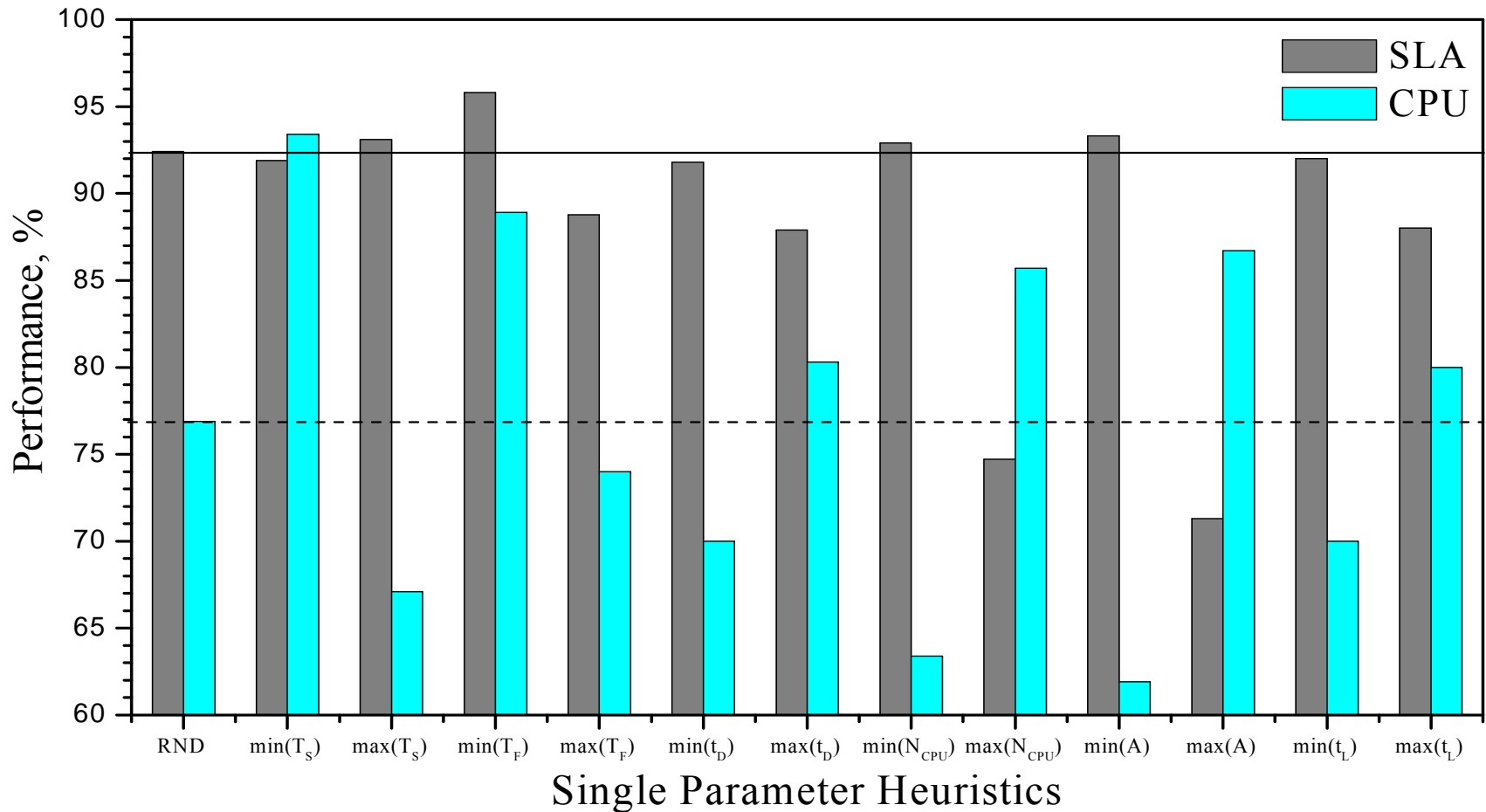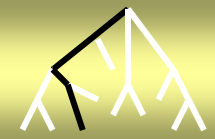5. If failed to find $N_{CPU}$ nodes, reject the request.
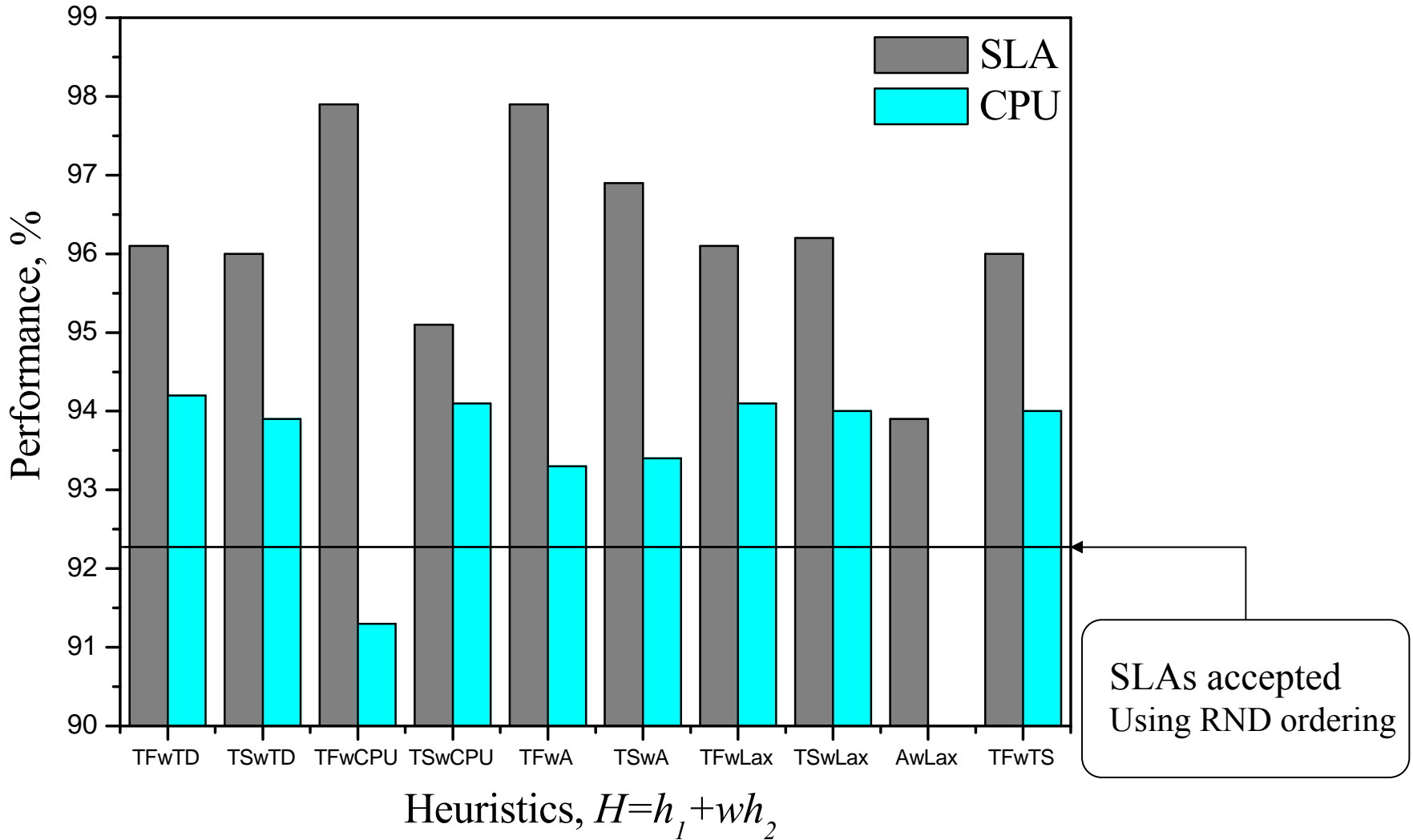
$T_S$

