

P4 Spring Workshop 2022

Architecture for multi-Terabit programmable networking functions

Petr Kastovsky, petr.kastovsky@intel.com

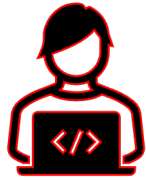
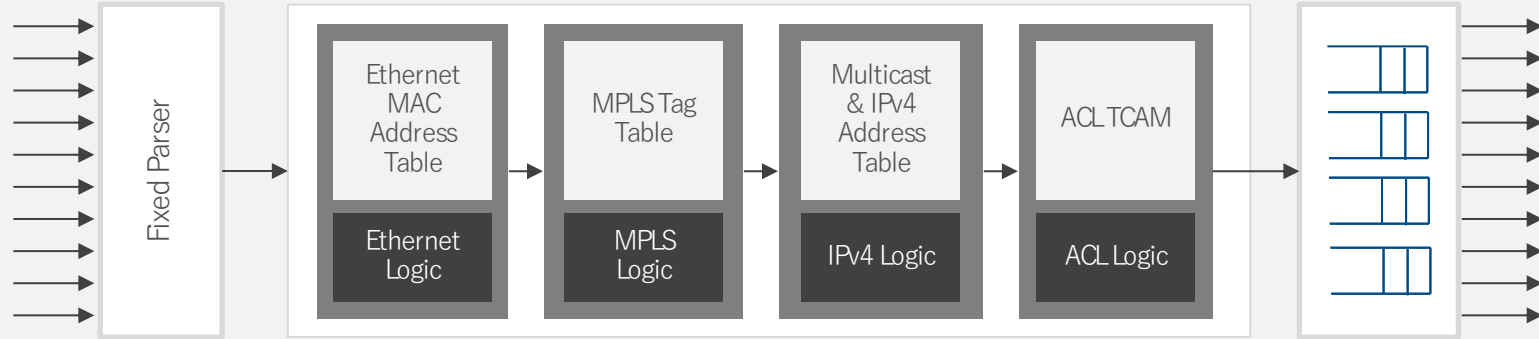
Georgios Nikolaidis, georgios.nikolaidis@intel.com



intel®

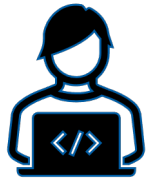
Fixed-function ASIC:

Features and table-sizes are **hard-coded**, not optimized, and often **wasted**.



This is how you process packets ...

New World of Network Design: Software-Defined

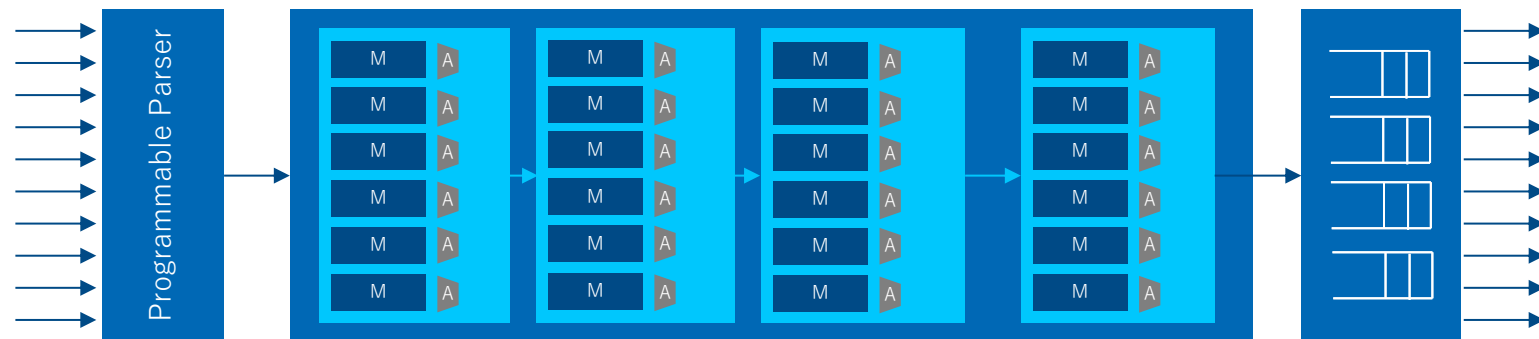


This is how I want to process packets ...

