



KubeCon



CloudNativeCon

Europe 2020

Virtual

What You Didn't Know About Ingress Controllers' Performance

Ismo Puustinen & Mikko Ylinen, Intel

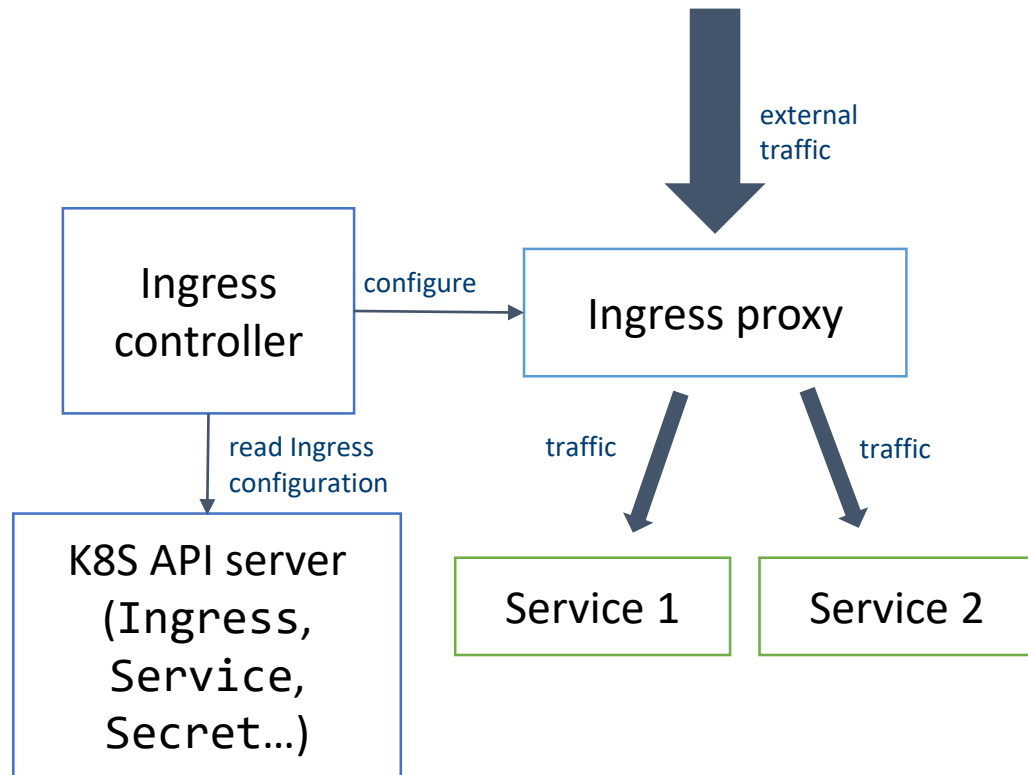
Agenda



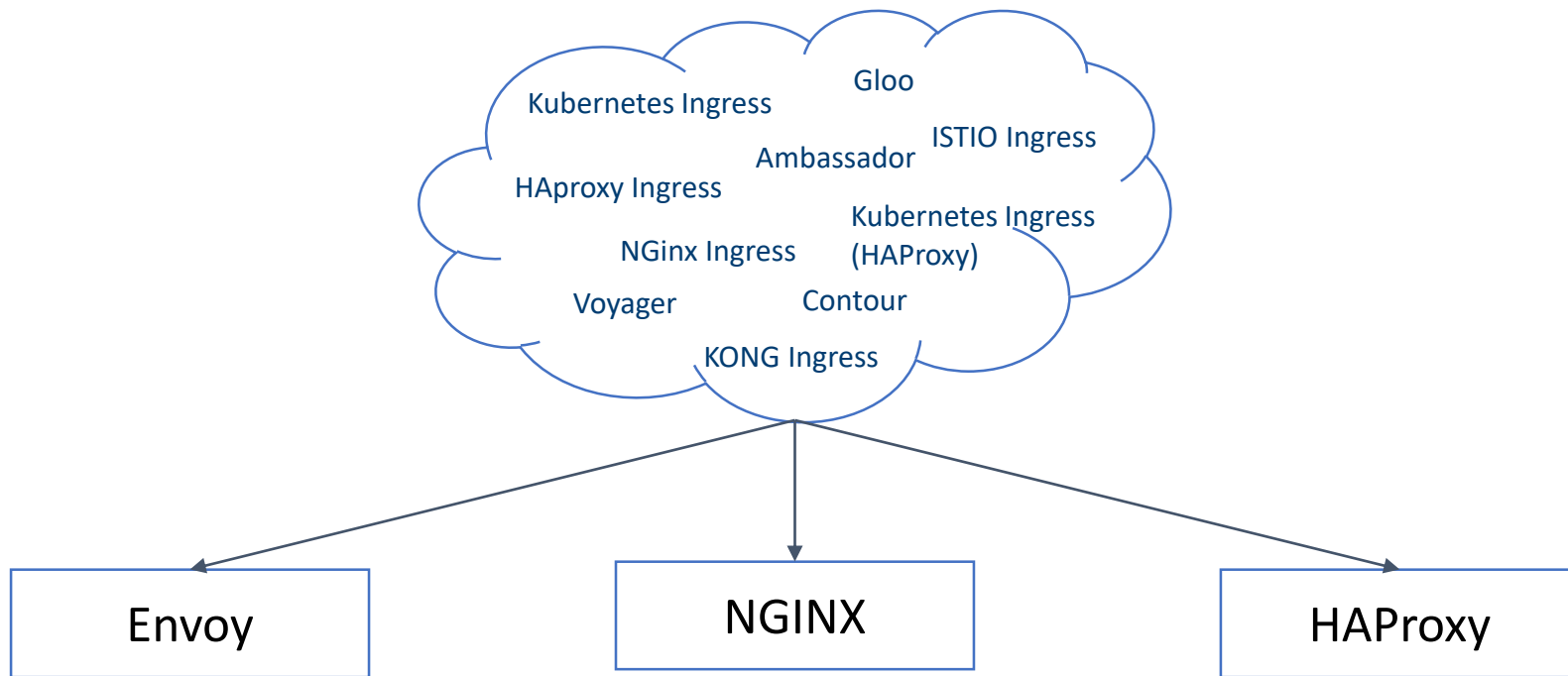
- Ingress overview and performance metrics
- Ingress controller optimization using
 - TLS offloading and asynchronous processing
 - Kubernetes CPU manager and alignment with CPU resources
- Call to Action

- What is Ingress?
 - Ingress exposes HTTP(s) routes to services within the cluster
 - “Route traffic coming to `https://example.com/foo/bar` to service foobar”
 - An Ingress controller is the implementation, at least one is required

- Typical Ingress setup
 - Ingress controller follows the Kubernetes objects and creates the route configuration
 - Ingress proxy reads the configuration and handles the actual traffic routing



Kubernetes Ingress Controllers



<https://kubernetes.io/docs/concepts/services-networking/ingress-controllers/#additional-controllers>

Performance Metrics

- Ingress proxy processing as performance bottleneck
 - **Bandwidth** (requests/sec)
 - **Latency** (ms/request, especially for tail)

```
time="2020-07-09T05:55:47Z" level=info msg="Test finished" i=43932 t=10.000065893s
