



Red Hat OpenShift Serverless 1.29

통합

OpenShift Serverless와 Service Mesh 및 Cost Mangement 서비스와 통합

Red Hat OpenShift Serverless 1.29 통합

OpenShift Serverless와 Service Mesh 및 Cost Mangement 서비스와 통합

법적 공지

Copyright © 2023 Red Hat, Inc.

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at

<http://creativecommons.org/licenses/by-sa/3.0/>

. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, the Red Hat logo, JBoss, OpenShift, Fedora, the Infinity logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux[®] is the registered trademark of Linus Torvalds in the United States and other countries.

Java[®] is a registered trademark of Oracle and/or its affiliates.

XFS[®] is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL[®] is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js[®] is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack[®] Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

초록

이 문서에서는 OpenShift Serverless와 Service Mesh를 통합하는 방법에 대한 정보를 제공합니다. 또한 비용 관리 서비스를 사용하여 비용을 이해하고 추적하는 데 도움이 되며 서버리스 애플리케이션에서 NVIDIA GPU 리소스를 사용하는 방법을 보여줍니다.

차례

1장. OPENSIFT SERVERLESS와 SERVICE MESH 통합	3
1.1. 사전 요구 사항	3
1.2. 수신 외부 트래픽을 암호화하기 위한 인증서 생성	3
1.3. OPENSIFT SERVERLESS와 SERVICE MESH 통합	4
1.4. MTLS와 함께 서비스 메시지를 사용할 때 KNATIVE SERVING 메트릭 활성화	8
1.5. KOURIER가 활성화된 경우 OPENSIFT SERVERLESS와 SERVICE MESH 통합	9
1.6. 서비스 메시에 보안 필터링을 사용하여 NET-ISTIO 메모리 사용량 개선	11
2장. SERVERLESS와 COST MANAGEMENT SERVICE 통합	13
2.1. 사전 요구 사항	13
2.2. 비용 관리 쿼리에 라벨 사용	13
2.3. 추가 리소스	13
3장. 서버리스 애플리케이션과 함께 NVIDIA GPU 리소스 사용	14
3.1. 서비스에 대한 GPU 요구 사항 지정	14
3.2. OPENSIFT CONTAINER PLATFORM의 추가 리소스	14

1장. OPENSIFT SERVERLESS와 SERVICE MESH 통합

OpenShift Serverless Operator는 Kourier를 Knative의 기본 수신으로 제공합니다. 그러나 Kourier가 활성화되어 있는지 여부에 관계없이 OpenShift Serverless에서 Service Mesh를 사용할 수 있습니다. Kourier disabled와 통합하면 mTLS 기능과 같이 Kourier 수신에서 지원하지 않는 추가 네트워킹 및 라우팅 옵션을 구성할 수 있습니다.



중요

OpenShift Serverless는 본 안내서에 명시되어 있는 Red Hat OpenShift Service Mesh 기능만 지원하며 문서화되지 않은 다른 기능은 지원하지 않습니다.

1.1. 사전 요구 사항

- 다음 절차의 예제에서는 도메인 **example.com**을 사용합니다. 이 도메인에 사용되는 예제 인증서는 하위 도메인 인증서에 서명하는 CA(인증 기관)로 사용됩니다. 배포에서 이 절차를 완료하고 확인하려면 널리 신뢰받는 공용 CA에서 서명한 인증서 또는 조직에서 제공하는 CA가 필요합니다. 도메인, 하위 도메인, CA에 따라 예제 명령을 조정해야 합니다.
- OpenShift Container Platform 클러스터의 도메인과 일치하도록 와일드카드 인증서를 구성해야 합니다. 예를 들어 OpenShift Container Platform 콘솔 주소가 <https://console-openshift-console.apps.openshift.example.com>인 경우 도메인이 ***.apps.openshift.example.com**이 되도록 와일드카드 인증서를 구성해야 합니다. 와일드카드 인증서 구성에 대한 자세한 내용은 수신되는 외부 트래픽을 암호화하기 위한 인증서 생성에 관한 다음의 내용을 참조하십시오.
- 기본 OpenShift Container Platform 클러스터 도메인의 하위 도메인이 아닌 도메인 이름을 사용하려면 해당 도메인에 대한 도메인 매핑을 설정해야 합니다. 자세한 내용은 [사용자 정의 도메인 매핑 생성에 대한 OpenShift Serverless 설명서](#)를 참조하십시오.

