# Solving (NP-Hard) Scheduling Problems with oVirt & OptaPlanner

Jason Brooks
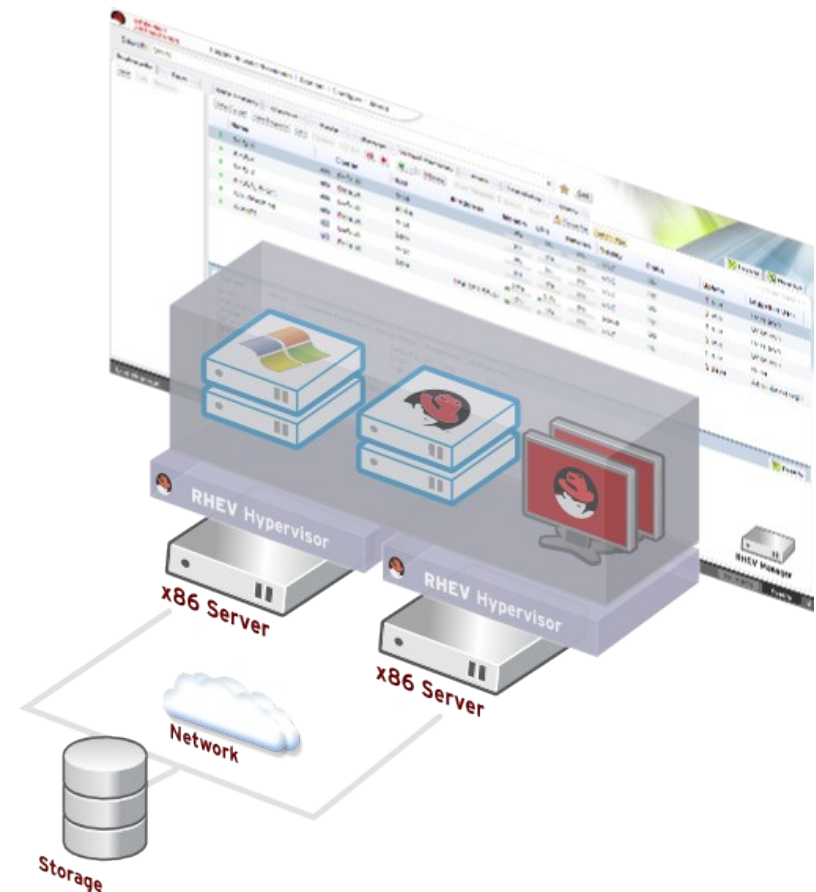Red Hat Open Source & Standards
SCALE13x, Feb 2015

