

TA24

DRS Deep Dive and Technology Preview of Distributed Power Management

Minwen Ji

Staff Engineer

VMware

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

VMWORLD 2007

This session may contain product features that are currently under development.

This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product.

Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind.

Technical feasibility and market demand will affect final delivery.

Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Overview

○ What DRS does

- DRS stands for Distributed Resource Scheduler
- Maps cluster-wide SLA to host-level SLA
- Recommends initial placement of VMs
- Recommends dynamic load balancing with VMotion

○ What DPM does

- DPM stands for Distributed Power Management
- Recommends powering off hosts to save energy
- Recommends powering on hosts when workload increases
- Recommends powering on hosts for certain events

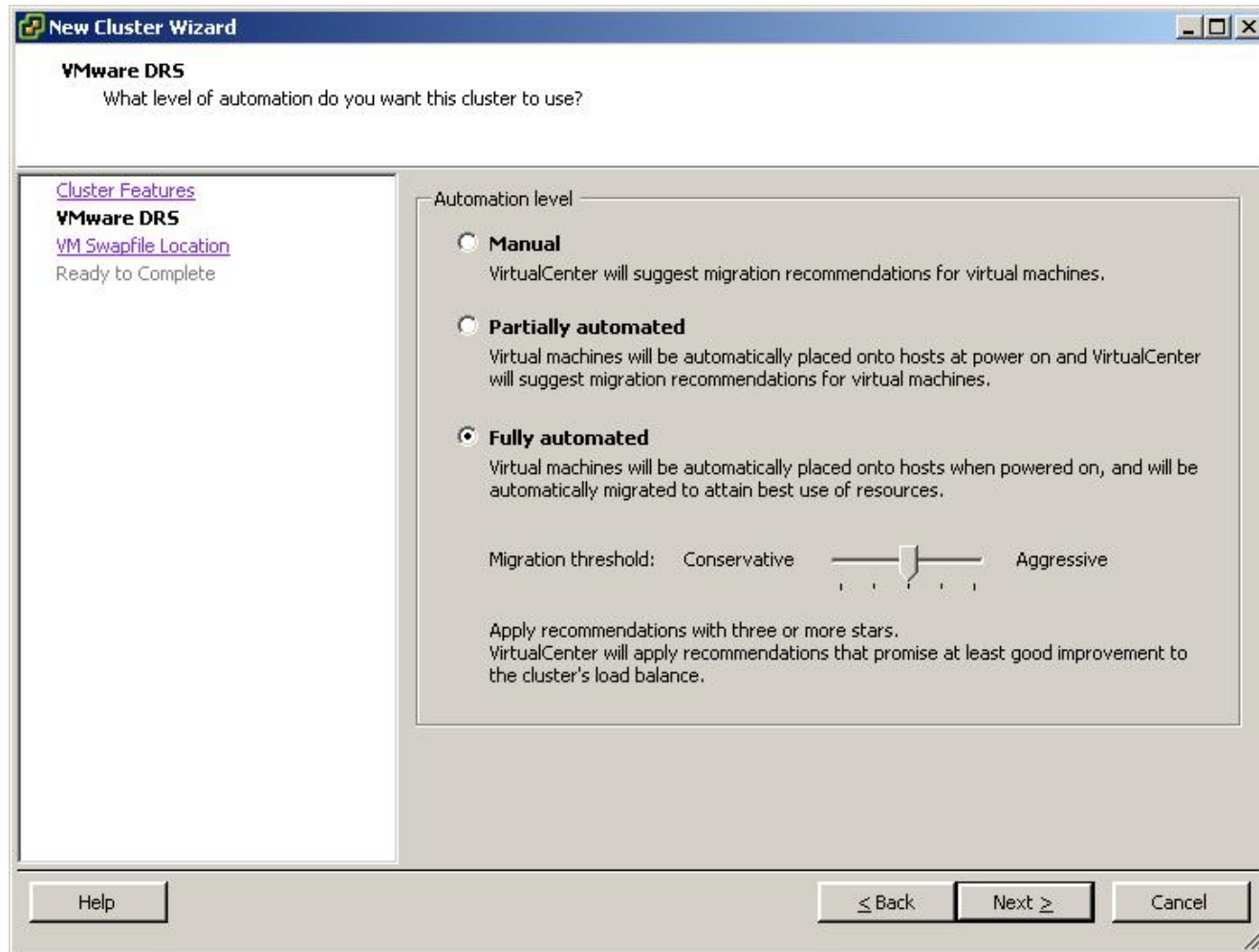
How to Influence DRS/DPM Decisions

- **Many ways to influence the decisions of DRS/DPM:**
- **Service Level Agreement**
 - > e.g., reservations, shares and limits
- **Migration threshold**
- **Automation level**
- **HA (High Availability) requirements**
- **Affinity / anti-affinity rules**
- **Some advanced options**

Outline

- **Create and configure a DRS cluster**
- Compound and dependent recommendations
- Enable DPM – power recommendations
- Power on VMs – host recommendations
- Maintain failover level for High Availability (HA)
- Keep VMs together or separate them – affinity rules

Configure DRS Clusters



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Migration Threshold

- **Reflects tolerance of cluster load imbalance**
 - > Aggressive: less tolerance and more migrations
 - > Conservative: more tolerance and fewer migrations
- **Internal load imbalance metric**
 - > Standard deviation of host load/capacity
- **Internal load imbalance threshold**
 - > Proportional to external threshold (stars) / sqrt (number of hosts)
 - Consistent across clusters with different numbers of hosts
 - Consistent across clusters with different numbers of VMs
 - Effective for a wide range of loads