```

```
data_received.....: 60 MB 6.0 MB/s
data_sent.....: 20 MB 2.0 MB/s
http_req_blocked.....: avg=13.25ms min=1.56ms med=10.65ms max=112.89ms p(90)=25.55ms p(95)=31.43ms
http_req_connecting.....: avg=157.05µs min=45.26µs med=157.05µs max=15.03ms p(90)=206.51µs p(95)=232.21µs
http_req_duration.....: avg=13.6ms min=217.88µs med=11.26ms max=76.48ms p(90)=27.64ms p(95)=33.39ms
http_req_receiving.....: avg=127.50µs min=15.84µs med=100.4µs max=41.77ms p(90)=149.96µs p(95)=170.55µs
http_req_sending.....: avg=137.3µs min=18µs med=66.64µs max=24.05ms p(90)=104.57µs p(95)=153.98µs
http_req_tls_handshaking...: avg=12.72ms min=1.43ms med=10.31ms max=75.83ms p(90)=24.87ms p(95)=30.42ms
http_req_waiting.....: avg=5.54ms min=111.57µs med=10.97ms max=76.33ms p(90)=27.37ms p(95)=33.14ms
http_reqs.....: 4394 4394.171045/s
```

Kubernetes

- Linux scheduling (QoS)
 - CFS, NUMA
- CPU manager policy
- Assignment of HW resources
- Ingress proxy configuration
 - Number of routes
 - Replicas
- Cluster topology (network)

Proxy

- Proxy threads and memory
 - In line with HW resources
 - Buffers
- Request content processing
 - TLS handshake
- Ingress proxy configuration
 - Upstream (cluster) load balancing
 - Timeouts, connection limits, ...

