



# Towards Distributed Adaptive Computing

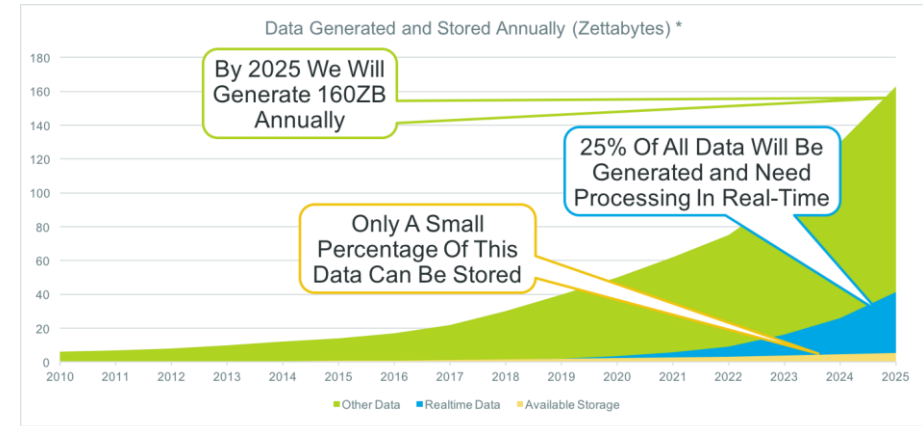
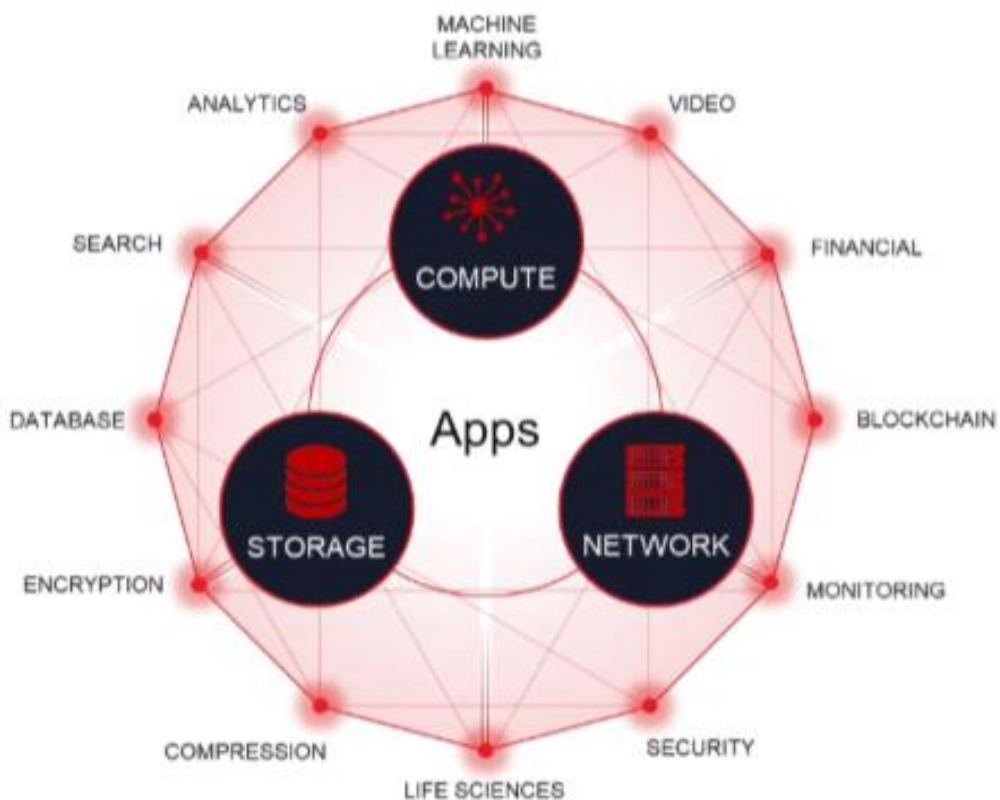
Chengchen Hu

Xilinx Labs

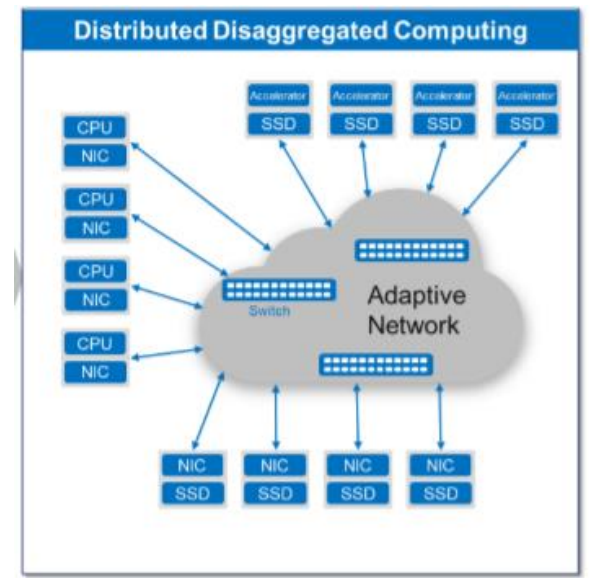
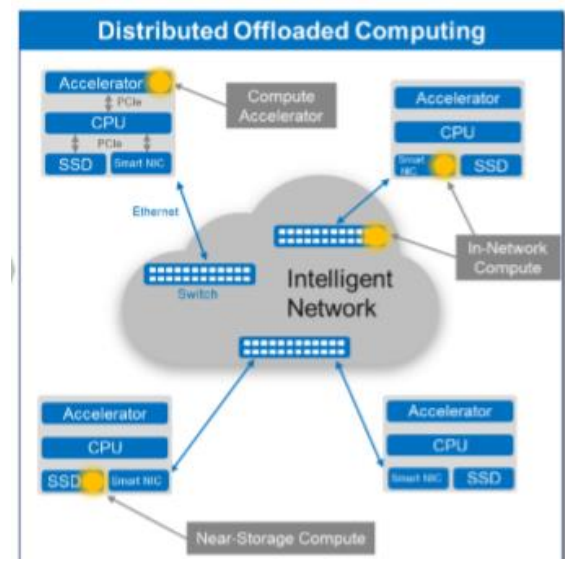
July 24, 2020