$T_F$

Results: Single Parameter Ordering

Order jobs in ascending ( *min(H)* ) or descending ( *max(H)* ) order of *H*.

Results: Two Parameter Ordering



SLAs accepted
Using RND ordering

## Results: Three Parameter Ordering



$h_0 = T_F,$
$h_1 = A,$
$h_2 = t_L,$

Results: Interesting Observations

## For best performance jobs must be always ordered by the lowest ⬇ or the highest ⬆ parameter first

| Pricing | $T_F$ | $T_S$ | $t_D$ | $t_L$ | $A$ | $N_{CPU}$ |
|---------|-------|-------|-------|-------|-----|-----------|
| ⊓ SLA | ⬇ | ⬇ | ⬇ | ⬇ | ⬇ | ⬇ |
| ⊓ CPU | ⬇ | ⬇ | ⬆ | ⬆ | ⬆ | ⬆ |

Other pricing policies were explored:

Viktor Yarmolenko, Rizos Sakellariou, "An Evaluation of Heuristics for SLA Based Parallel Job Scheduling", *Proceedings of the 3rd High Performance Grid Computing Workshop (HPGC)* (in conjunction with IPDPS 2006), Rhodes, Greece (April 2006), IEEE Computer Society Press

Higher level overview is here:

Rizos Sakellariou, Viktor Yarmolenko, "Job Scheduling on the Grid: Towards SLA-Based Scheduling", in L. Grandinetti (ed.), *High Performance Computing and Grids in Action"*, IOS Press, 2008

Workload: The Challenges

- The performance of the system is heavily dependant on the type of the workload

- There is no or very little data available about SLA workloads

- We need a tool that would allow to assemble a synthetic workload based on real logs, existing models, various other rules, and their combinations, easily integrated in a workload generator tool which would provide a workload stream with desired characteristics and behaviour.