Programmable ASIC:

Software determines headers, table sizes, and packet-processing functions.

Pipeline is fully workload-optimized.

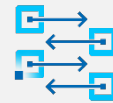


Innovate at the pace of software with new P4 Applications



```
reads {table int_table
}
ip.protocol;
}
actions {
  export_queue_latency;
}
```

```
actionadd_header(int_header);
modify_field(int_header.kind, TCP_OPTION_INT);
modify_field(int_header.len, TCP_OPTION_INT_LEN);
modify_field(int_header.sw_id, sw_id);
modify_field(int_header.q_latency,
  intrinsic_metadata.deq_timedelta);
add_to_field(tcp.dataOffset, 2);
add_to_field(ipv4.totalLen, 8);
subtract_from_field(ingress_metadata.tcpLength,
  12);
}
export_queue_latency (sw_id) {
```



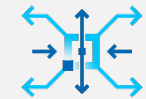
Enhanced switching



User plane functions



Physical to virtual



Broadband network gateway (BNG)



Security, DDoS detection



L4 load balancing



Tunnel gateways



Network packet broker (NPB)



Real-time telemetry



DNS caching



Machine learning



5G gateway

Alibaba Cloud Gateway XGW



- **Business challenge:** network traffic during a shopping festival can draw millions of shopper visits creating terabits per second of data traffic
- Scale and growth rate of Alibaba Cloud today makes horizontal scaling (adding servers) unsustainable
 - Increased CAPEX (hardware acquisition costs) and OPEX (maintenance, management, troubleshooting)
 - Heavy-hitter flows overloading single CPU core, need for headroom further increasing CAPEX & OPEX
- Technical challenge:
 - Limited total on-chip memory capacity storing the $O(1M)$ VXLAN routing table and $O(1M)$ VM-NC mapping table **while stateful tables and ultra-large in table entries, such as SNAT require $O(100M)$ entries**

Source: Tian Pan, Nianbing Yu, Chenhao Jia, Jianwen Pi, Liang Xu, Yisong Qiao, Zhiguo Li, Kun Liu, Jie Lu, Jianyuan Lu, Enge Song, Jiao Zhang, Tao Huang, Shunmin Zhu. 2021. Sailfish: Accelerating Cloud-Scale Multi-Tenant Multi-Service Gateways with Programmable Switches. In ACM SIGCOMM 2021 Conference (SIGCOMM '21), August 23–28, 2021, Virtual Event, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3452296.3472889>

Tofino Expanded Architecture: Intro

Complementing Tofino by FPGAs to enable 100x increase in table and buffer capacity

