

上节课我们一起学习了如何在 Prometheus Operator 下面自定义一个监控选项，以及自定义报警规则的使用。那么我们还能够直接使用前面课程中的自动发现功能吗？如果在我们的 Kubernetes 集群中有了很多的 Service/Pod，那么我们都需要一个一个的去建立一个对应的 ServiceMonitor 对象来进行监控吗？这样岂不是又变得麻烦起来了？

自动发现配置

为解决上面的问题，Prometheus Operator 为我们提供了一个额外的抓取配置的来解决这个问题，我们可以通过添加额外的配置来进行服务发现进行自动监控。和前面自定义的方式一样，我们想要在 Prometheus Operator 当中去自动发现并监控具有`prometheus.io/scrape=true`这个 annotations 的 Service，之前我们定义的 Prometheus 的配置如下：

```
1 - job_name: 'kubernetes-service-endpoints'
2   kubernetes_sd_configs:
3     - role: endpoints
4   relabel_configs:
5     - source_labels: [__meta_kubernetes_service_annotation_prometheus_io_scrape]
6       action: keep
7       regex: true
8     - source_labels: [__meta_kubernetes_service_annotation_prometheus_io_scheme]
9       action: replace
10      target_label: __scheme__
11      regex: (https?)
12     - source_labels: [__meta_kubernetes_service_annotation_prometheus_io_path]
13       action: replace
14      target_label: __metrics_path__
15      regex: (.+)
16     - source_labels: [__address__,
17       __meta_kubernetes_service_annotation_prometheus_io_port]
18       action: replace
19      target_label: __address__
20      regex: ([^:])(?::\d+)?;(\d+)
21      replacement: $1:$2
22     - action: labelmap
23       regex: __meta_kubernetes_service_label_(.+)
24     - source_labels: [__meta_kubernetes_namespace]
25       action: replace
26      target_label: kubernetes_namespace
27     - source_labels: [__meta_kubernetes_service_name]
```

```
27     action: replace  
28     target_label: kubernetes_name
```

如果你对上面这个配置还不是很熟悉的话，建议去查看下前面关于 [Kubernetes常用资源对象监控章节的介绍](#)，要想自动发现集群中的 Service，就需要我们在 Service 的 annotation 区域添加 `prometheus.io/scrape=true` 的声明，将上面文件直接保存为 `prometheus-additional.yaml`，然后通过这个文件创建一个对应的 Secret 对象：

```
1 $ kubectl create secret generic additional-configs --from-file=prometheus-  
additional.yaml -n monitoring  
2 secret "additional-configs" created
```

