

Cloud Computing and Hardware Accelerated Clouds

Ece Güran Schmidt

METU Electrical and Electronics
Engineering

Overview of the talk

- Part I: Cloud Computing
- Part II: Hardware Acceleration
- Part III: Hardware Accelerated Clouds
- Part IV: ACCLOUD Research Project



Part I: Cloud Computing

How the computing performed
(HW/OS/SW) is largely
irrelevant to the user.



Legacy Definitions of Cloud Computing

“A model for enabling, ubiquitous, convenient, **on-demand network access to a shared pool of configurable computing resources**”

- Resources: (networks, servers, storage, applications, and services)
- Can be rapidly provisioned and released with minimal management effort or service provider interaction.

<https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-145.pdf>

NIST

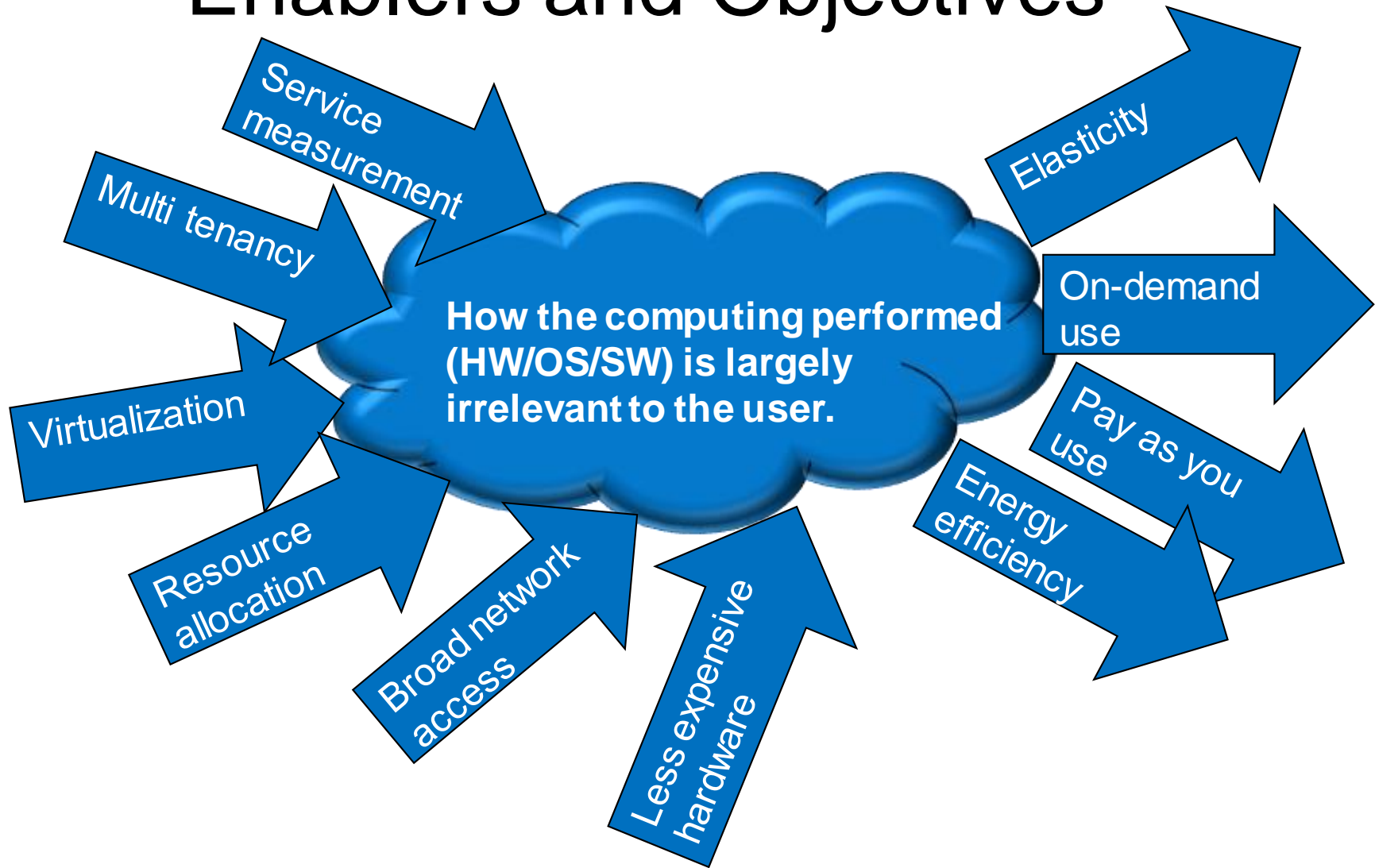
National Institute of
Standards and Technology
U.S. Department of Commerce