1.2. 수신 외부 트래픽을 암호화하기 위한 인증서 생성

기본적으로 Service Mesh mTLS 기능은 사이드카가 있는 수신 게이트웨이와 개별 Pod 간에 Service Mesh 자체의 트래픽만 보호합니다. OpenShift Container Platform 클러스터로 전달될 때 트래픽을 암호화하려면 OpenShift Serverless 및 Service Mesh 통합을 활성화하기 전에 인증서를 생성해야 합니다.

사전 요구 사항

- OpenShift Container Platform에 대한 클러스터 관리자 권한이 있거나 AWS 또는 OpenShift Dedicated의 Red Hat OpenShift Service에 대한 클러스터 또는 전용 관리자 권한이 있습니다.
- OpenShift Serverless Operator 및 Knative Serving이 설치되어 있습니다.
- OpenShift CLI(**oc**)를 설치합니다.
- 프로젝트를 생성했거나 애플리케이션 및 기타 워크로드를 생성할 수 있는 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.

절차

1. Knative 서비스의 인증서에 서명할 root 인증서 및 개인 키를 생성합니다.

```
$ openssl req -x509 -sha256 -nodes -days 365 -newkey rsa:2048 \
  -subj '/O=Example Inc./CN=example.com' \
  -keyout root.key \
```

```
-out root.crt
```

2. 와일드카드 인증서를 만듭니다.

```
$ openssl req -nodes -newkey rsa:2048 \
  -subj "/CN=*.apps.openshift.example.com/O=Example Inc." \
  -keyout wildcard.key \
  -out wildcard.csr
```

3. 와일드카드 인증서를 서명합니다.

```
$ openssl x509 -req -days 365 -set_serial 0 \
  -CA root.crt \
  -CAkey root.key \
  -in wildcard.csr \
  -out wildcard.crt
```

4. 와일드카드 인증서를 사용하여 시크릿을 생성합니다.

```
$ oc create -n istio-system secret tls wildcard-certs \
  --key=wildcard.key \
  --cert=wildcard.crt
```

이 인증서는 OpenShift Serverless를 Service Mesh와 통합할 때 생성된 게이트웨이에서 선택하여 수신 게이트웨이가 이 인증서와 트래픽을 제공하도록 합니다.

1.3. OPENSIFT SERVERLESS와 SERVICE MESH 통합

Kourier를 기본 수신으로 사용하지 않고 OpenShift Serverless와 Service Mesh를 통합할 수 있습니다. 이렇게 하려면 다음 절차를 완료하기 전에 Knative Serving 구성 요소를 설치하지 마십시오. 일반적인 Knative Serving 설치 절차에서는 **KnativeServing** CRD(사용자 정의 리소스 정의)를 생성하여 Knative Serving을 Service Mesh와 통합할 때 추가 단계가 필요합니다. 이 절차는 서비스 메시를 OpenShift Serverless 설치에 대한 기본값 및 수신만 통합하려는 경우에 유용할 수 있습니다.

사전 요구 사항

- OpenShift Container Platform에 대한 클러스터 관리자 권한이 있거나 AWS 또는 OpenShift Dedicated의 Red Hat OpenShift Service에 대한 클러스터 또는 전용 관리자 권한이 있습니다.
- 프로젝트를 생성했거나 애플리케이션 및 기타 워크로드를 생성할 수 있는 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.
- Red Hat OpenShift Service Mesh Operator를 설치하고 **istio-system** 네임스페이스에서 **ServiceMeshControlPlane** 리소스를 생성합니다. mTLS 기능을 사용하려면 **ServiceMeshControlPlane** 리소스의 **spec.security.dataPlane.mtls** 필드도 **true** 로 설정해야 합니다.