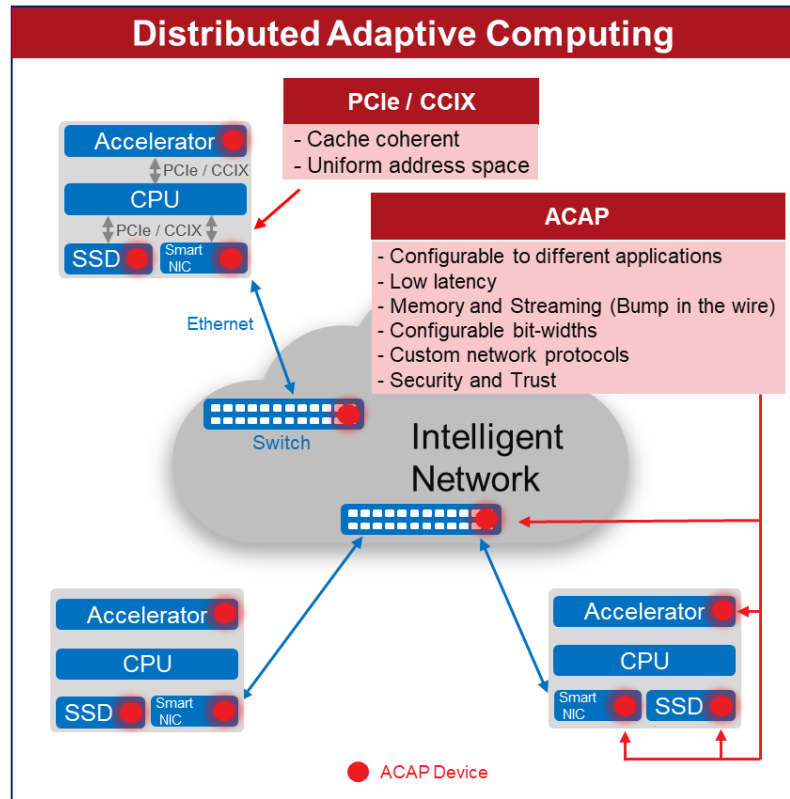
# Trends in DCs



\* Data Age 2025: The Evolution of Data to Life-Critical. An IDC White Paper, Sponsored by Seagate



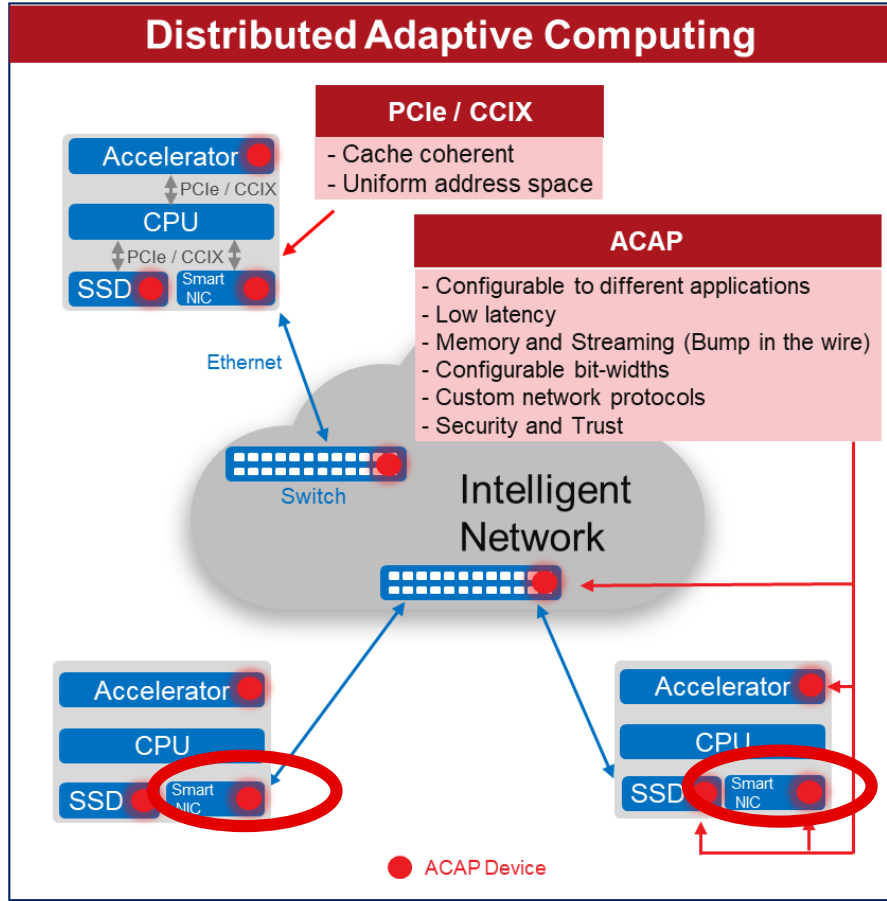
# Concept of Distributed Adaptive Computing



