# The Extension of Torque Scheduler Allowing the Use of Planning and Optimization Algorithms in Grids

**Václav Chlumský,** Dalibor Klusáček, Miroslav Ruda

`vchlumsky@cesnet.cz`

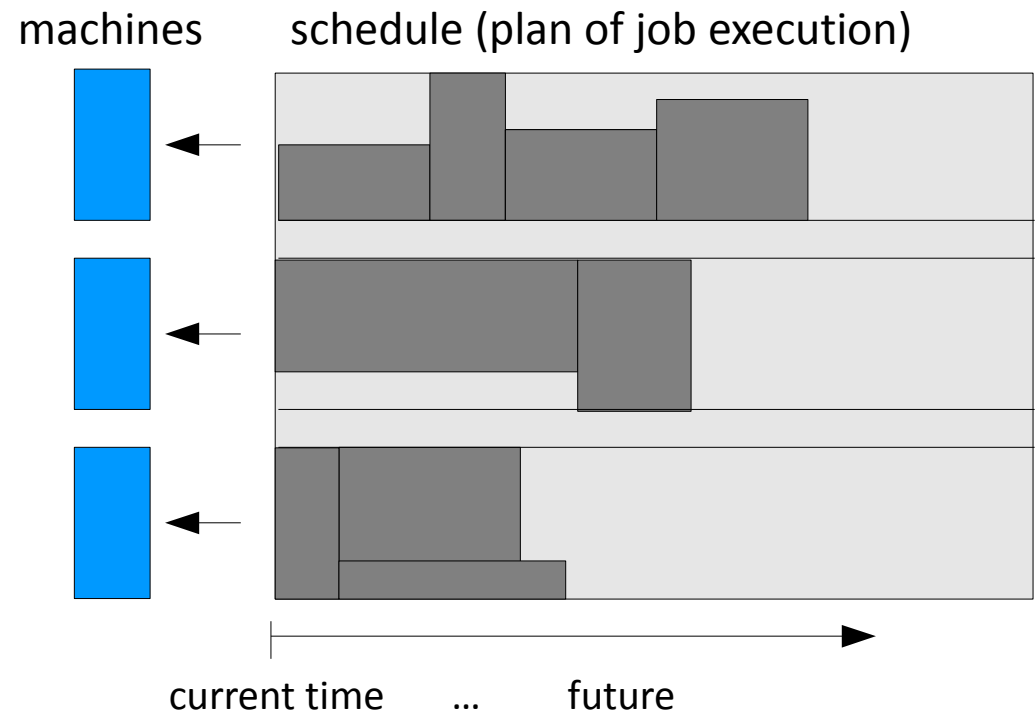CESNET z.s.p.o., Czech Republic

# Motivation

- **Efficient job scheduling in Grids is difficult task**

  - large, heterogeneous and dynamic system

- **State of the art scheduling approaches**

  - queue-based

  - robust and tolerable to dynamic environment

  - ad hoc decisions (at the very last moment)

  - limited predictability, planning and "self evaluation"

  - hard to use them when several objectives (goals) are to be followed simultaneously

# Schedule-based Solutions

- **Schedule represents plan of job execution**

  - straightforward planning

  - allows to predict behavior

  - easy **evaluation** wrt. selected optimization criteria

    - evaluation detects "problems"