중요

OpenShift Serverless를 Service Mesh와 함께 사용하는 것은 Red Hat OpenShift Service Mesh 버전 2.0.5 이상에서만 지원됩니다.

- OpenShift Serverless Operator를 설치합니다.

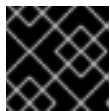
- OpenShift CLI(**oc**)를 설치합니다.

절차

1. Service Mesh와 통합할 네임스페이스를 **ServiceMeshMemberRoll** 오브젝트에 멤버로 추가합니다.

```
apiVersion: maistra.io/v1
kind: ServiceMeshMemberRoll
metadata:
  name: default
  namespace: istio-system
spec:
  members: 1
    - knative-serving
    - <namespace>
```

- 1** Service Mesh와 통합할 네임스페이스 목록입니다.



중요

이 네임스페이스 목록에 **knative-serving** 네임스페이스가 포함되어야 합니다.

2. **ServiceMeshMemberRoll** 리소스를 적용합니다.

```
$ oc apply -f <filename>
```

3. Service Mesh가 트래픽을 수락할 수 있도록 필요한 게이트웨이를 생성합니다.

HTTP를 사용하는 **knative-local-gateway** 오브젝트의 예

```
apiVersion: networking.istio.io/v1alpha3
kind: Gateway
metadata:
  name: knative-ingress-gateway
  namespace: knative-serving
spec:
  selector:
    istio: ingressgateway
  servers:
    - port:
        number: 443
        name: https
        protocol: HTTPS
      hosts:
        - "*"
      tls:
        mode: SIMPLE
        credentialName: <wildcard_certs> 1
  ---
apiVersion: networking.istio.io/v1alpha3
kind: Gateway
metadata:
```

```

name: knative-local-gateway
namespace: knative-serving
spec:
  selector:
    istio: ingressgateway
  servers:
  - port:
      number: 8081
      name: http
      protocol: HTTP 2
    hosts:
    - "*"
  ---
apiVersion: v1
kind: Service
metadata:
  name: knative-local-gateway
  namespace: istio-system
  labels:
    experimental.istio.io/disable-gateway-port-translation: "true"
spec:
  type: ClusterIP
  selector:
    istio: ingressgateway
  ports:
  - name: http2
    port: 80
    targetPort: 8081

```

1 와일드카드 인증서가 포함된 시크릿 이름을 추가합니다.

2 **knative-local-gateway**는 HTTP 트래픽을 제공합니다. HTTP를 사용하면 Service Mesh 외부에서 들어오는 트래픽이 표시되지만 **example.default.svc.cluster.local**과 같은 내부 호스트 이름을 사용하는 것은 암호화되지 않습니다. 다른 와일드카드 인증서 및 다른 **protocol** 사양을 사용하는 추가 게이트웨이를 생성하여 이 경로에 대한 암호화를 설정할 수 있습니다.

HTTPS를 사용하는 **knative-local-gateway** 오브젝트의 예

```

apiVersion: networking.istio.io/v1alpha3
kind: Gateway
metadata:
  name: knative-local-gateway
  namespace: knative-serving
spec:
  selector:
    istio: ingressgateway
  servers:
  - port:
      number: 443
      name: https
      protocol: HTTPS
    hosts:
    - "*"

