

# Ergonomics and efficiency of workflows on HPC clusters

Jakub Beránek

Supervisor  
**Ing. Jan Martinovič, Ph.D.**



IT4INNOVATIONS  
NATIONAL SUPERCOMPUTING  
CENTER

# Ergonomics and efficiency of workflows on HPC clusters

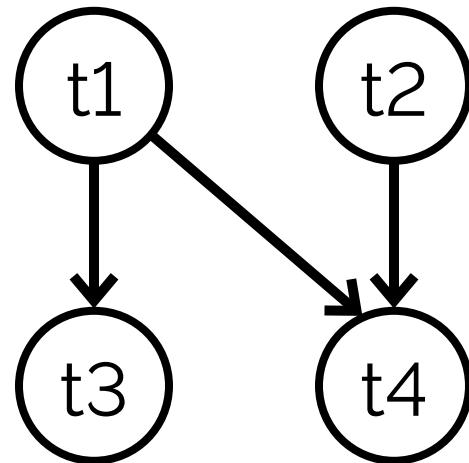
## **Goal**

Help users execute task workflows  
on HPC clusters in an easy & efficient way

# Ergonomics and efficiency of workflows on HPC clusters

## Goal

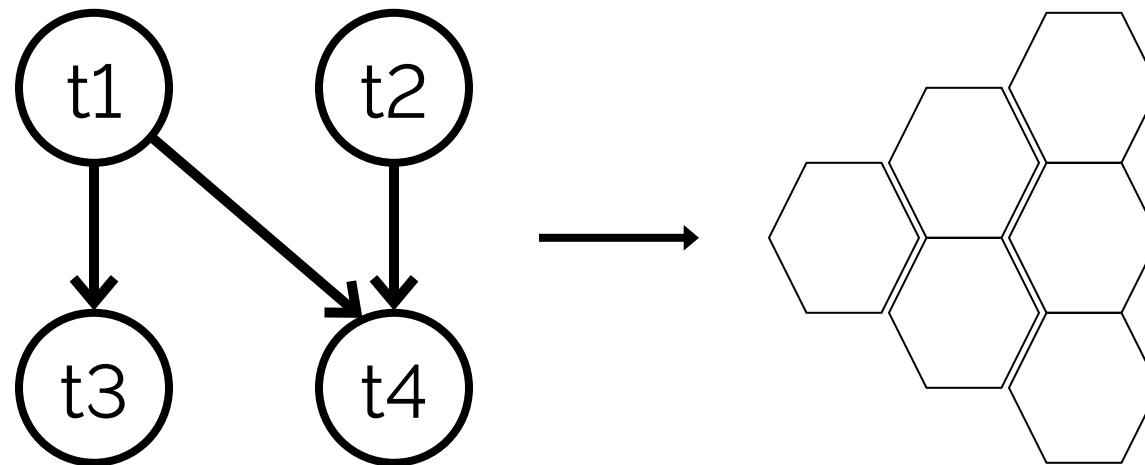
Help users execute task workflows  
on HPC clusters in an easy & efficient way



# Ergonomics and efficiency of workflows on HPC clusters

## Goal

Help users execute task workflows  
on HPC clusters in an easy & efficient way



# Objectives

1. Identify HPC workflow challenges

# Objectives

1. Identify HPC workflow challenges
2. Design approaches for overcoming them

# Objectives

1. Identify HPC workflow challenges
2. Design approaches for overcoming them
3. Implement them in a task runtime

# Objectives

1. Identify HPC workflow challenges
2. Design approaches for overcoming them
3. Implement them in a task runtime
4. Analyze results on real use-cases

# Motivation

Ergonomics and efficiency  
of **workflows on HPC clusters**

# Workflows

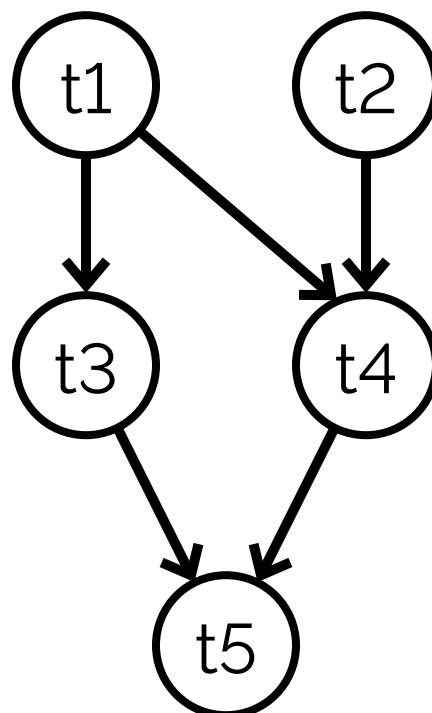
= task graphs, pipelines, task DAGs\*, ...

\* Directed Acyclic Graph

# Workflows

= task graphs, pipelines, task DAGs\*, ...

Popular programming model

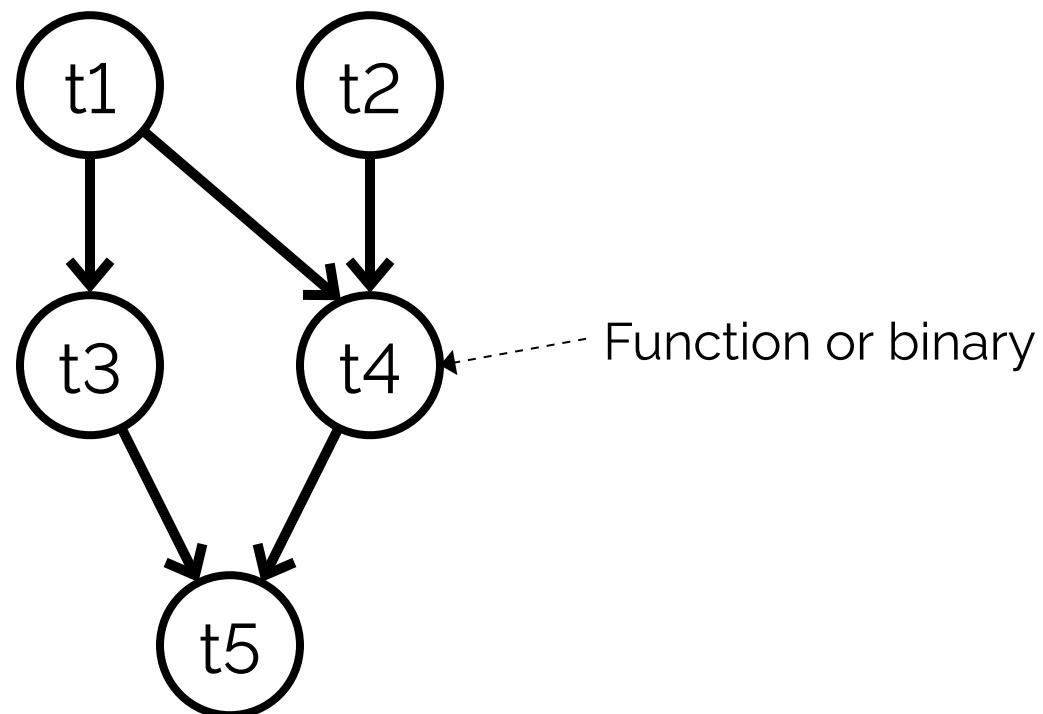


\* Directed Acyclic Graph

# Workflows

= task graphs, pipelines, task DAGs\*, ...

Popular programming model

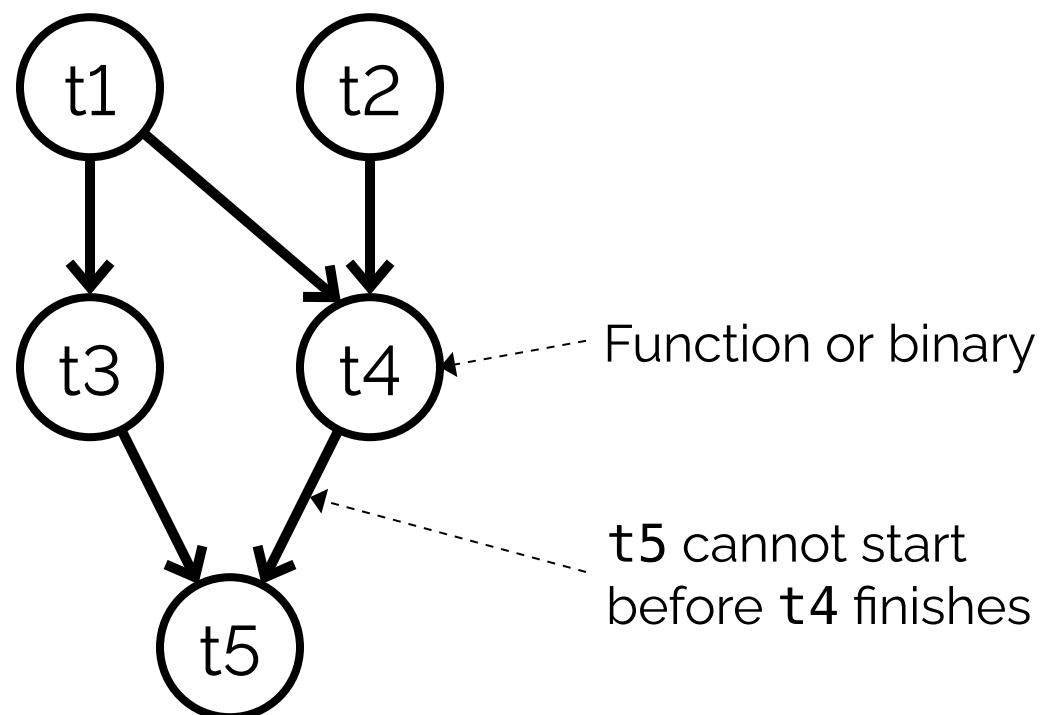


\* Directed Acyclic Graph

# Workflows

= task graphs, pipelines, task DAGs\*, ...

Popular programming model



\* Directed Acyclic Graph

# Main workflow benefits

- **Implicit parallelism**
  - Build a DAG vs. use MPI

# Main workflow benefits

- **Implicit parallelism**
  - Build a DAG vs. use MPI
  - Extracted by a task runtime

# Main workflow benefits

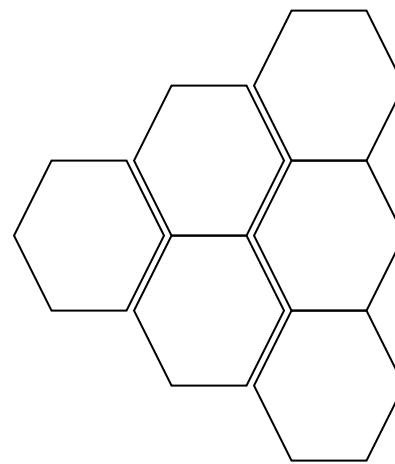
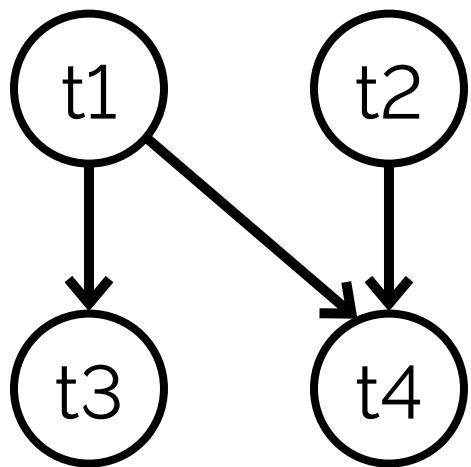
- **Implicit parallelism**
  - Build a DAG vs. use MPI
  - Extracted by a task runtime
- High-level description
  - Python/DSL vs. C/C++

# Main workflow benefits

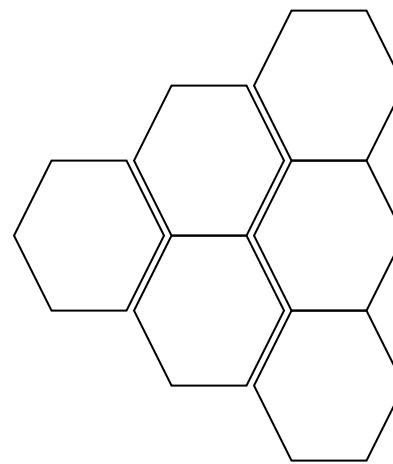
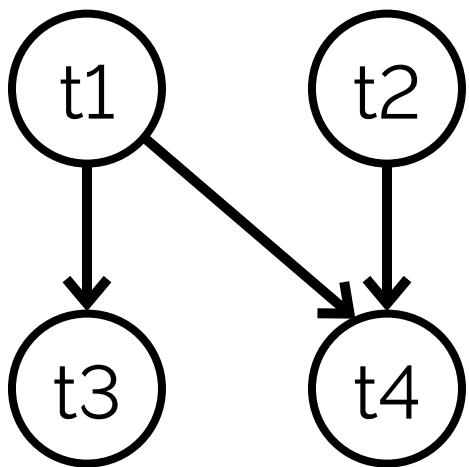
- **Implicit parallelism**
  - Build a DAG vs. use MPI
  - Extracted by a task runtime
- High-level description
  - Python/DSL vs. C/C++
- Portability
  - Use-cases, hardware

# Task graph challenges on HPC clusters

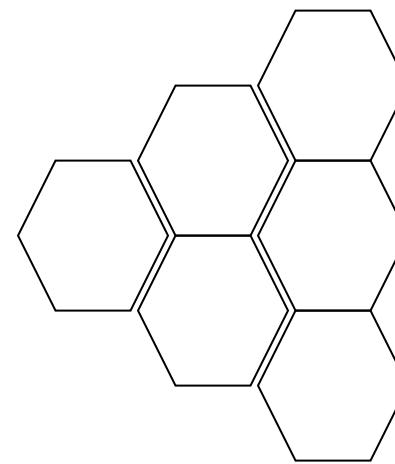
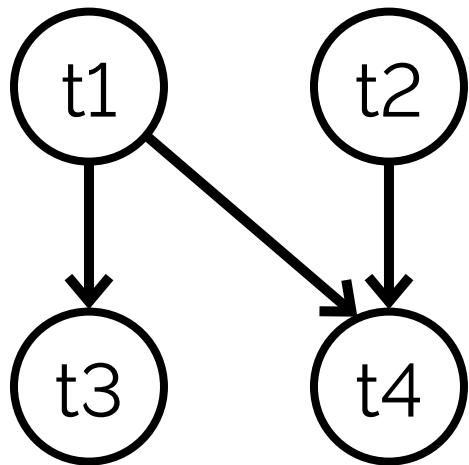
# Job manager



# Job manager



# Job manager

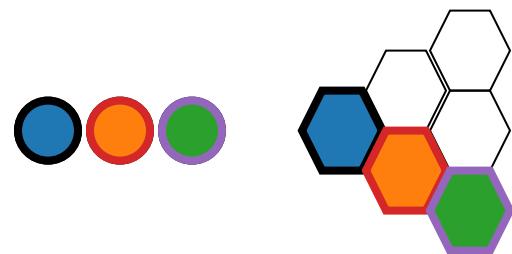


How to map tasks to PBS jobs?

# Submitting task graphs into PBS

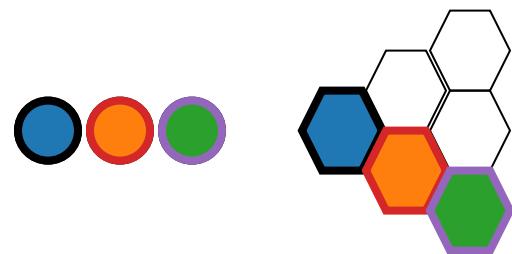
# Submitting task graphs into PBS

- Submit each task as a job



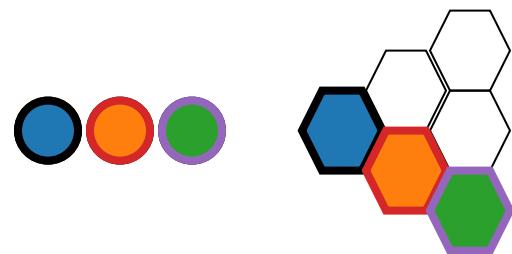
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)



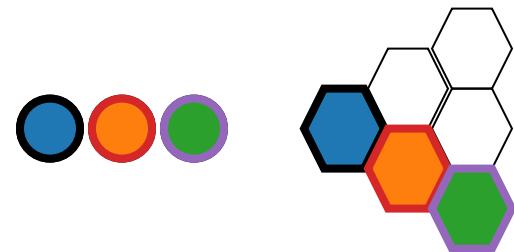
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits



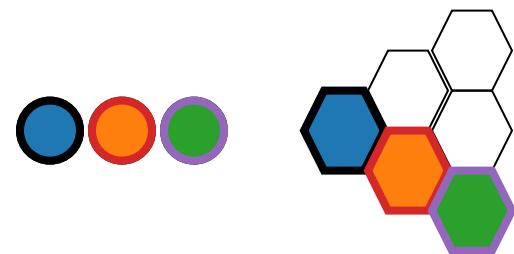
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits
  - **Node granularity**



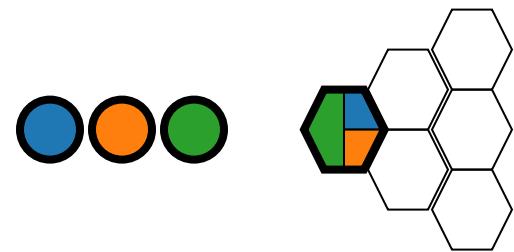
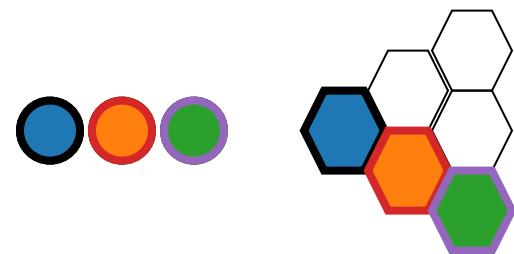
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits
  - **Node granularity**
  - Difficult with dependencies



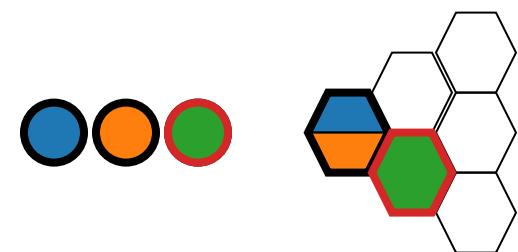
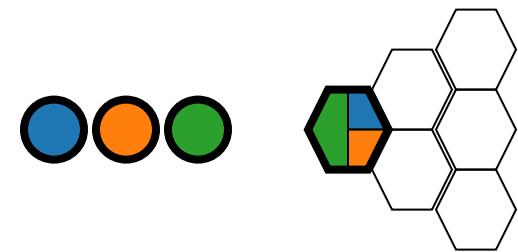
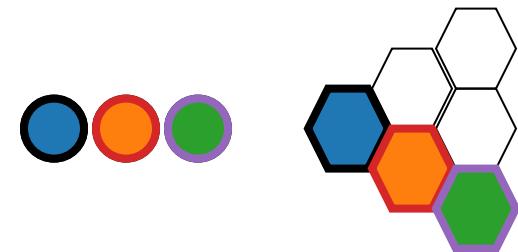
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits
  - **Node granularity**
  - Difficult with dependencies
- Submit task graph as a single job
  - Only for small-ish task graphs



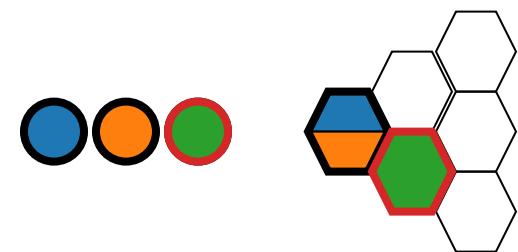
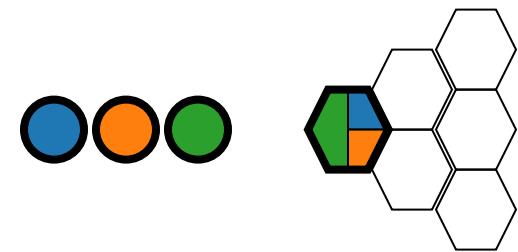
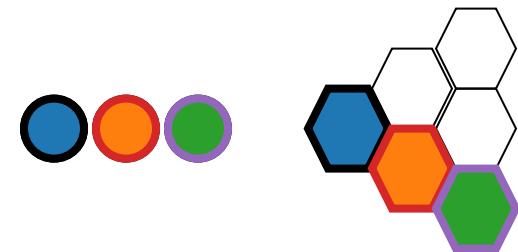
# Submitting task graphs into PBS

- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits
  - **Node granularity**
  - Difficult with dependencies
- Submit task graph as a single job
  - Only for small-ish task graphs
- Split task graph into multiple jobs
  - Lot of work

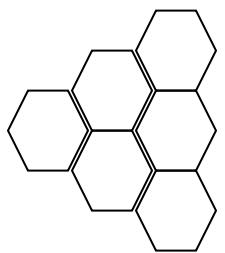


# Submitting task graphs into PBS

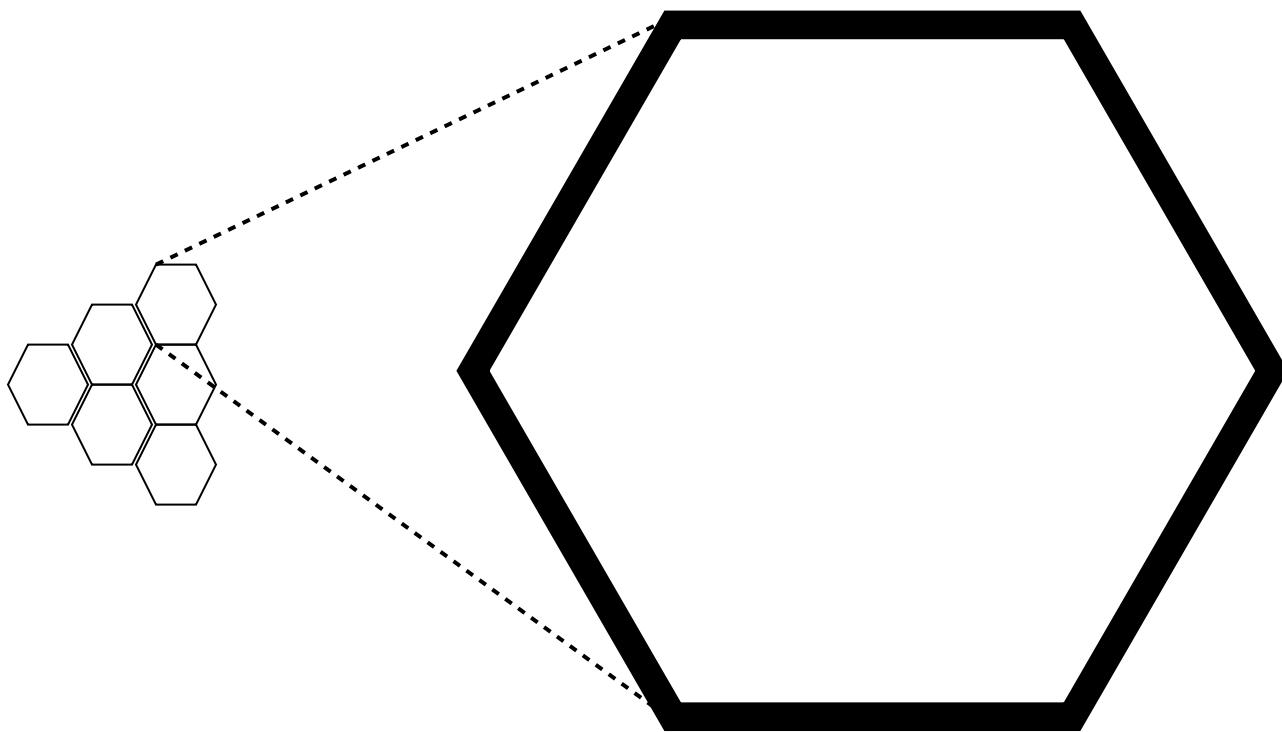
- Submit each task as a job
  - Massive overhead (millions of jobs)
  - Job count limits
  - **Node granularity**
  - Difficult with dependencies
- Submit task graph as a single job
  - Only for small-ish task graphs
- Split task graph into multiple jobs
  - Lot of work
  - No load balancing across jobs



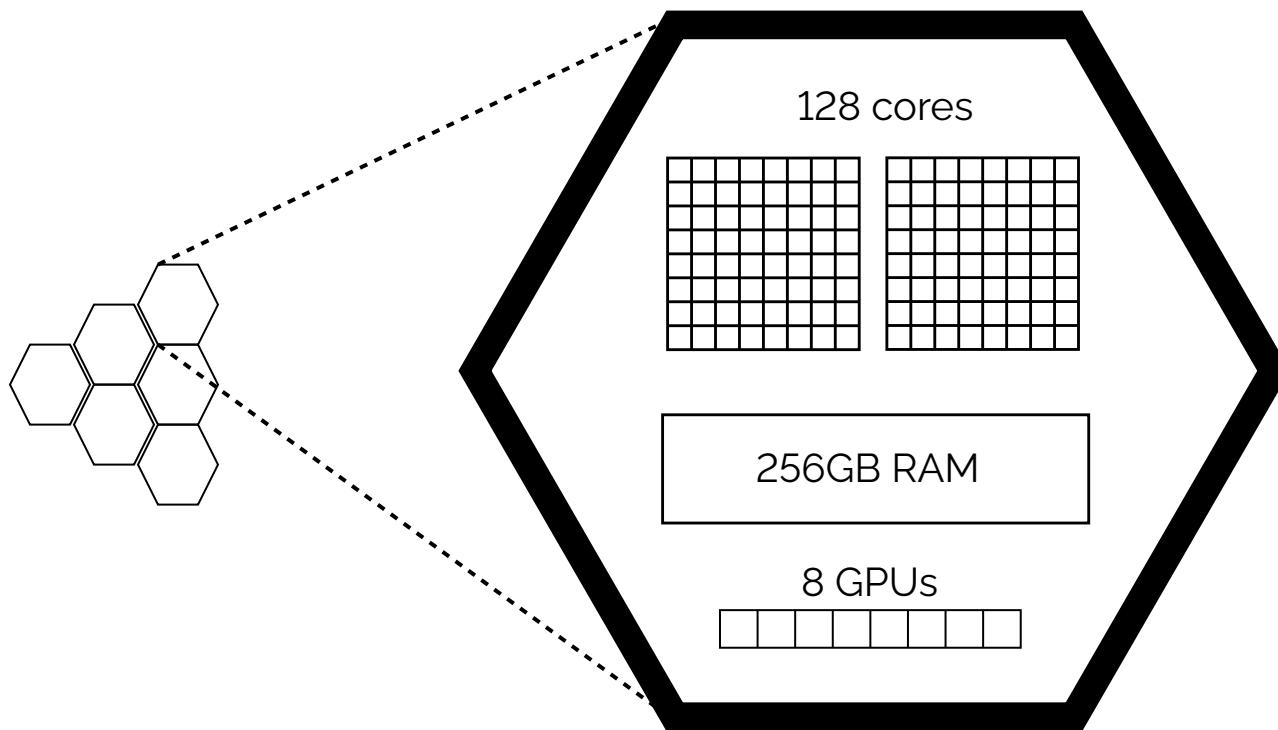
# Resource management



# Resource management



# Resource management



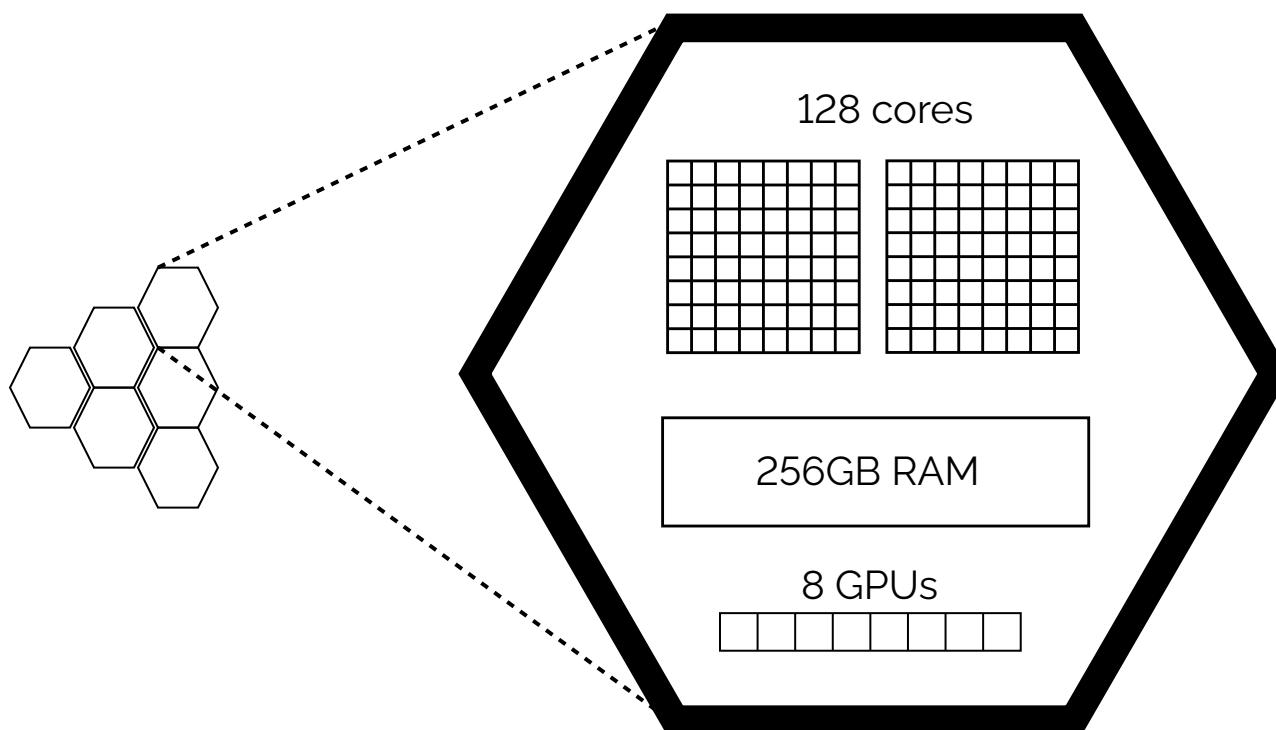
# Resource management

1 core

32 cores  
128 GB ram

16 cores  
2 gpu

32 cores  
in the same socket



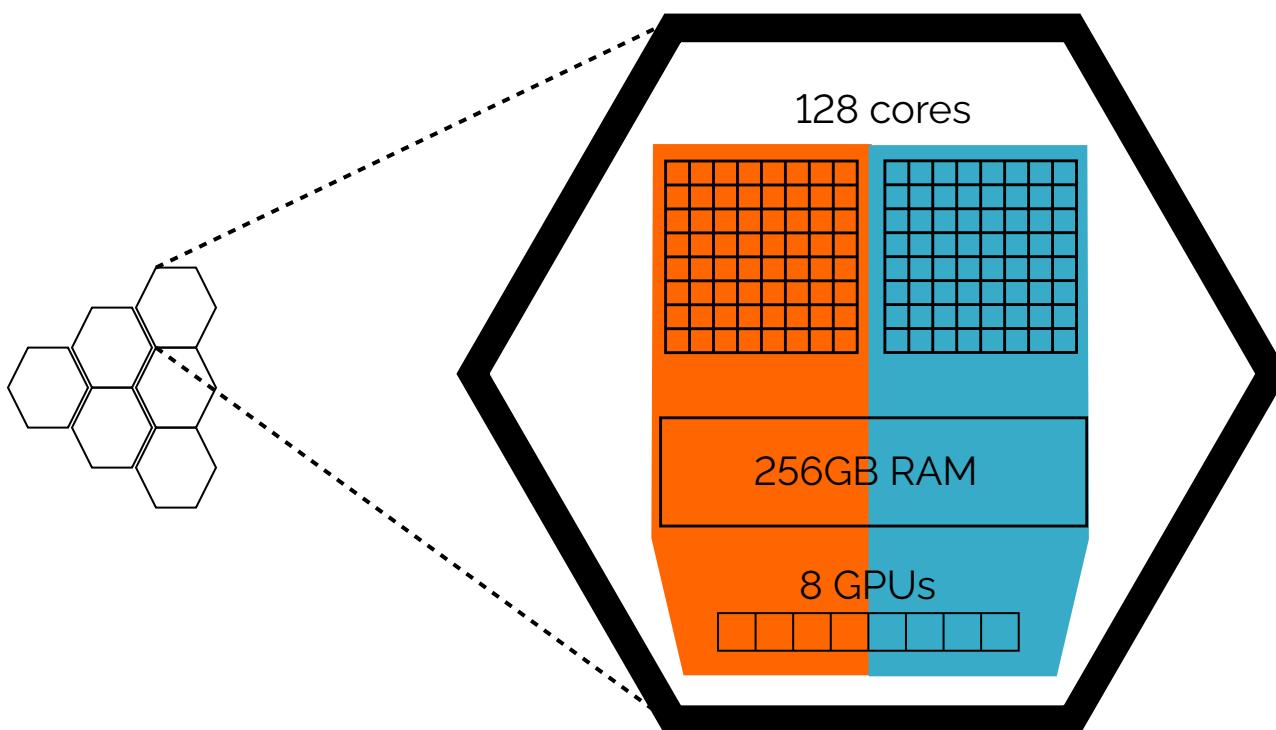
# Resource management

1 core

32 cores  
128 GB ram

16 cores  
2 gpu

32 cores  
in the same socket



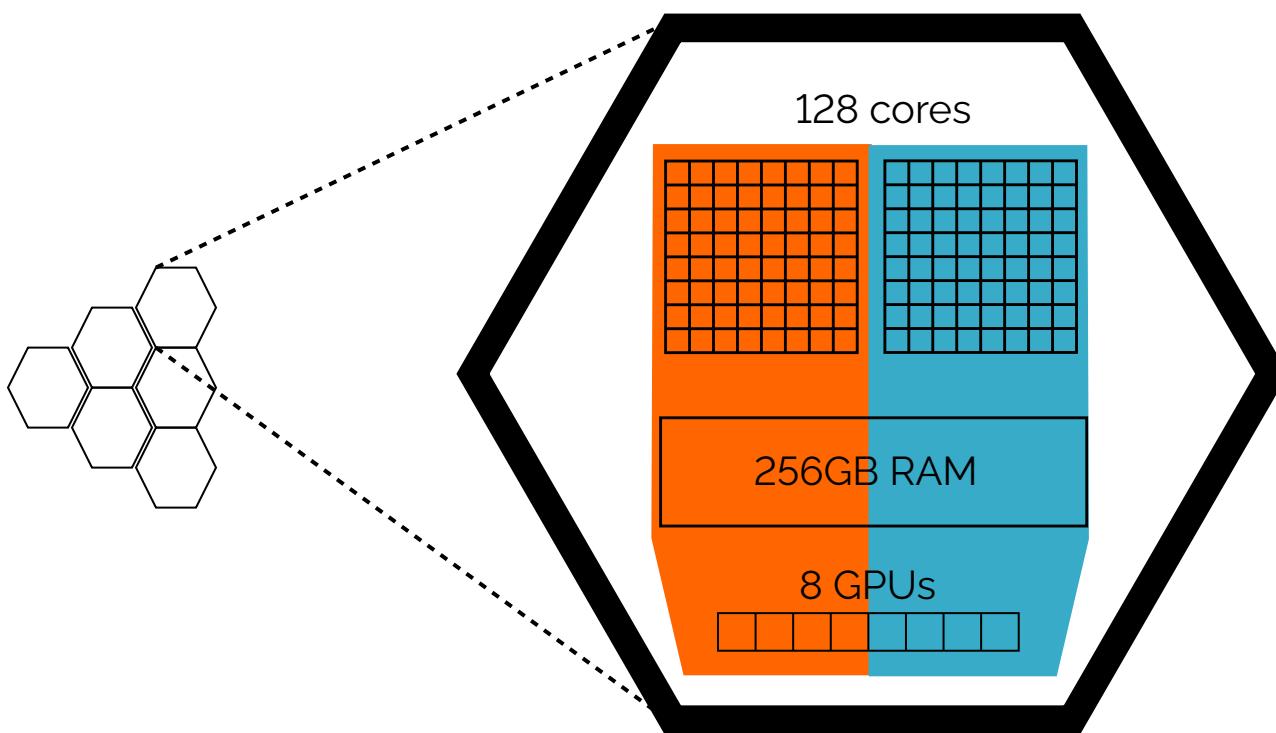
# Resource management

1 core

32 cores  
128 GB ram

16 cores  
2 gpu

32 cores  
in the same socket

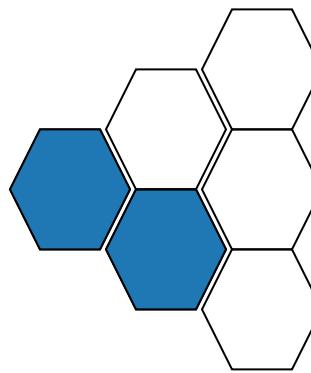
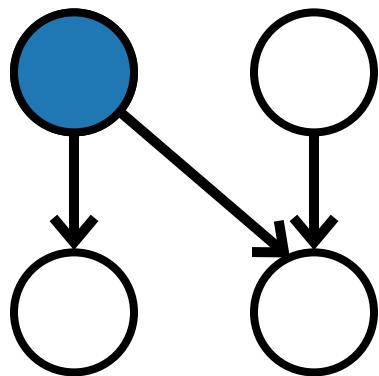


128 cores

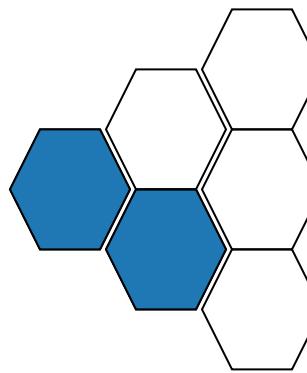
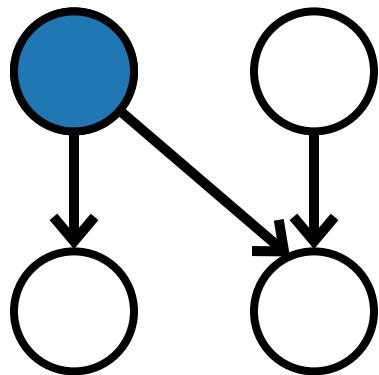
64 cores  
8 GPUs

32 cores  
24TB ram

# Multi-node tasks

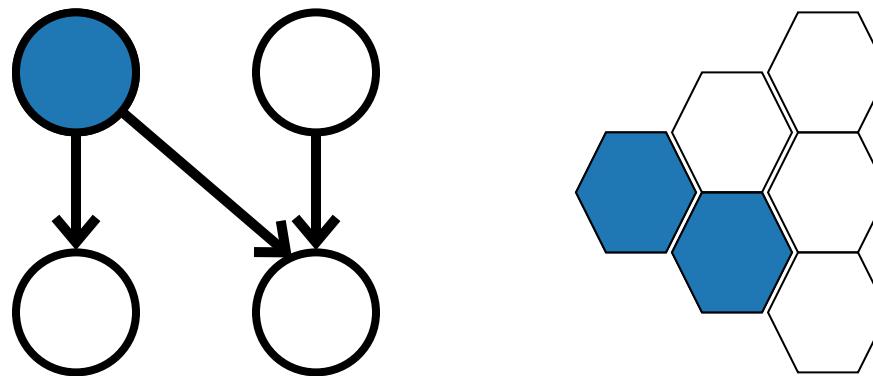


# Multi-node tasks



- MPI application within a task

# Multi-node tasks



- MPI application within a task
  - Complicates scheduling and data transfers

# Other workflow challenges on HPC

- Scalability

# Other workflow challenges on HPC

- Scalability
  - Millions of tasks, hundreds of nodes

# Other workflow challenges on HPC

- Scalability
  - Millions of tasks, hundreds of nodes
  - Communication bottleneck

# Other workflow challenges on HPC

- Scalability
  - Millions of tasks, hundreds of nodes
  - Communication bottleneck
  - Fault tolerance

# Other workflow challenges on HPC

- Scalability
  - Millions of tasks, hundreds of nodes
  - Communication bottleneck
- Fault tolerance
- Iterative computation

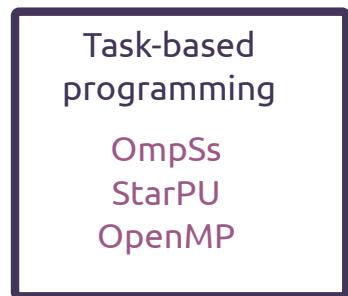
# Other workflow challenges on HPC

- Scalability
  - Millions of tasks, hundreds of nodes
  - Communication bottleneck
- Fault tolerance
- Iterative computation
- ...

