

# Distributed Global Scheduling in Datacenters

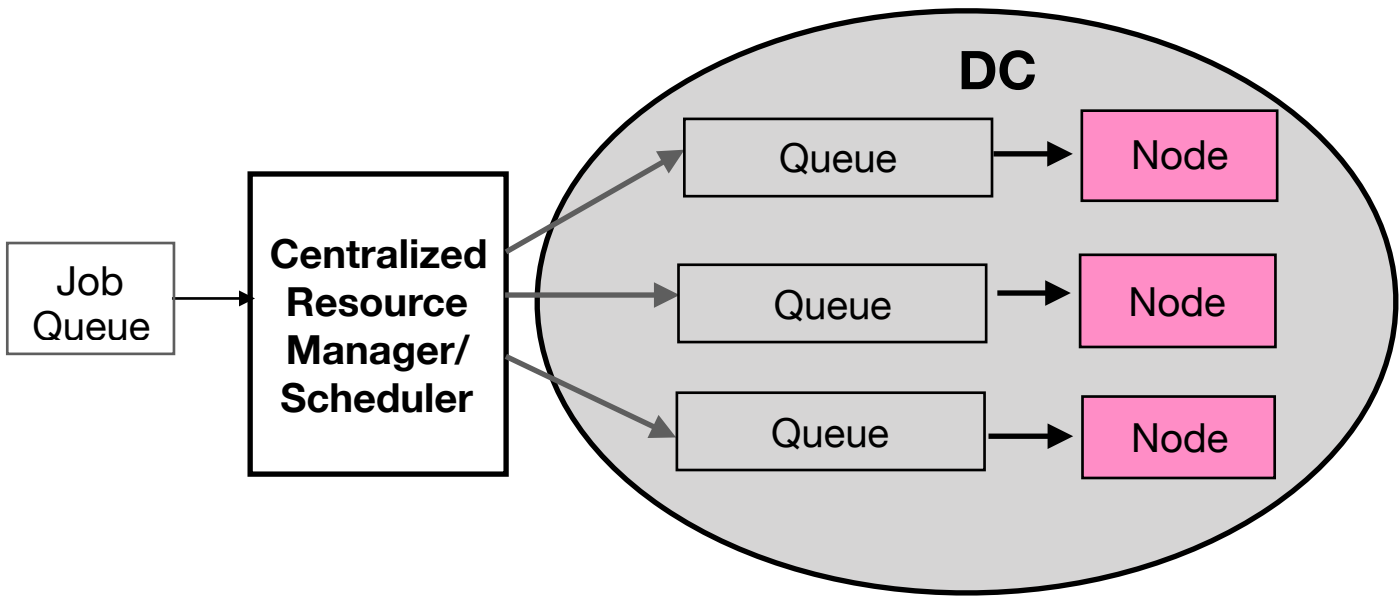
## State of the Art Scheduling

### Underutilised Datacenter Resources

- Azure**<sup>1</sup>
- ❖ 60% VMs have  $\leq 20\%$  CPU usage!
- Alibaba**<sup>2</sup>
- ❖ Average server CPU 50%
  - ❖ Memory  $\leq 60\%$
- Underutilisation is expensive!**<sup>3</sup>