“The applications delivered as **services over the Internet** and the hardware and systems software in the data centers that provide those services.”

[Armbrust, Michael, et al. "A view of cloud computing." *Communications of the ACM* 53.4 \(2010\): 50-58. \(11846 citations\)](#)

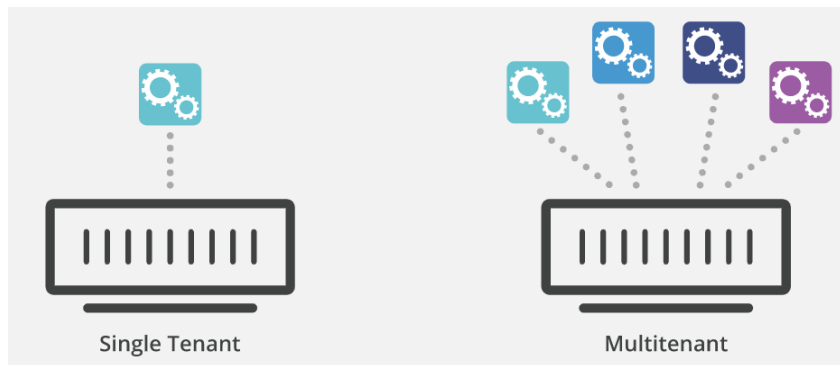


Enablers and Objectives



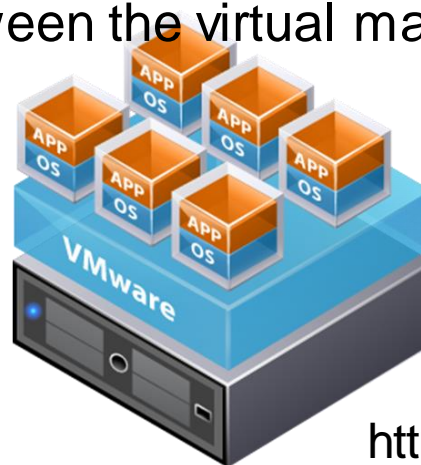
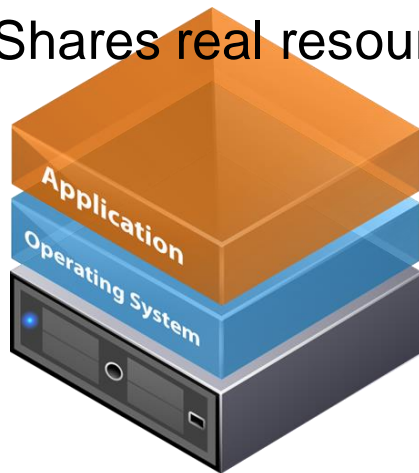
Multitenancy

- **Single-tenancy** is an architecture in which each customer has their own software instance; it requires a dedicated set of resources to fulfill the needs of just one organization
- **Multitenancy** is an architecture on which multiple customers (tenants) share the **same application**, running on the same operating system, on the same hardware, with the same data-storage mechanism.
- The distinction between the customers is achieved during application design, thus customers do not share or see each other's data.



