

Federator.ai User Guide

Federator.ai version 5.0 User Guide

ProphetStor Data Services, Inc.
830 Hillview Court, Suite 100
Milpitas, CA 95035 USA
Phone: 1.408.508.6255
Website: www.prophetstor.com

Copyright © 2020-2021 ProphetStor Data Services, Inc. All Rights Reserved.

Federator.ai® is a registered trademark of ProphetStor Data Services, Inc in the United States and other countries.

Kubernetes is a registered trademark of the Linux Foundation, and OpenShift is a trademark of Red Hat, Inc.

All other brand and product names are trademarks or registered trademarks of their respective owners.

12.30.2021

Contents

Overview	5
Terminology	5
Getting Started	7
Access the Federator.ai Portal.....	7
Setup Wizard	9
Kubernetes Cluster	9
VM Cluster	13
Federator.ai Administration Portal	15
Portal Sections	15
Portal Icons.....	15
Common Administration Portal Functions	16
Refresh Statistics	16
License Status	16
User Functions.....	17
Filters	17
Specify Time Range.....	17
Search/Sort Information in Tables.....	18
Show/Hide Metrics in Charts	18
Zoom In/Out of Charts.....	19
Dashboard	20
Cluster Workload Prediction	20
Cluster Charts	20
Node/Virtual Machine Chart.....	22
Namespace Chart (Kubernetes)	22
Application Workload Prediction (Kubernetes)	24
Application Charts	24
Controllers Chart	25
Application Insight.....	27
Topology	28
Application Correlation.....	29
Cluster Node Correlation	32
Resource Utilization.....	35

Resource Utilization - Controllers	35
Resource Utilization - Nodes	36
Resource Predictions	38
Resource Predictions - Controllers	38
Resource Predictions - Nodes	39
Application KPI Metrics.....	40
Cluster Overview - Cluster Health	41
Cluster Overview - Node Health (Kubernetes)	43
Predictions and Planning – Kubernetes or VM Resources	44
Managed Nodes Table (Kubernetes)	44
Managed VMs Table (VM)	45
Managed Containers Table (Kubernetes)	45
Workload Prediction Table and Workload Observation and Prediction Charts	45
Workload Prediction Table	46
Workload Observation and Prediction Charts	47
Utilization Analysis Charts	48
CPU and Memory Utilization Heatmap Charts.....	48
CPU and Memory Utilization Goals Charts	49
Autoscaling - HPA (Kubernetes)	50
CPU and Memory Charts	50
Autoscaling – Kafka Consumer (Kubernetes)	51
Number of Replicas Chart.....	51
Production Rate and Consumption Rate Chart	51
Consumer Lag Chart	52
Consumer Queue Latency Chart.....	52
CPU and Memory Observation Charts.....	52
Autoscaling – Ingress Upstream Services (Kubernetes).....	53
Number of Replicas Chart.....	53
HTTP Request Rate Chart	53
HTTP Response Error Rate Chart	54
Average Response Time Chart.....	54
Upstream Latency Chart.....	55
CPU and Memory Observation Charts.....	55
Cost Management – Cost Analysis.....	56

Cost Analysis Charts.....	56
Cost Analysis Summary Table	57
Cost Management – Cost Trends.....	58
Cost Trends Chart	58
Cost Trends Summary Table	59
Cost Management – Cost Optimization	60
Cost Optimization Chart	60
Cost Efficiency Charts	60
Cluster Cost Optimization.....	61
Cluster Cost Efficiency Table	61
Cluster Recommendations Table	61
Cluster Cost Optimization Details	62
Cluster Node Cost Optimization	63
Cluster Node Cost Efficiency Table	63
Cluster Node Recommendations Table.....	63
Cluster Node Cost Optimization Details.....	64
Namespace Cost Optimization	65
Namespace Cost Efficiency Table.....	65
Namespace Optimization Details	65
Application Cost Optimization	66
Application Cost Efficiency Table	66
Application Optimization Details	67
Cost – Multi-cloud Cost Analysis	68
Resource Utilization and Cost Efficiency Charts	70
Recommended Cluster Configuration	71
Configuration - Clusters	73
Kubernetes Clusters	73
Add a Kubernetes Cluster.....	74
Manage Kubernetes Clusters	76
VM Clusters	77
Add a VM Cluster.....	78
Manage VM Clusters	80
Configuration – Applications.....	81
Add an Application	82

Generic Kubernetes Application	83
Kafka Consumer Application	85
Ingress Application	87
Manage Applications	89
Configuration – Auto Provisioning	90
Auto Provisioning Scripts	91
Add a Profile	92
Manage Profiles	93
Configuration – System Settings.....	94
Admin Password	94
Metrics Data Source	94
Notification	95
Enable Notification	95
Licenses	97
Add a License.....	100
Manage Licenses	101
Price Books	102
Public Cloud Price Books	102
Custom Price Book	102
Events.....	103

Overview

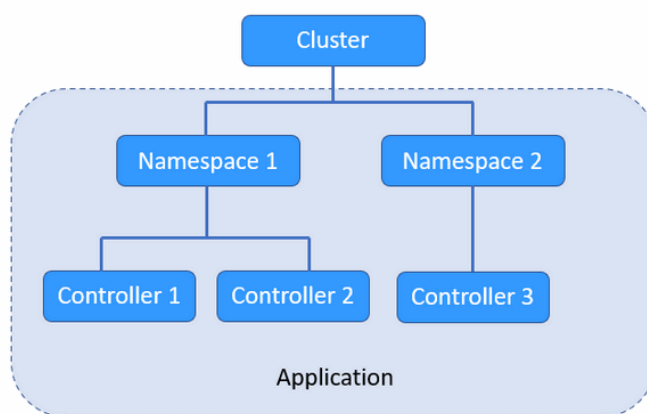
ProphetStor Federator.ai is an AI-based solution that helps enterprises manage and optimize resources for applications on Kubernetes and virtual machines (VMs) in VMware clusters. Using advanced machine learning algorithms to predict application workloads, Federator.ai offers:

- AI-based workload prediction for containerized applications in Kubernetes clusters as well as VMs in VMware clusters and Amazon Web Services (AWS) Elastic Compute Cloud (EC2)
- Resource recommendations based on workload prediction, application, Kubernetes, and other related metrics
- Automatic provisioning of CPU/memory for generic Kubernetes application controllers/namespaces
- Correlation and causality analysis of microservices/controllers of Kubernetes applications
- Automatic scaling of Kubernetes application containers, Kafka consumer groups, and Ingress upstream services
- Multicloud cost analysis and recommendations based on workload predictions for Kubernetes clusters and VM clusters
- Actual cost and potential savings based on recommendations for clusters, Kubernetes applications, VMs, and Kubernetes namespaces

If you have not installed Federator.ai yet, refer to your *Federator.ai Installation Guide* for information.

Terminology

Application – Defined by Federator.ai as a group of Kubernetes controllers that work together to serve tasks from the view of the end user. For example, an e-commerce web application consists of controllers for frontend and backend and can be considered as an application. An application is not a Kubernetes object.



Auto Provisioning – The ability to automatically deploy CPU and memory resource recommendations to controllers and namespaces of generic applications in Kubernetes clusters based on pre-defined profiles.

Autoscaling – In Kubernetes, the ability for the system to automatically increase or decrease containers/pods based on workload demands.

Auto Scaling group – (ASG) A collection of Amazon EC2 instances that are treated as a logical grouping for automatic scaling and management.

Container – An object that contains a software module with everything needed to run an application.

Controller – In Kubernetes, controllers are control loops that watch the state of your cluster, then make or request changes where needed. Each controller tries to move the current cluster state closer to the desired state. The types of controllers supported by Federator.ai are *Deployment* and *StatefulSet*. Additionally, Federator.ai supports *DeploymentConfig* controllers for OpenShift.

Cluster – A Kubernetes cluster with one or more nodes or a VM cluster with one or more VMs.

Deployment – A Deployment provides declarative updates for Pods and ReplicaSets. The user describes a desired state in a deployment and the deployment controller changes the actual state to the desired state at a controlled rate.

HPA – Horizontal Pod Autoscaling – In Kubernetes, the system automatically increases or decreases the number of containers/pods (replicas) based on the workload.

Namespace – Kubernetes supports multiple virtual clusters backed by the same physical cluster. These virtual clusters are called namespaces.

Microservice – Also known as controllers, microservices are independent, modular Kubernetes components that work together as a single application.

Node – In Kubernetes, nodes are server-like machines, such as a virtual machine running complete systems and multiple applications. There can be master nodes and worker nodes.

Pod – A group of one or more containers with shared storage/network resources and a specification for how to run the containers. Typically, one container runs in each pod.

Replica – A copy of a pod running for an application.

StatefulSet – A Kubernetes object that manages stateful applications. Unlike a Deployment, a StatefulSet maintains a sticky identity for each of its pods that remains the same across any rescheduling.

VM – A VMware virtual machine running on a physical *host* machine.

VM Cluster – A cluster with one or more VMs.

Related topics:

[Federator.ai Administration Portal](#)

[Configure Kubernetes Clusters](#)

[Configure VMWare Clusters](#)

[Configure Applications](#)

[Auto Provisioning](#)

Getting Started

After installation of Federator.ai, you must access the Federator.ai portal in order to configure your system.

Access the Federator.ai Portal

To access the Federator.ai administration portal, use the URL that is displayed at the end of the installation process.

You can also find the URL for the Federator.ai administration portal via the following methods:

Kubernetes

In a Kubernetes environment, use the `kubectl` command to find the administration portal service port number and node IP address.

```
# kubectl get svc -n federatorai |grep federatorai-dashboard-frontend-node-port
```

The output will look something like this:

```
federatorai-dashboard-frontend-node-port NodePort 10.103.181.133 <none> 9001:31012/TCP
```

Get the node's IP to access (INTERNAL-IP).

```
$kubectl get nodes -o wide
```

For example:

```
# kubectl get nodes -o wide
NAME      STATUS    ROLES    AGE   VERSION   INTERNAL-IP   EXTERNAL-IP   OS-IMAGE
KERNEL-VERSION   CONTAINER-RUNTIME
h7-130    Ready     master   35d   v1.18.5   172.31.7.130  <none>        CentOS Linux 7 (Core)
3.10.0-957.el7.x86_64  docker://19.3.13
h7-131    Ready     <none>   35d   v1.18.5   172.31.7.131  <none>        CentOS Linux 7 (Core)
3.10.0-957.el7.x86_64  docker://19.3.13
h7-132    Ready     <none>   35d   v1.18.5   172.31.7.132  <none>        CentOS Linux 7 (Core)
3.10.0-957.el7.x86_64  docker://19.3.13
h7-133    Ready     <none>   35d   v1.18.5   172.31.7.133  <none>        CentOS Linux 7 (Core)
3.10.0-957.el7.x86_64  docker://19.3.13
```

The URL will be `https://172.31.7.130:31012`.

OpenShift

In an OpenShift environment, use the `oc get route` command to find the URL.

```
# oc get route -n federatorai | grep federatorai-dashboard-frontend
```

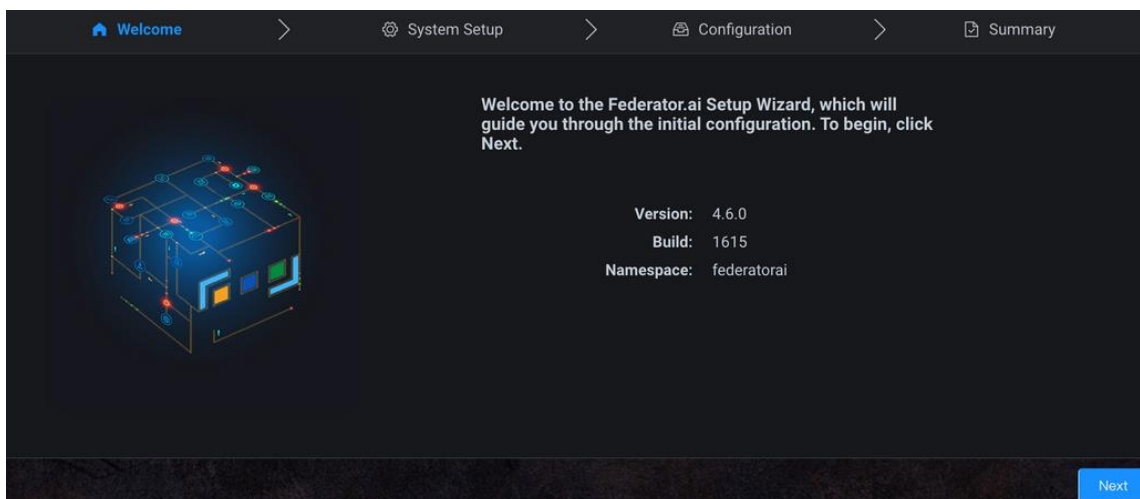
The output will look something like this:

```
federatorai-dashboard-frontend "federatorai-dashboard-frontend-federatorai.apps.ocp4.172-31-11-30.nip.io"
```

The URL will be `https://federatorai-dashboard-frontend-federatorai.apps.ocp4.172-31-11-30.nip.io`

Setup Wizard

The first time you log in after installation of Federator.ai, a setup wizard launches that allows you to configure a cluster that should be monitored by Federator.ai. You can add more clusters from the *Configuration* section of the portal after completing the setup wizard. Click *Next* to begin configuration.



Kubernetes Cluster

1. Set the administrator password and the name of a cluster to be monitored by Federator.ai, select *Kubernetes Cluster* and the source of metrics for this cluster, and specify authentication information. Then, click *Test Connection* to confirm that all information is correct. Click *Next* when the system can connect to the cluster.

The screenshot shows the 'System Setup' step of the Federator.ai Setup Wizard. The navigation bar at the top is the same as the previous screen, with 'System Setup' now active. The main content area is divided into two sections. The first section, titled 'Administrator account' with a user icon, contains the instruction 'Set the administrator password.' and two password input fields: 'New Password' and 'Confirm Password'. The second section, titled 'Monitoring & Metrics Data Source' with a database icon, contains the instruction 'Specify the name of a cluster to be monitored by Federator.ai and the source of metrics for this cluster.' This section includes a 'Cluster Name' input field, a 'Cluster Type' selection with 'Kubernetes Cluster' selected (radio button), and a 'Metrics Data Source' selection with 'Prometheus' selected (radio button). Below these, there is a checkbox for 'Federated Prometheus' which is checked, a 'Target Label' input field with the placeholder '<label-name><label-value>', a 'URL' input field, and a 'Token' input field. A blue 'Test Connection' button is located at the bottom right of the form. At the very bottom of the screen, there are 'Back' and 'Next' buttons.

You must set the administrator password in order to continue but your cluster can be configured later from the *Configuration* section of the portal after completing the setup wizard.

For the Prometheus open-source monitoring system, the URL is required but the token is optional for authentication. Specify if you are using Federation, which is a group of Prometheus servers that send metrics to a centralized Prometheus server. You will need to specify the target label of the centralized Prometheus server. The format is: <label-name>:<label-value> (e.g., clusterID:host-1).

The screenshot shows the Prometheus configuration interface. At the top, 'Cluster Type' has 'Kubernetes Cluster' selected. Below, 'Metrics Data Source' has 'Prometheus' selected. A checkbox for 'Federated Prometheus' is checked, and the 'Target Label' field contains 'clusterId:cluster-189'. The 'URL' field contains 'http://prometheus-operator-prometheus.monitoring.svc:9090'. There is an empty 'Token' field. A 'Test Connection' button is at the bottom right, along with 'Back' and 'Next' navigation buttons.

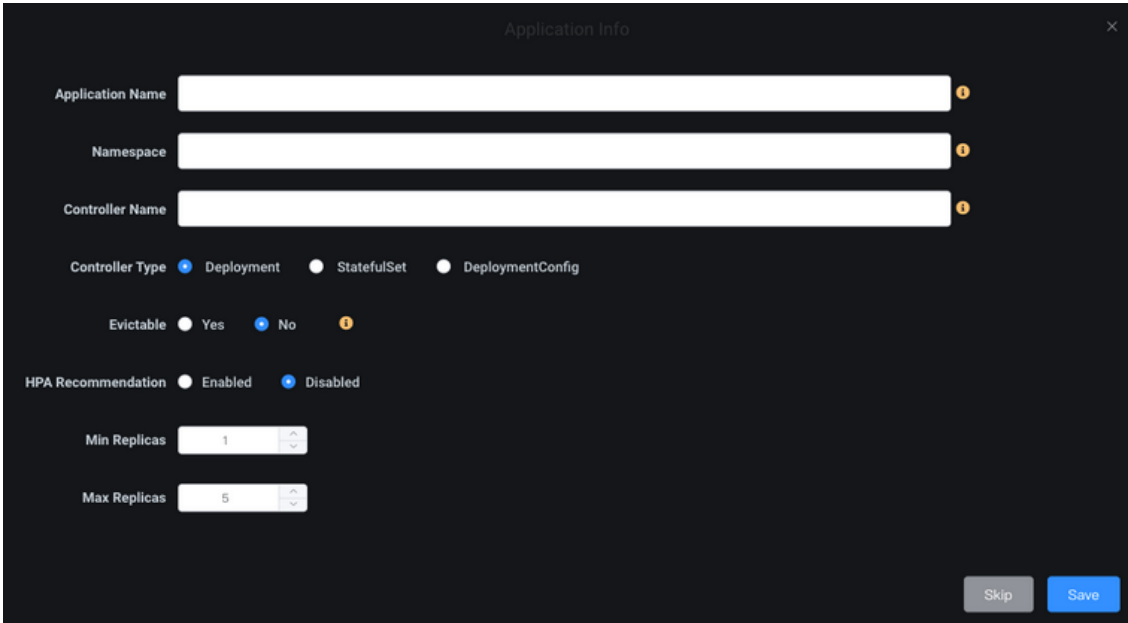
For Datadog, the API key and application key are required for authentication. If needed, you can click on the link to the Datadog website, which is included in the popup help text.

The screenshot shows the Datadog configuration interface. 'Cluster Type' is 'Kubernetes Cluster'. 'Metrics Data Source' has 'Datadog' selected. There are two required fields: 'API Key' and 'Application Key', both containing placeholder text. A 'Test Connection' button is at the bottom right, along with 'Back' and 'Next' navigation buttons.

For Sysdig, a URL and token are required for authentication. If needed, you can click on the link to the Sysdig website, which is included in the popup help text.

The screenshot shows the Sysdig configuration interface. 'Cluster Type' is 'Kubernetes Cluster'. 'Metrics Data Source' has 'Sysdig' selected. There are two required fields: 'URL' containing 'https://app.sysdigcloud.com' and 'Token' containing a placeholder. A 'Test Connection' button is at the bottom right, along with 'Back' and 'Next' navigation buttons.

2. Enter information about the first application you want to monitor. You can skip this step and configure applications to be monitored from the *Configuration* section of the portal.



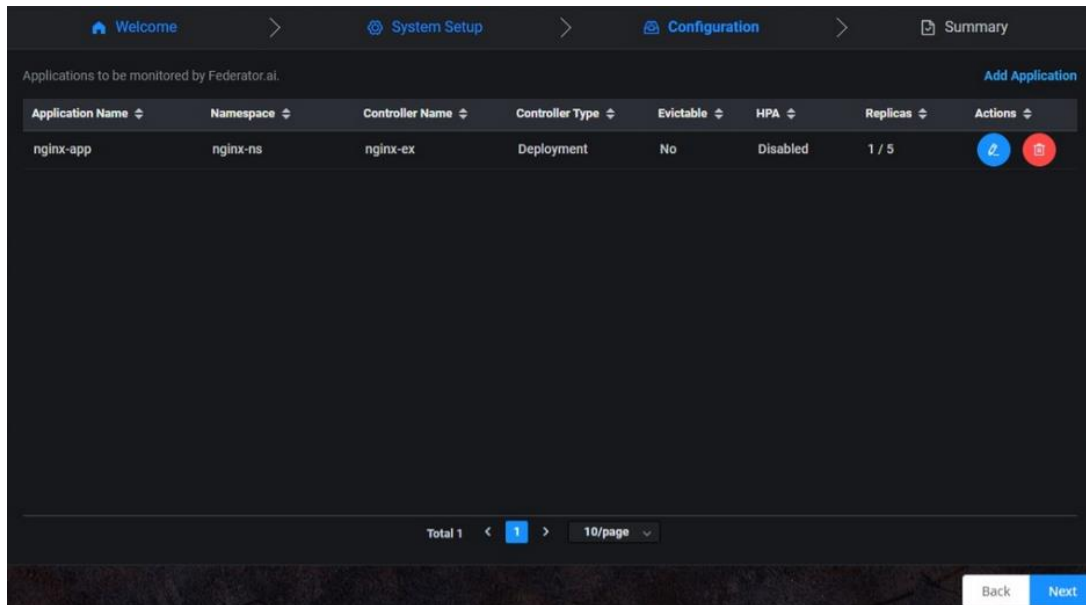
The screenshot shows a dark-themed 'Application Info' form. It contains the following fields and options:

- Application Name:** A text input field with a yellow information icon to its right.
- Namespace:** A text input field with a yellow information icon to its right.
- Controller Name:** A text input field with a yellow information icon to its right.
- Controller Type:** Radio buttons for ☒ Deployment, ☐ StatefulSet, and ☐ DeploymentConfig.
- Evictable:** Radio buttons for ☐ Yes and ☒ No, with a yellow information icon to the right.
- HPA Recommendation:** Radio buttons for ☐ Enabled and ☒ Disabled.
- Min Replicas:** A numeric input field with a value of 1 and up/down arrows.
- Max Replicas:** A numeric input field with a value of 5 and up/down arrows.
- Buttons:** 'Skip' and 'Save' buttons at the bottom right.

You can add applications now or from the *Configuration* section of the portal after completing the setup wizard.

- *Application Name* - The name of your application to be monitored by Federator.ai. An application is a group of one or more Kubernetes controllers that work together to serve tasks from the view of the end user; an application is not a Kubernetes object.
- *Namespace* – The Kubernetes namespace where the controller is deployed.
- *Controller Name* - The name of controller to be monitored.
- *Controller Type* – Supported controller types are *Deployment*, *StatefulSet*, and *DeploymentConfig* (OpenShift only).
- *Evictable* – Indicate if the controller can be interrupted if the node is shut down. Evictable controllers are good candidates to be deployed in Spot instances.
- *HPA Recommendation* – Indicate if you want to enable Horizontal Pod Autoscaling (HPA). When enabled, CPU and memory usage is monitored, and the number of pods is automatically increased/decreased based on the CPU/memory usage workload. HPA and Auto Provisioning are mutually exclusive; you can use HPA or auto provisioning, but not both.
- *Min/Max Replicas* – Specify the minimum and maximum number of pods when HPA is enabled.

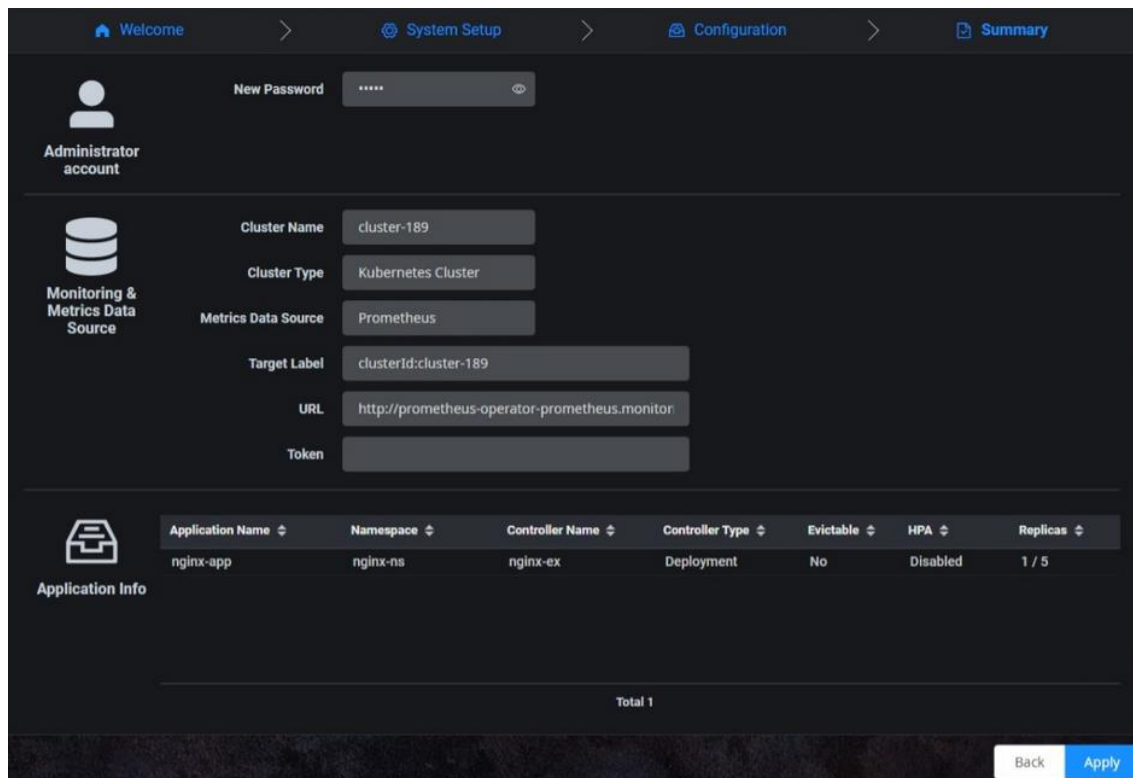
3. Click *Save* to view a list of applications that will be monitored by Federator.ai.



Click the + key to add an application.

Click the blue *Edit* icon to update application information or click the red *Delete* icon to remove an application.

- Click *Next* to view the *Summary* screen.



- Click *Apply* to apply all changes or click *Back* to edit information.

VM Cluster

1. Set the administrator password and the name of a cluster to be monitored by Federator.ai, select *VM Cluster* and *vCenter* or *AWS CloudWatch* as the source of metrics for this cluster, and specify connection information.

The screenshot shows the 'System Setup' step of the Federator.ai setup wizard. The interface is dark-themed. At the top, there's a navigation bar with 'Welcome', 'System Setup' (active), 'Configuration', and 'Summary'. Below the navigation bar, there are two main sections. The first section, 'Administrator account', has a user icon and a label 'Administrator account'. It contains two password fields: 'New Password' and 'Confirm Password', both with masked characters (*****). The second section, 'Monitoring & Metrics Data Source', has a database icon and a label 'Monitoring & Metrics Data Source'. It contains several fields: 'Cluster Name' (value: cluster-189), 'Cluster Type' (radio buttons for 'Kubernetes Cluster' and 'VM Cluster', with 'VM Cluster' selected), 'Metrics Data Source' (radio buttons for 'vCenter' and 'AWS CloudWatch', with 'vCenter' selected), 'vCenter' (value: 172.31.2.189), 'Login ID' (value: vsphere.local/administrator), 'Password' (masked), and 'Cluster Path' (value: Datacenter/172.31.2.189). There are 'Test Connection' and 'Back' buttons at the bottom right, and a 'Next' button at the bottom right.

You must set the administrator password in order to continue but your cluster can be configured later from the *Configuration* section of the portal after completing the setup wizard.

The cluster name must have a maximum of 253 lowercase characters, "-", or "." allowed. The name must start and end with an alphanumeric character.

For vCenter, specify the vCenter IP address (you can have multiple vCenters in your system), login ID, password, and the path to the cluster, within vCenter. If needed, you can click on the link to the vCenter website, which is included in the popup help text.

For AWS CloudWatch, specify the region of Amazon AWS S3 service, the AWS Identity and Access Management key ID (16 to 128 bytes), and the secret access key of the key ID that is used for access. Note: The CloudWatch agent must be installed on the EC2 node in order to use this data source.

Federator.ai Administration Portal

The Federator.ai administration portal displays the overall health of each cluster, as well as application workload and resource recommendations. Information is presented in tables and charts.


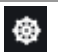
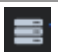









Portal Sections

The portal is separated into the following sections:

- **Dashboard** – Overall system information, including the number of monitored resources, as well as cluster and Kubernetes application workload predictions and recommendations.
- **Application Insight** - Provides statistical analysis and predictions based on the correlation between resource usage and application workload for Kubernetes.
- **Cluster Overview** - Cluster and Kubernetes node health information, including CPU utilization, memory utilization, disk capacity, and VM network throughput.
- **Predictions and Planning** – Forecasting tools, including actual CPU and memory usage observations, predicted workload usage, utilization analysis, and recommendations for Kubernetes and VM resources.
- **Autoscaling** - (Kubernetes) Displays HPA recommendations for controllers enabled with autoscaling and autoscaling predictions for each Kafka individual topic/consumer group and Ingress upstream service
- **Cost** – cost analysis, cost trends, and cost optimization for resources at different levels, including Kubernetes and VM clusters and nodes, as well as Kubernetes namespaces and applications.
- **Configuration** – Configuration of clusters, Kubernetes applications and controllers/consumer groups, as well as system configuration, including resetting the admin password, metrics data source, system notifications, licensing, and price books.
- **Events** – System events that have occurred.

Portal Icons

To make it easy to distinguish between cluster types and providers, the following icons are used throughout the portal:

Icon	Meaning						
	On-premises provider						
	Kubernetes clusters						
	VM clusters. The cluster type is: <table><tr><td></td><td>AWS CloudWatch VM cluster configured with AWS Auto Scaling groups.</td></tr><tr><td></td><td>AWS CloudWatch VM cluster with individual VMs.</td></tr><tr><td></td><td>vCenter cluster.</td></tr></table>		AWS CloudWatch VM cluster configured with AWS Auto Scaling groups.		AWS CloudWatch VM cluster with individual VMs.		vCenter cluster.
	AWS CloudWatch VM cluster configured with AWS Auto Scaling groups.						
	AWS CloudWatch VM cluster with individual VMs.						
	vCenter cluster.						

Common Administration Portal Functions

The administration portal presents information in tables and charts. At the top right of each portal page, you can do the following:

- Refresh statistics
- Check license status
- Get technical support contact information
- View Federator.ai product documentation
- Display the product software version
- Log out



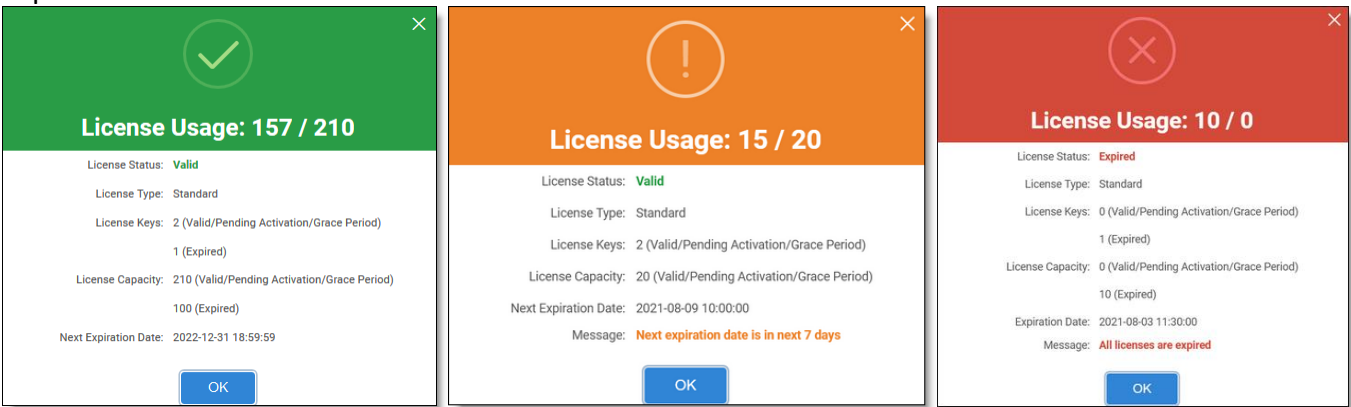
Refresh Statistics

By default, Federator.ai information is refreshed every five minutes. To change the interval, click the drop-down at the top right and select a 1, 5, 15, or 30-minute interval.

