

# K8s 污点和容忍度

AUTHOR: 彭玲 TIME: 2022/5/12

## K8s 污点和容忍度

污点和容忍度

给 node 加污点

定义 Pod 的容忍度

效果 (effect)

Pod 示例

应用案例

污点 (Taints)

容忍度 (Tolerations) 配置文件

## 污点和容忍度

**节点亲和性** 是 Pod 的一种属性，它使 Pod 被吸引到一类特定的节点（这可能出于一种偏好，也可能是硬性要求）。**污点** (Taint) 则相反——它使节点能够排斥一类特定的 Pod。

### 给 node 加污点

使用命令 `kubectl taint` 给节点增加一个污点。比如，给节点 `node1` 增加一个污点，它的键名是 `key1`，键值是 `value1`，效果是 `NoSchedule`。这表示只有拥有和这个污点相匹配的容忍度的 Pod 才能够被分配到 `node1` 这个节点。

```
1 | kubectl taint nodes node1 key1=value1:NoSchedule # NoSchedule 为 effect(效果)
```

### 定义 Pod 的容忍度

您可以在 PodSpec 中定义 Pod 的容忍度。Pod 的容忍度 (tolerations) 配置中，`operator` 的默认值是 `Equal`。

一个容忍度和一个污点相“匹配”是指它们有一样的键名和效果，并且：

- 如果 `operator` 是 `Exists`（此时容忍度不能指定 `value`），或者
- 如果 `operator` 是 `Equal`，则它们的 `value` 应该相等。

存在两种特殊情况：

- 如果一个容忍度的 `key` 为空且 `operator` 为 `Exists`，表示这个容忍度与任意的 `key`、`value` 和 `effect` 都匹配，即这个容忍度能容忍任意 taint。
- 如果 `effect` 为空，则可以与所有键名 `key1` 的效果相匹配。

下面两个容忍度均与上面例子中使用 `kubectl taint` 命令创建的污点相匹配，因此如果一个 Pod 拥有其中的任何一个容忍度都能够被分配到 `node1`：

```
1 tolerations:
2 - key: "key1"
3   operator: "Equal"
4   value: "value1"
5   effect: "NoSchedule"
6
7 # 或者
8 tolerations:
9 - key: "key1"
10  operator: "Exists"
11  effect: "NoSchedule"
```

## 效果 (effect)

K8s 处理多个污点和容忍度的过程就像一个过滤器：

从一个节点的所有污点开始遍历，过滤掉那些 Pod 中存在与之相匹配的容忍度的污点。余下未被过滤的污点的 effect 值决定了 Pod 是否会被分配到该节点，特别是以下情况：

- 如果未被过滤的污点中存在至少一个 effect 值为 `NoSchedule` 的污点，则 Kubernetes 不会将 Pod 分配到该节点。
- 如果未被过滤的污点中不存在 effect 值为 `NoSchedule` 的污点，但是存在 effect 值为 `PreferNoSchedule` 的污点，则 Kubernetes 会尝试不将 Pod 分配到该节点。
- 如果未被过滤的污点中存在至少一个 effect 值为 `NoExecute` 的污点，则 Kubernetes 不会将 Pod 分配到该节点（如果 Pod 还未在节点上运行），或者将 Pod 从该节点驱逐（如果 Pod 已经在节点上运行）。

## Pod 示例

下面是一个使用了容忍度的 Pod：

```
1 # pods/pod-with-toleration.yaml
2 apiVersion: v1
3 kind: Pod
4 metadata:
5   name: nginx
6   labels:
7     env: test
8 spec:
9   containers:
10  - name: nginx
11    image: nginx
12    imagePullPolicy: IfNotPresent
13  tolerations: # 容忍度
14  - key: "example-key"
15    operator: "Exists"
16    effect: "NoSchedule"
```