Virtualization

- First idea: technology that allowed a host operating system (such as Linux) to execute one or more client operating systems (e.g., Windows).
- Hypervisor: A program
 - synthesizes virtual computing environments
 - virtual NIC, BIOS, sound adapter, and video
 - Shares real resources between the virtual machines.



<https://www.vmware.com/>

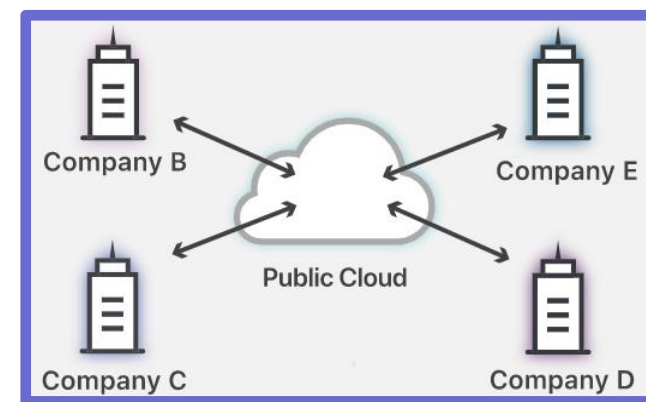
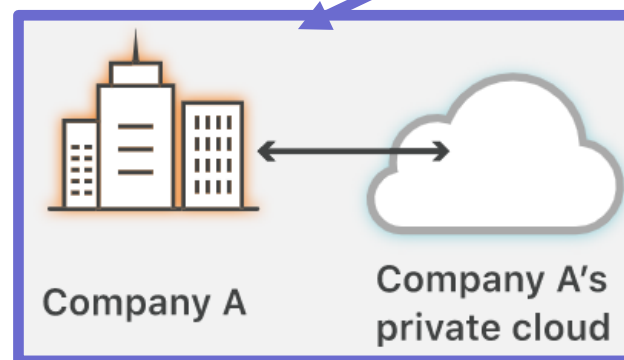
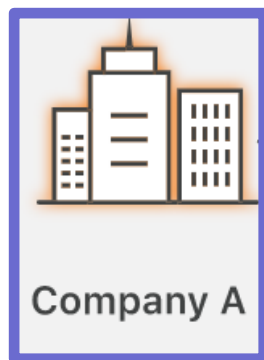
Traditional Architecture

Virtual Architecture



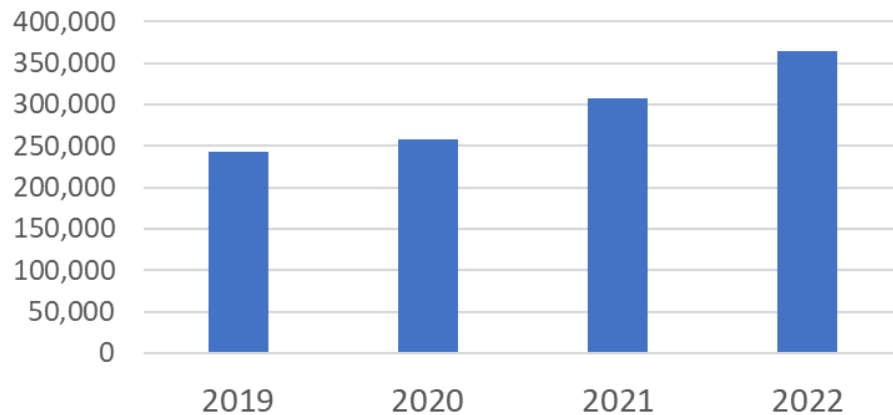
Computing Models

| | On-Premise (Traditional IT) | Private Cloud | Public Cloud |
|-----------------------------|--|--|---|
| HW Infrastructure Ownership | Company | Company | Service Provider |
| Software Ownership | Company | Company | Service Provider/Company |
| Maintenance | Company | Company | Service Provider/Company |
| Resource allocation | Fixed allocation to Company and applications | Fixed allocation to Company Elastic allocation to the applications. | Elastic allocation to Company and applications. |