To force a refresh, click *Refresh Now* where the current interval is displayed.

License Status

Click the *License* icon at the top right of the dashboard to see Federator.ai license information, including license status and type, number and type of license keys, licensed capacity and usage, as well as license expiration. When the icon is green, all licenses are valid. Orange indicates a situation that requires attentions, such as a license is near expiration, a license is in a grace period, or the number of monitored resources exceeds the license limit. Red requires immediate attention because it indicates a license has expired.

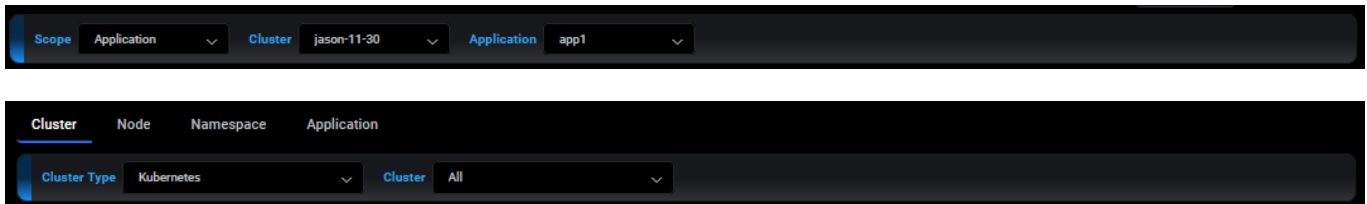


User Functions

Click the *User* icon to contact technical support, view the Federator.ai product documentation, display the product software version, or log out from the system.

Filters

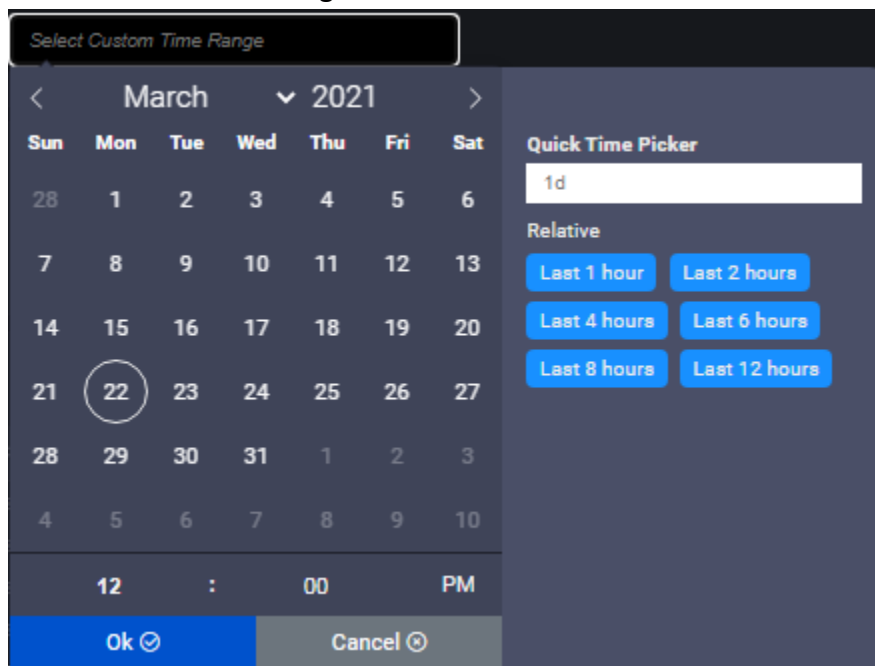
A panel appears on most pages that allows you to filter the display of information. Depending upon the page and cluster type, you may be able to select a scope/resource type, cluster, application, namespace, controller, time range, etc.



Specify Time Range

When a chart allows you to specify a time range to display data, you can select a predefined time frame (e.g., last 1 hour, last 24 hours) from the drop-down box or you can specify a custom time range via one of the following methods:

- Use the calendar to select the start and end dates.
- Specify the number of hours (e.g., 5h), days (e.g., 5d), weeks (e.g., 3w), or months (e.g., 6m) in the *Quick Time Picker* box.
- Click a predefined relative time range.



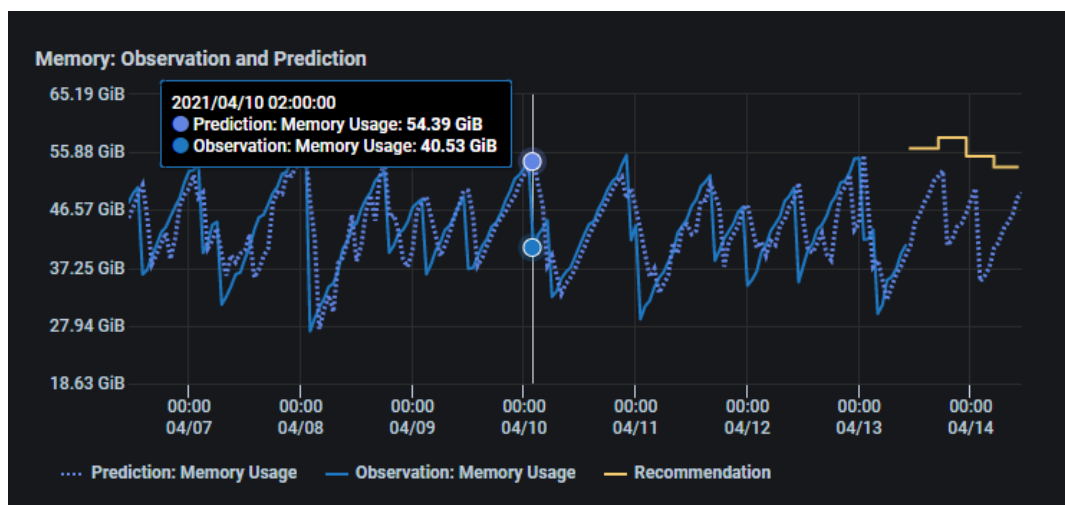
Search/Sort Information in Tables

Click a column heading to sort the entire list based upon values in that column. The blue highlighted triangle or inverted triangle on the column indicates the direction of the sort. To search, type a name or value in the *Search* box; clear the search field to return to the full view of the list. As soon as you start typing, only those items that have matching text are displayed. You can also determine how many rows to show per page (5, 10, or 20).

Managed Containers				
<input type="text" value="Search"/>				
Container ↕	Project (Names...	Application ↕	Pod ↕	Node ↕
my-nginx	nginx2	nginx2-jason-6-...	my-nginx-6f97b...	h6-182
my-nginx	nginx1	alamedascaler-...	my-nginx-6c99d...	h6-182
my-nginx	nginx1	alamedascaler-...	my-nginx-6c99d...	h6-182
my-nginx	nginx1	alamedascaler-...	my-nginx-6c99d...	h6-182
my-nginx	nginx1	alamedascaler-...	my-nginx-6c99d...	h6-182
Total 6 < 1 2 > 5/page ▾				

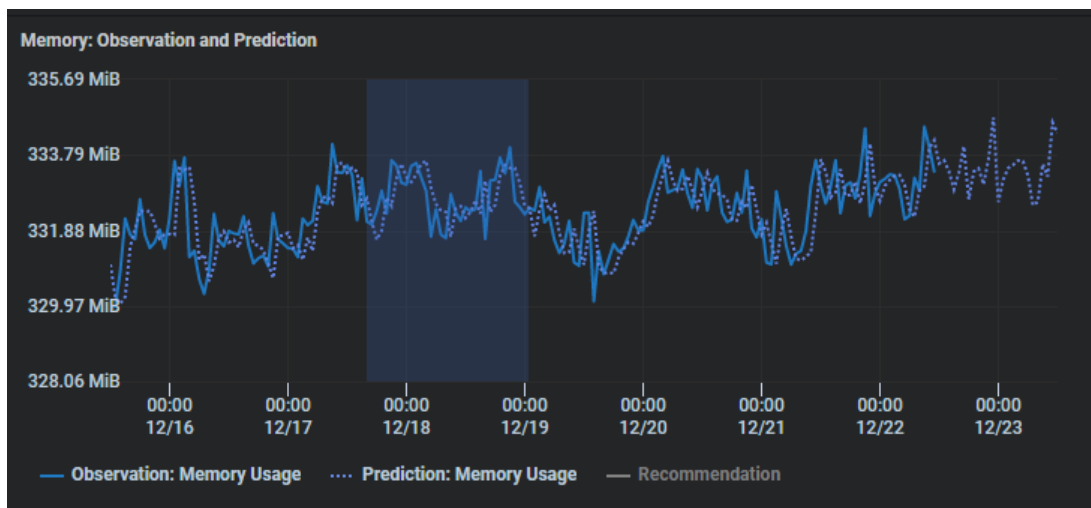
Show/Hide Metrics in Charts

Click anywhere on a chart to see values for a specific point in time. Highlight or click on the key at the bottom of the chart to show/hide individual metrics.

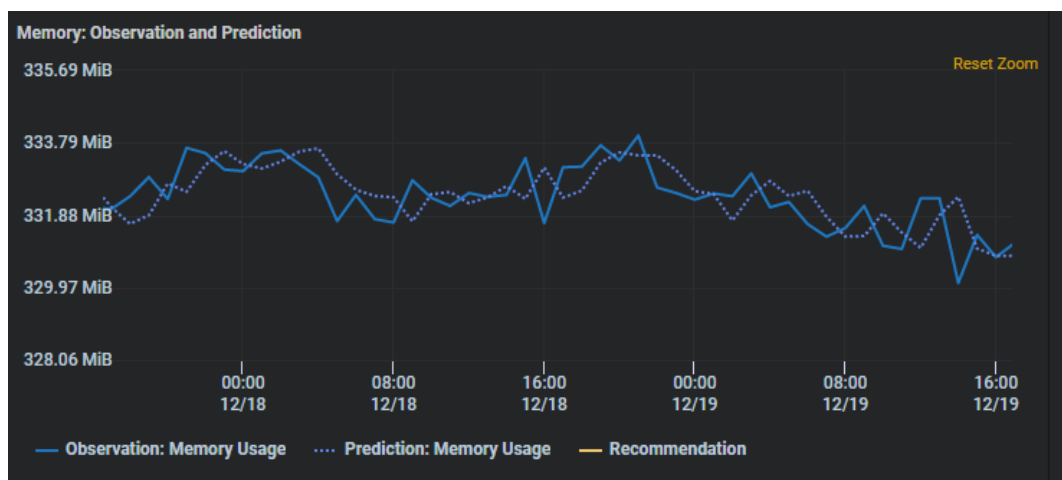


Zoom In/Out of Charts

Click and drag a pointer in a chart to zoom into a specific time frame of interest.



Click *Reset Zoom* to return to the original time frame.



Related topics:

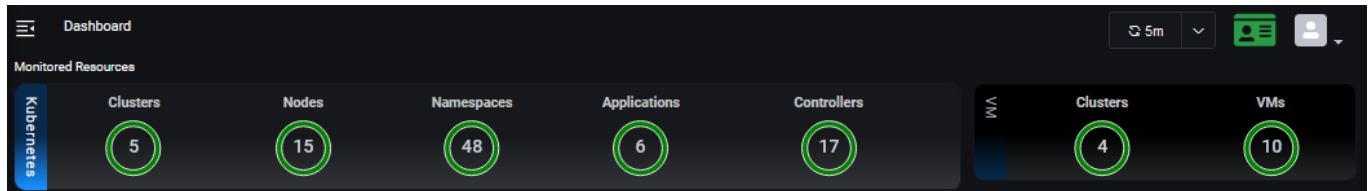
[Common Administration Portal Functions](#)

[Dashboard](#)

[Licenses](#)

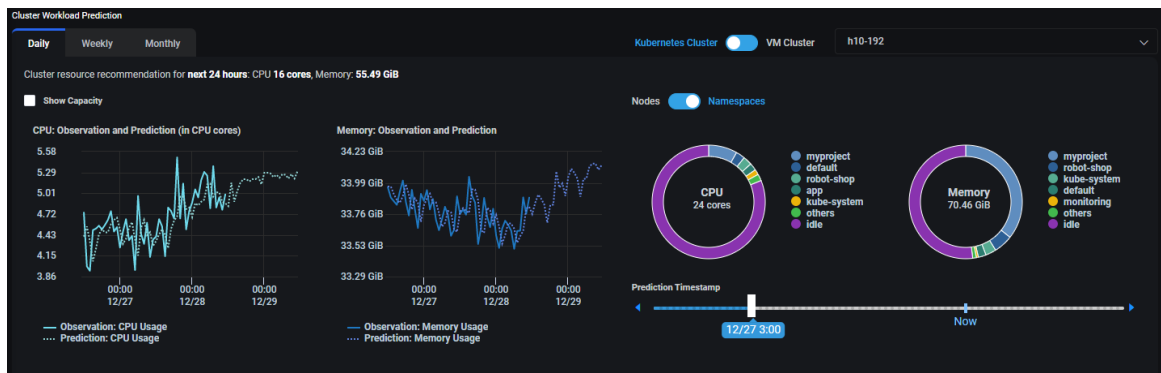
Dashboard

The *Dashboard* displays the number of monitored resources in Kubernetes clusters and VM clusters, as well as cluster and Kubernetes application workload predictions and recommendations.



Cluster Workload Prediction

Toggle between *Kubernetes Cluster* and *VM Cluster*, select a cluster from the drop-down list, and select the time frame (daily, weekly, or monthly) to display CPU and memory observations, predictions, and recommendations.



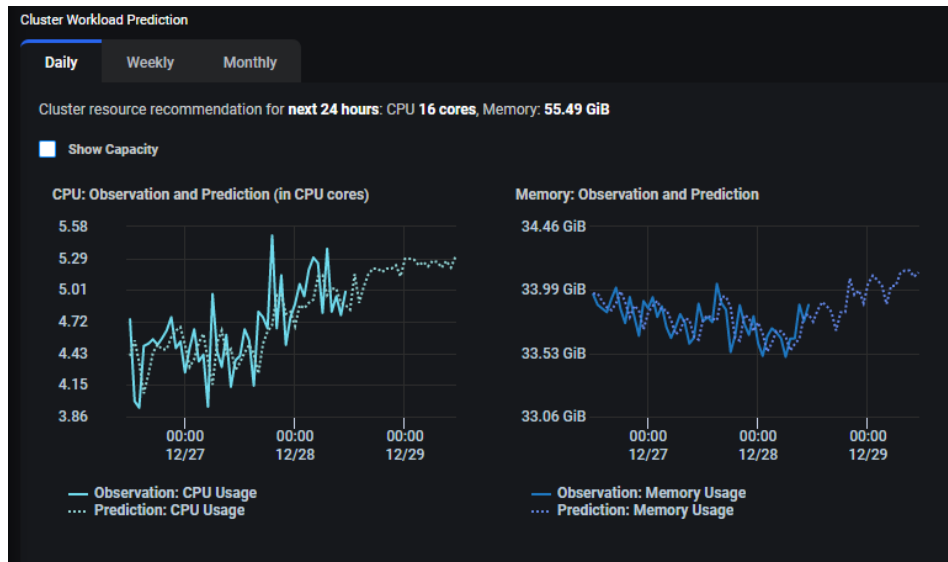
The first two charts display information for the selected cluster; the chart on the right toggles the display of information for nodes and namespaces of the selected Kubernetes cluster. For a VM cluster, the chart on the right displays information for the VMs in the selected cluster.

The text above the charts summarizes the CPU and memory recommendations for the next 24 hours (daily), 7 days (weekly), or 30 days (monthly).

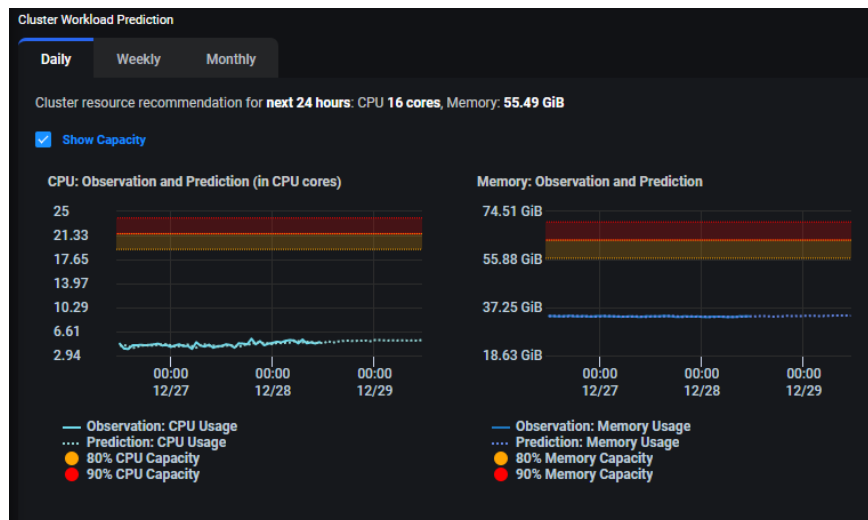
Cluster Charts

The cluster charts display CPU and memory observations and predictions for the cluster.

The solid lines represent the observed actual usage while the dotted lines show the historical and future predicted usage. Click anywhere on the charts to see values for a specific point in time. This will adjust the slider in the Nodes/Namespaces chart accordingly.



Check *Show Capacity* to see the maximum CPU and memory usage limits for the cluster. Orange represents 80-90% and red represents 90-100%. This is a useful way to see if the utilization of resources is approaching the overall cluster capacity.

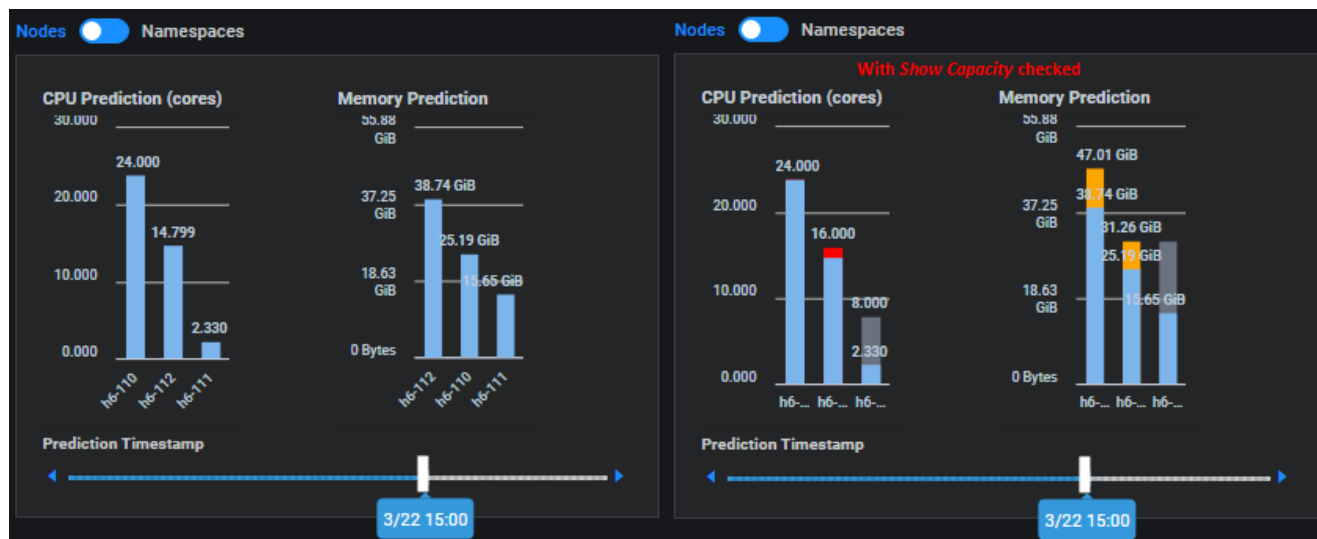


Node/Virtual Machine Chart

This chart displays CPU and memory observations and predictions for each member of the cluster. For Kubernetes, toggle to *Nodes* to display this chart.

The slider always starts at the current time but allows you to select any day/hour. Slide to the left of *Now* for historical usage; slide right for future predictions.

If *Show Capacity* was selected for the cluster, the chart will show if the node's utilization of resources is approaching the maximum CPU and memory usage limits. Orange represents less than 20% availability and red represents less than 10% availability.

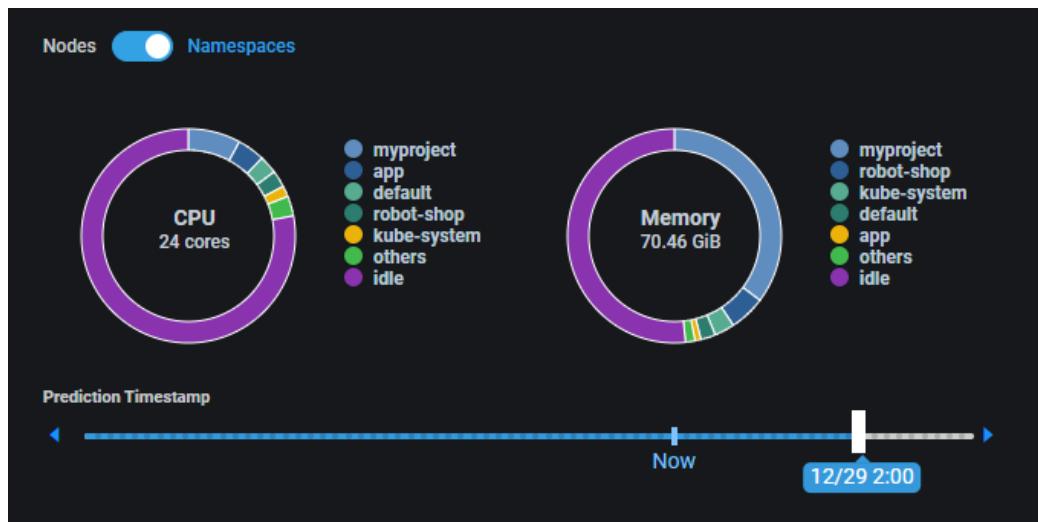


Namespace Chart (Kubernetes)

Toggle to *Namespaces* to display CPU and memory predictions for each namespace and for when the system is idle. Move your cursor over each section to show the usage or predicted usage by a specific namespace and the overall percentage used or predicted to use by the namespace or when the resource is idle.

Highlighting the amount of predicted idle (unused) resources provides a useful way to determine where your cluster is over-provisioned and can help you balance the resource allocation within the cluster.

The slider always starts at the current time but allows you to select any day/hour. Slide to the left of *Now* for historical predictions; slide right for future predictions.



Application Workload Prediction (Kubernetes)

Select an application from the drop-down list and select the time frame (daily, weekly, or monthly) to display CPU and memory observations and predictions.

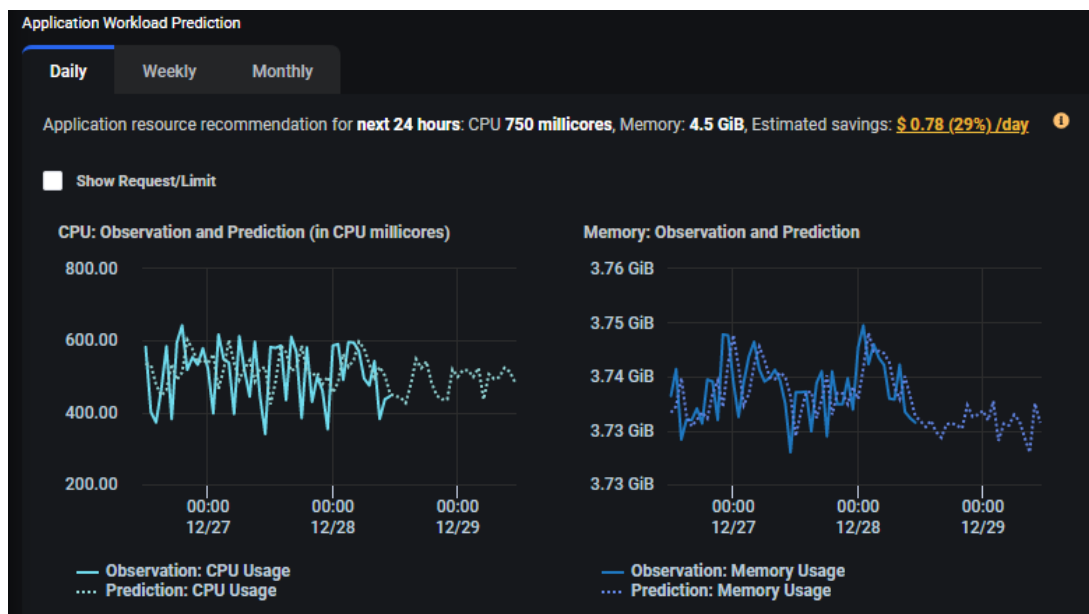
The first two charts display information for the selected application; the chart on the right displays information for controllers of the application.

The text above the charts summarizes the CPU and memory recommendations for the next 24 hours (daily), 7 days (weekly), or 30 days (monthly). It may also show estimated savings based on the system recommendations. You can link to the *Cost Optimization/Application* page for more details.

Application Charts

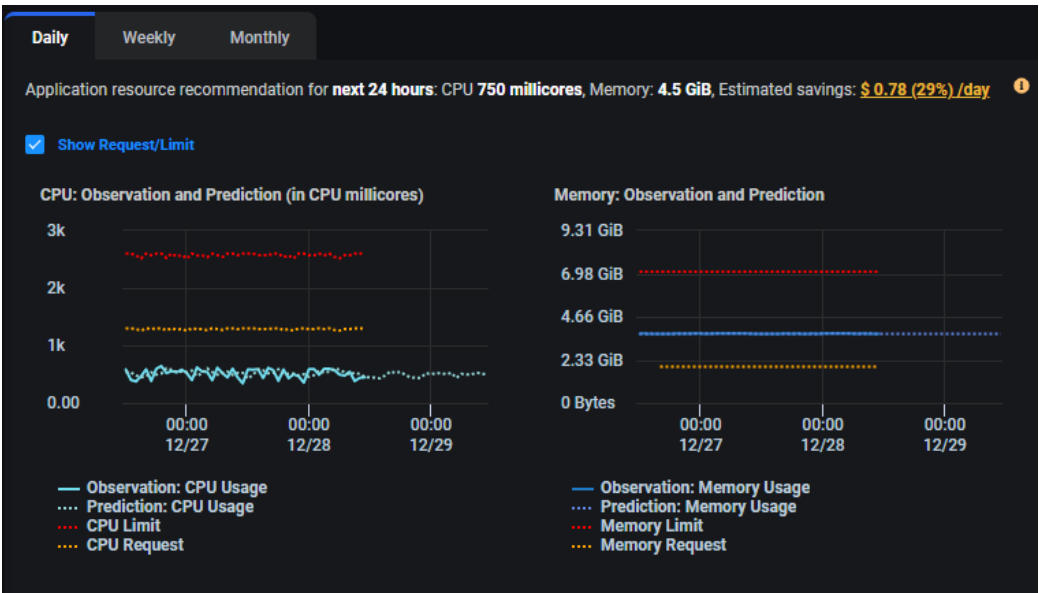
The application charts display CPU and memory observations, predictions, and recommendations for the application.

The solid lines represent the observed actual usage while the dotted lines show the historical and future predicted usage. Click anywhere on the charts to see values for a specific point in time. This will adjust the slider in the Controllers chart accordingly.



Select *Show Request/Limit* to see your application’s CPU and memory usage limits in Kubernetes. Orange represents your resource request and red represents the resource limit.

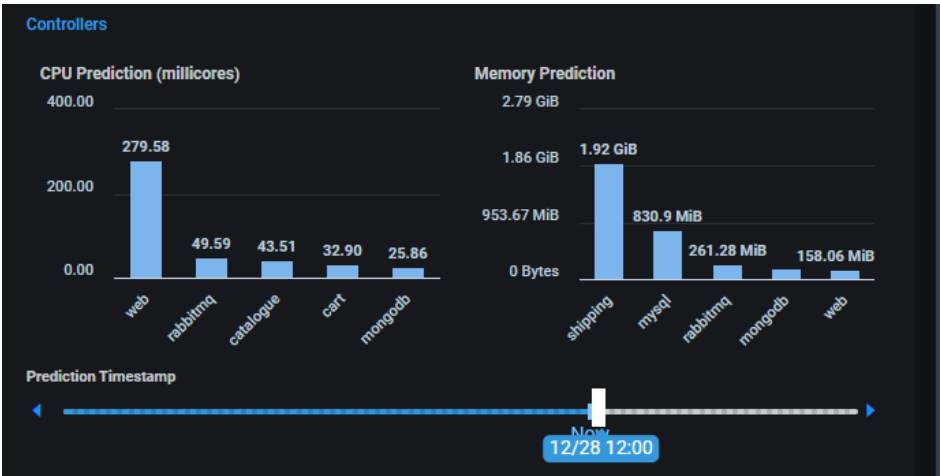
Generally, memory is a hard limit, but CPU is more *stretchable*. This is a useful way to see if you are over-provisioned or under-provisioned for your application.



Controllers Chart

The chart displays CPU and memory observations and predictions for each controller defined in a generic application.

The slider always starts at the current time but allows you to select any day/hour. Slide to the left of *Now* for historical predictions; slide right for future predictions.



Related topics:

[Common Administration Portal Functions](#)

[Refresh Statistics](#)

[License Status](#)

[User Functions](#)

[Search/Sort Information in Tables](#)

[Show/ Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Filters](#)

Application Insight

Each Kubernetes application typically consists of multiple microservices. Also known as controllers, microservices are independent, modular components that work together as a single application. For example, a web shopping application might include microservices for a catalog, database, shopping cart, payment, shipping, etc.

The overall application workload can impact each of the microservices differently. Using our example, a highly advertised sale during the holiday season may drive a larger than usual number of customers to the shopping website, greatly impacting the shopping cart, payment, and shipping services as well as the catalog and database.

Understanding how individual microservices/controllers are impacted by external factors that affect the primary workload is very important in order to prevent over- or under-provisioning of resources. For example, how will a 50% increase in primary workload increase the CPU or memory usage for each microservice/controller?

Federator.ai provides insight to individual controllers and cluster nodes to predict the resource usage of each microservice and make system recommendations based on application workload and the impact to each microservice.

