Vantage
App Note

# Using Vantage Load Balancing

**This App Note applies to Vantage 7.0, 7.0 UP1, 7.0 UP2, & later**

**Note:** This guide is written for video professionals who are familiar with using Vantage. To implement applications in Vantage, you should know how to create workflows and submit jobs. If you aren't familiar with Vantage, we suggest that you review the *Vantage User's Guide* and *Vantage Domain Management Guide* as needed.

telestream

# Synopsis

Load balancing offers significant benefits to enhance the operational speed of your Vantage installation. Using the built-in CPU-aware task scheduling or the optional Advanced Task Scheduling, you can ensure that tasks are adjusted in priority in a way that maximizes processing efficiency on a single machine or an array. Your workflows will execute more quickly, reducing the time it takes to transcode media.

The following three methods of load balancing are explained in this application note:

- Session-based Load Balancing
- Cost-based Load Balancing (legacy)
- Task Based Load Balancing

The Advanced Task Scheduling license enables all three selections, so you can choose among them (they are mutually exclusive). Without the Advanced Task Scheduling license, only the Session Based is available as the default.

## Key Concepts

- Load balancing spreads tasks across multiple servers using rules you specify.
- Session-based Load Balancing runs by default on array machines to set session limits on services. Session limits prevent service overload by limiting the maximum number of actions a service can run.
- Cost-based Load Balancing, a legacy method, lets you assign costs to specific actions and target usage levels (maximum total costs) for specific services. This limits the number of actions that can run on a service, and enables actions to be shifted to the same service on other machine to balance workload.
- Task-based Load Balancing provides a single set of metrics shared across all services for scheduling tasks and balancing workload on a single machine or all the machines in a domain.
- Load balancing can operate on a machine only if...
  – The Vantage service being balanced is installed on that machine.
  – An Advanced Task Scheduling license is available for that machine.
- Lightspeed Bundle licenses handle load balancing differently—one entire policy operates together and only affects transcoding licenses. Any Flip, Multiscreen Flip, or IPTV VOD Flip that uses Lightspeed will only work on a Lightspeed Server.
- Task scheduling looks at available session limits and CPU usage to determine which machine is least busy and assigns services accordingly.
- Task Routing using Run On rules lets you configure an action to run on a particular server.
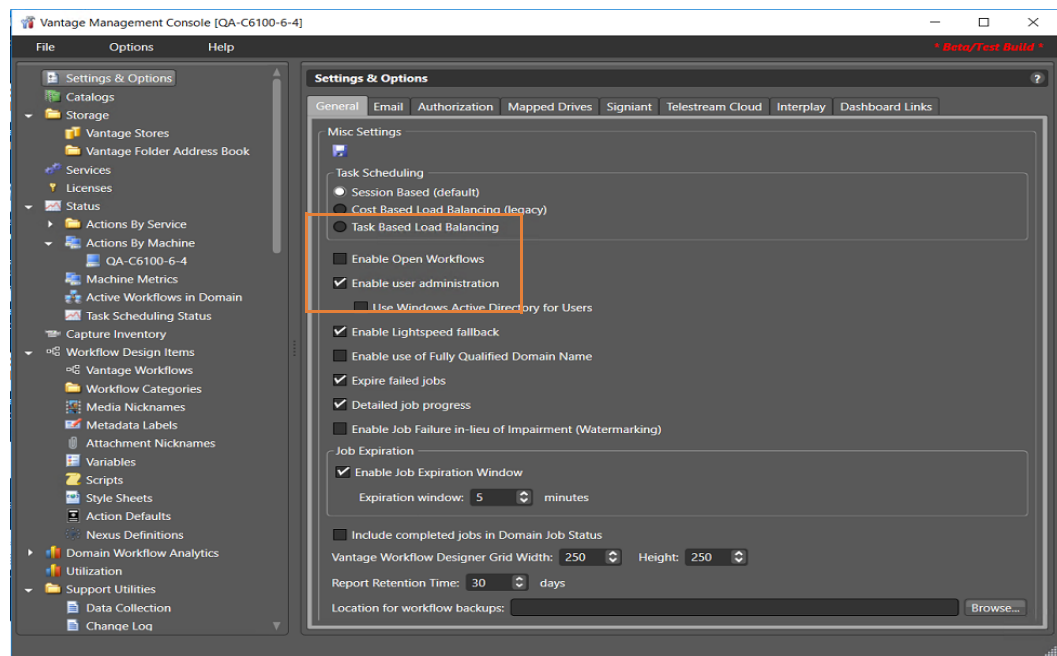- Priority variables let you assign higher or lower priorities to specific actions.