Economy

Cloud Service Revenue Forecast (2020)
in Millions of \$



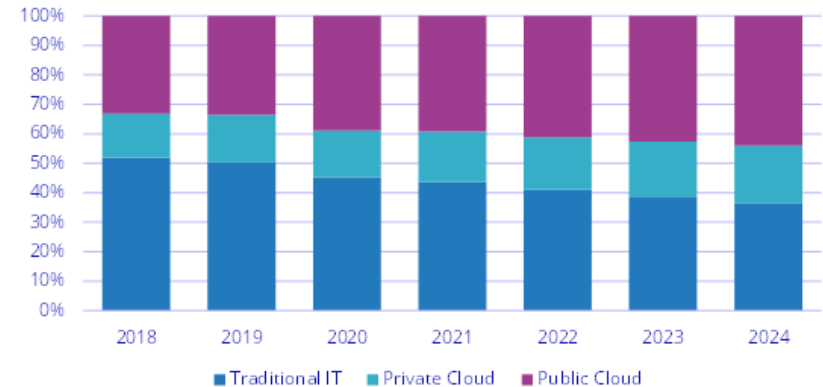
Worldwide Public Cloud Service Revenue Forecast (Millions of U.S. Dollars)

<https://www.gartner.com/>

Date: 2020-07-23



Worldwide Cloud IT Infrastructure Market Forecast by
Deployment Type, 2018- 2024 (shares based on Value)