# What Is oVirt?

Large scale, centralized management for server and desktop virtualization

Based on leading performance, scalability and security infrastructure technologies
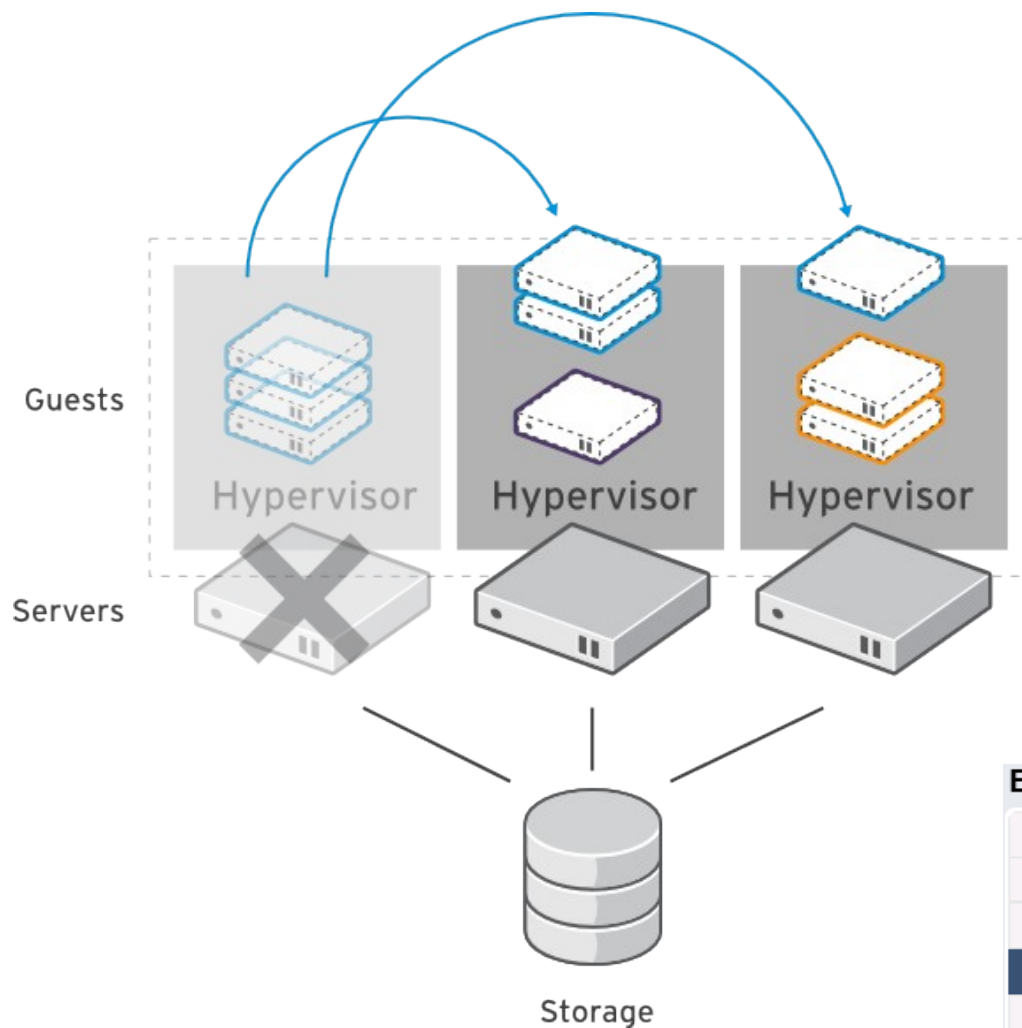