# Overview

As Vantage has evolved, three methods of load balancing have been developed:

- Session-based Load Balancing
- Cost-based Load Balancing (legacy)
- Task Based Load Balancing

You select among the three methods, which are mutually exclusive, in the Vantage Management Console *Settings and Options* panel, General tab. Click the method you prefer to select it and deselect the other two.

The three load balancing methods are described below.



## Session-Based Load Balancing

Session-based load balancing, which is included with Vantage and does not require a license, allows you to distribute the processing of actions across multiple servers that support the same service, up to the configured session limit for each server.
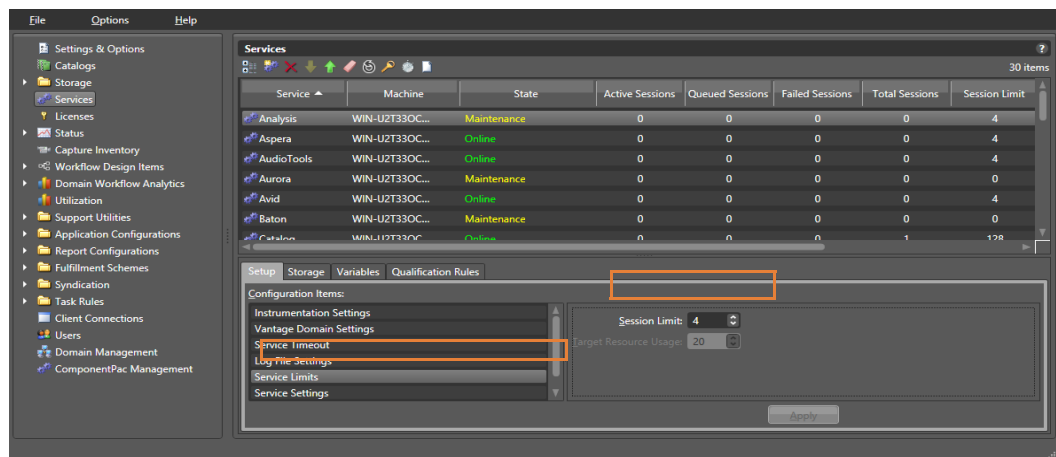
**Note:** A Lightspeed K20 Server can run a maximum of six concurrent Lightspeed jobs regardless of the capacity setting. (K80 Lightspeed servers have no such concurrency limit.) The default session limit setting for each transcode service is four. This helps ensure the server retains enough GPU memory to efficiently process all jobs. If Lightspeed jobs fail because they exceed memory limits, try a lower setting.

This load balancing is done by taking the CPU capacity and current utilization of each machine into account, as well as ensuring that any workflow Run On rules are observed.

telestream

(Run On rules allow tasks to be routed to specific machines. They are configured in the Management Console > Services > Variables and are set in the Workflow Designer by right-clicking on actions and selecting Run On Rules.)

When the configured session limit is reached on all servers, actions are queued until the workload drops below the configured session limit. (An exception to this is that some actions support the *pause for priority* feature in which an action pauses, yielding to a higher priority action of the same type.)

You set the Session Limit in the Vantage Management Console by selecting Services and selecting the service you want to change. The Session Limit is listed in the Setup panel below the services list. To change the limit for a service, you must first right-click and select Enter Maintenance Mode, which enables the Session Limit selection. After changing the Session Limit, right-click the service and select Exit Maintenance Mode.



## Cost Based Load Balancing (Legacy)

Cost based load balancing is a legacy method that has been improved upon by Task Based Load Balancing. Cost based load balancing requires a Vantage Array license and the Advanced Task Scheduling license and allows you to do the following:

- Assign a default cost for each type of action used in workflows
- Define a target resource usage level for each service on a server
- Override the default resource cost in workflows as needed

---

**Note:** Session limits are ignored when cost based load balancing is enabled, and the target resource usage level is ignored when session-based load balancing is enabled.

---

To determine its available resources, each service monitors the total cost of all actions it is processing and compares that cost to the configured target resource usage level for that service. A service can accept an action for processing if the available resource usage is at least half the cost of the action.

