

K8s StatefulSets 部署有状态应用

AUTHOR: 彭玲 TIME: 2022/2/24

K8s StatefulSets 部署有状态应用

StatefulSets 工作负载资源

StatefulSet 示例

前置条件

yaml 配置清单

创建

验证

StatefulSet 中的 Pod

检查 Pod 的顺序索引

使用稳定的网络身份标识

写入稳定的存储

动态提供 PV

Nginx 服务器文件

StatefulSet 扩容/缩容

扩容

缩容

清理工作

1. 顺序终止 Pod

2. 删除 PVC 和 PV

Pod 管理策略

OrderedReady Pod 管理

Parallel Pod 管理

StatefulSets 工作负载资源

StatefulSet 是用来管理有状态应用的工作负载 API 对象。

StatefulSet 用来管理某 Pod 集合的部署和扩缩，并保证这些 pod 的顺序和唯一性。StatefulSet 为这些 Pod 提供持久存储和持久标识符。StatefulSets 对于需要满足以下一个或多个需求的应用程序很有价值：

- 稳定的、唯一的网络标识符。
- 稳定的、持久的存储。
- 有序的、优雅的部署和缩放。
- 有序的、自动的滚动更新。

StatefulSet 示例

对于一个拥有 n 个副本的 StatefulSet，Pod 被部署时是按照 $\{0..n-1\}$ 的序号顺序创建的。

前置条件

确认集群中默认的 StorageClass 是否 **由 PersistentVolume Provisioner 提供** or **管理员 预先提供**:

```
1 anxin@node38:~$ kubectl get sc
2 NAME                                PROVISIONER                                RECLAIMPOLICY  VOLUMEBINDINGMODE
3 ALLOWVOLUMEEXPANSION  AGE
4 localhostpath         kubernetes.io/no-provisioner  Retain         Immediate
   false                47d
4 nfs-storage (default) fuseim.pri/ifs                 Retain         Immediate
   false                27h
```

yaml 配置清单

下面例子中，创建了一个 `nginx` 无头服务，用来发布 StatefulSet 对象 `web` 中 Pod 的 IP 地址。

```
1 anxin@node38:~/pengling/k8s/controller-demo$ vi statefulset-web.yaml
2
3 apiVersion: v1
4 kind: Service
5 metadata:
6   name: nginx
7   namespace: julin
8   labels:
9     app: nginx
10 spec:
11   ports:
12     - port: 80
13     name: web
14   clusterIP: None
15   selector:
16     app: nginx
17 ---
18 apiVersion: apps/v1
19 kind: StatefulSet
20 metadata:
21   name: web
22   namespace: julin
23 spec:
24   serviceName: "nginx"
25   replicas: 2
26   selector:
27     matchLabels:
28       app: nginx
29   template:
30     metadata:
31       labels:
32         app: nginx
33     spec:
34       containers:
35         - name: nginx
36           image: nginx
37           ports:
38             - containerPort: 80
39             name: web
40           volumeMounts: # 保证了 /usr/share/nginx/html 文件夹由一个 PV 支持
```

```

41     - name: www
42       mountPath: /usr/share/nginx/html
43   volumeClaimTemplates:
44   - metadata: # 未指定 julin 命名空间
45     name: www
46     spec:
47       accessModes: [ "ReadWriteOnce" ]
48       resources:
49         requests:
50           storage: 1Gi

```

创建

下面打开两个终端窗口。

- 在第一个终端中，查看 StatefulSet 的 Pods 的创建情况。

```

1 anxin@node38:~$ kubectl get pods -w -l app=nginx -n julin
2
3 NAME          READY   STATUS              RESTARTS   AGE
4 web-0         0/1    Pending            0          0s # web-0
5 web-0         0/1    Pending            0          0s
6 web-0         0/1    Pending            0          1s
7 web-0         0/1    ContainerCreating  0          2s
8 web-0         1/1    Running            0          10s
9 web-1         0/1    Pending            0          0s # web-1
10 web-1         0/1    Pending            0          0s
11 web-1         0/1    Pending            0          2s
12 web-1         0/1    ContainerCreating  0          2s
13 web-1         1/1    Running            0          12s

