

A Reference Architecture for Datacenter Scheduling

Design, Validation, and Experiments



Georgios Andreadis, Laurens Versluis,
Fabian Mastenbroek, Alexandru Iosup

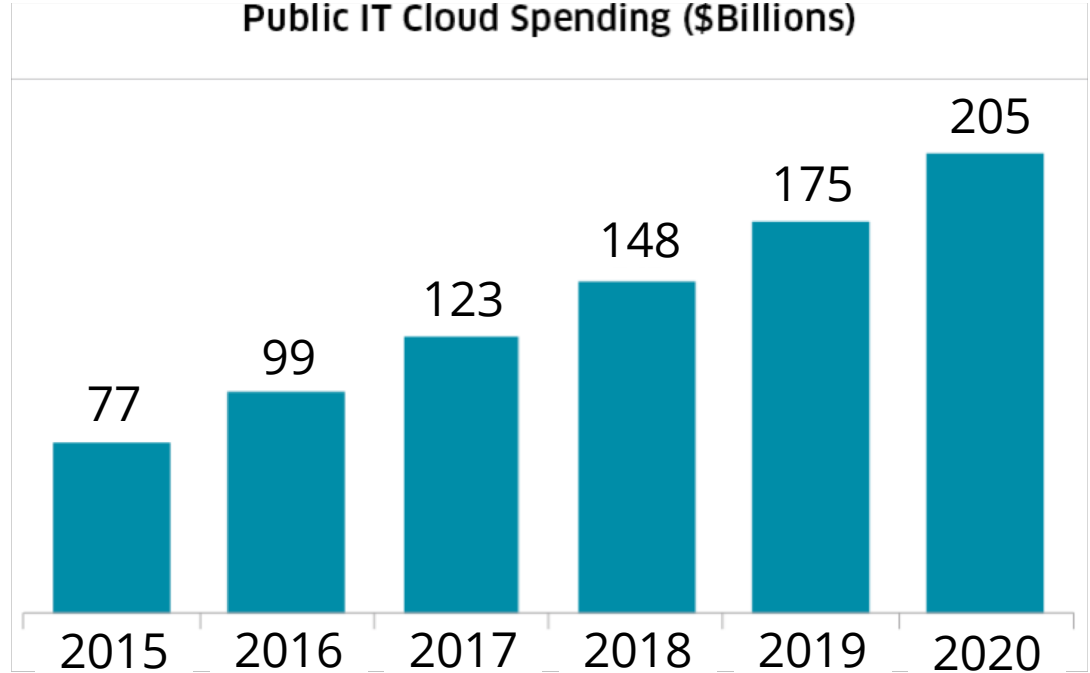
g.andreadis@atlarge-research.com

Delft University of Technology, Vrije Universiteit Amsterdam

Delft > The Netherlands > Europe

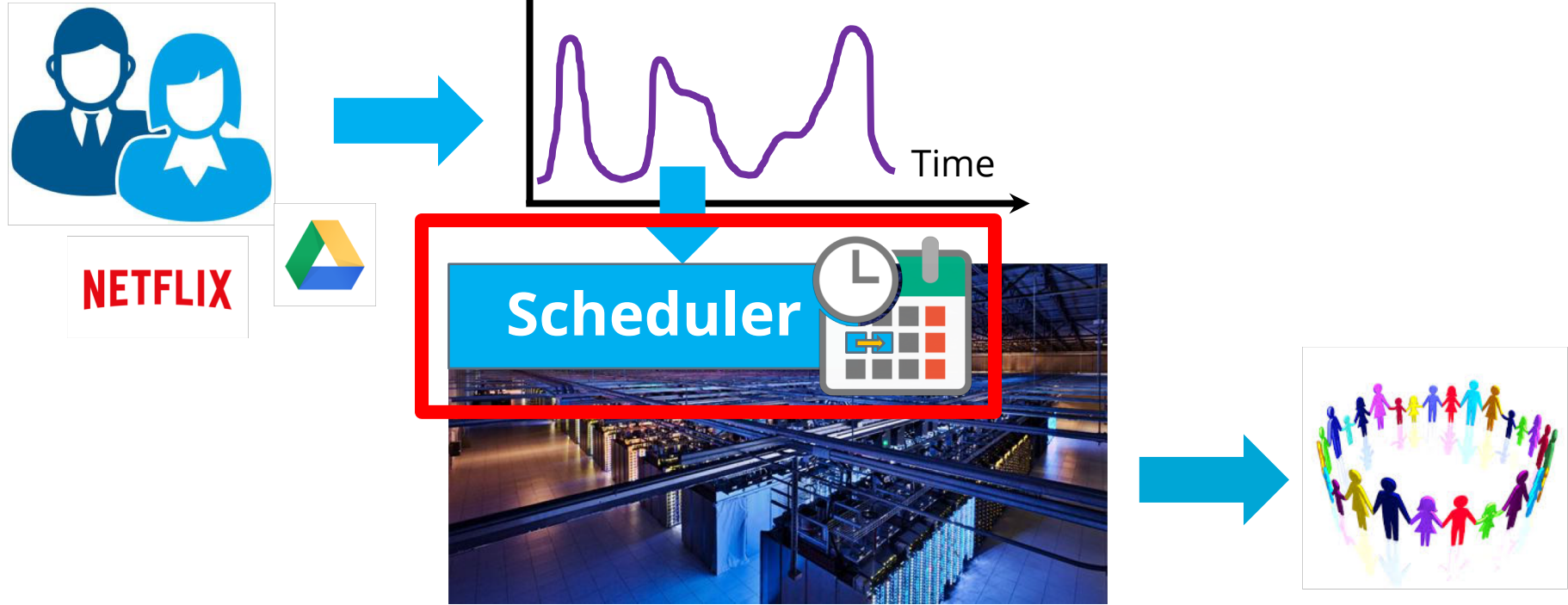


Cloud Infrastructure is Vital to Digital Society



Source: <https://business.nasdaq.com/marketinsite/2017/Cloud-Computing-Industry-Report-and-Investment-Case.html>, 2018-11-02

