

# Breaking Limits: Terabit Speeds on a single CPU server

**Marco Chiesa** 

**KTH Royal Institute of Technology** 



# Breaking Limits: Terabit Speeds on a single CPU server

Marco Chiesa

KTH Royal Institute of Technology

#### Joint work with





Mariano Hamid Scazzariello Ghasemirahni



Farshin



Dejan

Kostić



Gerald Q. Maguire Jr

Alireza



#### Massive amount of data to be processed at **cloud/edge networks**

### **External** internet-facing

- users, sensors, services
- terabits per second (Tbps)
- millions of connections per second

Internal machine-to-machine

- big-data, AI/ML, etc
- thousands of servers
- petabits per second (Pbps) of traffic



Processing massive amounts of traffic at high rates is hard!



Network Functions Virtualization is an **essential** architectural paradigm of today's data centers



Network Functions Virtualization is an **essential** architectural paradigm of today's data centers

NFs are deployed on commodity hardware rather than specialized one



**Cost-Effective** 



**Flexible** 



Scalable



NFs are deployed on commodity hardware rather than specialized one



Lower performance compared to specialized hardware



























**None** of the existing frameworks can handle such loads!





# Can we scale complex, stateful packet processing chains on commodity hardware to over 100 millions packets per second?



Batch Processing

Minimize Memory Accesses

**Minimize Sync. Overheads** 



> Batch Processing

Minimize Memory Accesses

**Minimize Sync. Overheads** 





VPPBatch Processing

**Minimize Memory Accesses** 

Minimize Sync. Overheads





NF



Brables **software prefetching** of data



✓ VPP
 Batch Processing > Minimize Memory Accesses Minimize Sync. Overheads
 ⓓ Hash tables are the essential data structure in stateful NFs

Require at least two random memory accesses





VPP
 Batch Processing > Minimize Memory Accesses Minimize Sync. Overheads
 Hash tables are the essential data structure in stateful NFs

Require at least two random memory accesses













VPP
 FastClick
 Batch Processing
 Minimize Memory Accesses
 Minimize Sync. Overheads

How to minimize memory accesses?



Aggregate States







✓ FastClick

VPP **Batch Processing** Minimize Memory Accesses

How to scale NFs on multiple cores?



> Minimize Sync. Overheads

Shared-Nothing Model



✓ FastClick

Minimize Memory Accesses

VPP Batch Processing

How to **scale** NFs on **multiple cores**?

High Performance

Per-flow Consistency

Not applicable in all scenarios!



> Minimize Sync. Overheads

Shared-Nothing Model



#### (i) Efficient Inter-Core Load Balancing

RSS **struggles** to efficiently distribute packets among cores with **adversarial loads** 





### (i) Efficient Inter-Core Load Balancing

- RSS **struggles** to efficiently distribute packets among cores with **adversarial loads**
- Monitor and rebalance flows among cores (e.g., Dyssect, RSS++)





### (i) Efficient Inter-Core Load Balancing

- RSS **struggles** to efficiently distribute packets among cores with **adversarial loads**
- Monitor and rebalance flows among cores (e.g., Dyssect, RSS++)





### (i) Efficient Inter-Core Load Balancing

- RSS **struggles** to efficiently distribute packets among cores with **adversarial loads**
- Monitor and rebalance flows among cores (e.g., Dyssect, RSS++)





(i) Efficient Inter-Core Load Balancing

(ii) Stateful NFs with Diverse Flow Definitions

- 6
- The definition of state is **not exclusively** the 5-tuple!
- What happens if they are **chained together**?

Load Balancer	Src IP	Dst IP	Src Port	Dst Port	Protocol	
Policer	Src IP	Dst IP	Src Port	Dst Port	Protocol	
Port Scan Detector	Src IP	Dst IP	Src Port	Dst Port	Protocol	



(i) Efficient Inter-Core Load Balancing

(ii) Stateful NFs with Diverse Flow Definitions



- The definition of state is **not exclusively** the 5-tuple!
- What happens if they are **chained together**?