# State of the art: task duration

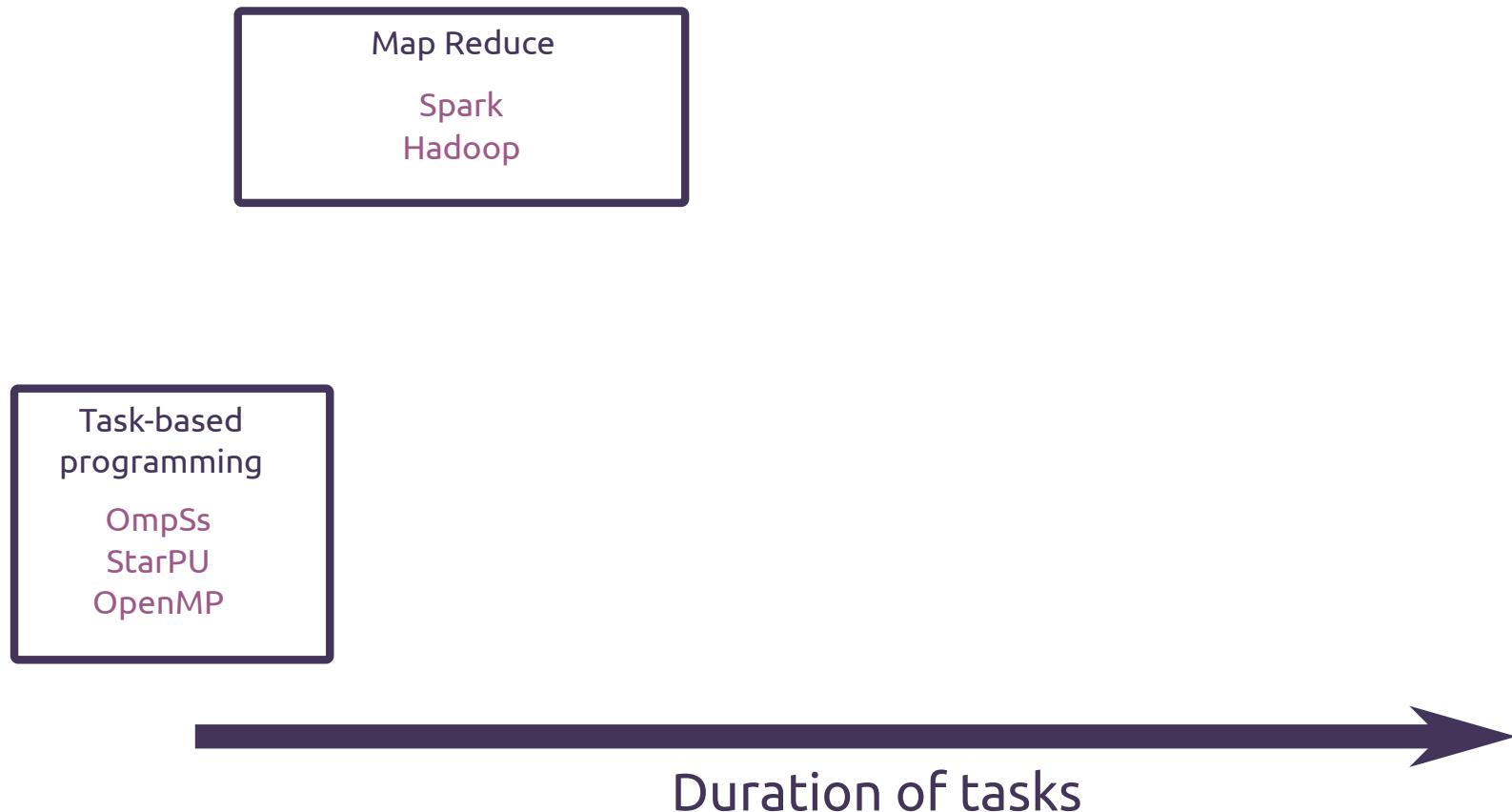


# State of the art: task duration

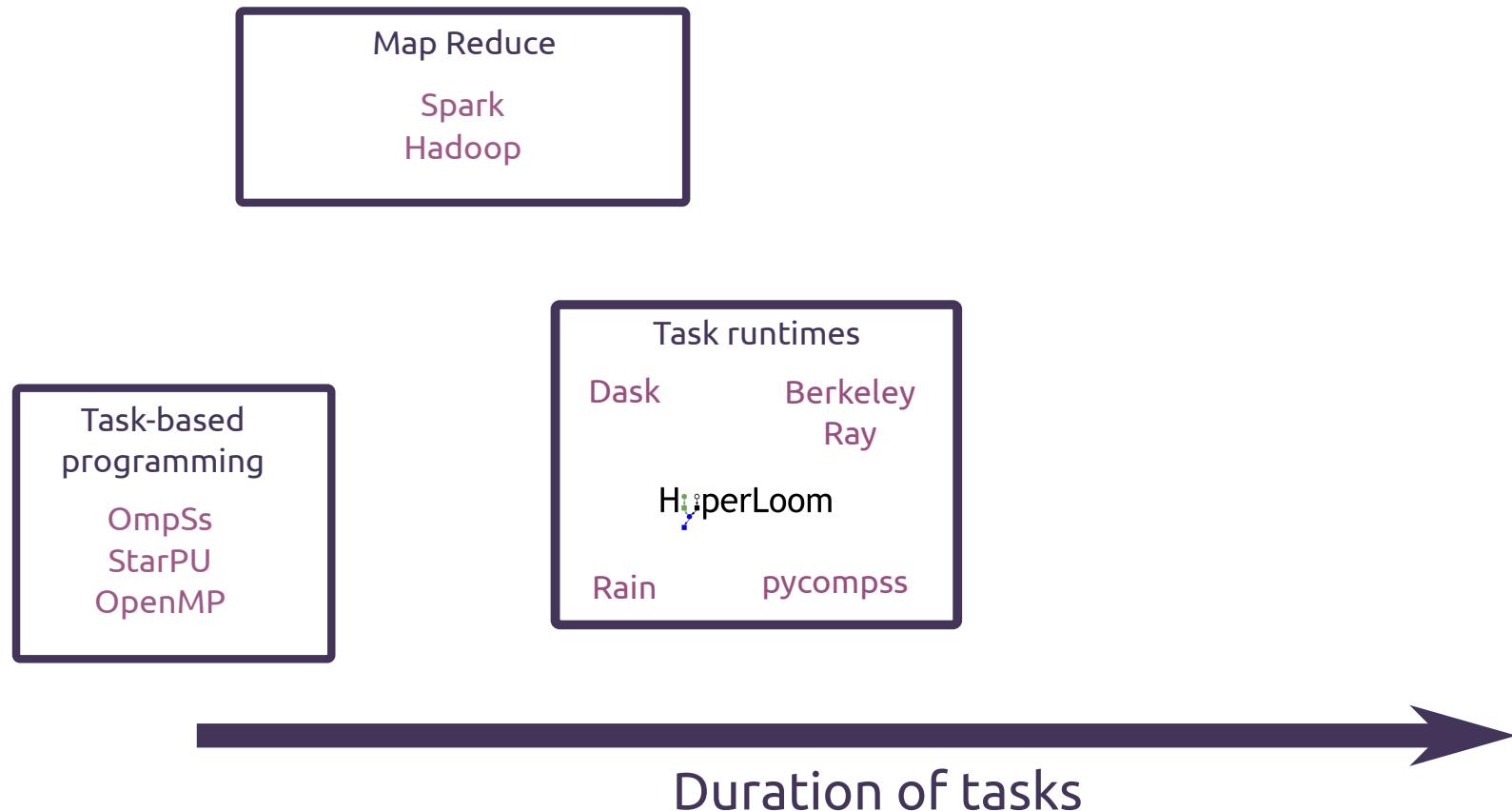


Duration of tasks

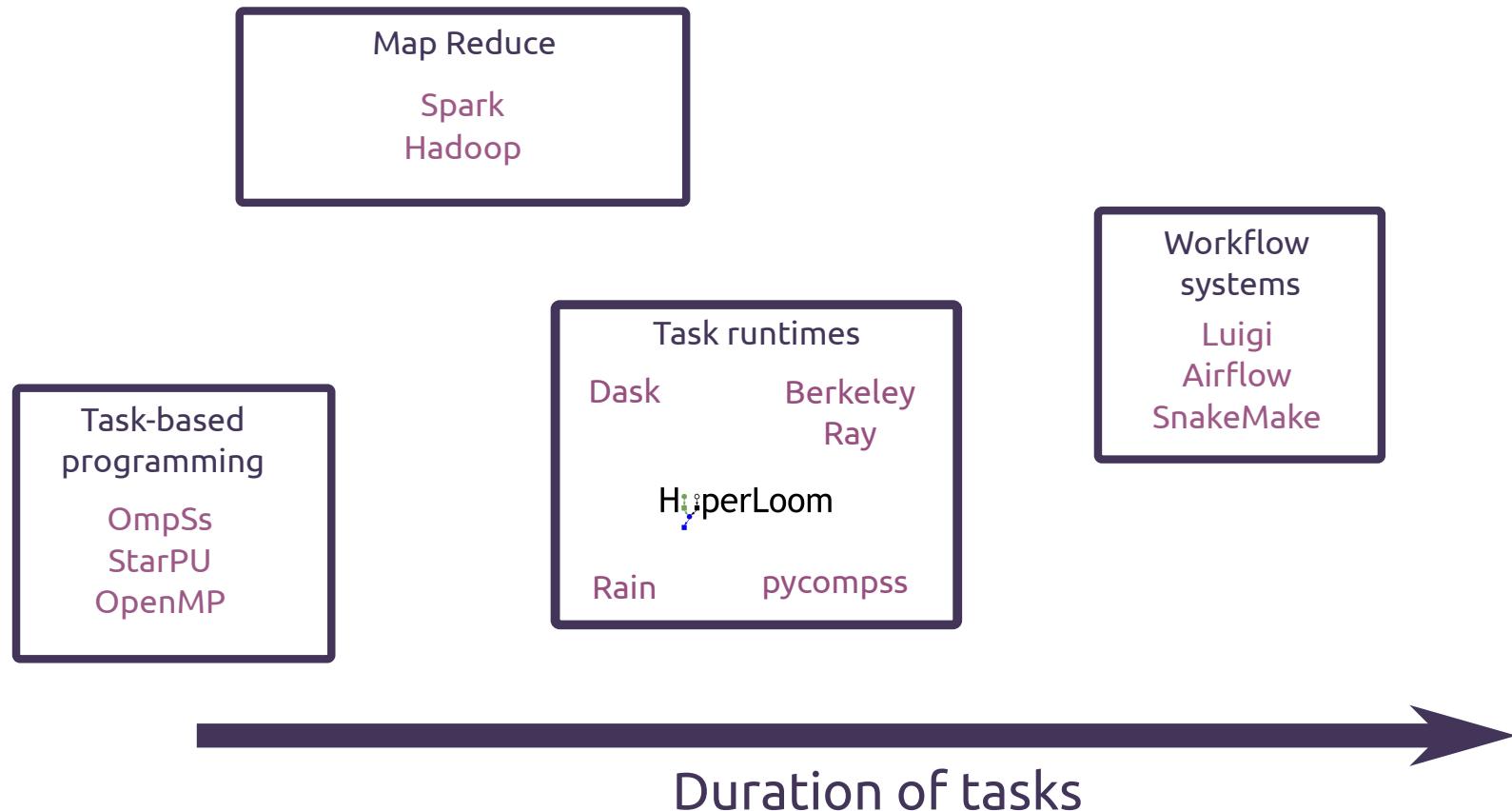
# State of the art: task duration



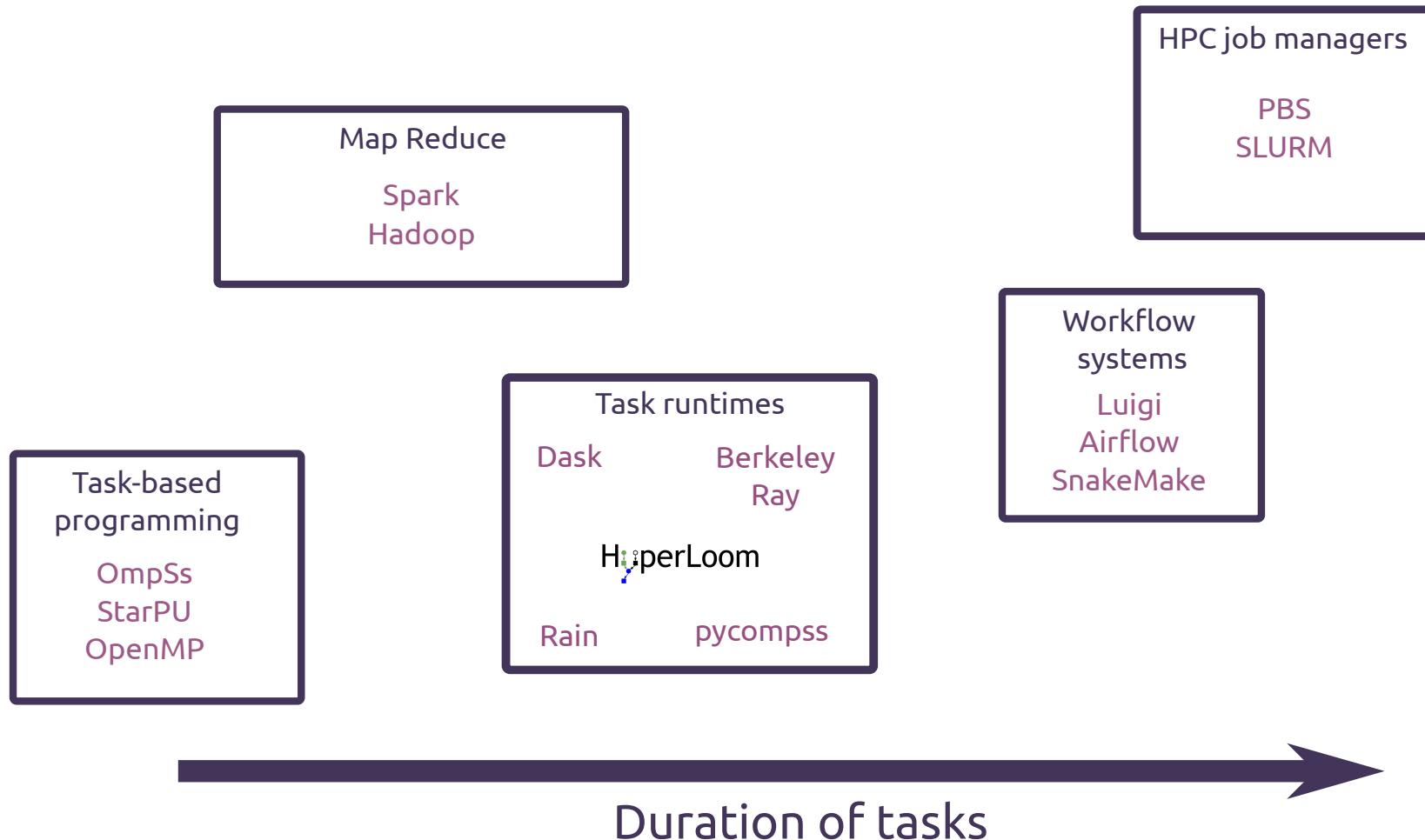
# State of the art: task duration



# State of the art: task duration



# State of the art: task duration

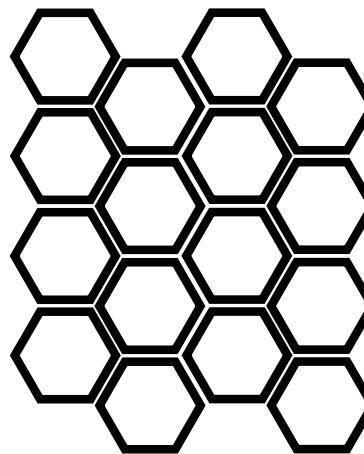


# State of the art: job manager



Login nodes

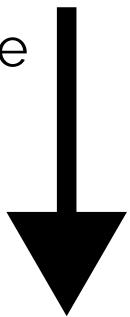
PBS/Slurm



Compute nodes

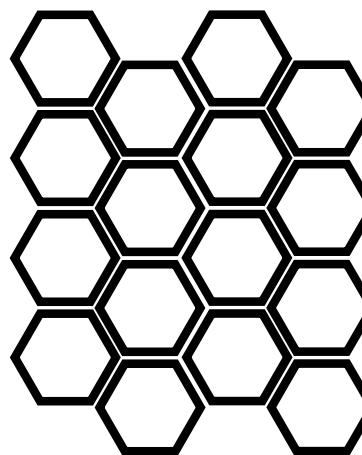
# State of the art: job manager

SnakeMake  
Nextflow  
Airflow



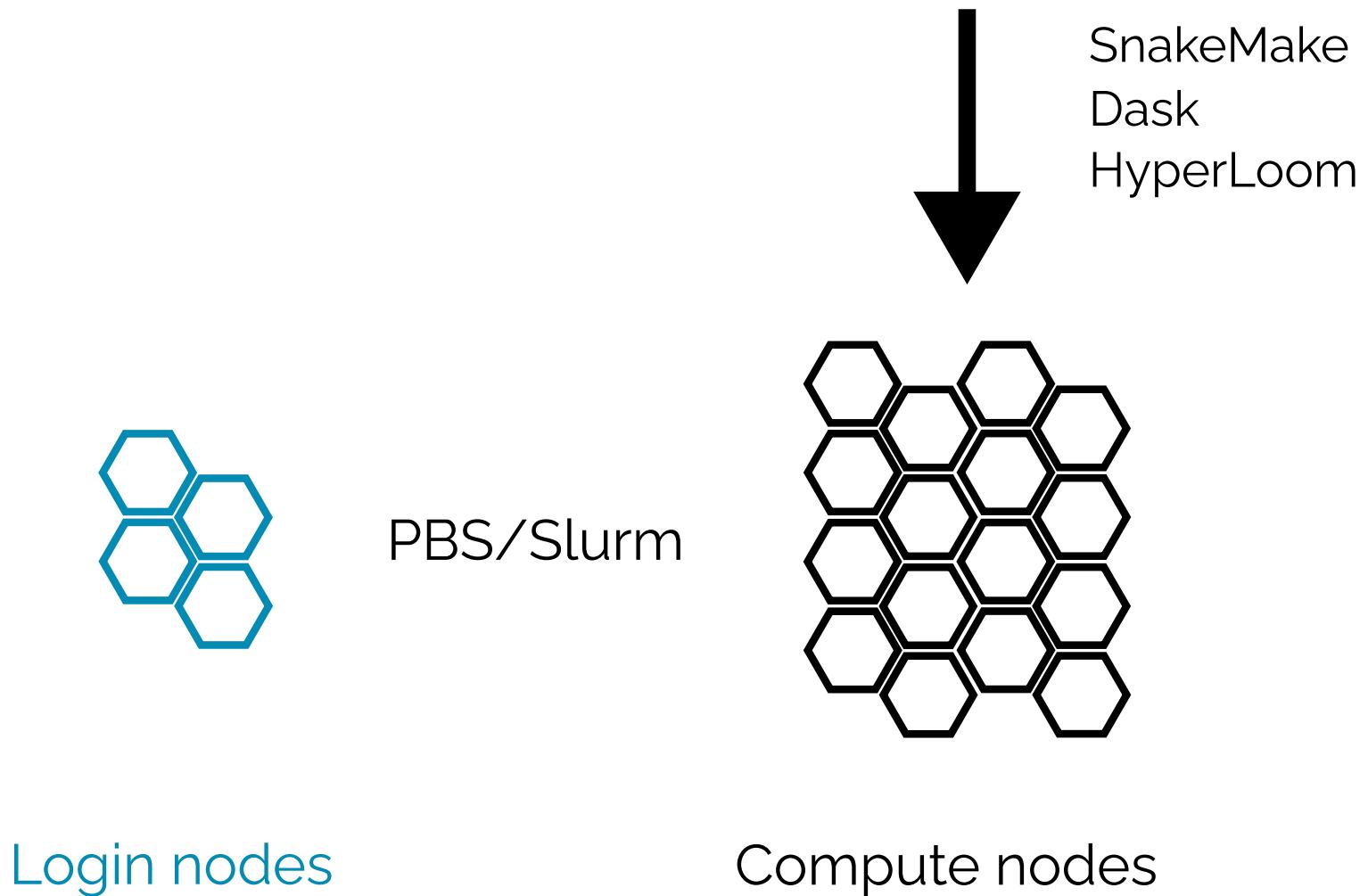
Login nodes

PBS/Slurm



Compute nodes

# State of the art: job manager

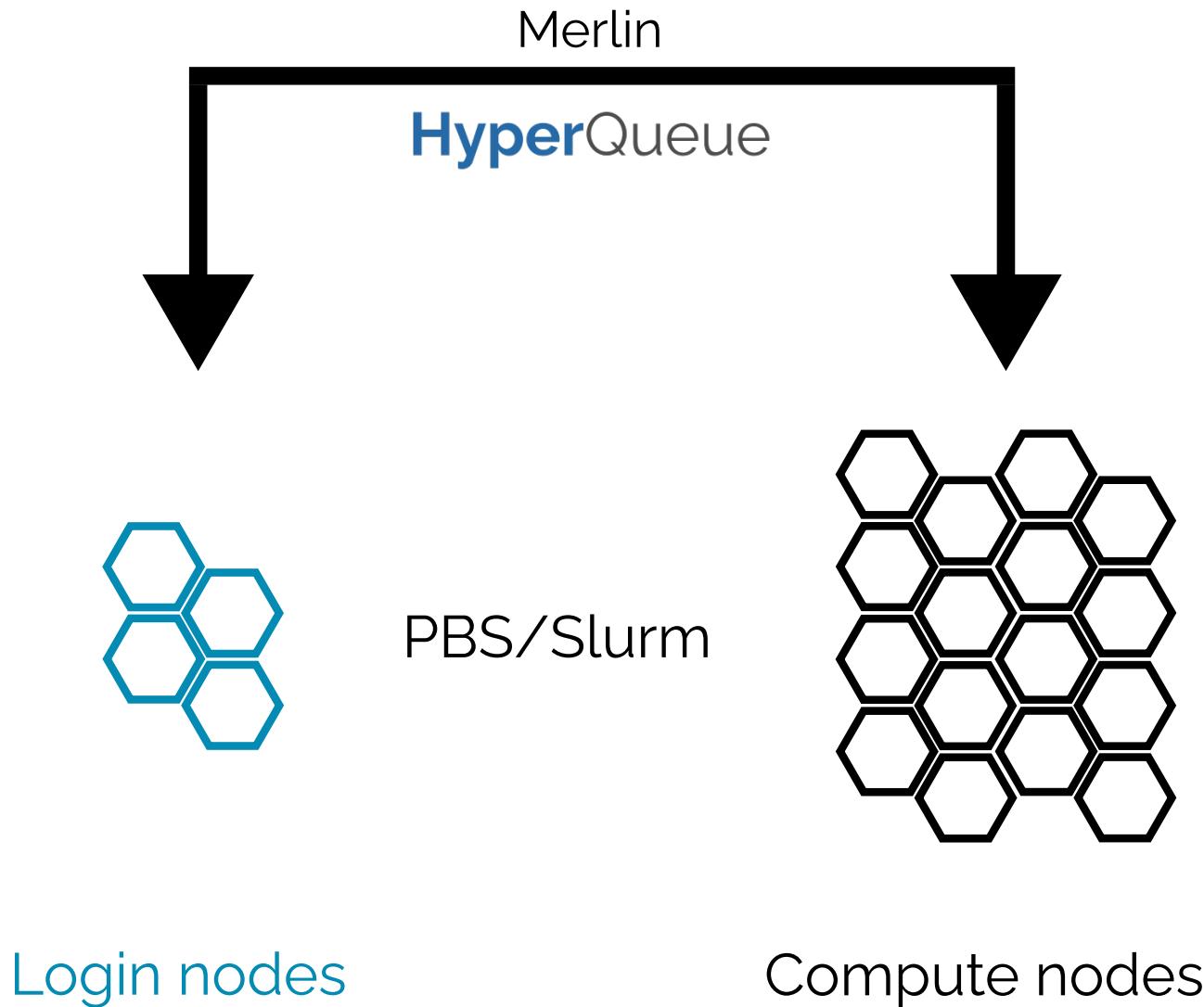


Login nodes

PBS/Slurm

Compute nodes

# State of the art: job manager



# **State of the art: resources & multi-node**

- Resource requirements: ✓ (simple)

# State of the art: resources & multi-node

- Resource requirements: ✓ (simple)
- Multi-node tasks: x

# Ergonomics and **efficiency** of workflows on HPC clusters



**Abstract**

Task graphs provide a simple way to describe scientific workflow (sets of tasks with dependencies) that can be executed on both HPC clusters and in the cloud. An important aspect of workflow execution is the choice of scheduling algorithm. Many scheduling metrics have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research, which is freely available online. It contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of the effect of various factors on the performance of workflow schedulers. The effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocations or partially unknown task dependencies, are indicated that several factors used by popular works might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that the performance of some well-known scheduling algorithms, which are often neglected, can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsb.cz

Stanislav Böhm  
stanislaw.bohm@vsb.cz

Vojtěch Cima  
vojtěch.cima@vsb.cz

<sup>1</sup> IT Innovations, VSB – Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022



# Analysis of workflow schedulers in simulated distributed environments

Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)



**Abstract** Task graphs provide a simple way to describe scientific workflow sets of tasks with dependencies that can be executed on both HPC clusters and in the cloud. An important aspect of workflow management is the scheduler algorithm. Many scheduling metrics have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research. The environment is a Java application which contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of workflow scheduling performance. Our results show that there is a large effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocations or partially unknown task dependencies. Our results indicate that several challenges in real-world workflows might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that the performance of some well-known scheduling algorithms, which are often neglected, can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsb.cz  
Stanislav Böhm  
stanislaw.bohm@vsb.cz

Vojtěch Cima  
vojtech.cima@vsb.cz

<sup>1</sup> IT Innovations, VSB – Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022

Springer

# Analysis of workflow schedulers in simulated distributed environments

Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)

- Compare scheduler performance



**Abstract**  
Task graphs provide a simple way to describe scientific workflow (sets of tasks with dependencies) that can be executed on both HPC clusters and in the cloud. An important aspect of workflow execution is the choice of scheduling algorithm. Many scheduling metrics have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research. The environment is freely available, which contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of workflow scheduling performance. Our results show a significant effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocations or partially unknown task dependencies. Our results indicate that several factors used by previous works might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that the performance of some well-known algorithms, which are often neglected can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsbr.cz  
Stanislav Böhm  
stanislaw.bohm@vsbr.cz

Vojtěch Cima  
vojtech.cima@vsbr.cz

<sup>1</sup> IT Innovations, VSBR – Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022

Springer

# Analysis of workflow schedulers in simulated distributed environments

Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)

- Compare scheduler performance
- Analyze neglected factors