Challenge: Manage Resources Efficiently & Fully Automated



Scheduling is Hard

“30—70% scheduler decisions
incorrect in datacenters”

Source: IEEE Computer '15

“current schedulers not efficient
for many users, diverse services”

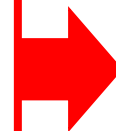
Source: Dutch industry,
CCGRID '15

“new schedulers not used in
datacenters, fear of failure”

Source: Euro-Par '16



Need **Smarter**
Schedulers



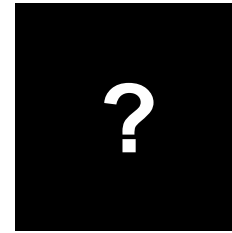
Need Scheduler
Reproducibility

Previous Approaches vs. Our Solution

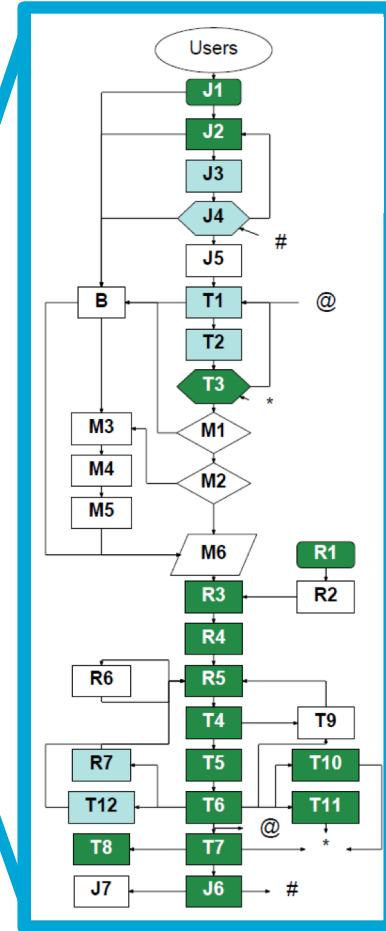
- Current models stay at black-box level
- Scheduler do many things
- Difficult to understand and compare
- Need for a common language

In this work:

1. Design a **reference architecture**
2. **Map** existing schedulers
3. Conduct **experiments**



Scheduler



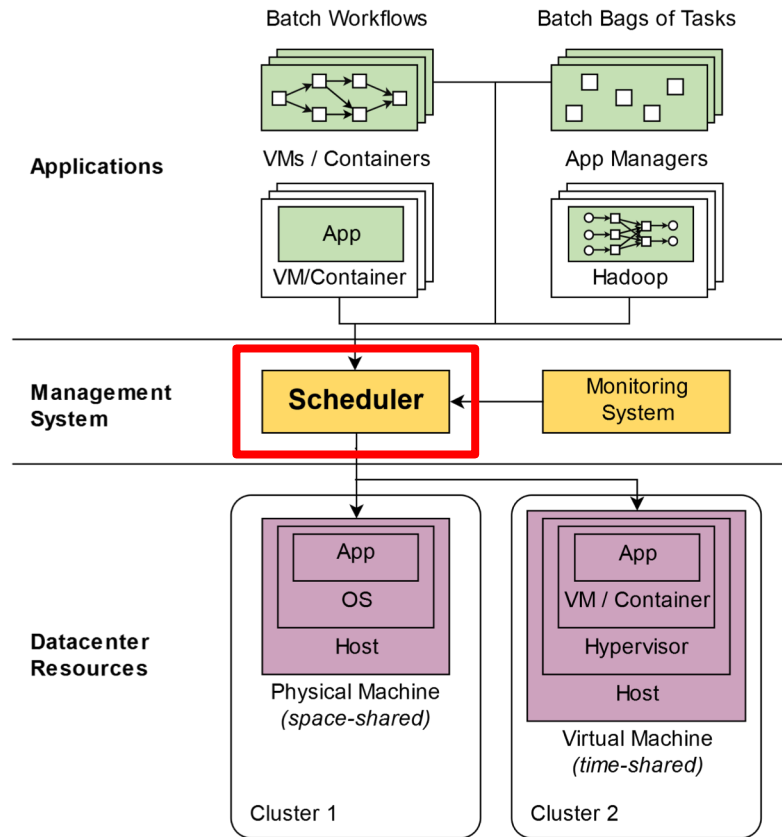
Closely Related Work

- Scheduling architectures
 - Schopf (10 steps when grid-scheduling): Main inspiration
 - Surveys (Rodriguez et al., Singh et al.): More coarse-grained and lacking features
- Architectures of cloud systems
 - Complements NIST Cloud and Big Data RAs
 - Fits in the ISO/IEC/IEEE 42010:2011 standard
- Large-scale software architectures
 - Rozanski et al., Bass et al.: Theory of Software Architectures



System Model

- Workflow-compliant workloads
- Users can specify requirements
- Physical and virtual resources
- The scheduler...
 - Allocates
 - Provisions
 - Replicates
 - Migrates
 - ...