- Novel workloads require new metrics or ways that characterise them (job rate, or throughput may not longer be an adequate description of a workload)

## Workload Generator Framework: The Concept (paper in preparation)



N of CPUs

Time Constraints

Job Arrival Time

Other Characteristics

Job Duration

Job Rate

Price

**SLA Workload Framework**

**Generator**

SLA ID: 0001
SLA Birth: 12:00
SLO: job duration = 3:30h
SLO: CPUCount = 16
SLO: time constraints =
 {12:00-21:00}
SLO: Price = £1

SLA ID: 0002
SLA Birth: 12:01
SLO: …
SLO: …
SLO: …
SLO: …

SLA ID: 0003
SLA Birth: 12:01
SLO: …
SLO: …
SLO: …
SLO: …

…

## Modelling Variation in User Demand: Related Work

Active Time

Quiet Time

Calzarossa & Serrazi Model (1985) models the daily cycle of job rate, with peaks in the morning and afternoon, and a drop at lunchtime.

Extended to include a night time, currently *const*, but can be modelled by any function

The Feitelson Model (1996):

**Degrees of Parallelism:** harmonic distribution of order 1.5, with the most popular sizes of powers of two.
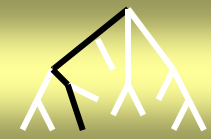**Repeated Executions:** the same job is likely to be executed again, $n^{-2.5}$, where n is the job duration.
**Arrival Process:** used in this model is Poisson.
**Correlation of Runtime with Parallelism**

Modelling SLA Constraints: Paper in Preparation

CPU
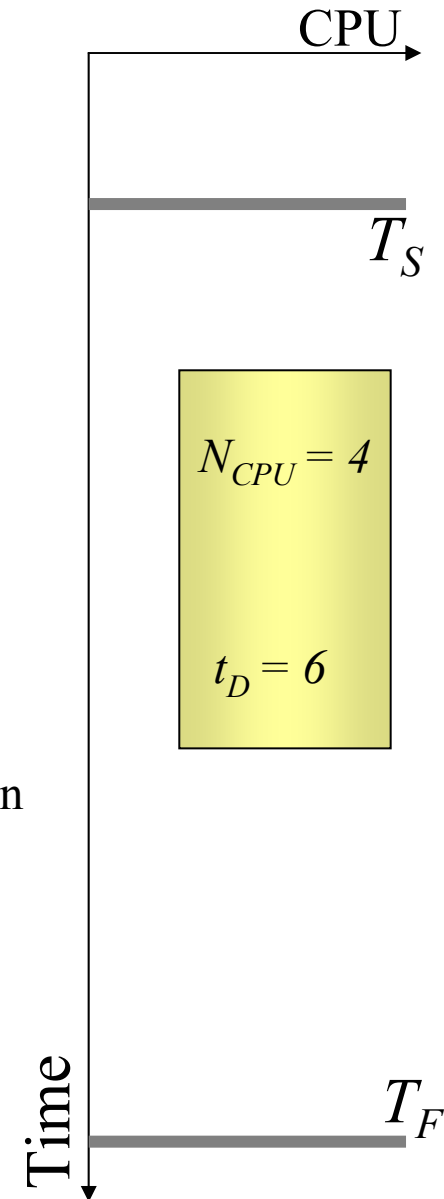
The constraint is represented as tightness:
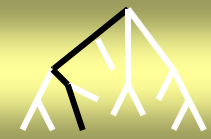
$$t_T = \frac{t_D}{T_F - T_S}$$

The distribution of tightness in the workload is described by a specific function $f_{TT}()$