```

- 在另一个终端中，创建定义在 `web.yaml` 中的 Headless Service 和 StatefulSet。

```

1 anxin@node38:~/pengling/k8s/controller-demo$ kubectl apply -f statefulset-
  web.yaml
2 service/nginx created
3 statefulset.apps/web created

```

验证

获取 `nginx` Service 和 `web` StatefulSet 来验证是否成功的创建了它们。

```

1 anxin@node38:~$ kubectl get service nginx -n julin
2 NAME      TYPE          CLUSTER-IP   EXTERNAL-IP   PORT(S)    AGE
3 nginx     ClusterIP    None         <none>        80/TCP     3m42s

```

获取 `web` StatefulSet，以验证两者均已成功创建：

```

1 anxin@node38:~$ kubectl get statefulset web -n julin
2 NAME      READY   AGE
3 web       2/2     3m46s

```

注意：

在 `web-0` Pod 处于 `Running` 和 `Ready` 状态后 `web-1` Pod 才会被启动。

例如，该示例中应该创建 2 个 Pod，但下面的结果显示：`web-0` 为 `Pending` 异常状态，所以 `web-1` 未被创建。

```
1 anxin@node38:~$ kubectl get po -l app=nginx
2 NAME      READY   STATUS    RESTARTS   AGE
3 web-0     0/1    Pending  0           23m
```

在 `web-0` Pod 正常启动后，`web-1` Pod 会被启动：

```
1 anxin@node38:~$ kubectl get po -l app=nginx -n julin
2 NAME      READY   STATUS    RESTARTS   AGE
3 web-0     1/1    Running   0           4m43s
4 web-1     1/1    Running   0           4m31s
```

StatefulSet 中的 Pod

检查 Pod 的顺序索引

```
1 anxin@node38:~$ kubectl get pods -l app=nginx -n julin
2 NAME      READY   STATUS    RESTARTS   AGE
3 web-0     1/1    Running   0           12m
4 web-1     1/1    Running   0           11m
```

使用稳定的网络身份标识

每个 Pod 都拥有一个基于其顺序索引的稳定的主机名。使用 `kubectl exec` 在每个 Pod 中执行 `hostname` 命令。

```
1 anxin@node38:~$ for i in 0 1; do kubectl exec "web-$i" -n julin -- sh -c
  'hostname'; done
2 web-0
3 web-1
```

使用 `kubectl run` 运行一个提供 `nslookup` 命令的容器。通过对 Pod 的主机名执行 `nslookup`，你可以检查他们在集群内部的 DNS 地址。

```
1 kubectl run -i --tty --image busybox:1.28 dns-test --restart=Never --rm
```

这将启动一个新的 shell。在新 shell 中，运行：

```
1 nslookup web-0.nginx # <pod_name>.<StatefulSet.spec.serviceName>.<namespace>
```

Headless Service 的 CNAME 指向 SRV 记录（记录每个 Running 和 Ready 状态的 Pod）。SRV 记录指向一个包含 Pod IP 地址的记录表项。完整过程如下：

```
1 anxin@node38:~$ kubectl run -it --image busybox:1.28 dns-test --
  restart=Never --rm
2 If you don't see a command prompt, try pressing enter.
3 / #
4 / # nslookup web-0.nginx.julin
5 Server:      10.96.0.10
```

```

6 Address 1: 10.96.0.10 kube-dns.kube-system.svc.cluster.local
7
8 Name:      web-0.nginx.julin
9 Address 1: 10.244.1.52 web-0.nginx.julin.svc.cluster.local
10 / #
11 / # nslookup web-1.nginx.julin
12 Server:   10.96.0.10
13 Address 1: 10.96.0.10 kube-dns.kube-system.svc.cluster.local
14
15 Name:      web-1.nginx.julin
16 Address 1: 10.244.2.44 web-1.nginx.julin.svc.cluster.local

```

写入稳定的存储

动态提供 PV

获取 `web-0` 和 `web-1` 的 PersistentVolumeClaims。

```

1 anxin@node38:~$ kubectl get pvc -l app=nginx -n julin
2 NAME          STATUS    VOLUME                                     CAPACITY  ACCESS MODES
3 STORAGECLASS  AGE
4 www-web-0     Bound    pvc-7f3d201c-847f-4272-ae52-3c4132bf00a8  1Gi       RWO
5               nfs-storage  48m
6 www-web-1     Bound    pvc-ccf60a35-1a29-476c-a722-bdc16f4fb249  1Gi       RWO
7               nfs-storage  47m