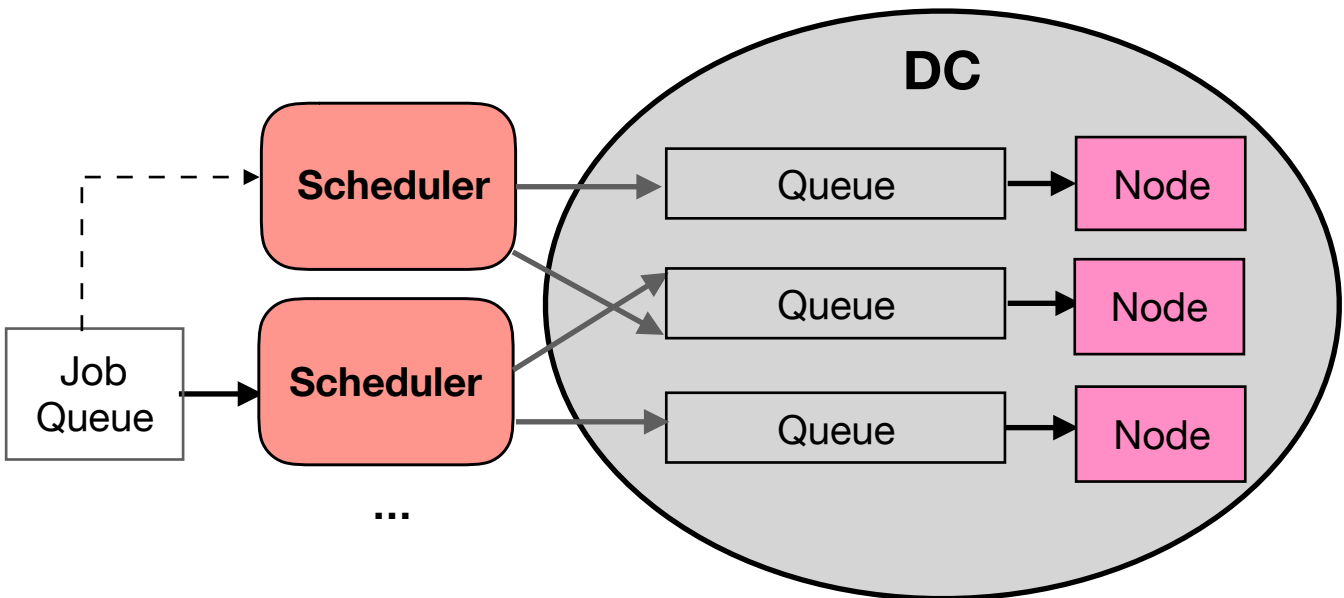
<sup>1</sup>[Resource Central, SOSP'17]  
<sup>2</sup>[https://github.com/alibaba/clusterdata]  
<sup>3</sup>[Scalable system scheduling for HPC and big data, JPDC, 17]

### Centralised



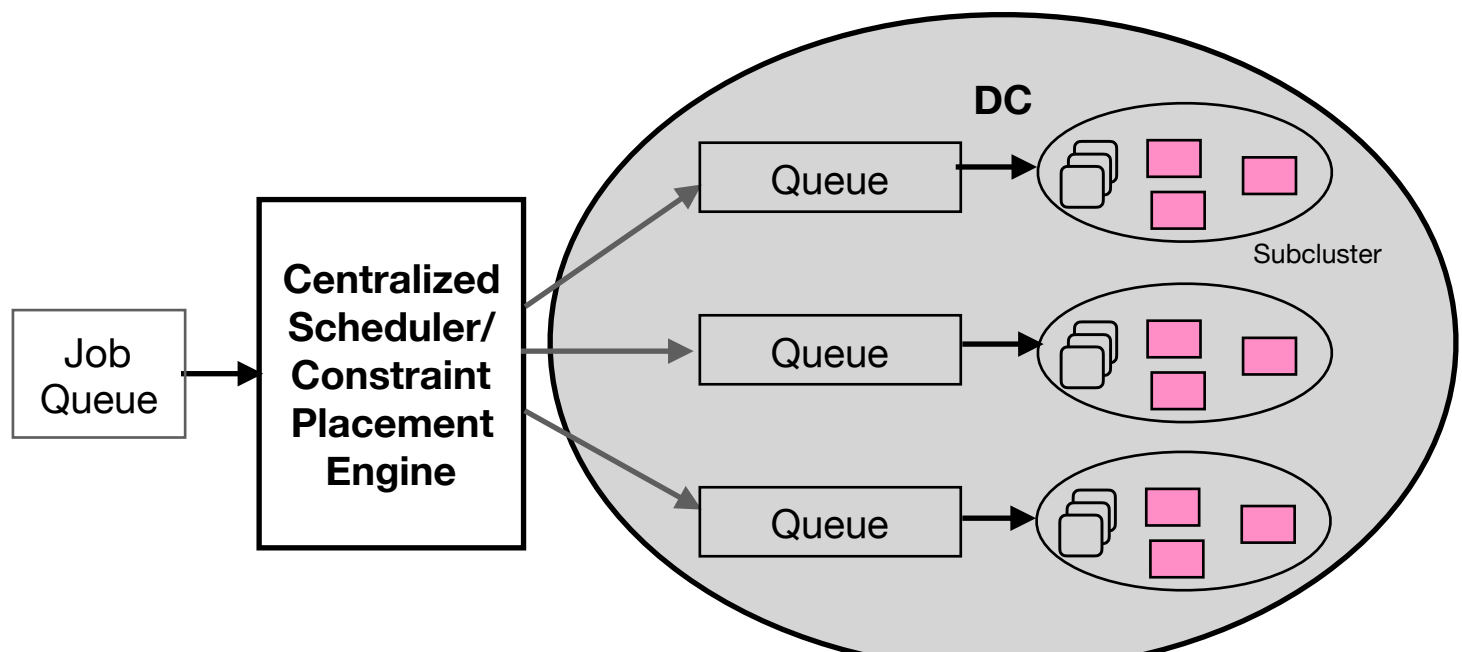
- Examples - Mesos [NSDI'11], Yarn, Apollo [OSDI'14]
- ✓ Global resource view
  - ✗ Scheduler can be a bottleneck
  - ✗ Delayed, high volumes of resource updates