Load Balancing Algorithm

While (load imbalance metric > threshold) {

Select the best migration in the current state;

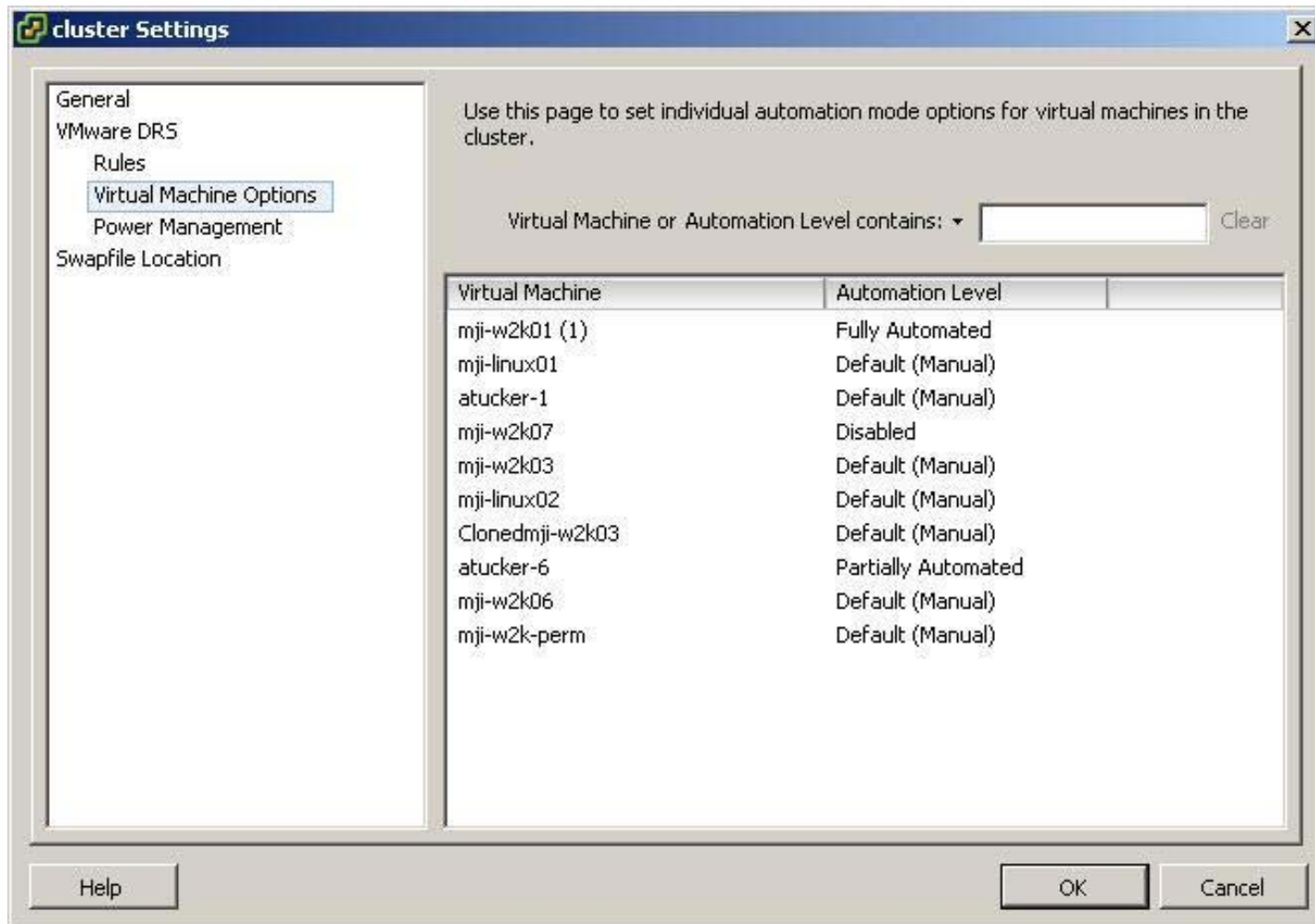
If no good migration is found, stop;

Else recommend the best migration;

Update cluster to the state after the recommended migration is executed;

}

VM Automation Level



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

VM Automation Levels

○ Manual VMs

- > User needs migrate privilege on manual VMs in order to migrate them.