**Adaptable NIC:**  
a new type of SmartNIC as host network device

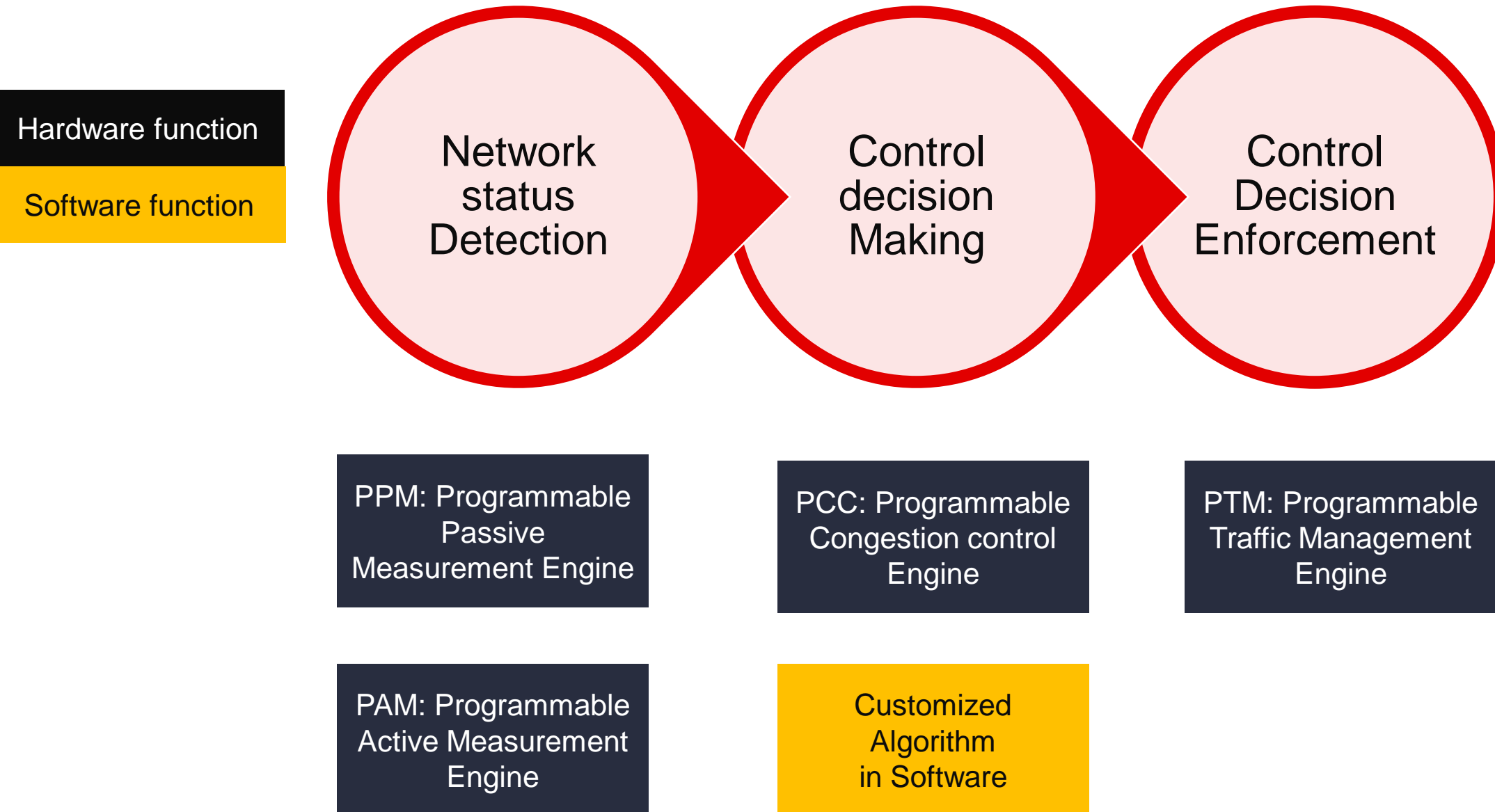
**Adaptable Switch:**  
Programmability in intermediate network nodes

# Adaptable NIC



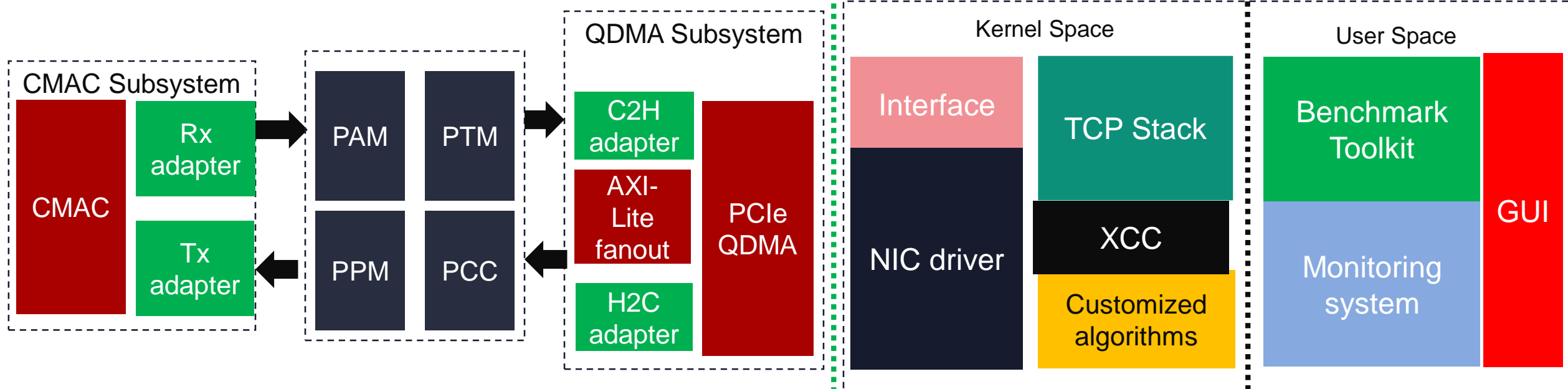
	Description	Features
<b>Type 1</b>	Basic Connectivity NIC	<ul style="list-style-type: none"> <li>Basic offloads, simple virtualization</li> </ul>
<b>Type 2</b>	SmartNIC for Network Acceleration	<ul style="list-style-type: none"> <li>Crypto, <u>vSwitch</u>, custom tunneling</li> </ul>
<b>Type 3</b>	SmartNIC for Network + Compute + Storage Acceleration	<ul style="list-style-type: none"> <li>Machine Learning, video transcoding</li> <li>Database Analytics, smart storage</li> </ul>
<b>Adaptable NIC</b>	<u>SmartNIC</u> as standalone managed network node	<ul style="list-style-type: none"> <li>Domain Specific Programmable Engines offloaded</li> </ul>

# Why it is adaptable: HW+SW Abstraction



# Adaptable NIC prototype

Alveo  
U250



## ► NIC Hardware

- Network measurement: PPM(passive), PAM(active)
- Network Congestion Control: PCC
- Network Traffic Shaping: PTM

## > NIC Software at host server

- > Network Congestion Control: XCC
- > Network Monitoring: Monitoring System, GUI
- > Network Benchmark Generation Toolkit

# Story of Transport Congestion Control

## ► Congestion Control on HW

- TCP variants on HW (Chelsio).
- Socket API (TCP offload engine).
- RMDA(iWARP).
- RoCE using DCQCN (Microsoft).
- HPCC (Alibaba),FPGA based.