**Abstract**  
Task graphs provide a simple way to describe scientific workflow (sets of tasks with dependencies) that can be executed on both HPC clusters and in the cloud. An important aspect of workflow execution is the scheduling algorithm. Many scheduling approaches have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research. The environment is a Java application which contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of workflow scheduling. Our results show that there is a significant effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocations or partially unknown task dependencies. Our results indicate that several factors used by previous works might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that the performance of some well-known scheduling algorithms, which are often neglected, can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsb.cz

Stanislav Böhm  
stanislaw.bohm@vsb.cz

Vojtěch Cima  
vojtech.cima@vsb.cz

<sup>1</sup> IT Innovations, VSB – Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022



# Analysis of workflow schedulers in simulated distributed environments

Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)

- Compare scheduler performance
- Analyze neglected factors
  - Delay between scheduling



**Abstract**  
Task graphs provide a simple way to describe scientific workflows (sets of tasks with dependencies) that can be executed on both HPC clusters and in the cloud. An important aspect of workflow execution is the scheduling algorithm. Many scheduling approaches have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research. The environment is a Java application which contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of workflow scheduling performance. Our results show a significant effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocation or partially unknown task durations. Our results indicate that several factors used by previous works might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that knowledge about task durations is crucial for scheduler algorithms which are often neglected can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsbr.cz  
Stanislav Böhm  
stanislaw.bohm@vsbr.cz

Vojtěch Cima  
vojtech.cima@vsbr.cz

<sup>1</sup> IT Innovations, VSBR – Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022

Springer

# Analysis of workflow schedulers in simulated distributed environments

Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)

- Compare scheduler performance
- Analyze neglected factors
  - Delay between scheduling
  - Knowledge about task durations



**Abstract**  
Task graphs provide a simple way to describe scientific workflow (sets of tasks with dependencies) that can be executed on both HPC clusters and in the cloud. An important aspect of workflow execution is the scheduling algorithm. Many scheduling metrics have been proposed in existing works; nevertheless, they are often tested in oversimplified environments. We provide an extensible simulation environment for workflow scheduling research. The environment is open-source, which contains implementations of various scheduling algorithms and is open-sourced, in order to be fully reproducible. We use this environment to perform a comprehensive analysis of workflow scheduling. We show that there is a significant effect of scheduling challenges that have so far been mostly neglected, such as delays between scheduler invocation or partially unknown task durations. Our results indicate that several factors used by previous works might produce results that are off by an order of magnitude in comparison to a more realistic model. Additionally, we show that knowledge about task durations and network models, which are often neglected, can have a large effect on the scheduler's performance, and they should thus be described in great detail to enable proper evaluation.

**Keywords** Distributed computing · DAG scheduling · Task Scheduling · Network models

Jakub Beránek, Stanislav Böhm and Vojtěch Cima these authors contributed equally to this work.

✉ Jakub Beránek  
jakub.beranek@vsbr.cz  
Stanislav Böhm  
stanislaw.bohm@vsbr.cz

Vojtěch Cima  
vojtech.cima@vsbr.cz

<sup>1</sup> IT Innovations, VSBR - Technical University of Ostrava, Ostrava, Czech Republic

Published online: 14 April 2022

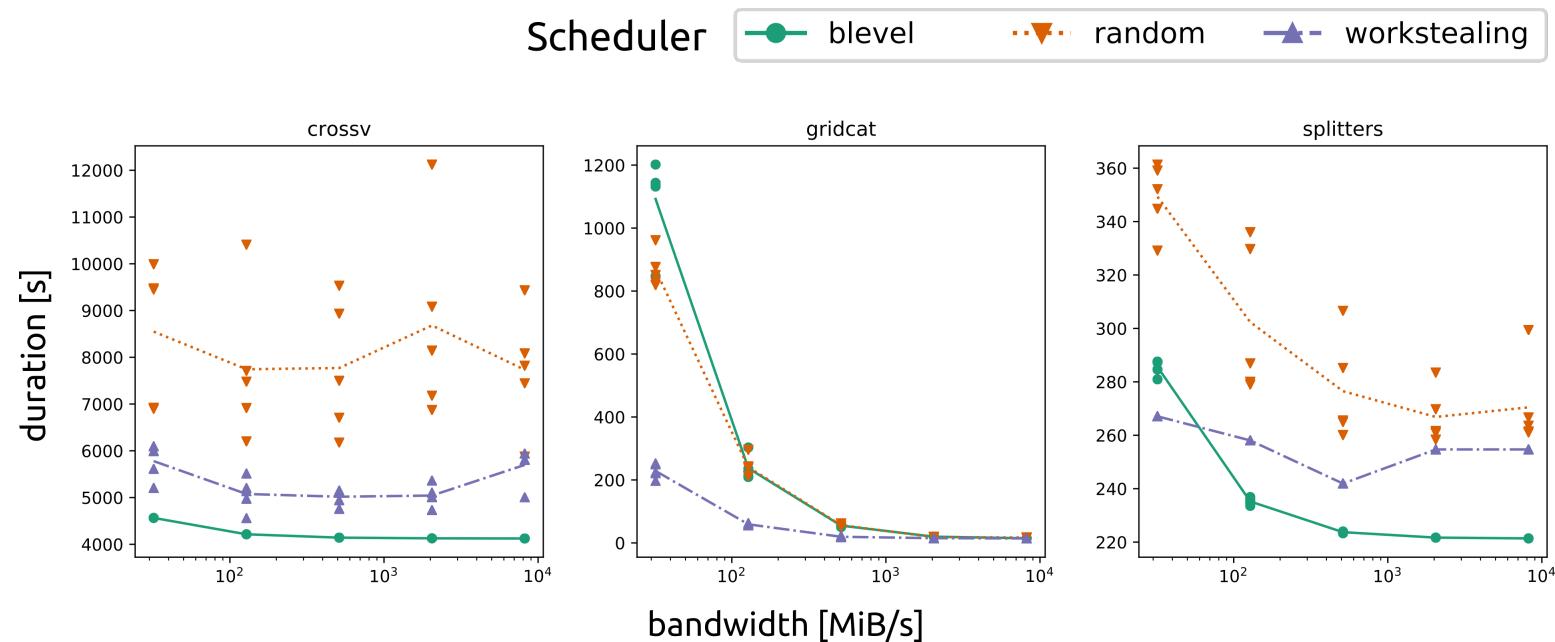


# Analysis of workflow schedulers in simulated distributed environments

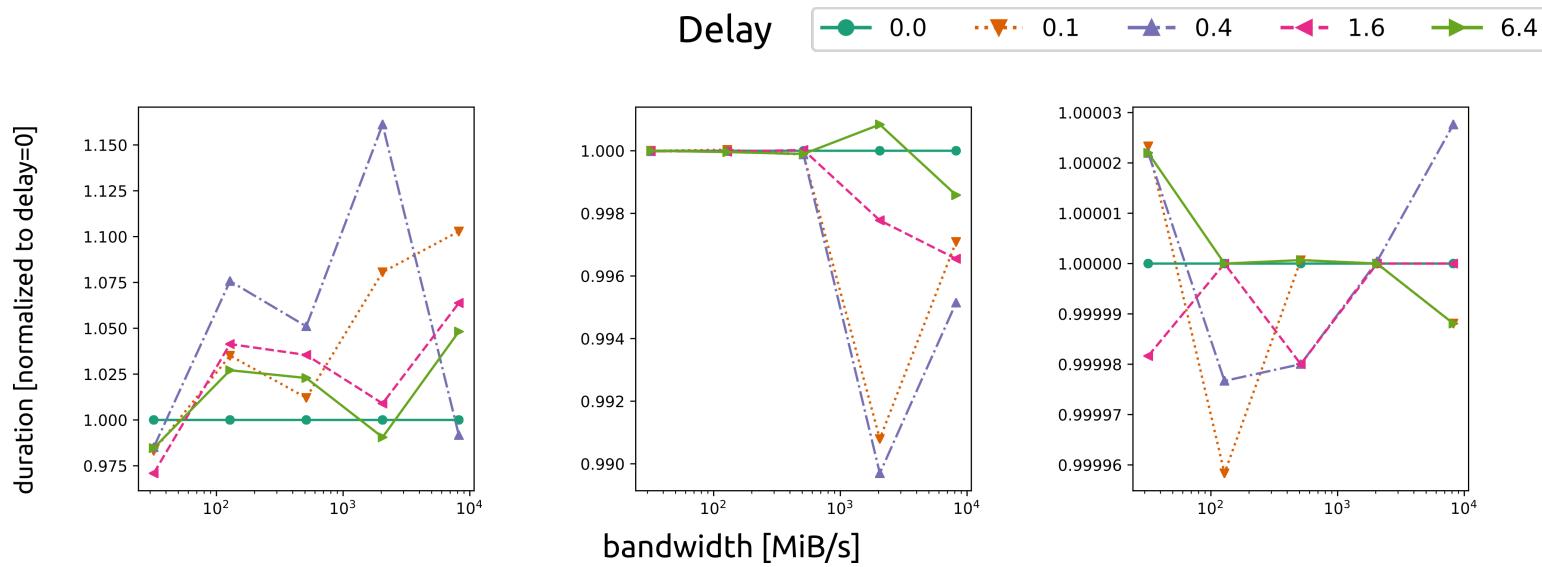
Jakub Beránek, Stanislav Böhm, Vojtěch Cima  
(The Journal of Supercomputing 2022)

- Compare scheduler performance
- Analyze neglected factors
  - Delay between scheduling
  - Knowledge about task durations
  - Network model

# Runtime vs. network speed



# Scheduling delay effect



# Outcome

- Most competitive schedulers

# Outcome

- Most competitive schedulers
  - Work-stealing (used in Dask)

# Outcome

- Most competitive schedulers
  - Work-stealing (used in Dask)
  - B-level (basic heuristic)

# Outcome

- Most competitive schedulers
  - Work-stealing (used in Dask)
  - B-level (basic heuristic)
- Implementation details matter a lot!