○ Automatic VMs

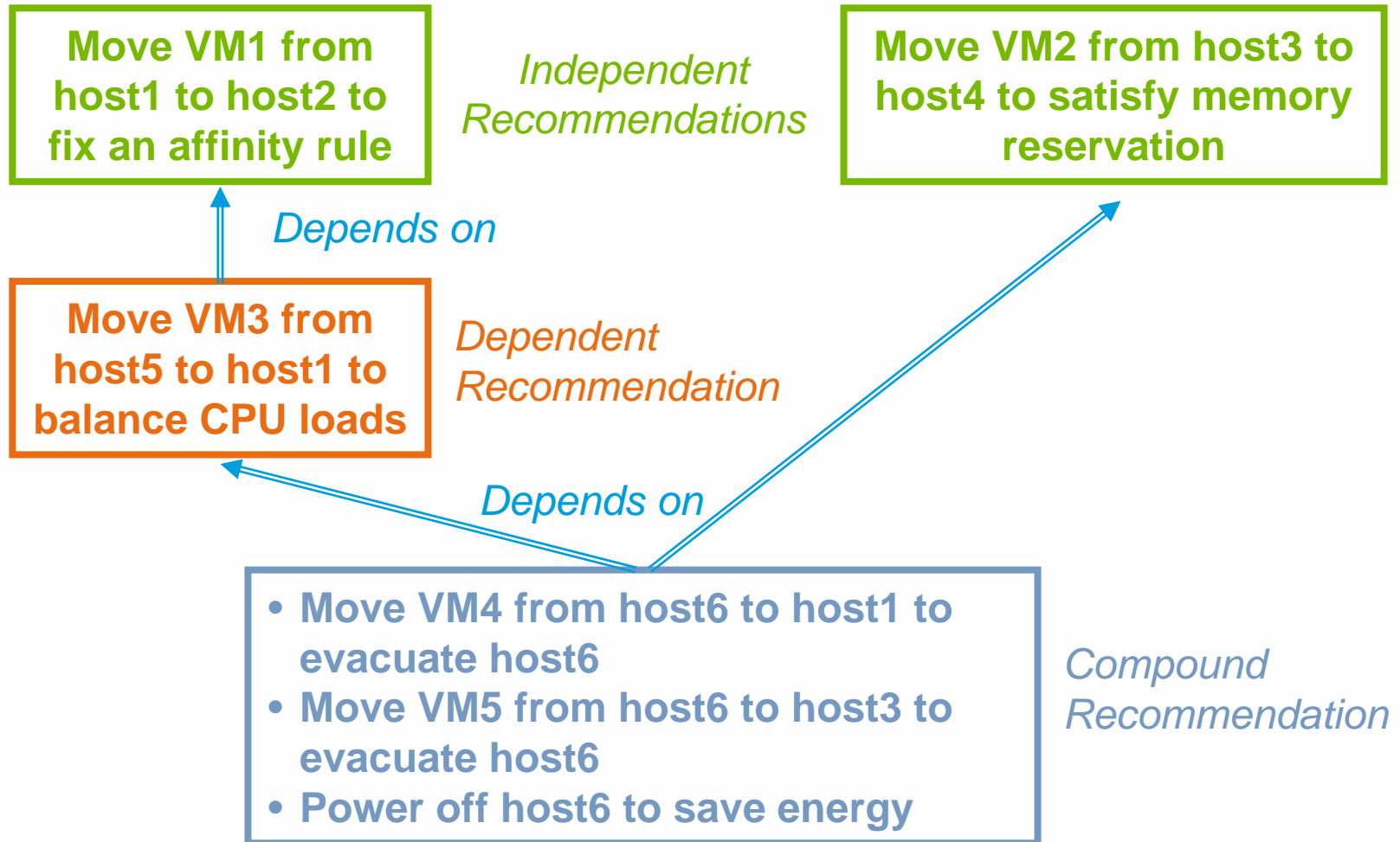
- > DRS can migrate automatic VMs without user's privilege or approval.

○ In selecting migration candidates, DRS prefers automatic VMs to manual ones

Outline

- Create and configure a DRS cluster
- **Compound and dependent recommendations**
- Enable DPM – power recommendations
- Power on VMs – host recommendations
- Maintain failover level for High Availability (HA)
- Keep VMs together or separate them – affinity rules

Workflow of Recommendations



Attributes of Recommendations

○ Automatic vs. manual recommendations

- > A recommendation is manual if and only if any VM or host involved in it has a manual automation level

○ Compound recommendations

- > Have multiple actions, all of which must be executed atomically

○ Independent vs. dependent recommendations

- > Independent ones can start execution as soon as approved
- > Dependent ones cannot start execution until the ones they depend on complete execution

Dependent Recommendations

- **Recommendation A depends on B if and only if**
 - > A is moving a VM into a host that B is moving a VM out of, or
 - > A is moving a VM into a host that B is powering on.
- **If A is executed before B completes, A may fail due to**
 - > Admission failure on host
 - > Invalid host state
 - > Other reasons ...

Compound Recommendations





- **A compound recommendation has multiple actions A, B, C, ...**
- **All actions share the same goal (or reason) and the same star rating**
- **All actions must be executed in order to achieve the goal**
- **Execution of subset may leave the cluster in a worse state than before the execution**

Compound Recommendations

cluster

Getting Started Summary Virtual Machines Hosts **DRS Recommendations** Resource Allocation Performance Tasks & Events Alarms

DRS Recommendations:

Priority	Recommendation	Reason	Apply
★★★★★	Migrate mji-w2k07 from drm104.eng.vmware.com to dr...	Balance CPU reservations	 <input checked="" type="checkbox"/>
★★★★★	Migrate mji-w2k01 (1) from drm103.eng.vmware.com t...	Balance CPU reservations	
★★★★★	Migrate atucker-1 from drm104.eng.vmware.com to dr...	Balance CPU reservations	 <input checked="" type="checkbox"/>
★★★★★	Migrate mji-w2k06 from drm103.eng.vmware.com to dr...	Balance CPU reservations	

Override suggested DRS recommendations

Generate Recommendations Apply Recommendations

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Dependent Recommendations (1)

cluster

Getting Started | Summary | Virtual Machines | Hosts | **DRS Recommendations** | Resource Allocation | Performance