eXtra large tables

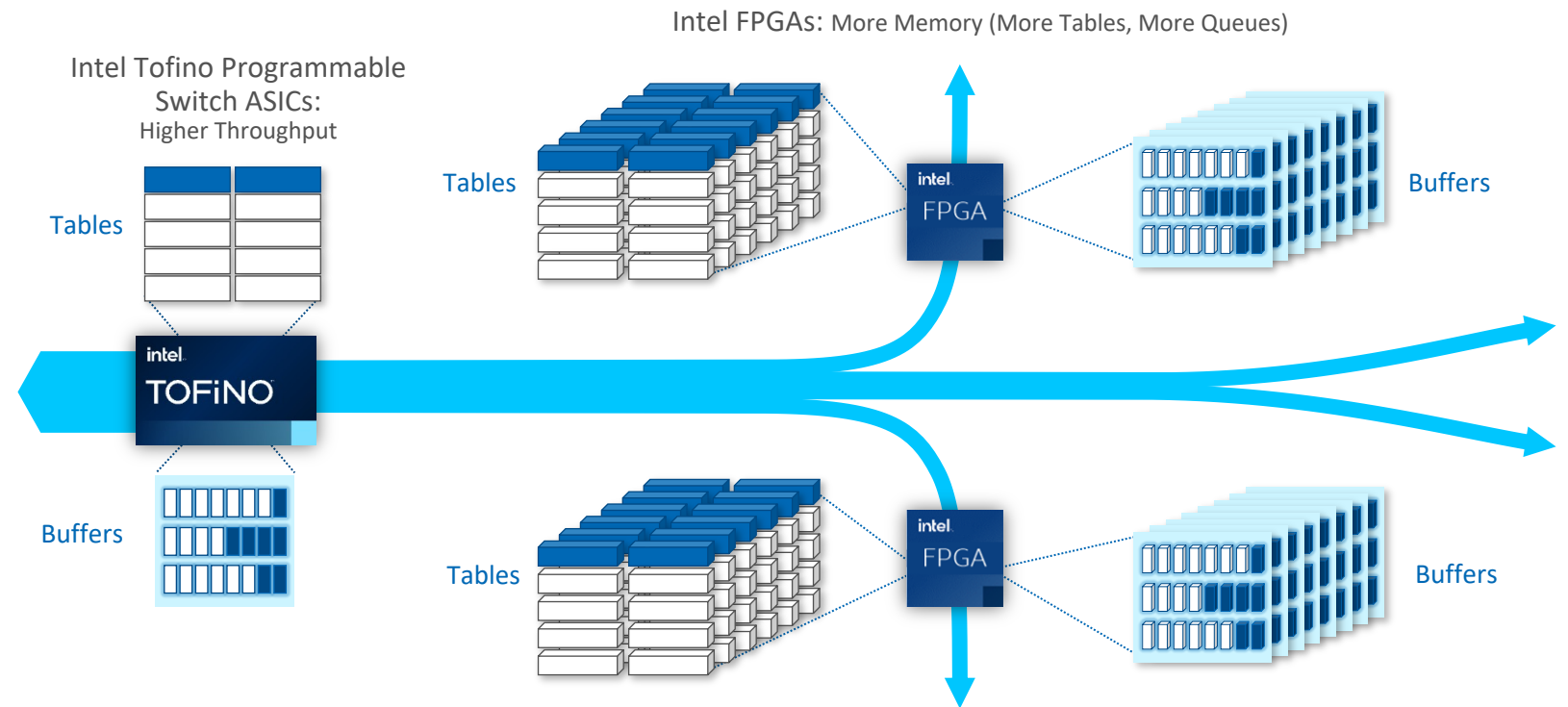
up to 100s of millions of entries

- CSP: cloud gateway (L4 LB, firewall, VxLAN, NAT)
- CoSP: carrier grade NAT, NPB, IPv6 NAT, 5G metro router etc.