Current state of grid is such that this function is a discrete distribution with non-zero values for:
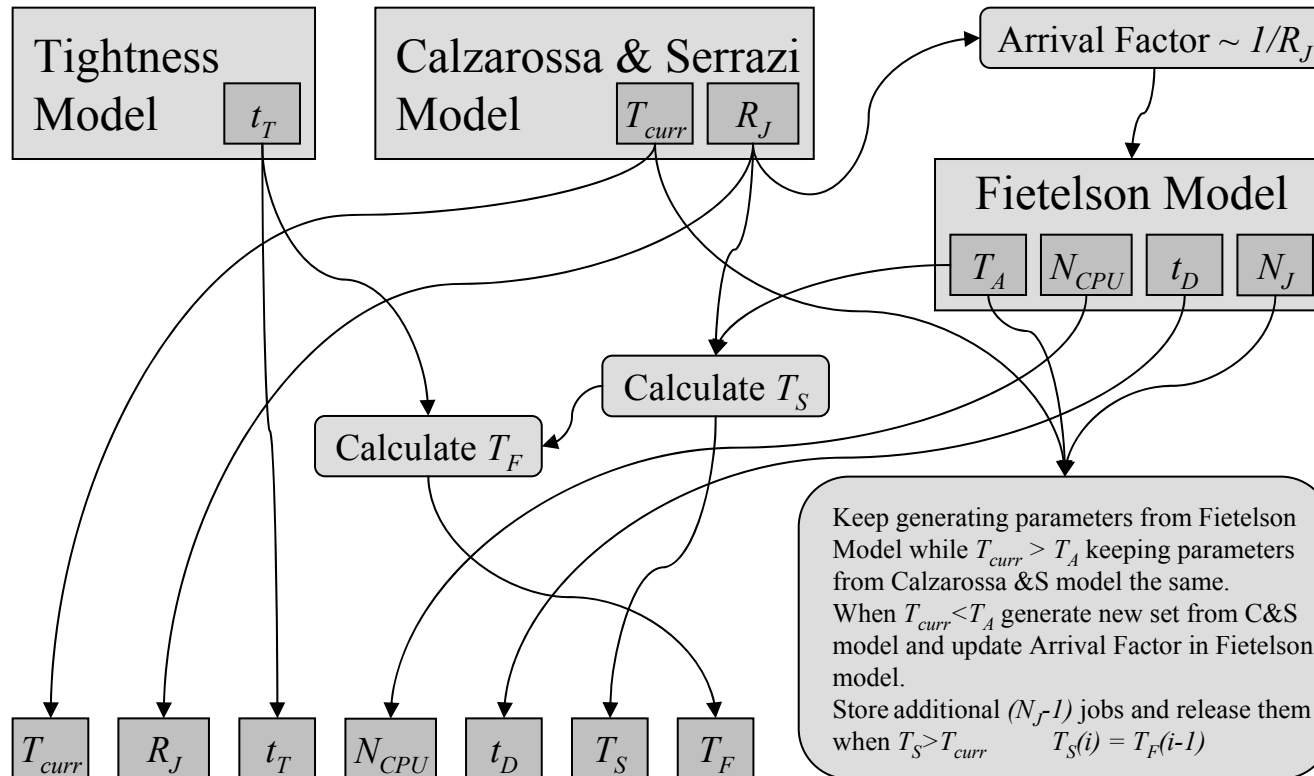
$t_T = 0$ – no constraints, deadlines

$t_T = 1$ – Advance Reservation

$T_S$

$N_{CPU} = 4$

$t_D = 6$

Time

$T_F$

# Workload Generator Framework: Example (paper in preparation)



WebApp Demo available from http://www.gridscheduling.org
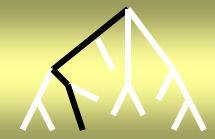
Workload Generator Framework: What Metrics (paper in preparation)



Example of a new metric: Load Density of the SLA workload
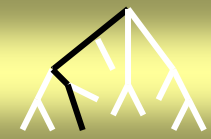
## Workload Generator Framework: Output

- Dynamic SLAs for DAGs
- Dynamic SLAs for dynamic workflows
- SLA workload models and generator
- Scheduling Heuristics
- Dynamic SLAs for market models
- SLAs with notion of trust

Thanks!

In preparation: Viktor Yarmolenko, Rizos Sakellariou, "A Framework for SLA Workload Generation"

In preparation: Viktor Yarmolenko, Rizos Sakellariou, "An SLA Workload Model for Time Constraints"

Rizos Sakellariou, Viktor Yarmolenko, "Job Scheduling on the Grid: Towards SLA-Based Scheduling", in L. Grandinetti (ed.), *High Performance Computing and Grids in Action,* IOS Press, 2008

Viktor Yarmolenko, Rizos Sakellariou, "Towards Increased Expressiveness in Service Level Agreements", *Concurrency and Computation: Practice and Experience,* vol.19, 1975-1990 (2007)

Viktor Yarmolenko, Rizos Sakellariou, "An Evaluation of Heuristics for SLA Based Parallel Job Scheduling", *Proceedings of the 3rd High Performance Grid Computing Workshop (HPGC)* (in conjunction with IPDPS 2006), Rhodes, Greece (April 2006), IEEE Computer Society Press

Viktor Yarmolenko, Rizos Sakellariou, Djamila Ouelhadj, Jonathan M Garibaldi, "SLA Based Job Scheduling: A Case Study on Policies for Negotiation With Resources", *Proceedings of the All Hands Meeting AHM'05*, Nottingham, UK (September 2005)

Keep checking *http://www.gridscheduling.org*