The *Application Insight* section includes pages for:

- Topology – Displays how microservices/controllers interact and the top five CPU and memory users.
- Application correlation – Displays the correlation between controllers of an application.
- Cluster node correlation – Displays the correlation between controllers and nodes on which an application runs of an application.
- Resource utilization – Displays the resource breakdown for each metric for controllers and nodes.
- Resource predictions – Displays resource predictions and recommendations for the primary workload for each controller and node.
- Application KPI Metrics – Displays metrics from each controller of the application, including common metrics, such as CPU, memory, and network traffic, as well as application-specific KPI metrics, if configured.

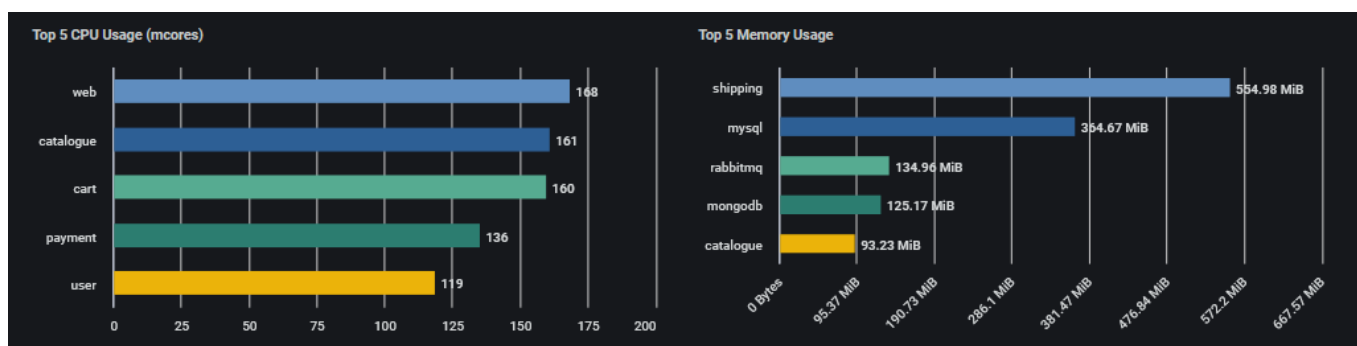
At the top of each page, you will see the controller considered the primary indicator of application load and the workload metric being used. The metric can be CPU, memory, number of network bytes received, number of network bytes transmitted, or an application-specific metric for a controller. Information for metrics from each controller is compared against the primary. The primary controller and the workload metric are set from the *Configuration - Applications* page.

Topology

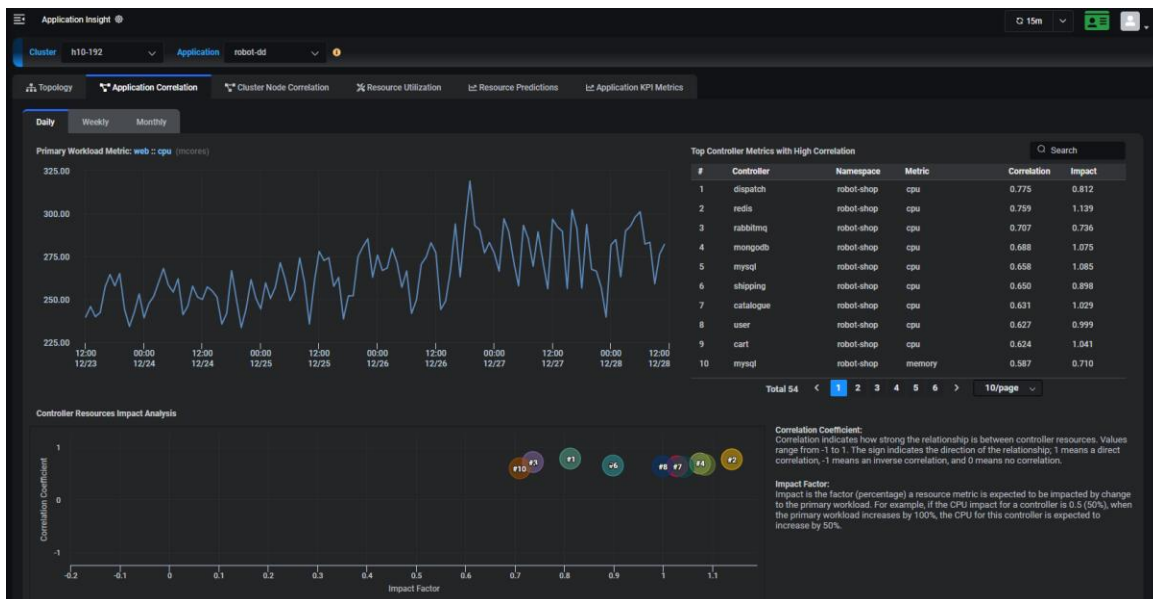


The *Topology* page shows the microservices/controllers of an application and which nodes they are running on. Here, you will see the namespace, controllers, pods, and nodes for a specific application. The number next to each pod tells you how many pods are used for that function.

The *Top 5 CPU/Memory Usage* charts display the current top five controllers for CPU and memory usage. Data is calculated every hour.



Application Correlation



The *Application Correlation* page shows the resource usage correlation between controllers and the primary workload metric of an application. There can be a large number of controller/metrics combinations to analyze because there are at least four metrics (CPU, memory, number of network bytes received, number of network bytes transmitted), and n number of controllers. Therefore, there are $4 \times n$ possible controller/metric combinations.

The top chart displays values of the primary workload metric during the specified time frame:

- Daily – Displays the last five days.
- Weekly – Displays the last 30 days.
- Monthly – Displays the last four months.

Click anywhere on the chart to see values for a specific point in time.

The *Top Controller Correlation* table displays the controller metrics in the order of the correlation to the primary workload metric. This chart identifies which resource metrics of controllers have the highest correlation to the primary workload and how much of an impact there is to these metrics when the primary workload metrics increase or decrease. For each controller in the table, you will see the associated namespace, metric, correlation, and impact. Correlation indicates how strong the relationship is between controller resource metrics or other performance metrics and the primary workload metric. It is a measure of whether resource usage or performance of a controller go up/down similarly to the primary workload metric. Values range from -1 to 1, where 0 means no correlation. The sign indicates the direction of the relationship:

- 1 = Direct correlation. If the primary controller/metric usage goes up, this controller/metric usage is predicted to go up.
- -1 = Inverse correlation. If the primary controller/metric usage goes up, this controller/metric usage is predicted to go down.

Impact is the factor (percentage) a resource metric is expected to be impacted by change to the primary workload. It measures how much the increase/decrease of the workload metric impacts the increase/decrease of resource usage or a performance metric for a controller. For example, if the CPU impact for a controller is 0.5 (50%), when the primary workload increases by 100%, the CPU for this controller is expected to increase by 50%. Note that the impact factor of a metric of a controller is only meaningful when there is a stronger correlation between the metric and the primary workload metric. The impact factor of a metric is set to 0 if the correlation coefficient between this metric and the primary workload metric is between 0.5 and -0.5.

For a quick view on the primary workload’s impact to a specific controller (e.g., *redis*), type the name of the controller in the search box. You will see the correlation coefficient and impact factor of each controller metric versus the primary workload metric.

Top Controller Metrics with High Correlation

redis

#	Controller	Namespace	Metric	Correlation	Impact
2	redis	robot-shop	cpu	0.759	1.139
12	redis	robot-shop	redis_info_latency	0.543	0.786
13	redis	robot-shop	redis_keyspace_hits	0.507	0.950
15	redis	robot-shop	redis_commands_pro...	0.507	0.950
21	redis	robot-shop	network_receive_bytes	0.398	0.000
37	redis	robot-shop	network_transmit_byt...	0.336	0.000
39	redis	robot-shop	memory	0.294	0.000
43	redis	robot-shop	redis_mem_fragment...	-0.248	0.000
48	redis	robot-shop	redis_connected_clien...	-0.088	0.000

Total 9

< 1 >

10/page

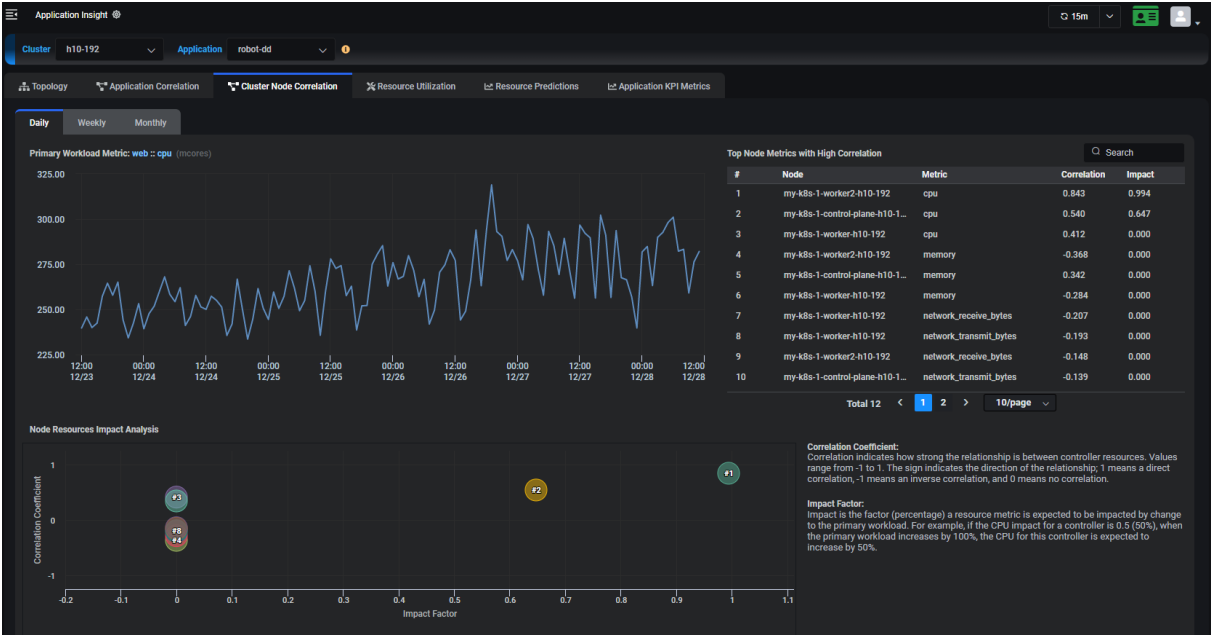
The *Controller Resources Impact Analysis* chart is a visualization of the data that appears in the *Top Controller Metrics with High Correlation* table. This chart displays the metrics shown in the current page of *Top Controller Correlation Table*. If there are more than 10 controller metrics in the table, click a new page number to view the next 10 controller metrics in the chart.

The *Controller Metric Correlation* charts provide a visualization of the top five controllers metrics with the highest correlation to the primary workload metric. The blue line is the primary workload metric while a different color is used for the controller/metric being compared. While the distance between lines may not be very pronounced, the relationship between lines (one goes up, the other goes up/down) reveals how the resource usage or performance metric of a controller is tied together. You can show additional comparisons of controller metrics with the primary workload metric by selecting a different controller metric from the drop-down menu and adding it to the list.

The range for each metric is displayed on the sides of the chart and may be different for each controller metric.



Cluster Node Correlation



The *Cluster Node Correlation* page shows the resource usage correlation between cluster node resource usage and the primary workload metric of an application. There can be a large number of node/metrics combinations to analyze because there are at least four metrics (CPU, memory, number of network bytes received, number of network bytes transmitted), and n number of nodes. Therefore, there are $4 \times n$ possible node/metric combinations.

The top chart displays actual usage for the primary controller during the specified time frame:

- Daily – Displays the last five days.
- Weekly – Displays the last 30 days.
- Monthly – Displays the last four months.

Click anywhere on the chart to see values for a specific point in time.

The *Top Node Correlation* table displays the top node resource usage metrics with the highest correlation to the primary workload metric. This chart identifies which node resources get impacted the most when primary workload metric is higher/lower. For each node in the table, you will see the associated metric, correlation, and impact.

Correlation indicates how strong the relationship is between node resources and the primary workload metric of an application. It is a measure of whether resource usage of a node goes up/down similarly to the primary workload metric. Values range from -1 to 1, where 0 means no correlation. The sign indicates the direction of the relationship:

- 1 = Direct correlation. If the primary workload metric usage goes up, this node’s resource usage metric is predicted to go up.
- -1 = Inverse correlation. If the primary workload metric usage goes up, this node’s resource usage metric is predicted to go down.

Impact is the factor (percentage) a resource metric is expected to be impacted by change to the primary workload metric. It measures how much the increase/decrease of the workload metric impacts the increase/decrease of resource usage for a node. For example, if the CPU impact for a node is 1.05 (105%), when the primary workload increases by 100%, the CPU for this node is expected to increase by 105%. Note that the impact factor of a resource metric of a node is only meaningful when there is a stronger correlation between the metric and the primary workload metric. The impact factor of a metric is set to 0 if the correlation coefficient between this metric and the primary workload metric is between 0.5 and -0.5.

For a quick view on the primary workload’s impact to a specific cluster node, type the name of the cluster node in the search box. You will see the correlation coefficient and impact factor of each cluster node resource metric versus the primary workload metric.

Top Node Metrics with High Correlation

worker2

#	Node	Metric	Correlation	Impact
1	my-k8s-1-worker2-h10-192	cpu	0.843	0.994
4	my-k8s-1-worker2-h10-192	memory	-0.368	0.000
9	my-k8s-1-worker2-h10-192	network_receive_bytes	-0.148	0.000
12	my-k8s-1-worker2-h10-192	network_transmit_bytes	0.062	0.000

Total 4 < 1 > 10/page

The *Node Resources Impact Analysis* chart is a visualization of the data that appears in the *Top Node Correlation* table.

The *Node Metric Correlation* charts provide a visualization of the top five node metrics with the highest correlation to the primary workload metric of an application. The blue line is the primary workload metric while a different color is used for the node/metric being compared. While the distance between lines may not be very pronounced, the relationship between lines (one goes up, the other goes up/down) reveals how the node resource usage is tied together.

The range for each metric is displayed on the sides of the chart and may be different for each node metric.



Resource Utilization

The *Resource Utilization* page displays the application resource breakdown for each metric (CPU, memory, number of network bytes received, number of network bytes transmitted) by each controller or by each cluster node for the specified time frame:

- Daily – Displays the last five days.
- Weekly – Displays the last 30 days.
- Monthly – Displays the last four months.

Click anywhere on a chart to see values for a specific point in time.

Resource Utilization - Controllers

These tables show the breakdown of resource usage by each controller. Data can be filtered by cluster node to show the breakdown of resource usage by controllers running on that node.



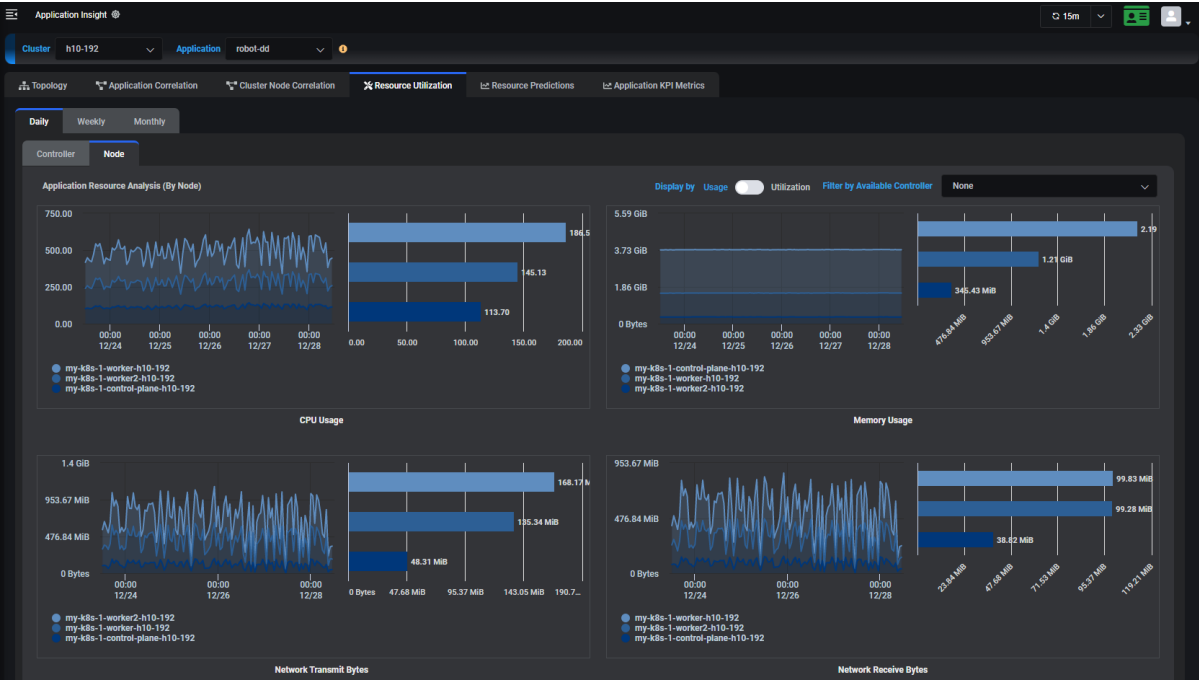
For each metric, there are two charts. The chart on the left displays the usage per controller over time. Click anywhere to see the values for all controllers at that time. You can click on the key at the bottom of the chart to show/hide individual controllers.

The donut chart on the right displays the total amount used in the last hour, six hours, or 24 hours in Daily/Weekly/Monthly view, respectively. Move your cursor over each section to show the usage (value and percentage) by a specific controller.

Resource Utilization - Nodes

These tables show the breakdown of resource usage by each node where an application runs. Data can be filtered by controller to show the resource usage of the controller on each node. You can also toggle to display data by actual usage or utilization (compared to capacity).

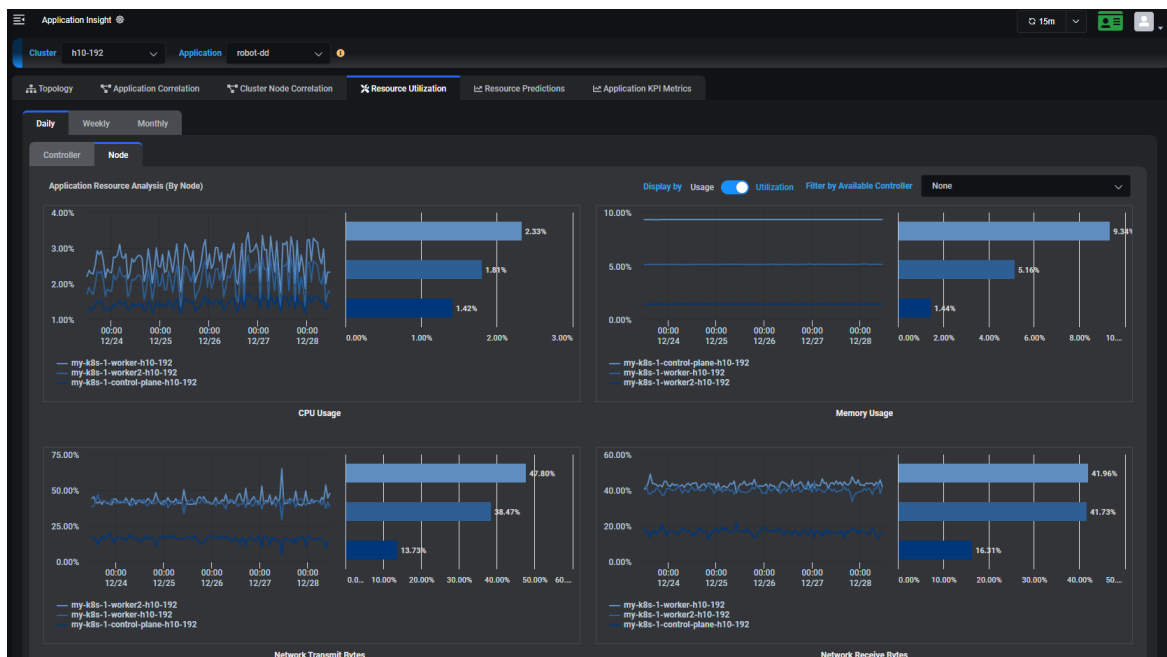
Usage



For each metric, there are two charts. The chart on the left displays the resource usage per node over time. Click anywhere to see the values for all nodes at that time. You can click on the key at the bottom of the chart to show/hide individual nodes.

The chart on the right displays the total amount used by the application on each node in the last hour, six hours, or 24 hours in Daily/Weekly/Monthly view, respectively.

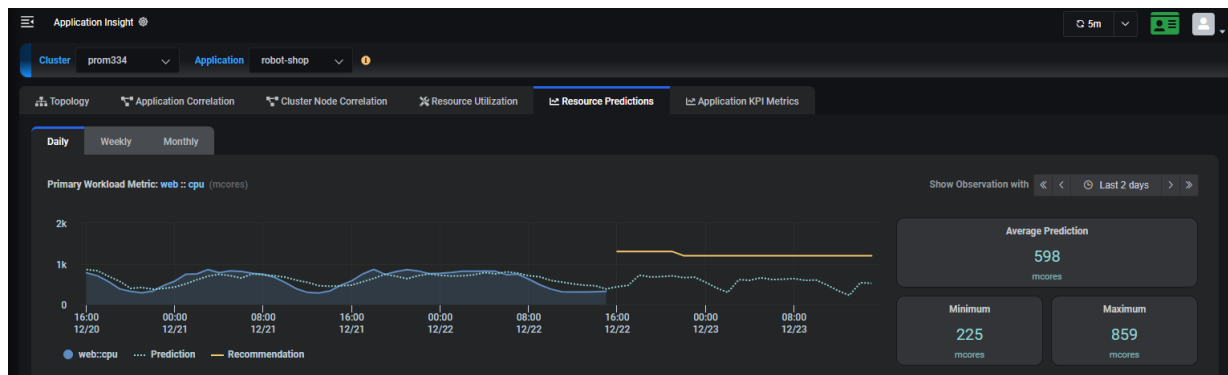
Utilization



For each metric, there are two charts. The chart on the left displays the utilization per node by the application over time. Click anywhere to see the values for all nodes at that time. You can click on the key at the bottom of the chart to show/hide individual nodes.

The chart on the right displays the percentage used by the application on each node in the last hour, six hours, or 24 hours in Daily/Weekly/Monthly view, respectively. For example, if the *Network Received Bytes* metric shows 55%, then 55% of traffic for this application went through this node.

Resource Predictions



The *Resource Predictions* page displays resource predictions and recommendations, along with average, minimum, and maximums for the primary workload metric for the specified time frame:

- **Daily** – Displays data and predictions for seven days but can be changed to display one day, two days, seven days, 14 days, 28 days, or a custom time frame.
- **Weekly** – Displays data and predictions for four weeks but can be changed to display one week, two weeks, four weeks, eight weeks, 12 weeks, or a custom time frame.
- **Monthly** – Displays data and predictions for three months but can be changed to display one month, two months, three months, six months, nine months, or a custom time frame.

Click anywhere on the chart to see values for a specific point in time.

Resource Predictions - Controllers

You can display a birds-eye view of usage, predictions, and recommendations for each metric (CPU, memory, number of network bytes received, number of network bytes transmitted) per controller for the specified time frame.

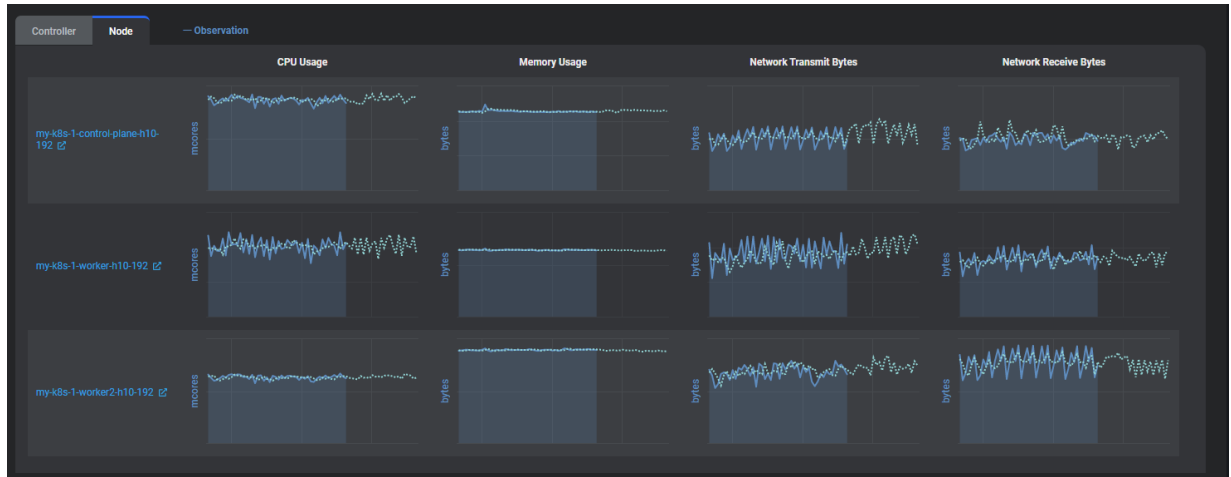
Click anywhere on the chart to see values for a specific point in time.



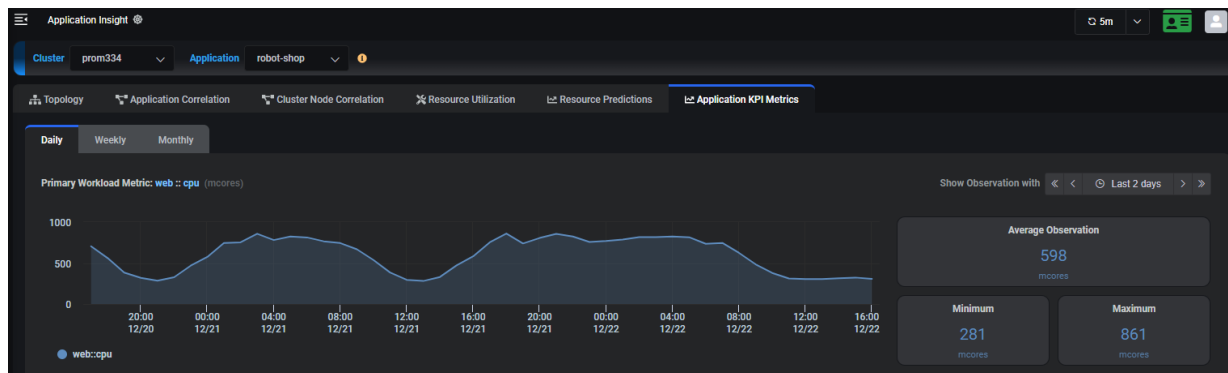
Resource Predictions - Nodes

You can display a birds-eye view of usage, predictions, and recommendations for each metric (CPU, memory, number of network bytes received, number of network bytes transmitted) per node for the specified time frame.

Click anywhere on the chart to see values for a specific point in time.



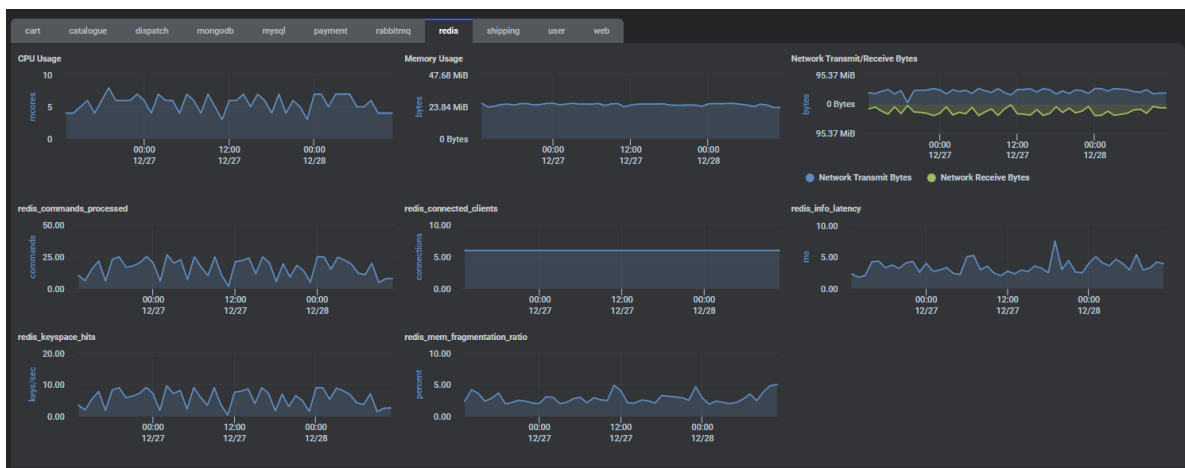
Application KPI Metrics



The *Application KPI Metrics* page displays common resource metrics from each controller of an application, as well as application-specific metrics, if configured. It also displays the average, minimum, and maximum for the primary workload metric for the specified time frame:

- Daily – Displays data and predictions for seven days but can be changed to display one day, two days, seven days, 14 days, 28 days, or a custom time frame.
- Weekly – Displays data and predictions for four weeks but can be changed to display one week, two weeks, four weeks, eight weeks, 12 weeks, or a custom time frame.
- Monthly – Displays data and predictions for three months but can be changed to display one month, two months, three months, six months, nine months, or a custom time frame.

Click anywhere on the chart to see values for a specific point in time.

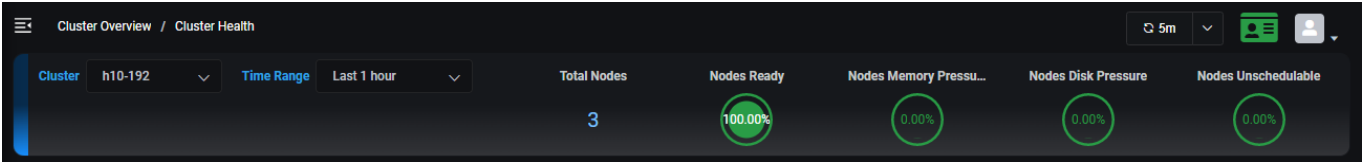


The bottom chart displays metrics for individual controllers. Click the name of the controller to see all KPI metrics for that controller.

Cluster Overview - Cluster Health

The *Cluster Health* page displays actual usage observations about the nodes/VMs in a cluster. Usage for the last 1, 2, 4, 6, or 12 hours can be displayed and can further be filtered by selecting a range of dates.

For the selected cluster, you will see the total number of nodes/VMs and the percentage that are ready. For Kubernetes, you will also see the percentage of nodes that are under memory or disk pressure, as well as the percentage of nodes that are not schedulable. Memory pressure and disk pressure are defined by Kubernetes. Refer to <https://kubernetes.io/docs/tasks/administer-cluster/out-of-resource/#node-conditions> for more information.



The table below displays the current configuration for each cluster node, including Kubernetes role (master, worker), instance types being used at your cloud provider or operating system (OS) for your local cluster, cloud provider region, number of CPUs, memory size, storage size, and node status.