[illegible]

TLS handshake

HTTP parsing

Connection teardown

Data writing

Connection creation

- Two separate parts: handshake (asymmetric crypto) and data transfer (symmetric crypto)
 - Symmetric crypto is pretty fast
 - OpenSSL, BoringSSL, GnuTLS support AES instructions which make it even faster
 - **Asymmetric crypto is slow (ECDHE-RSA cipher suites)**
 - Receiving new HTTPs connections is where the CPU cycles are spent!

Faster TLS Handshakes?

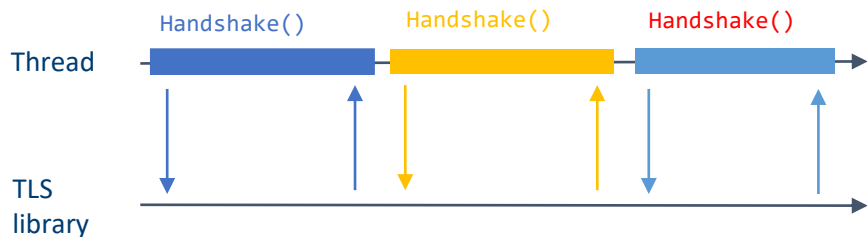
- Crypto parameters (TLS1.3, ECDSA vs. RSA, ECDHE vs. RSA...) – feasible?
- Crypto operations can be accelerated by offloading them to an external accelerator
 - Individual operations or vectorization (wait until several operations are received)
 - Different TLS libraries have different methods for offloading crypto operations to external providers:

OpenSSL	OpenSSL Engines
BoringSSL	Private key operations
GnuTLS	Cryptography provider layer

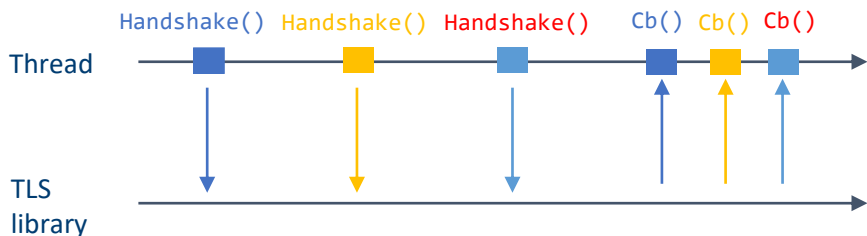
- It's common that accelerating a single crypto operation isn't much faster, but many operations can be done in parallel: **asynchronous handshake processing**

Sync vs. Async TLS

Synchronous TLS



Asynchronous TLS



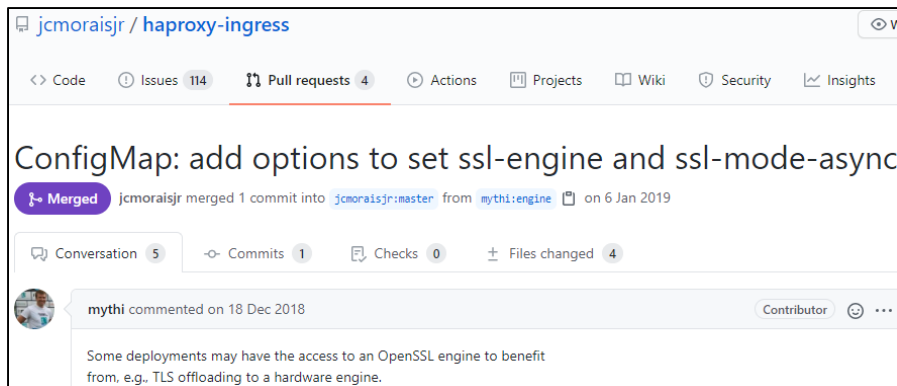
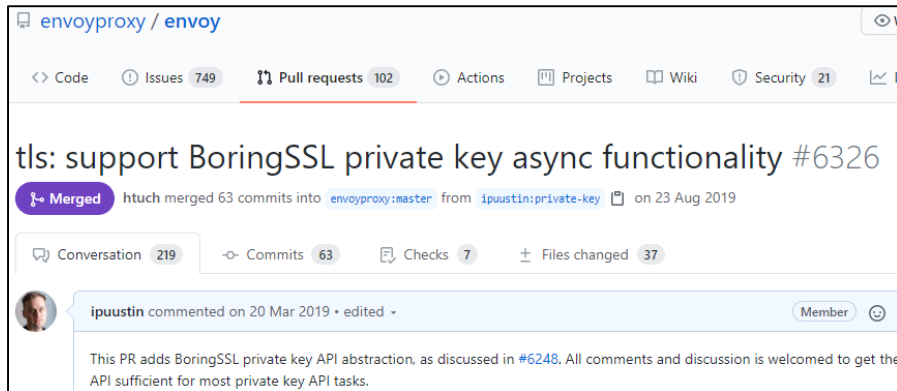
Asynchronous TLS

