

通過基於 DSCP 的網路優先級加速 Kubernetes 環境中的分散式 Ray 應用 Accelerating Distributed Ray Applications in Kubernetes Through DSCP-based Network Prioritization

指導教授：莊坤達

專題成員：林宥呈、張羿軒

開發工具：VSCode, Golang, Kubernetes

測試環境：x86_64 *2, AGX ORIN *2、Switch

EPS121

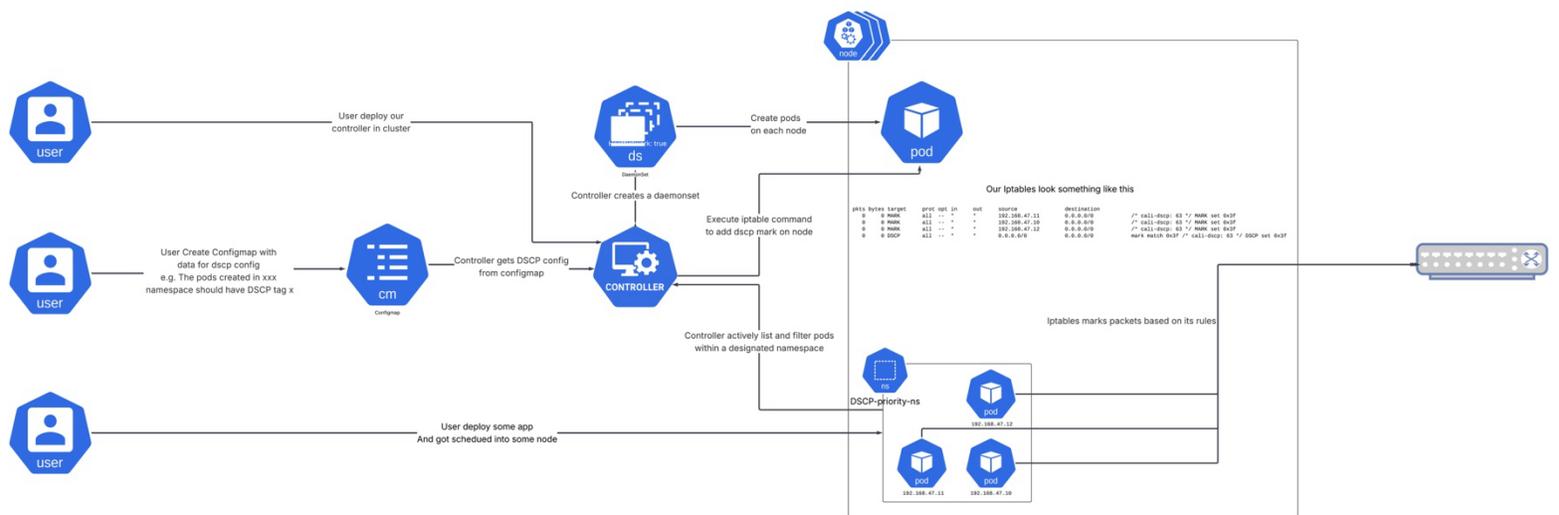
一、簡介：

在分散式機器學習的環境中，由於節點之間會有頻繁的網路傳輸，大量的模型訓練和參數同步過程會產生密集的網路傳輸需求。當多個訓練任務同時進行時，網路壅塞成為一個嚴重的問題，可能導致關鍵數據包被丟棄，影響整體訓練效率和模型收斂性。

為了解決這個問題，我們開發了一個專門的 Kubernetes Controller，能夠基於 namespace 自動為網路封包設置 DSCP (Differentiated Services Code Point) 標記。通過這種方式，我們可以為不同優先級的網路流量分配適當的服務品質 (QoS) 等級。

這個控制器可以在任何支援 DSCP 的網路環境中使用，通過 iptables mangle 表來實現封包標記。使用者可以根據業務需求，為不同的 Kubernetes namespace 配置不同的 DSCP 值，確保高優先級工作負載（如關鍵的模型訓練任務）的網路封包能夠獲得優先處理，降低在網路壅塞時被丟棄的機率並提高該任務效率。