Managed VMs

Name	Role	Instance Type	Region	vCPU	Memory Size	Storage Size	Status
ocp4-qd7hn-master-0	master	m5.4xlarge	us-west-1a	16	62.91 GiB	115.83 GiB	Ready
ocp4-qd7hn-worker-0-wwnc2	worker	m5.4xlarge	us-west-1a	16	62.91 GiB	116.32 GiB	Ready
ocp4-qd7hn-worker-0-v9pbg	worker	m5.4xlarge	us-west-1a	16	62.91 GiB	116.32 GiB	Ready
ocp4-qd7hn-worker-0-2p4nn	worker	m5.4xlarge	us-west-1a	16	62.91 GiB	116.32 GiB	Ready

Total 4

Kubernetes cluster

Managed VMs

Name	OS Type	vCPU	Memory Size	Storage Size	Status
h4-137	CentOS 7 (64-bit)	32	64 GiB	630 GiB	Ready
h4-146	CentOS 7 (64-bit)	8	16 GiB	60 GiB	Ready
h4-145	CentOS 7 (64-bit)	8	16 GiB	60 GiB	Ready
h4-144	CentOS 7 (64-bit)	8	16 GiB	60 GiB	Ready
h4-143	CentOS 7 (64-bit)	8	16 GiB	60 GiB	Ready

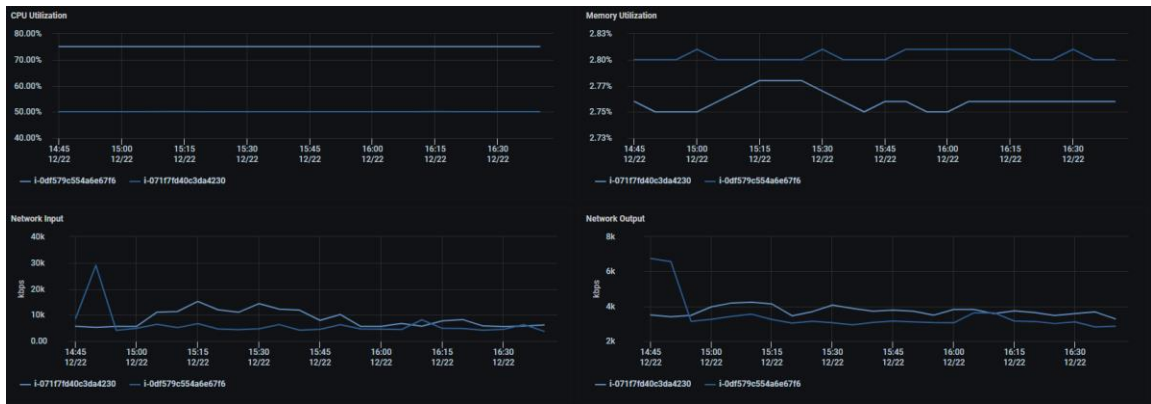
Total 9 < 1 2 > 5/page

VM cluster

CPU utilization, memory utilization, disk capacity, and disk IO utilization, network transmit and receive bytes charts are displayed for each Kubernetes node.



CPU utilization, memory utilization, and network input and output charts are displayed for each VM.



Related topics:

- [Search/Sort Information in Tables](#)
- [Show/Hide Information in Charts](#)
- [Zoom In/Out of Charts](#)
- [Terminology](#)

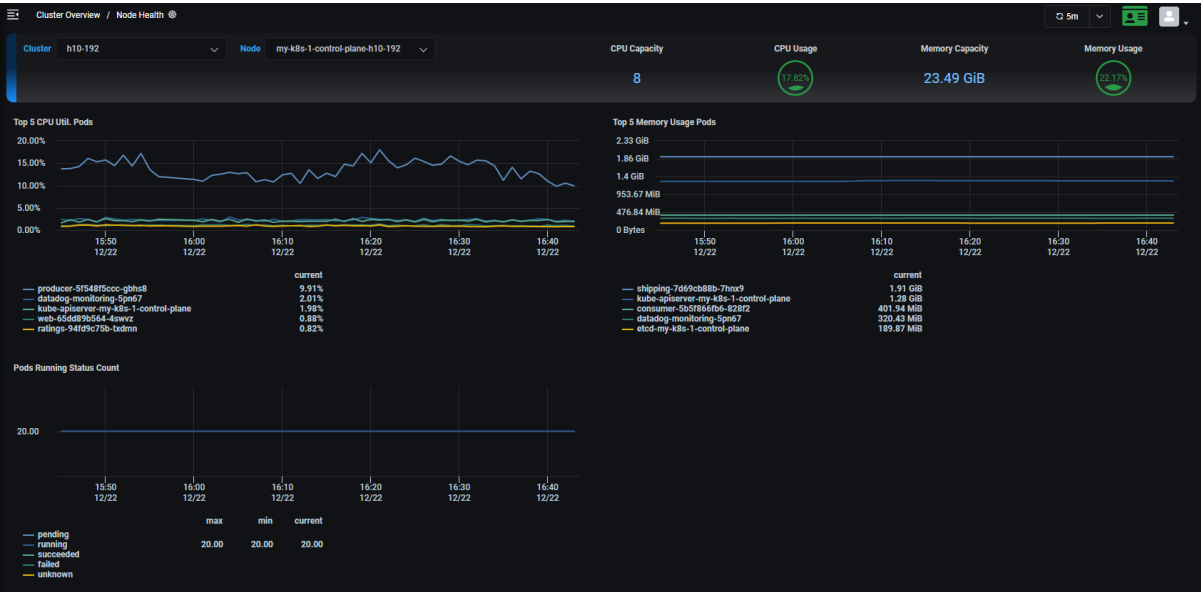
Cluster Overview - Node Health (Kubernetes)

The *Node Health* page displays actual usage observations about each node in a Kubernetes cluster. Select which cluster and node to display.

For the selected node, you will see the total CPU capacity and usage as well as memory capacity and usage.

The *Top 5* charts below display the CPU utilization of the top five pods on each node and the memory usage of the top five pods.

The *Pods Running Status Count* chart displays the number of pods running, the minimum and maximum number of pods that can run, along with the status of each pod.



Related topics:

[Search/Sort Information in Tables](#)

[Show/ Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Terminology](#)

Predictions and Planning – Kubernetes or VM Resources

The *Kubernetes Resources* page displays actual CPU and memory usage observations, predicted usage, and recommendations for resources in a Kubernetes cluster.

For Kubernetes, Federator.ai monitors resource usage for monitored clusters, nodes, namespaces, user-defined applications, and controllers and provides workload prediction, recommendations, and utilization analysis at each level.

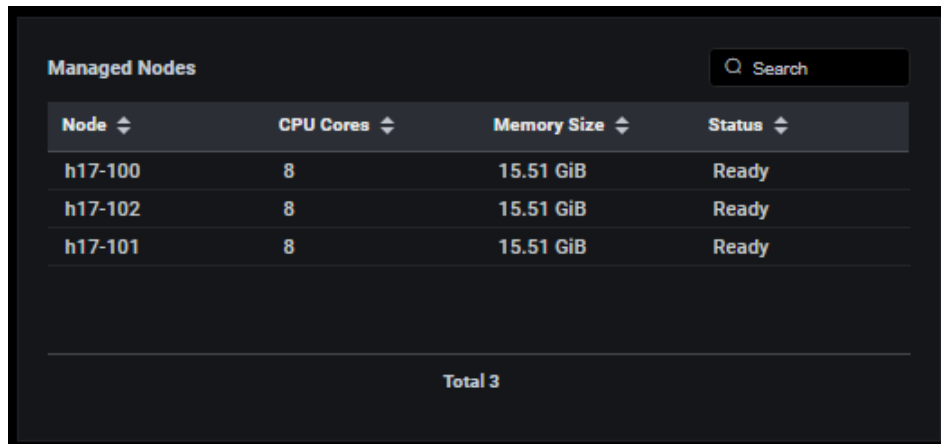
For VM, Federator.ai monitors resource usage for monitored clusters and VMs and provides workload prediction, recommendations, and utilization analysis at each level.

With the analysis and recommendations, you can decide if a resource is over-provisioned (wasting resources), or if it is under-provisioned and will not sustain an increased workload.

Use filters at the top of the page to select the level of information you want to display. When you select a Kubernetes namespace, the namespace status will be displayed.

Managed Nodes Table (Kubernetes)

This table shows the list of nodes in the selected cluster with the number of CPU cores, the size of memory, and the status for each node.



Node ↕	CPU Cores ↕	Memory Size ↕	Status ↕
h17-100	8	15.51 GiB	Ready
h17-102	8	15.51 GiB	Ready
h17-101	8	15.51 GiB	Ready
Total 3			

Managed VMs Table (VM)

This table shows the list of VMs in the selected cluster with OS type, the number of CPU cores, the size of memory, and the status for each VM.

Managed VMs				
Q Search				
Node ↕	OS Type ↕	CPU Cores ↕	Memory Size ↕	Status ↕
h4-137	CentOS 7 (64-bit)	32	64 GiB	Ready
h4-148	CentOS 7 (64-bit)	8	16 GiB	Ready
h4-147	CentOS 7 (64-bit)	8	16 GiB	Ready
h4-146	CentOS 7 (64-bit)	8	16 GiB	Ready
h4-145	CentOS 7 (64-bit)	8	16 GiB	Ready
Total 11 < 1 2 3 > 5/page				

Managed Containers Table (Kubernetes)

Based on the selected scope (cluster, node, namespace, application, or controller), this table lists the containers for the scope. Each container is listed along with its namespace, application, Kubernetes pod name, and the node where this container runs.

Managed Containers				
Q Search				
Container ↕	Namespac...	App ↕	Pod ↕	Node ↕
alameda-ai-e...	federatorai	multiple	alameda-...	h6-110
kafka	myproject	multiple	my-clust...	h6-110
kafka	myproject	multiple	my-clust...	h6-111
kafka	myproject	multiple	my-clust...	h6-112
tls-sidecar	myproject	multiple	my-clust...	h6-110
Total 7 < 1 2 > 5/page				

Workload Prediction Table and Workload Observation and Prediction Charts

The *Workload Prediction* table displays daily, weekly, and monthly predictive CPU and memory data.

The *Workload Observation and Prediction* charts display observed actual usage for the selected time as well as predictive CPU and memory data:

- Daily – Predicts CPU and memory usage every hour for the next 24 hours.
- Weekly – Predicts CPU and memory usage every 6 hours for the next 7 days.
- Monthly – Predicts CPU and memory usage every day for the next 30 days.

Use the *Time Range* field to set a custom time period for observed CPU and memory usage.

Workload Prediction Table

This table displays average/minimum/maximum CPU and memory usage and recommendations for the upcoming time selected - 24 hours (daily), 7 days (weekly), 30 days (monthly).

VM:

Daily

Weekly

Monthly

Time Range

Last 24 hours

Workload prediction for next 24 hours (from 5/20 18:00 ~5/21 18:00)

Estimated Savings: \$ 5.764(36%) /day ⓘ

Average CPU	Minimum CPU	Maximum CPU	Recommended C...	Average Memory	Minimum Memory	Maximum Memory	Recommended M...
5.09 cores	5.08 cores	5.11 cores	8.00 cores	1.59 GiB	1.58 GiB	1.59 GiB	6 GiB

Kubernetes:

Daily

Weekly

Monthly

Time Range

Last 24 hours

Workload prediction for next 24 hours (from 5/20 18:00 ~5/21 18:00)

Average CPU

19.05

cores

Minimum CPU

16.20

cores

Maximum CPU

20.94

cores

Recommended C...

24.00

cores

Average Memory

30.54 GiB

Minimum Memory

30.07 GiB

Maximum Memory

30.89 GiB

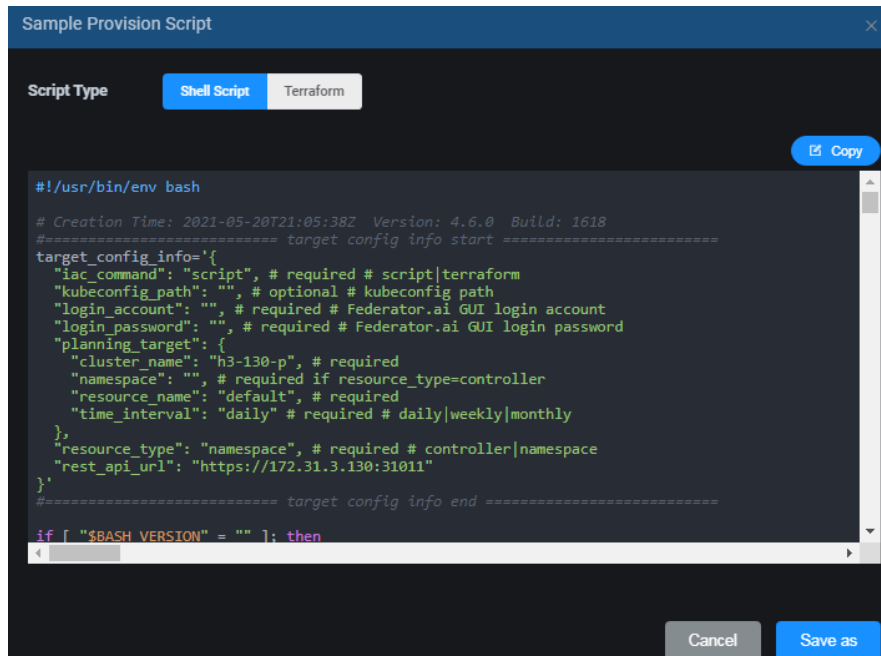
Recommended M...

40 GiB

If you are displaying information for a Kubernetes namespace and the status is anything but *Monitoring*, this section will provide more information. For example, you will see the message, "Workload prediction is not configured for this namespace" or "Not enough information for predictions" for newly added, monitored namespaces.

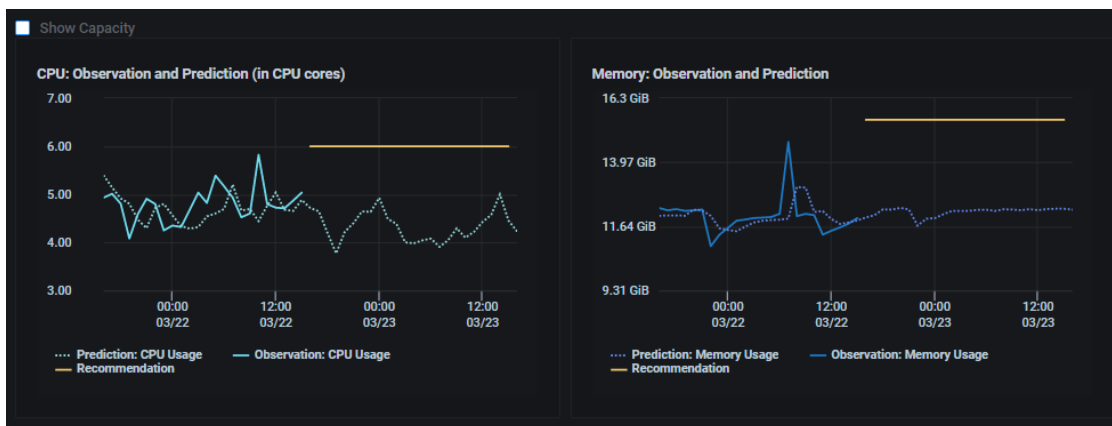
For Kubernetes namespaces and controllers, Federator.ai provides a resource provisioning script that can be used to automatically apply the recommended CPU/memory for the namespace or controller. If an auto provisioning profile is assigned to a namespace or a controller, the resource provisioning script uses recommendations set by the auto provisioning profile. Otherwise, the resource provisioning script uses system recommendations based on the time frame you are viewing (daily/weekly/monthly). For remote Kubernetes clusters, you can copy a resource provisioning script to the remote cluster in order to run auto provisioning.

Workload prediction for next 24 hours (from 5/20 18:00 ~5/21 18:00)								Resource Provision Script
Average CPU	Minimum CPU	Maximum CPU	Recommended C...	Average Memory	Minimum Memory	Maximum Memory	Recommended M...	
19.05 millicores	16.20 millicores	20.94 millicores	24.00 millicores	30.54 GiB	30.07 GiB	30.89 GiB	40 GiB	



You can copy the script and run it in the Kubernetes cluster where the controller or the namespace is located. The script queries Federator.ai for the most recent recommendations and applies them to the controller or the namespace. Refer to [Auto Provisioning Scripts](#) for more information.

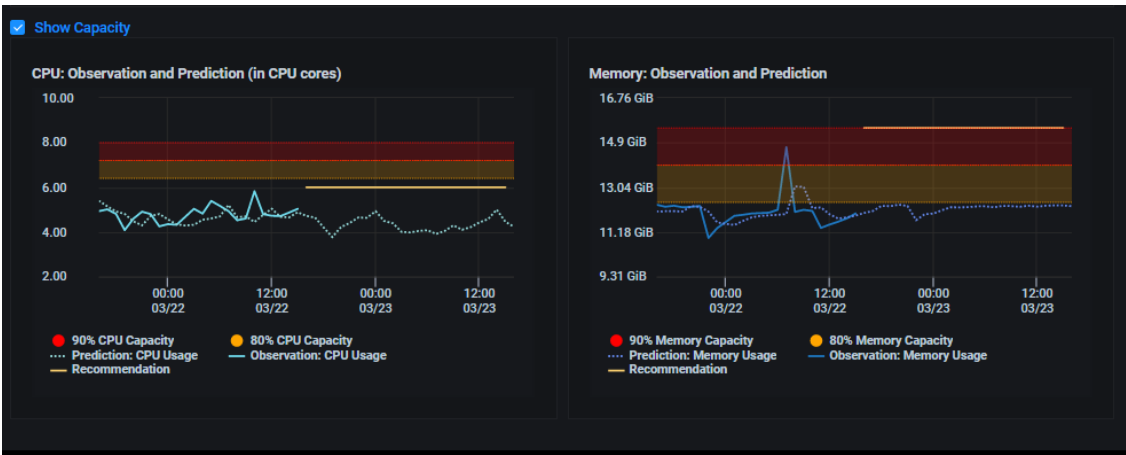
Workload Observation and Prediction Charts



These charts display CPU and memory observations and predictions for all resources specified in the *Filter* panel.

- The solid line represents the observed actual usage.
- The dotted green line represents the past and future predicted usage.
- The solid yellow line represents the recommended usage, which can help you from over-provisioning resources.
- If the selected scope is a Kubernetes node, the solid line represents the node's total CPU and memory. A big difference between total resources and actual and predicted usage can indicate that you are over-provisioned. A small difference between total resources and actual and predicted usage can indicate that you might be under-provisioned.

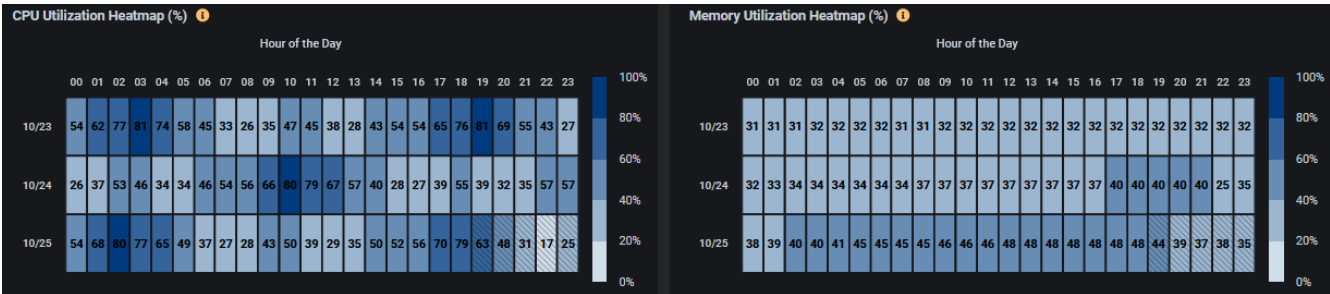
Check *Show Capacity* to see the maximum CPU and memory usage limits for the cluster, node, or VM. Orange represents 80-90% and red represents 90-100%. This is a useful way to see if the utilization of resources is approaching the overall capacity.



Utilization Analysis Charts

The *Utilization Analysis* charts display daily, weekly, and monthly CPU and memory utilization data for Kubernetes clusters, nodes, applications, and controllers and for VM clusters and VMs. Use filters to select the resources and the *Day/Week/Month* field to select a time period. For example, if *Weekly* is selected, use the *Week* field to select a different calendar week.

CPU and Memory Utilization Heatmap Charts



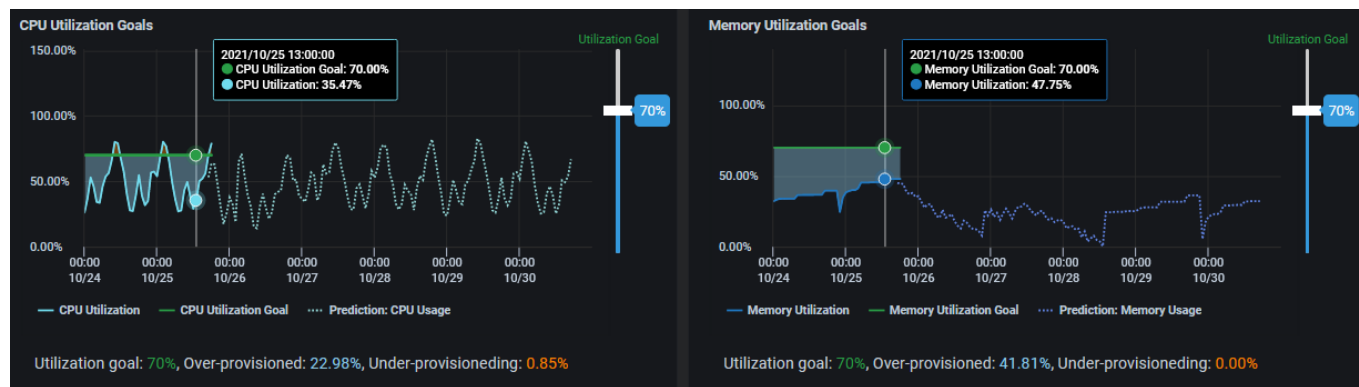
These charts display the actual and predicted CPU and memory usage percentage for all resources specified using filters for the selected time frame:

- Daily – Displays usage every hour for the last three days.
- Weekly – Displays usage each hour for a calendar week.
- Monthly – Displays usage every day for a calendar month.

For clusters and nodes, the percentage is calculated by actual usage divided by capacity. For applications and controllers, the percentage is calculated by actual usage divided by the requested (minimum) CPU/memory or the limit (maximum) CPU/memory.

The color gradient illustrates the percentage range, making it easy to see periods of high and low usage. Boxes with diagonal lines represent future predicted utilization.

CPU and Memory Utilization Goals Charts



These interactive charts display target goals along with actual and predicted CPU and memory usage for all resources specified by filters for the selected time frame.

- The blue line represents the actual CPU or memory utilization.
- The green line represents your utilization goal.
- The dotted blue line represents future predicted utilization.
- Gray areas represent periods of time when your resources were over-provisioned (wasted utilization).
- Orange areas represent periods of time when your resources were under-provisioned.

By comparing actual usage to your utilization goals, you can easily see where you are over- or under-provisioned, enabling you to adjust your cloud resources for more efficient usage. For example, if you see times when actual usage is consistently much lower than your utilization goals for a cluster, you may want to deploy additional applications in that cluster.

You can adjust your target utilization goals by using the slider on the right side of the chart.

Related topics:

[Terminology](#)

[Search/Sort Information in Tables](#)

[Show/Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Cluster Configuration](#)

[Auto Provisioning](#)

Autoscaling - HPA (Kubernetes)

The *HPA Recommendation* page displays usage and recommendation information about replicas for selected controllers. These controllers must be enabled with autoscaling during configuration. Refer to the *Configuration - Applications* section for information about how to enable autoscaling for a controller. When autoscaling is enabled, CPU and memory usage is monitored, and the number of pods is increased/decreased based on the workload. An autoscaled pod is called a *replica*.

Use the correct namespace, application, and controller name to see the history of autoscaling for a specific controller. The historical number of replicas and CPU/memory usage for the last 1, 2, 4, 6, or 12 hours can be displayed and can further be filtered by selecting a range of dates.

When the system was configured, the requested (minimum) CPU/memory and the limit (maximum) CPU/memory were set. This page shows the number of recommended and current replicas along with the total CPU/memory limit/request being used by all current replicas.

CPU and Memory Charts

These charts display CPU and memory usage and recommendations.

- The blue line represents the observed actual usage.
- The yellow line represents the limits after autoscaling.



Related topics:

[Common Administration Portal Functions](#)

[Terminology](#)

[Configuration](#)

[Search/Sort Information in Tables](#)

[Show/Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

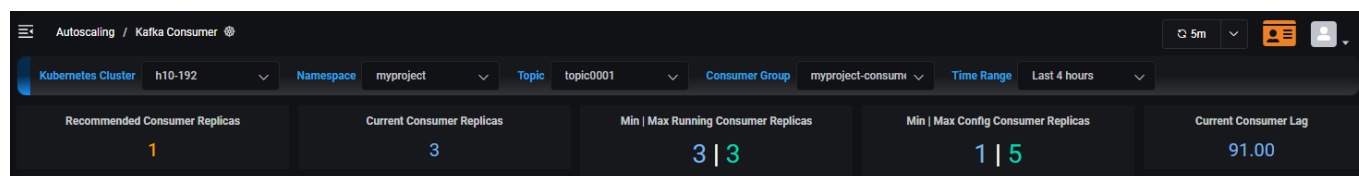
[Configure Applications](#)

Autoscaling – Kafka Consumer (Kubernetes)

The traditional method of autoscaling, based on consumer CPU/memory usage, is not enough to achieve desired performance goals for Kafka. Kafka production and consumption of messages (including message arrival rate and consumer lag) offer a better indicator of workload and performance, and is monitored by Federator.ai. Using message production rate predictions, Federator.ai autoscales the number of Kafka consumer pods to fit the workload and optimize performance.

If you have configured Federator.ai to monitor and autoscale Kafka consumers, the *Kafka Consumer* page displays predictions for the message production rate and Federator.ai scales Kafka consumer replicas to satisfy the workload. You can configure multiple Kafka topics and consumer groups. Federator.ai will predict and autoscale consumers for each individual topic/consumer group. Refer to the [Add an Application](#) section for information about how to configure Kafka consumer monitoring and autoscaling.

Use filters at the top of the page to select the correct cluster, namespace, topic, and consumer group for the Kafka consumer being monitored. The number of replicas and Kafka message production/consumption rate and the prediction of message production rate for the last 1, 2, 4, 6, or 12 hours can be displayed and can further be filtered by selecting a range of dates.



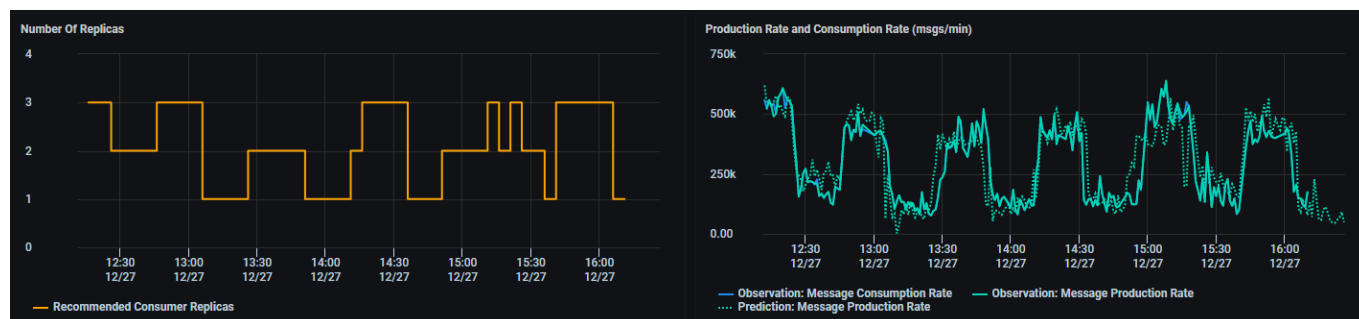
The number of recommended and current consumer replicas is displayed along with the minimum and maximum number of running and configured replicas and the current consumer lag.

Number of Replicas Chart

The *Number of Replicas* chart displays the Kafka consumer replicas for the specified time range as a result of recommendations from Federator.ai.

Production Rate and Consumption Rate Chart

The *Production Rate and Consumption Rate* chart displays the actual observed message production rate (solid green line) and consumption/processed rate (solid blue line) for the specified time range and the historical and future predicted rate (dotted green line).

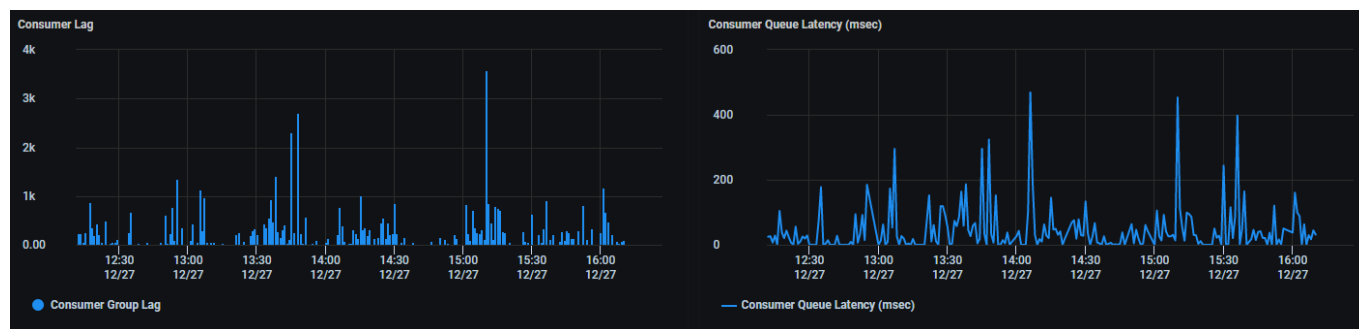


Consumer Lag Chart

The *Consumer Lag* chart displays the number of messages in the Kafka brokers that are yet to be processed by the Kafka consumers for the selected topic and consumer group and time range. This represents messages in the queue waiting to be processed.

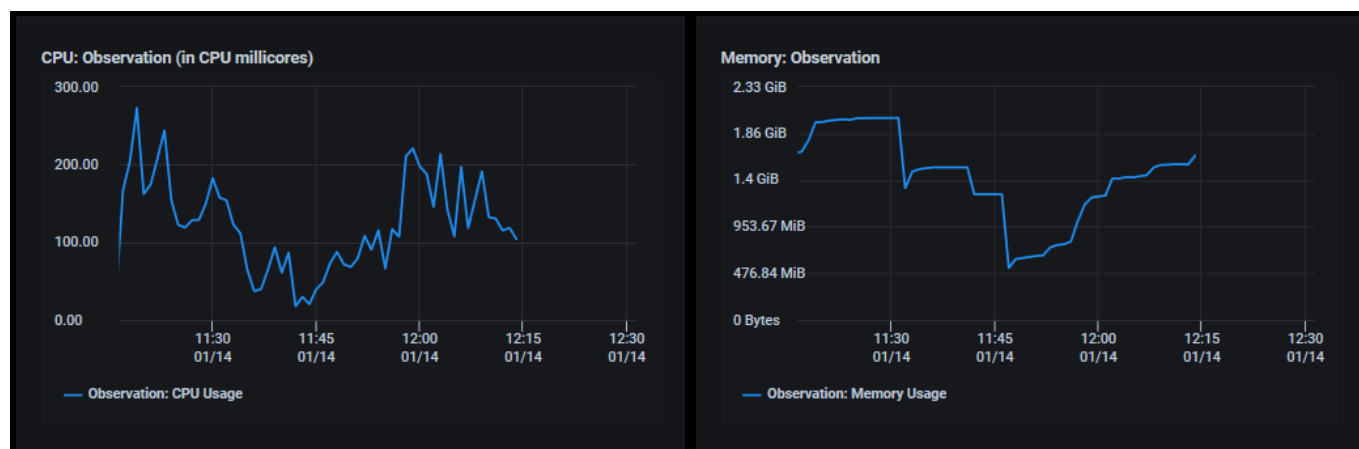
Consumer Queue Latency Chart

The *Consumer Queue Latency* chart displays how long it takes for each message to be processed for the selected time range.