Unless priority is used, actions are assigned to services in the order in which they are processed, which prevents low cost actions from consuming all service resources and

starving higher-cost actions. Each action is forwarded to the service with the most available resources. If the action is forwarded to a service that does not have resources to immediately process the action, the action is queued for later processing.

To implement cost based load balancing, enable the feature (as described in *Implementing the Cost-Based Option*) and monitor the following:

- The server performance using server monitoring tools.

- The Services panel in the Vantage Management Console, which displays a list of services, including the current and target resource usage levels. For more information, see *Service Resource Utilization*.

If the server is under- or over-loaded, you might adjust the default target resource usage level for each service as described in *Service Resource Utilization*.

To adjust the default resource cost of an action, search for *Managing Actions* in the Domain Management Guide. For information on overriding an action's default resource cost in a workflow, see the *Vantage User Guide*.

## Task Based Load Balancing

Vantage Task Based load balancing is Vantage's most advanced load balancing capability and augments the existing Vantage Session Based load balancing and the Vantage Cost Based load balancing. Task Based load balancing ensures that all Vantage servers remain fully utilized, maximizing server efficiency and increasing workflow throughput. Task Based load balancing also helps avoid the situation where some Vantage servers are over utilized and others are under-utilized.

The older session-based and cost based methods discussed previously are limited to a single service on a single machine. These methods can prevent a service from becoming overloaded, but they cannot effectively balance between different types of actions on a single machine, and they do not offer a global way of managing sessions across multiple machines.

Task based load balancing is a separately licensed option that provides a single set of metrics shared across all services for scheduling tasks and balancing workload on a single machine or all the machines in a domain.

These are the key features of task based load balancing:

- Each metric can be edited and has machine scope (across all services on the machine) or global scope (across all services on all machines).

- Each metric has two inputs: Cost (set for each action) and Capacity (set per machine, or globally). You can edit both on a per-machine basis (for per-machine metrics) or for the entire Vantage Domain (for global metrics).

- A default metric called System Resources is pre-configured to provide a per-machine metric. This metric cannot be deleted and must remain per-machine. However, you may set the capacity for this metric on each machine, and you may remove the metric from an individual action.

- You can create and name your own new metrics.

telestream

- By default, each metric has machine scope, but you can edit this after creation.
- The default capacity is 100 "units" for each machine, allowing each cost to indicate the resource percentage an action will use. The capacity can be changed.
- You can choose which metrics are relevant to an action, and you can set the cost for each action.
- When executing a task:
  - Vantage evaluates all metrics that you attach to the action. If a metric is not attached to an action, it is not considered.
  - Vantage tracks the aggregate cost of all executing actions on a machine, across all machine-based metrics, and will track the aggregate cost of executing actions across the entire system for global metrics.
  - Vantage only executes an action when the aggregate cost of executing current actions, plus the new task, does not exceed the capacity of any of the metrics attached to the task.

Vantage will never execute a task when the aggregate cost of executing actions, including the new task, would exceed the capacity of any of the metrics attached to the task.

- When tasks are queued, they may be pre-empted only by higher priority actions:
  - Certain high priority tasks will pause an appropriate number of low priority tasks (executing on the same service on the same machine) to free up sufficient resources across all metrics, to allow that high priority task to execute. If doing so would exceed the maximum number of paused jobs specified by the user, then the high priority task will wait until a different set of tasks are present which may be paused. Pauses will not occur across multiple machines; they occur only for actions hosted by the same service.

The new task based load balancing capability allows you to define custom task scheduling rules which allow jobs to be controlled/allocated based on criteria on a given machine or across an entire domain.

Vantage task based load balancing is enabled via the settings panel in the Vantage Management Console. As with other system level settings, when task based load balancing is enabled (or changed), services must be restarted to implement the new settings.

Task based load balancing includes capabilities that replace and extend the Cost based load balancing mode. The new capabilities replace the existing action cost mechanisms, and are used in-conjunction with the existing service capacity to provide more granularity on when and where an action may run.

# Choosing a Load Balancing Method

The three methods of load balancing provide different benefits, so the choice of method depends on your particular workflow processing pattern. The following discussions explain the different benefits and setup of each load balancing method.