  - **optimization of schedule**

    - to fix the problem

    - to improve the quality

    - using advanced scheduling techniques such as (meta)heuristics

machines     schedule (plan of job execution)
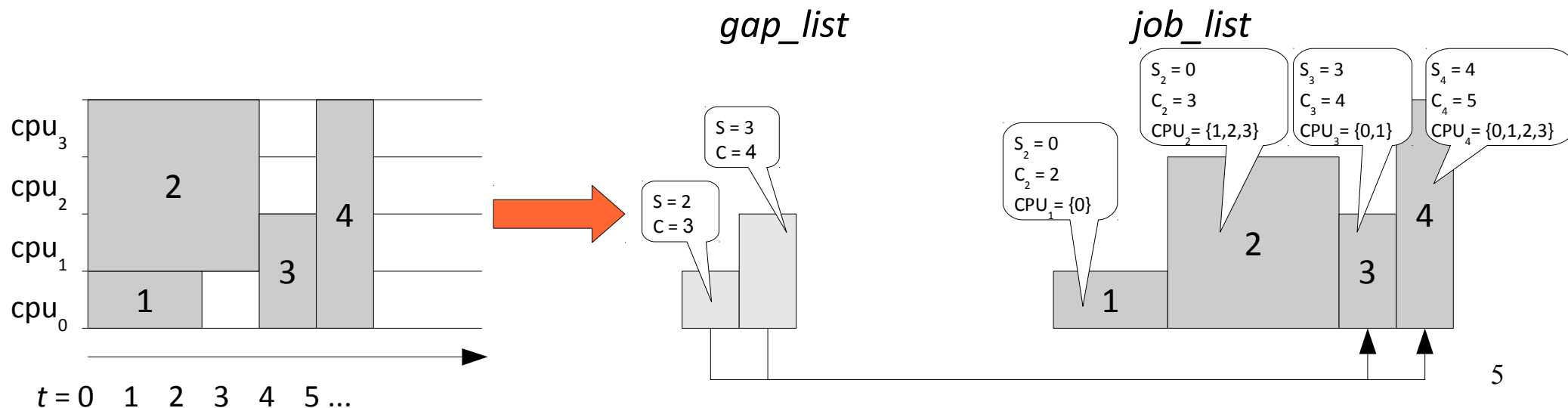
current time     ...     future

# Main Contribution of Our Work

- **Implementation of schedule-based solution in Torque Resource Manager**

- Main features

  - schedule data structure in the `pbs_sched` module
    - planning and prediction
    - when and where jobs will be executed

  - the use of **optimization algorithms**
    - subject to selected optimization criteria
    - **local search-based methods**
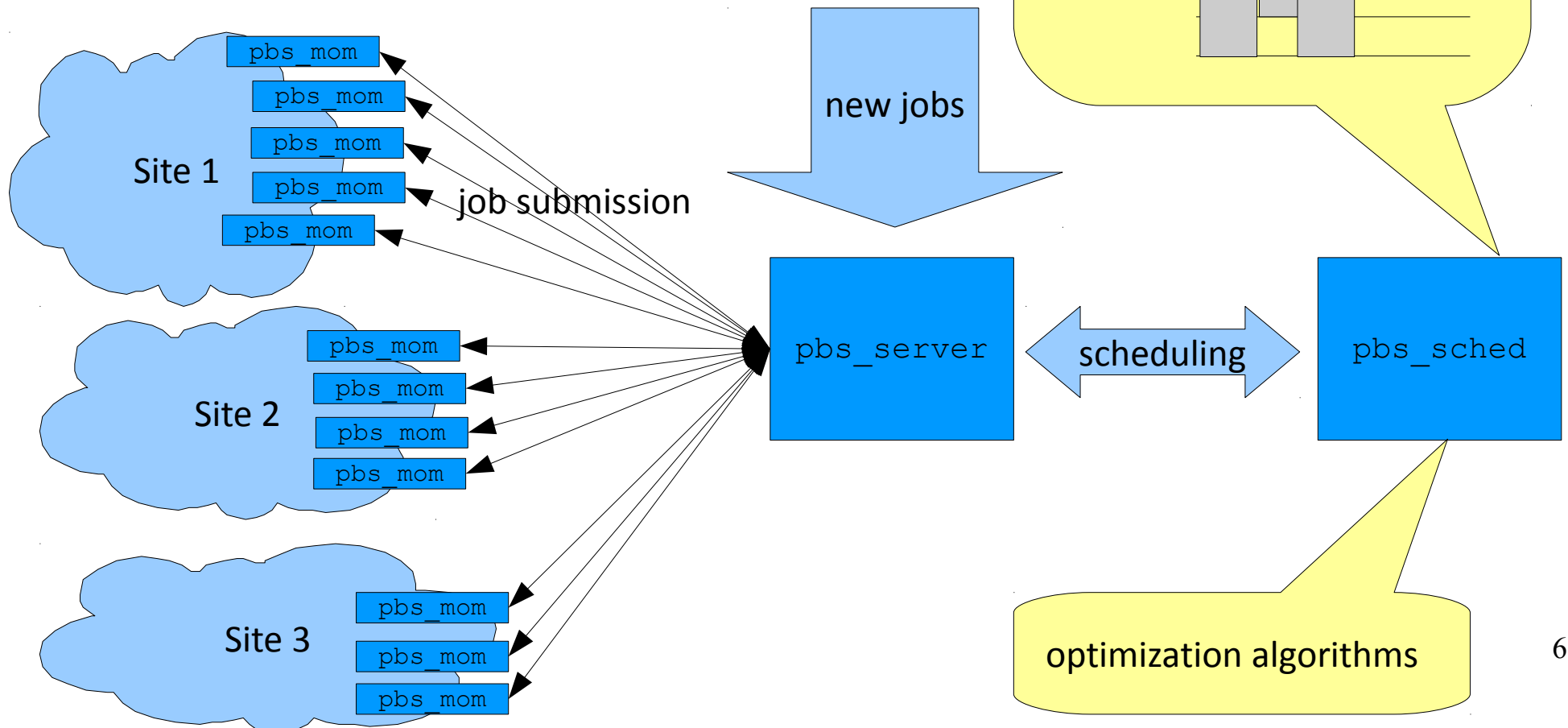    - periodically improving the quality of schedule