CPU and Memory Observation Charts

The *CPU Observation* and *Memory Observation* charts display observed actual usage for the selected time range. This reflects the resources used as a result of autoscaling consumer pods.



Related topics:

[Common Administration Portal Functions](#)

[Terminology](#)

[Show/Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Setup Wizard](#)

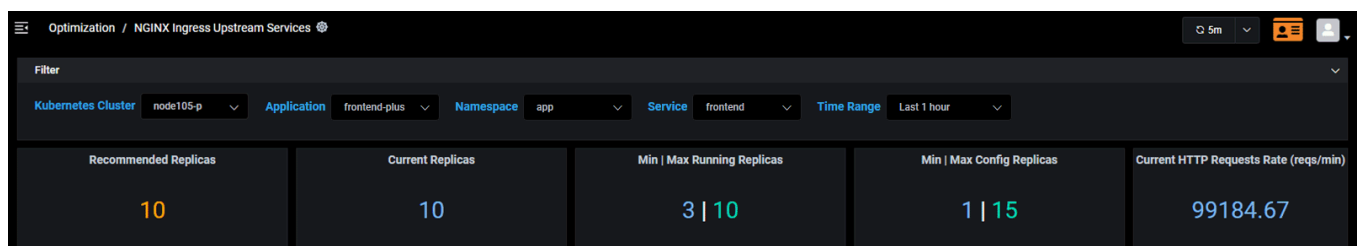
[Configure Applications](#)

Autoscaling – Ingress Upstream Services (Kubernetes)

The traditional method of autoscaling, based on CPU/memory usage, is not enough to achieve desired performance goals for Ingress services. Key performance goals include the ability of Ingress to forward requests to upstream services with minimal response time and errors; this provides a better indicator of workload and performance, and is monitored by Federator.ai. Using HTTP request rate predictions, Federator.ai autoscales the number of services to fit the workload and optimize performance.

If you have configured Federator.ai to monitor and autoscale Ingress upstream services, the *Ingress Upstream Services* page displays predictions for the HTTP request rate and Federator.ai scales replicas to satisfy the workload. You can configure multiple upstream services. Federator.ai will predict and autoscale for each individual service. Refer to the [Add an Application](#) section for information about how to configure Ingress upstream services for autoscaling.

Use filters at the top of the page to select the correct cluster, namespace, application, and service. The information and predictions for the last 1, 2, 4, 6, or 12 hours can be displayed and can further be filtered by selecting a range of dates.



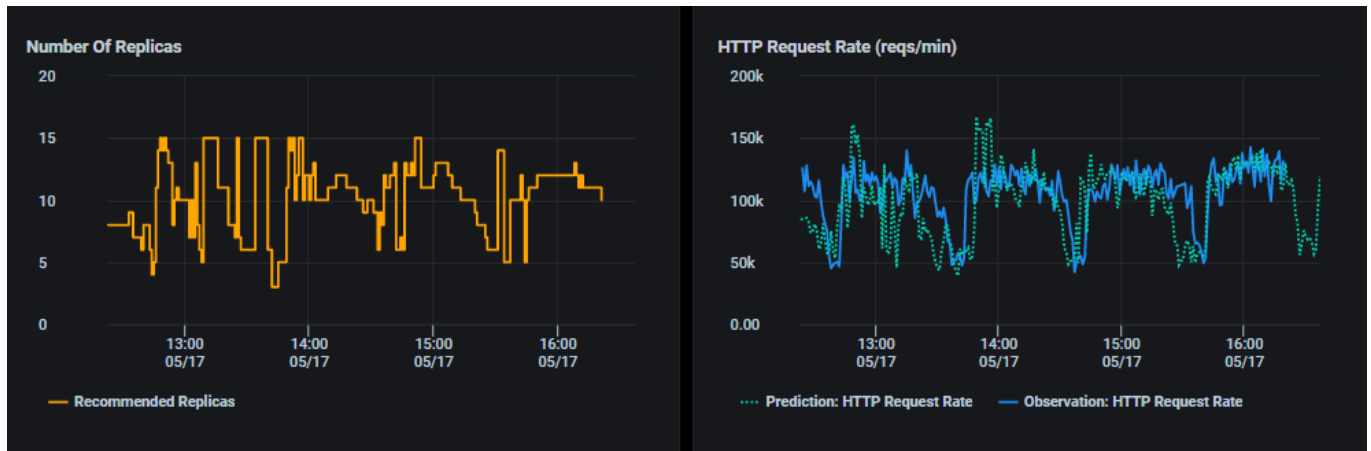
The number of recommended and current replicas is displayed along with the minimum and maximum number of running and configured replicas and the current HTTP request rate.

Number of Replicas Chart

The *Number of Replicas* chart displays the number of replicas of an upstream service for the specified time range as a result of recommendations from Federator.ai.

HTTP Request Rate Chart

The *HTTP Request Rate* chart displays the actual observed HTTP request rate (solid blue line) for the specified time range and the historical and future predicted rate (dotted green line).

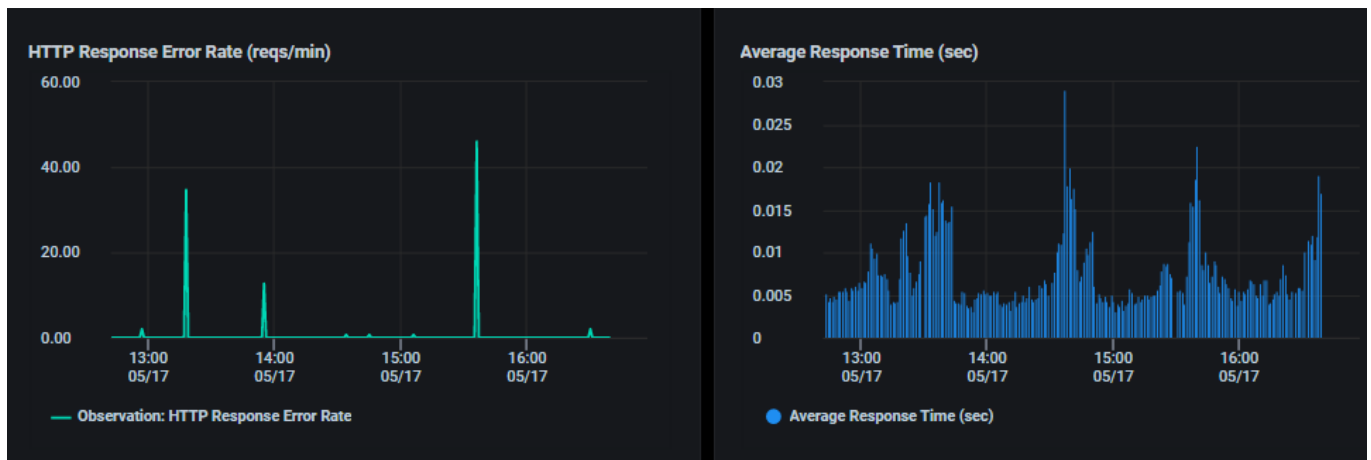


HTTP Response Error Rate Chart

The *HTTP Response Error Rate* chart displays the actual observed number of 5xx server errors for the specified time range. The goal is a zero-error rate.

Average Response Time Chart

The *Average Response Time* chart displays the average amount of time, in seconds, it takes for HTTP requests to be processed for the specified time range.



Upstream Latency Chart

The *Upstream Latency* chart displays the actual average delay for a request to upstream, in milliseconds, for the specified time range.

CPU and Memory Observation Charts

The *CPU Observation* and *Memory Observation* charts display observed actual usage of the upstream service for the selected time range. This reflects the services used as a result of autoscaling.



Related topics:

[Common Administration Portal Functions](#)

[Terminology](#)

[Show/Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Setup Wizard](#)

[Configure Applications](#)

Cost Management – Cost Analysis

The *Cost Analysis* page displays daily, weekly, and monthly costs and predictions. By default, information is displayed for all clusters but can be changed to display one or more Kubernetes or VM clusters or nodes, or one or more Kubernetes namespaces or applications.

Cost Analysis Charts



The left chart displays actual costs and predictions for the specified time frame:



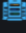

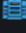

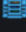

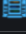

- Daily – Displays costs for each day in the last seven days but can be changed to the last 14 or 28 days, or to a custom time range.
- Weekly – Displays costs for each week in the last four weeks but can be changed to the last eight or 12 weeks, or to a custom time range
- Monthly – Displays costs for the last three months but can be changed to the last six or nine months, or to a custom time range.

Click anywhere on the chart to see values for a specific point in time. Highlight or click on the key at the bottom of the chart to show/hide individual metrics (i.e., clusters).

The right chart displays the total cost for all clusters. Highlight a section of the chart to see values for a specific cluster.

Cost Analysis Summary Table

Cluster Cost (12/20 ~ 12/26) Search

Cluster	Provider	Cluster Type	Cluster Nodes	CPU (cores)	Memory (GiB)	Cost
h10-192			3	24	70.46	\$397.33
h11-180			3	24	62.29	\$42.81
h11-180-p			3	24	62.29	\$355.74
h17-100-p			3	24	93.78	\$414.75
h17-100-s			3	24	93.78	\$414.75

Total 8 < 1 2 > 5/page

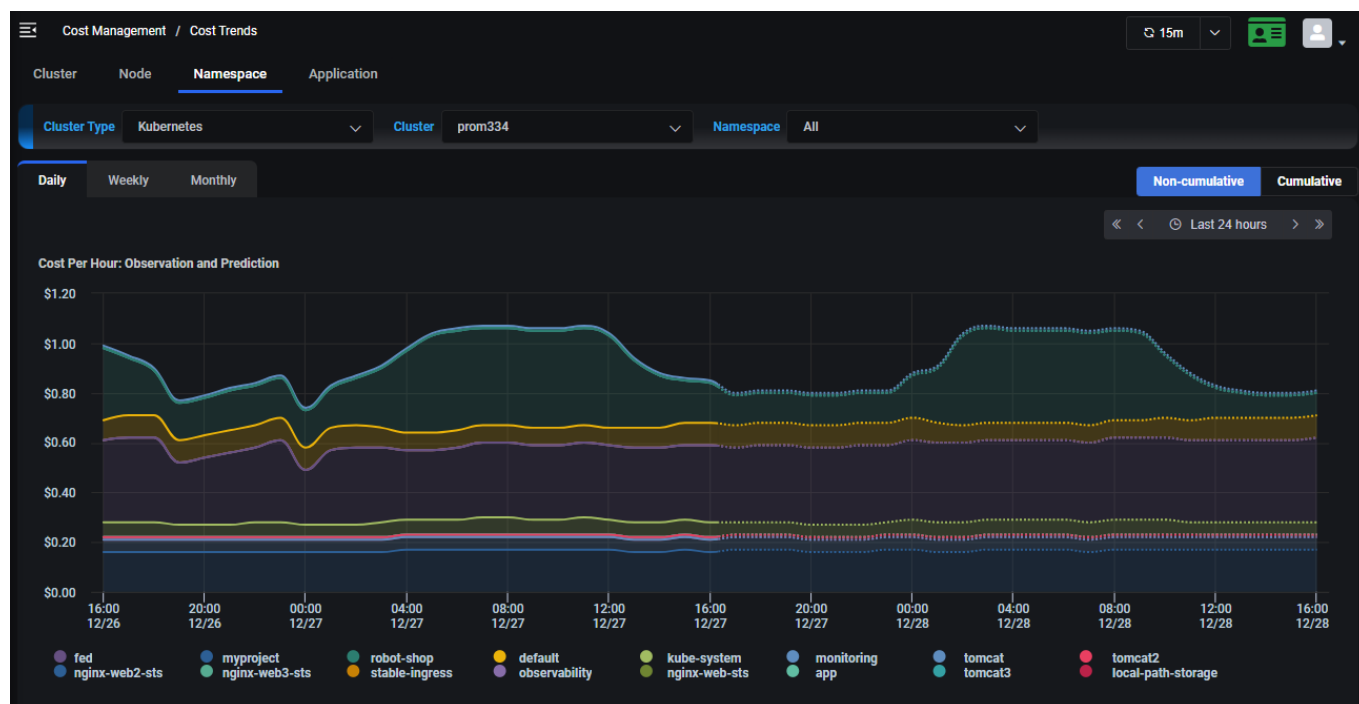
Summary information includes:

- Cluster – Cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), number of cluster nodes, CPU cores and memory capacity allocated for the cluster during the specified time frame, cost for each cluster for the specified time frame.
- Node – Node name, cluster name, provider, cluster type, instance type used at your cloud provider, CPU cores and memory capacity allocated for the node during the specified time frame, cost for each node for the specified time frame.
- Namespace – Kubernetes namespace name, cluster name, provider, cluster type, number of CPU millicores and amount of memory used during the specified time frame, cost for each namespace for the specified time frame.
- Application – Kubernetes application name, cluster name, provider, cluster type, namespace, CPU millicore and memory requests/limits, cost for each application for the specified time frame.

Cost Management – Cost Trends

The *Cost Trends* page displays daily, weekly, and monthly costs and predictions based on your expected workload. By default, information is displayed for all clusters but can be changed to display one or more Kubernetes or VM clusters or nodes, or one or more Kubernetes namespaces or applications. Note that typically, the costs for cluster capacity will be consistent, as will the costs for node capacity as they do not fluctuate much unless a cluster node is added or removed. For namespaces and applications, the cost trends are based on the actual usage of resources.

Cost Trends Chart



The chart displays actual costs and predictions for the specified time frame:





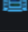



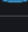
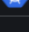
- Daily – Displays hourly costs for the last 24 hours and predictions for the next 24 hours but can be changed to display costs for the last two or four days, or to a custom time range.
- Weekly – Displays hourly costs for the last seven days and predictions for the next seven days but can be changed to display costs for the last two or four weeks, or to a custom time range.
- Monthly – Displays daily costs for the last 30 days and predictions for the next month but can be changed to display costs for the last two, three, or four months, or to a custom time range.

Click the *Cumulative* button to view trends cumulatively through the end of the selected time period.

Click anywhere on the chart to see values for a specific point in time. Click on the key at the bottom of the chart to show/hide individual metrics (i.e., namespaces).

Cost Trends Summary Table

Namespace Cost 12/26 16:00 ~ 12/27 16:00 Search

Namespace	Cluster	Provider	Cluster Type	CPU (mcores)	Memory (MiB)	Avg Daily Cost	Pred. Daily Cost
app	prom334			0	10 - 10	N/A	N/A
default	prom334			819 - 1143	1.29 GiB - 1.36 GiB	\$1.95	\$1.97
fed	prom334			2503 - 3498	5.42 GiB - 14.13 GiB	\$7.14	\$7.71
kube-node-lease	prom334			N/A	N/A	N/A	N/A
kube-public	prom334			N/A	N/A	N/A	N/A

Total 18 < 1 2 3 4 > 5/page

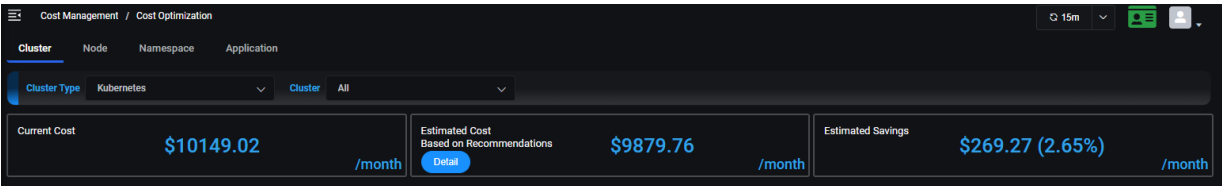
Summary information includes:

- Cluster – Cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), number of cluster nodes, CPU cores and memory capacity allocated for the cluster during the specified time frame, average and predicted cost for each cluster for the specified time frame (costs are typically the same, unless a cluster node is added or removed).
- Node – Node name, cluster name, provider, cluster type, instance type used at your cloud provider, CPU cores and memory capacity allocated for the node during the specified time frame, average and predicted cost for each node for the specified time frame (costs are typically the same, unless a node is added or removed).
- Namespace – Kubernetes namespace name, cluster name, provider, cluster type, amount CPU/memory used during the specified timeframe, average and predicted cost for each namespace for the specified time frame.
- Application – Kubernetes application name, cluster name, provider, cluster type, CPU millicore and memory requests/limits, average and predicted cost for each application for the specified time frame.

Cost Management – Cost Optimization

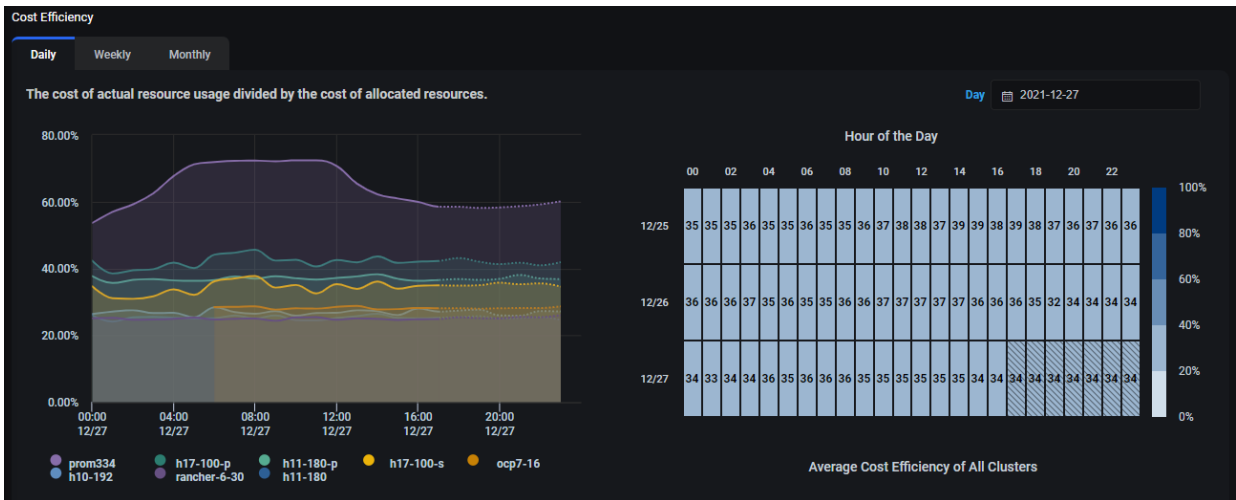
The *Cost Optimization* page displays daily, weekly, and monthly recommendations and potential savings based on your expected workload. By default, information is displayed for all clusters but can be changed to display one or more Kubernetes or VM clusters or nodes, or one or more Kubernetes namespaces or applications.

Cost Optimization Chart



This chart displays your current monthly cost with your existing resource configuration, the estimated monthly cost based on recommendations, and the potential monthly savings. Click the *Detail* link to view the recommendations (further down the page) to see how these savings are calculated. The savings typically come from reducing idle resources. The information continually refreshes itself as new data becomes available.

Cost Efficiency Charts



The *Cost Efficiency* charts display the actual cost of resource usage divided by the cost of allocated resources for the specified time frame as well as cost efficiency predictions going forward. Information is displayed in a time-series graph and heatmap for each cluster/node/namespace/application or an average for all clusters/nodes/namespaces/applications, depending upon what was selected. The lower the percentage, the more wasteful your configuration. The goal is to get to as close to 100% cost efficiency without performance risk as possible.

- **Daily** – Displays hourly cost efficiency from the start of today and predictions going forward for the day but can be changed to display cost efficiency for a different day.
- **Weekly** – Displays hourly cost efficiency for the current week and predictions going forward for the week but can be changed to display cost efficiency for a different week.

- Monthly – Displays daily cost efficiency for the current month and predictions going forward for the month but can be changed to display cost efficiency for a different month.

Click anywhere on either chart to see values for a specific point in time. Click on the key at the bottom of the graph to show/hide individual metrics (i.e., clusters).

The color gradient in the heatmap illustrates the percentage range, making it easy to see periods of high and low usage. Boxes with diagonal lines represent future predicted cost efficiency values.

Cluster Cost Optimization

Cluster Cost Efficiency Table

Cluster	Provider	Cluster Type	Cluster Nodes	CPU (cores)	Memory (GiB)	Cost/day	Cost	Cost Based on Us...	Cost Efficiency
h10-192			3	24	70	\$56.76	\$40.21	\$10.81	26.89%
h11-180			3	24	62	\$54.07	\$38.30	\$9.69	25.31%
h11-180-p			3	24	62	\$50.82	\$36.00	\$13.31	36.98%
h17-100-p			3	24	94	\$59.25	\$41.97	\$17.62	41.99%
h17-100-s			3	24	94	\$59.25	\$41.97	\$14.35	34.20%

Total 8 < 1 2 > 5/page

Cluster Cost Efficiency information includes cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), number of cluster nodes, CPU cores and memory capacity allocated for the cluster during the specified time frame, cost for each cluster for the specified time frame, cost from the actual cluster usage, and the cost efficiency percentage (the cost from actual cluster usage divided by the allocated cost) for the specified time frame.

Cluster Recommendations Table

Cluster	Provider	Cluster Type	Recomm. Cluster ...	Recomm. CPU (co...	Recomm. Memor...	Est. Cost/day	Est. Savings/day...	Est. Cost Efficienc...	
h10-192			2	16	55	\$37.90	\$18.86 (33.23%)	40.79%	View Details
h11-180			1	8	32	\$18.98	\$35.09 (64.89%)	73.80%	View Details
h11-180-p			3	18	55	\$39.94	\$10.88 (21.41%)	46.90%	View Details
h17-100-p			2	18	40	\$36.10	\$23.15 (39.07%)	70.05%	View Details
h17-100-s			3	14	32	\$28.74	\$30.51 (51.49%)	70.99%	View Details

Total 8 < 1 2 > 5/page

The *Cluster Recommendations* table shows the changes you can make to save money. For example, it may show that your cost efficiency is 60%, meaning you are only using 60% of your currently allocated resources, and based on your predicted workload, four nodes are sufficient instead of the five current ones.

The table displays cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), recommended number of cluster nodes, CPU cores, and memory, as well as the

estimated cost, savings, and cost efficiency percentage (the cost from predicted cluster usage divided by the allocated cost) per day/week/month for the cluster if the recommendations are followed. Click the *View Details* link to view the recommendations to see how these savings are calculated.

Cluster Cost Optimization Details



The *Cost Optimization* details page for a cluster shows the current configuration vs. the recommended number of cluster nodes, CPU cores, memory, and specific CPU/memory of each instance for the specified time frame (day/week/month).

Average actual CPU and memory usage is displayed along with the cost savings if the recommendations are followed.

Current CPU and memory capacity, recommended CPU and memory capacity, along with observed and predicted CPU and memory usage are also shown for Kubernetes clusters and AWS CloudWatch VM clusters configured with AWS ASG.

Cluster Node Cost Optimization

Cluster Node Cost Efficiency Table

Cluster Node Cost Efficiency 12/26 04:00 ~ 12/27 16:00

Node Name	Cluster	Provider	Cluster Type	Instance Type	CPU (cores)	Memory (GiB)	Cost/week	Cost	Cost Based on ...	Cost Efficiency
h6-30	rancher-6-30				8	31	\$156.58	\$39.15	\$8.75	22.36%
h6-31	rancher-6-30				8	31	\$133.57	\$33.39	\$5.92	17.74%
h6-32	rancher-6-30				8	31	\$133.57	\$33.39	\$3.64	10.91%
h6-33	rancher-6-30				8	31	\$133.57	\$33.39	\$17.09	51.18%
Total 4										

Cluster Node Cost Efficiency information includes node name, cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), instance type used at your cloud provider, CPU cores and memory capacity allocated for the node during the specified time frame, cost for each node and cost based on actual usage for the specified time frame, and the cost efficiency percentage (the cost from actual node usage divided by the allocated cost) for the specified time frame.

Cluster Node Recommendations Table

Node Recommendations for Next Week

Node Name	Cluster	Provider	Cluster Type	Instance Type	Recomm. CPU	Recomm. Mem	Est. Cost/week	Est. Savings/w	Est. Cost Effi
h6-30	rancher-6-30				8	31	\$156.58	\$344.54 (61.82%)	66.43%
new_worker_nod...	rancher-6-30				4	8	\$56.18		

View Details

The *Cluster Node Recommendations* table displays cluster name, provider (AWS, Azure, Google, or on-premises), cluster type (Kubernetes, VM), instance type, recommended CPU cores and memory, as well as the estimated cost, savings, and cost efficiency percentage (the cost from predicted usage divided by the allocated cost) per day/week/month if the recommendations are followed. Note that since containers are deployed to different cluster nodes by Kubernetes, there is no individual recommendation for a cluster node. Instead, the cost optimization and recommendations are applied to the entire cluster. Click the *View Details* link to view the recommendations to see how these savings are calculated.

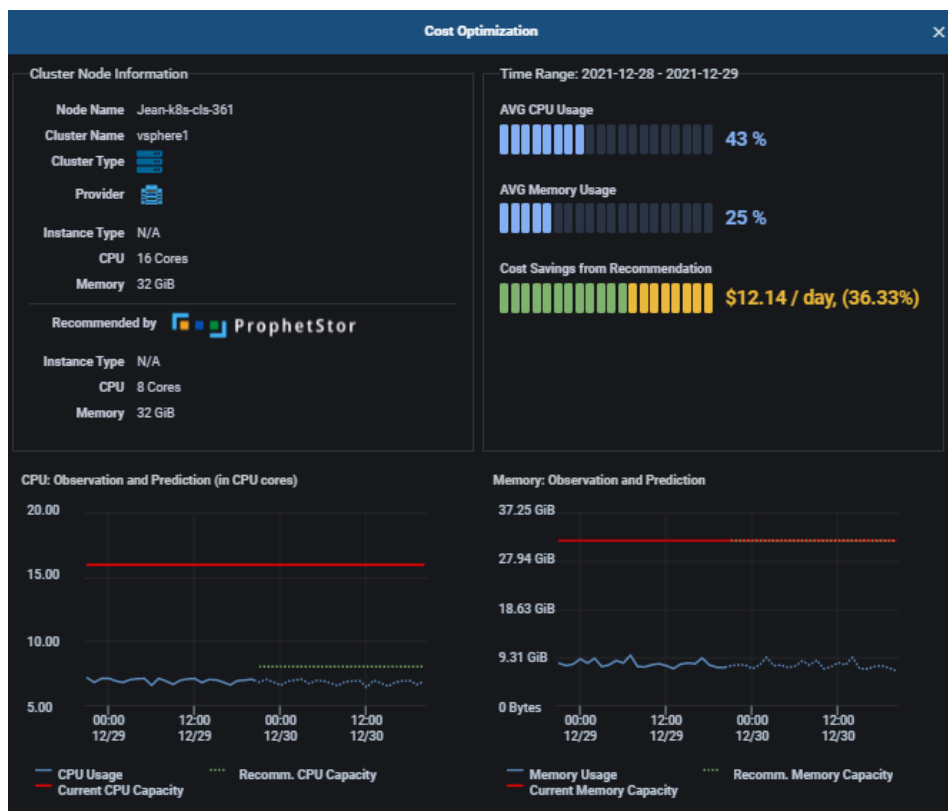
This table shows the changes you can make to save money. However, the information presented is based on the cluster type:

- Kubernetes – Since containers are deployed to different cluster nodes by Kubernetes, there is no individual recommendation of individual cluster node. Recommendations for the optimal number of nodes in Kubernetes clusters are based on the cluster, even if an individual node is selected for display. Recommendations for CPU and memory are per node. For example, the system may recommend fewer nodes for the cluster but more CPU and memory for each node.

- AWS CloudWatch VM clusters configured with AWS ASG – Similar to a Kubernetes cluster, recommendations for the optimal number of nodes are based on the cluster, even if an individual node is selected for display. Recommendations for CPU and memory are per node.
- AWS CloudWatch VM clusters with individual VMs or VMware VM cluster – Each VM node will have its own recommendations for optimal number of CPU cores and memory.

Node Name	Cluster	Provider	Cluster Type	Instance Type	Recomm. CPU	Recomm. Mem.	Est. Cost/week	Est. Savings/w...	Est. Cost Effci...	
Jean-k8s-cls-361	vsphere1				8	32	\$148.94	\$84.97	59.50%	View Details
Jean-k8s-cls-362	vsphere1				8	32	\$135.13	\$84.97	44.04%	View Details
Jean-k8s-cls-363	vsphere1				4	16	\$69.87	\$150.23	51.67%	View Details
Jean-k8s-cls-364	vsphere1				8	32	\$135.13	\$84.97	55.03%	View Details
Jean-k8s-cls-365	vsphere1				4	16	\$83.67	\$42.49	34.56%	View Details
Total 5										

Cluster Node Cost Optimization Details



The *Cost Optimization* details page shows the current configuration vs. the recommended instance type, number of CPU cores and memory for the specified time frame (day/week/month) of a specific node.

For Kubernetes and AWS CloudWatch VM clusters configured with AWS Auto Scaling groups, recommendations are displayed for the entire cluster. For AWS CloudWatch VM clusters with individual VMs or VMware VM clusters, recommendations are for each VM node.

Cost from usage, cost from recommendations, and potential savings are also displayed. Click *Cumulative* to view information cumulatively through the end of the selected time period.