## Benefits of Session-based Load Balancing

Session-based load balancing offers good, basic load balancing as a starting point for Vantage Array systems, and it comes built-in and operating with the Array license. You can make adjustments to this automatic system by setting the session limits for each service. If your systems handle a modest volume and no very large jobs, this may be all the balancing you need. It's a convenient set it and forget it system. However, if your typical workflow consists of multiple transcode actions that are a combination of very small/very low CPU (resource) utilization and very large/very high resource utilization, using cost-based or task-based load balancing will allow you to design a policy that provides optimal machine loading.

### Adjusting Session-based Load Balancing

The Session-based feature is operational by default with the Array license. You can adjust this automatic system by setting the session limits for each service:

1. Click the service in the Vantage Management Console > Services panel.

2. In the lower panel, click the Setup tab > Service Limits selection > Session Limit number.

3. Increase or decrease the number of concurrent jobs (Session Limit) you want to permit the service to handle.

4. Observe performance and adjust Session Limits as required.

In addition to the performance information in the VMC Services panel, you can consult the VMC Domain Analytics panels, which provide extensive performance information about specific actions and workflows. If you see actions associated with a particular service causing bottlenecks, you can adjust the session limit lower for that service and observe the system to see if performance improves. The goal is to set the highest session limit that provides acceptable overall performance.

## Benefits of Cost-based Load Balancing

Cost-based load balancing offers considerable control and flexibility in balancing a Vantage Array for optimum, efficient operation. Dynamic task scheduling helps ensure very large jobs or many small jobs do not overwhelm the system. You can also assign processing costs independently to services and actions, allowing you to finely tune processing resource usage. There are many benefits to this option, and it is easy to set up, but it requires some time to monitor and adjust.

## Cost-based Examples

In the VMC Services panel, you can use the target resource usage setting to limit the total action cost that a service will simultaneously process. For example, if you configure the Transcode service for a target resource usage of 16 and you configure the default Flip action resource cost to 4, the Transcode service supports a maximum of 4 simultaneous Flip actions.

Some services support multiple actions, and you can set the action resource cost for each action to a different value to define the relative usage levels. For example, the Transport service supports the Copy, Delete, Deploy, and Move actions, which are assigned default resource costs of 2, 1, 1, and 2, respectively. Based on the default values, a Transport service can support twice as many Delete and Deploy actions as Copy and Move actions.

Different codecs require different levels of processing resources. Using Workflow Designer, you can override the default resource cost for specific actions within a workflow. For example, you might configure a Transcode service to target resource usage of 16, specify a resource cost of 1 for a 3GP Flip action, and a resource cost of 8 for an MPEG-2 HD Flip action. During load balancing, Vantage will ensure that:

- No more than two MPEG-2 HD jobs execute simultaneously
- Up to 16 3GP jobs can execute simultaneously
- One MPEG-2 HD and eight 3GP jobs can execute simultaneously in ideal situations
- MPEG-2 jobs are not starved by 3GP jobs which only require one slot

**Note:** Resource costs are integer values that have an ordinal relationship to one another. You can implement any scale you want in your domain. Greater scales enable finer-grained control.

## Implementing the Cost-Based Option

To implement cost-based load balancing, check the feature in VMC > Settings and Options > Details panel > General tab. Then monitor the following:

- Server performance using server monitoring tools.
- Service current and target resource usage levels in the VMC Services panel.
- Action Analytics and Workflow Analytics in the VMC Domain Workflow Analytics panel.

As you observe performance and draw conclusions about which resources need more or less processing power, you can adjust the following settings.

**Note:** A Lightspeed Server can run only four concurrent Lightspeed jobs at a time regardless of the cost/capacity setting. This limitation helps ensure the server retains enough GPU memory to efficiently process all jobs.

## Service Resource Utilization

If the server is under- or over-loaded when processing a particular service, you can adjust the Target Resource Usage to change the allowed usage level for that Vantage service:

1. Right-click the service in the VMC Settings panel, and place the service in maintenance mode.

2. Select the Setup tab.

3. Select Service Limits, and locate the Target Resource Usage setting.

4. Set the usage number higher if the server is underloaded, lower if overloaded when processing this service.

5. Right-click the service and exit from maintenance mode.

## Action Resource Utilization

As you observe action resource utilization, if you determine that a particular action should be set to consume more or less processing power, you can adjust that action in the VMC Action Defaults panel, where you can define the Resource Cost. A higher resource cost means that fewer actions can run at the same time, which reduces demand on server resources. A lower resource cost allows more actions to run simultaneously.