### Decentralised



- Example - Sparrow [NSDI'14]
- ✓ Fast and simple
  - ✗ Unsuitable for long running jobs
  - ✗ Not globally optimal

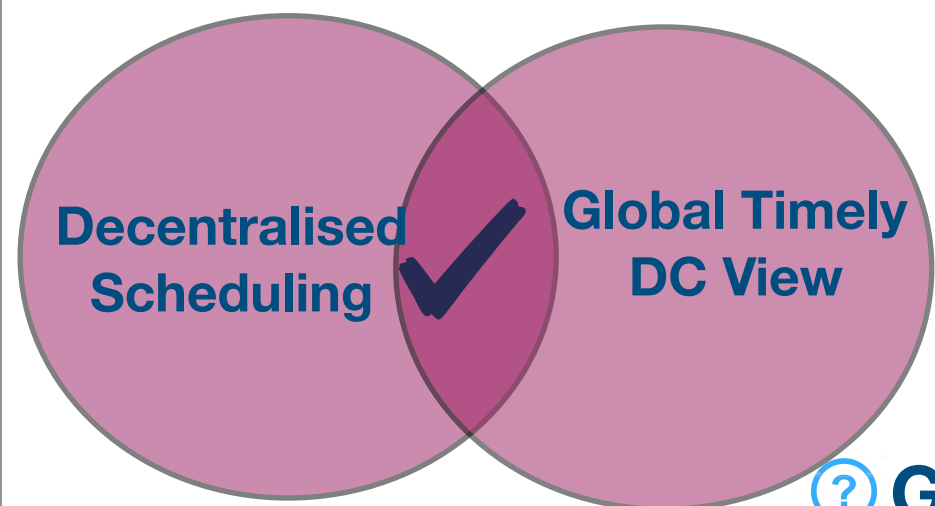
### Hybrid



- Example - Hydra [NSDI'19], Medea [EuroSys'18], Borg [EuroSys'15]
- ✓ Policy-driven job/task placement
  - ✓ Less node information traffic
  - ✗ Centralised or decentralised components