# Outcome

- Most competitive schedulers
  - Work-stealing (used in Dask)
  - B-level (basic heuristic)
- Implementation details matter a lot!
- Open source scheduler simulator ESTEE
  - [github.com/it4innovations/estee](https://github.com/it4innovations/estee)









# Dask vs. RSDS: work-stealing scheduler

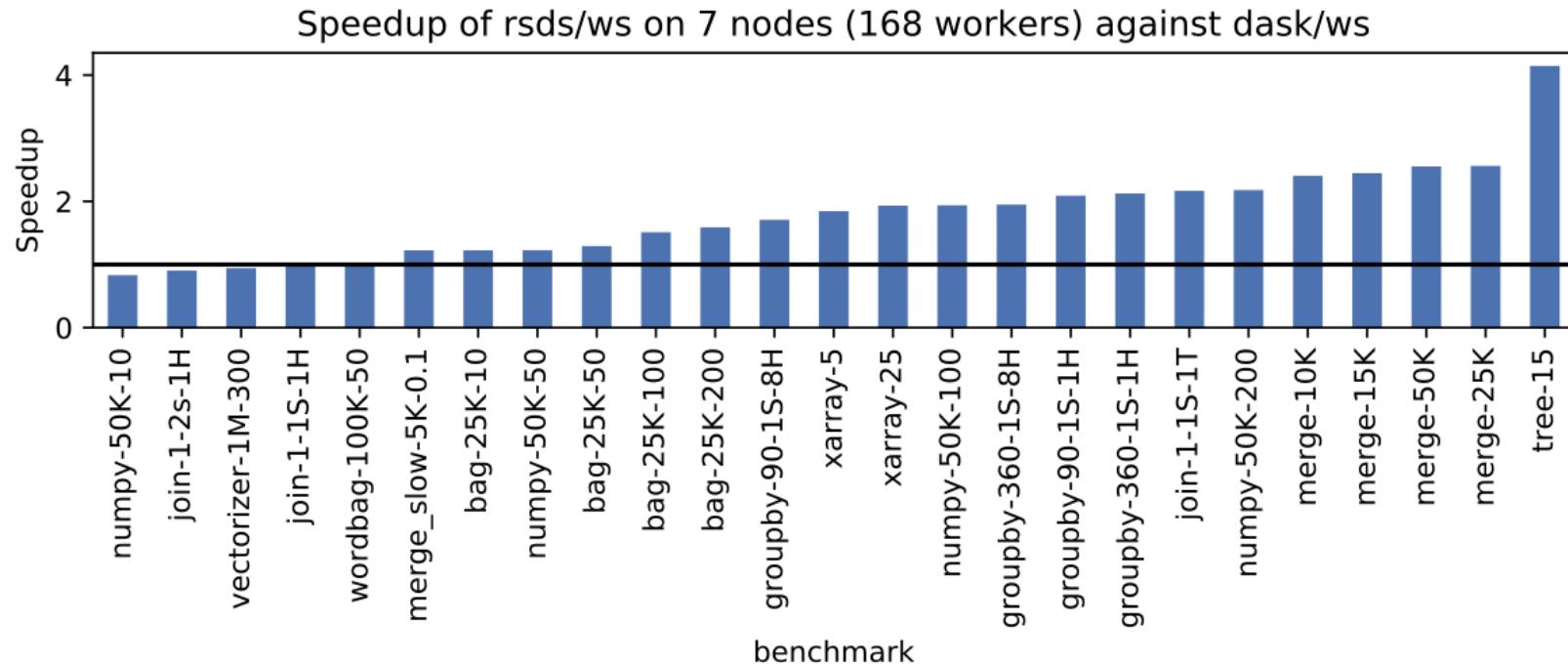
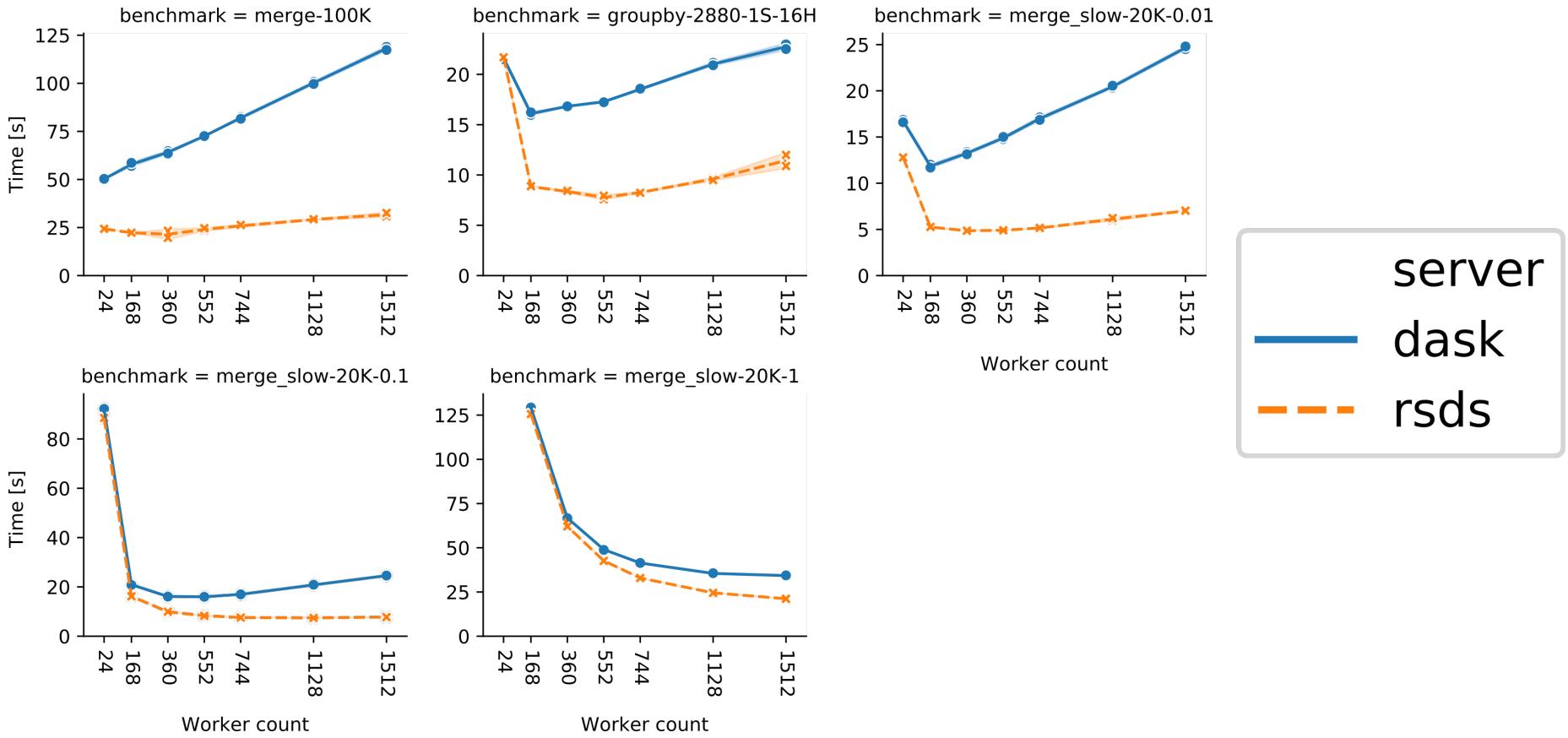


Fig. 3. Speedup of RSDS/ws scheduler; DASK/ws is baseline.

# Dask vs. RSDS: scaling on Salomon



# Outcome

- Runtime efficiency as important as scheduling

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!
  - Dask scaled poorly on HPC

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!
- Dask scaled poorly on HPC
  - Caused by inefficient runtime

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!
- Dask scaled poorly on HPC
  - Caused by inefficient runtime
  - <100 ms tasks problematic

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!
- Dask scaled poorly on HPC
  - Caused by inefficient runtime
  - <100 ms tasks problematic
- RSDS: open source alternative to Dask's server

# Outcome

- Runtime efficiency as important as scheduling
  - Even a random scheduler can be competitive!
- Dask scaled poorly on HPC
  - Caused by inefficient runtime
  - <100 ms tasks problematic
- RSDS: open source alternative to Dask's server
  - Backward-compatible with existing Dask programs
  - [github.com/it4innovations/rsds](https://github.com/it4innovations/rsds)

# **Ergonomics** and efficiency of workflows on HPC clusters

# HyperQueue

[github.com/it4innovations/hyperqueue](https://github.com/it4innovations/hyperqueue)

Task runtime designed for HPC

# HyperQueue

[github.com/it4innovations/hyperqueue](https://github.com/it4innovations/hyperqueue)

Task runtime designed for HPC



IT4INNOVATIONS  
NATIONAL SUPERCOMPUTING  
CENTER



**LUMI**

# HyperQueue

[github.com/it4innovations/hyperqueue](https://github.com/it4innovations/hyperqueue)

Task runtime designed for HPC



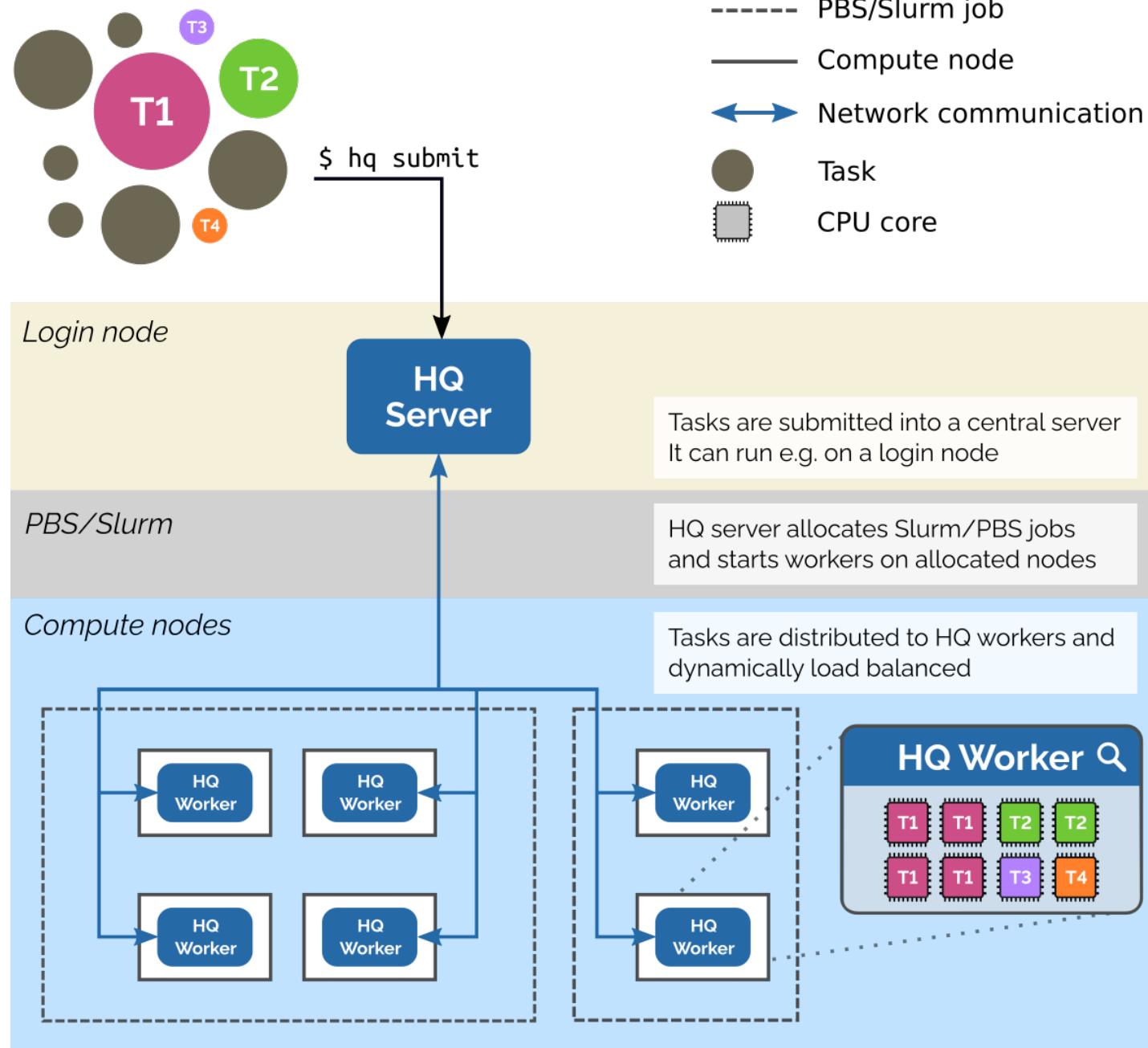
IT4INNOVATIONS  
NATIONAL SUPERCOMPUTING  
CENTER



LUMI

Team effort @ IT4I  
(primary contributors: Ada Böhm & me)

# Architecture



# Command-line interface

Start server

```
(login1) $ hq server start
```

# Command-line interface

Start server

```
(login1) $ hq server start
```

Submit tasks

```
(login1) $ hq submit -- ./my-program --foo=bar
```

# Command-line interface

Start server

```
(login1) $ hq server start
```

Submit tasks

```
(login1) $ hq submit -- ./my-program --foo=bar
```

Provide computational resources

```
(login1) $ hq alloc add pbs --time-limit 1h -- -qqexp
```

# Command-line interface

Start server

```
(login1) $ hq server start
```

Submit tasks

```
(login1) $ hq submit -- ./my-program --foo=bar
```

Provide computational resources

```
(login1) $ hq alloc add pbs --time-limit 1h -- -qqexp
```

or

```
(node1) $ hq worker start
```

# Python API

```
def preprocess(path):
    # preprocess data

path = "/tmp/foo"

job = Job()
t1 = job.function(preprocess, args=(path, ))
job.program(["simulate", "--path", path], deps=[t1])
client.submit(job)
```

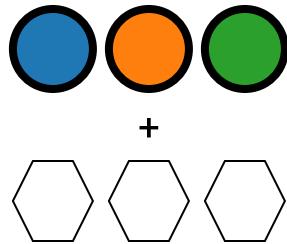
# **Challenge: PBS/Slurm**

Disentangle computation description from resources

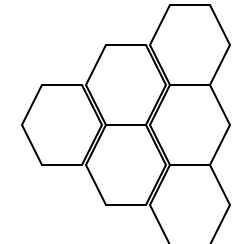
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



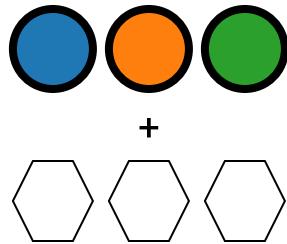
→ PBS →



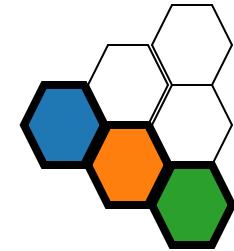
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



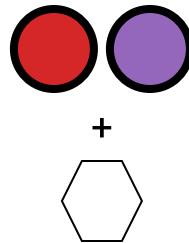
→ PBS →



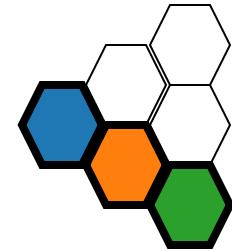
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



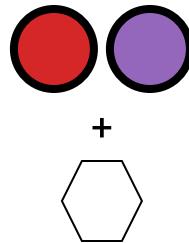
→ PBS →



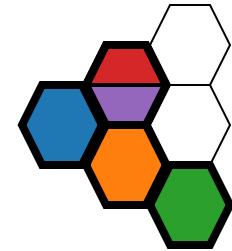
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



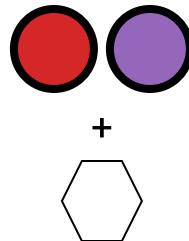
→ PBS →



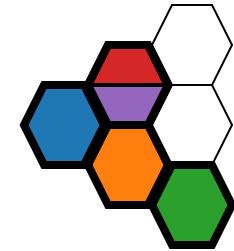
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



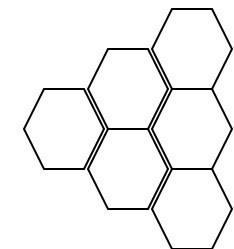
→ PBS →



PBS + HQ



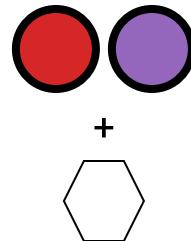
→ HyperQueue → PBS →



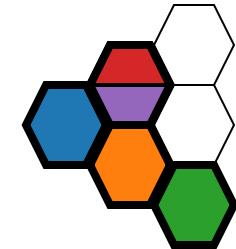
# Challenge: PBS/Slurm

Disentangle computation description from resources

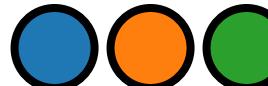
PBS only



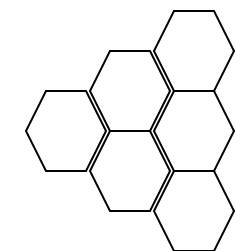
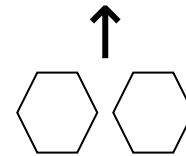
→ PBS →



PBS + HQ



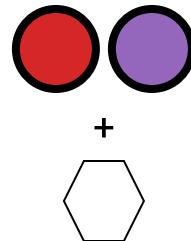
→ HyperQueue → PBS →



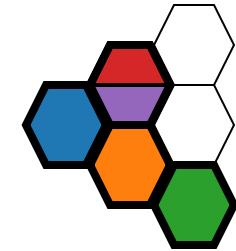
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



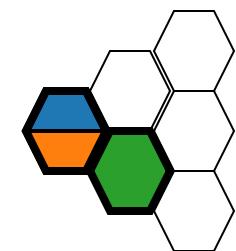
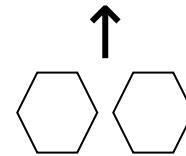
→ PBS →



PBS + HQ



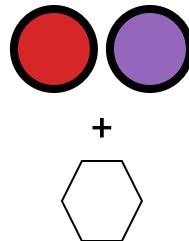
→ HyperQueue → PBS →



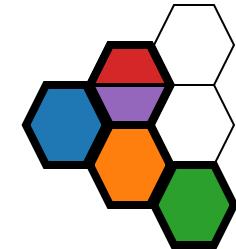
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



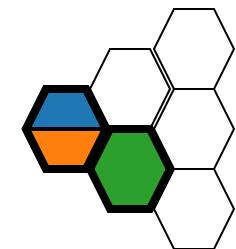
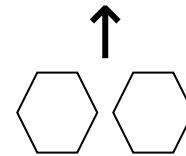
→ PBS →



PBS + HQ



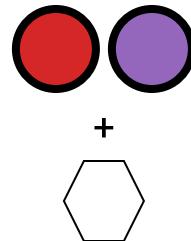
→ HyperQueue → PBS →



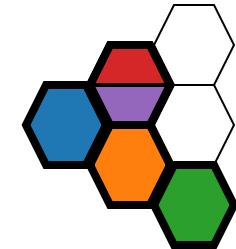
# Challenge: PBS/Slurm

Disentangle computation description from resources

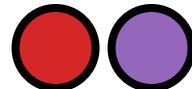
PBS only



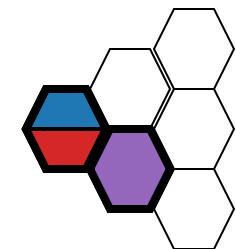
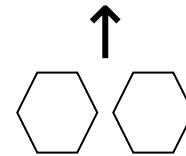
→ PBS →



PBS + HQ



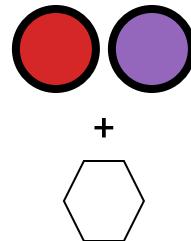
→ HyperQueue → PBS →



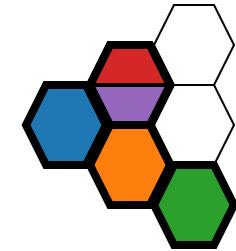
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



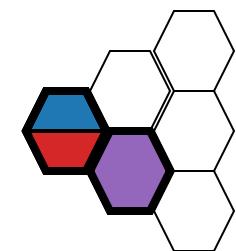
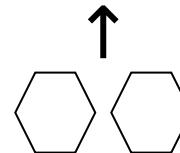
→ PBS →



PBS + HQ



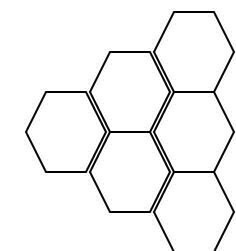
→ HyperQueue → PBS →



PBS + HQ  
+autoalloc



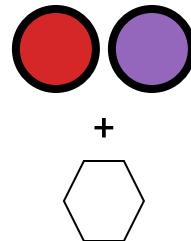
→ HyperQueue → PBS →



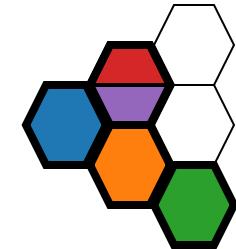
# Challenge: PBS/Slurm

Disentangle computation description from resources

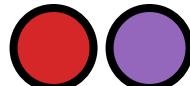
PBS only



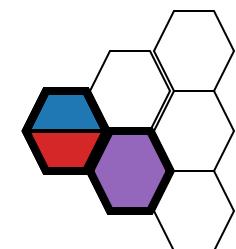
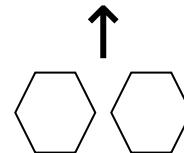
→ PBS →



PBS + HQ



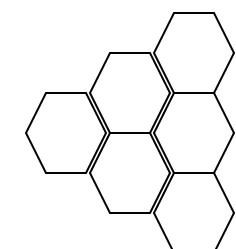
→ HyperQueue → PBS →



PBS + HQ  
+autoalloc



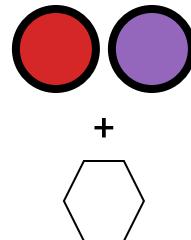
→ HyperQueue → PBS →



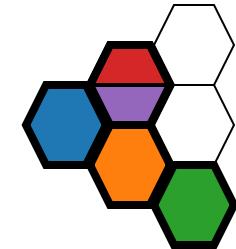
# Challenge: PBS/Slurm

Disentangle computation description from resources

PBS only



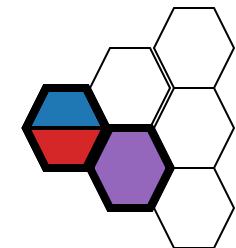
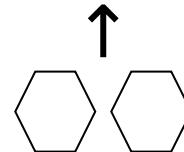
→ PBS →