- Handshake function returns immediately, callback is called when the handshake is done
 - Crypto is accelerated in the background while Ingress proxy server can do other things
- Sub-problems:
 1. How to get the Ingress proxy servers to support asynchronous handshakes?
 2. How to get accelerator resources to Kubernetes Ingress Controllers?

Async TLS Support In Ingress

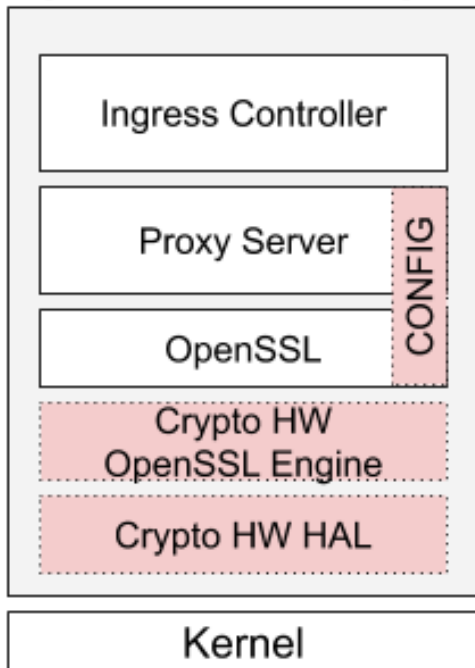


1. Add asynchronous TLS support to underlying proxy servers
2. Add support to Ingress controllers' configuration for getting asynchronous TLS used in the proxy



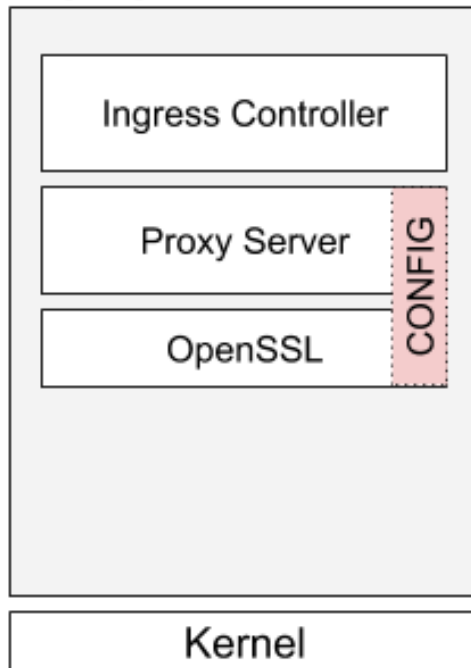
Async TLS Offloading

Ingress Controller w/ Engine



vs.