Design Validation Experiments

Requirements and Goals

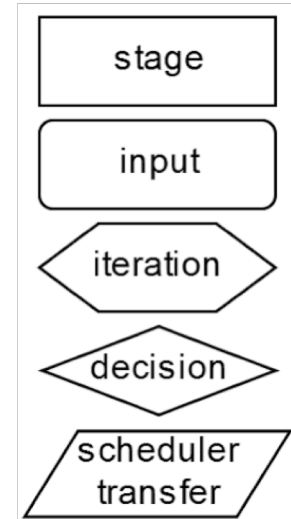
- Validity
 - Accurately represents the field
 - Verify through mapping and peer-review
- Usefulness
 - Its utility to stakeholders
 - Empirical results from mapping and experiments



Specification: Our Workflow-based Model for Datacenter Scheduling

- 33 stages with control and data flows between them
- 5 modes of operation
- Support hierarchy of schedulers
- Stages can be divided into 4 groups:
 - Job processing (**J**)
 - Task processing (**T**)
 - Scheduler management (**M**)
 - Resource management (**R**)

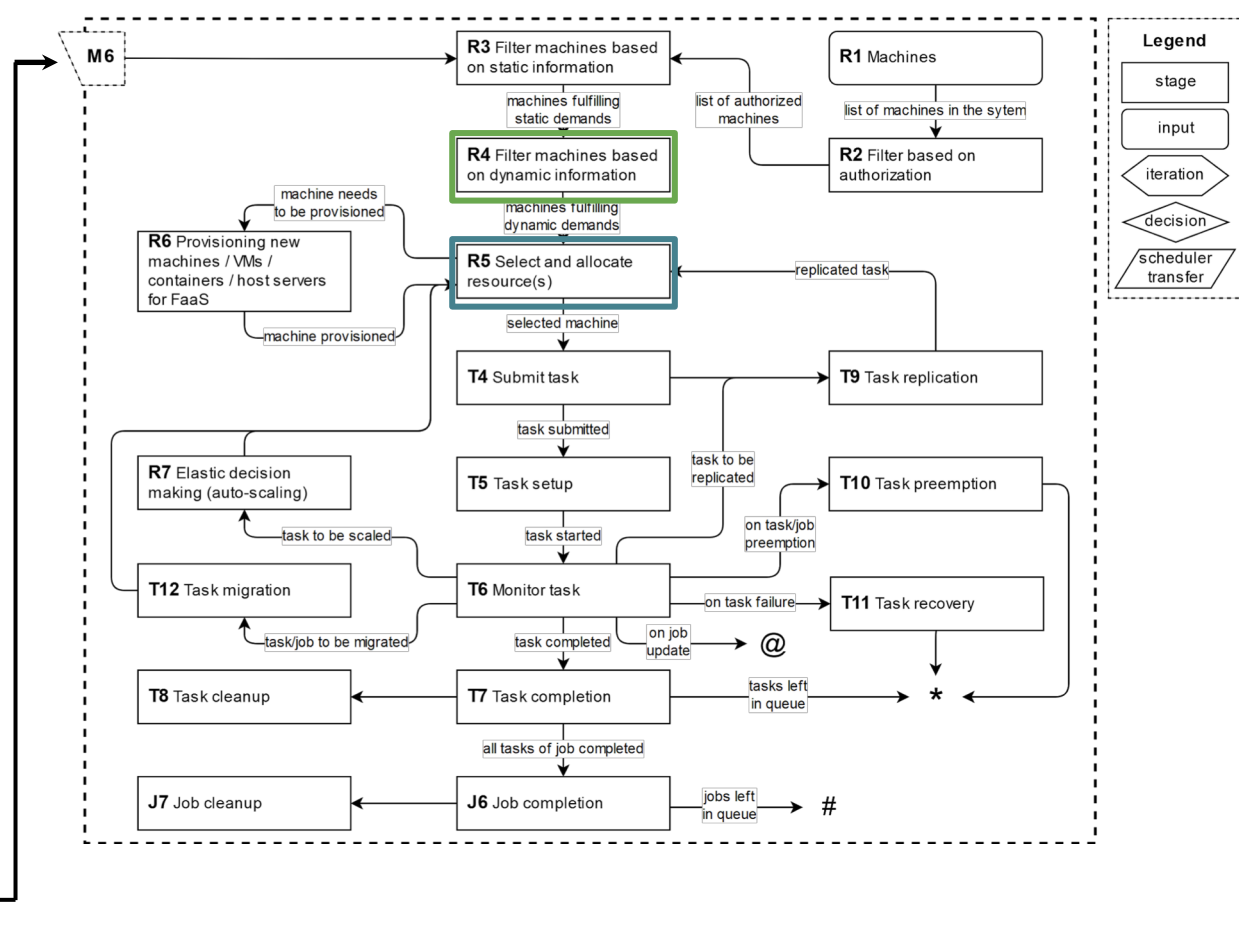
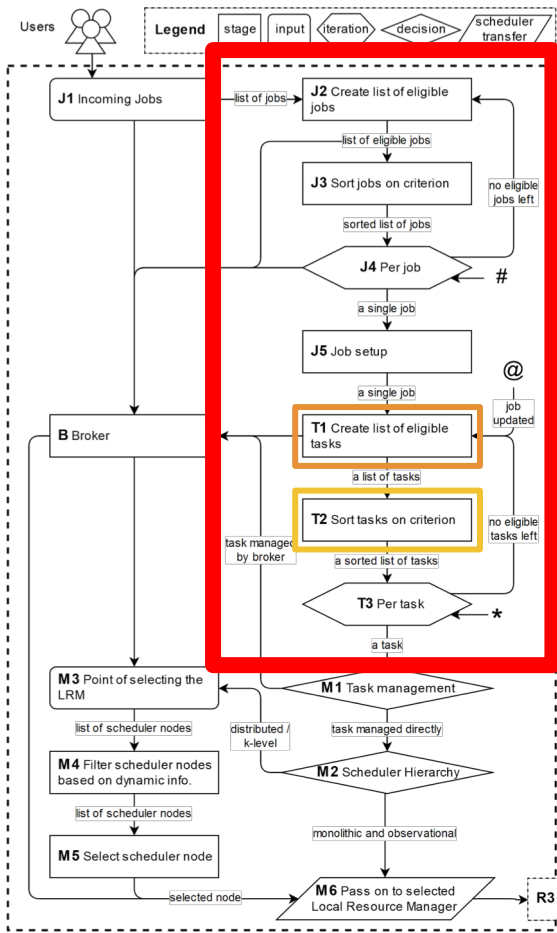
The 5 modes of operation
for scheduling stages