Provide an open source alternative to vCenter/vSphere
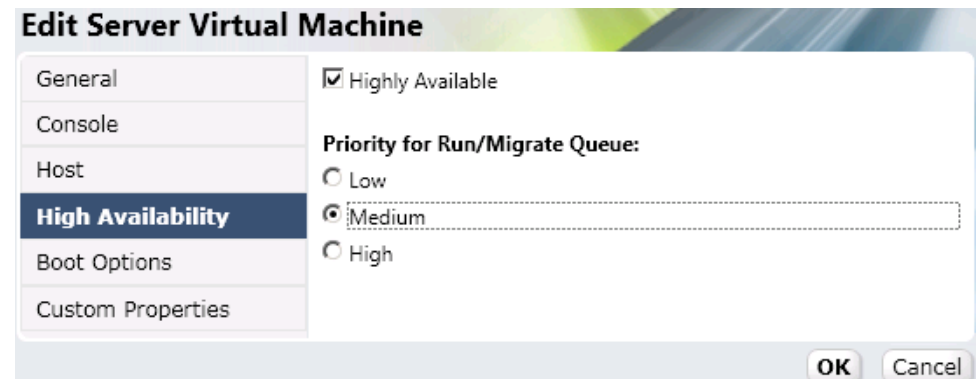
Focus on KVM for best integration/performance
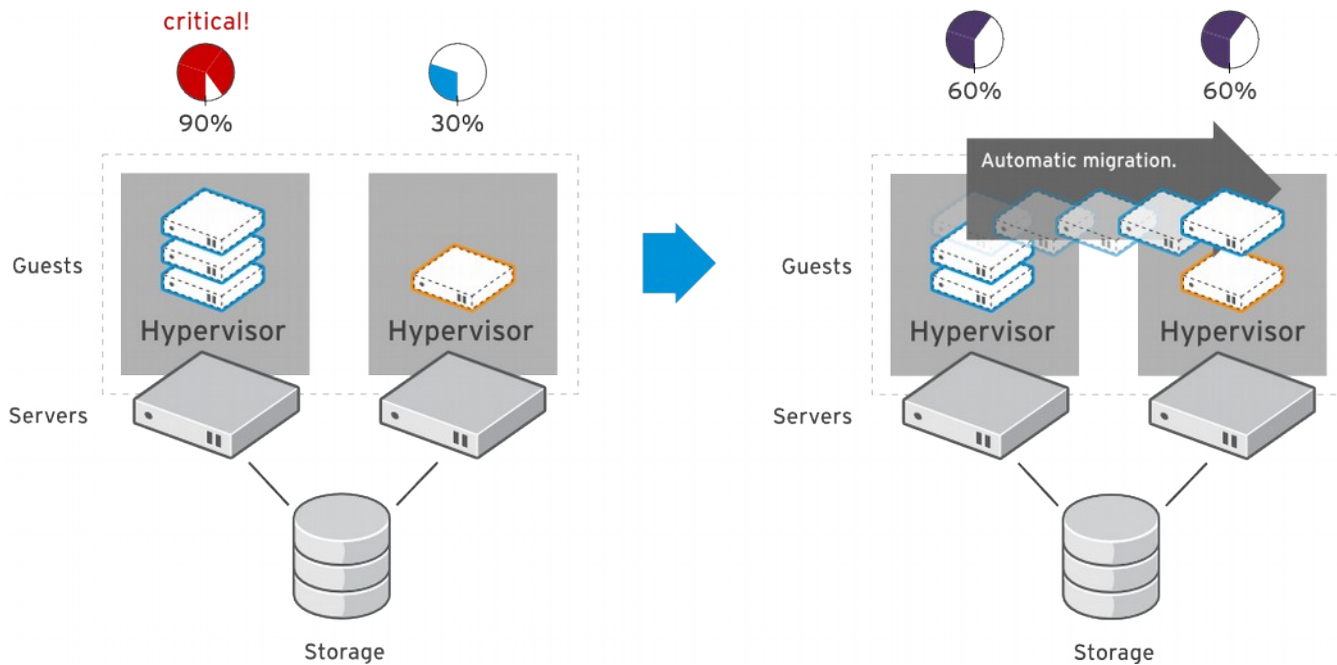
Focus on ease of use/deployment