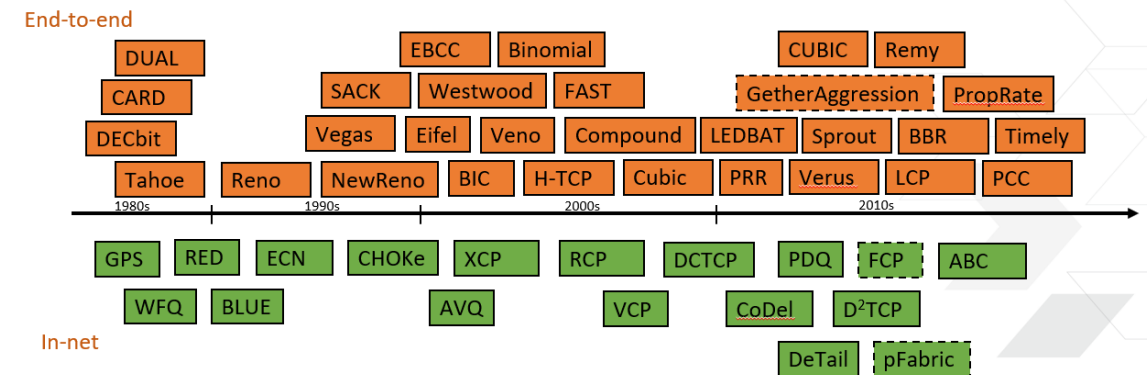
“needed to modify the data delivery algorithm to avoid livelocks in their network but had to rely on the NIC vendor to make that change.”

## ► Congestion Control keeps evolving

- Environment changes
- Application requirement changes

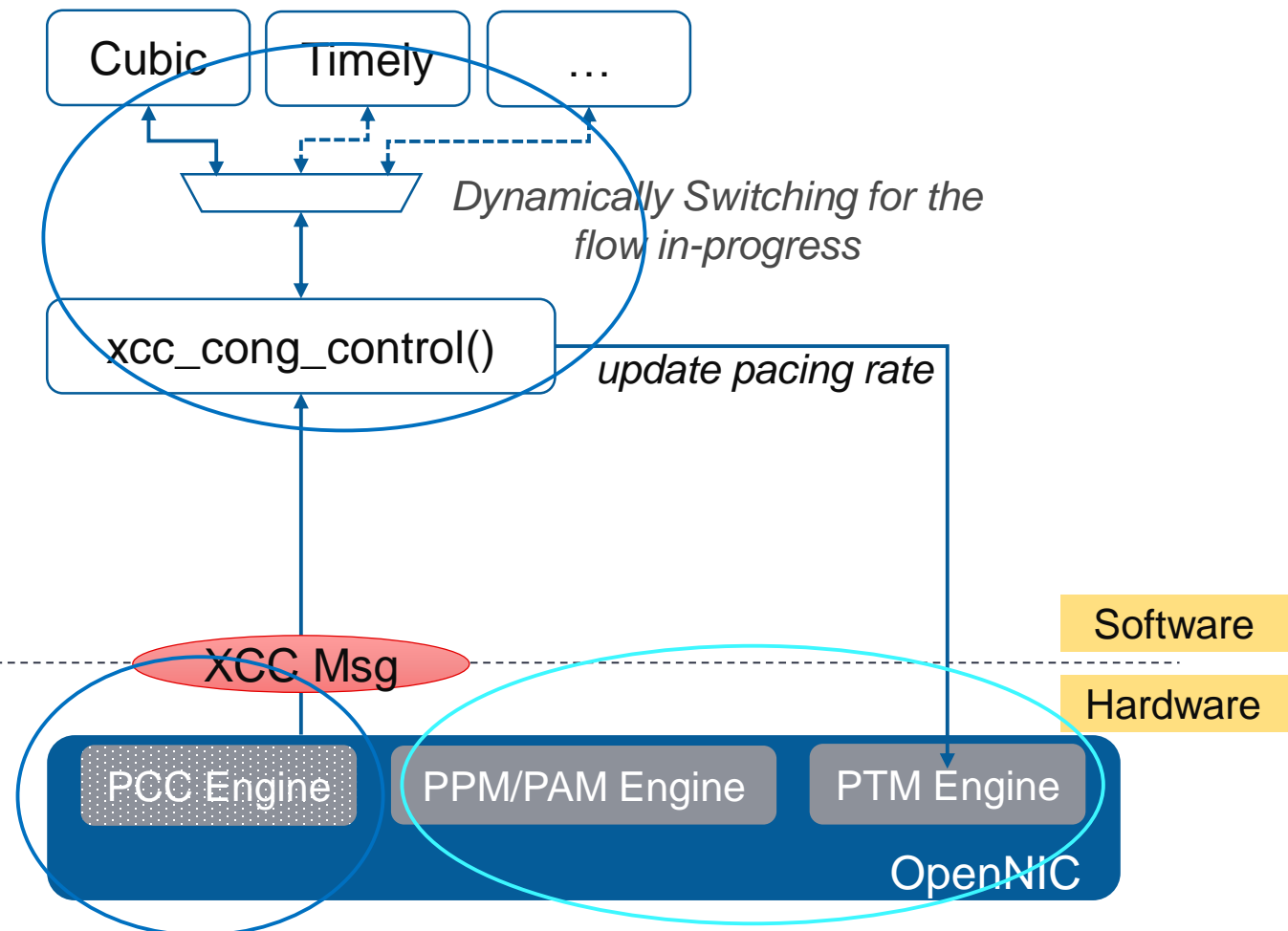
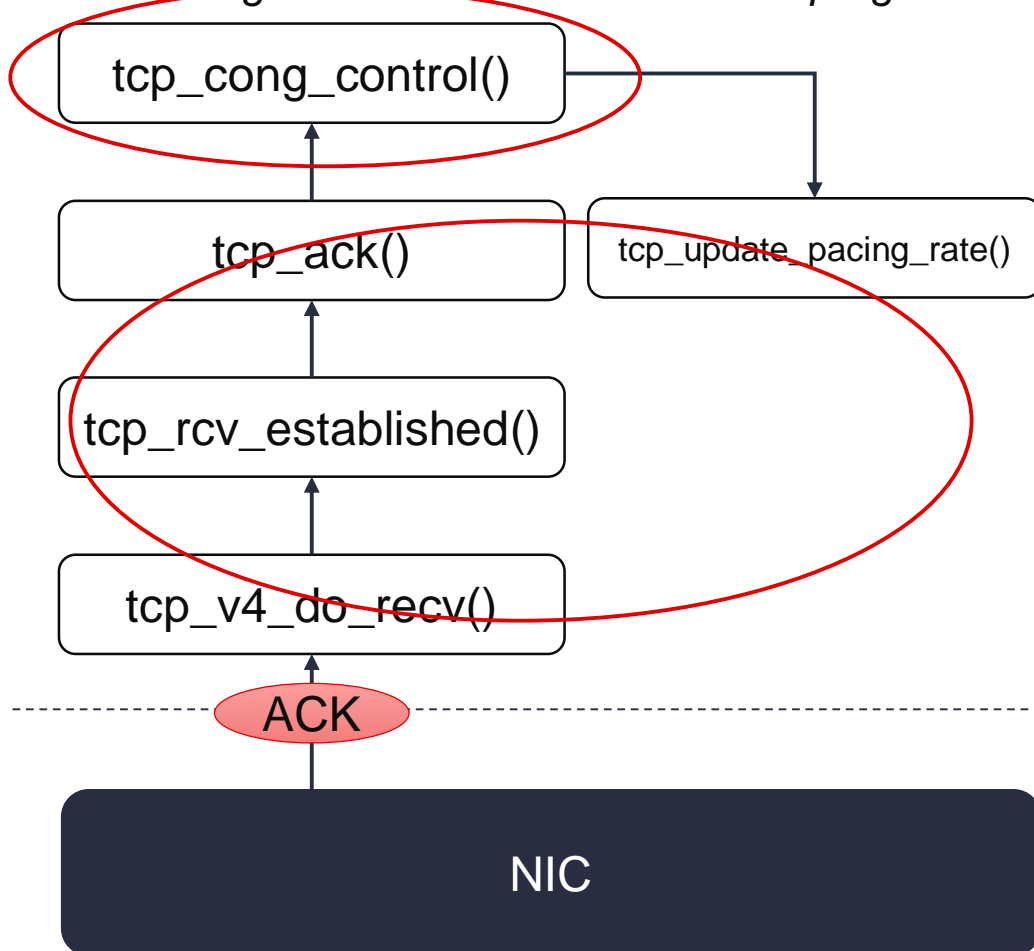
## ► Programmable CC on FPGA