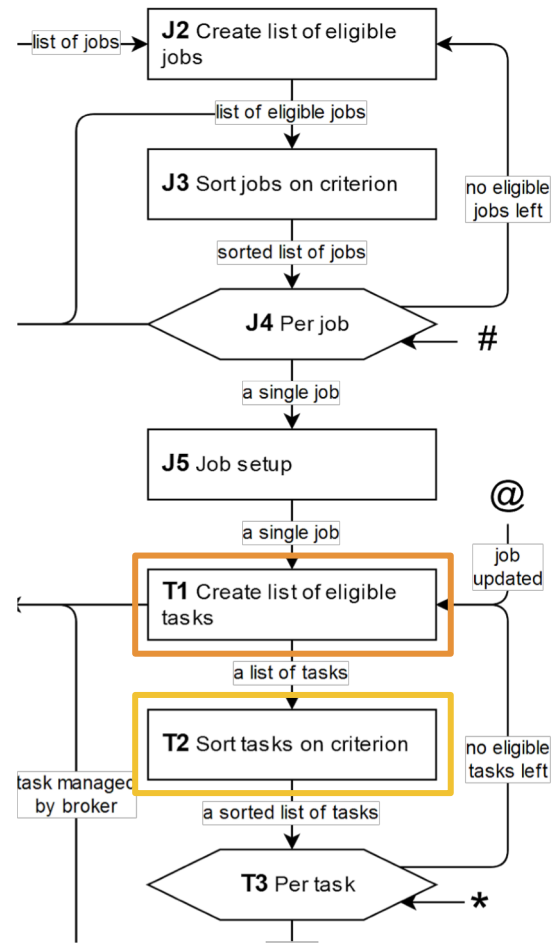
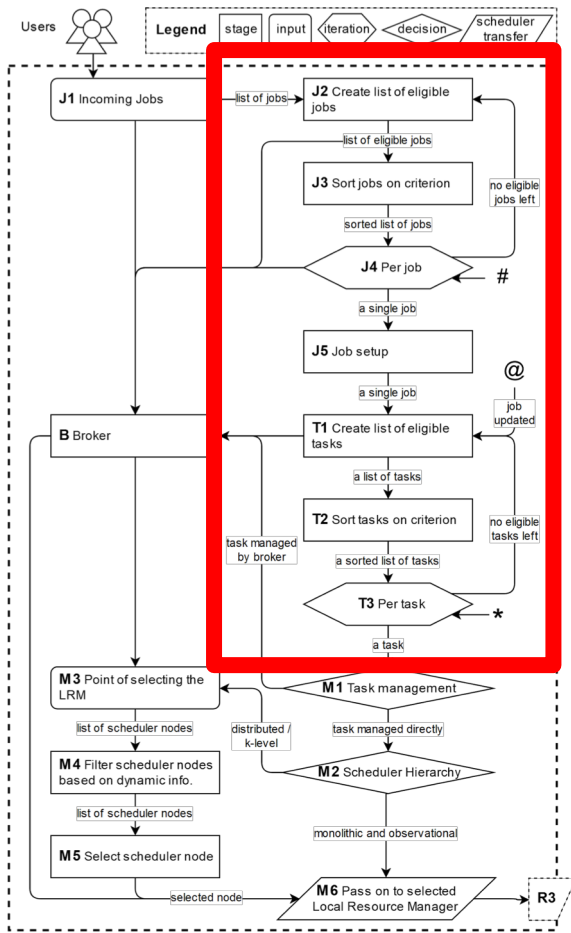