Namespace Cost Optimization

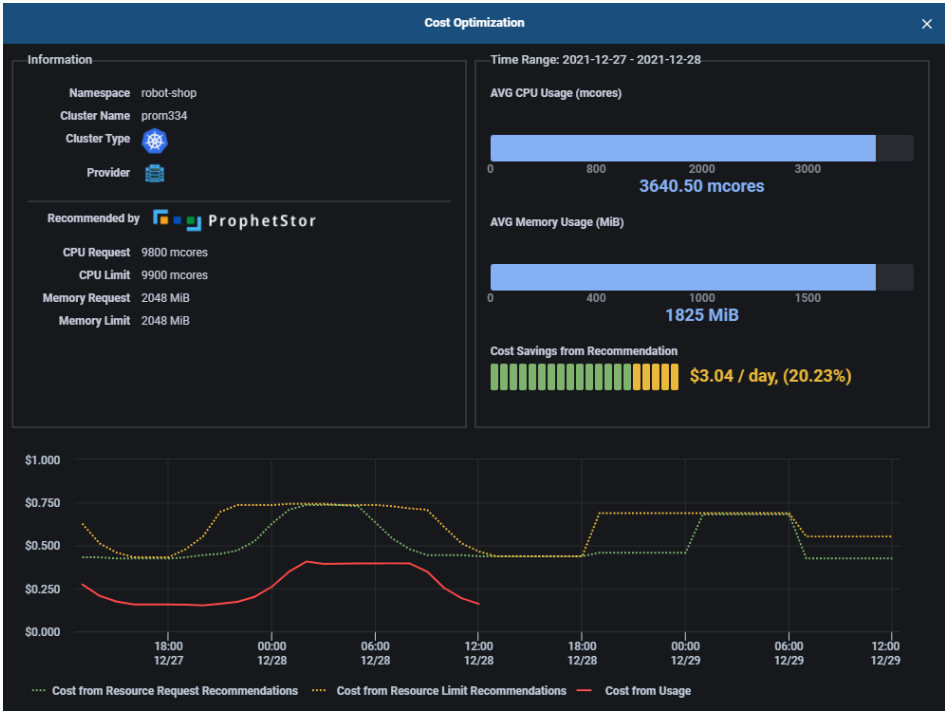
Namespace Cost Efficiency Table

Namespace	Cluster	CPU (mcores)	Memory (MiB)	CPU Request/Limit Quota(mcores)	Memory Request/Limit Quota(MiB)	Cost/day	Cost	Cost Based on Usa...	Cost Efficiency
tomcat	prom334	10 - 11	850 - 850	N/A / N/A	N/A / N/A	N/A	N/A	\$0.11	N/A
tomcat2	prom334	3 - 5	367 - 469	200 / 200	448 / 576	\$0.45	\$0.24	\$0.05	18.64%
tomcat3	prom334	1	63	100 / 500	100 / 500	\$0.92	\$0.50	\$0.01	1.45%

Namespace Cost Efficiency information includes Kubernetes namespace, cluster, number of CPU millicores and amount of memory used during the specified time frame, CPU and memory request/limit quota, as well as the daily/weekly/monthly cost based on resource quota, current up-to-date cost based on quota, cost based on actual usage, and cost efficiency percentage (the cost from actual resource usage by this namespace divided by the namespace resource quota) per day/week/month for the namespace.

Namespace Recommendation information includes Kubernetes namespace, cluster, recommended CPU and memory request/limit quota, as well as the estimated cost, estimated savings, and estimated cost efficiency percentage based on recommended resource quota per day/week/month for the namespace. Click the *View Details* link to view the recommendations to see how these savings are calculated.

Namespace Optimization Details



The *Namespace Optimization* details page shows the current configuration vs. the recommended amount of CPU and memory requests/limits quota for the specified time frame (day/week/month).

Average actual CPU and memory usage is displayed along with the cost savings if the recommendations are followed. Cost from usage and cost from request/limit recommendations are also displayed. Note that the calculation of cost efficiency and estimated savings requires namespace resource quota information. If the information is not available, the current cost efficiency and estimated cost savings will not be available.

Application Cost Optimization

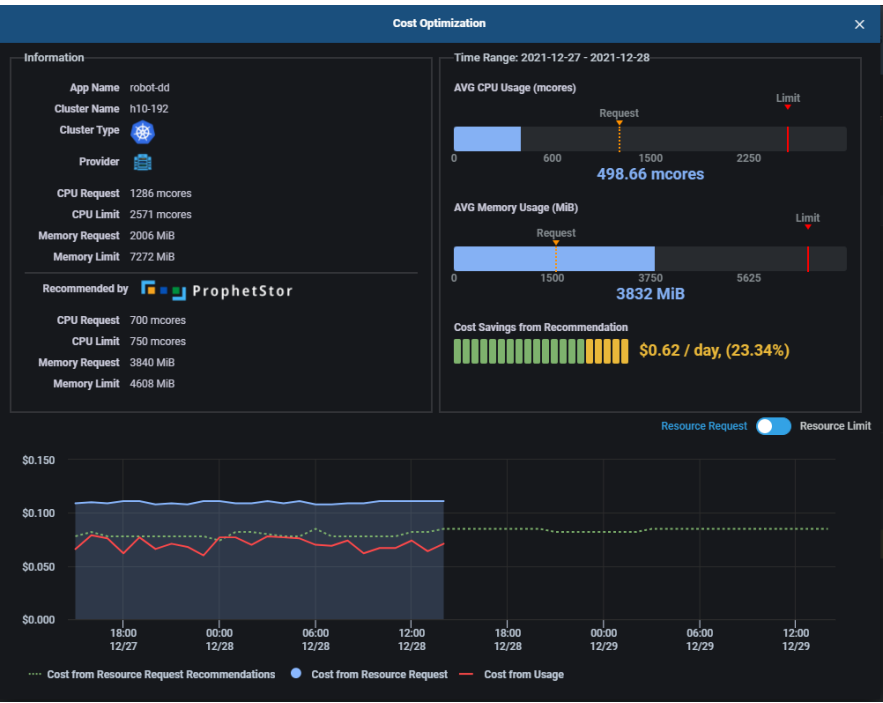
Application Cost Efficiency Table

Application Cost Efficiency 12/28 00:00 ~ 12/28 13:00									
App Name	Cluster	CPU (mcores)	Memory (MiB)	CPU Request/Limit (mcores)	Memory Request/Limit (MiB)	Cost/day	Cost	Cost Based on Usage	Cost Efficiency
consumer	rancher-6-30	52 - 77	1549 - 2068	223 / 221	5000 / 5000	\$1.43	\$0.87	\$0.28	32.14%
fed-630	rancher-6-30	2143 - 2814	10653 - 11665	2000 / N/A	1000 / N/A	\$6.52	\$3.70	\$3.70	100.00%
ingress1	rancher-6-30	82 - 88	44 - 45	625 / 943	80 / 120	\$1.05	\$0.64	\$0.09	13.52%
Total 3									

Cost Efficiency information includes Kubernetes application name, cluster, CPU and memory usage during the time period, CPU and memory requests/limits, as well as the daily/weekly/monthly cost, up-to-date cost, actual cost based on usage, and cost efficiency percentage (the cost from actual application usage divided by the allocated cost) per day/week/month.

Application Recommendations information includes Kubernetes application name, cluster, recommended CPU and memory requests/limits, as well as the estimated cost, savings, and cost efficiency percentage (the cost from actual application usage divided by the allocated cost) per day/week/month for the application if the recommendations are followed. Each application has its own recommendations. Click the *View Details* link to view the recommendations to see how these savings are calculated.

Application Optimization Details



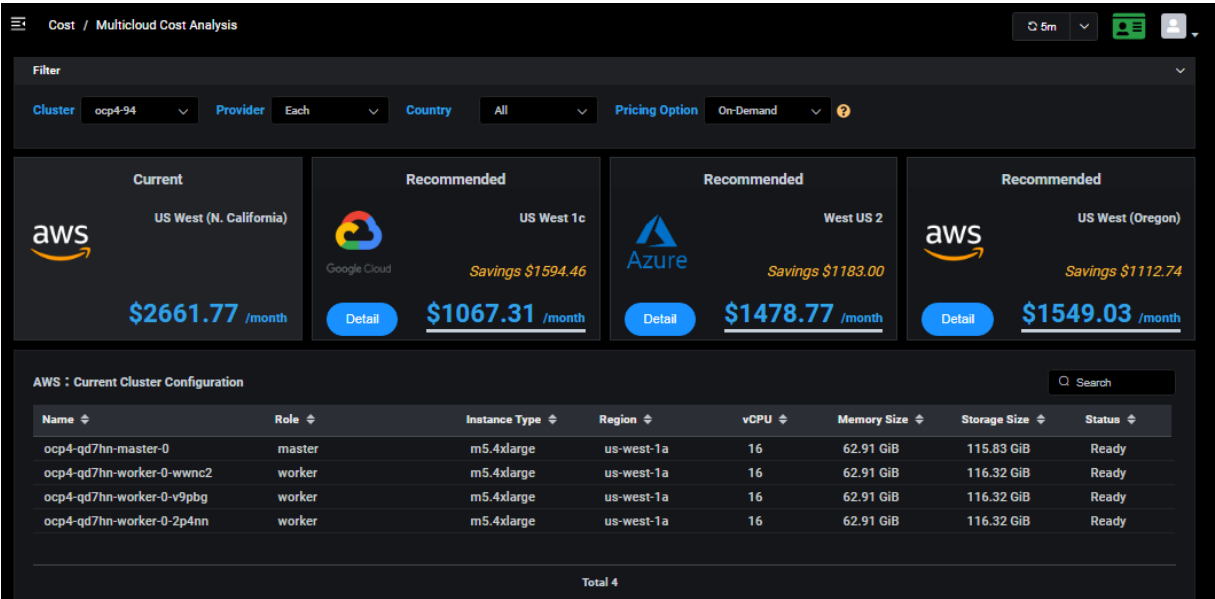
The *Application Optimization* details page shows the current configuration vs. the recommended amount of CPU and memory requests/limits for the specified time frame (day/week/month).

Average actual CPU and memory usage is displayed along with the cost savings if the recommendations are followed.

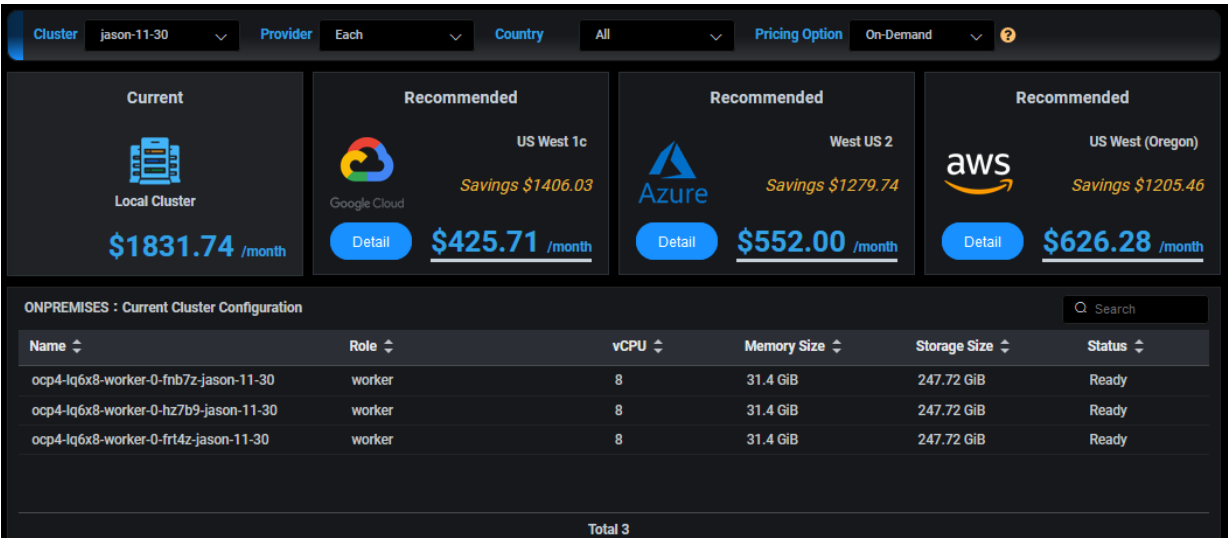
Cost from usage, cost from current resource requests/limits, and cost from recommendations are also displayed. You can switch between *Resource Request* and *Resource Limit* to see the savings comparison.

Cost – Multi-cloud Cost Analysis

The *Multi-cloud Cost Analysis* page calculates the amount you are spending with your current cloud provider (if applicable) or calculates what you are spending for your local cluster (based on your custom price book) and outlines your estimated cost for doing business with a variety of cloud providers based on fees charged by each provider and recommended changes to your current resource configuration. Federator.ai monitors and analyzes the overall CPU/memory usage of a cluster and predicts the usage for the next 30 days. Based on the past usage and the forecast, Federator.ai recommends the right instance types of cluster nodes for the workload with the lowest cost from each cloud provider.



Spending with current cloud provider



Spending for local cluster

You can filter the data that appears on the whole page by selecting:

- Cluster
- Provider – Display data for the:
 - Three lowest costs among all providers.
 - Lowest cost from each provider.
 - Three lowest costs from a selected provider.
- Country – All or a specific country.
- Pricing option - The following pricing options are available:
 - Reserved - Stationary workloads will be served with Reserved instances and temporary increased workloads will be served with On-Demand instances.
 - On-Demand - All workloads will be served with On-Demand instances.
 - Spot – (Kubernetes only) Evictable workloads will be served with Spot instances while the other workloads will be served with On-Demand instances.
 - Spot + Reserved - (Kubernetes only) Evictable workloads will be served with Spot instances while the other workloads will be served with Reserved and On-Demand instances.

The *Current* box displays what you are currently paying (locally or to your cloud provider) with your existing resource configuration. The current configuration for each cluster node, including Kubernetes role (master, worker), instance types being used at your cloud provider, cloud provider region, number of CPUs, memory size, and storage size, is displayed in the table below.

The *Recommended* boxes display potential savings with different cloud providers or different regions for your current provider. Click the *Detail* button to see how these savings are calculated. The savings typically come from lower provider fees, as well as from reducing idle resources and using more cost-effective instance types with the provider. The information continually refreshes itself as new data becomes available.

While it is difficult to move between cloud providers, this information can help you reduce costs by changing instance types, reducing idle resources, and selecting a different region from your current provider.

Even if you are not currently using a cloud provider, this information can show you what your options are if you move to the cloud.

Resource Utilization and Cost Efficiency Charts

The *Resource Utilization* charts display CPU and memory utilization information for all nodes/VMs in the selected cluster based on your daily, weekly, or monthly workload.



The solid lines represent the observed actual usage while the dotted lines show the historical and future predicted usage. The predicted usage is used in the analysis for cluster configuration recommendations.

If there are evictable Kubernetes applications in this cluster, you will also see their actual and predicted usage.

The *Cost Efficiency* chart displays the actual cost of resource usage compared to the cost of allocated resources for the selected cluster based on your daily, weekly, or monthly workload.

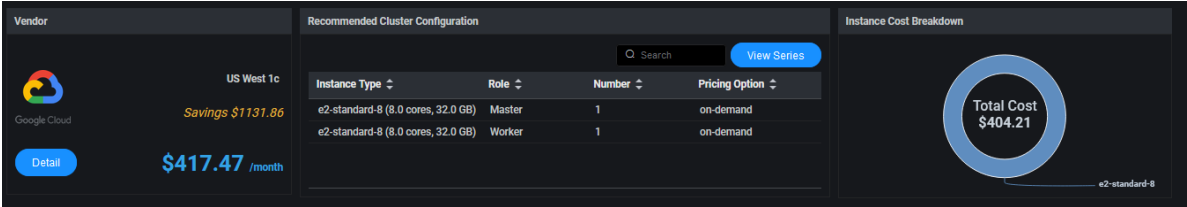


The color gradient illustrates the percentage range. Boxes with diagonal lines represent future predicted cost efficiency.

By comparing actual usage costs to the cost of allocated resources, you can easily see where you are over-provisioned, enabling you to adjust your resources for better cost efficiency. For example, if you see that the percentage is consistently lower than expected (indicating low cost efficiency), you may want to reduce size of the cluster or allocate less CPU/memory.

Recommended Cluster Configuration

The recommended configuration with each provider is displayed based on your daily, weekly, or monthly workload prediction (e.g., if *Daily* is selected for the *Resource Utilization* charts, the cost and saving are calculated based on the workload prediction for the next 24 hours). This includes the recommended instance types (scroll down in the table to see all instance types), as well as the number of Kubernetes master and worker servers for the selected pricing option.



The *Instance Cost Breakdown* chart shows the cost per instance type for the specified time frame.

Click the *View Series* button to view additional, more fine-grained, savings information. For daily predictions, data is displayed every hour for the next 24 hours. For weekly predictions, data is displayed every six hours for the next seven days. For monthly predictions, data is displayed every day for the next 31 days.



The *Cluster Potential Cost Savings* chart displays how much you are predicted to save by following the recommended configuration.

The *Recommended Cluster Configuration: Number of Instances* charts display the recommended number of on-demand, reserved, or spot instances (for Kubernetes), depending on the pricing option selected.

The recommendations for instance type include the number and type for future workloads at specific times and can help you determine how many resources to reserve. Place your cursor over one of the bars to see the recommended number of instances.

In the case of a VM cluster, the recommended cluster configuration will display the recommended instance type for each individual VM.

Related topics:

[VM Cost Analysis](#)

[Application Cost Analysis](#)

[Cost Allocation Kubernetes Namespaces](#)

[Common Administration Portal Functions](#)

[Terminology](#)

[Search/Sort Information in Tables](#)

[Show/Hide Information in Charts](#)

[Zoom In/Out of Charts](#)

[Price Books](#)

Configuration - Clusters

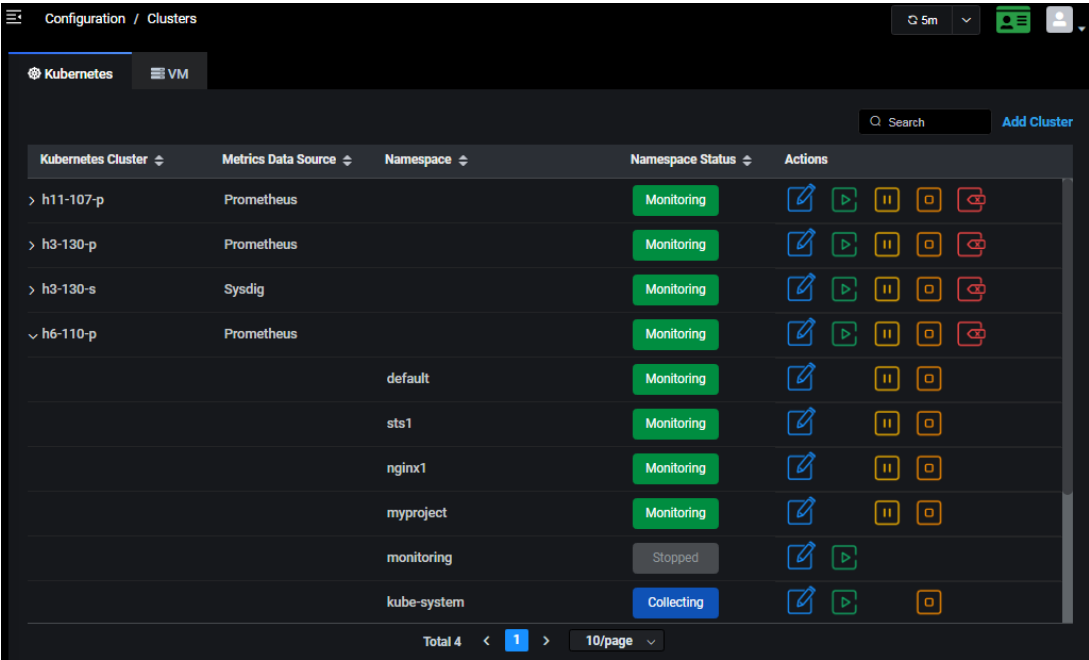
The *Clusters* page displays all Kubernetes or VM clusters being monitored by Federator.ai. You can add, manage, and remove clusters from this page.

Kubernetes Clusters




Select the *Kubernetes* tab to see the list of monitored Kubernetes clusters. Expand a cluster to see the namespaces for a cluster.



You can also see namespace status for the cluster and individual namespaces. The status will be *Monitoring* when the system is collecting metrics and providing workload predictions. The status will be *Collecting* when all the namespaces in the cluster are in collecting state. When a namespace is in collecting state, the system will still collect metrics, but prediction tasks are paused. When a namespace is added to an existing cluster, it will collect data but will not be automatically monitored until it is manually set for monitoring.

A status of *Stopped* means there is no collecting of metrics and no predictions. A cluster is *Stopped* if all of its namespaces are stopped.



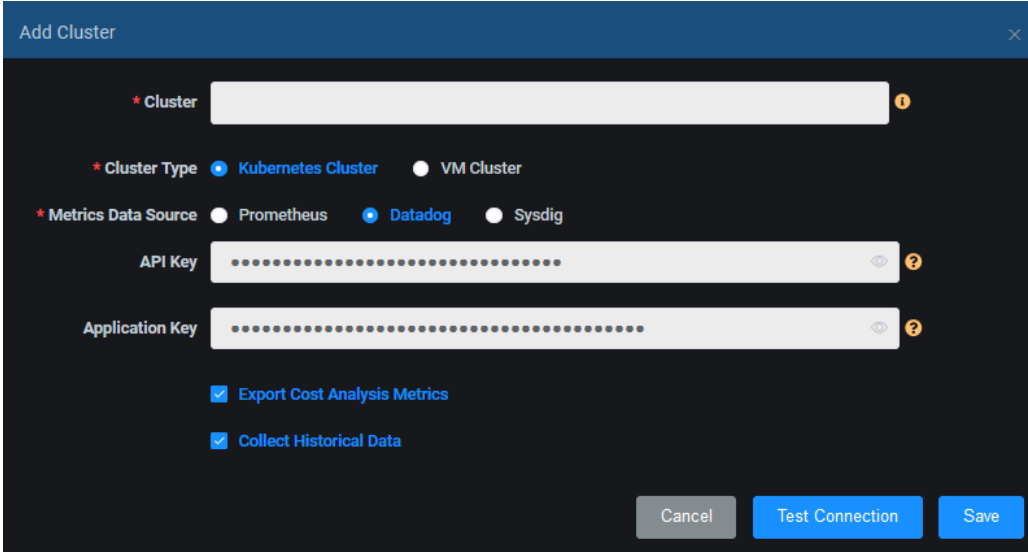
In addition to adding a cluster, you can perform the following functions for existing clusters and namespaces:

Icon	Function
	Edit settings for clusters and namespaces.
	Start monitoring and predictions for all namespaces or a specific namespace.
	Pause monitoring and predictions for all namespaces or a specific namespace.

	Stop collecting metrics and making predictions for all namespaces or a specific namespace.
	Remove a cluster.

Add a Kubernetes Cluster

1. On the *Configuration / Clusters* page, click *Add Cluster*.



Cluster – Specify the name of a cluster to be managed. There is a maximum of 253 lowercase characters, "-", or "." allowed. The name must start and end with an alphanumeric character.

Cluster Type – By default, *Kubernetes Cluster* is selected for you.

Metrics Data Source – Select the source of metrics for this cluster.

Prometheus Federation – If *Prometheus* is your data source, specify if you are using Federation, which is a group of Prometheus servers that send metrics to a centralized Prometheus server. You will need to specify the target label of the centralized Prometheus server. The format is: <label-name>:<label-value> (e.g., clusterID:host-1).

API Key/Application Key – For Datadog, the API key and application key are required for authentication. By default, they are set as Datadog API key and application key from the *Metrics Data Source* tab under *Configuration / System Settings*. Each cluster can use a different API key and application key.

URL/Token – For Sysdig, a URL and token are required for authentication. For the Prometheus open-source monitoring system, the URL is required but the token is optional. By default, they are set as Sysdig URL/token or Prometheus URL from the *Metrics Data Source* tab under *Configuration / System Settings*. Each cluster can use different values.

- Authenticate Prometheus by using basic authentication with a username and password.

Use the following command to generate the token:

```
# echo -n "<username>:<password>" | base64
```

Refer to the following for information about securing the Prometheus API using basic authentication (Basic Auth): <https://prometheus.io/docs/guides/basic-auth/>

- Authenticate Prometheus by using a service account token in OpenShift:

Use the following commands to get the service account name for Prometheus:

```
# oc get prometheus -n openshift-monitoring
NAME AGE
k8s 169d
# oc get prometheus -n openshift-monitoring k8s -oyaml | grep serviceAccount
serviceAccountName: prometheus-k8s
```

Use the following command to get the token for the Prometheus service account:

```
# oc serviceaccounts get-token prometheus-k8s -n openshift-monitoring
```

Export Cost Analysis Metrics – Specify if you want cost analysis metrics to be automatically exported to Datadog, so that the information can appear in Datadog’s user interface.

Collect Historical Data – Specify if you want the system to collect up to three months' worth of historical data for existing nodes and namespaces in this cluster. This will enable weekly and monthly predictions, recommendations, and cost analysis for newly added clusters without waiting to collect weeks’ or months’ worth of data. If less than three months’ worth of data exists, the system will collect the maximum data that is available. Typically, collection of historical data will complete in about 2-3 hours. If you need to pause collection, you can edit cluster settings.

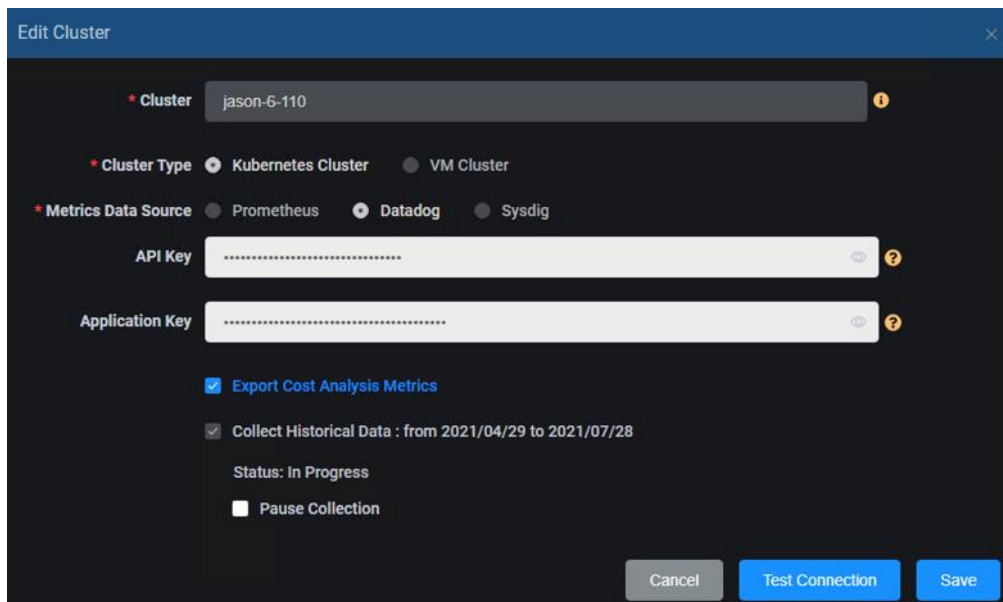
Note that Federator.ai uses the APIs provided by your metrics data source (Datadog, Sysdig) to access historical data. The data source imposes limits on how many calls can be made to their service per hour. If the rate limit is too low, the queries for historical data may exceed the limit and the API will return an error. You will need to contact your metrics data service provider to raise your API rate limit.

2. Click *Test Connection* to confirm that all information is correct.
3. Click *Save* when the system can connect to the cluster.

Manage Kubernetes Clusters

You can do the following from the *Configuration / Clusters* page:

- Edit cluster settings – Click the *Edit Cluster* icon to manage export of cost analysis metrics (Datadog), configure Prometheus Federation (Prometheus), or test the connection to the cluster. For historical data collection, you can:
 - Start the collection of up to three months' worth of historical data (from the current time) for existing nodes and namespaces in the cluster, if it was not enabled when the cluster was added.
 - See the status of data collection, including the time period collected.
 - Pause/resume historical data collection that is in progress.



- Edit namespace settings – Change monitoring status and configure auto provisioning for the namespace. To do this, click the *Edit Namespace* icon.
- Start monitoring and prediction for all namespaces in the cluster or a specific namespace. To do this, click the *Start Monitoring and Predictions* icon.
- Pause monitoring and prediction for all namespaces or a specific namespace. To do this, click the *Pause Monitoring and Predictions* icon for a cluster or for a namespace.
- Stop collecting metrics and making predictions for all namespaces or a specific namespace. To do this, click the *Stop Collecting Metrics and Predictions* icon for a cluster or for a namespace.
- Remove a cluster that does not have any applications configured. To do this, click the *Remove Cluster* icon.

Related topic:

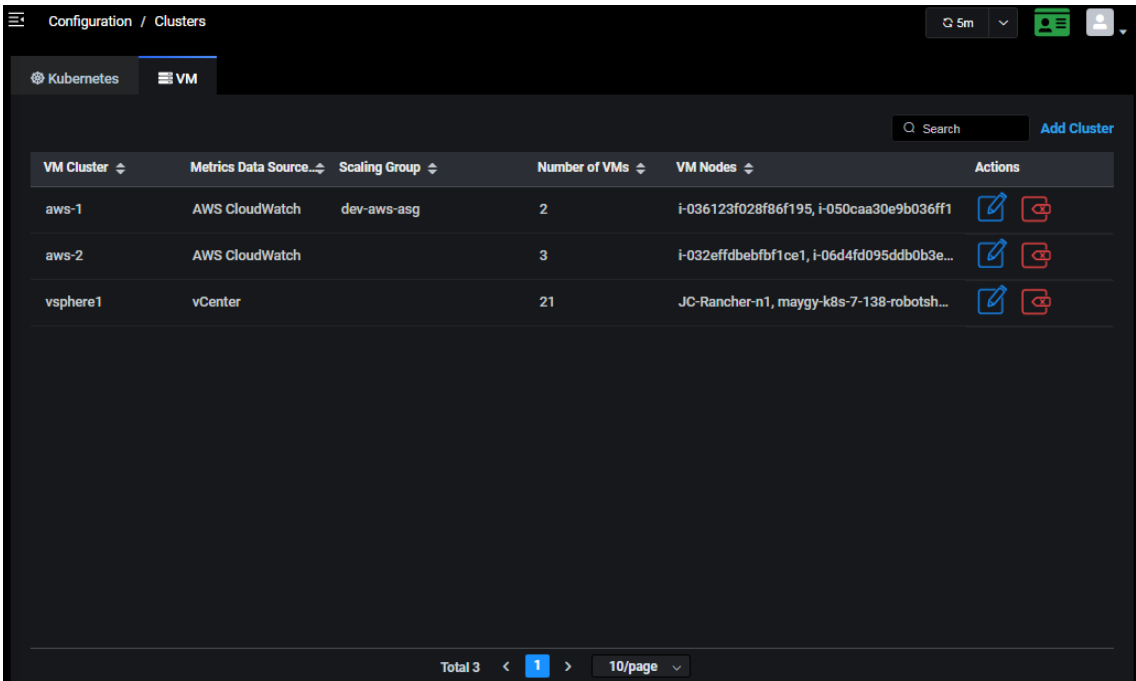
[Terminology](#)

[Search/Sort Information in Tables](#)



[Configure Applications](#)

VM Clusters

Select the *VM* tab to see the list of monitored VM clusters and the number and names of the virtual machines (VMs) in each cluster.