```

```

tls:
  mode: SIMPLE
  credentialName: <wildcard_certs>

```

4. **Gateway** 리소스를 적용합니다.

```
$ oc apply -f <filename>
```

5. 다음 **KnativeServing** CRD(사용자 정의 리소스 정의)를 생성하여 Knative Serving을 설치하여 Istio 통합을 활성화합니다.

```

apiVersion: operator.knative.dev/v1beta1
kind: KnativeServing
metadata:
  name: knative-serving
  namespace: knative-serving
spec:
  ingress:
    istio:
      enabled: true ①
  deployments: ②
  - name: activator
    annotations:
      "sidecar.istio.io/inject": "true"
      "sidecar.istio.io/rewriteAppHTTPProbers": "true"
  - name: autoscaler
    annotations:
      "sidecar.istio.io/inject": "true"
      "sidecar.istio.io/rewriteAppHTTPProbers": "true"

```

- ① Istio 통합을 활성화합니다.
- ② Knative Serving 데이터 플레인 Pod에 사이드카 삽입을 활성화합니다.

6. **KnativeServing** 리소스를 적용합니다.

```
$ oc apply -f <filename>
```

7. 사이드카 삽입이 활성화되고 패스쓰루(pass-through) 경로를 사용하는 Knative 서비스를 생성합니다.

```

apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: <service_name>
  namespace: <namespace> ①
  annotations:
    serving.knative.openshift.io/enablePassthrough: "true" ②
spec:
  template:
    metadata:
      annotations:
        sidecar.istio.io/inject: "true" ③

```

```
sidecar.istio.io/rewriteAppHTTPProbers: "true"
spec:
  containers:
  - image: <image_url>
```

- 1 Service Mesh 멤버 룰에 포함된 네임스페이스입니다.
- 2 생성한 인증서가 수신 게이트웨이를 통해 직접 제공되도록 Knative Serving에 OpenShift Container Platform 패스스루 지원 경로를 생성하도록 지시합니다.
- 3 Service Mesh 사이드카를 Knative 서비스 Pod에 삽입합니다.

8. **Service** 리소스를 적용합니다.

```
$ oc apply -f <filename>
```

검증

- CA에서 신뢰하는 보안 연결을 사용하여 서버리스 애플리케이션에 액세스합니다.

```
$ curl --cacert root.crt <service_url>
```

명령 예

```
$ curl --cacert root.crt https://hello-default.apps.openshift.example.com
```

출력 예

```
Hello Openshift!
```

1.4. mTLS와 함께 서비스 메시를 사용할 때 KNATIVE SERVING 메트릭 활성화

Service Mesh가 mTLS를 사용하여 활성화되면 Service Mesh가 Prometheus가 메트릭 스크랩을 허용하지 않기 때문에 기본적으로 Knative Serving에 대한 메트릭이 비활성화됩니다. 이 섹션에서는 서비스 메시 및 mTLS를 사용할 때 Knative Serving 메트릭을 활성화하는 방법을 보여줍니다.

사전 요구 사항

- OpenShift Serverless Operator 및 Knative Serving이 클러스터에 설치되어 있습니다.
- mTLS 기능이 활성화된 Red Hat OpenShift Service Mesh를 설치했습니다.
- OpenShift Container Platform에 대한 클러스터 관리자 권한이 있거나 AWS 또는 OpenShift Dedicated의 Red Hat OpenShift Service에 대한 클러스터 또는 전용 관리자 권한이 있습니다.
- OpenShift CLI(**oc**)를 설치합니다.
- 프로젝트를 생성했거나 애플리케이션 및 기타 워크로드를 생성할 수 있는 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.

절차

1. Knative Serving 사용자 정의 리소스(CR)의 **observability** 사양에서 **prometheus**를 **metrics.backend-destination**으로 지정합니다.

```

apiVersion: operator.knative.dev/v1beta1
kind: KnativeService
metadata:
  name: knative-serving
spec:
  config:
    observability:
      metrics.backend-destination: "prometheus"
  ...

```

이 단계에서는 메트릭이 기본적으로 비활성화되지 않습니다.