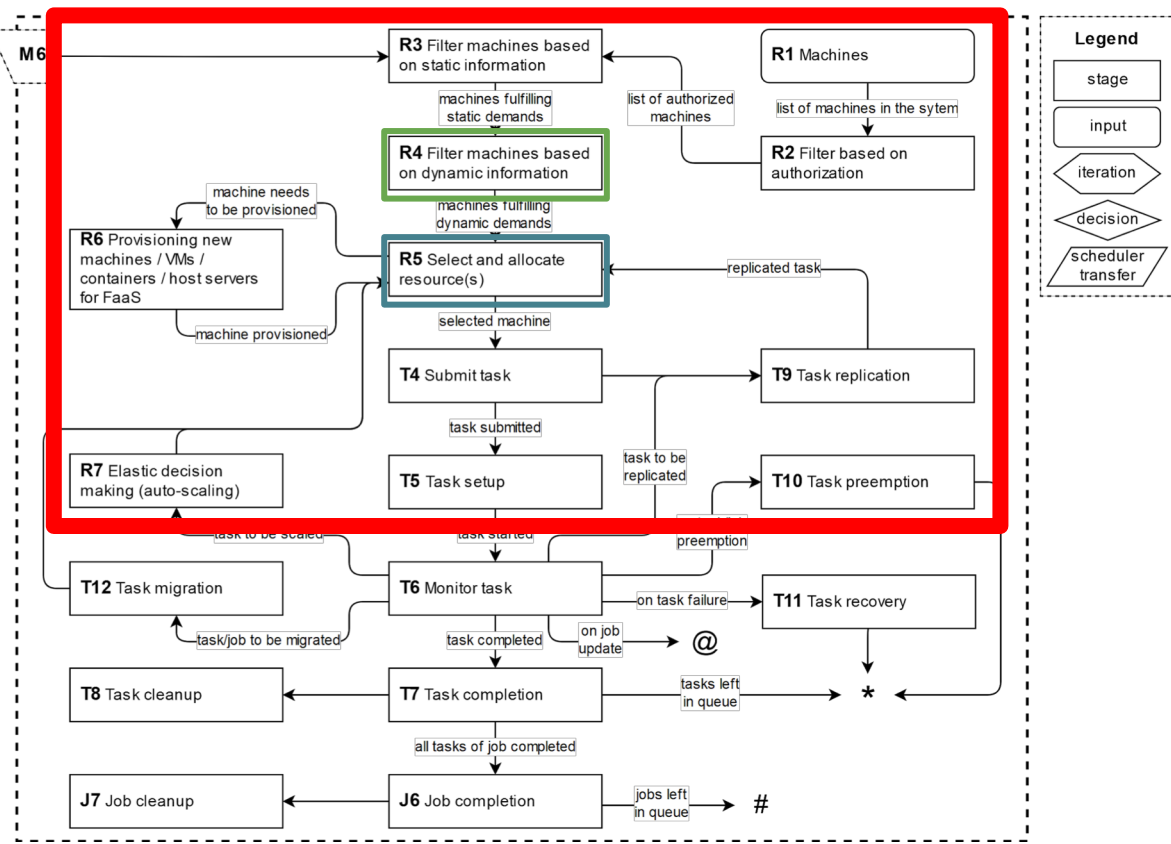
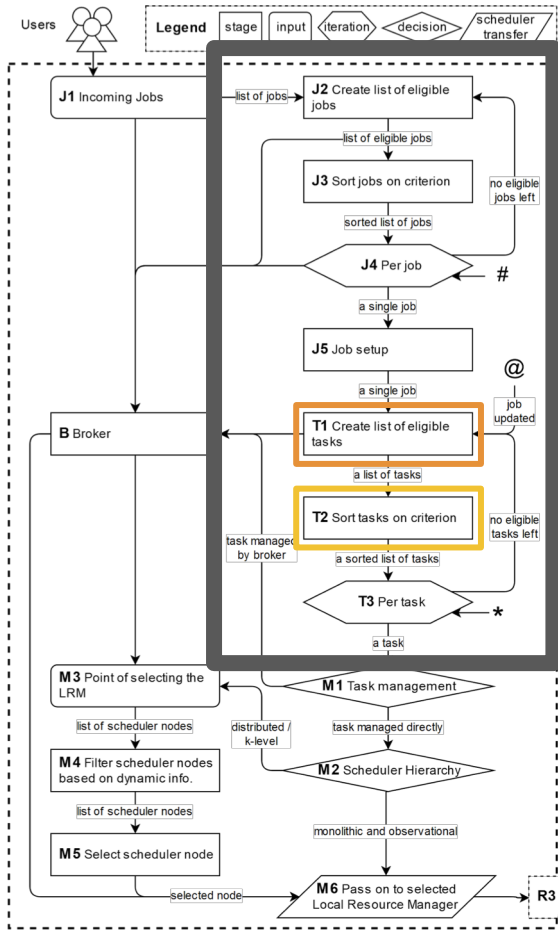
Process of Designing the Reference Architecture