1. In VMC Action Defaults, click the action to highlight it.

2. Set the Resource Cost higher or lower as needed.

3. You may also want to change the Automated Retry Settings. To consume less processing power, set retry to fewer times (or No Retry) after a processing failure.

## Recommended Starting Service Settings

The Vantage default Target Resource Usage (TRU) and Resource Cost (RC) settings may not be optimum for a typical Vantage Array. You may want to start your Cost-based Load Balancing implementation with the following settings and then adjust by trial and error.

| Vantage Service or Action | TRU or RC Setting |
| --- | --- |
| Transcode service | 100 |
| Flip action | 25 |
| Catalog service | 128 |
| Communicate service | 128 |
| Metadata service | 64 |
| Monitor service | 64 |

telestream

# Benefits of Task-based Load Balancing

Task Based load balancing is the most efficient form of Vantage system load balancing. It provides the ability to balance workload across services and servers in a Vantage Array, as opposed to Session-based and Cost-based load balancing, which only balance workload within a single service.
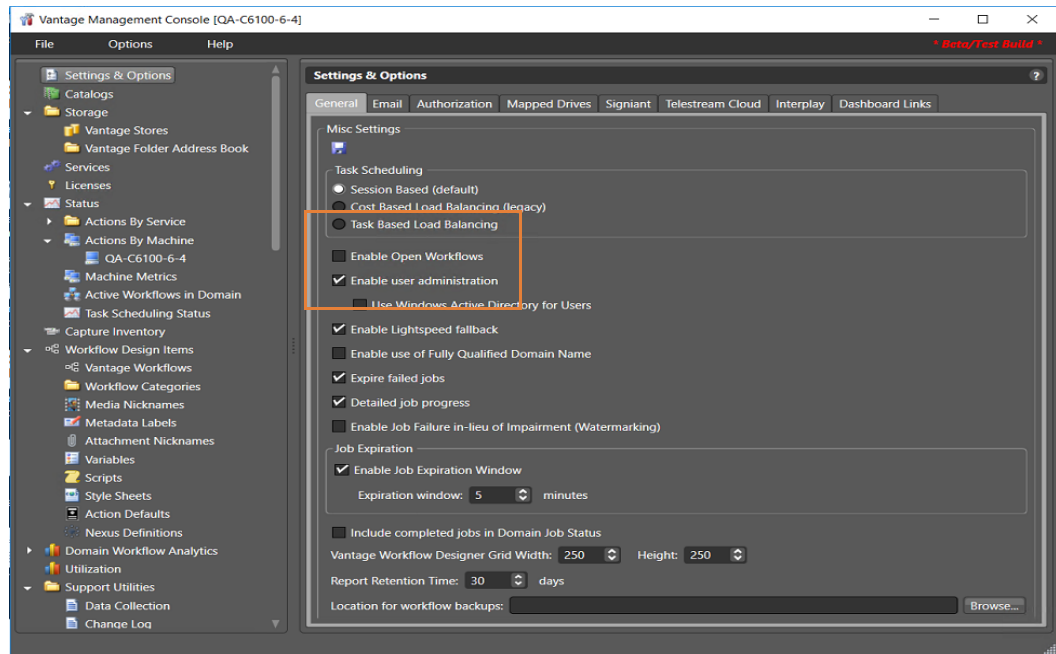
For example, Task Based load balancing allows you to balance the workload between CPU-intensive actions such as a Flip action, a Multiscreen action, an IPTV action and an Analyze action. Each of these actions is managed by a separate Vantage service yet Task-based load balancing allows for tuning of the execution to maximize server utilization and minimize processing time. Task-based load balancing uses the concept of Capacity for the servers (or entire Vantage Array) and Cost for the Actions that need to be processed.

Rules are created to define Capacities and Costs and are assigned at a default level for all servers and services respectively. Additionally, the default Cost assigned to an action can be overridden within the Workflow that utilizes the Action. When jobs run, total Costs are constantly monitored for each server and service to ensure they never exceed specified Capacities. Jobs are shifted as necessary among the resources to maintain the load balance.

Not only can Rules be created to control the balancing of Actions across services, but also Rules can be created that control access to physical characteristics of a Vantage Array. Access to a storage sub-system, for example, can be controlled through Task Based load balancing to ensure that the storage sub-system is not oversaturated with too many processes.   Likewise, if you need to send files via FTP to a remote server, and the remote server only allows a single FTP connection, Rules defined in Task-based load balancing allow this to be accomplished.