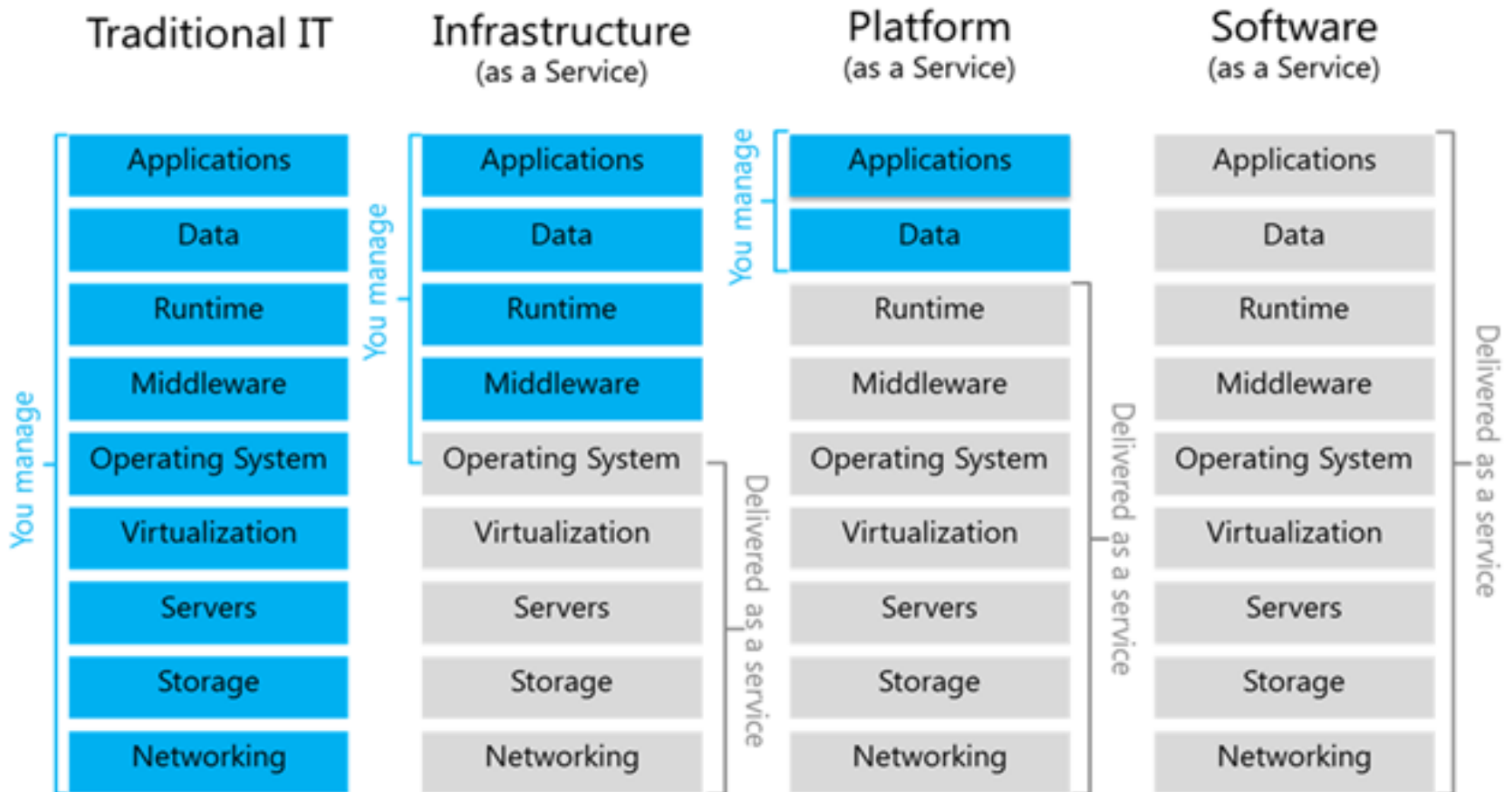


Source: IDC 2020

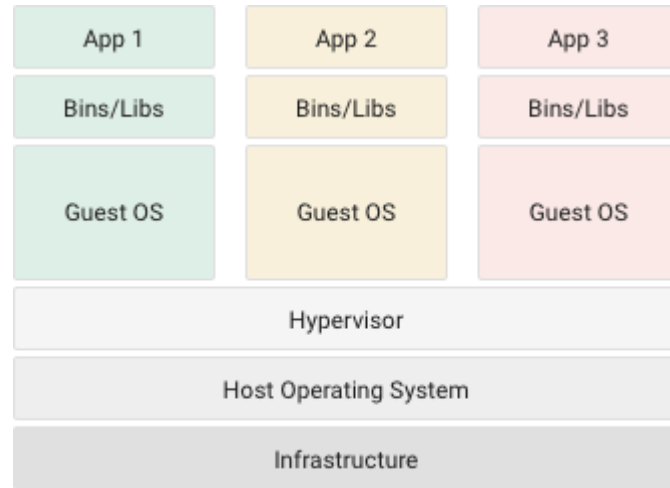
<https://www.idc.com/getdoc.jsp?containerId=prUS46895020>



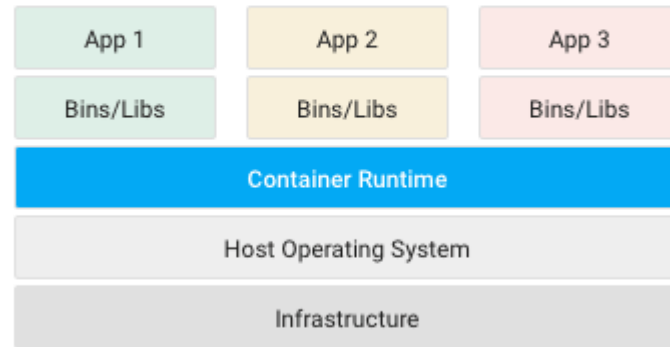
Cloud Services (Legacy Breakdown)