PBS + HQ



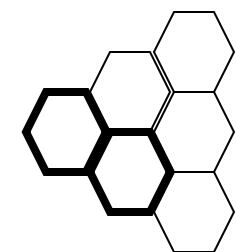
→ HyperQueue → PBS →



PBS + HQ  
+autoalloc



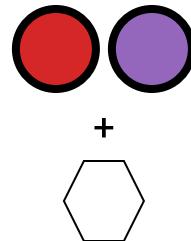
→ HyperQueue → PBS →



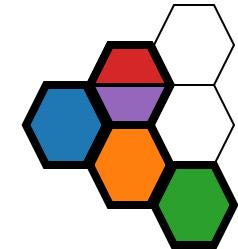
# Challenge: PBS/Slurm

Disentangle computation description from resources

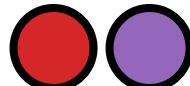
PBS only



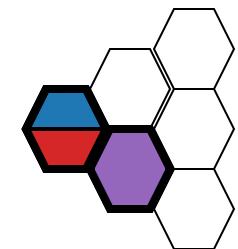
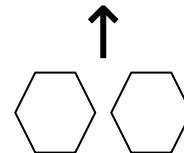
→ PBS →



PBS + HQ



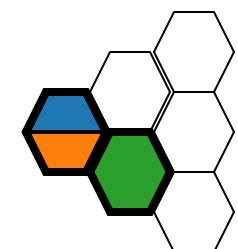
→ HyperQueue → PBS →



PBS + HQ  
+autoalloc



→ HyperQueue → PBS →



# Automatic allocation

- Automatic submit of PBS/Slurm jobs

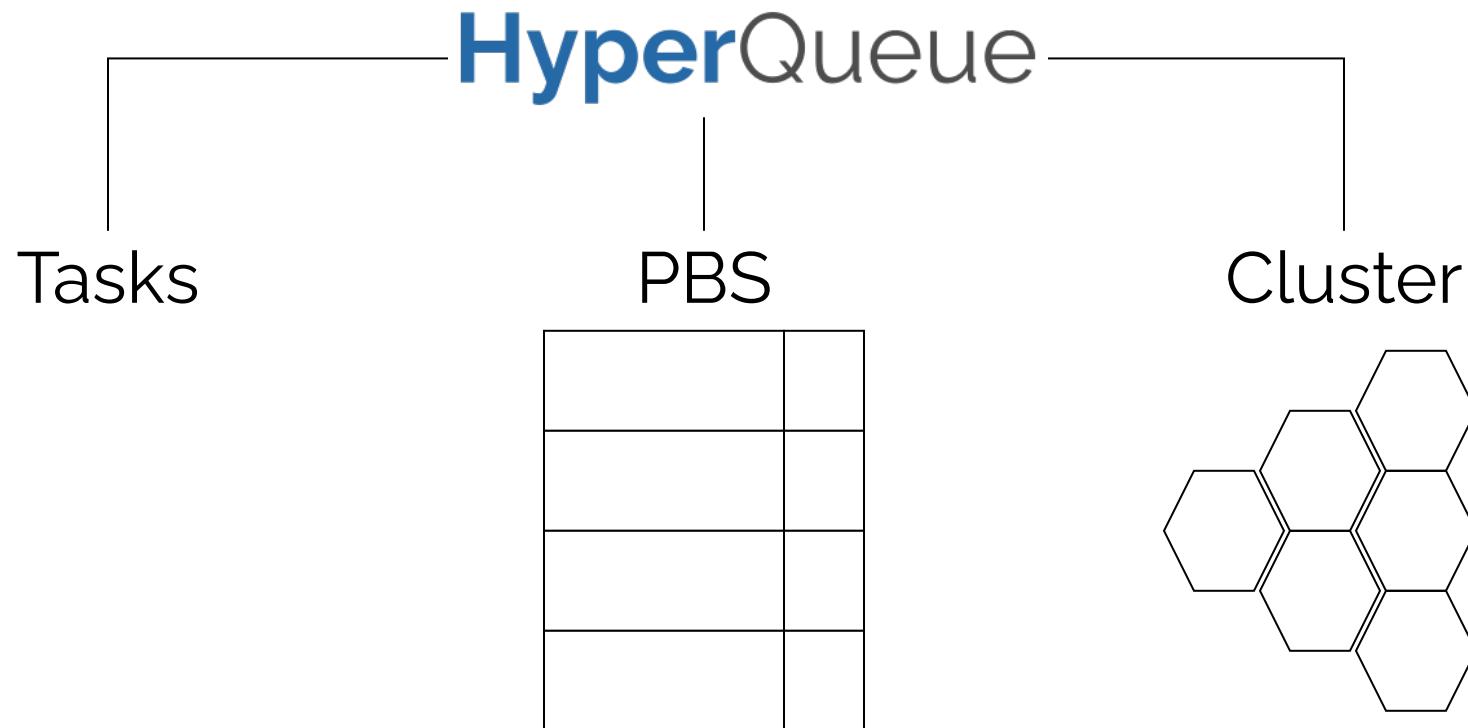
# Automatic allocation

- Automatic submit of PBS/Slurm jobs
- Based on computational demand (submitted tasks)

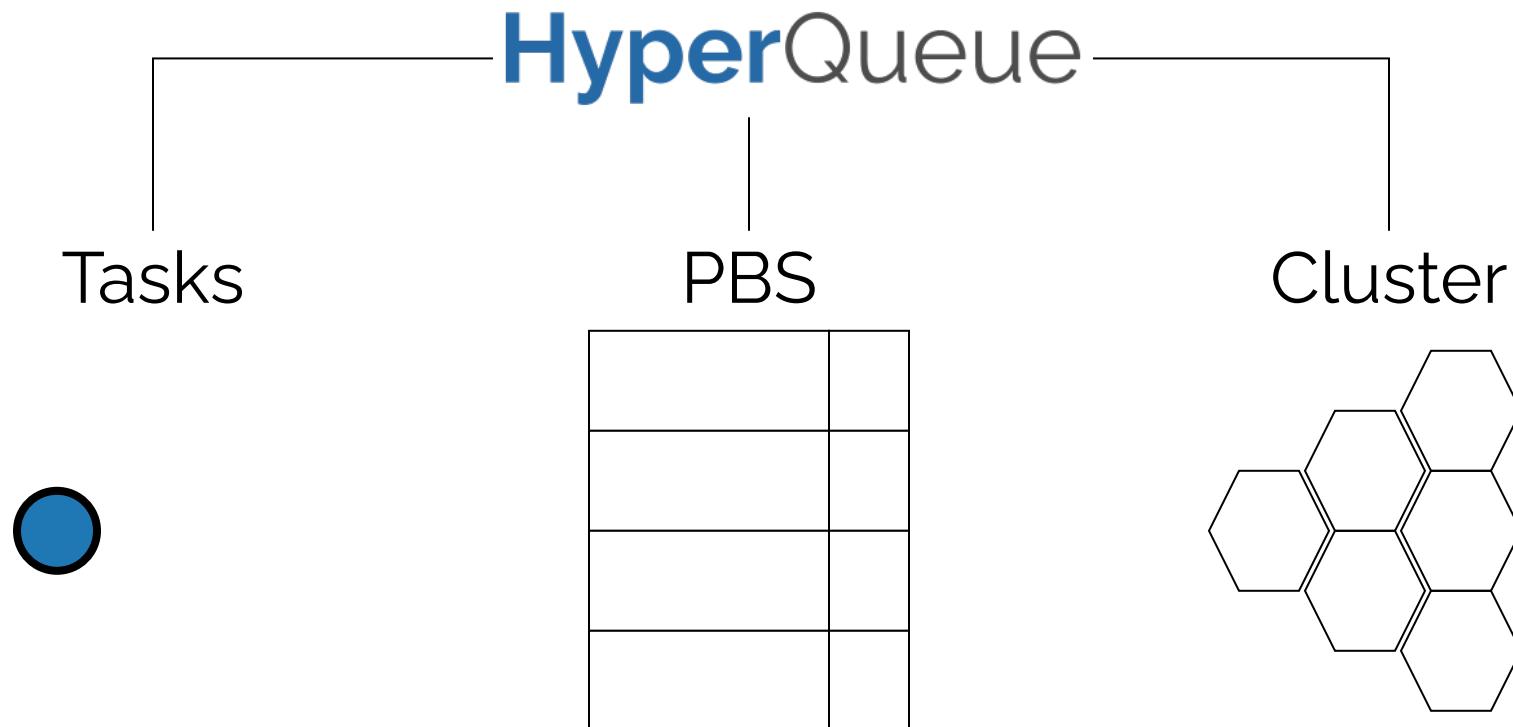
# Automatic allocation

- Automatic submit of PBS/Slurm jobs
- Based on computational demand (submitted tasks)
- To be further improved and analysed