- TONIC (Princeton & MIT), NSDI 2020
- Limited programmability provided



# Mapping CC to Adaptable NIC

*The CC algorithm is fixed for the flow in-progress*

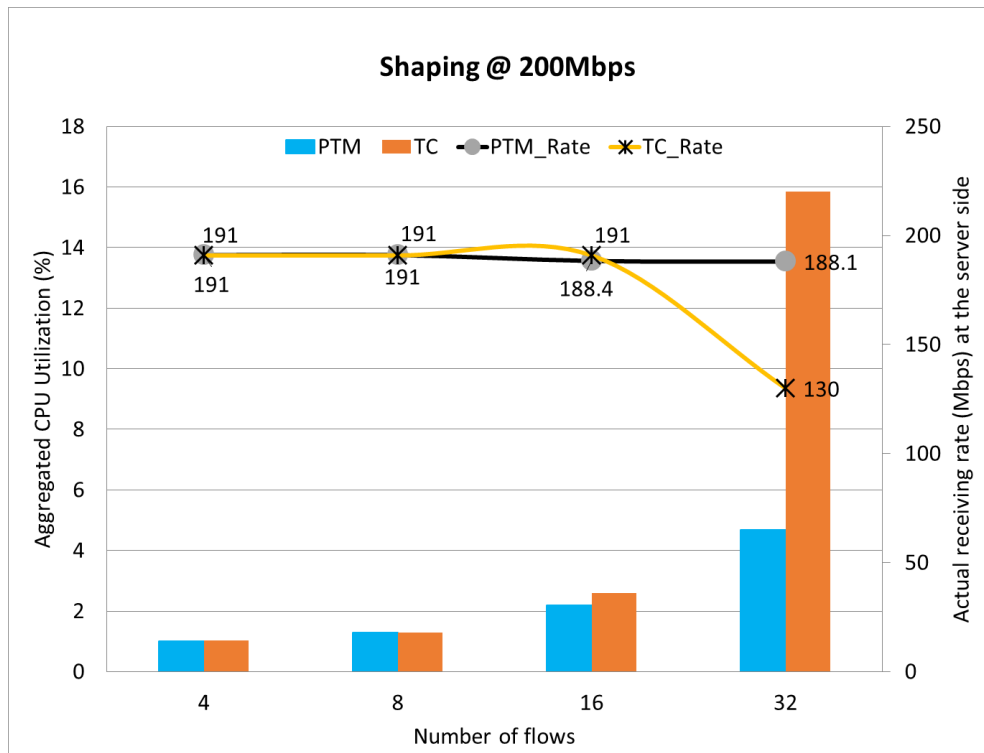


*Parts of the **xcc\_cong\_control()** could be further moved to the **PCC Engine***



# Use Case: Public Cloud

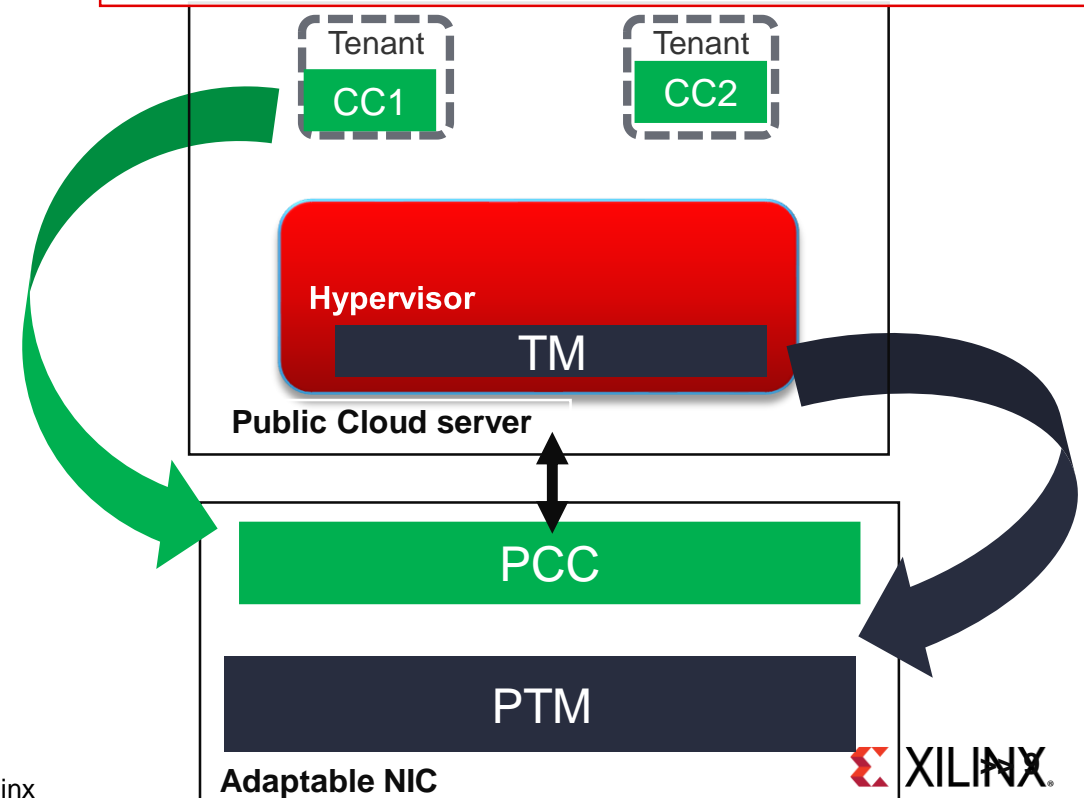
- > Tenants Share Machine in Public Cloud
- > Each tenant can customize CC policy for service optimization, with the support of PCC in Adaptable NIC
- > Cloud operator define scheduling/policing/shaping in hypervisor supported by PTM in Adaptable NIC
- > QoS guarantee on both ingress and egress