Containers



Virtual Machines

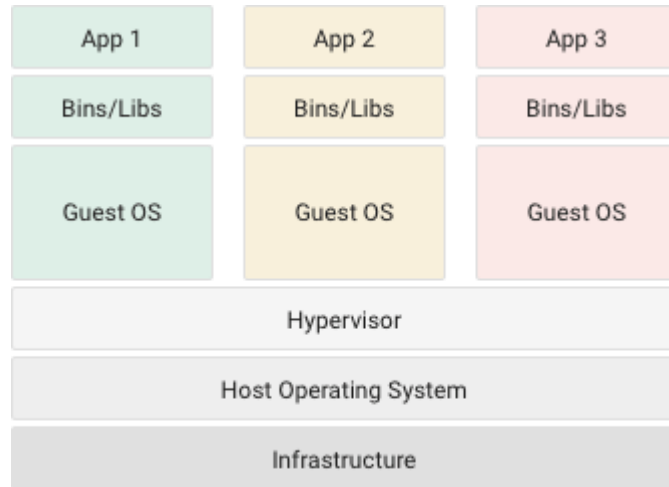


Containers

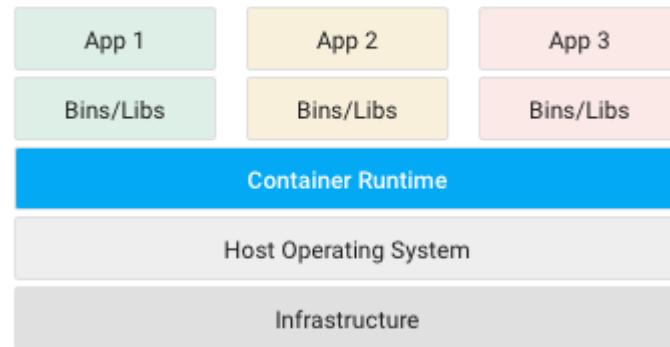
- an executable unit of software in which application code is packaged, along with its libraries and dependencies
- There are many container formats available. Docker is a popular, open-source container format.

<https://cloud.google.com/containers>

Containers



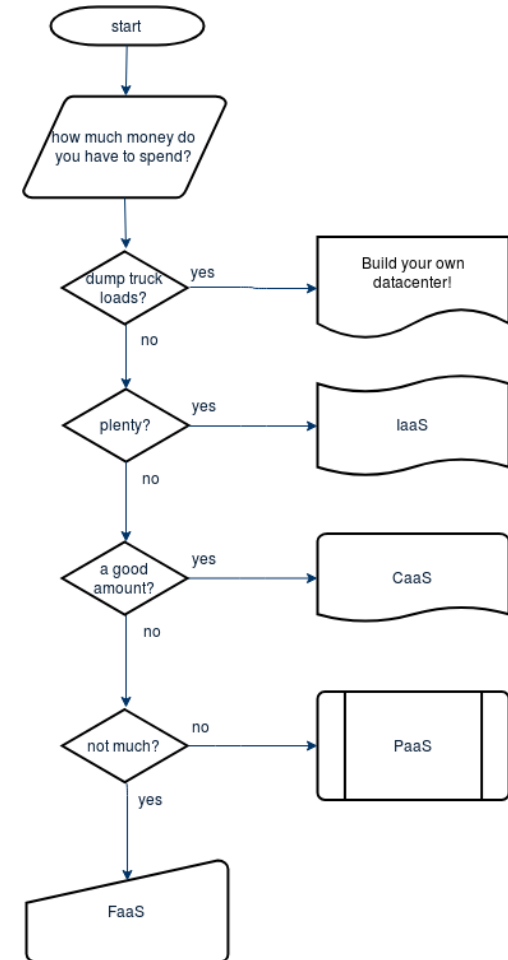
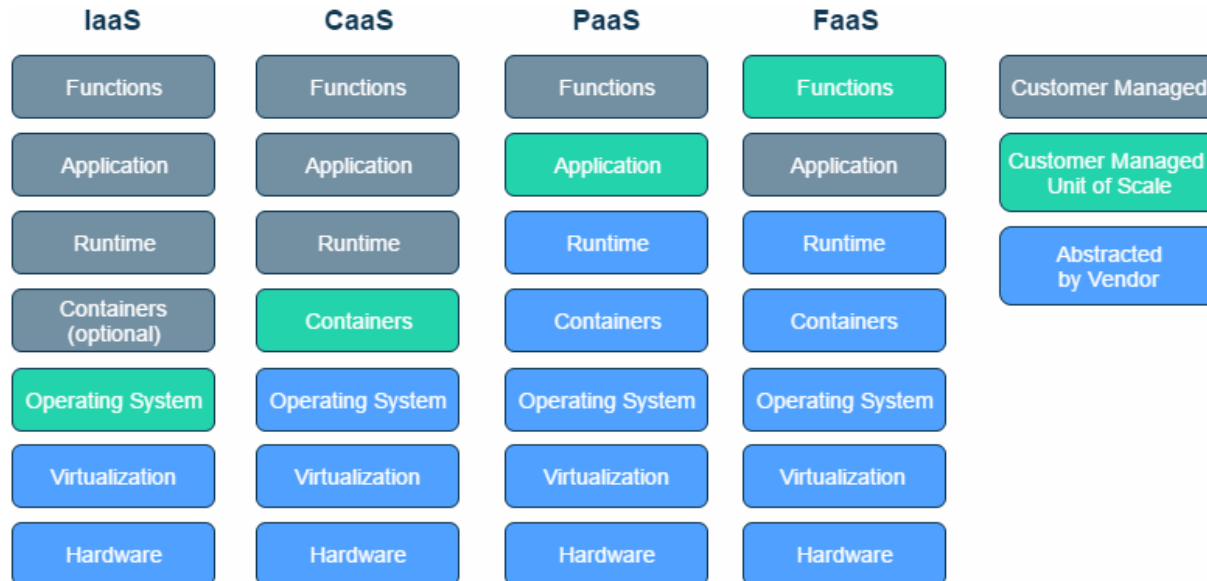
Virtual Machines