DRS Recommendations:

Priority	Recommendation	Reason	Apply
*****	Migrate mji-w2k06 from...	Satisfy anti-affinity rule	<input type="checkbox"/>
****	Migrate mji-w2k01 (1) fro...	Balance average CPU loads	<input type="checkbox"/>

Override suggested DRS recommendations

Generate Recommendations | Apply Recommendations

Recent Tasks

Name	Target	Status	Initiated by	Time	Start Time	Complete Time
Refresh Recommendation...	cluster	Completed	vcadmin	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM
Power On Virtual Mach...	mji-w2k01 (1)	Completed	vcadmin	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:40 PM
Get Recommendation	cluster	Completed	vcadmin	6/21/2007 1:07:26 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Dependent Recommendations (1)

The screenshot shows the VMware Virtual Infrastructure Client interface. The left pane displays the 'Hosts & Clusters' tree with a 'cluster' selected. The main pane shows the 'DRS Recommendations' tab for the selected cluster. A dialog box titled 'Prerequisite Recommendations Selected' is open, displaying an information icon and the message: 'The selected recommendation is dependent on one or more other recommendations. These prerequisite recommendations have also been selected.' Below the message is a checkbox labeled 'Do not show this message again.' and an 'OK' button. The background table shows the following DRS Recommendations:

Priority	Recommendation	Reason	Apply
*****	Migrate mji-w2k06 from...	Satisfy anti-affinity rule	<input type="checkbox"/>
****	Migrate mji-w2k01 (1) fro...	Balance average CPU loads	<input type="checkbox"/>

At the bottom of the window, the 'Recent Tasks' pane shows a table of completed tasks:

Name	Target	Status	Initiated by	Time	Start Time	Complete Time
Refresh Recommendation...	cluster	Completed	vcadmin	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM
Power On Virtual Mach...	mji-w2k01 (1)	Completed	vcadmin	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:40 PM
Get Recommendation	cluster	Completed	vcadmin	6/21/2007 1:07:26 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Dependent Recommendations (1)

cluster

Getting Started | Summary | Virtual Machines | Hosts | **DRS Recommendations** | Resource Allocation | Performance

DRS Recommendations:

Priority	Recommendation	Reason	Apply
*****	Migrate mji-w2k06 from...	Satisfy anti-affinity rule	<input checked="" type="checkbox"/>
****	Migrate mji-w2k01 (1) fro...	Balance average CPU loads	<input checked="" type="checkbox"/>

Override suggested DRS recommendations

Generate Recommendations | Apply Recommendations

Recent Tasks

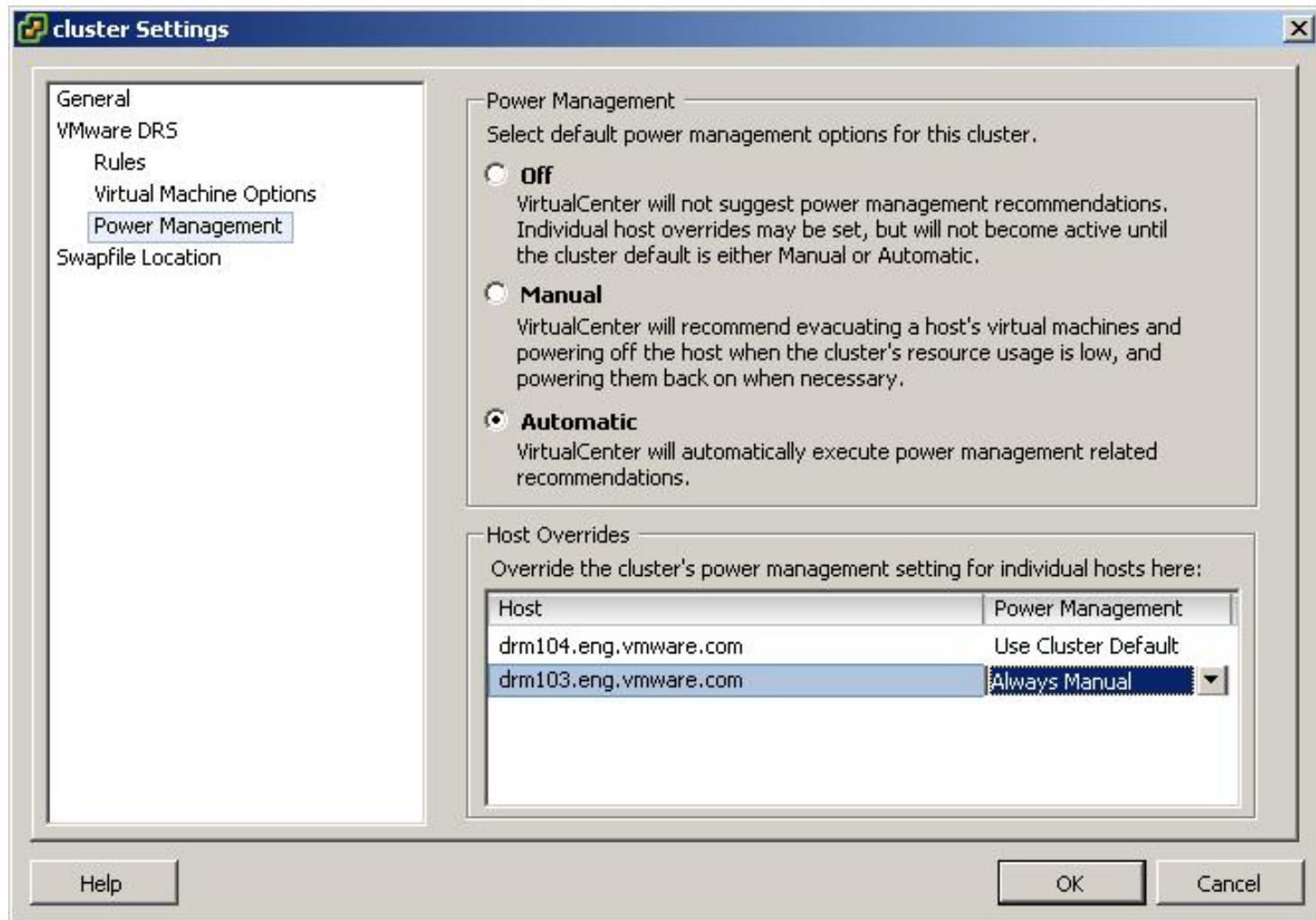
Name	Target	Status	Initiated by	Time	Start Time	Complete Time
Refresh Recommendation...	cluster	Completed	vcadmin	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM	6/21/2007 1:07:41 PM
Power On Virtual Mach...	mji-w2k01 (1)	Completed	vcadmin	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:40 PM
Get Recommendation	cluster	Completed	vcadmin	6/21/2007 1:07:26 PM	6/21/2007 1:07:27 PM	6/21/2007 1:07:27 PM