# Virt & Cloud Scheduling

- Running a new VM

- Selecting migration destination

- Load balancing

# High Availability



- Build a highly available enterprise infrastructure

- Continually monitor host systems and virtual machines

- Automatically restart virtual machines in case of host failure

  - Restart virtual machine on another node in the cluster

- Use live migration to "fail-back" a VM to it's original host when the server is restored
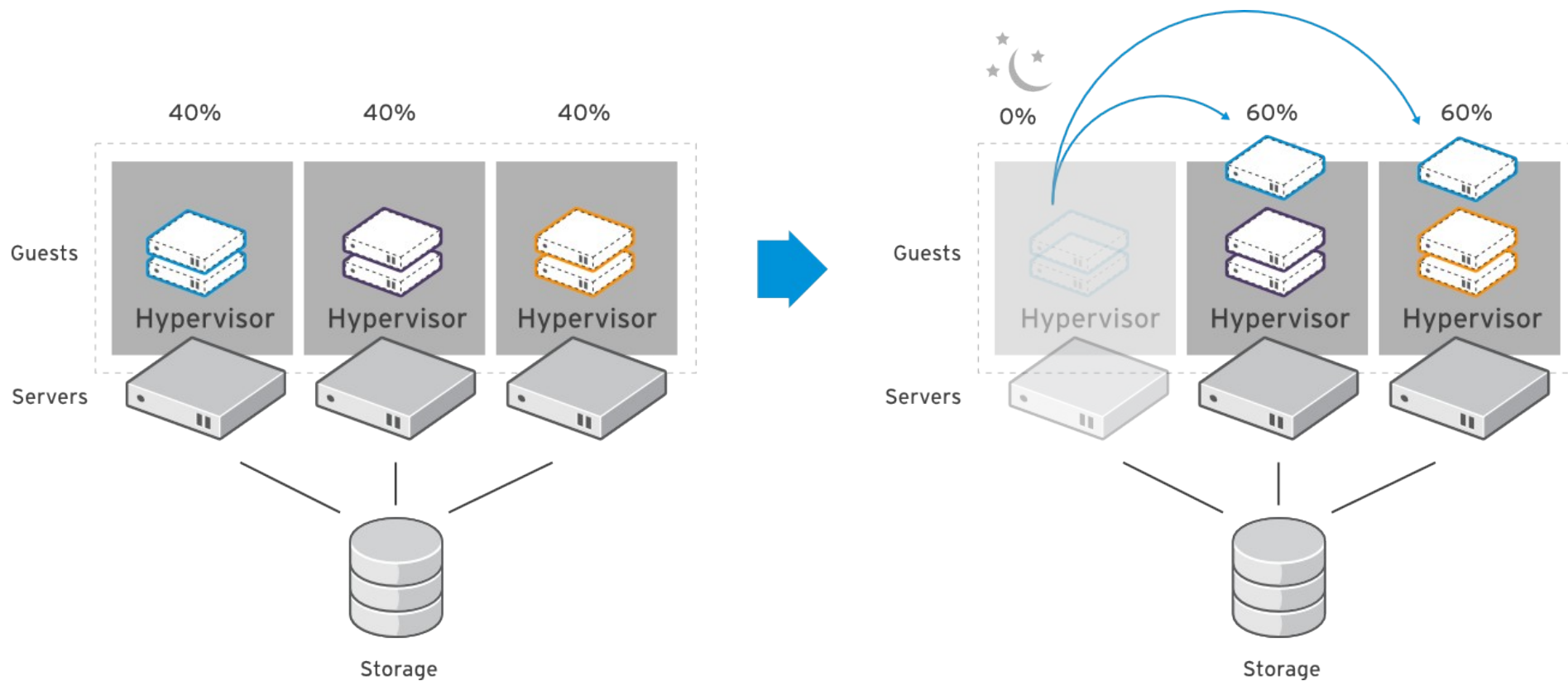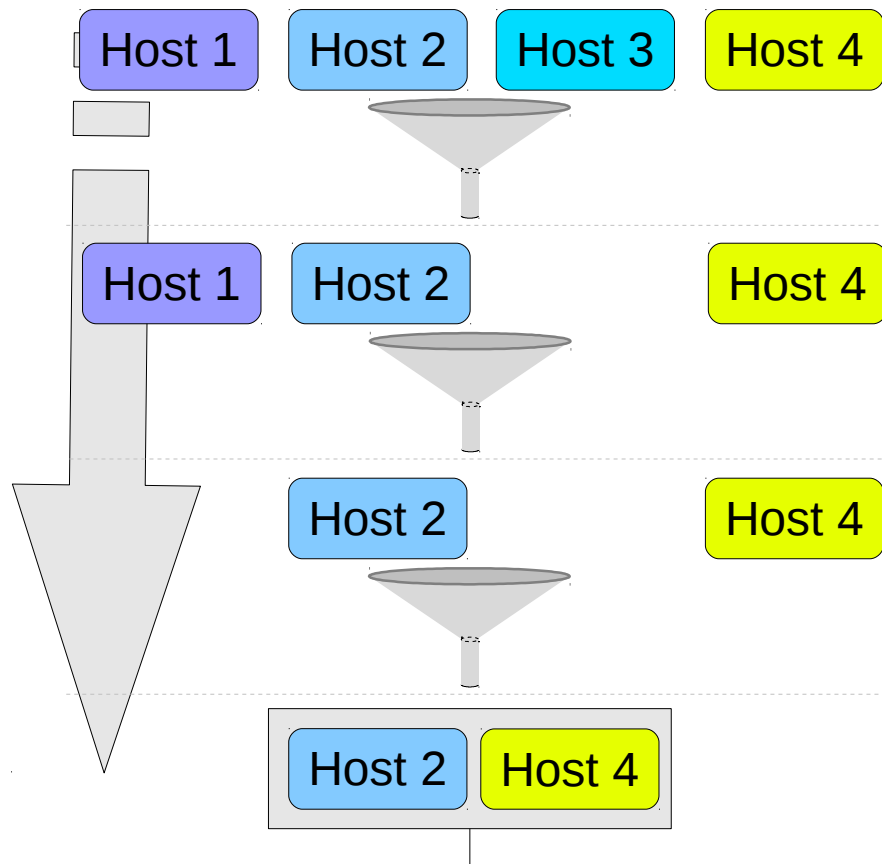
# System Scheduler



- Dynamically balance workloads in the data center.