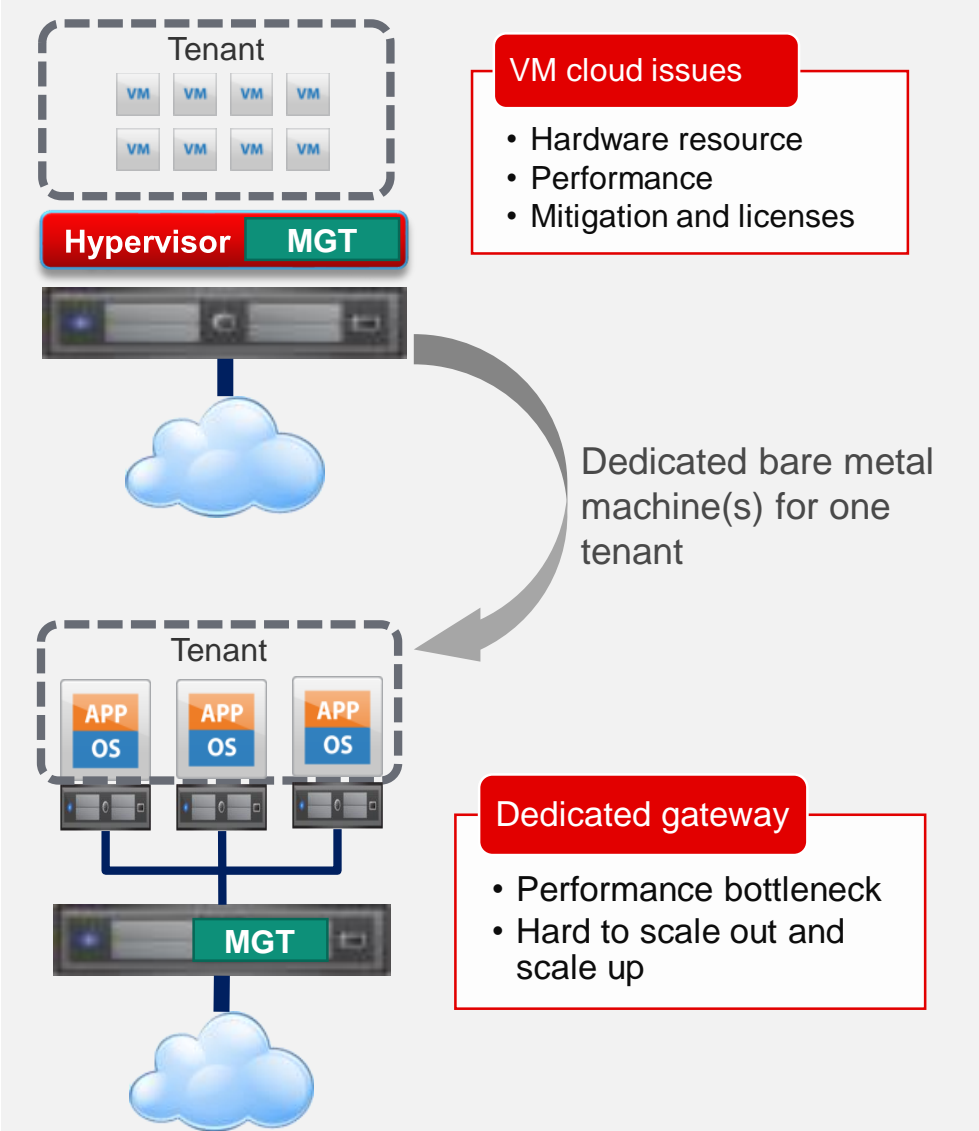
© Copyright 2020 Xilinx

## Management in Adaptable NIC

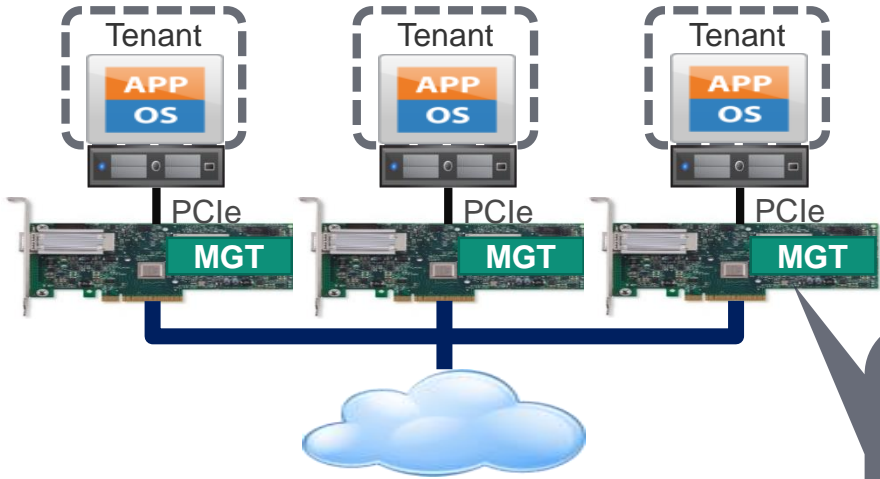
- Less CPU, more income
- Fine-grained bandwidth control and guarantee



# Use Case: Bare Metal Cloud



>> 10 >> 10

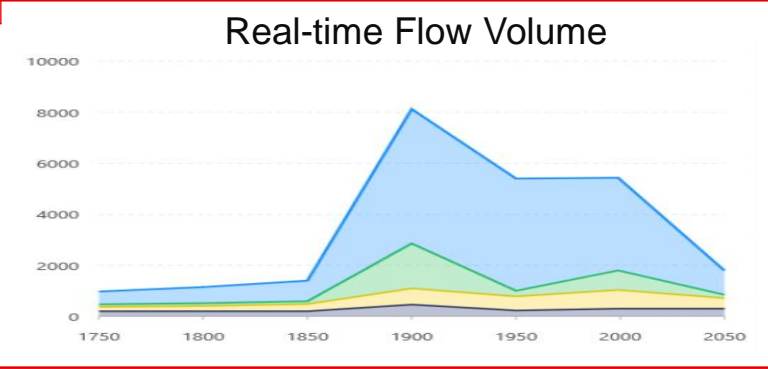


- Measurements
- Usage Statistics
  - Anomaly Detection
  - Heavy jitter
  - Active Probes
  - .....