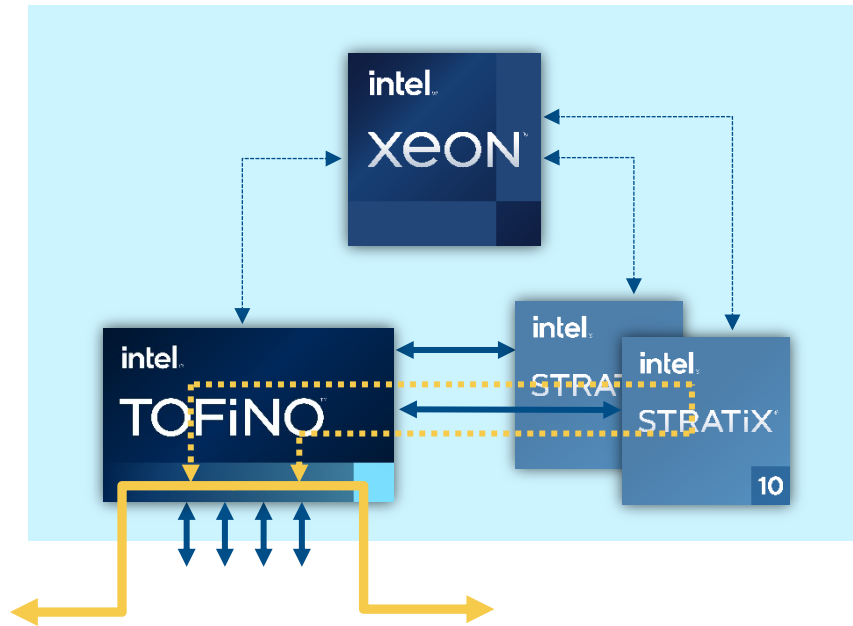
eXtra large buffers

up to 10s of GBs of buffers

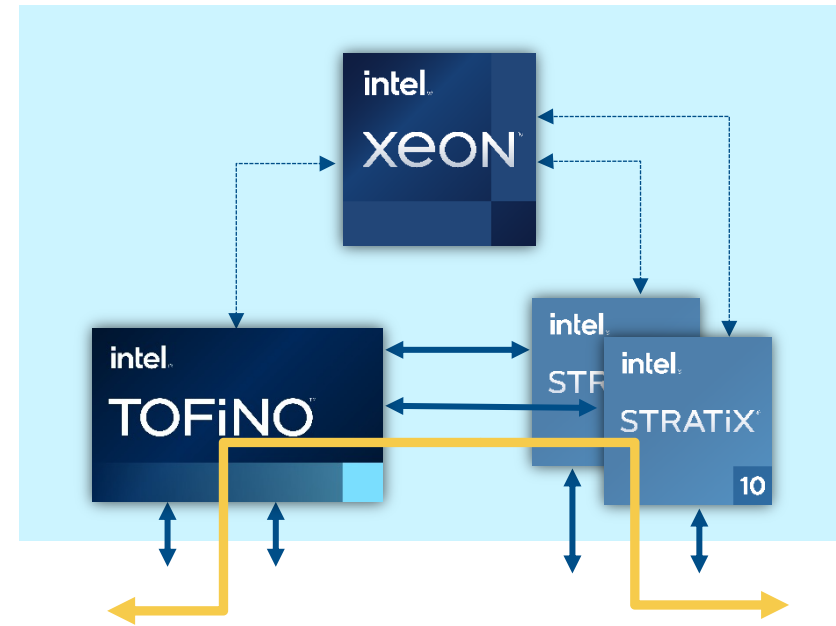
- CoSP: telco gateway BNG/5G UPF/AGF, 5G metro router, NFV acceleration



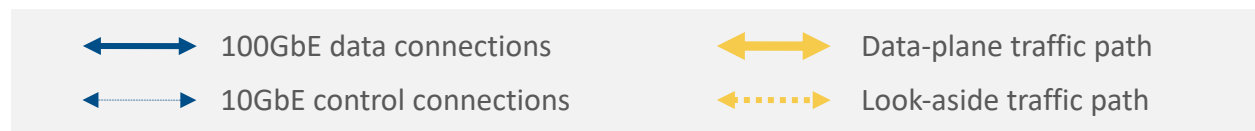
Architecture Implementations



FPGA Look-Aside to Switch ASIC



FPGA Inline with Switch ASIC



Hardware Form factors

Switch + FPGA SmartNIC

- Tofino-based switch
- FPGA SmartNIC cards in a separate server
 - Intel® FPGA PAC N3000
 - N5010 (LC)