Orig. Ingress Controller



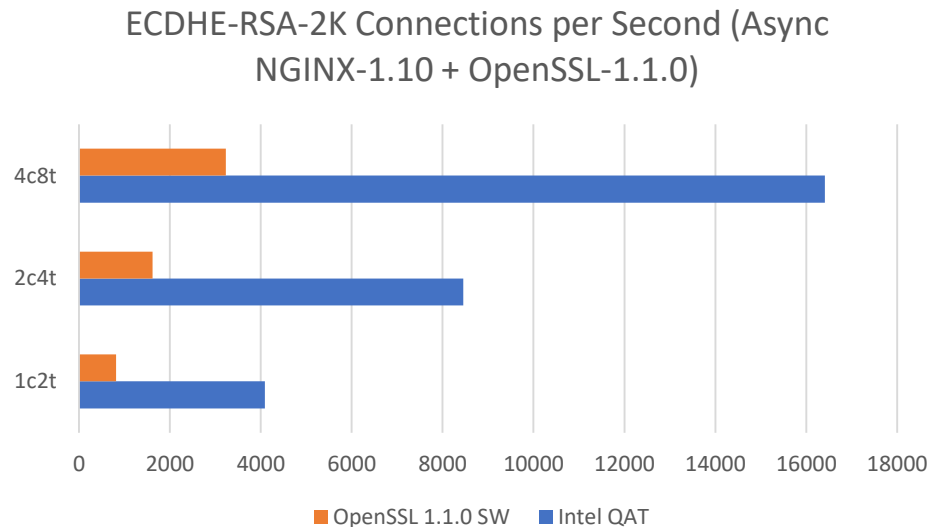
TLS Offloading Status

	NGINX	HAProxy	Envoy	Envoy-OpenSSL
Crypto offload	OpenSSL Engine	OpenSSL Engine	BoringSSL Private Key Provider	OpenSSL Engine
Asynchronous Processing	No*	Handshakes	Handshakes	Handshakes
Rebuild to use Offloading	No	No	Yes	No
Ingress Configuration Supported	(OPENSSL_CONF override)	Haproxy-ingress, Kubernetes (Haproxy) Ingress		(OPENSSL_CONF override)

* `Asynch_mode_nginx`: Handshakes and bulk crypto

When To Offload?

- Identify if the ingress controller proxy is **CPU-bound**
 - Adding more cores increases throughput and/or reduces latency
- Find out if you have a lot of incoming **new** HTTPs connections
 - Upstream connections don't matter, because the TLS connections are reused
- The biggest benefit from the acceleration is when the Ingress proxy is running on a **low CPU core count** (1-8 cores)

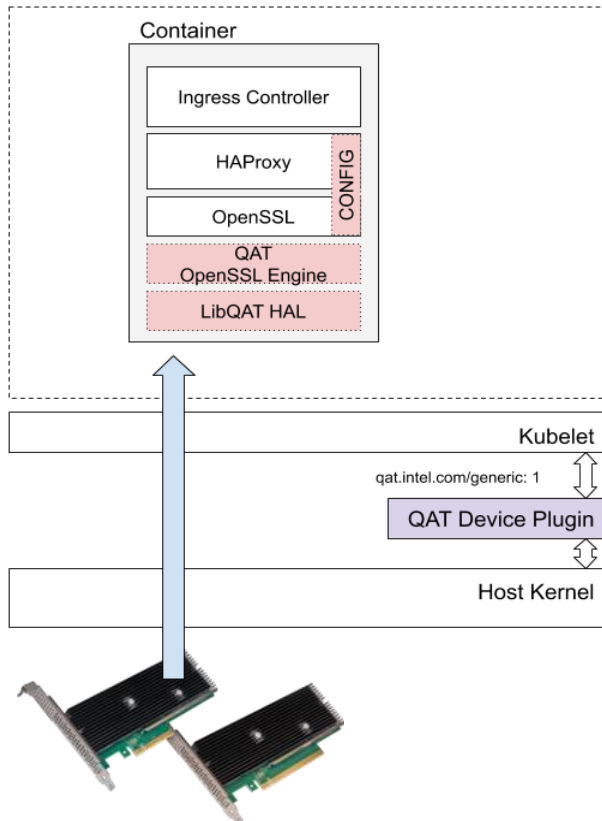


Reference: <https://01.org/sites/default/files/downloads/intelr-quickassist-technology/337003-001-intelquickassisttechnologyandopenssl-110.pdf>