注意我们所有的操作都在 Prometheus Operator 源码 `contrib/kube-prometheus/manifests/` 目录下面。

创建完成后，会将上面配置信息进行 base64 编码后作为 `prometheus-additional.yaml` 这个 key 对应的值存在：

```
1 $ kubectl get secret additional-configs -n monitoring -o yaml  
2 apiVersion: v1  
3 data:  
4   prometheus-additional.yaml:  
LSBqb2JfbmFtZTogJ2t1YmVybmV0ZXMtC2VydmljZS1lbnRwb21udHMnCiAga3ViZXJuZXRLc19zZF9jb25maWd  
z0gogIC0gcm9sZTogZW5kcG9pbnRzCiAgcmVsYWI1bF9jb25maWdz0gogIC0gc291cmNlX2xhYmVsczogW19fbW  
V0YV9rdWJ1cm51dGVzX3N1cnZpY2VfYW5ub3RhdG1vb19wcm9tZXRoZXVzX21vX3NjcmFwZV0KICAgIGFjdG1vb  
joga2V1cAogICAgnVnZXG6IHRydWUKICAtIHNVdXJjZV9sYWJ1bHM6IFtfX211dGFfa3ViZXJuZXRLc19zZXJ2  
aWNlX2Fubm90YXRpb25fcHJvbW0aGV1c19pb19zY2h1bWdCiAgICBhY3Rpb246IHJ1cGxhY2UKICAgIHRhcmd  
1dF9sYWJ1bDogX19zY2h1bWVfxwogICAgnVnZXG6ICChodHRwcz8pCiAgLSBzb3VyY2VfbGFizWxz0iBbX19tZX  
RhX2t1YmVybmV0ZXNfc2VydmljZV9hb5vdGF0aW9uX3Byb211dGhldXNfaW9fcGF0aF0KICAgIGFjdG1vbjogc  
mVwbGFjZQogICAgdGFyZ2V0X2xhYmVsOibfx211dHjpY3NfcGF0aF9fcAgICByZwdleDogKC4rKQogIC0gc291  
cmNlX2xhYmVsczogW19fYWRkcmVzc19fLCBfx211dGFfa3ViZXJuZXRLc19zZXJ2aWNlX2Fubm90YXRpb25fcHJ  
vbW0aGV1c19pb19wb3J0XQogICAgnYWN0aW9u0iByZXBsYWN1CiAgICB0YXJnZXrbfGFizWw6IF9fYWRkcmVzc1  
9fcAgICByZwdleDogKfTe0l0rKsg/OjpcZCspPzs0XGQrKQogICAgnVwbGFjZw1lbnQ6ICQx0iQyCiAgLSBhY  
3Rpb246IGxhYmVsbfWciAgICByZwdleDogX19tZXrhX2t1YmVybmV0ZXNfc2VydmljZV9sYWJ1bF8oLispCiAg  
LSBzb3VyY2VfbGFizWxz0iBbX19tZXrhX2t1YmVybmV0ZXNfbmFtZXnwYWN1XQogICAgnYWN0aW9u0iByZXBsYWN  
1CiAgICB0YXJnZXrbfGFizWw6IGt1YmVybmV0ZXNfbmFtZXnwYWN1CiAgLSBzb3VyY2VfbGFizWxz0iBbX19tZX  
RhX2t1YmVybmV0ZXNfc2VydmljZV9uYw11XQogICAgnYWN0aW9u0iByZXBsYWN1CiAgICB0YXJnZXrbfGFizWw6I  
Gt1YmVybmV0ZXNfbmFtZQo=  
5 kind: Secret  
6 metadata:  
7   creationTimestamp: 2018-12-20T14:50:35Z  
8   name: additional-configs  
9   namespace: monitoring
```

```
10 resourceVersion: "41814998"
11 selfLink: /api/v1/namespaces/monitoring/secrets/additional-configs
12 uid: 9bbe22c5-0466-11e9-a777-525400db4df7
13 type: Opaque
```

然后我们只需要在声明 prometheus 的资源对象文件中添加上这个额外的配置：(prometheus-prometheus.yaml)

```
1 apiVersion: monitoring.coreos.com/v1
2 kind: Prometheus
3 metadata:
4   labels:
5     prometheus: k8s
6   name: k8s
7   namespace: monitoring
8 spec:
9   alerting:
10    alertmanagers:
11      - name: alertmanager-main
12        namespace: monitoring
13        port: web
14   baseImage: quay.io/prometheus/prometheus
15   nodeSelector:
16     beta.kubernetes.io/os: linux
17   replicas: 2
18   secrets:
19     - etcd-certs
20   resources:
21     requests:
22       memory: 400Mi
23   ruleSelector:
24     matchLabels:
25       prometheus: k8s
26       role: alert-rules
27   securityContext:
28     fsGroup: 2000
29     runAsNonRoot: true
30     runAsUser: 1000
```

```
31 additionalScrapeConfigs:  
32   name: additional-configs  
33   key: prometheus-additional.yaml  
34   serviceAccountName: prometheus-k8s  
35   serviceMonitorNamespaceSelector: {}  
36   serviceMonitorSelector: {}  
37   version: v2.5.0
```

添加完成后，直接更新 prometheus 这个 CRD 资源对象：

```
1 $ kubectl apply -f prometheus-prometheus.yaml  
2 prometheus.monitoring.coreos.com "k8s" configured
```

隔一小会儿，可以前往 Prometheus 的 Dashboard 中查看配置是否生效：

在 Prometheus Dashboard 的配置页面下面我们可以看到已经有了对应的的配置信息了，但是我们切换到 targets 页面下面却并没有发现对应的监控任务，查看 Prometheus 的 Pod 日志：

```
1 $ kubectl logs -f prometheus-k8s-0 prometheus -n monitoring  
2 level=error ts=2018-12-20T15:14:06.772903214Z caller=main.go:240  
component=k8s_client_runtime  
err="github.com/prometheus/prometheus/discovery/kubernetes/kubernetes.go:302: Failed  
to list *v1.Pod: pods is forbidden: User \"system:serviceaccount:monitoring:prometheus-  
k8s\" cannot list pods at the cluster scope"  
3 level=error ts=2018-12-20T15:14:06.773096875Z caller=main.go:240  
component=k8s_client_runtime  
err="github.com/prometheus/prometheus/discovery/kubernetes/kubernetes.go:301: Failed  
to list *v1.Service: services is forbidden: User  
\"system:serviceaccount:monitoring:prometheus-k8s\" cannot list services at the  
cluster scope"  
4 level=error ts=2018-12-20T15:14:06.773212629Z caller=main.go:240  
component=k8s_client_runtime  
err="github.com/prometheus/prometheus/discovery/kubernetes/kubernetes.go:300: Failed  
to list *v1.Endpoints: endpoints is forbidden: User  
\"system:serviceaccount:monitoring:prometheus-k8s\" cannot list endpoints at the  
cluster scope"  
5 .....
```