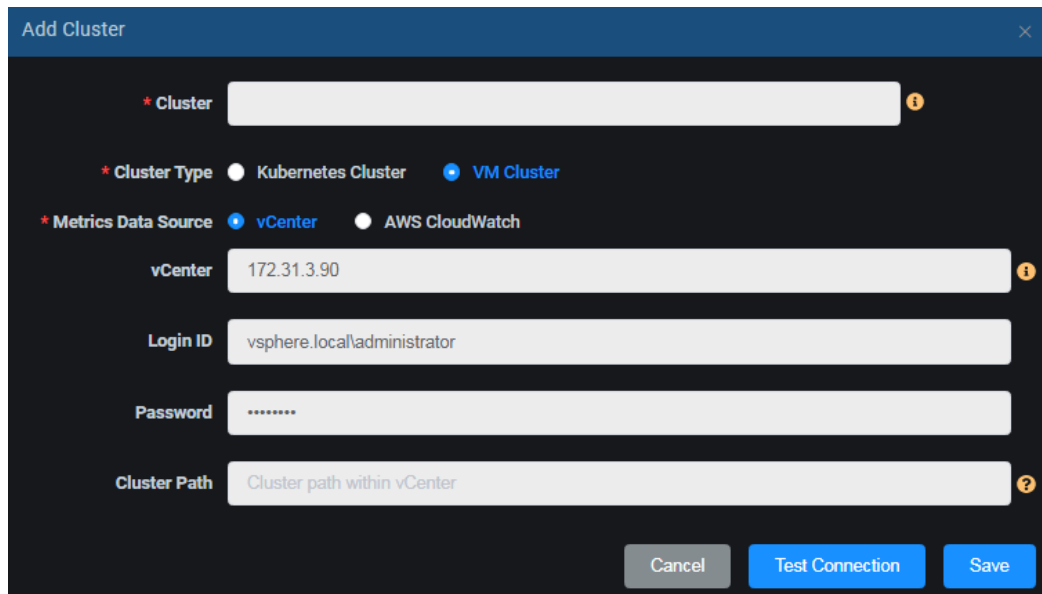


In addition to adding a cluster, you can perform the following functions for existing clusters:

Icon	Function
	Add/remove the cluster nodes to monitor.
	Remove a cluster.

Add a VM Cluster

1. On the *Configuration / Clusters* page, click *Add Cluster*.



Cluster – Specify the name of the VM cluster to be managed. There is a maximum of 253 lowercase characters, "-", or "." allowed. The name must start and end with an alphanumeric character.

Cluster Type – By default, *VM Cluster* is selected for you.

Metrics Data Source – Select the source of metrics for this cluster:

vCenter

vCenter – Specify the vCenter IP address. You can have multiple vCenters in your system.

Login ID and Password – Specify the login credentials.

Cluster Path – Specify the path to the cluster, within vCenter. If needed, you can click on the link to the vCenter website, which is included in the popup help text.

AWS CloudWatch

Note: The CloudWatch agent must be installed on the EC2 node in order to use this data source.

Region - Specify the region of Amazon AWS EC2 service.

Access Key ID - Specify the access key ID of an IAM user (16 to 128 bytes).

Secret Access Key - Specify the secret access key of the key ID that is used for access.

Collect Historical Data – Specify if you want the system to collect up to three months' worth of historical data for VMs in this cluster. This will enable weekly and monthly predictions, recommendations, and cost analysis for newly added clusters without waiting to collect weeks' or months' worth of data. If less than three months' worth of data exists, the system will collect the maximum data that is available. Typically, collection of historical data will complete in about 2-3 hours. If you need to pause collection, you can edit cluster settings.

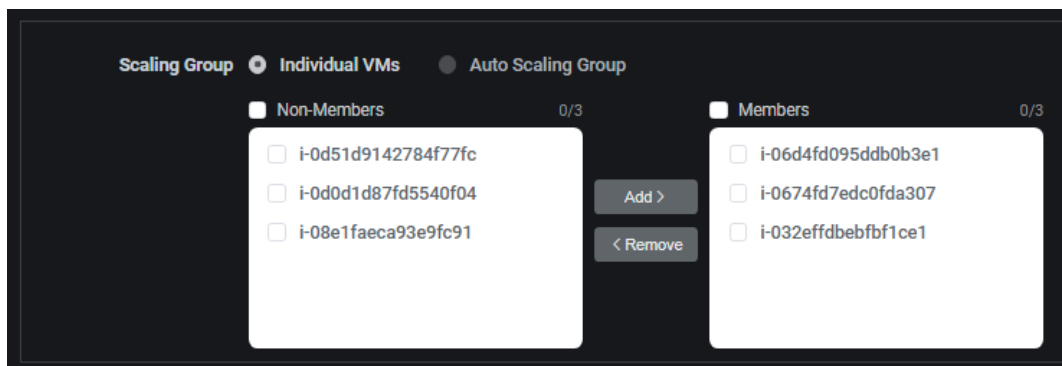
2. Click *Test Connection* to confirm that all information is correct.
3. Click *Save* when the system can connect to the cluster.

The system will discover the EC2 VMs and auto scaling groups in this region, which may take a few minutes to complete. Use the *Edit Cluster* function to select which nodes that you want to monitor.

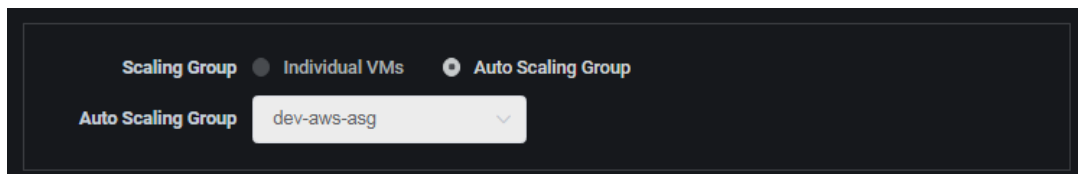
Manage VM Clusters

You can do the following from the *Configuration / Clusters* page:

- Add/remove nodes or auto scaling groups to monitor. To do this, click the *Edit Cluster* icon. Select which members of a VM cluster to monitor.



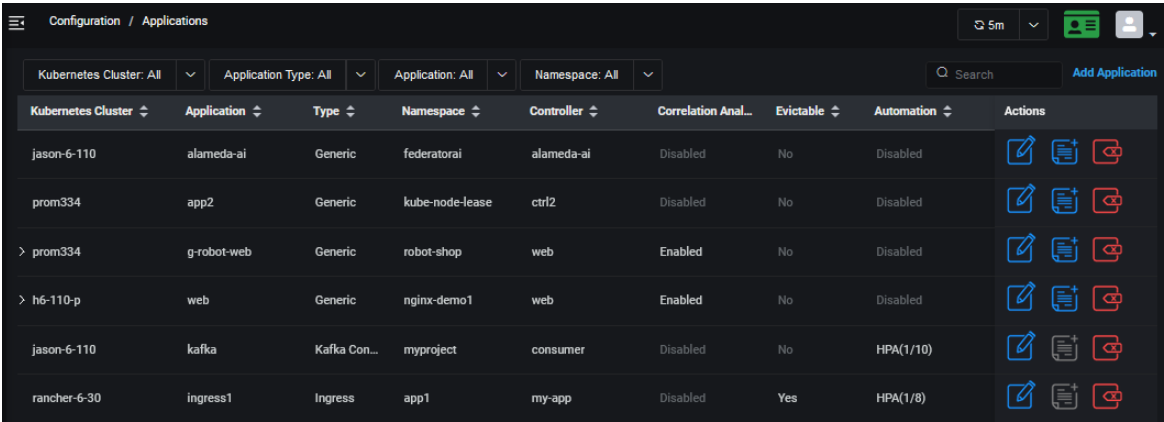
If you have Auto Scaling groups configured in AWS, select a group to monitor.






- Remove a cluster. To do this, click the *Remove Cluster* icon.

Configuration – Applications

The *Applications* page displays all Kubernetes applications being monitored by Federator.ai and allows you to manage them.



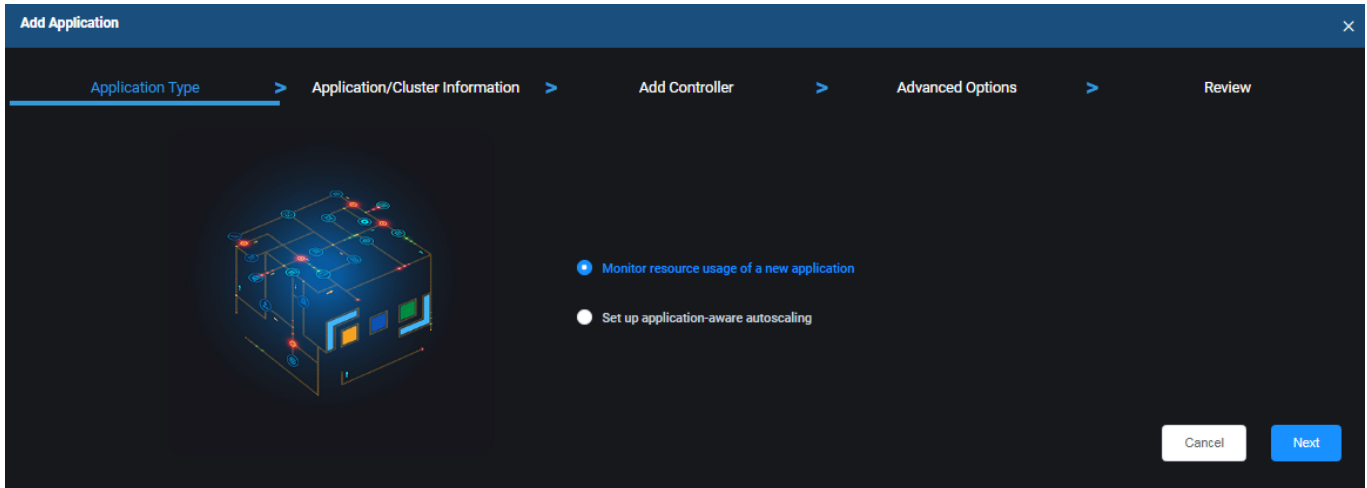
In addition to adding an application, you can perform the following functions for existing applications:

Icon	Function
	Edit an application.
	View the resource provisioning script.
	Remove an application.

Add an Application

Applications are configured via a wizard.

1. On the *Configuration / Applications* page, click *Add Application* and specify the type of application you are adding.



Depending upon the type of application you are adding, the system will be monitored and can optionally provide autoscaling.

The traditional method of autoscaling, based on consumer CPU/memory usage, is not enough to achieve desired performance goals for certain applications. For example, Kafka production and consumption of messages offer a better indicator of workload and performance. Using message production rate predictions, Federator.ai provides application-aware autoscaling of Kafka consumer pods to fit the workload and optimize performance.

As another example, Ingress services use key performance goals, including the ability of Ingress to forward requests to upstream services with minimal response time and errors, as an indicator of workload and performance. Using HTTP request rate predictions, Federator.ai provides application-aware autoscaling of Ingress services to fit the workload and optimize performance.

For other Kubernetes applications, CPU and memory usage is monitored and Federator.ai autoscales the number of pods based on workload and proactive usage predictions.

For Kafka and Ingress applications, select *Set up application-aware autoscaling*. For all other Kubernetes applications, select *Monitor resource usage of a new application*.

2. Click *Next* to continue.

Generic Kubernetes Application

1. Specify an application name and Kubernetes cluster and then click *Next* to continue.

The screenshot shows the 'Add Application' wizard with the 'Application/Cluster Information' step selected. The 'Application Name' field is empty and has a red asterisk and an information icon. The 'Kubernetes Cluster' dropdown is set to 'Select' and also has a red asterisk. At the bottom are 'Back', 'Cancel', and 'Next' buttons.

Application Name - Name of the application you want to manage. The name must be a maximum of 253 characters, contain only lowercase alphanumeric characters, “-”, or “.”, and start and end with an alphanumeric character. Note, an application is not a native Kubernetes object. You will define the controllers that are part of your application later.

Kubernetes Cluster – Select an existing cluster where this application resides.

2. Configure controllers.

The screenshot shows the 'Add Application' wizard at the 'Add Controller' step. On the left, a table titled 'Current Controllers' shows one controller: 'Deployment' in the 'Namespace' column. The 'Add Controller' button is in the top right of the table. On the right, the 'Namespace' dropdown is set to 'Select'. The 'Controller Type' is set to 'Deployment'. The 'Controller Name' dropdown is set to 'Select'. The 'Controller Workload Evictable' toggle is set to 'No'. The 'Provision Profile' dropdown is set to 'Select'. The 'Automation' section has 'No Automation' selected, with 'Auto Provisioning by Profile' and 'Horizontal Pod Autoscaling' as options. At the bottom are 'Back', 'Cancel', and 'Next' buttons.

Namespace – Select the Kubernetes namespace where the controller is deployed.

Controller Type – Supported controller types are *Deployment*, *StatefulSet*, and *DeploymentConfig* (OpenShift only).

Controller Name – Specify the name of controller to be monitored.

Controller Workload Evictable – Indicate if the controller can be interrupted if the node is shut down. Evictable controllers are good candidates to be deployed in Spot instances.

Provision Profile – Indicate if you want to select a provision profile for this application. If needed, you can even create a profile. If Federator.ai is installed in the same Kubernetes cluster as the application, this profile will be used to automatically apply resource recommendations. For remote clusters, you can click the script icon and copy the resource provisioning script for a profile to the

remote cluster in order to run auto provisioning. If you do not select a profile, you can click the script icon and copy the generic provisioning script provided by the system. When a resource provisioning script is run in a Kubernetes cluster, it queries Federator.ai for the most recent recommendations for this controller and applies the resource recommendations. Refer to [Auto Provisioning Scripts](#) for more information.

Automation: No Automation – Indicate if Federator.ai should monitor resource usage only.

Automation: Auto Provisioning by Profile - Indicate if you want to use auto provisioning based on the selected profile. This option can only be selected if Federator.ai is installed in this Kubernetes cluster.

Automation: Horizontal Pod Autoscaling – Indicate if you want to enable Horizontal Pod Autoscaling (HPA). When enabled, the number of pods is automatically increased/decreased based on the CPU/memory usage workload. HPA and Auto Provisioning are mutually exclusive; you can use HPA or auto provisioning, but not both.

Mix/Max Replicas – If *Horizontal Pod Autoscaling* is selected, specify the minimum and maximum number of pods.

3. Click *Add Controller* and repeat the previous step to add more controllers. When done, click *Next* to continue.
4. Configure advanced options. When done, click *Next* to continue.

Add Application

Application Type > Application/Cluster Information > **Add Controller (2/2)** > Advanced Options > Review

Collect Historical Data

Enable ☒ Yes ☐ No

Enable to collect application historical metrics to build prediction model.

Application Correlation Analysis

Enable ☒ Yes ☐ No

Enable to compute the correlation and causality among application metrics if they exist.

Primary Workload Metric

Controller Metric

By default, when Application Correlation Analysis is enabled, CPU, memory, network transmit bytes and network receive bytes are automatically collected. You can select additional application specific metrics if applicable.

Controller Name	Application Specific Metrics
monitoring/cart	<input type="text" value="Select"/>
monitoring/catalogue	<input type="text" value="Select"/>

Back Cancel Next

Collect Historical Data – Specify if you want the system to collect up to three months' worth of historical data for existing controllers of this application. This will enable weekly and monthly predictions, recommendations, and cost analysis for newly added applications. If less than three months' worth of data exists, the system will collect the maximum data that is available. Typically,

collection of historical data will complete in about 2-3 hours. If you need to pause collection, you can edit application settings.

Note that Federator.ai uses the APIs provided by your metrics data source (Datadog, Sysdig) to access historical data. The data source imposes limits on how many calls can be made to their service per hour. If the rate limit is too low, the queries for historical data may exceed the limit and the API will return an error. You will need to contact your metrics data service provider to raise your API rate limit.

Application Correlation Analysis – Specify if you want the system to provide statistical analysis and predictions based on the correlation between resource usage and application workload. Information will appear on the *Application Insight* pages.

Primary Workload Metrics – If you enabled *Correlation Analysis*, select which controller and which metric from this controller is the primary indicator of application load. The metric can be CPU, memory, number of network bytes received or transmitted.

Application-specific Metrics – Select additional application-specific metrics, if applicable. For example, if you are using mongoDB, mySQL, NGINX, RabbitMQ, or Redis, you can select available metrics for these applications.

5. Review all information and click *Create* to create the application.

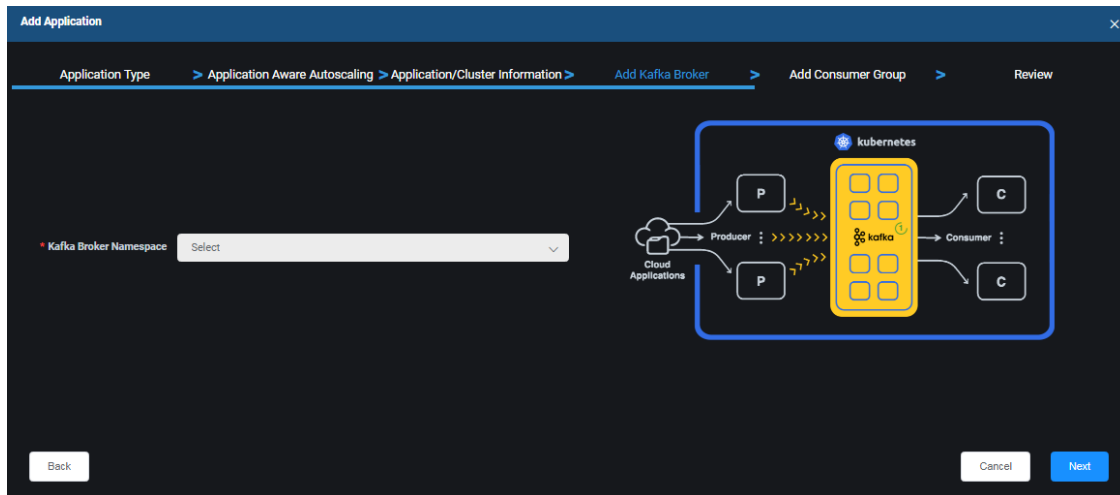
Kafka Consumer Application

1. Select *Configure Autoscaling for Kafka Consumers* and then click *Next* to continue.
2. Specify an application name and Kubernetes cluster and then click *Next* to continue.

Application Name - Name of the application you want to manage. The name must be a maximum of 253 characters, contain only lowercase alphanumeric characters, “-”, or “.”, and start and end with an alphanumeric character. Note, an application is not a native Kubernetes object. You will define the controllers that are part of your application later.

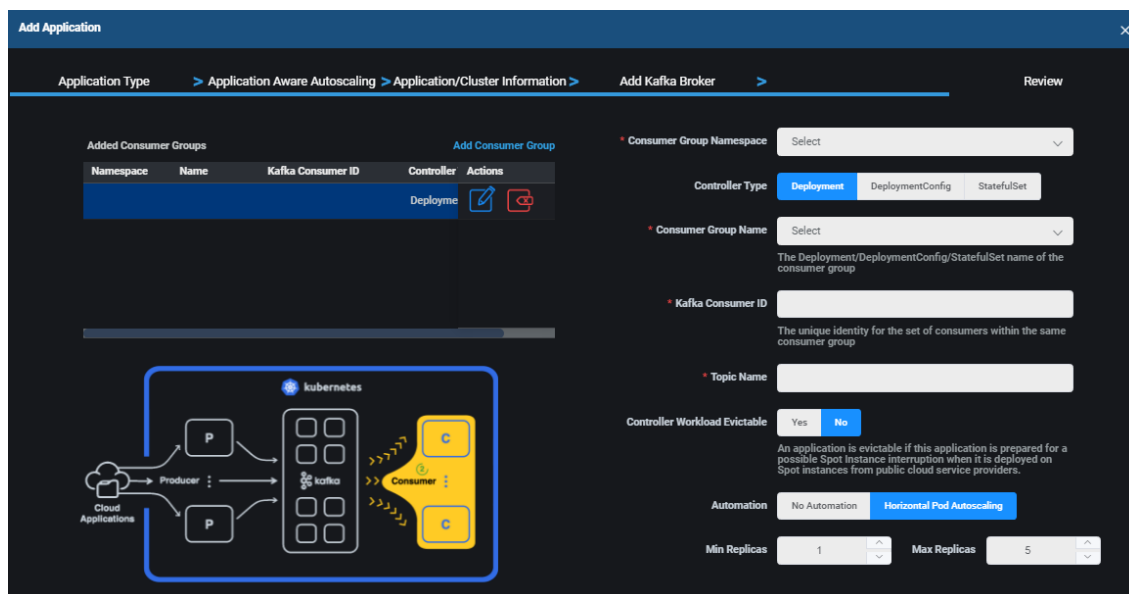
Kubernetes Cluster – Select an existing cluster where this application resides.

3. Specify the Kafka broker namespace and then click *Next* to continue.



Kafka Broker Namespace – Select the namespace for the Kafka broker that receives and stores messages from producers and allows consumers to fetch the messages.

4. Configure Kafka consumer groups.



Consumer Group Namespace – Select the namespace of the consumer group used for checking the offset on topics and processing messages.

Controller Type – Specify the controller type of the consumer group. Supported controller types are *Deployment*, *DeploymentConfig* (OpenShift only), and *StatefulSet*.

Consumer Group Name – Specify the name of the Kafka consumer group used for checking the offset on topics and processing messages.

Kafka Consumer ID – Specify the unique ID of the Kafka consumer group for the set of consumers within the same consumer group.

Topic Name – Indicate the topic (where records are stored and published) that will be processed by consumers.

Controller Workload Evictable – Indicate if the controller can be interrupted if the node is shut down. Evictable controllers are good candidates to be deployed in Spot instances.

Automation: No Automation – Indicate if Federator.ai should monitor the message production rate only and not autoscale Kafka consumer pods.

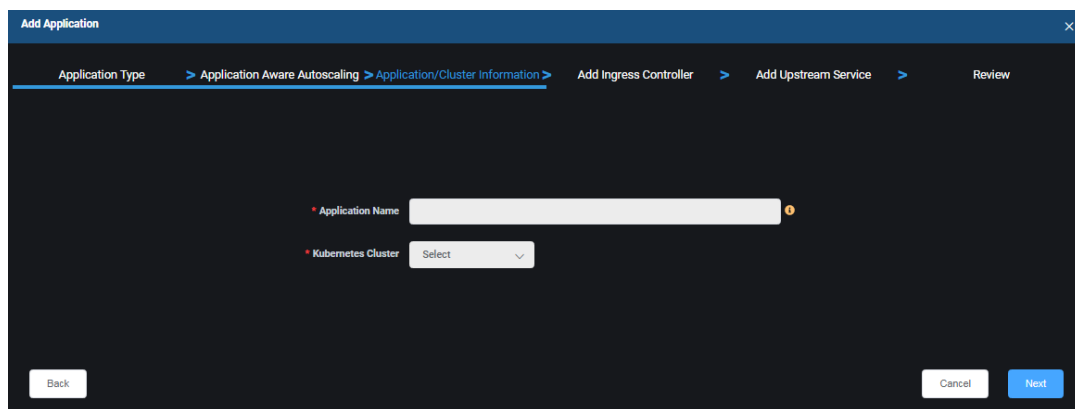
Automation: Horizontal Pod Autoscaling – Indicate if you want to enable Horizontal Pod Autoscaling (HPA). When enabled, the Kafka message production rate is monitored, and the number of Kafka consumer pods is automatically increased/decreased based on the Kafka message production rate.

Mix/Max Replicas – Specify the minimum and maximum number of consumers to be autoscaled.

5. Click *Add Consumer Group* and repeat the previous step to add more consumer groups. When done, click *Next* to continue.
6. Review all information and click *Create* to create the application.

Ingress Application

1. Select *Configure Autoscaling for Nginx Ingress Controller Upstream* and then click *Next* to continue.
2. Specify an application name and Kubernetes cluster and then click *Next* to continue.



Application Name - Name of the application you want to manage. The name must be a maximum of 253 characters, contain only lowercase alphanumeric characters, “-”, or “.”, and start and end with an alphanumeric character. Note, an application is not a native Kubernetes object. You will define the controllers that are part of your application later.

Kubernetes Cluster – Select an existing cluster where this application resides.

3. Configure the Ingress controller and then click *Next* to continue.

Ingress Controller Type – Select the controller type. Ingress NGINX is the community-supported version and NGINX Ingress Controller is a commercial version by F5. The NGINX Ingress Controller option is only available when you select a Kubernetes cluster that uses Prometheus as the metrics data source.

Ingress Controller Namespace – Select the namespace of the Ingress controller.

Ingress Controller Name – Specify the deployment name of the Ingress controller.

4. Configure upstream services.

An *upstream* service is provided by a group of upstream servers that receive HTTP requests from an Ingress controller. You must identify the Ingress controller and upstream HTTP service that you want to scale.

Upstream Namespace – Specify the namespace of the upstream service.

Upstream Service Name – Specify the name of the upstream service to be scaled.

Controller Type – Specify the Ingress controller type. Supported controller types are *Deployment*, *DeploymentConfig* (OpenShift only), and *StatefulSet*.

Upstream Controller Name – Specify the deployment name of the controller that is used to scale the upstream service.

Controller Workload Evictable – Indicate if the controller can be interrupted if the node is shut down. Evictable controllers are good candidates to be deployed in Spot instances.

Automation: No Automation – Indicate if Federator.ai should monitor the HTTP services only and not autoscale them.

Automation: Horizontal Pod Autoscaling – Indicate if you want to enable Horizontal Pod Autoscaling (HPA). When enabled, HTTP services are monitored, and the number of services is automatically increased/decreased based on the workload.

Mix/Max Replicas – Specify the minimum and maximum number of services to be autoscaled.

5. Click *Add Upstream* and repeat the previous step to add more upstream services. When done, click *Next* to continue.
6. Review all information and click *Create* to create the application.

Manage Applications

You can do the following from the *Applications* page:

- Edit application settings – Using the application wizard, you can add, edit, or remove a controller (generic application), consumer group (Kafka application), or upstream HTTP service (Ingress application) as well as enable/disable correlation analysis for a generic application. For historical data collection of generic applications, you can:
 - Start the collection of up to three months' worth of historical data (from the current time) for existing controllers of the application, if it was not enabled when the application was added.
 - See the status of data collection, including the time period collected.
 - Pause/resume historical data collection that is in progress.
- View the resource provisioning script associated with a generic controller. If no script is associated, you will see the generic system script. A script can be copied to a remote Kubernetes cluster in order to run auto provisioning for that controller. To do this, click the *Resource Provisioning Script* icon.
- Remove an application. To do this, click the *Remove Application* icon and confirm the removal.

Related topics:

[Auto Provisioning](#)

[Terminology](#)

[Search/Sort Information in Tables](#)

[Application Insight](#)

Configuration – Auto Provisioning

Federator.ai predicts CPU and memory usage for each application controller and application namespace in Kubernetes clusters and makes recommendations for the optimal amount of resources. Auto provisioning can automatically deploy resource recommendations to controllers and namespaces for generic applications based on a pre-defined profile.

An auto provisioning profile defines the conditions under which the resource recommendations will be automatically applied. It defines which recommendations to use (daily, weekly, or monthly), any adjustments to make on top of the system recommendations, and the schedule for when the resource recommendations should be applied.

If Federator.ai is installed in the same Kubernetes cluster as the application, you can assign auto provisioning profiles to controllers via the *Configuration / Applications* page or assign profiles to namespaces via the *Configuration / Clusters* page.

For remote clusters, you can copy a resource provisioning script to the remote cluster in order to run auto provisioning. Refer to [Auto Provisioning Scripts](#) below.

Note that auto provisioning and Horizontal Pod Autoscaling (HPA) are mutually exclusive; you can use HPA or auto provisioning, but not both.

The *Auto Provisioning* page displays all the existing profiles and allows you to add, edit, and remove profiles.



Auto Prov...	Recommendation	Additional Headroom		Schedule	Used By	Actions
p1	Daily	CPU: Large	Mem: Large	Hourly (4 hour at 00)	yell-db, alameda-ai, alameda-ai +3	 
p1-week	Weekly	CPU: Medium	Mem: Medium	Hourly (2 hour at 15)	federator-influxdb	 
p2	Daily	CPU: Small	Mem: Small	Hourly (1 hour at 00)	my-app	 
m1	Monthly	CPU: 500 m	Mem: Large	Hourly (1 hour at 00)	app1, testc	 

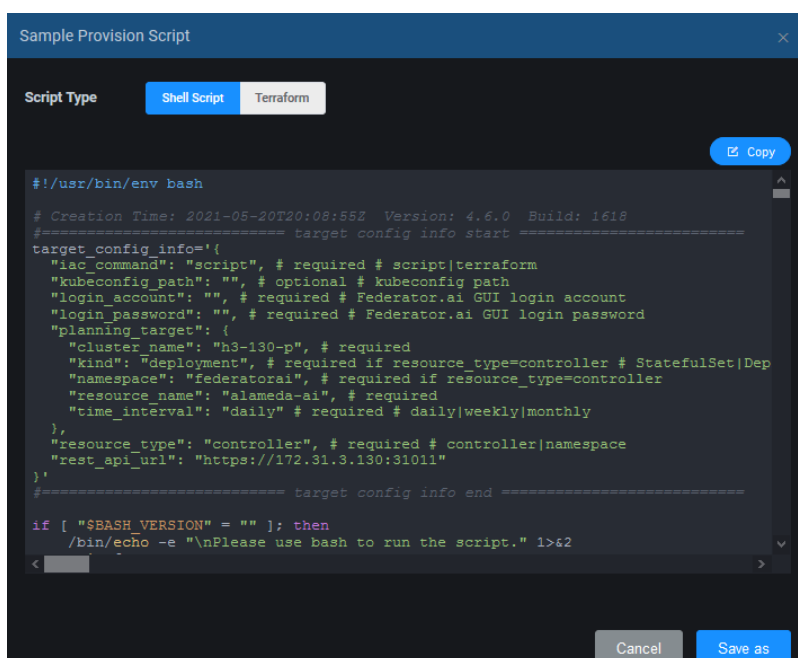
For each profile, you will see the frequency of the recommendations, CPU and memory adjustments, auto provisioning schedule, and which controllers and namespaces are using the profile. Purple represents a controller and blue represents a namespace.

Auto Provisioning Scripts

When you create an auto provisioning profile, the system generates a *resource provisioning script* that contains all the conditions in the profile.

For remote clusters, you can copy a resource provisioning script to the remote cluster in order to run auto provisioning. You can use a script associated with an auto provisioning profile or you can use the generic provisioning script provided by the system. This generic script uses system recommendations and does not have any adjustments or boundary (min/max) settings. When a resource provisioning script is run in a Kubernetes cluster, it queries Federator.ai for the most recent recommendations and applies them to a controller or a namespace.

You can find the scripts, save the scripts locally, and copy these scripts via the following pages: *Configuration / Applications*, *Configuration / Clusters* (when you edit a namespace), or *Planning / Kubernetes Workload Prediction* (when you are viewing information for controllers or namespaces).



```
#!/usr/bin/env bash

# Creation Time: 2021-05-20T20:08:55Z Version: 4.6.0 Build: 1618
##### target config info start #####
target_config_info='{
  "iac_command": "script", # required # script|terraform
  "kubeconfig_path": "", # optional # kubeconfig path
  "login_account": "", # required # Federator.ai GUI login account
  "login_password": "", # required # Federator.ai GUI login password
  "planning_target": {
    "cluster_name": "h3-130-p", # required
    "kind": "deployment", # required if resource_type=controller # StatefulSet|Dep
    "namespace": "federatorai", # required if resource_type=controller
    "resource_name": "alameda-ai", # required
    "time_interval": "daily" # required # daily|weekly|monthly
  },
  "resource_type": "controller", # required # controller|namespace
  "rest_api_url": "https://172.31.3.130:31011"
}'
##### target config info end #####

if [ "$BASH_VERSION" = "" ]; then
  /bin/echo -e "\nPlease use bash to run the script." 1>&2
fi
```

You can select a shell script to be run in a Kubernetes cluster where the application is run. For Terraform integration, you can select the Terraform script.