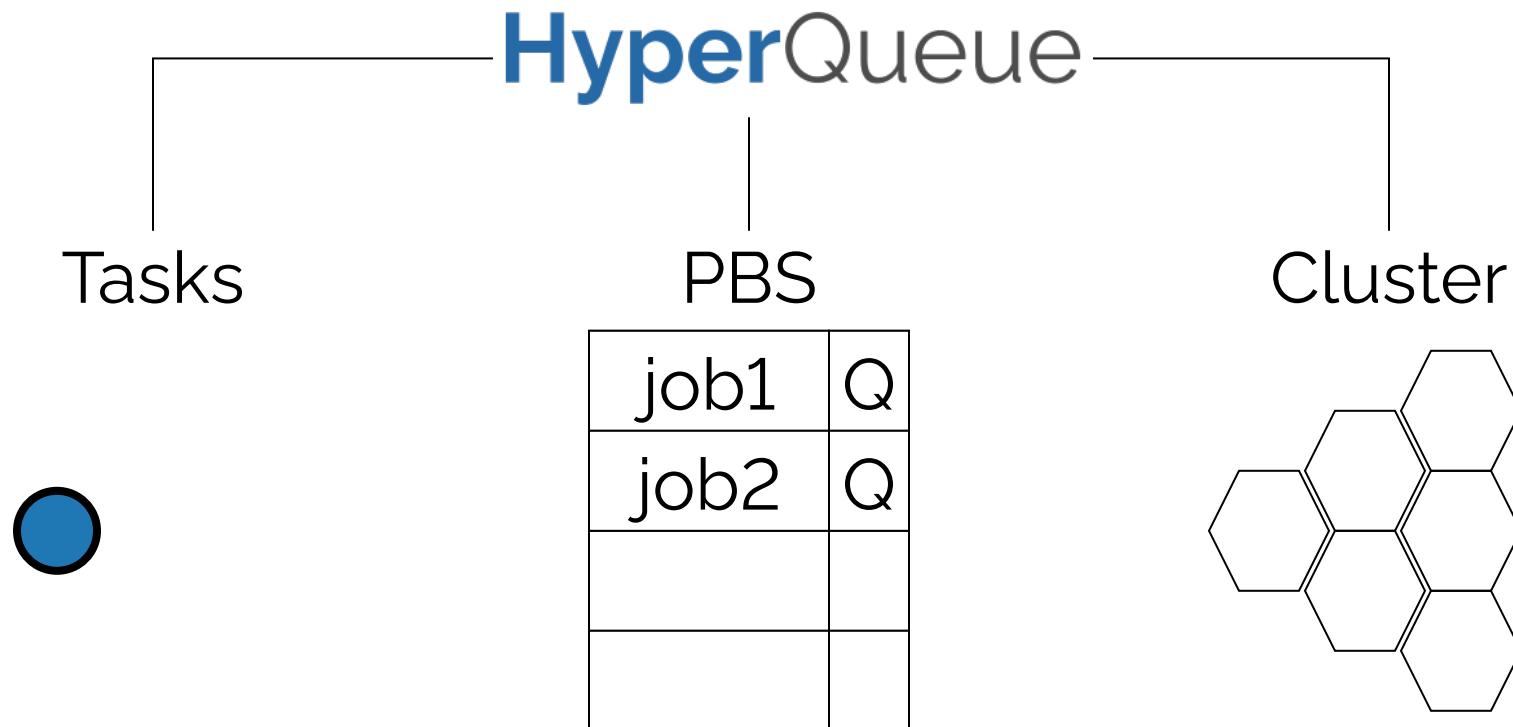
# Automatic allocation



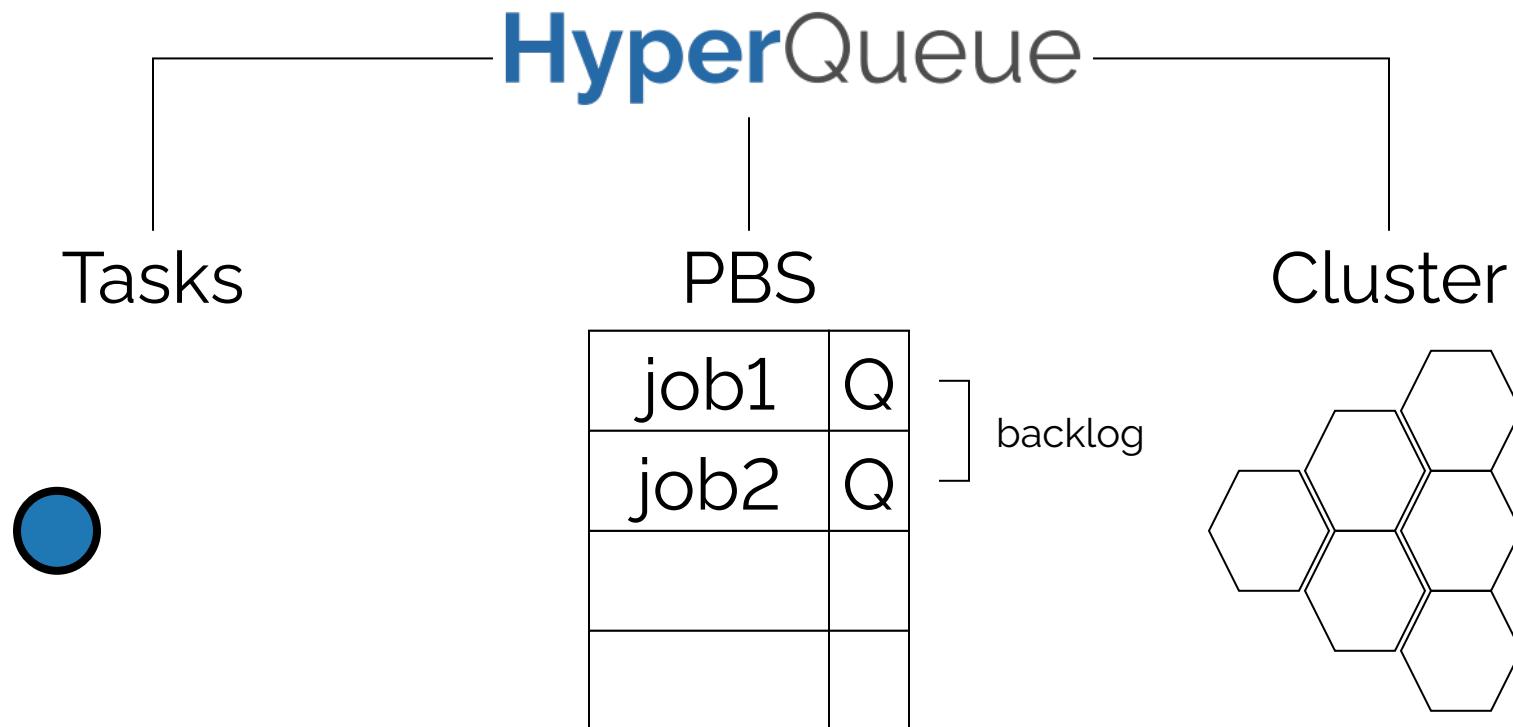
# Automatic allocation



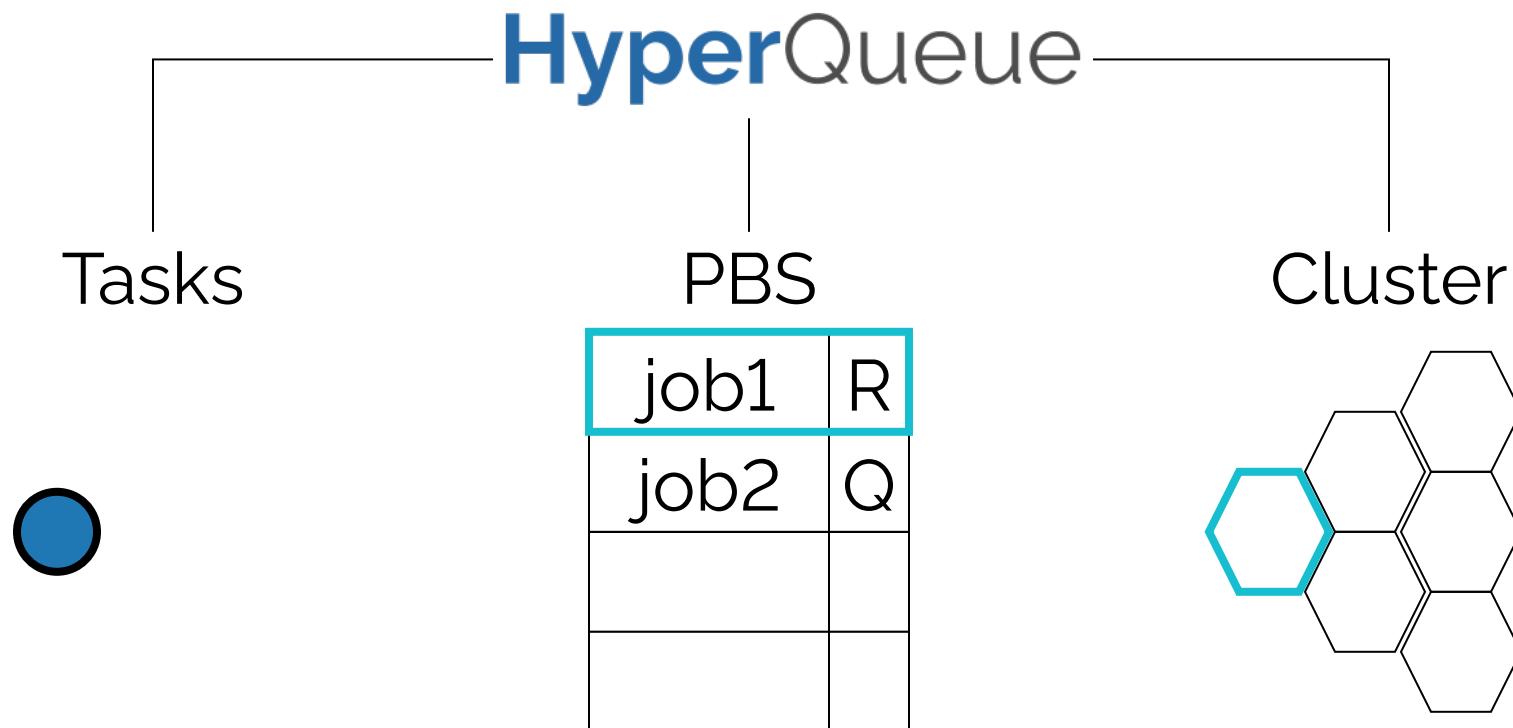
# Automatic allocation



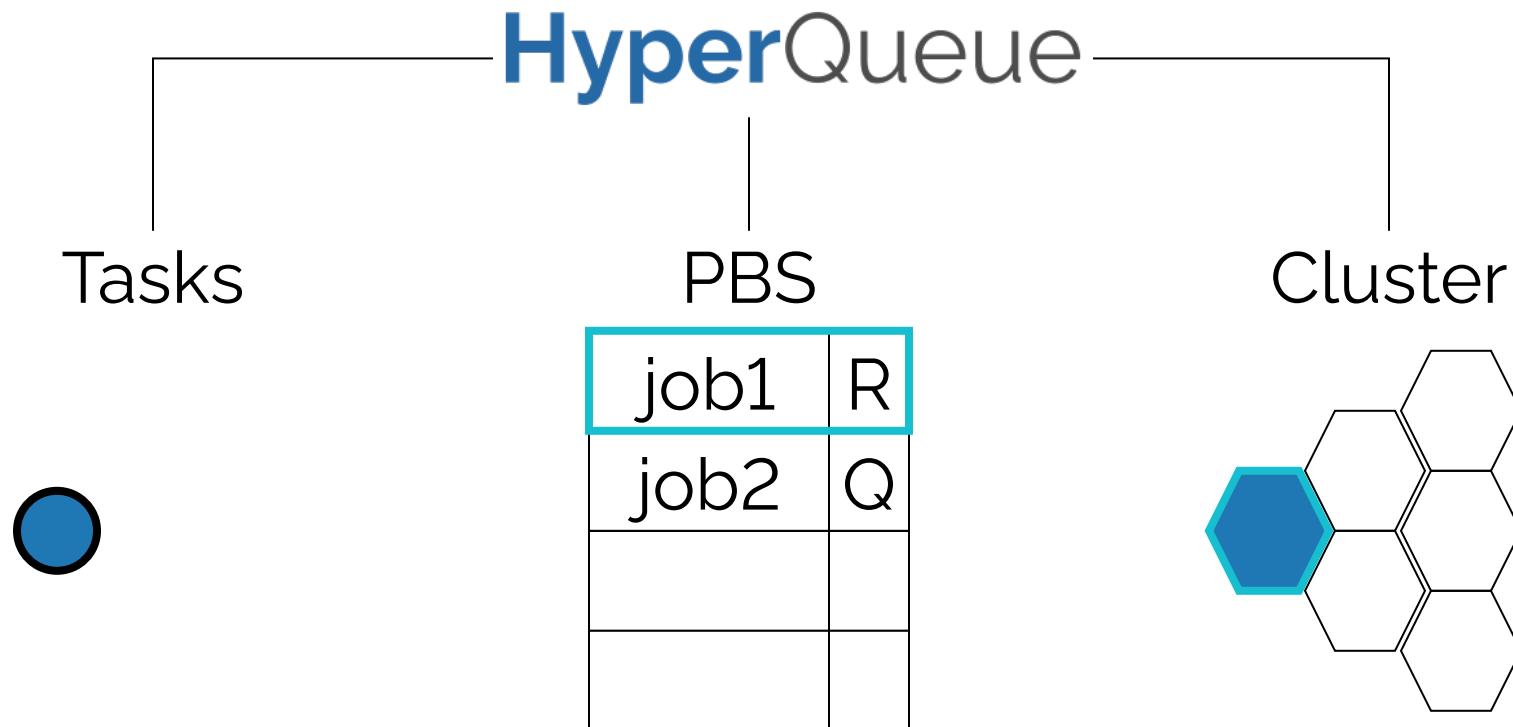
# Automatic allocation



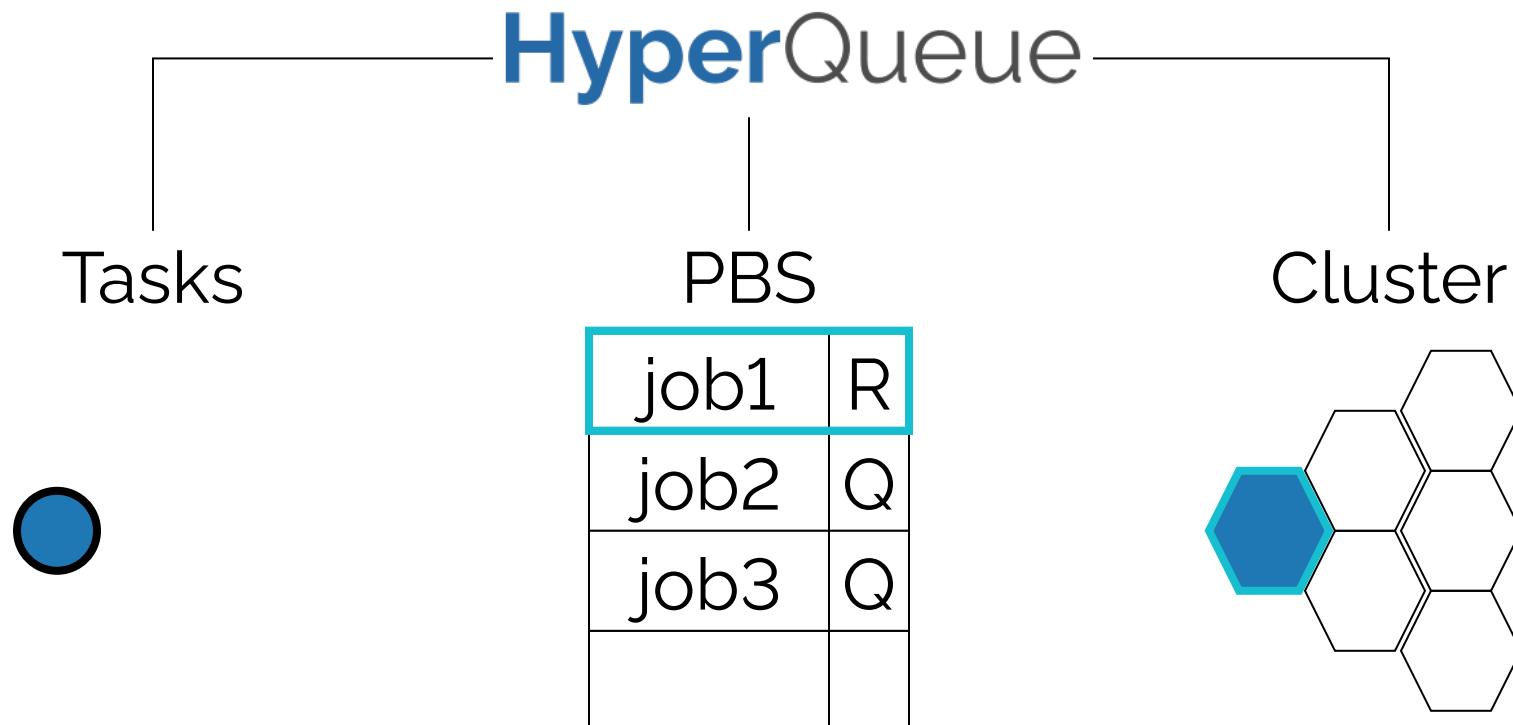
# Automatic allocation



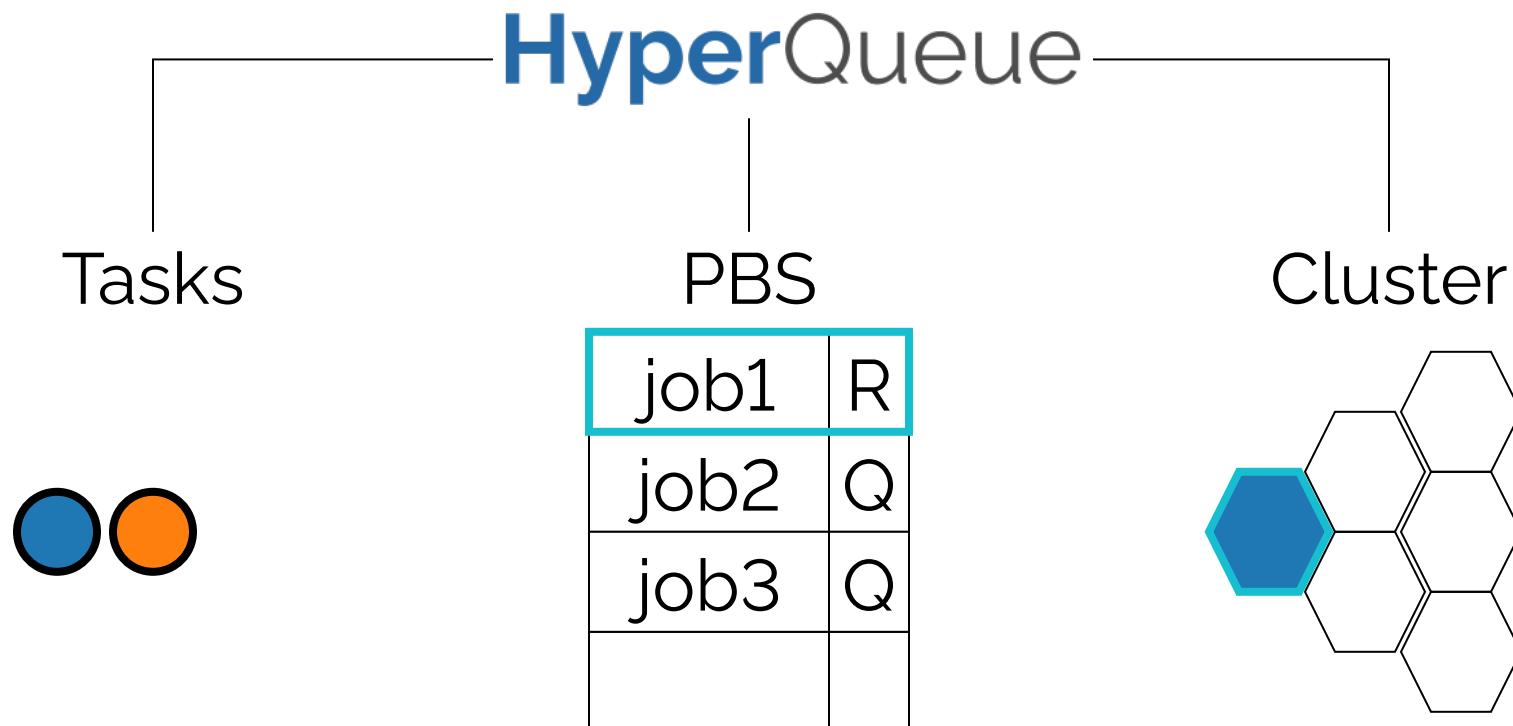
# Automatic allocation



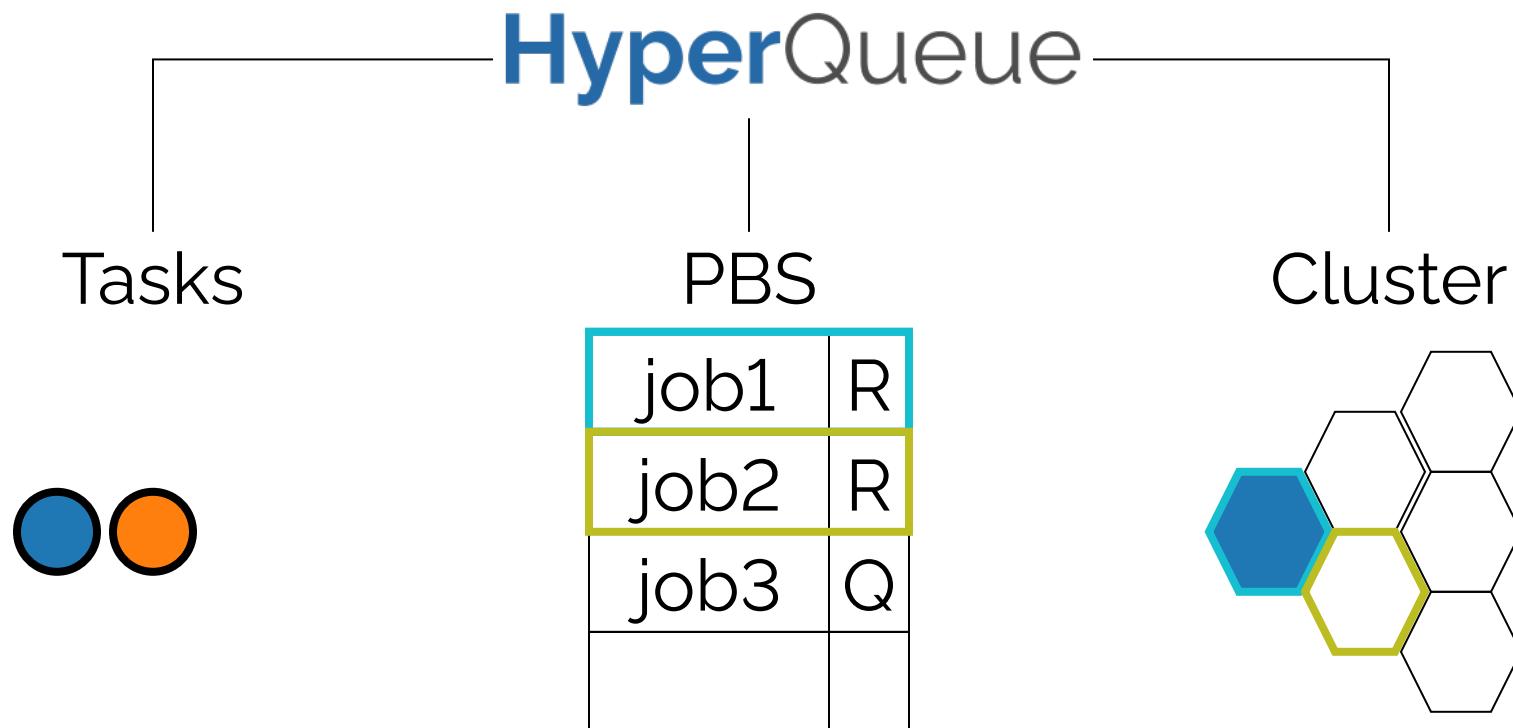
# Automatic allocation



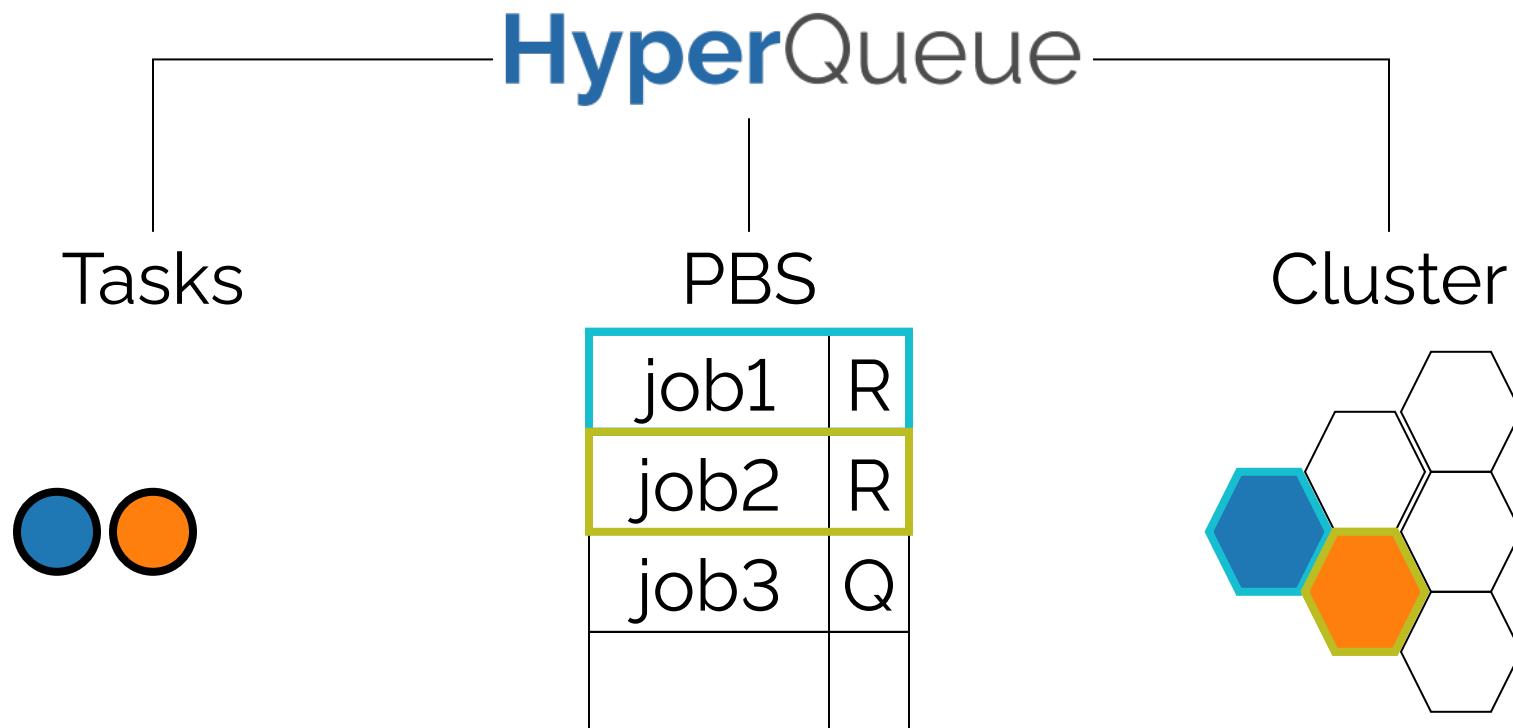
# Automatic allocation



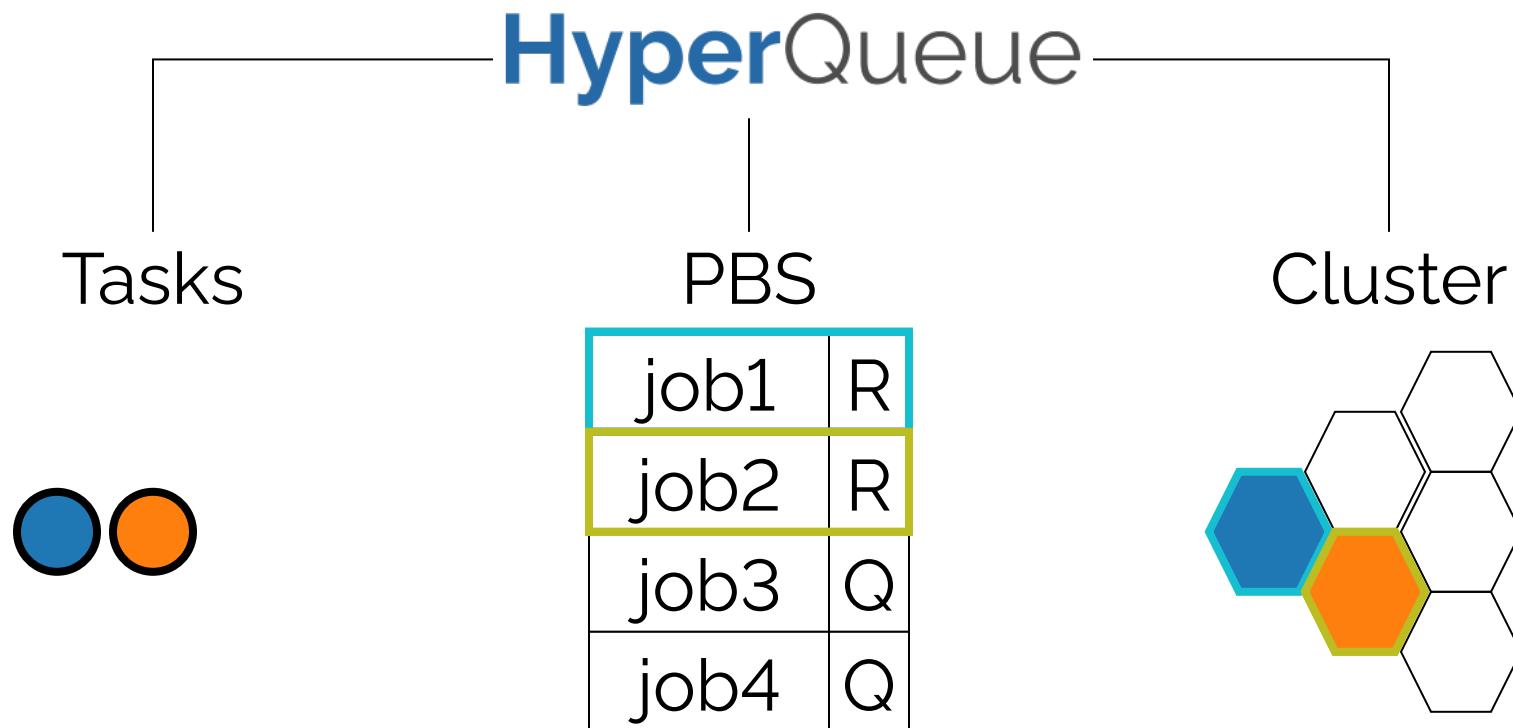
# Automatic allocation



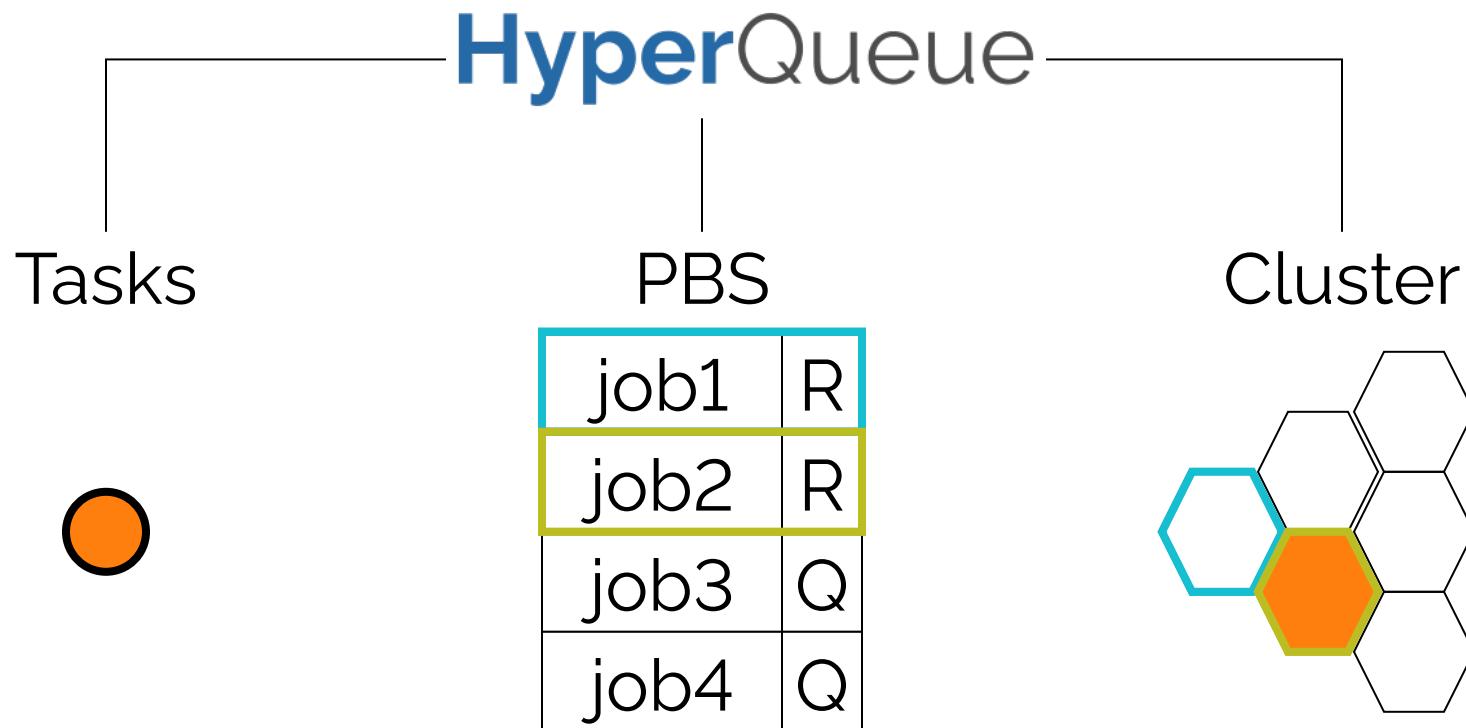
# Automatic allocation



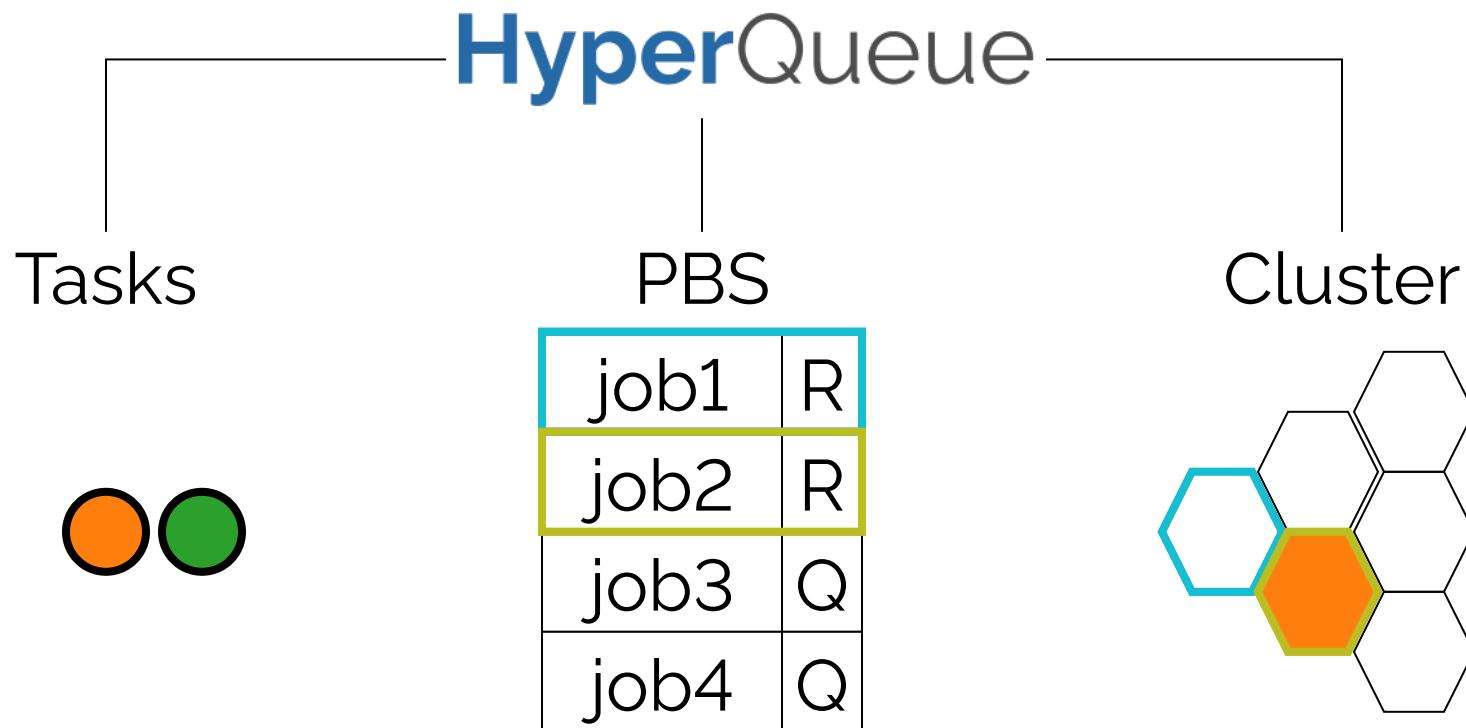
# Automatic allocation



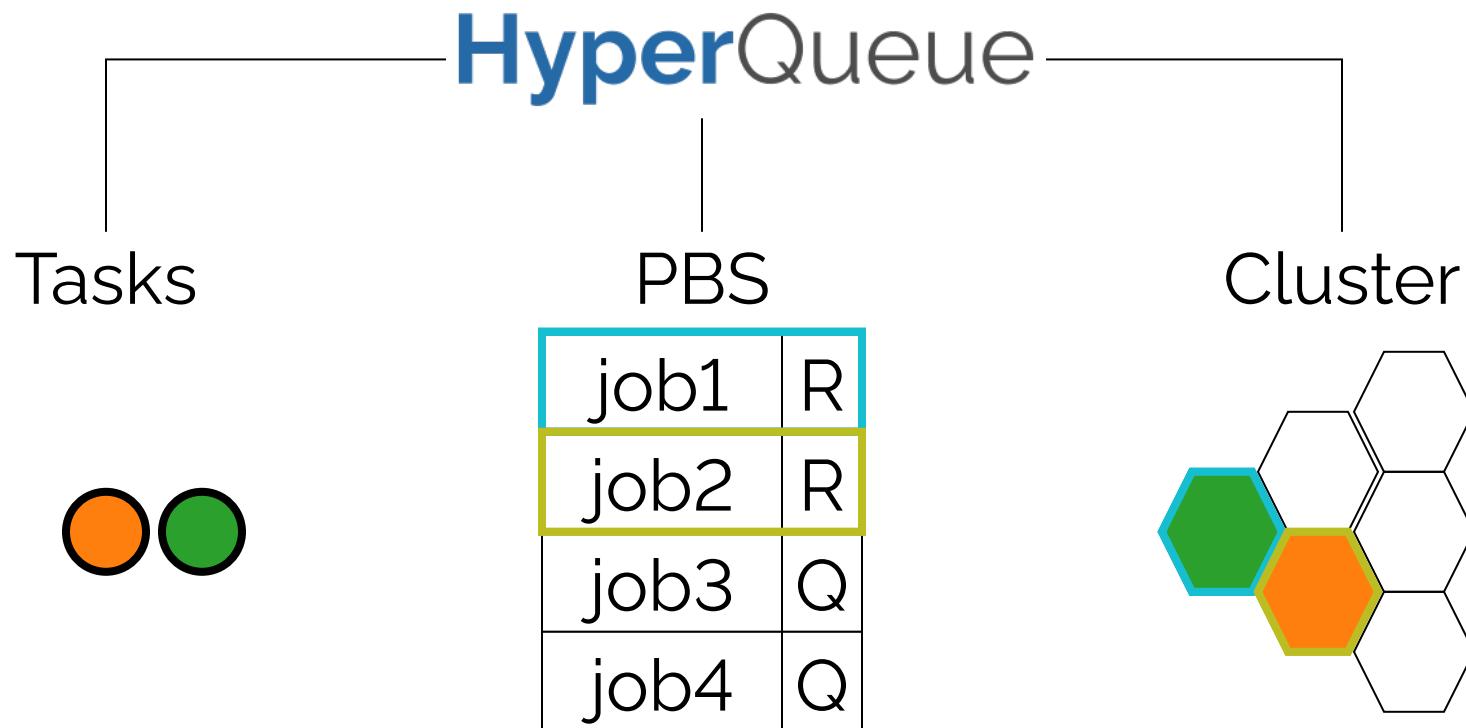
# Automatic allocation



# Automatic allocation



# Automatic allocation



# Challenge: fine-grained resources

- Arbitrary resources per task and worker

# Challenge: fine-grained resources

- Arbitrary resources per task and worker
- Resource-aware scheduler

# Challenge: fine-grained resources

- Arbitrary resources per task and worker
- Resource-aware scheduler
- NUMA support

# Challenge: fine-grained resources

- Arbitrary resources per task and worker
- Resource-aware scheduler
- NUMA support

```
$ hq worker start --resource gpus=range(1-3)
```

# Challenge: fine-grained resources

- Arbitrary resources per task and worker
- Resource-aware scheduler
- NUMA support

```
$ hq worker start --resource gpus=range(1-3)
```

```
$ hq submit --cpus=4 --resources gpus=2 ...
```

# Challenge: multi-node tasks

- Initial support for multi-node tasks

# Challenge: multi-node tasks

- Initial support for multi-node tasks

```
$ hq submit --nodes=4 ...
```

# **Challenge: efficiency**

- Low overhead per task (~0.1ms)

# Challenge: efficiency

- Low overhead per task (~0.1ms)
- I/O streaming: `stdout/stderr` over network

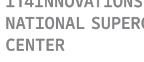
# HyperQueue in the wild

- H2020 EU projects
  -  **LIGATE** Complex workflows (GROMACS + LiGen)
  -  **EVEREST** Heterogeneous nodes (traffic simulator)
  -  **ACROSS** Multi-node tasks

# HyperQueue in the wild

- H2020 EU projects
  -  **LIGATE** Complex workflows (GROMACS + LiGen)
  -  **EVEREST** Heterogeneous nodes (traffic simulator)
  -  **ACROSS** Multi-node tasks
- HPC centers
  -  **VSB TECHNICAL UNIVERSITY OF OSTRAVA** | **IT4INNOVATIONS NATIONAL SUPERCOMPUTING CENTER**
  -  **CINECA**
  -  **LUMI**

# HyperQueue in the wild

- H2020 EU projects
  -  **LIGATE** Complex workflows (GROMACS + LiGen)
  -  **EVEREST** Heterogeneous nodes (traffic simulator)
  -  **ACROSS** Multi-node tasks
- Tool integrations
  -  **nextflow**
  -  **AiiDA**
  - ERT
- HPC centers
  -  **VSB TECHNICAL UNIVERSITY OF OSTRAVA** |  **IT4INNOVATIONS NATIONAL SUPERCOMPUTING CENTER**
  -  **CINECA**
  -  **LUMI**

# HyperQueue in the wild

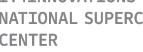
- H2020 EU projects

-  **LIGATE** Complex workflows (GROMACS + LiGen)
-  **EVEREST** Heterogeneous nodes (traffic simulator)
-  **ACROSS** Multi-node tasks

- Tool integrations

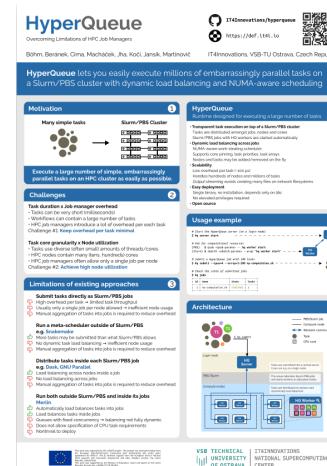
-  **nextflow**
-  **AiiDA**
- **ERT**

- HPC centers

-  **VSB TECHNICAL UNIVERSITY OF OSTRAVA** |  **IT4INNOVATIONS**  
NATIONAL SUPERCOMPUTING CENTER
-  **CINECA**
-  **LUMI**

- Research

- **ATLAS experiment (CERN)**  
ARC-CE+HyperQueue based submission system of ATLAS jobs for the Karolina HPC
- **Supercomputing'21 poster**



# Recap: HyperQueue

- **Ergonomics**
  - Job manager bypass
  - Fine-grained resource requirements
  - Multi-node tasks
  - ...

# Recap: HyperQueue

- **Ergonomics**
  - Job manager bypass
  - Fine-grained resource requirements
  - Multi-node tasks
  - ...
- **Efficiency**
  - Low overhead per task
  - Efficient scheduler
  - ...

# Objectives + plan

- Identify HPC workflow challenges ✓

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓
  - **HyperQueue performance study** •

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓