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Outline

- Create and configure a DRS cluster
- Compound and dependent recommendations
- **Enable DPM – power recommendations**
- Power on VMs – host recommendations
- Maintain failover level for High Availability (HA)
- Keep VMs together or separate them – affinity rules

Enable DPM (Distributed Power Management)



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

DRS/DPM Interactions

○ **Conflicting goals ?**

- DRS balances load to satisfy SLAs
- DPM optimizes for power consumption

○ **High level interaction**

- DPM determines how many hosts to be powered on
- DRS balances load among powered-on hosts

○ **Detailed interactions**

- DRS rebalances with DPM-recommended power actions in what-if simulations
- DPM evaluates potentials of power actions based on DRS rebalancing results
- Final power actions and VMotions are recommended to user

Power Management Algorithm Overview

- **Goal: try to maintain per-host load/capacity ratio within a target range (centered at 63%), while satisfying all DRS constraints.**
- **After DRS load balancing, divide hosts into highly and lightly loaded groups.**
 - Consider power-on for highly loaded group and power-off for lightly loaded group.
 - Run load balancing simulation for each power-on candidate
 - No VM can move to the new host => do not power it on
 - Run evacuation simulation for each power-off candidate
 - Host cannot be evacuated => do not power it off

Power-on and Power-off Decisions

- **Compute a weighted average load/capacity ratio for each group, i.e., R_h and R_l .**
 - > Consider longer load history for power-off decisions (20 minutes) than for power-on decisions (5 minutes) – conservative in powering off and responsive in powering on.
- **Power on for highly loaded group**
 - > If $R_h > \text{highThreshold}$, power on additional hosts, until $R_h \leq \text{highThreshold}$.
- **Power off for lightly loaded group**
 - > If $R_l < \text{lowThreshold}$, power off existing hosts, until $R_l \geq \text{lowThreshold}$.
 - > Take into account load variation, power operation overhead, and power-performance trade-off in powering off decisions.
- **All parameters are tunable online through advanced options.**

Power Off To Save Energy

The screenshot shows the VMware vSphere interface for a cluster. The 'DRS Recommendations' tab is active, displaying a table of recommendations. The first recommendation is to migrate VM 'atucker-1' from host 'drm104.eng.vmware.com' to another host, with a reason of 'Power off host for power savings'. The second recommendation is to power off host 'drm104.eng.vmware.com' for the same reason. Both recommendations have a priority of five stars. At the bottom, there is a checkbox for 'Override suggested DRS recommendations' and two buttons: 'Generate Recommendations' and 'Apply Recommendations'.

Priority	Recommendation	Reason	Apply
★★★★★	Migrate atucker-1 from drm104.eng.vmware.com to dr...	Power off host for power savings	<input checked="" type="checkbox"/>
★★★★★	Power off host drm104.eng.vmware.com	Power off host for power savings	

Override suggested DRS recommendations

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Power On To Handle Increased Load

The screenshot shows the vSphere DRS Recommendations interface for a cluster. The breadcrumb navigation includes: Getting Started, Summary, Virtual Machines, Hosts, DRS Recommendations, Resource Allocation, Performance, Tasks & Events, Alarms, and Page navigation. The main section is titled "DRS Recommendations:" and contains a table with the following data:

Priority	Recommendation	Reason	Apply
*****	 Power on host drm104.eng.vmware.com	Increase capacity	<input checked="" type="checkbox"/>