Add a Profile

1. On the *Configuration / Auto Provisioning* page, click *Add Profile*.

The screenshot shows the 'Add Auto Provision Profile' dialog box. It includes a title bar, a close button, and several configuration sections. The 'Profile Name' section has a text input field. The 'Recommendation' section has three buttons: 'Daily' (selected), 'Weekly', and 'Monthly'. The 'Additional Adjustments' section contains 'Extra CPU Headroom' with buttons 'None', 'Small', 'Medium', 'Large', and 'Custom'. Below this is 'Allocation Constraints' for CPU, with checkboxes for 'Minimum CPU' and 'Maximum CPU', each followed by a text input and a unit dropdown (mcores). The 'Extra Memory Headroom' section has buttons 'None', 'Small', 'Medium', 'Large', and 'Custom'. Below this is 'Allocation Constraints' for Memory, with checkboxes for 'Minimum Memory' and 'Maximum Memory', each followed by a text input and a unit dropdown (MiB). The 'Trigger Condition' section has a checked checkbox and text 'When reducing the resources, only apply if the reduction is more than 10 % of currently configured.'. The 'Schedule' section has a 'Select' dropdown. At the bottom right are 'Cancel' and 'Save' buttons.

Profile Name – Specify a name for the profile.

Recommendation – Specify which system recommendations to use (daily, weekly, or monthly).

Adjustments: Extra Headroom – If desired, specify any adjustments to make on top of the system recommendations for CPU and memory. *Small* means that the adjustment is 10% more than the recommendation, *Medium* means 20%, and *Large* means 30%. You can also specify a custom adjustment (millicores or percent for CPU; MB, GB, or percent for memory).

Adjustments: Allocation Constraints – If desired, specify minimum and maximum limits for CPU and memory. Resources will not be deployed if above or below these boundaries.

Trigger Condition – Recommendations may trigger a reduction of resources. If desired, specify a percentage to limit reduction of resources that are currently configured. The application will be restarted when resources are reduced. Therefore, you should be careful not to make the difference too small, causing frequent application restarts.

Schedule – Specify when the resource recommendations should be applied. If you are using *Daily* recommendations, you may want to apply changes hourly, daily, or automatically at midnight. For *Weekly* recommendations, you may want to apply changes hourly, daily, weekly, or automatically (12:00 a.m. Sunday). For *Monthly* recommendations, you may want to apply changes daily, weekly, monthly, or automatically (12:00 a.m. on the first day of the month). Note that all times are local.

2. Click *Save* when you are done.

Manage Profiles

You can do the following from the *Configuration / Auto Provisioning* page:

- Edit a profile. To do this, click the *Edit Profile* icon.
- Remove a profile. You can only remove a profile if it is not being used by a namespace or controller. To do this, click the *Remove Profile* icon.

Related topic:

[Applications](#)

[Terminology](#)

[Search/Sort Information in Tables](#)

[Configure Applications](#)

Configuration – System Settings

The *System Settings* page has tabs that allow you to:

- Change the admin password.
- Update metrics data source information.
- Set system notification.
- Manage the system license.
- Set the policy/update price books.

Admin Password

To access this page, select *Configuration, System Settings, Admin Password* tab. The *Admin Password* page allows you to change the admin password.

You must know the current password and *New Password* must match *Confirm Password*.

Metrics Data Source

The *Metrics Data Source* page allows you to set the authentication values that are needed by clusters to access metrics from different data sources. To access this page, select *Configuration, System Settings, Metrics Data Source* tab. You will need to select a cluster type, data source, and cluster.

For Kubernetes clusters, the available data sources are Prometheus, Datadog, and Sysdig.

For VM clusters, the available data sources are vCenter and AWS CloudWatch.

- For Datadog, the *API Key* and *Application Key* are required for authentication.
- For Sysdig, a *URL* and *Token* are required for authentication.
- For the Prometheus open-source monitoring system, the *URL* is required but the *Token* is optional.
- For vCenter, a *Login ID* and *Password* are required for authentication to the specified vCenter. A *Cluster Path* is needed for identifying VMs to be managed.
- For AWS CloudWatch, the *Region*, *Access Key ID*, and *Secret Access Key* are required for authentication.

When you are done, click *Test Connection* to confirm that all information is correct.

Note that for VMware, you can only modify cluster information if there are no VMs being monitored for the cluster.

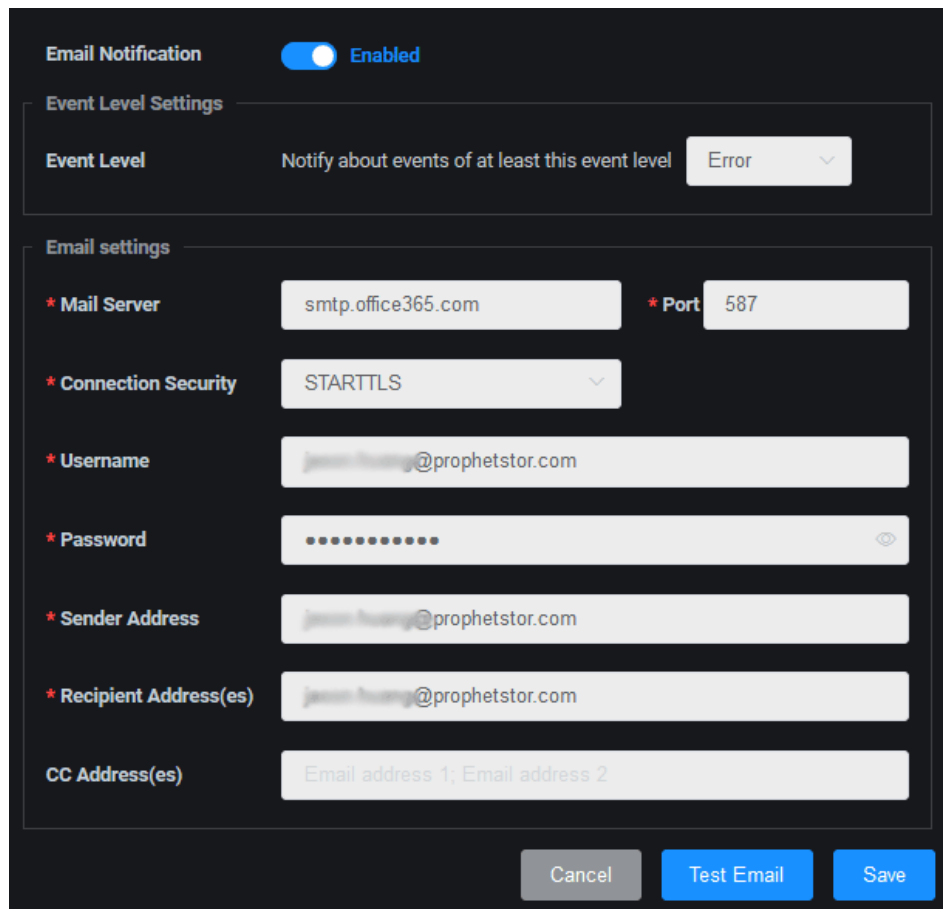
Notification

The *Notification* page allows you to configure email notifications to administrators when system errors and fatal issues occur. To access this page, select *Configuration, System Settings, Notification* tab.

Enable Notification

Follow the steps below to enable notification:

1. Toggle the *Email Notification* icon to *Enabled*.



The screenshot shows the 'Email Notification' configuration window. At the top, there is a toggle switch labeled 'Email Notification' which is currently turned on, indicated by a blue circle and the word 'Enabled'. Below this is the 'Event Level Settings' section, which includes a label 'Event Level' and a text input field containing 'Notify about events of at least this event level'. To the right of this text is a dropdown menu currently set to 'Error'. The 'Email settings' section follows, containing several fields: '* Mail Server' with the value 'smtp.office365.com', '* Port' with the value '587', '* Connection Security' with a dropdown set to 'STARTTLS', '* Username' with a masked email address '@prophetstor.com', '* Password' with a masked password field and an eye icon, '* Sender Address' with a masked email address '@prophetstor.com', '* Recipient Address(es)' with a masked email address '@prophetstor.com', and 'CC Address(es)' with a placeholder 'Email address 1; Email address 2'. At the bottom right, there are three buttons: 'Cancel', 'Test Email', and 'Save'.

Event Level – Select the minimum event level that should trigger notification. Higher levels will also trigger an email. For example, if you select *Error*, fatal events will also trigger an email.

Mail Server - Specify the mail server that should be used to send notification emails.

Port - Specify the mail server port that should be used.

Connection Security – Specify the protocol used by the mail server to secure email transmissions.

Username/Password - Specify the user account that will be used to log into the mail server.

Sender Address - Specify the email account that will be used in the “From” field of emails that are sent.

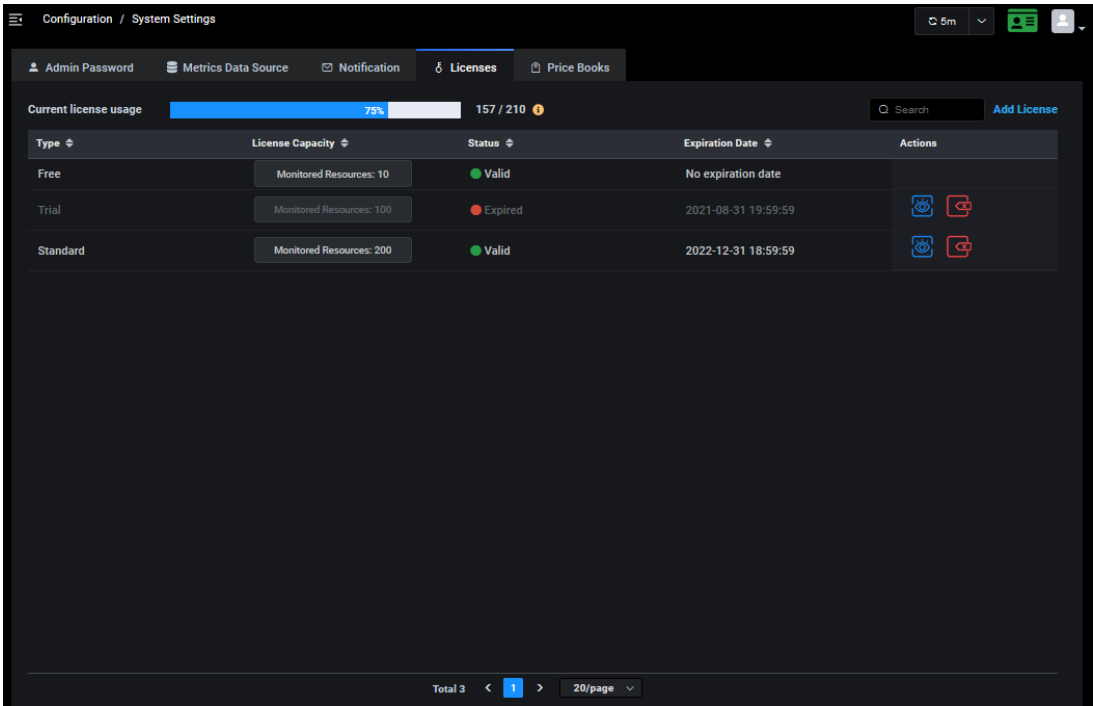
Recipient Address(es) - Specify the email address of the account that will receive emails. This will be used in the “To” field of emails. Separate multiple email addresses with semicolons.

CC Address(es) - Specify any other email accounts that should receive emails. Separate multiple email addresses with semicolons.

2. Click *Test Email* to confirm that all information is correct.
3. Once the test emails are received, click *Save*.

Licenses

To access this page, select *Configuration, System Settings, Licenses* tab. The *Licenses* page displays your current system licenses, including license type, status, and expiration date. It also shows what capacity is included in your license and your current usage.



The *Current license usage* graph displays the percentage of licensed resources being monitored as well as the number of resources being monitored and the total for which you are licensed. For Kubernetes clusters, licensed resources include nodes, namespaces in “Monitoring” state, and configured controllers. For VM clusters, licensed resources include VMs.

There are several license types, *Free*, *Trial*, and *Standard*. By default, a *Free* license with 10 resources is automatically applied when Federator.ai is installed. A *Trial* license may be provided during Federator.ai evaluation. Once a *Standard* license applied, the *Trial* license expires. The number of licensed resources is cumulative and includes the total number *Free* and *Standard* resources; for trials, it includes the total number of *Free* and *Trial* resources.

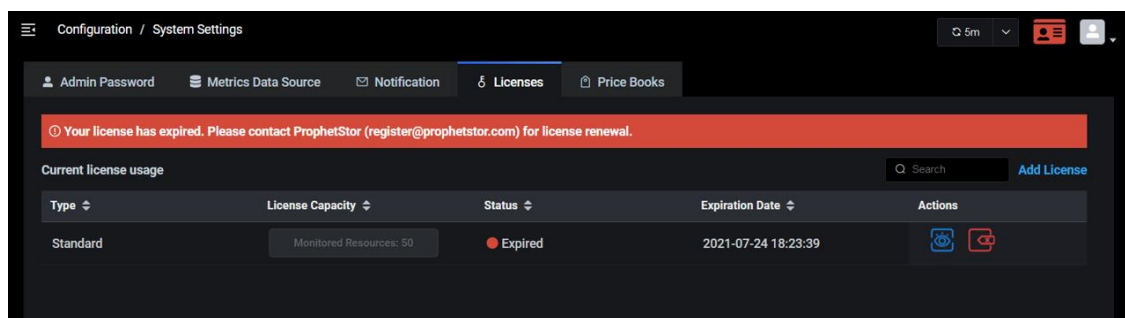
If you reach the number of licensed resources, there is a 30-day grace period, at which time the system prevents you from adding resources and some product functions, such as predictions, recommendations, cost analysis, and autoscaling, will stop. However, data collection will continue until additional license capacity is purchased.

A *Standard* license must be activated within 30 days after install, however a license with a status of *Pending Activation* will still make predictions and recommendations during that period. The status will be *Valid* once the license is activated.

The *Expiration Date* shows when the license expires. If the expiration of a license results in the system exceeding the license limit, predictions, recommendations, cost analysis, and autoscaling will not work.

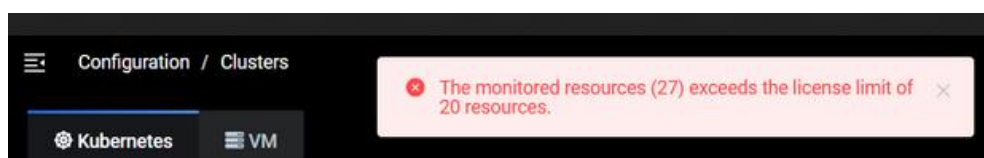
License Expiration

- Free licenses – Do not expire.
- Trial licenses – Expire on the expiration date displayed. A warning will be displayed when a trial license is expiring within 7 days. There is no grace period for an expired trial license. A trial license will be immediately expired when a standard license is added to the system.
- Standard licenses – Expire on the expiration date displayed. A warning will be displayed when a standard license is expiring within 7 days and will continue during the 30-day grace period after expiration. If a standard license expires and a remaining license does not have enough capacity, predictions and recommendations will stop.

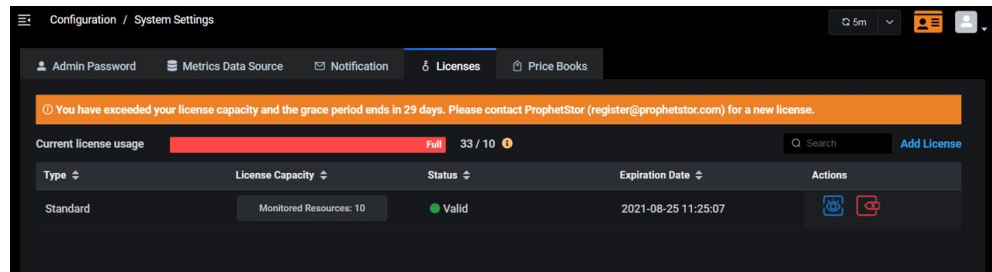


License Limits and Grace Periods

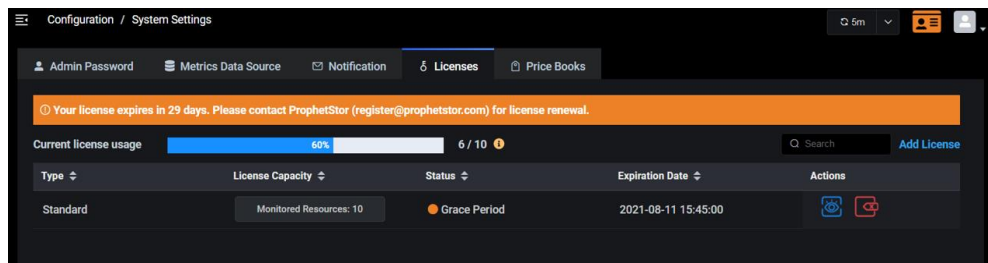
- In most cases, no new monitored resources can be added when there is not enough license capacity.





- If new cluster nodes are added to a cluster and there is not enough license capacity for new cluster nodes, there will be a 30-day grace period. After the grace period, predictions and recommendations will stop.



- There is a 30-day grace period after expiration of a standard license.



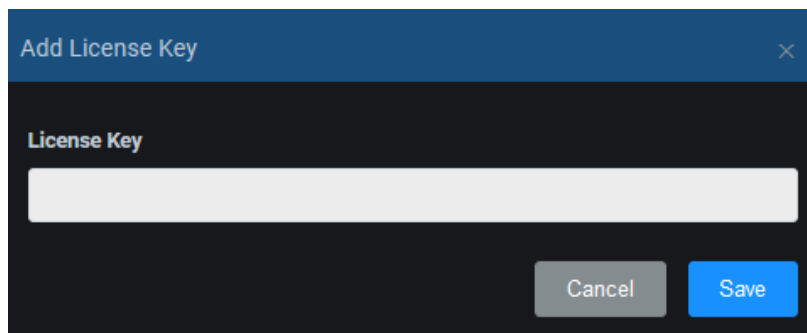
In addition to adding and activating licenses, you can perform the following functions from the *Licenses* page:

Icon	Function
	Show license key.
	Remove a license key.

Add a License

Follow the steps below to add and register a license:

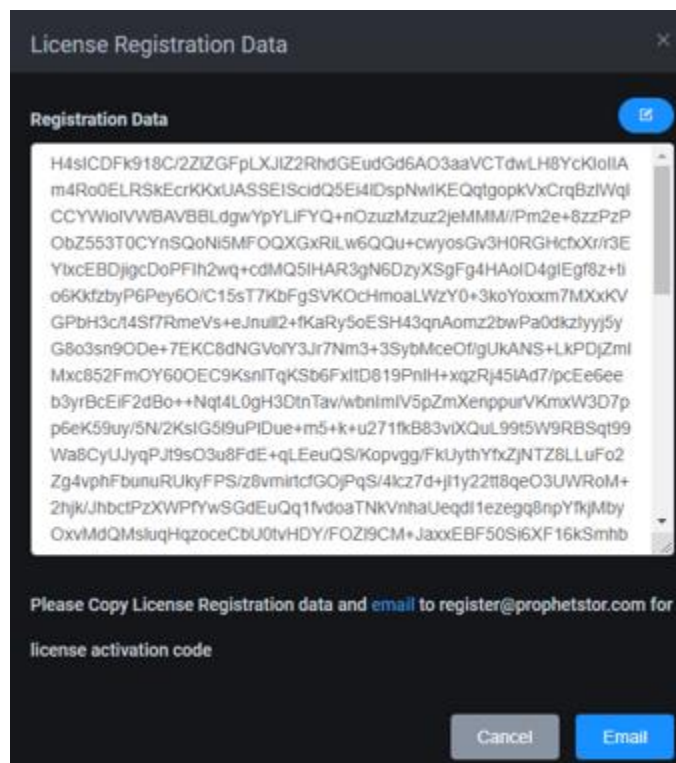
1. Click the *Add License* icon and enter the license key.

A dialog box titled "Add License Key" with a close button (X) in the top right corner. It features a text input field labeled "License Key" and two buttons at the bottom: "Cancel" and "Save".

2. Click *Add*.

A trial license is valid immediately after it is added. A standard license needs to be registered in order to be activated.

3. Email the registration data to register@prophetstor.com for a license activation code.

A dialog box titled "License Registration Data" with a close button (X) in the top right corner. It contains a text area with a copy icon (two overlapping squares) in the top right corner. Below the text area, there is a message: "Please Copy License Registration data and email to register@prophetstor.com for license activation code". At the bottom, there are two buttons: "Cancel" and "Email".

Registration Data

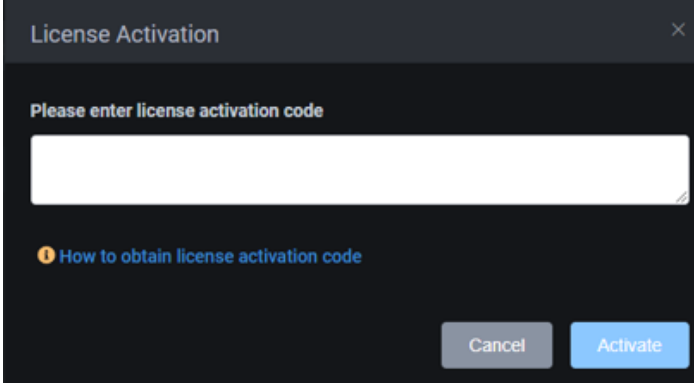
H4sICDFk918C/2ZIZGFpLXJIZ2RhGEudGd6AO3aaVCTdwLH8YcKlOIA
m4Ro0ELRSkEcrKKxUASSEIScidQ5Ei4IDspNwiKEQgtgopkVxCrQbzIWql
CCYWoIVVBAVB8LdgwYpYLIFYQ+nOzuzMzuz2jeMMM//Pm2e+8zzPzP
ObZ553T0CYnSQoNi5MFOQXGxRILw6QQu+cwyosGv3H0RGHctfXr/r3E
YlxcEBDjgcDoPFih2wq+cdMQ5IHAR3gN6DzyXSgFg4HAoID4glEgf8z+ti
o6KkftzbyP6Pey6O/C15sT7KbFgSVKOcHmoaLWzY0+3koYoxxm7MXXdKV
GPbH3c/4Sf7RmeVs+eJnull2+fkARy5oESH43qnAomz2bwPa0dkzlyy5y
G8o3sn9ODE+7EKC8dNGVofY3Jr7Nm3+3SybMceOfIgUKANS+LkPDjZml
Mxc852FmOY60OEC9KsnITqKSb6FxtID819PnlH+qxzRj45iAd7/pcEe6ee
b3yrBcEiF2dBo++Nqt4L0gH3DtnTav/wbnimIV5pZmXenppurVKmxW3D7p
p6eK59uy/5N/2KsIG5I9uPiDue+m5+k+u271fkB83viXQuL99t5W9RBSqt99
Wa8CyUJyqPJt9sO3u8FdE+qLEuQS/Kopvgg/FkUythYxZJNTZ8LLuFo2
Zg4vphFbunuRUkyFPS/z8vmirtcfGOjPqS/4kcz7d+jf1y22t8qeO3UWRoM+
2hjk/JhbctPzXWPFYwSGdEuQq1fvdaTNkVnhaUeqd11ezegq8npYtkjMby
OxvMdmQmsluqHqzocCbU0tvHDY/FOZi9CM+JaxxEbF50Si6XF16kSmhb

Please Copy License Registration data and email to register@prophetstor.com for
license activation code

Click the *Email* button to launch your email program. Click the *Copy* icon to copy the registration text and paste it into your email before sending.

The license status will now say *Pending Activation*. It must be activated within 30 days.

4. Once you have received the license activation code, click the *Activate License* icon and paste the activation code.

A dark-themed dialog box titled "License Activation" with a close button (X) in the top right corner. Inside the dialog, there is a text prompt "Please enter license activation code" above a large white text input field. Below the input field is a link with an information icon and the text "How to obtain license activation code". At the bottom right of the dialog are two buttons: a gray "Cancel" button and a blue "Activate" button.

5. Click *Activate*.

Manage Licenses

You can do the following from the *Licenses* page:

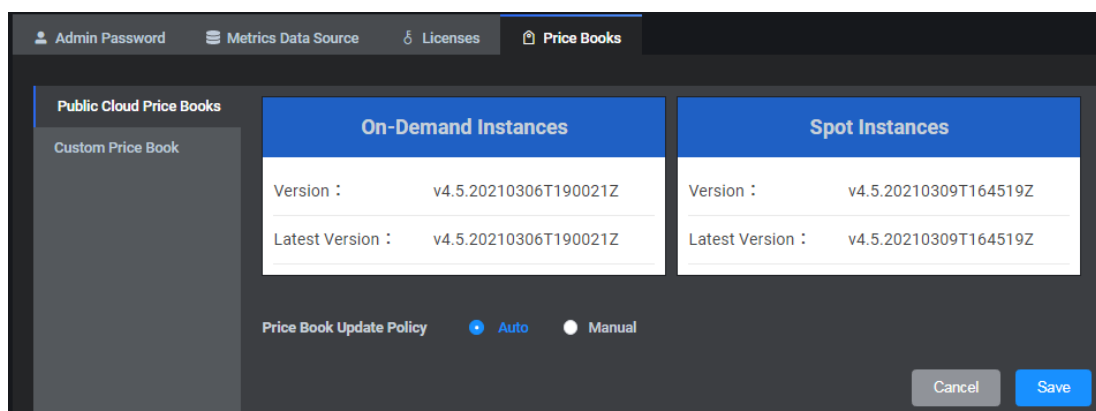
- View the key for your license. To do this, click the *Show License Key* icon.
- Remove the license. To do this, click the *Remove License Key* icon and confirm the removal.

Price Books

To access this page, select *Configuration, System Settings, Price Books* tab.

Public Cloud Price Books

The *Public Cloud Price Books* page displays the versions of the on-demand and spot instance price books that have been collected from the following cloud providers: Amazon Web Services (AWS), Google Cloud, and Microsoft Azure.



The screenshot shows the 'Price Books' configuration page. The top navigation bar includes 'Admin Password', 'Metrics Data Source', 'Licenses', and 'Price Books'. The left sidebar has 'Public Cloud Price Books' and 'Custom Price Book'. The main content area is divided into two columns: 'On-Demand Instances' and 'Spot Instances'. Each column displays 'Version' and 'Latest Version' as v4.5.20210306T190021Z and v4.5.20210309T164519Z respectively. Below these columns is a 'Price Book Update Policy' section with radio buttons for 'Auto' (selected) and 'Manual'. At the bottom right are 'Cancel' and 'Save' buttons.

The version numbers represent the ProphetStor version of the compiled price books.

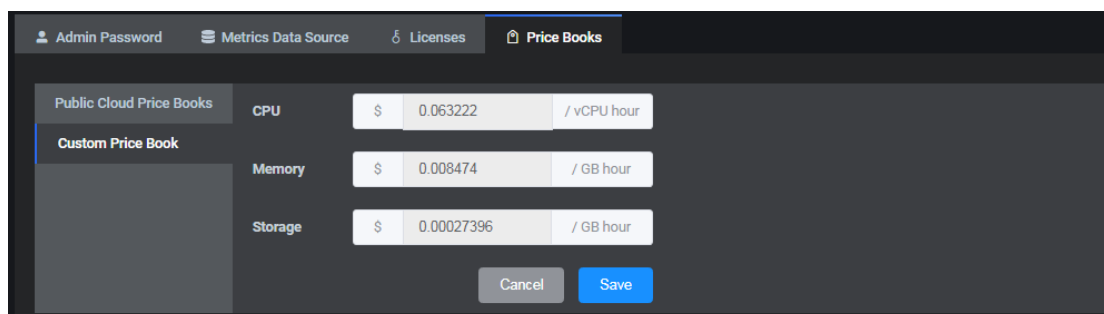
The *Price Book Update Policy* allows you to specify how you want the price books updated. *Auto* regularly checks availability and automatically downloads the latest version, which requires Internet access.

If *Manual* is selected and the *Latest Version* is different from *Version*, Federator.ai is aware of a newer price book but it has not been applied. Click the *Update Price Books* button to update.

Note that the browser needs Internet access in order to display the *Latest Version*. If your computer does not have access, the data must be pushed to system.

Custom Price Book

The *Custom Price Book* page allows you to define your hourly operating costs for CPU, memory, and storage and use these numbers for calculating costs/savings instead of using cloud provider pricing.



The screenshot shows the 'Custom Price Book' configuration page. The top navigation bar is the same as the previous page. The left sidebar has 'Public Cloud Price Books' and 'Custom Price Book'. The main content area has three rows: 'CPU' with a value of \$ 0.063222 / vCPU hour, 'Memory' with a value of \$ 0.008474 / GB hour, and 'Storage' with a value of \$ 0.00027396 / GB hour. At the bottom are 'Cancel' and 'Save' buttons.

When determining your hourly costs, be sure to include electricity, cooling/heating, labor, hardware, etc.

Events

The *Events* page displays all system events that have occurred.

There are five levels of events:

- Fatal - Issues that may stop the system from operating properly.
- Error - Indicates that a failure has occurred.
- Warning - Indicates that something occurred that may require maintenance or corrective action; however, the system is still operational.
- Info - Day-to-day activities, which require no action.
- Debug - Detailed activities used for troubleshooting.

You can filter by the event level, event type, and the time frame to display. For example, if you select the warning level, all warnings, errors, and fatal events will be displayed for the specified time period.

To specify a custom range, select *Custom* under *Time Range* and then specify a date range.

The screenshot displays the 'Events' page with a dark theme. At the top, there's a 'Filter' section with three dropdowns: 'Event Level' (set to 'Info'), 'Event Type' (set to 'All'), and 'Time Range' (set to 'Last 24 hours'). Below the filter is a summary bar with five columns: 'Fatal' (0), 'Error' (0), 'Warning' (0), 'Info' (4), and 'Debug' (0). The main section is a table of events with columns: 'Time', 'Level', 'Cluster', 'Resource', 'Namespace', 'Event Type', and 'Message'. The table contains four rows of event data. At the bottom, there is a pagination bar showing 'Total 4', a page indicator '1', and a dropdown for '20/page'.

Time	Level	Cluster	Resource	Namespace	Event Type	Message
2021-10-28 13:05:25	Info	Federator.ai	Price Book			Updated spot price book from v4.7.20211028T105428Z to v4.7.20211028T165232Z successfully.
2021-10-28 07:07:37	Info	Federator.ai	Price Book			Updated spot price book from v4.7.20211028T045241Z to v4.7.20211028T105428Z successfully.
2021-10-28 01:07:54	Info	Federator.ai	Price Book			Updated spot price book from v4.7.20211027T225233Z to v4.7.20211028T045241Z successfully.
2021-10-27 19:03:52	Info	Federator.ai	Price Book			Updated spot price book from v4.7.20211027T172147Z to v4.7.20211027T225233Z successfully.

At the bottom of the page, the total number messages being displayed is shown, along with navigation to other pages. You can also determine how many events to show per page.

Related topics:

[Terminology](#)

[Search/Sort Information in Tables](#)