- Automatically live migrate virtual machines based on resources

- Define custom policies for distribution of virtual machines

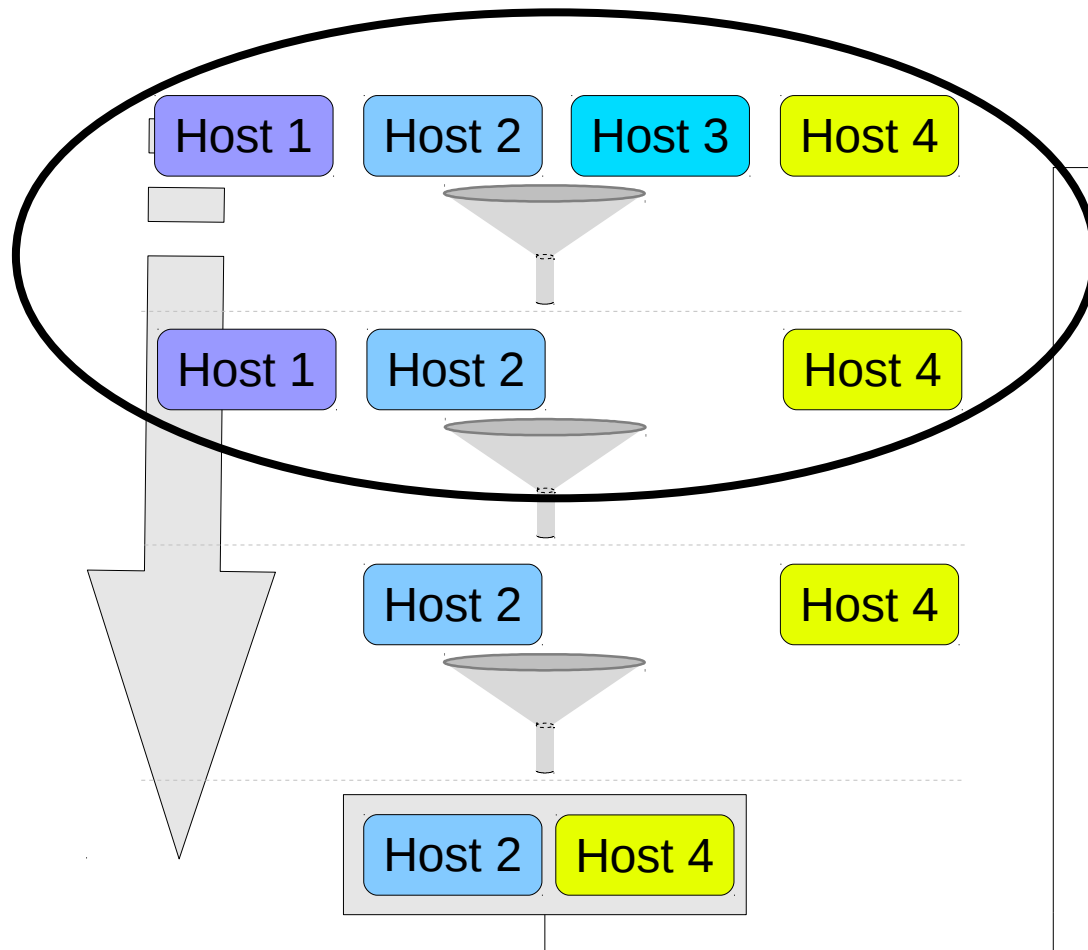Maintain consistent resource usage across the enterprise data center

Define policies to optimize workload on a fewer number of servers during "off-peak" hours

# Scheduling in oVirt
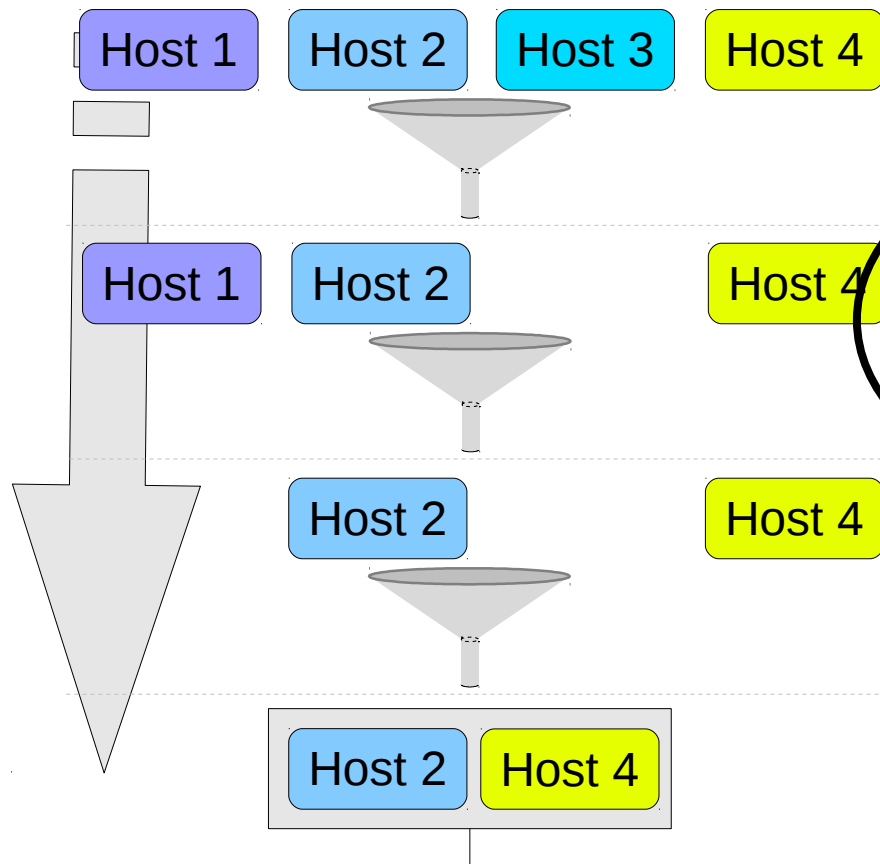
- 



| | func 1 | func 2 | sum |
|---|---|---|---|
| Factor | 5 | 2 | |
| Host 2 | 10 | 2 | 54 |
| Host 4 | 3 | 12 | **39*** |