2. Prometheus 네임스페이스의 트래픽을 허용하려면 다음 네트워크 정책을 적용합니다.

```

apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: allow-from-openshift-monitoring-ns
  namespace: knative-serving
spec:
  ingress:
    - from:
      - namespaceSelector:
          matchLabels:
            name: "openshift-monitoring"
      podSelector: {}
  ...

```

3. 다음 사양을 포함하도록 **istio-system** 네임스페이스에서 기본 서비스 메시 컨트롤 플레인을 수정하고 다시 적용합니다.

```

...
spec:
  proxy:
    networking:
      trafficControl:
        inbound:
          excludedPorts:
            - 8444
  ...

```

1.5. KOURIER가 활성화된 경우 OPENSIFT SERVERLESS와 SERVICE MESH 통합

Kourier가 이미 활성화되어 있어도 OpenShift Serverless에서 Service Mesh를 사용할 수 있습니다. 이 절에서는 Kourier를 사용하여 Knative Serving을 이미 설치했지만 나중에 서비스 메시 통합을 추가하도록 결정하는 경우 유용할 수 있습니다.

사전 요구 사항

- OpenShift Container Platform에 대한 클러스터 관리자 권한이 있거나 AWS 또는 OpenShift Dedicated의 Red Hat OpenShift Service에 대한 클러스터 또는 전용 관리자 권한이 있습니다.
- 프로젝트를 생성했거나 애플리케이션 및 기타 워크로드를 생성할 수 있는 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.
- OpenShift CLI(**oc**)를 설치합니다.
- 클러스터에 OpenShift Serverless Operator 및 Knative Serving을 설치합니다.
- Red Hat OpenShift Service Mesh를 설치합니다. OpenShift Serverless with Service Mesh and Kourier는 Red Hat OpenShift Service Mesh 버전 1.x 및 2.x에서 모두 사용이 지원되고 있습니다.

절차

1. Service Mesh와 통합할 네임스페이스를 **ServiceMeshMemberRoll** 오브젝트에 멤버로 추가합니다.

```

apiVersion: maistra.io/v1
kind: ServiceMeshMemberRoll
metadata:
  name: default
  namespace: istio-system
spec:
  members:
    - <namespace> ❶
  ...

```

❶ Service Mesh와 통합할 네임스페이스 목록입니다.

2. **ServiceMeshMemberRoll** 리소스를 적용합니다.

```

$ oc apply -f <filename>

```

3. Knative 시스템 Pod에서 Knative 서비스로의 트래픽 흐름을 허용하는 네트워크 정책을 생성합니다.

- a. Service Mesh와 통합할 각 네임스페이스에 **NetworkPolicy** 리소스를 생성합니다.

```

apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: allow-from-serving-system-namespace
  namespace: <namespace> ❶
spec:
  ingress:
    - from:
      - namespaceSelector:
          matchLabels:
            knative.openshift.io/part-of: "openshift-serverless"
  podSelector: {}
  policyTypes:
    - Ingress
  ...

```

1 Service Mesh와 통합할 네임스페이스를 추가합니다.



참고

knative.openshift.io/part-of: "openshift-serverless" 라벨이 OpenShift Serverless 1.22.0에 추가되었습니다. OpenShift Serverless 1.21.1 또는 이전 버전을 사용하는 경우 **knative.openshift.io/part-of** 레이블을 **knative-serving** 및 **knative-serving-ingress** 네임스페이스에 추가합니다.