# Configuring Task-based Load Balancing

**1.** Install the Advanced Task Scheduling license via the Vantage Management Console > Licenses panel.

**2.** Enable Task Based Load Balancing in the Vantage Management Console > Settings & Options > General panel.
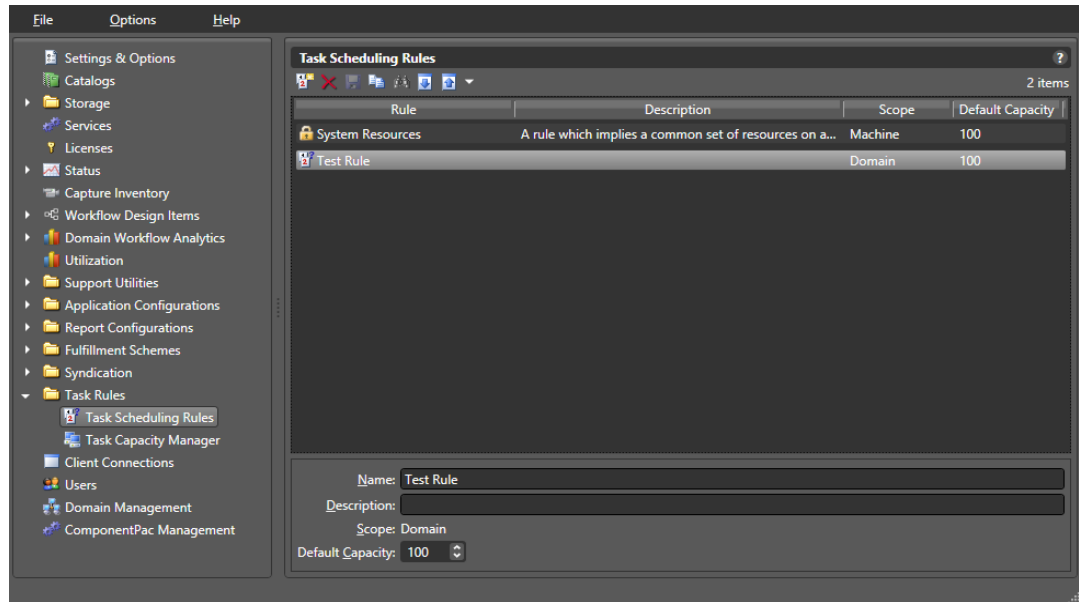


**3.** Put all services into Maintenance mode.

**4.** Set Service and Session limits for each service in the Vantage Management Console > Services panel. This limits how many actions can run within a given service to prevent service overload.
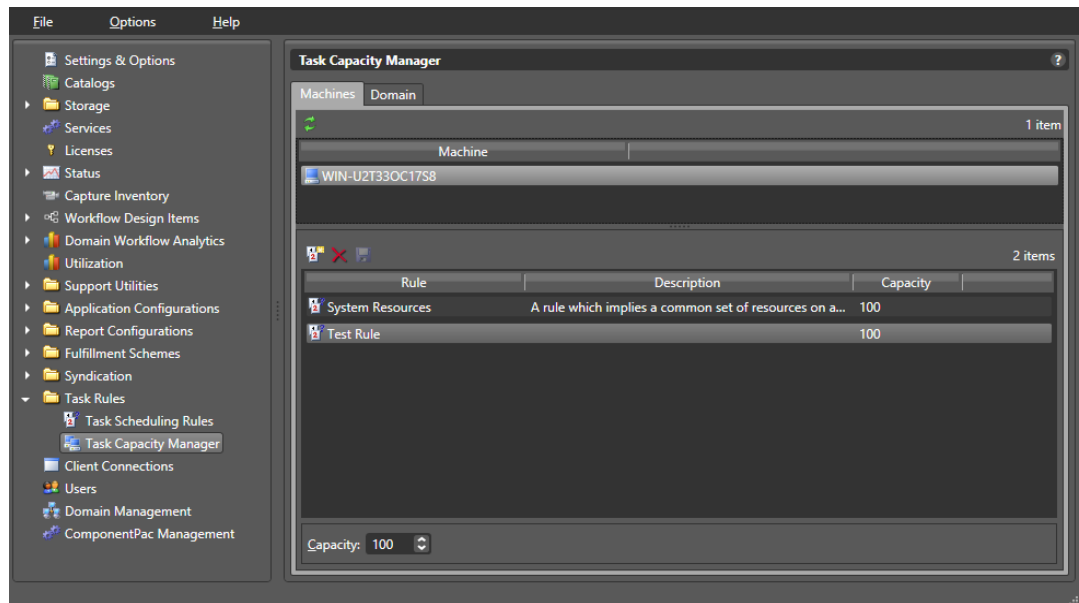