**\*Host 4 sum: 3\*5+12\*2 = 39**

| | func 1 | func 2 | sum |
|---|---|---|---|
| Factor | 5 | 2 | |
| Host 2 | 10 | 2 | 54 |
| Host 4 | 3 | 12 | **39*** |

**\*Host 4 sum: 3*5+12*2 = 39**

# Weights



| | func 1 | func 2 | sum |
|---|---|---|---|
| Factor | 5 | 2 | |
| Host 2 | 10 | 2 | 54 |
| Host 4 | 3 | 12 | **39*** |

**\*Host 4 sum: 3\*5+12\*2 = 39**

# Balancers

- Triggers a scheduled task to determine which VM needs to be migrated to one of under-utilized hosts
- A single load balancing logic is allowed per cluster

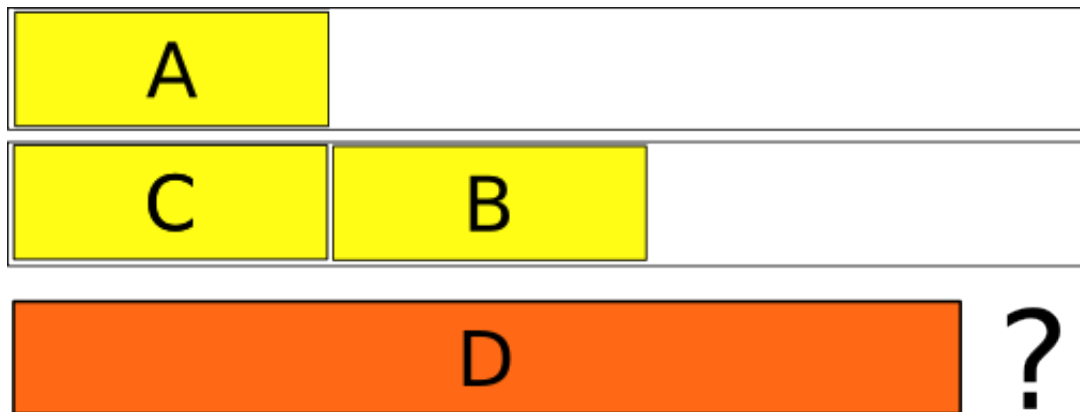# Filters, Weights, Balancers

# External Scheduler

- External service written in python and run as a separate process from the engine

- External service provides:

  - Engine safety

  - Should allow additional languages

  - Future option of scheduling as a service

**Filter Modules**  Drag or use context menu to make changes

Enabled Filters

CPU

Network

(EXT) max_vms

**Weights Modules**  Drag or use context menu to make changes

Enabled Weights & Factors

(EXT) even_vm_distribution

**Load Balancer**

vm_balance  (EXT)

# Optimizer Goals

- Better load balancing

- Configurable by existing cluster policy

- Separate machine to protect ovirt-engine

- Starting a VM that can't be placed directly
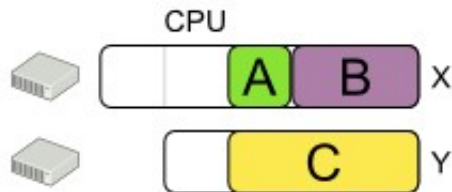  - Space needs to be created first

- Defined by set of machines and set of processes

- Each machine has some resources (CPU, RAM, ...)

- Each process requires resources

- NP-complete (variant of bin packing)

  - Easy to verify a given solution to a problem in reasonable time.

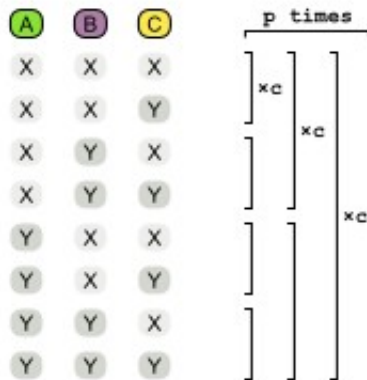  - There is no silver bullet to find the optimal solution of a problem in reasonable time (*).