Server switch

- Integrated platform
- Tofino, FPGA, and CPU in one box

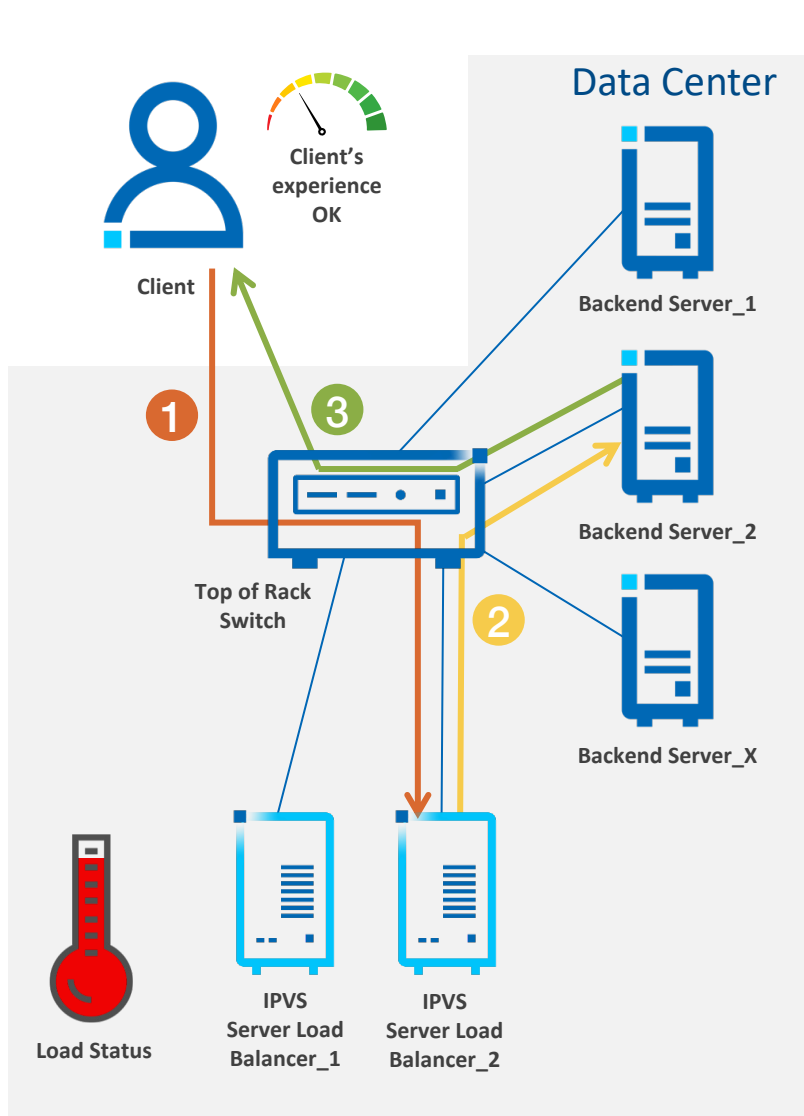


Chassis platform

- Modules with Tofino, FPGA, and CPU



Efficient High-Bandwidth Load Balancing



BENEFITS

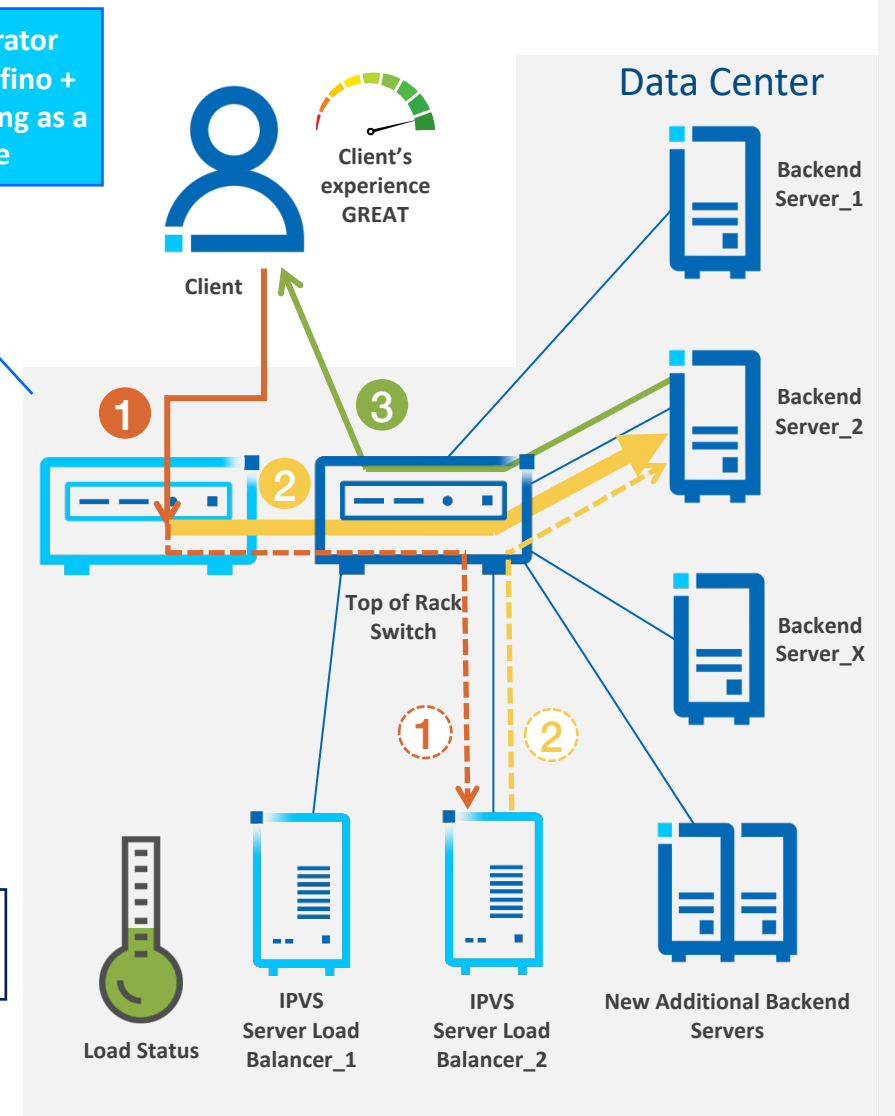