_							Separate hash table per
Load Balancer	Src IP	Dst IP	Src Port	Dst Port	Protocol	←	flow definition!
Policer	Src IP	Dst IP	Src Port	Dst Port	Protocol	←	
			-	-	-		
Port Scan Detector	Src IP	Dst IP	Src Port	Dst Port	Protocol	$\leftarrow$	



(i) Efficient Inter-Core Load Balancing

(ii) Stateful NFs with Diverse Flow Definitions



- The definition of state is **not exclusively** the 5-tuple!
- What happens if they are **chained together**?

						_	Sepa	rate ha	i <b>sh table</b> pe	r
Load Balancer	Src IP	Dst IP	Src Port	Dst Port	Protocol		flow definition!			
						_				_
Policer	Src IP	Dst IP	Src Port	Dst Port	Protocol	<b>←</b>	States are <b>shared</b>			
						among all cores				
Port Scan Detector	Src IP	Dst IP	Src Port	Dst Port	Protocol	$\longleftarrow$				
						-	A Performan degradatio		ormance radation!	



#### VPP Batch Processing

✓ FastClick
Minimize Memory Accesses

**?** Minimize Sync. Overheads



**Overlooking** one of these principles **hinders** Stateful NF frameworks to process >100Mpps



We **must** support the three principles **together**!





#### Batch Processing

Minimize Memory Accesses

**Minimize Sync. Overheads** 



**Batch Processing** 

#### Minimize Memory Accesses

#### **Minimize Sync. Overheads**



**Optimized** state aggregation



Bulk Lookup for all packets in the batch





**Batch Processing** 

#### **Minimize Memory Accesses**

**Minimize Sync. Overheads** 

**Optimized** state aggregation



Bulk Lookup for all packets in the batch






# Minimize Memory Accesses

### **Minimize Sync. Overheads**





# Minimize Memory Accesses

#### **Minimize Sync. Overheads**

Each NF prefetches its state data at the right time









# Minimize Memory Accesses

# **Minimize Sync. Overheads**

Only considers 5-tuple flows definition!





Minimize Memory Accesses

**Minimize Sync. Overheads** 





**Minimize Memory Accesses** 

**Minimize Sync. Overheads** 





Minimize Memory Accesses

**Minimize Sync. Overheads** 



47



Minimize Memory Accesses

**Minimize Sync. Overheads** 



48



Minimize Memory Accesses

**Minimize Sync. Overheads** 



49



# Minimize Memory Accesses

# **Minimize Sync. Overheads**





# Minimize Memory Accesses

# **Minimize Sync. Overheads**





# Minimize Memory Accesses

# **Minimize Sync. Overheads**













# Minimize Memory Accesses

# **Minimize Sync. Overheads**





# Minimize Memory Accesses

# **Minimize Sync. Overheads**

















#### **Batch Processing Minimize Memory Accesses Minimize Sync. Overheads** How can we **reduce** the **inevitable** sync. overhead with different flow definitions? Only a single expensive insert/lookup into/from shared HT Minimize sync. overheads Core 0 NF 2 (Dst IP) NF 1 (Src IP) **Auxiliary HT** (5-tuple) **%** } ₩0 **.....** Shared Shared Shared Shared HT HT **States** States **Auxiliary HT** Core 1 (5-tuple)



# Minimize Memory Accesses

### **Minimize Sync. Overheads**

How can we **reduce** the **inevitable** sync. overhead with different flow definitions?

Only a single expensive insert/lookup into/from shared HT





# **Evaluation**









# 16 CPU cores Synthetic 64-B trace w/ 2 million flows





# 16 CPU cores Synthetic 64-B trace w/ 2 million flows





16 CPU cores Synthetic 64-B trace w/ 2 million flows



#### FAJITA is the first system to enable >100Mpps stateful NF chains on a single socket!

200 176.6 Throughput (Mpps) 150 2.43x 100 72.5 70.8 52.9 50 0 FastClick Dyssect VPP FAJITA Framework

**Shared Nothing** 

(Flow Statistics Counter + Rate Limiter + Load Balancer)

#### Shared (Src IP Counter + Policer + Load Balancer)





Does FAJITA Improve End-to-End Latency?



Load Balancer + Flow Statistics Counter



Does FAJITA Improve End-to-End Latency?