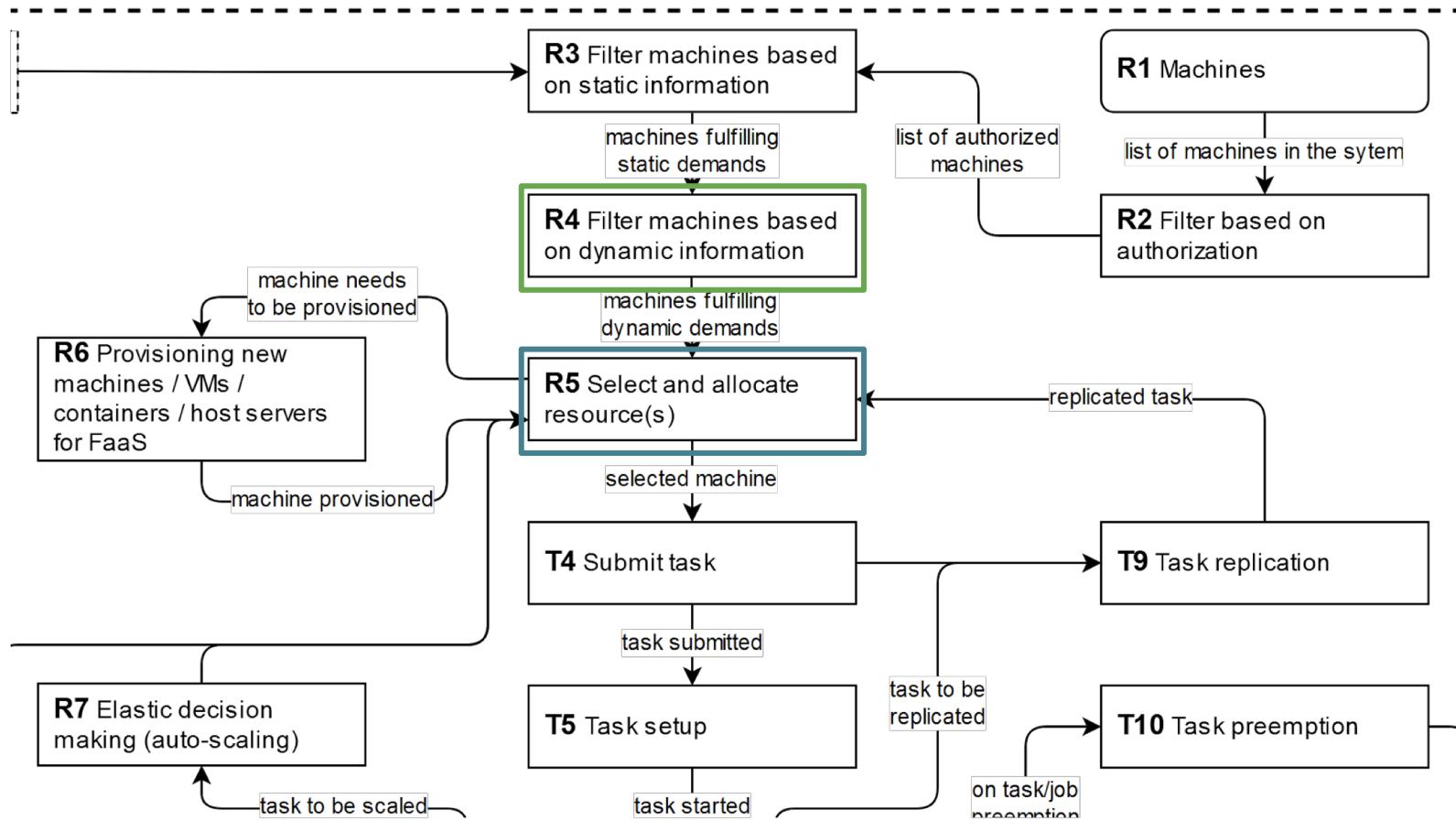
- Goal: Comprehensiveness
- Publications selected:
 - Scheduling core analyzed and presented
 - Highly cited or deployed in a large real-world environment
- Leading principles in this process:
 - Components with Clearly Distinct Responsibilities
 - Separation of Mechanism from Policy











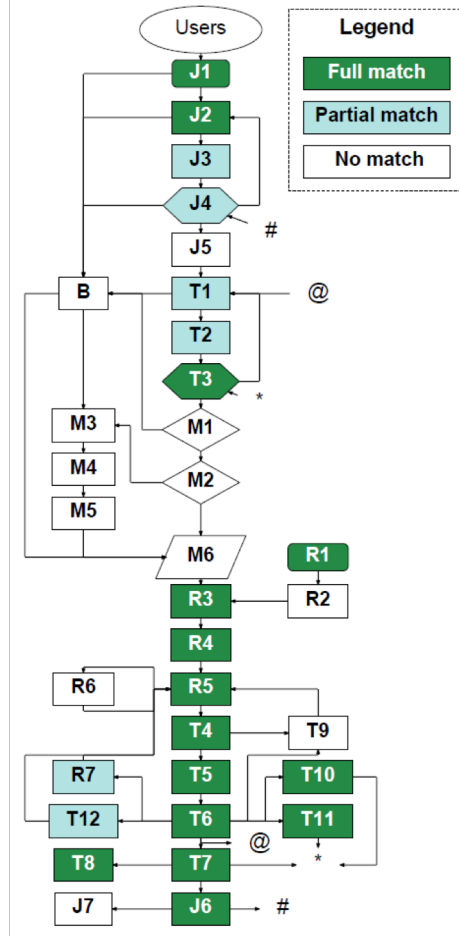
Design Validation Experiments

Mapping Schedulers to the Model

- 14 scheduler publications analyzed
- For each stage: 3 possible types of a 'match'
 - Full match:** Stage is described in detail
 - Partial match:** Detail on the 'how' is lacking
 - No match:** Remaining cases **Underspecification**
- Here, a Borg-like scheduler, given no open-source