Calculate the size of the search space

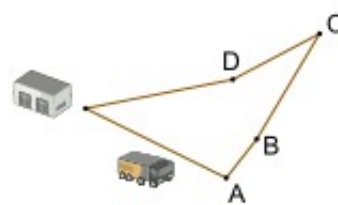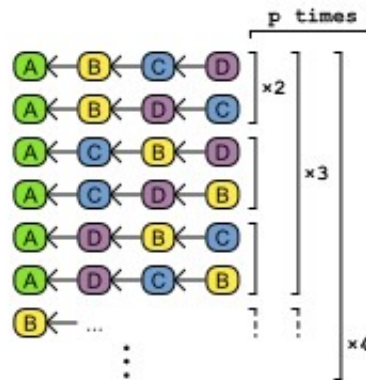Given a Solution model, how many different combinations can it represent?
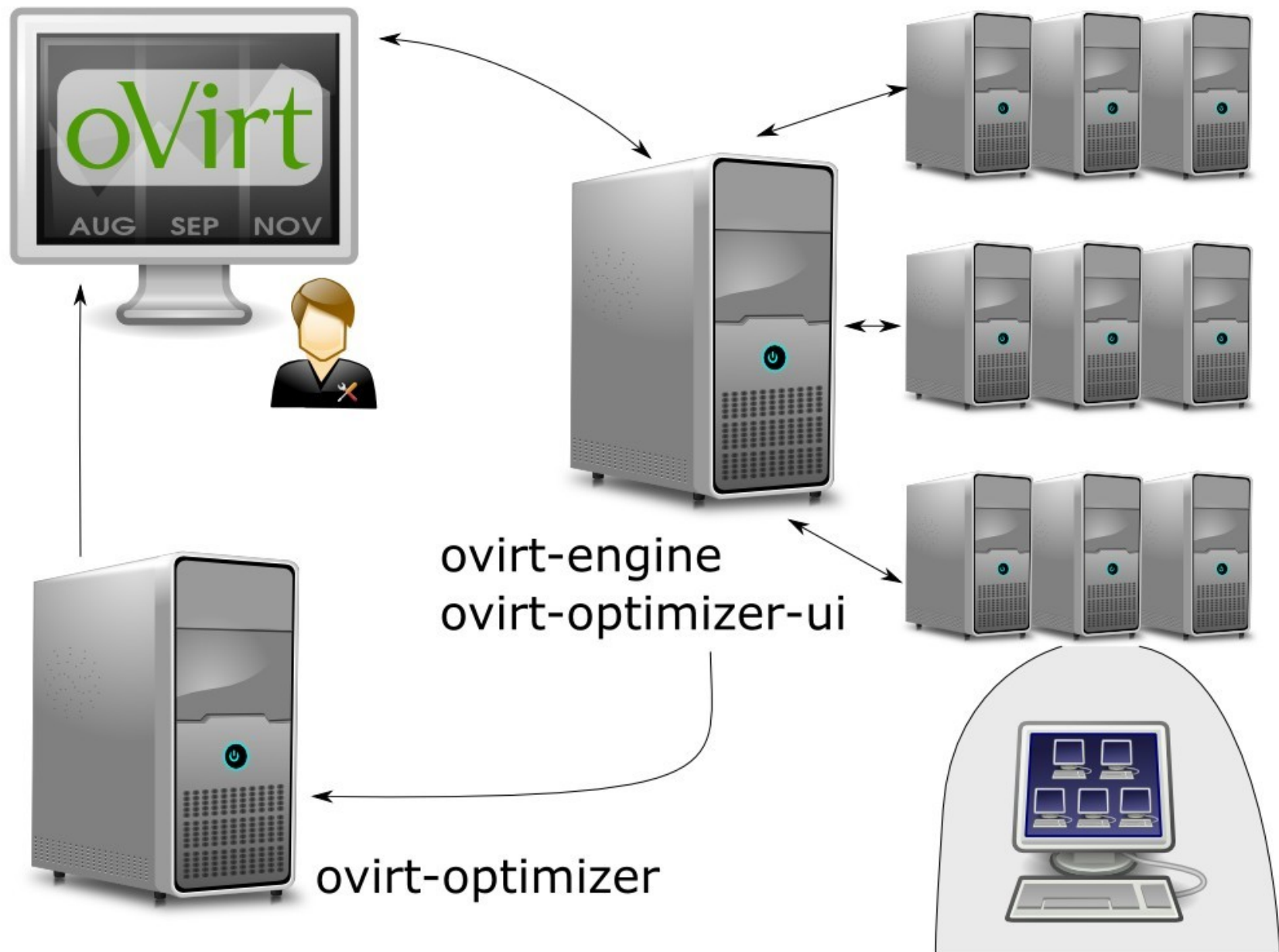
- optaplanner.org