**5.** Take all Vantage services out of Maintenance mode to restart them.

**6.** Create machine and domain rules in the Vantage Management Console > Task Rules > Task Scheduling Rules. You can create multiple rules to govern machine and domain capacities under different operating circumstances. Next you will select which capacity rules to apply to particular machines and domains.



**7.** Set the capacity of machines and domains by applying rules to them in the Vantage Management Console > Task Rules > Task Capacity Manager.

**8.** Apply default costs and rules to actions using the Vantage Management Console > Workflow Design Items > Action Defaults. Use the Configure Task Scheduling button to set default values and rules to apply to each action.



**9.** When creating workflows in the Workflow Designer, you can use the action defaults, or you can customize costs and rules for specific instances. To customize an action, right-click the action, and select Resource Cost > Configure Task Scheduling.

telestream

**10.** When running workflows, you can monitor workflow performance in the Workflow Designer Status panels, such as Domain Job Status. You can monitor task scheduling rule performance in the Management Console Task Scheduling Status panel.

# Other Load Balancing Options

Some other options are available to assist with load balancing: Task Routing (Run On Rules) and the Priority Variable.

## Task Routing (Run On Rules)

Run On Rules provide another method of balancing workload. You can route actions with specialized needs to the service running on the servers in your array best suited to handle them. To implement this as you can design a workflow, define a True/False variable (certain other variables can also be used), and provide it a value of TRUE to indicate that this instance of the service can satisfy the requirement. Next, you assign the variable to the service or services that qualify, in the Vantage Management Console > Services > Variables.

For example, if you are designing an HD transcoding workflow, you can specify that the SD-to-HD Flip action must be executed on a server with sufficient power to handle HD media transcoding. For this example, the server is named Vantage_HD_Transcoder.

**1.** In Workflow Designer, open your HD workflow, right-click on the Flip action, and select Run On Rules to display the Run On Rules dialog.

**2.** Click the Add Run On Rules button in the toolbar at the top, create an HD_Capable variable, and set its value to TRUE.

**3.** Next, click OK to create this variable, and click OK to save the Run On Rule.

**4.** In the Vantage Management Console, select Services > Variables and add the HD_Capable variable to the Transcode Service running on the Vantage_HD_Transoder server.

**5.** Now, return to Workflow Designer and open your HD workflow again. Right-click on your Transcode action and select Run On Rules to display the Run On Rules dialog.

**6.** Click Check Services to verify that the Vantage_HD_Transcoder now responds correctly to queries about being HD-capable.

**Note:** When a variable is added to a service, that variable is added to the job if the service executes it. Be careful about which variables you use for Run On rules, and which are used for decision-making; generally you will want to use different variables.

Run On Rules only analyze variables explicitly set by a service. They do not analyze variables already set in a job. This ensures that Run On rules only execute actions on a service that explicitly sets them.

# The Priority Variable

Another way to balance resources is to assign the Priority Variable to particular actions in your workflows. Using this variable, you can assign higher or lower priority to specific actions during workflow execution.

1. In the Workflow Designer, right-click the action in your workflow whose priority you want to set.

2. Select Add Variables from the context menu to display the Add Variables dialog.

3. Click the Add Variable(s) button in the toolbar to display the Select Variables dialog.

4. Scroll through the list (or press the first letter of the variable—in this case P) and select the Priority variable—click OK to add it to this action. Now, the Priority variable is attached to the workflow.

5. You can change the value of the variable at the bottom of the Add Variables panel—higher numbers indicate higher priority, lower numbers indicate lower priority. Enter a large number, such as 100, to give the action highest priority, or enter a small number, such as 1 to give it a low priority.

Now when the workflow executes, the action you selected will run at the priority level you specified. High priority actions can even automatically cause low priority actions to pause, if necessary, so that the high priority action can complete first. These actions increase your control over the amount of processing power you devote to a particular action.