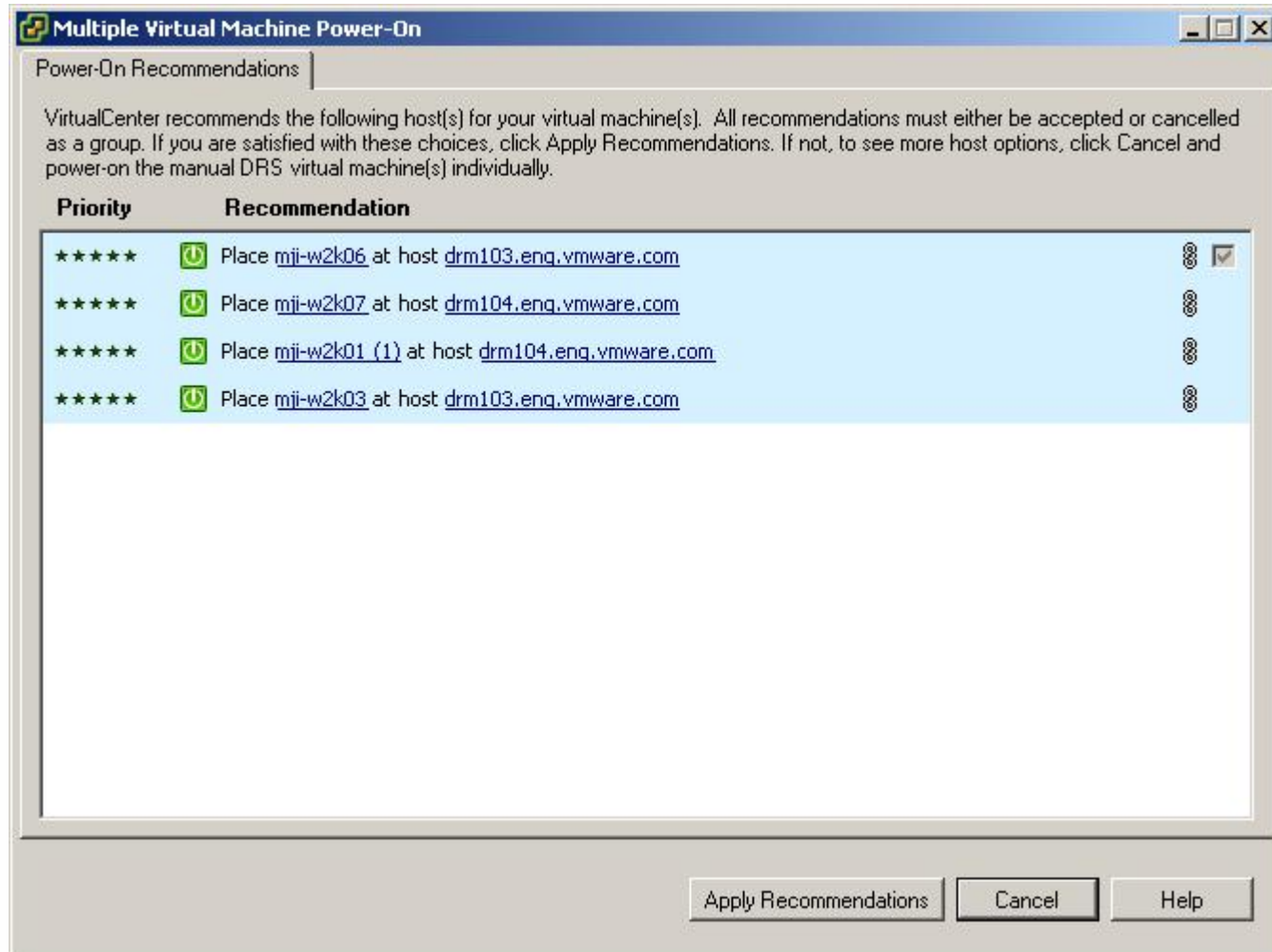
At the bottom of the interface, there is a checkbox labeled "Override suggested DRS recommendations" which is currently unchecked. To the right of this checkbox are two buttons: "Generate Recommendations" and "Apply Recommendations".

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Outline

- Create and configure a DRS cluster
- Compound and dependent recommendations
- Enable DPM – power recommendations
- **Power on VMs – host recommendations**
- Maintain failover level for High Availability (HA)
- Keep VMs together or separate them – affinity rules

Power On Multiple VMs



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Power-on With Prerequisite



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Host Recommendations for VM Power-On

- **May power on a single VM or multiple VMs simultaneously**
 - > In case of multiple power-on, a compound recommendation per DRS cluster
- **May have multiple choices**
- **Each recommendation has at least one VM power-on action, and optionally prerequisite actions**
 - > Migration
 - > Host power-on (if DPM is enabled)

Advantages of Multiple Power-on

○ Better performance

- DRS is invoked only once for a multiple power-on, but is invoked N times for N single power-ons.

○ Better placement decision

- If powered on in the wrong order, later VMs may fail in single power-on.
- In multiple power-on, DRS places VMs in descending order of constraints, to avoid failures in later stage.