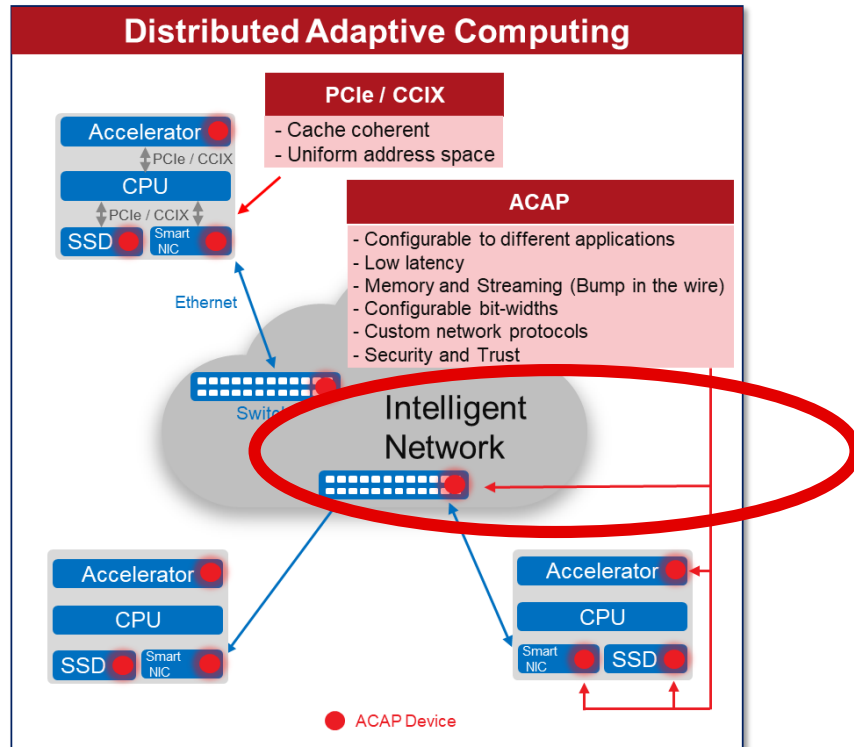
IP Address	Port	Alarm
192.168.1.12	22	True
192.168.1.17	1433	True
192.168.1.18	80	True

Flow ID	Heavy-changer
2	True

IP Address	Port Scanner
192.168.1.8	False
192.168.1.5	True

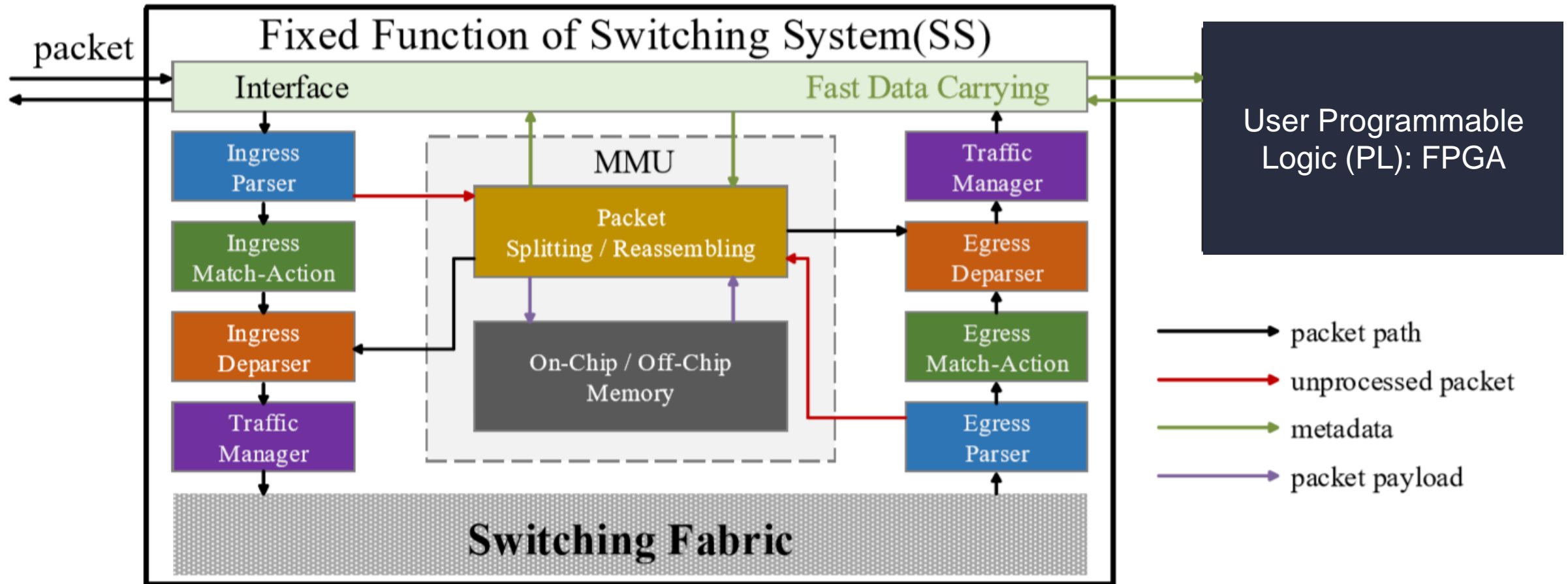


# Adaptable Switch



- ▶ Push adaptivity to switch for new opportunities.
  - In-network compute, e.g. distributed ML
  - In-network storage, e.g., cached streaming video
- ▶ Step 1: Two-chip approach
  - Combined software stacks
  - 'Flow-through' processing architecture
  - Switch chip and Xilinx chip on mother board
- ▶ Step 2: Integrated chip
  - Integrated software stack
  - 'Lookaside' processing architecture
  - Integrated switching and processing