- Optimization engine

- Many search algorithms

- Uses Drools Rule Language (DRL) for scoring

  – drools.org

# Probabilistic approach

- **Random search**
  - Randomly generate a candidate solution
  - Evaluate and assign a score
  - Accept if better than the current
  - Rinse and repeat
- **Smarter than random**
  - Simulated annealing – closer and closer neighbors
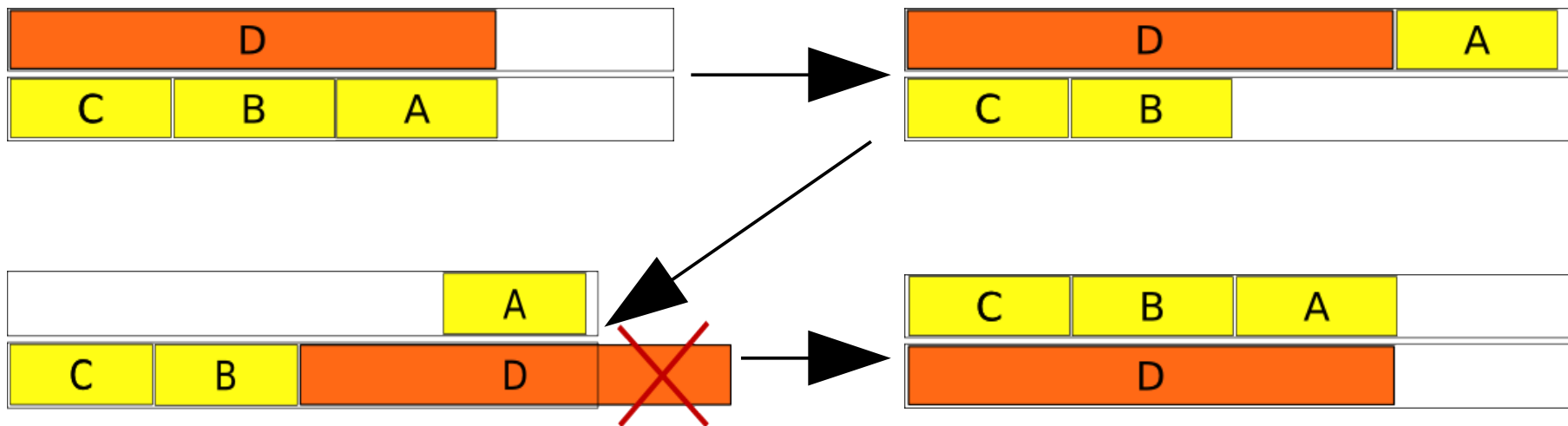  - Tabu search – do not repeat mistakes

# OptaPlanner and oVirt



ovirt-engine
ovirt-optimizer-ui

ovirt-optimizer

- oVirt's Java-based policy units converted to DRL-based rules in order to honor admin-set filters and weights

  - not all policy units yet available through API

    - hosted engine score filters

    - CPU load-based balancing

- cluster info periodically acquired by the optimizer over oVirt's REST API, converted, and fed to the OptaPlanner's fact database

- performance is improved by caching all rule matches

- All previous facts and rules are then used together by the OptaPlanner solver engine to compute the result.

- The optimizer service keeps running and improving the solution.

- When something in the cluster changes, the facts update and the solver resumes using the current best solution as a base point.

# Optimization steps

- Number of steps limited

- Slower to converge than simple "get me the optimum"

- Hard constraint check for each intermediate state

- Soft constraint check for the final situation only

# Web admin integration

SCALE 13x, Feb 2015

# Looking Ahead

- Tighter integration with BRMS

- Full automation of the optimization

  – Using the optimizer instead of the internal scheduler in oVirt engine

- Support for more Policy Units

  – Custom DRL rules

  – Units coming from external scheduler

- Long term cooperation potential

  – OpenStack Gantt

  – Kubernetes

  – Mesos

# Questions?

http://www.ovirt.org
jbrooks@redhat.com
@jasonbrooks
jbrooks on OFTC & Freenode IRC