💡 Overall: Schedulers tend to be **underspecified**

Mapping of the Borg Scheduler to the RA



More on the Mapping Process

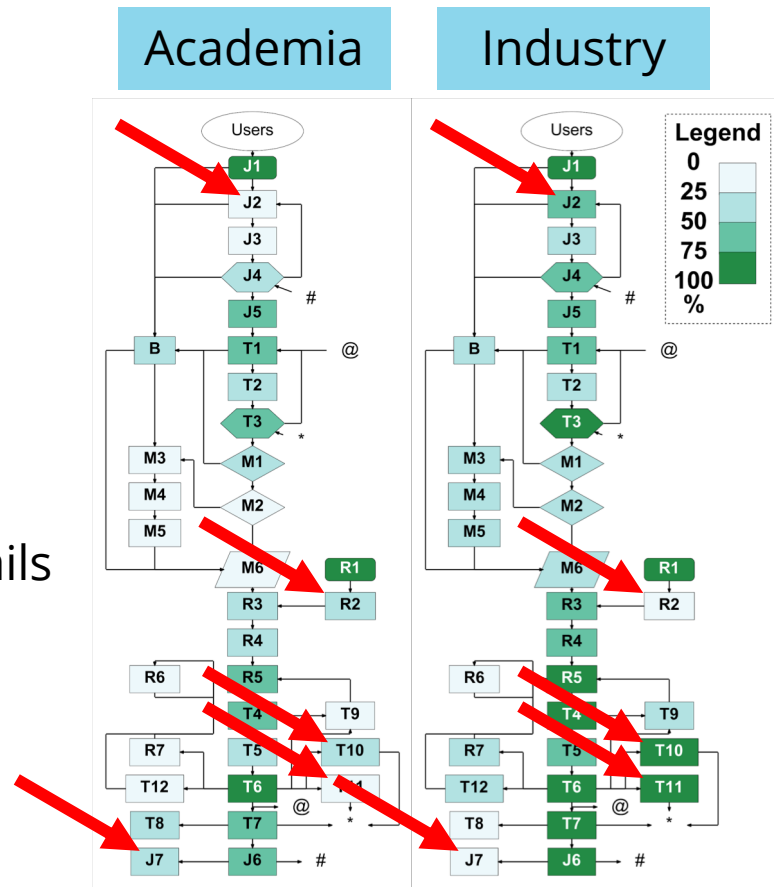
- Subset of the publications used for design
- Mapping done in parallel by two reviewers
- Small set of publications used as calibration
 - No significant difference in interpretation

Discussion on mapping schedulers with
different levels of expertise:
bit.ly/sc18-scheduling



Head-to-Head: Academia vs. Industry

- Comparison through **heat maps**
- Specified more by industry
 - T11: Task recovery
 - T10: Task preemption
 - J2: Job filtering
 - → More attention to the technical details
- Specified more by academia
 - R2: Filtering resources based on authorization
 - J7: Job cleanup

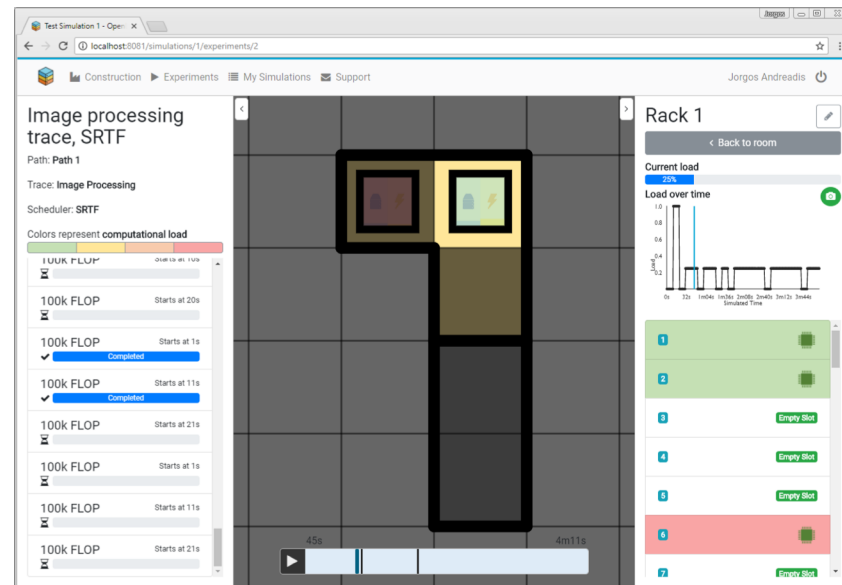


Design Validation Experiments

- Q1:** Impact of underspecification on performance?
Q2: Impact of underspecification on real-world performance?

Experimental Setup

- Prototype implemented in **OpenDC**
- Selected subset of stages made configurable
- Real-world engineering and industrial traces
- Standard DC topology



Impact of Underspecification on Performance

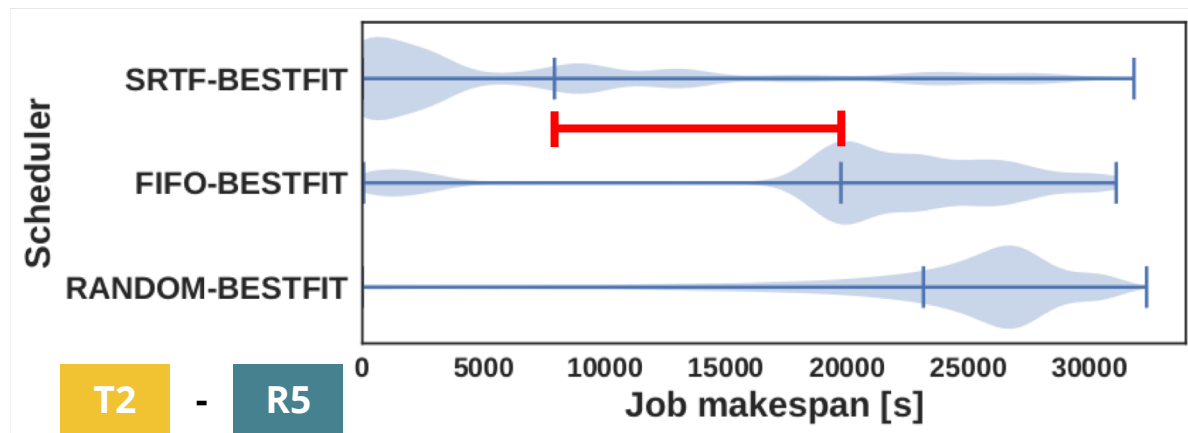
- Selected two underspecified stages in a Borg-like system
- Compare various metrics of different configurations