Examples



HAProxy Hardware Acceleration



ConfigMap

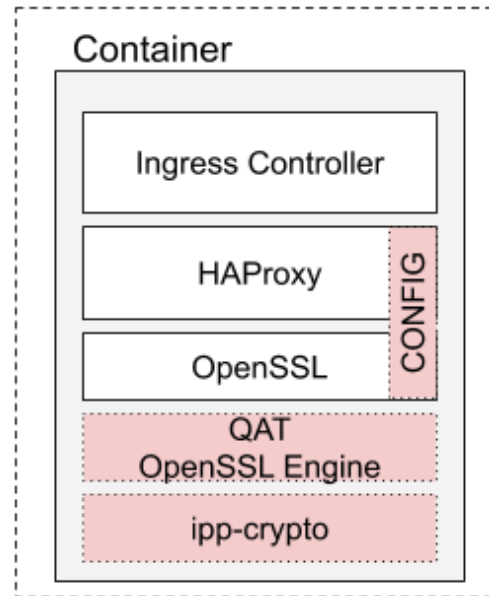
```
config-snippet: |
  ssl-engine qat
  ssl-mode-async
```

Custom Deployment

```
containers:
  ...
  - name: demo-container-1
    image: custom-haproxy-ingress:2.0
    resources:
      limits:
        qat.intel.com/generic: 1
```

HAProxy RSA Multibuffer

- RSA Multibuffer and asynchronous OpenSSL engine
 - queue up-to eight RSA operations and process them in parallel
 - ipp-crypto's multi-buffer
 - AVX512 IFMA (e.g., Ice Lake CPUs) enabled nodes



ConfigMap

```
config-snippet: |  
  ssl-engine qat  
  ssl-mode-async
```

Custom Deployment (node-feature-discovery labeled nodes):

```
...  
spec:  
  template:  
    spec:  
      nodeSelector:  
        feature.node.kubernetes.io/cpu-cpuid.AVX512IFMA: "true"
```

Performance Tuning Areas

Kubernetes

- Linux scheduling (QoS)
 - CFS, NUMA
- CPU manager policy
- Assignment of HW resources
- Ingress proxy configuration
 - Number of routes
 - Replicas
- Cluster topology (network)

Proxy

- Proxy threads and memory
 - In line with HW resources
 - Buffers
- Request content processing
 - TLS handshake
- Ingress proxy configuration
 - Upstream (cluster) load balancing
 - Timeouts, connection limits, ...

Kubernetes CPU Manager



- Ingress proxy performance scales almost linearly with CPU cores
- Kubelet CPU Manager “static” policy can be used to assign full cores for a container in a pod with Guaranteed QoS class
 - “Full cores” means that other containers will not be scheduled to the cpuset belonging to the container. Also, CPU manager selects cores which are topologically close, meaning threads on the same physical core, cores on the same socket, and so on.
- Problem: how to choose the hard limits? What if Ingress controller memory usage increases when changing the Ingress configuration?

QoS class is a pod-level property.

Guaranteed: all containers in the pod need to have both request and limit set to same value for all native resources.

Burstable: at least one container needs to have a native resource request or limit set.

Best-effort: everything else.

- The performance impact depends on other workloads running on the same node
 - Difficult to quantify! Test this in your own environment.
- Performance scales both to larger and smaller workloads
 - Large workloads will get CPU cores with good topology
 - Small workloads benefit from less processor cache misses and immunity to CFS throttling

- The Ingress proxies' resource use can be configured
 - Worker threads, memory buffers, connections per route, ...
- Proxies typically try to autodetect the amount of resources available: for bare metal and virtual machines this works fine
 - Problem: when running in container, the CPU core count auto-detection doesn't necessarily work at all, because the proxies aren't cgroup-aware: the worker thread count is easily much bigger than intended
 - Solution: set the parameters manually or use more advanced configuration values
 - Example: Envoy option `--cpuset-threads` adjusts the worker thread count to match with the Envoy cpuset size and not the number of cores in the system

- Users: analyze your workloads to see where the bottlenecks are
- Ingress controller developers: ensure the following config knobs are there:
 - Easy customization of crypto offload and async processing modes
 - Easy customization of non-native resources
 - Easy customization of nodeAffinity labels
- Ingress proxy developers: switch to async TLS and integrate mechanisms for enabling custom TLS handshake



KubeCon



CloudNativeCon

Europe 2020



...

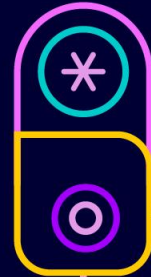


Virtual



KEEP CLOUD NATIVE

CONNECTED



...