knative-serving 네임스페이스에 라벨을 추가합니다.

```
$ oc label namespace knative-serving knative.openshift.io/part-of=openshift-serverless
```

knative-serving-ingress 네임스페이스에 라벨을 추가합니다.

```
$ oc label namespace knative-serving-ingress knative.openshift.io/part-of=openshift-serverless
```

b. NetworkPolicy 리소스를 적용합니다.

```
$ oc apply -f <filename>
```

1.6. 서비스 메시에 보안 필터링을 사용하여 NET-ISTIO 메모리 사용량 개선

기본적으로 Kubernetes **client-go** 라이브러리의 **정보** 관리자 구현은 특정 유형의 모든 리소스를 가져옵니다. 이로 인해 많은 리소스가 사용 가능할 때 상당한 오버헤드가 발생할 수 있으므로 메모리 누수로 인해 Knative **net-istio** Ingress 컨트롤러가 대규모 클러스터에서 실패할 수 있습니다. 그러나 컨트롤러에서 Knative 관련 시크릿만 가져올 수 있는 Knative **net-istio** 수신 컨트롤러에 필터링 메커니즘을 사용할 수 있습니다. **KnativeServing** CR(사용자 정의 리소스)에 주석을 추가하여 이 메커니즘을 활성화할 수 있습니다.



중요

시크릿 필터링을 활성화하면 모든 시크릿에 **networking.internal.knative.dev/certificate-uid: "<id>"** 로 레이블이 지정됩니다. 그러지 않으면 Knative Serving에서 이를 탐지하지 않아 오류가 발생합니다. 새 시크릿과 기존 보안의 라벨을 모두 지정해야 합니다.

사전 요구 사항

- OpenShift Container Platform에 대한 클러스터 관리자 권한이 있거나 AWS 또는 OpenShift Dedicated의 Red Hat OpenShift Service에 대한 클러스터 또는 전용 관리자 권한이 있습니다.
- 프로젝트를 생성했거나 애플리케이션 및 기타 워크로드를 생성할 수 있는 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.
- Red Hat OpenShift Service Mesh를 설치합니다. OpenShift Serverless with Service Mesh는 Red Hat OpenShift Service Mesh 버전 2.0.5 이상에서만 지원됩니다.
- OpenShift Serverless Operator 및 Knative Serving을 설치합니다.
- OpenShift CLI(**oc**)를 설치합니다.

절차

- 서버리스.[openshift.io/enable-secret-informer-filtering](https://serverless.openshift.io/enable-secret-informer-filtering) 주석을 **KnativeServing** CR에 추가합니다.

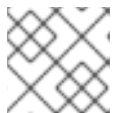
KnativeServing CR의 예

```

apiVersion: operator.knative.dev/v1beta1
kind: KnativeServing
metadata:
  name: knative-serving
  namespace: knative-serving
  annotations:
    serverless.openshift.io/enable-secret-informer-filtering: "true" 1
spec:
  ingress:
    istio:
      enabled: true
  deployments:
    - annotations:
        sidecar.istio.io/inject: "true"
        sidecar.istio.io/rewriteAppHTTPProbers: "true"
      name: activator
    - annotations:
        sidecar.istio.io/inject: "true"
        sidecar.istio.io/rewriteAppHTTPProbers: "true"
      name: autoscaler

```

- 1 이 주석을 추가하면 환경 변수 **ENABLE_SECRET_INFORMER_FILTER_BY_CERT_UID=true** 를 **net-istio** 컨트롤러 Pod에 삽입합니다.



참고

배포를 재정의하여 다른 값을 설정하면 이 주석은 무시됩니다.

2장. SERVERLESS와 COST MANAGEMENT SERVICE 통합

Cost Management 는 클라우드 및 컨테이너 비용을 더 잘 이해하고 추적할 수 있는 OpenShift Container Platform 서비스입니다. 오픈 소스 **Koku** 프로젝트를 기반으로 합니다.

2.1. 사전 요구 사항

- 클러스터 관리자 권한이 있어야 합니다.
- 비용 관리를 설정하고 **OpenShift Container Platform** 소스를 추가했습니다.

2.2. 비용 관리 쿼리에 라벨 사용

비용 관리에서 태그라고도 하는 레이블은 노드, 네임스페이스 또는 Pod에 적용할 수 있습니다. 각 레이블은 키와 값 쌍입니다. 여러 레이블을 결합하여 보고서를 생성할 수 있습니다. **Red Hat 하이브리드 콘솔** 을 사용하여 비용에 대한 보고서에 액세스할 수 있습니다.