  - **HyperQueue performance study** •
- Analyze ergonomics

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓
  - **HyperQueue performance study** •
- Analyze ergonomics
  - HyperQueue ✓

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓
  - **HyperQueue performance study** •
- Analyze ergonomics
  - HyperQueue ✓
  - **Automatic allocation analysis** •

# Objectives + plan

- Identify HPC workflow challenges ✓
- Analyze efficiency
  - Schedulers ✓
  - Dask ✓
  - **HyperQueue performance study** •
- Analyze ergonomics
  - HyperQueue ✓
  - **Automatic allocation analysis** •
- Prepare HyperQueue publication

# Publications related to thesis

- **Analysis of workflow schedulers in simulated distributed environments**  
*Jakub Beránek, Stanislav Böhm, Vojtěch Cima*  
The Journal of Supercomputing 2022
- **Runtime vs Scheduler: Analyzing Dask's Overheads**  
*Stanislav Böhm, Jakub Beránek*  
IEEE/ACM Workflows in Support of Large-Scale Science (WORKS) 2020

# Publications related to thesis

- **Analysis of workflow schedulers in simulated distributed environments**  
*Jakub Beránek, Stanislav Böhm, Vojtěch Cima*  
The Journal of Supercomputing 2022
- **Runtime vs Scheduler: Analyzing Dask's Overheads**  
*Stanislav Böhm, Jakub Beránek*  
IEEE/ACM Workflows in Support of Large-Scale Science (WORKS) 2020

(collaboration with ETH Zurich)

- **Network-Accelerated Non-Contiguous Memory Transfer**  
S. Di Girolamo, K. Taranov, A. Kurth, M. Schaffner, T. Schneider, *J. Beránek*, M. Besta,  
L. Benini, D. Roweth, T. Hoefler  
*SC (International Conference for High Performance Computing, Networking, Storage and Analysis) 2019*
- **A RISC-V in-Network Accelerator for Flexible High-Performance Low-Power Packet Processing**  
S. Di Girolamo, A. Kurth, A. Calotoiu, T. Benz, T. Schneider, *J. Beránek*, L. Benini, T. Hoefler  
*ISCA (International Symposium on Computer architecture) 2021*

# Publications unrelated to thesis

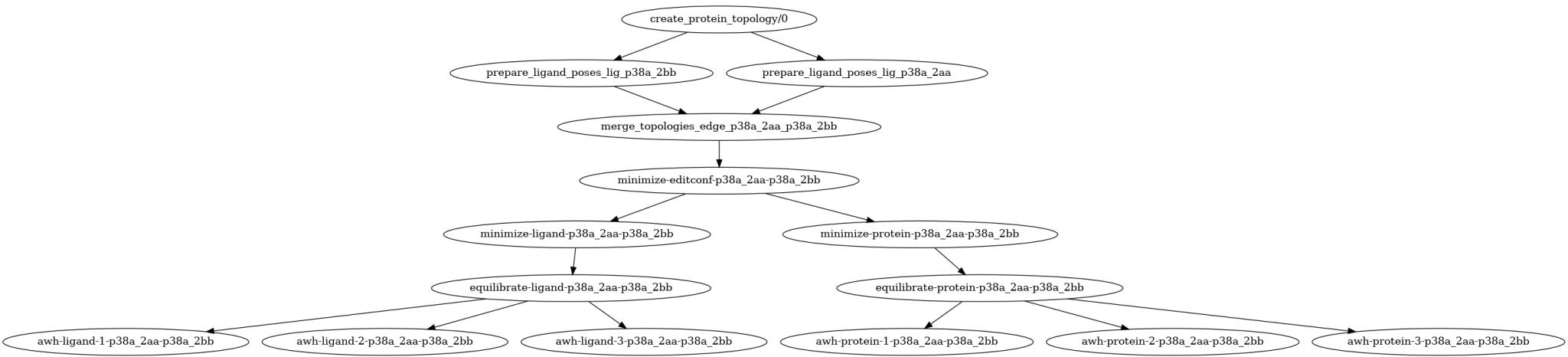
(collaboration with ETH Zurich)

- **SISA: Set-Centric Instruction Set Architecture for Graph Mining on Processing-in-Memory System**  
M. Besta, R. Kanakagiri, G. Kwasniewski, R. Ausavarungnirun, J. Beránek, K. Kanellopoulos, K. Janda, Z. Vonarburg-Shmaria, L. Gianinazzi, I. Stefan, J. G. Luna, J. Golinowski, M. Copik, L. Kapp-Schwoerer, S. Di Girolamo, N. Blach, M. Konieczny, O. Mutlu, T. Hoefler  
*IEEE/ACM MICRO (International Symposium on Microarchitecture) 2021*
- **GraphMineSuite: Enabling High-Performance and Programmable Graph Mining Algorithms with Set Algebra**  
M. Besta, Z. Vonarburg-Shmaria, Y. Schaffner, L. Schwarz, G. Kwasniewski, L. Gianinazzi, J. Beránek, K. Janda, T. Holenstein, S. Leisinger, P. Tatkowski, E. Ozdemir, A. Balla, M. Copik, P. Lindenberger, M. Konieczny, O. Mutlu, T. Hoefler  
*PVLDB 2021*
- **Streaming Message Interface: High-Performance Distributed Memory Programming on Reconfigurable Hardware**  
T. De Matteis, J. de Fine Licht, J. Beránek, T. Hoefler  
*SC (International Conference for High Performance Computing, Networking, Storage and Analysis) 2019*
- **Haydi: Rapid Prototyping and Combinatorial Objects**  
Stanislav Böhm, Jakub Beránek, Martin Šurkovský  
*Foundations of Information and Knowledge Systems 2018*
- **Alternative Paths Reordering Using Probabilistic Time-Dependent Routing**  
M. Golasowski, J. Beránek, M. Šurkovský, L. Rapant, D. Szturcová, J. Martinovič, Kateřina Slaninová  
*Advances in Networked-based Information Systems 2020*
- **A Distributed Environment for Traffic Navigation Systems**  
J. Martinovič, M. Golasowski, K. Slaninová, J. Beránek, M. Šurkovský, L. Rapant, D. Szturcová, R. Cmar  
*Complex, Intelligent, and Software Intensive Systems 2020*

# Thank you for your attention

Slides made with <https://github.com/spirali/elsie>

# GROMACS + LiGEN LIGATE workflow



# Projects and publications

- ANTAREX: distributed system for traffic simulator
- EVEREST: traffic simulator + HQ
- LIGATE: porting GROMACS + LiGEN workflow to HQ
- ETH internship: 5 publications (2 at SC'19)