💡 Different policy combinations lead to different results

T2: Sort tasks

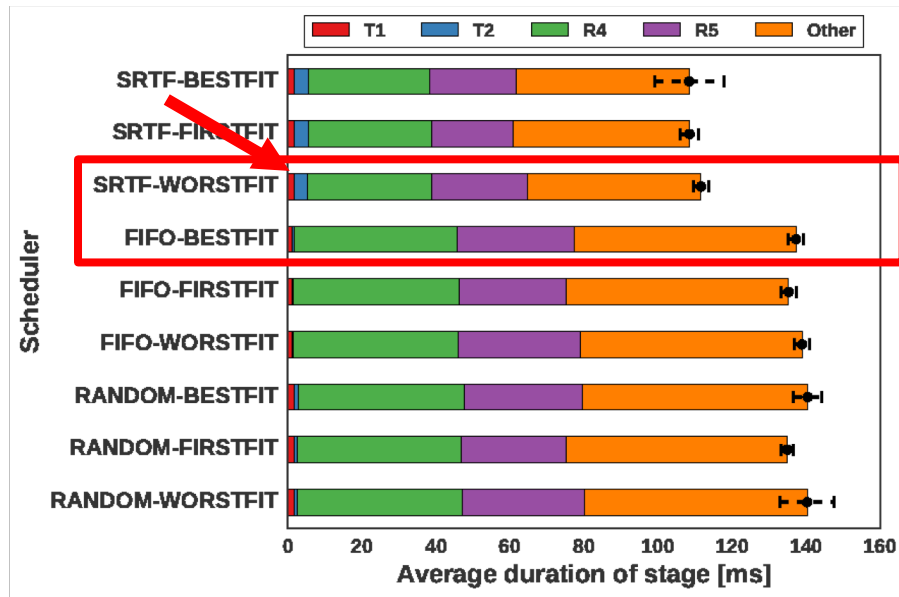
R5: Select resources

Significant difference in performance!



Stage Complexity at Runtime

- Instrumented simulation run
- Measure time spent in each stage
- Observation: SRTF leads to longer sorting times, but shortest total duration

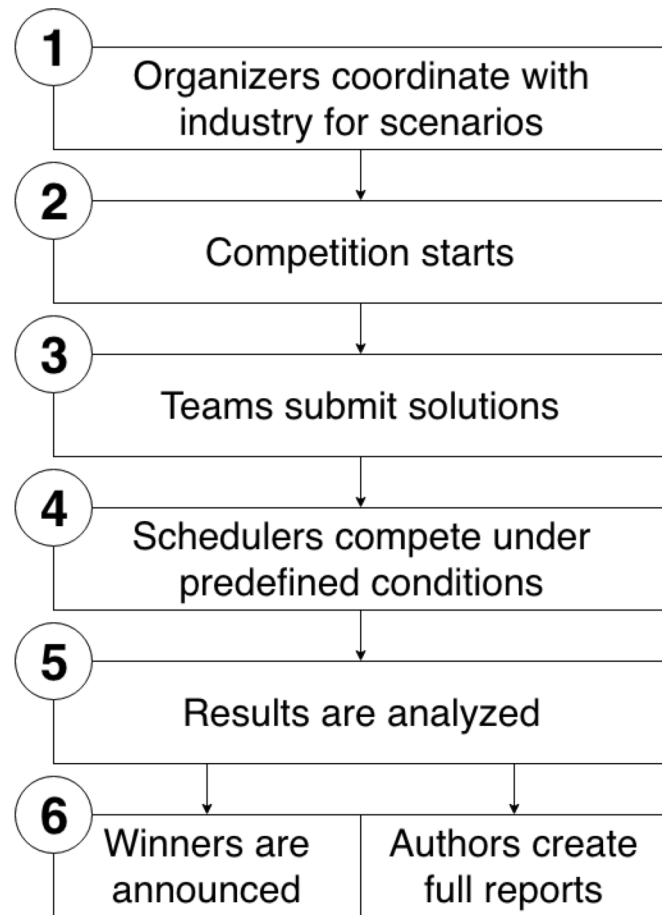


Stage-policy can have a non-trivial impact on stage durations

Future Directions

- Extended validation
- API/Language for scheduler design
- Dialogue with scheduler designers
- A global competition of schedulers
 - Future edition of JSSPP

We're collecting use cases
for the community:
bit.ly/jsspp-cfp



Speculation: Scheduler Design Space Exploration

- Systematic exploration of stage implementation combinations
- Using real-world stage policies
- “Bruteforcing” scheduler design

Will the results differ from human designs?

Take-Home Message

- Schedulers are complex and difficult to compare
- They remain underspecified
- Important implications for reproducibility, DevOps
- A guideline for new scheduler publications
- Conceptual model can help design, analyze, and improve schedulers

Learn more about OpenDC:
opendc.org

Read the Technical Report:
bit.ly/sc18-scheduling

Fill in the short survey:
bit.ly/sc18-scheduling-survey

Georgios Andreadis

g.andreadis@atlarge-research.com
atlarge-research.com/gandreadis