레이블은 노드에서 네임스페이스로, 네임스페이스에서 Pod로 상속됩니다. 그러나 이미 리소스에 있는 경우 레이블을 재정의하지 않습니다. 예를 들어 Knative 서비스에는 기본 **app=<revision_name>** 라벨이 있습니다.

Knative 서비스 기본 라벨의 예

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  name: example-service
spec:
  ...
  labels:
    app: <revision_name>
  ...
```

app=my-domain과 같은 네임스페이스 레이블을 정의하는 경우 비용 관리 서비스는 **app=my-domain** 태그를 사용하여 애플리케이션을 쿼리할 때 **app= my-domain** 태그를 사용하여 애플리케이션을 쿼리할 때 Knative 서비스에서 발생하는 비용을 고려하지 않습니다. 이 태그가 있는 Knative 서비스의 비용은 **app=<revision_name>** 태그에서 쿼리해야 합니다.

2.3. 추가 리소스

- [소스에 대한 태그 지정 구성](#)
- [Cost Explorer](#)를 사용하여 비용을 시각화합니다.

3장. 서버리스 애플리케이션과 함께 NVIDIA GPU 리소스 사용

NVIDIA는 OpenShift Container Platform에서 GPU 리소스 사용을 지원합니다. [OpenShift Container Platform에서 GPU 리소스를 설정하는 방법에 대한 자세한 내용은 OpenShift의 GPU Operator 를 참조하십시오.](#)

3.1. 서비스에 대한 GPU 요구 사항 지정

OpenShift Container Platform 클러스터에 GPU 리소스를 활성화하면 Knative(**kn**) CLI를 사용하여 Knative 서비스에 대한 GPU 요구 사항을 지정할 수 있습니다.

사전 요구 사항

- OpenShift Serverless Operator, Knative Serving 및 Knative Eventing이 클러스터에 설치되어 있습니다.
- Knative(**kn**) CLI가 설치되어 있습니다.
- OpenShift Container Platform 클러스터에 GPU 리소스가 활성화됩니다.
- 프로젝트를 생성했거나 OpenShift Container Platform에서 애플리케이션 및 기타 워크로드를 생성하는 데 적절한 역할 및 권한이 있는 프로젝트에 액세스할 수 있습니다.



참고

OpenShift Container Platform 또는 OpenShift Dedicated의 IBM zSystems 및 IBM Power에서는 NVIDIA GPU 리소스 사용이 지원되지 않습니다.

절차

1. Knative 서비스를 생성하고 **--limit nvidia.com/gpu=1** 플래그를 사용하여 GPU 리소스 요구 사항 제한을 **1**로 설정합니다.

```
$ kn service create hello --image <service-image> --limit nvidia.com/gpu=1
```

GPU 리소스 요구 사항 제한이 **1**이면 서비스의 전용 GPU 리소스가 1개임을 나타냅니다. 서비스에서는 GPU 리소스를 공유하지 않습니다. GPU 리소스가 필요한 기타 서비스는 GPU 리소스를 더 이상 사용하지 않을 때까지 기다려야 합니다.

또한 GPU가 1개로 제한되면 GPU 리소스를 2개 이상 사용하는 애플리케이션이 제한됩니다. 서비스에서 GPU 리소스를 1개 이상 요청하는 경우 GPU 리소스 요구 사항을 충족할 수 있는 노드에 배포됩니다.

2. 선택 사항: 기존 서비스의 경우 **--limit nvidia.com/gpu=3** 플래그를 사용하여 GPU 리소스 요구 사항 제한을 **3**으로 변경할 수 있습니다.

```
$ kn service update hello --limit nvidia.com/gpu=3
```

3.2. OPENSIFT CONTAINER PLATFORM의 추가 리소스

- [확장 리소스의 리소스 할당량 설정](#)