可以看到有很多错误日志出现，都是 `xxx is forbidden`，这说明是 RBAC 权限的问题，通过 prometheus 资源对象的配置可以知道 Prometheus 绑定了一个名为 prometheus-k8s 的

ServiceAccount 对象，而这个对象绑定的是一个名为 prometheus-k8s 的 ClusterRole：
(prometheus-clusterRole.yaml)

```
1 apiVersion: rbac.authorization.k8s.io/v1
2 kind: ClusterRole
3 metadata:
4   name: prometheus-k8s
5 rules:
6   - apiGroups:
7     - ""
8     resources:
9       - nodes/metrics
10    verbs:
11      - get
12   - nonResourceURLs:
13     - /metrics
14    verbs:
15      - get
```

上面的权限规则中我们可以看到明显没有对 Service 或者 Pod 的 list 权限，所以报错了，要解决这个问题，我们只需要添加上需要的权限即可：

```
1 apiVersion: rbac.authorization.k8s.io/v1
2 kind: ClusterRole
3 metadata:
4   name: prometheus-k8s
5 rules:
6   - apiGroups:
7     - ""
8     resources:
9       - nodes
10      - services
11      - endpoints
12      - pods
13      - nodes/proxy
14    verbs:
15      - get
```

```
16   - list
17   - watch
18 - apiGroups:
19   - ""
20   resources:
21   - configmaps
22   - nodes/metrics
23   verbs:
24   - get
25 - nonResourceURLs:
26   - /metrics
27   verbs:
28   - get
```

更新上面的 ClusterRole 这个资源对象，然后重建下 Prometheus 的所有 Pod，正常就可以看到 targets 页面下面有 kubernetes-service-endpoints 这个监控任务了：

我们这里自动监控了两个 Service，第一个就是我们之前创建的 Redis 的服务，我们在 Redis Service 中有两个特殊的 annotations：

```
1 annotations:
2   prometheus.io/scrape: "true"
3   prometheus.io/port: "9121"
```

所以被自动发现了，当然我们也可以用同样的方式去配置 Pod、Ingress 这些资源对象的自动发现。

数据持久化

上面我们在修改完权限的时候，重启了 Prometheus 的 Pod，如果我们仔细观察的话会发现我们之前采集的数据已经没有了，这是因为我们通过 prometheus 这个 CRD 创建的 Prometheus 并没有做数据的持久化，我们可以直接查看生成的 Prometheus Pod 的挂载情况就清楚了：

```
1 $ kubectl get pod prometheus-k8s-0 -n monitoring -o yaml
2 .....
3   volumeMounts:
4     - mountPath: /etc/prometheus/config_out
5       name: config-out
```

```
6     readOnly: true
7     - mountPath: /prometheus
8       name: prometheus-k8s-db
9     .....
10    volumes:
11    .....
12    - emptyDir: {}
13      name: prometheus-k8s-db
14    .....
```

我们可以看到 Prometheus 的数据目录 /prometheus 实际上是通过 emptyDir 进行挂载的，我们知道 emptyDir 挂载的数据的生命周期和 Pod 生命周期一致的，所以如果 Pod 挂掉了，数据也就丢失了，这也就是为什么我们重建 Pod 后之前的数据就没有了的原因，对应线上的监控数据肯定需要做数据的持久化的，同样的 prometheus 这个 CRD 资源也为我们提供了数据持久化的配置方法，由于我们的 Prometheus 最终是通过 Statefulset 控制器进行部署的，所以我们这里需要通过 storageclass 来做数据持久化，首先创建一个 StorageClass 对象：

```
1 apiVersion: storage.k8s.io/v1
2 kind: StorageClass
3 metadata:
4   name: prometheus-data-db
5 provisioner: fuseim.pri/ifs
```