# Architecture of Adaptable Switch



## Use case: Programmable Networking

- ▶ Event Driven processing
  - NDP – trigger congestion control logic by tracking buffer occupancy
- ▶ Stateful processing
  - Firewall
- ▶ Algorithmic capacity
  - Measurement – Statistical calculation

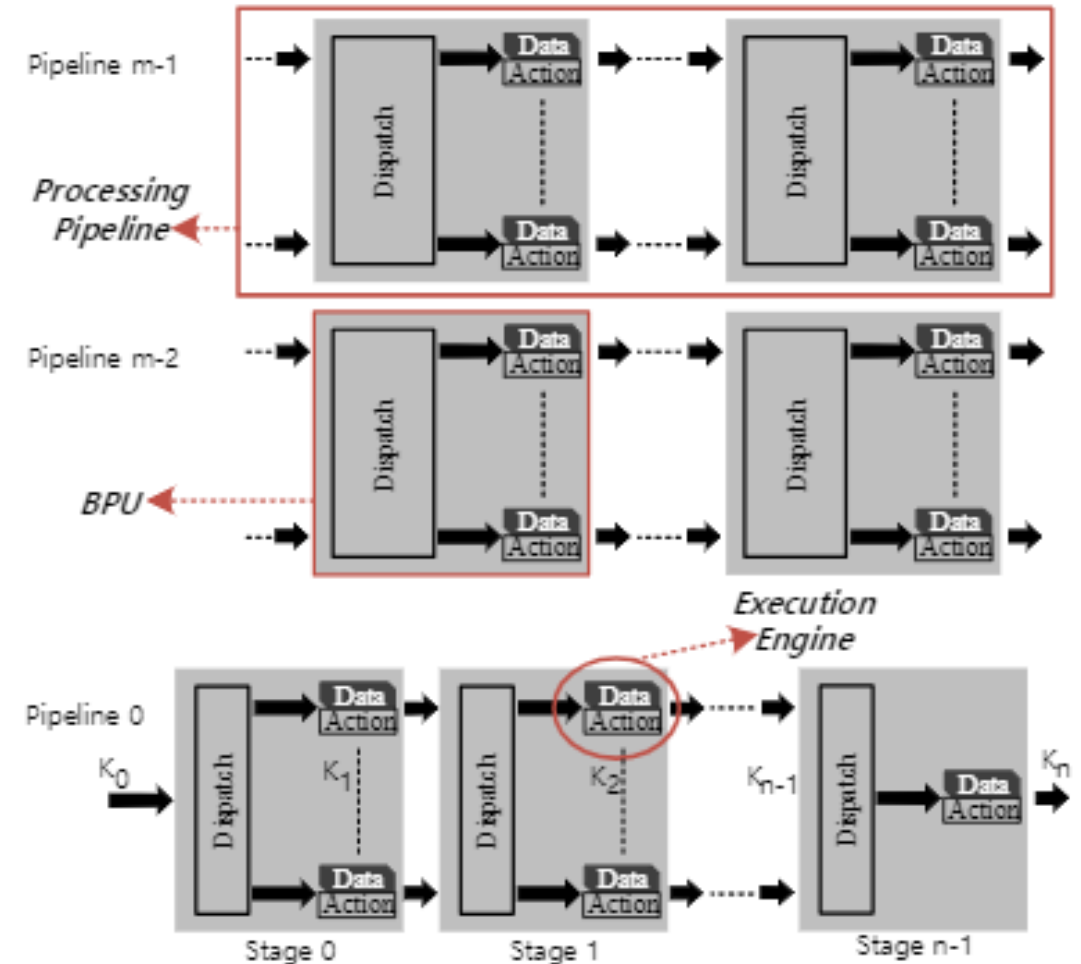
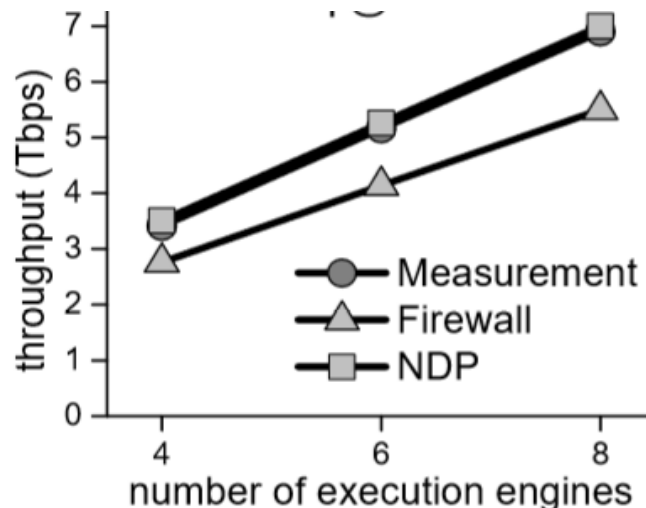
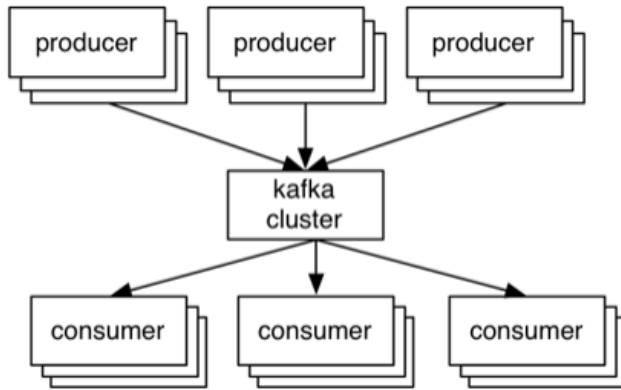


Fig. 3. A General Model of Parallel processing Pipeline in PL.

# Use case: In-Networking computing

- ▶ Kafka is a data streaming platform to manage/process data from multiple input streams
- ▶ How to enable a Kafka cluster in an Adaptable Switch
  - Add disks to Adaptable Switch attached to FPGA/MPSoC
  - FPGA/MPSoC enables brokers and zookeeper in the Kafka cluster
  - Filters out data to be kept in Kafka and offload to ACAP/FPGA/MPSoC for data



# Summary

- ▶ Distributed Adaptive Computing is flexible to enable user specific functions
  - Requirements vary user-by-user
  - Adaptable NIC/Switch provide new programming ways to deploy application accelerations
- ▶ User cases show value for cloud data center
  - Free server CPU for more users/tenants to increase revenue
  - Increase scalability (e.g., remove gateway for Bare Metal cloud)
  - Innovate with new functions, proprietary protocols and customized processing.



---

# Thank You