# Schedule Representation

- *job_list*: a list of objects that represent job-to-machine mapping in time

- *gap_list*: a list of gaps that represent unused CPU time slots

  - speeds up common operations

    - e.g., when trying to backfill some job

    - only the gap list is traversed

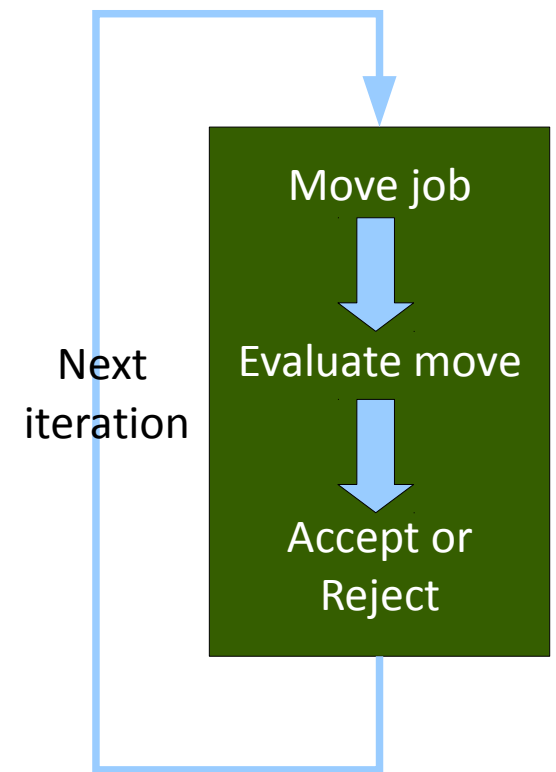    - start times of previous jobs remain guaranteed



*gap_list*

*job_list*

$S = 3$
$C = 4$

$S = 2$
$C = 3$

$S_2 = 0$
$C_2 = 2$
$CPU_1 = \{0\}$

$S_2 = 0$
$C_2 = 3$
$CPU_2 = \{1,2,3\}$

$S_3 = 3$
$C_3 = 4$
$CPU_3 = \{0,1\}$

$S_4 = 4$
$C_4 = 5$
$CPU_4 = \{0,1,2,3\}$

cpu$_3$ cpu$_2$ cpu$_1$ cpu$_0$

$t = 0$  1  2  3  4  5 ...