Note that when *Add Variables* is used to set priority, the priority is also applied to downstream actions. This can be overwritten by using *Add Variables* again on a downstream action to change or reset the priority to 0.

Also note that priority can be manually adjusted at run time when a job is in process by right-clicking the action and selecting *Priority....*

Some actions will pause lower priority actions in order to run. This can be disabled on a per service basis in Vantage Management Console -> Services.

telestream

# Easing into Load Balancing

The nice thing about Task Based Load Balancing is that it is designed so you can slowly phase into using the new capability. Begin by enabling task based load balancing:
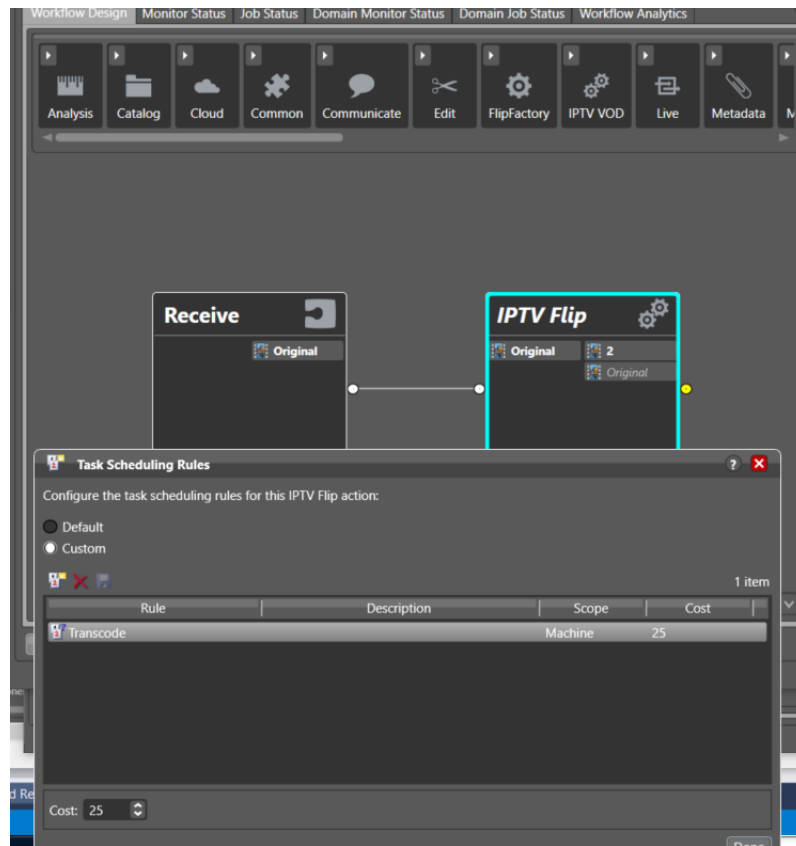


If you do nothing else, the system will behave just as it did in Session Based load balancing. Next you can begin to define rules and slowly start to use them in workflows.

You may want to start with a machine level rule: Transcode

Give this a capacity relative to the capability of the machine. Start with a value of 100, which can be refined later:

Create a test workflow and assign the Transcoding rule to the various transcoder-based actions in this workflow (Flip, Multiscreen, and so on).



You can then submit jobs to this workflow to gauge the behavior of the jobs running in Task based load balancing mode.

If you want to jump in and touch everything at one time, you could assign the default value for a certain action:

Any actions of this type (in this case Flip) would automatically run with this rule:



All of your workflows would run and enforce this rule simply by having the default set on the Flip action.

Finally, if there is any problem, all you have to do is to switch back to Session Based:



This will restore things to the way they were (after a service restart of course). Any settings and rules will be preserved, but you will be back in the session-based mode (just as you were when you started). This makes it very easy to get started with this function and roll it back if you encounter something that does not work as you.

# Conclusion

Vantage Arrays provide you exceptional flexibility in balancing server workload. Tools available to maximize workflow processing throughput and server efficiency include:

- Session-based automatic load balancing
- Cost-based configurable service and action load balancing
- Task-based load balancing of services and actions across machines or domains
- Task routing to services on particular servers
- Priority variable assignment (with automatic pause and resume) in workflow actions

Some of the above methods can be used together to fine tune your Vantage Array for optimum operation. Contact Telestream Support for more information on which load balancing tools are best for your particular application.

# Copyright and Trademark Notice

telestream