- IPVS less utilized so available for other tasks
- Latency goes down
- Revenue goes up



Server L4 LB Accelerator (SLBA) built on the Tofino + FPGA architecture acting as a connection cache

- 1 Client's request
- 2 Client's request load balanced to the right backend server
- 3 Content delivered to the client

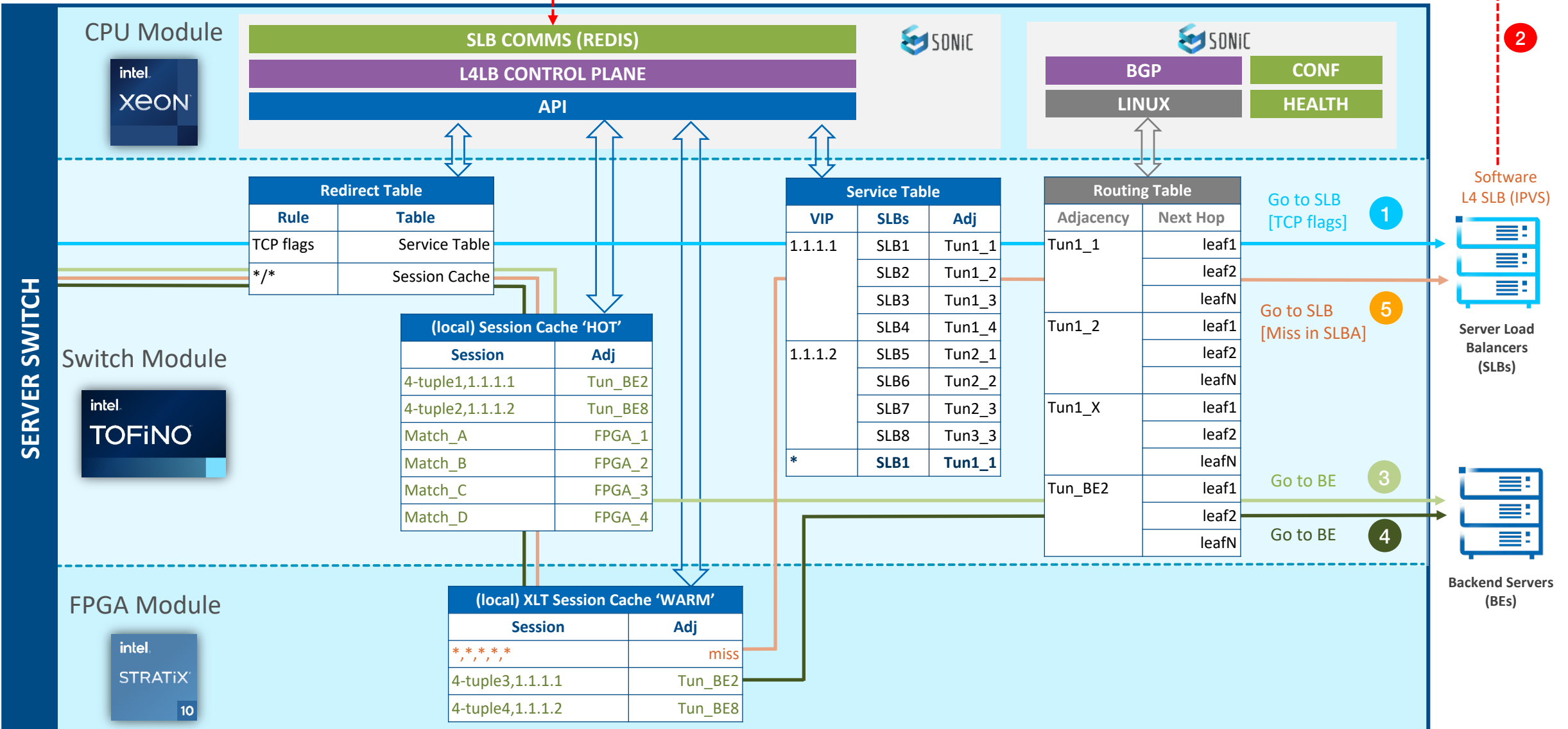
- 1 Client's request
- 1 Client's request forwarded to software SLB
- 2 Client's request load balanced to the right backend server
- 2 Accelerated fast path load balancing by SLBA
- 3 Content delivered to the client