这里我们声明一个 StorageClass 对象，其中 provisioner=fuseim.pri/ifs，则是因为我们集群中使用的是 nfs 作为存储后端，而前面我们课程中创建的 nfs-client-provisioner 中指定的 PROVISIONER_NAME 就为 fuseim.pri/ifs，这个名字不能随便更改，将该文件保存为 prometheus-storageclass.yaml:

```
1 $ kubectl create -f prometheus-storageclass.yaml
2 storageclass.storage.k8s.io "prometheus-data-db" created
```

然后在 prometheus 的 CRD 资源对象中添加如下配置：

```
1 storage:
2   volumeClaimTemplate:
```

```
3   spec:
4     storageClassName: prometheus-data-db
5   resources:
6     requests:
7       storage: 10Gi
```

注意这里的 storageClassName 名字为上面我们创建的 StorageClass 对象名称，然后更新 prometheus 这个 CRD 资源。更新完成后会自动生成两个 PVC 和 PV 资源对象：

```
1 $ kubectl get pvc -n monitoring
2 NAME                                     STATUS      VOLUME
3   CAPACITY   ACCESS MODES   STORAGECLASS        AGE
4   prometheus-k8s-db-prometheus-k8s-0    Bound      pvc-0cc03d41-047a-11e9-a777-
5   525400db4df7   10Gi        RWO          prometheus-data-db   8m
6   prometheus-k8s-db-prometheus-k8s-1    Bound      pvc-1938de6b-047b-11e9-a777-
7   525400db4df7   10Gi        RWO          prometheus-data-db   1m
8
9 $ kubectl get pv
10 NAME                                    CAPACITY   ACCESS MODES   RECLAIM POLICY
11   STATUS      CLAIM
12   REASON     AGE
13   STORAGECLASS
14   pvc-0cc03d41-047a-11e9-a777-525400db4df7   10Gi        RWO          Delete
15   Bound      monitoring/prometheus-k8s-db-prometheus-k8s-0   prometheus-data-db
16   2m
17   pvc-1938de6b-047b-11e9-a777-525400db4df7   10Gi        RWO          Delete
18   Bound      monitoring/prometheus-k8s-db-prometheus-k8s-1   prometheus-data-db
19   1m
```

现在我们再去看 Prometheus Pod 的数据目录就可以看到是关联到一个 PVC 对象上了。

```
1 $ kubectl get pod prometheus-k8s-0 -n monitoring -o yaml
2 .....
3   volumeMounts:
4     - mountPath: /etc/prometheus/config_out
5       name: config-out
6       readOnly: true
7     - mountPath: /prometheus
8       name: prometheus-k8s-db
9 .....
10  volumes:
11    .....
```

```
12 - name: prometheus-k8s-db
13   persistentVolumeClaim:
14     claimName: prometheus-k8s-db-prometheus-k8s-0
15 . . . . .
```

现在即使我们的 Pod 挂掉了，数据也不会丢失了，最后，下面是我们 Prometheus Operator 系列课程中最终的创建资源清单文件，更多的信息可以在<https://github.com/cnzych/kubernetes-learning> 下面查看。

```
1 apiVersion: monitoring.coreos.com/v1
2 kind: Prometheus
3 metadata:
4   labels:
5     prometheus: k8s
6   name: k8s
7   namespace: monitoring
8 spec:
9   alerting:
10    alertmanagers:
11      - name: alertmanager-main
12        namespace: monitoring
13        port: web
14   storage:
15     volumeClaimTemplate:
16       spec:
17         storageClassName: prometheus-data-db
18         resources:
19           requests:
20             storage: 10Gi
21   baseImage: quay.io/prometheus/prometheus
22   nodeSelector:
23     beta.kubernetes.io/os: linux
24   replicas: 2
25   secrets:
26     - etcd-certs
27   additionalScrapeConfigs:
28     name: additional-configs
29     key: prometheus-additional.yaml
```

```
30 resources:
31   requests:
32     memory: 400Mi
33   ruleSelector:
34     matchLabels:
35       prometheus: k8s
36       role: alert-rules
37   securityContext:
38     fsGroup: 2000
39     runAsNonRoot: true
40     runAsUser: 1000
41   serviceAccountName: prometheus-k8s
42   serviceMonitorNamespaceSelector: {}
43   serviceMonitorSelector: {}
44 version: v2.5.0
```

到这里 Prometheus Operator 系列教程就告一段落了，大家还有什么问题可以到微信群里面继续交流，接下来会和大家介绍 Kubernetes 日志收集方便的知识点