Load Balancer + Flow Statistics Counter



### Does FAJITA Improve End-to-End Latency?

# FAJITA performs stateful packet processing in half the time of other frameworks







8 CPU cores – Synthetic Trace w/ Elephants Only

Flow Statistics Counter – 100ms Monitoring Interval





8 CPU cores – Synthetic Trace w/ Elephants Only

Flow Statistics Counter – 100ms Monitoring Interval





8 CPU cores – Synthetic Trace w/ Elephants Only

Flow Statistics Counter – 100ms Monitoring Interval



# FAJITA achieves better performance **solely relying on RSS**

Load imbalance **only** depends on the number of flows, that are usually **>2K at high rates**!







The **FAJITA** system:

- incorporates the three essential principles for stateful packet processing together
- minimizes memory access overheads by exploiting batching & software prefetching
- alleviates overheads of accessing shared data structures by introducing auxiliary HTs
- achieves performance 2x higher than existing stateful packet processors
- demonstrates that dynamic inter-core load balancing is **detrimental** at high rates
- only processes packet headers

How to split packets?



State-of-the-art **packet-splitting** solutions **increase** <u>shallow</u> NFs throughput




State-of-the-art packet-splitting solutions increase NFs throughput



#### **Improve NF Performance**

Better CPU cache exploitation



#### Free up Bandwidth

Reduce the number of dedicated NF servers Lower energy consumption



#### **Complex Deployment**

Require end-host modifications Require available shared resources



State-of-the-art packet-splitting solutions increase NFs throughput























State-of-the-art packet-splitting solutions increase NFs throughput

Storing Payloads in the On-NIC Memory (nicmem)

Storing Payloads on <u>Shared Resources</u> (Ribosome) 🗫 💋 🍟

















State-of-the-art packet-splitting solutions increase NFs throughput

Storing Payloads in the <u>On-NIC Memory</u> (nicmem) **See** 

Storing Payloads on Shared Resources (Ribosome) 🗫 🗡

#### Store on the <u>ASIC SRAM</u>

• ideal, but how? no API to store payloads!



# Can we store the entire payload within the switch without requiring external devices?



**Core Idea: A Queue-Based Packet Storage** 

#### Exploit the switch shared buffer to store payloads



#### **Exploit the switch shared buffer to store payloads**





#### **Exploit the switch shared buffer to store payloads**







Enqueued packets **implicitly** control the buffer memory



Hold payloads in the port queues until the processed header from the NF is returned!



Would it work for real?



































Q

**Congest** the egress queues to delay payloads release by the time required for NF processing





**Congest** the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready




**Congest** the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues





**Congest** the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues



Stop the queue with Priority-Flow-Control (PFC)



Q

**Congest** the egress queues to delay payloads release by the time required for NF processing



Recirculate payloads until headers are ready



Create copies of payloads to fill up the queues





## **Evalution of QueueMem (combined with FAJITA)**

Network function: Per-flow counter + Load balancer + Per-flow rate limiter



## FAJITA + QueueMem conclusions

The **FAJITA** system:

- incorporates the three essential principles for stateful packet processing together
- minimizes memory access overheads by exploiting batching & software prefetching
- alleviates overheads of accessing shared data structures by introducing auxiliary HTs
- achieves performance 2x higher than existing stateful packet processors
- demonstrates that dynamic inter-core load balancing is detrimental at high rates
- only processes packet headers

The **QueueMem** system (under submission):

- supports packet splitting entirely on the switch
- reduces power consumption
- relies on PFC to limit buffer usage

## FAJITA + QueueMem conclusions

The **FAJITA** system:

- incorporates the three essential principles for stateful packet processing together
- minimizes memory access overheads by exploiting batching & software prefetching
- alleviates overheads of accessing shared data structures by introducing auxiliary HTs
- achieves performance 2x higher than existing stateful packet processors
- demonstrates that dynamic inter-core load balancing is detrimental at high rates
- only processes packet headers

The **QueueMem** (under submission):

- supports packet splitting entirely on the switch
- reduces power consumption
- relies on PFC to limit buffer usage