L4 Server Load-Balancer Accelerator Data-Plane Architecture

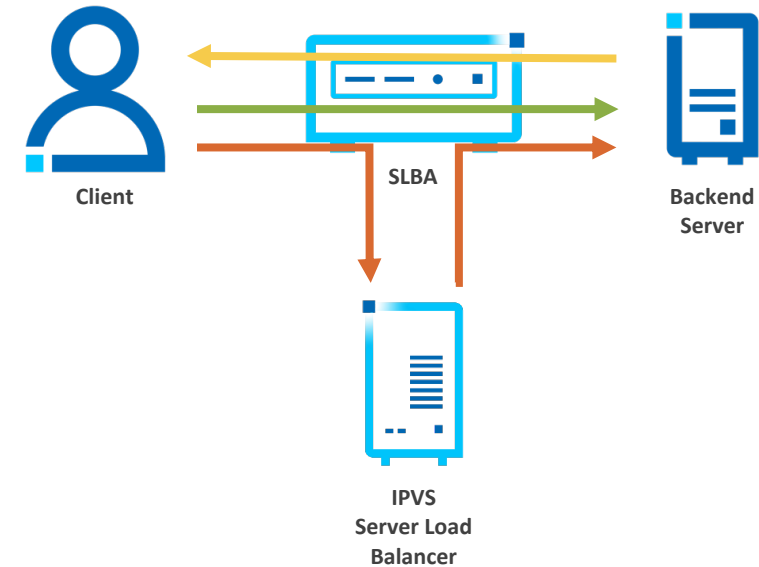


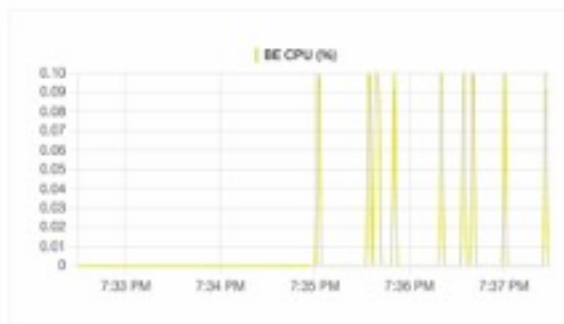
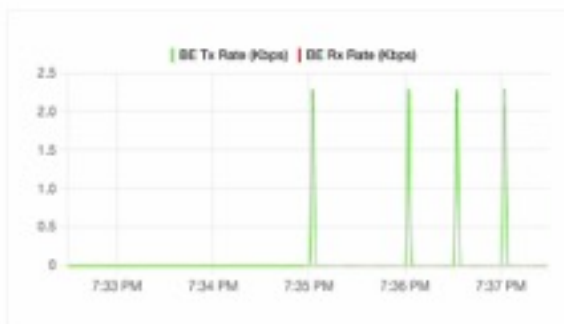
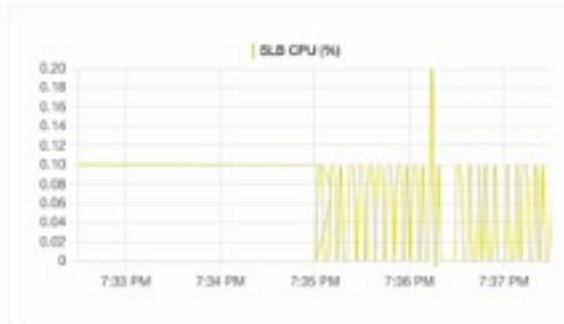
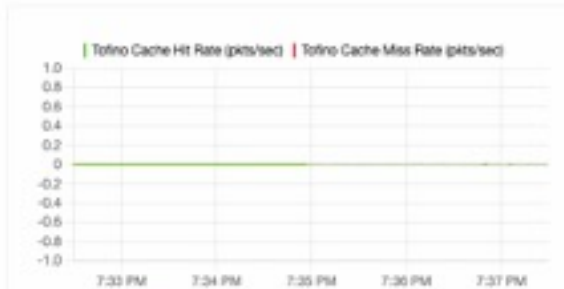
Add session message: client 5-tuple, BE VRF+IP+port



SLBA demo setup

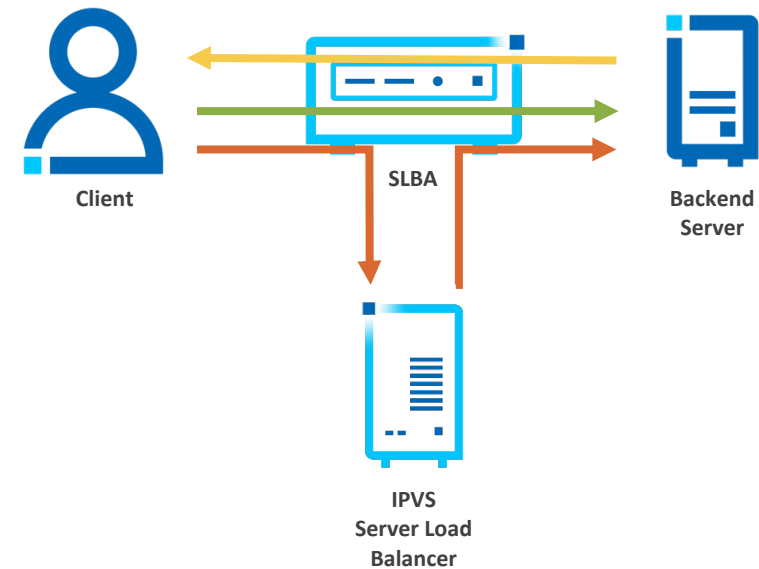
- One client, one SLB and one BE, 100GbE, 2x22 cores each
- Client initiates progressively 1000 iperf3 flows
- First run with SLBA disabled, second run with SLBA enabled





SynProxy demo setup

- 8k legitimate requests
- 6k attacker requests periodically
- Initially SynFlood detection is off, so all connections are accepted
- We then turn on SynFlood detection and SynProxy mitigation, and attacker traffic is dropped



intel®