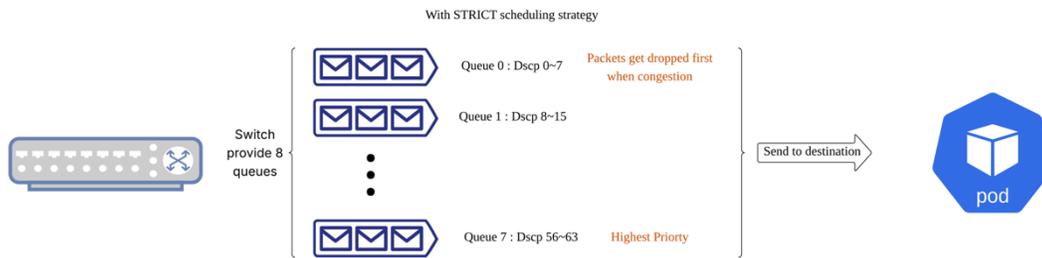
二、系統架構



我們實作了一個 k8s 控制器，當控制器被部署後會在集群內建立一個 daemonset，此 daemonset 負責確保在所有節點上，一定運行著我們的一個 DSCP Pod，當這些 pods 某些因素消失時會將他們自動復原。

透過這些 pods，我們便可以控制各個 node 的 iptables rules，使得我們目標的 app 中的 pod 所產生的封包能夠在被送出 node 之前被加上 DSCP 標籤。

當 Network switch 收到帶有 DSCP 標籤的封包時，會根據預先在 switch 上設定好的優先權數值，決定要放入哪個佇列中等待，而優先權最大的佇列會最先被送出，優先權越低的封包越容易被丟棄。



透過自定義優先權，我們就能夠擁有更高的彈性決定運行在擁擠程度很高的集群中哪種應用最重要，這能顯著提升指定應用在分散式系統中的效率，例如分散式機器學習。

三、測試結果：

我們利用 iperf3 進行測試，另有一台用作模擬 congestion 的機器，負責發送無用的封包佔滿頻寬，下面的測試結果是未使用我們的 controller 時，可以看到重傳封包的次數為 517，頻寬速度為 15.5 Mbits/s。

```
root@iperf3-client-priority:/# iperf3 -c iperf3-server-priority.dscp-priority.svc.cluster.local -p 5202 -t 5 -b 20G
Connecting to host iperf3-server-priority.dscp-priority.svc.cluster.local, port 5202
[ ID] Interval           Transfer             Bitrate             Retr
[  5]  0.00-5.00   sec   9.25 MBytes    15.5 Mbits/sec    517           sender
[  5]  0.00-5.02   sec   6.52 MBytes    10.9 Mbits/sec                receiver
```

以下是有使用我們的 controller，傳出的封包有加上 DSCP 標籤=63，可以看到重傳封包的次數為 184，頻寬速度為 944 Mbits/s。雖然實際機器學習的情況需要再實驗，但我們能看到，有加上 DSCP 標籤的組別速度相差了約 60 倍。

```
root@iperf3-client-priority:/# iperf3 -c iperf3-server-priority.dscp-priority.svc.cluster.local -p 5202 -t 5 -b 20G
Connecting to host iperf3-server-priority.dscp-priority.svc.cluster.local, port 5202
[ ID] Interval           Transfer             Bitrate             Retr
[  5]  0.00-5.00   sec   563 MBytes    944 Mbits/sec    184           sender
[  5]  0.00-5.00   sec   560 MBytes    939 Mbits/sec                receiver
Ethernet42 Last cached time was 2025-05-26T17:30:05.071666
```

另列出 Switch 中佇列的情況，未加上 DSCP 標籤的封包會被放入佇列 UC0 中，有加上 DSCP 標籤=63 的封包會被放入佇列 UC7 中，可以看到佇列 UC0 有相當高的封包丟棄量，而佇列 UC7 是 0，這顯示我們的 DSCP 控制器發揮了效果，使我們的應用獲得了較高的優先權。

Port	TxQ	Counter/pkts	Counter/bytes	Drop/pkts	Drop/bytes
Ethernet42	UC0	110,044,021	162,895,083,608	30,071	44,981,565
Ethernet42	UC7	1,378,088	2,022,900,011	0	0