## Proposed Direction

### Global Scheduling at Node Level



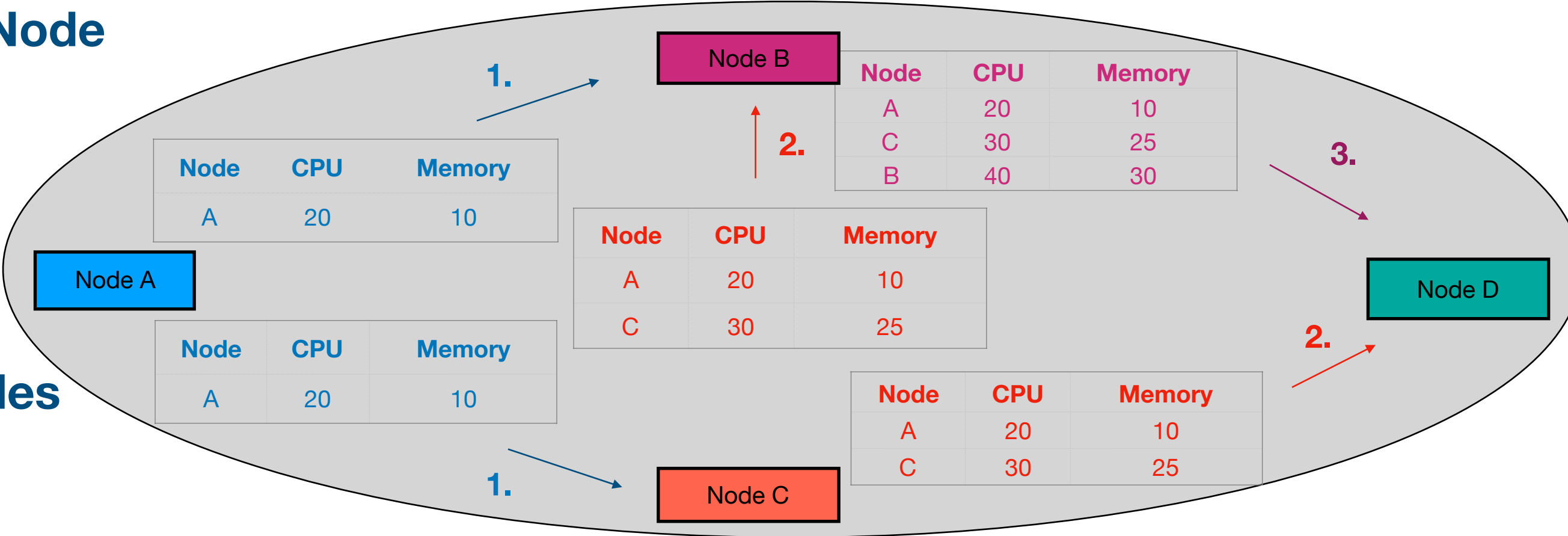
### Challenges

- ? Good for long and short jobs
- ? Volume and frequency of updates
- ? Time from local to global view

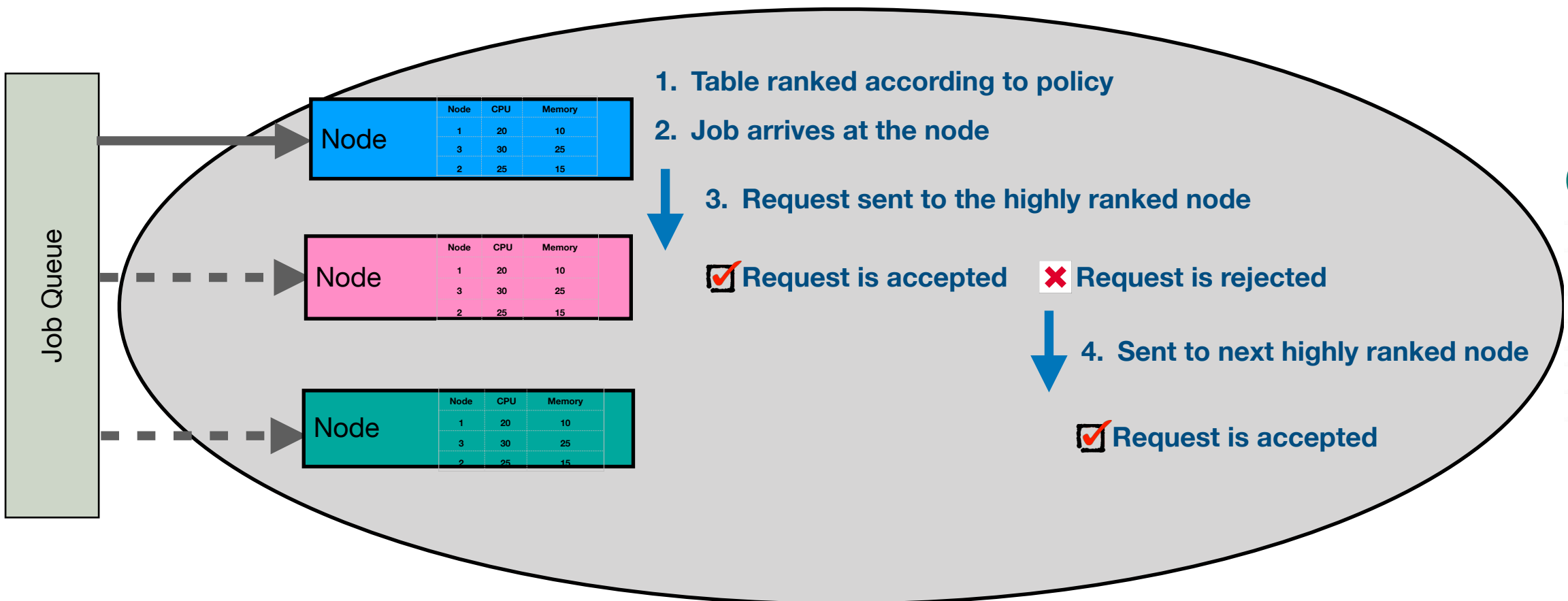
### Up-to-Date Global View at each Node

#### Inspired by routing protocols

- ✓ BGP, OSPF, ...
- ✓ Resource data propagation
- ✓ Global view convergence
- ✓ Same ranking policy across nodes



### Intra-DC Scheduling



### Scheduling Using "Up-to-Date" Global View

#### Challenges

- ? Collision avoidance
- ? Minimise inter-DC traffic
- ? Minimise scheduling time

### Inter-DC Scheduling