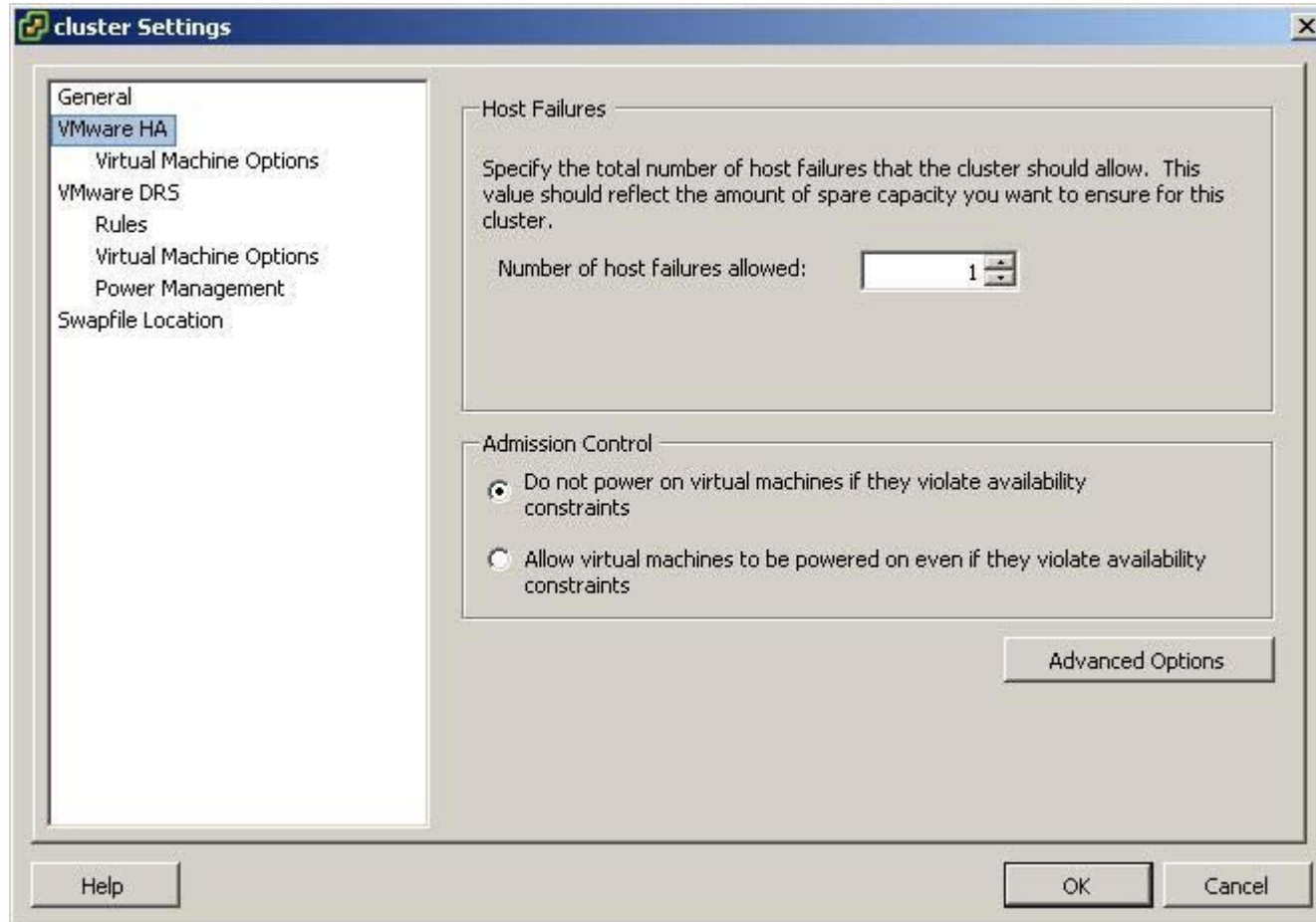
○ Ideal for Virtual Desktop environments

- All VMs need to be powered on at 9am when people come to work

Outline

- Create and configure a DRS cluster
- Compound and dependent recommendations
- Enable DPM – power recommendations
- Power on VMs – host recommendations
- **Maintain failover level for High Availability (HA)**
- Keep VMs together or separate them – affinity rules

Enable HA (High Availability)



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

HA/DPM Interactions

○ HA does **VM admission control** to satisfy the constraints:

- > **SlotSize = Largest VM reservation**
- > **NumSlots \geq Capacity of N largest hosts / SlotSize**
- > **NumVms \leq ClusterCapacity / SlotSize – NumSlots**

DPM does **host exit control** to satisfy the same constraints:

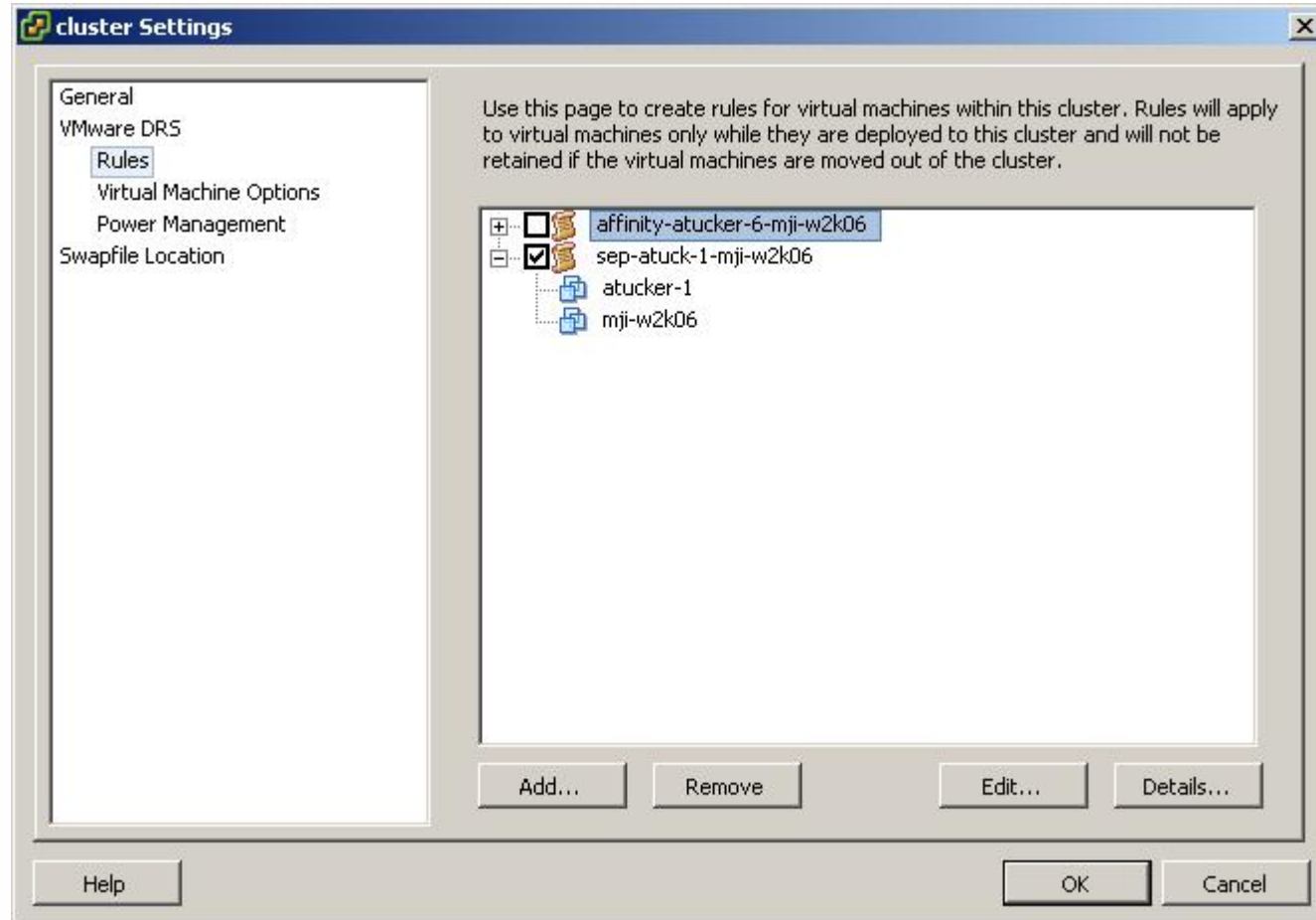
- > **SlotSize = Largest VM reservation**
- > **NumSlots \geq Capacity of N largest hosts / SlotSize**
- > **ClusterCapacity \geq (NumVms + NumSlots) * SlotSize**



Outline

- Create and configure a DRS cluster
- Compound and dependent recommendations
- Enable DPM – power recommendations
- Power on VMs – host recommendations
- Maintain failover level for High Availability (HA)
- **Keep VMs together or separate them – affinity rules**

Rules Overview



This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Rules Details

Virtual Machine Rule

Give the new rule a name and choose its type from the menu below. Then, select the virtual machines to which this rule will apply.

Name
affinity-atucker-6-mji-w2k06

Type
Keep Virtual Machines Together

Virtual Machines
atucker-6
mji-w2k06

Add... Remove

OK Cancel

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Use Cases and Constraints

○ Affinity rules

- > VMs communicating a lot with each other
- > VMs sharing a portgroup on a virtual switch (LabManager)
- > Constraint: at most one affinity rule per VM

○ Anti-affinity rules

- > Redundant VMs for fault tolerance
- > Constraint: at most two VMs per anti-affinity rule (may be lifted in the future)

Enforcing Affinity Rules

- **Each DRS invocation runs in a number of steps:**
 - First, try to enforce SLA, affinity rules and HA requirement.
 - Second, optimize for load balancing and power consumption.
 - Must not violate a new rule in enforcing another.
 - Must not violate a new rule in optimization.
 - Better error reporting is needed when rules cannot be enforced.

Summary

○ **DRS control knobs**

- Migration threshold, automation levels, affinity rules, etc.

○ **DRS recommendations**

- Compound, dependent, initial placement, power, etc.

○ **DRS/DPM interactions**

- DRS balances load among hosts that DPM powers on

○ **HA/DPM interactions**

- HA specifies spare capacity to maintained by DPM

Backup slides

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

Enforce Affinity Rules

Part of DRS placement algorithm:

For each affinity rule {

 If (rule violated) {

 For each host in cluster {

 If (can have all VMs in the rule) {

 Move all VMs in the rule to the host;

 Done with this rule;

 }

 }

 }

}

Enforce Anti-Affinity Rules

Part of DRS placement algorithm:

```
For each VM {  
  If (violates an anti-affinity rule) {  
    If (can move to another host) {  
      Move it to another host;  
    }  
  }  
}
```

Questions?

TA24

**DRS Deep Dive and Technology
Preview of Distributed Power
Management**

Minwen Ji

VMware

This session may contain product features that are currently under development. This session/overview of the new technology represents no commitment from VMware to deliver these features in any generally available product. Features are subject to change and must not be included in contracts, purchase orders, or sales agreements of any kind. Technical feasibility and market demand will affect final delivery. Pricing and packaging for any new technologies or features discussed or presented have not been determined.

VMWORLD 2007



VMWORLD 2007

EMBRACING YOUR VIRTUAL WORLD

BREAKOUT SESSION