5

# Torque Resource Manager

- Modification of `pbs_sched` module

  - **schedule structure**

  - **optimization algorithms**

  - related methods

job schedule

pbs_mom
pbs_mom
pbs_mom
pbs_mom

Site 1

pbs_mom

job submission

new jobs

pbs_mom
pbs_mom
pbs_mom
pbs_mom

Site 2

pbs_server

scheduling

pbs_sched

pbs_mom
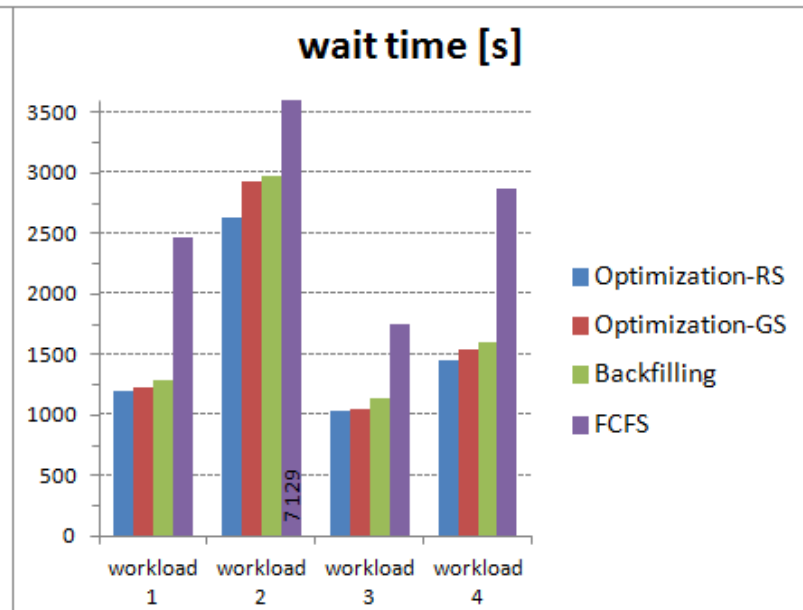pbs_mom

Site 3
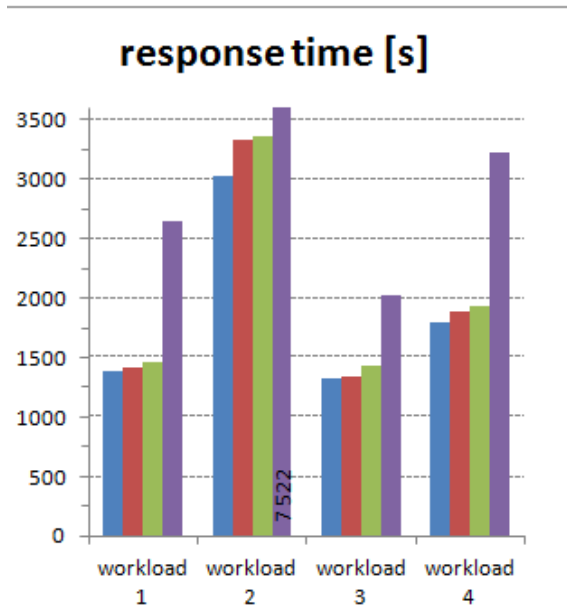
pbs_mom

optimization algorithms

6

# Preliminary Results

- 4 problem instances using real workloads from Czech NGI MetaCentrum

- **2 iterative Local Search-based optimization algorithms**

  - "conservative" gap-search algorithm

    - move random jobs into existing gaps

  - random-search algorithm

    - move random jobs into random positions

  - applied objective – **minimize avg. slowdown**

  - measured criteria

    - avg. slowdown

    - avg. response time

    - avg. wait time

- **Compared wrt. FCFS and Backfilling**

Next iteration → Move job → Evaluate move → Accept or Reject

# Preliminary Results

# Current and Future Work

- Further testing

  - implementation of multi-criteria objective function

- Fairness related problems

  - *job-to-job* and *user-to-user* fairness

  - proposal and integration of proper objective functions

- Deployment of such a solution in a production environment

- The use of runtime-prediction techniques

  - when estimates are very inaccurate/missing

- Preparation of a publicly available software package

  - Torque with the schedule-based extension of `pbs_sched` module