## 应用案例

---

# 污点 (Taints)

```
anxinyun@anxinyun-m1:~$ kubectl describe node jump-01 | grep -A 5 'Taints'
Taints:
  nginx=true:NoSchedule
  Unschedulable:
    false
Lease:
  HolderIdentity: jump-01
  AcquireTime: <unset>
  RenewTime: Mon, 09 May 2022 17:03:35 +0800
anxinyun@anxinyun-m1:~$ kubectl describe node jump-02 | grep -A 5 'Taints'
Taints:
  nginx=true:NoSchedule
  Unschedulable:
    false
Lease:
  HolderIdentity: jump-02
  AcquireTime: <unset>
  RenewTime: Mon, 09 May 2022 17:03:42 +0800
anxinyun@anxinyun-m1:~$ kubectl describe node jump-03 | grep -A 5 'Taints'
Taints:
  nginx=true:NoSchedule
  Unschedulable:
    false
Lease:
  HolderIdentity: jump-03
  AcquireTime: <unset>
  RenewTime: Mon, 09 May 2022 17:03:51 +0800
anxinyun@anxinyun-m1:~$
anxinyun@anxinyun-m1:~$
anxinyun@anxinyun-m1:~$ kubectl get po -A | grep jump
anxinyun@anxinyun-m1:~$ kubectl get po -A -o wide | grep jump
ingress-nginx      ingress-nginx-controller-vh7k5      1/1      Running      0      30d      10.8.40.120      jump-01      <none>      <none>
kube-system        kube-flannel-ds-amd64-69529        1/1      Running      0      30d      10.8.40.120      jump-01      <none>      <none>
kube-system        kube-flannel-ds-amd64-mxqg8        1/1      Running      0      30d      10.8.40.124      jump-02      <none>      <none>
kube-system        kube-flannel-ds-amd64-x8jzs        1/1      Running      0      30d      10.8.40.124      jump-03      <none>      <none>
kube-system        kube-proxy-nxptt                    8      412d      10.8.40.124      jump-03      <none>      <none>
kube-system        kube-proxy-qwmn9                    8      412d      10.8.40.123      jump-02      <none>      <none>
kube-system        kube-proxy-wxdg                      1/1      Running      11     480d      10.8.40.120      jump-01      <none>      <none>
kubernetes-monitoring-system  node-exporter-mbkpp                 2/2      Running      12     480d      10.8.40.120      jump-01      <none>      <none>
kubernetes-monitoring-system  node-exporter-mln6                  2/2      Running      8      412d      10.8.40.123      jump-02      <none>      <none>
kubernetes-monitoring-system  node-exporter-zfkxc                 2/2      Running      10     412d      10.8.40.124      jump-03      <none>      <none>
ops                            emqx-66879d5655-zn4qg              1/1      Running      0      30d      10.96.13.23      jump-01      <none>      <none>
ops                            emqx-broker-0                       0/1      Terminating 327    161d     10.96.14.109     jump-02      <none>      <none>
ops                            emqx-broker-1                       0/1      Terminating 1      160d     10.96.12.13      jump-03      <none>      <none>
```

# 容忍度 (Tolerations) 配置文件

```
1 kind: DaemonSet
2 apiVersion: apps/v1
3 metadata:
4   name: ingress-nginx-controller
5   namespace: ingress-nginx
6   labels:
7     app.kubernetes.io/component: controller
8     app.kubernetes.io/instance: ingress-nginx
9     app.kubernetes.io/name: ingress-nginx
10 spec:
11   selector:
12     matchLabels:
13       app.kubernetes.io/component: controller
14       app.kubernetes.io/instance: ingress-nginx
15       app.kubernetes.io/name: ingress-nginx
16   template:
17     metadata:
18       creationTimestamp: null
19     labels:
20       app.kubernetes.io/component: controller
21       app.kubernetes.io/instance: ingress-nginx
22       app.kubernetes.io/name: ingress-nginx
23   spec:
24     containers:
25       - name: controller
26         image: 'k8s.gcr.io/ingress-nginx/controller:v0.40.2'
27         args:
28           - /nginx-ingress-controller
29           - '--election-id=ingress-controller-leader'
30           - '--ingress-class=nginx'
31           - '--configmap=$(POD_NAMESPACE)/ingress-nginx-controller'
32           - '--tcp-services-configmap=$(POD_NAMESPACE)/tcp-services'
33           - '--udp-services-configmap=$(POD_NAMESPACE)/udp-services'
34           - '--validating-webhook=:8443'
35           - '--validating-webhook-
certificate=/usr/local/certificates/cert '
36           - '--validating-webhook-key=/usr/local/certificates/key'
37     ports:
```

```
38     - name: http
39       hostPort: 80
40       containerPort: 80
41       protocol: TCP
42     - name: https
43       hostPort: 443
44       containerPort: 443
45       protocol: TCP
46     - name: webhook
47       hostPort: 8443
48       containerPort: 8443
49       protocol: TCP
50   env:
51     - name: POD_NAME
52       valueFrom:
53         fieldRef:
54           apiVersion: v1
55           fieldPath: metadata.name
56     - name: POD_NAMESPACE
57       valueFrom:
58         fieldRef:
59           apiVersion: v1
60           fieldPath: metadata.namespace
61     - name: LD_PRELOAD
62       value: /usr/local/lib/libmimalloc.so
63   livenessProbe:
64     httpGet:
65       path: /healthz
66       port: 10254
67       scheme: HTTP
68     initialDelaySeconds: 10
69     timeoutSeconds: 1
70     periodSeconds: 10
71     successThreshold: 1
72     failureThreshold: 5
73   readinessProbe:
74     httpGet:
75       path: /healthz
76       port: 10254
77       scheme: HTTP
78     initialDelaySeconds: 10
79     timeoutSeconds: 1
80     periodSeconds: 10
81     successThreshold: 1
82     failureThreshold: 3
83   lifecycle:
84     preStop:
85       exec:
86         command:
87         - /wait-shutdown
88   securityContext:
89     capabilities:
90       add:
91         - NET_BIND_SERVICE
92       drop:
93         - ALL
94     runAsUser: 101
95     allowPrivilegeEscalation: true
```

```
96     nodeSelector:
97         ingress: nginx
98     hostNetwork: true
99     tolerations: # 容忍度
100         - key: nginx
101           operator: Exists
102
```