```

StatefulSet 控制器创建了 2 个 PVC，绑定到 2 个 PersistentVolumes (动态提供 PV，所有的 PV 都是自动创建和绑定的)。

```

1 anxin@node38:~$ kubectl get pv
2 NAME          RECLAIM    POLICY    STATUS
3 pvc-7f3d201c-847f-4272-ae52-3c4132bf00a8  Bound      julin/www-web-0  nfs-
4 storage
5 pvc-ccf60a35-1a29-476c-a722-bdc16f4fb249  Bound      julin/www-web-1  nfs-
6 storage

```

Nginx 服务器文件

NGINX web 服务器默认会加载位于 `/usr/share/nginx/html/index.html` 的 index 文件。

所有 web 服务器始终使用它们的主机名提供服务。

即使 `web-0` 和 `web-1` 被重新调度的了，但它们仍然继续监听各自的主机名，因为 PV 被重新挂载到了各自的 `volumeMount` 上。不管 `web-0` 和 `web-1` 被调度到了哪个节点上，它们的 PV 将会被挂载到合适的挂载点上。

```

1 # 1. 将 Pod 的主机名写入它们的 `index.html` 文件
2 anxin@node38:~$ for i in 0 1; do kubectl exec "web-$i" -n julin -- sh -c
3 'echo "$(hostname)" > /usr/share/nginx/html/index.html'; done
4 # 2. 验证 NGINX web 服务器使用该主机名提供服务
5 anxin@node38:~$ for i in 0 1; do kubectl exec -it "web-$i" -n julin -- curl
6 http://localhost/; done
7 web-0
8 web-1

```

StatefulSet 扩容/缩容

扩容/缩容 StatefulSet 指增加或减少它的副本数。这通过更新 `replicas` 字段完成。你可以使用 `kubectl scale` 或者 `kubectl patch` 来扩容/缩容一个 StatefulSet。

扩容

使用 `kubectl scale` 扩展副本数为 5。

```
1 anxin@node38:~$ kubectl scale sts web --replicas=5 -n julin
2 statefulset.apps/web scaled
```

缩容

使用 `kubectl patch` 将 StatefulSet 缩容回 3 个副本。

```
1 anxin@node38:~$ kubectl patch sts web -p '{"spec":{"replicas":3}}' -n julin
2 statefulset.apps/web patched
```

清理工作

1. 顺序终止 Pod

控制器会按照与 Pod 序号索引相反的顺序每次删除一个 Pod。在删除下一个 Pod 前会等待上一个被完全关闭。

2. 删除 PVC 和 PV

当删除 StatefulSet 的 Pod 时，挂载到 StatefulSet 的 Pod 的 PersistentVolumes 不会被删除。

这种删除行为由 StatefulSet 缩容引起时，也是一样的。

Pod 管理策略

对于某些分布式系统来说，StatefulSet 的顺序性保证是不必要和/或者不应该的。这些系统仅仅要求唯一性和身份标志。

为了解决这个问题，在 Kubernetes 1.7 中我们针对 StatefulSet API 对象引入了

`.spec.podManagementPolicy`。

`.spec.podManagementPolicy` 选项仅影响扩缩操作的行为，更新不受影响。

OrderedReady Pod 管理

`OrderedReady` pod 管理策略是 StatefulSets 的默认选项。它告诉 StatefulSet 控制器遵循上文展示的 `顺序性` 保证。

Parallel Pod 管理

Parallel pod 管理策略告诉 StatefulSet 控制器 并行 终止所有 Pod，在启动或终止另一个 Pod 前，不必等待这些 Pod 变成 Running 和 Ready 或者完全终止状态。例如，下面的例子：

```
1  apiVersion: v1
2  kind: Service
3  metadata:
4    name: nginx
5    ...
6  ---
7  apiVersion: apps/v1
8  kind: StatefulSet
9  metadata:
10   name: web
11   namespace: julin
12  spec:
13   serviceName: "nginx"
14   podManagementPolicy: "Parallel" # Pod 管理策略 (The default policy is
   `OrderedReady`)
15   ...
```