Containers

- Similar to virtual machines:
 - Package the application together with libraries and other dependencies,
 - Providing isolated environments for running the software
- Different than virtual machines:
 - multiple containers run atop the OS kernel directly.
 - more lightweight: they share the OS kernel, start much faster, and use a fraction of the memory compared to booting an entire OS.

Cloud Services (More Contemporary)



<https://serverless.zone/abstracting-the-back-end-with-faas-e5e80e837362>

<https://developer.ibm.com/articles/when-to-use-iaas-faas-paas-and-caas/>

CaaS/FaaS

- Containers as a service (CaaS): Allows users to upload, organize, start, stop, scale and otherwise manage containers, applications and clusters.
- FaaS (Function-as-a-Service)
 - an event-driven computing execution model that runs in stateless containers
 - infrastructure is usually metered on-demand, primarily through an event-driven execution model, so it's there when you need it but it doesn't require any server processes to be running constantly in the background, like platform-as-a-service (PaaS) would.



Cloud Data Center

- Traditional Data Center

- a single physical facility with all hardware infrastructure and equipment
- Houses all data and applications

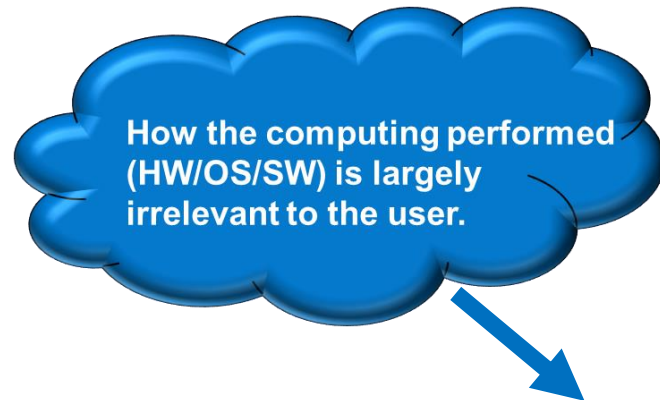
<https://www.cisco.com/c/en/us/solutions/data-center-virtualization/what-is-a-data-center.html#~types-of-data-centers>

- Cloud Data Center

- it's all online!
- **cloud servers host data and applications**
- Data automatically gets fragmented and duplicated across various locations for secure storage.



Part II: Hardware Accelerators

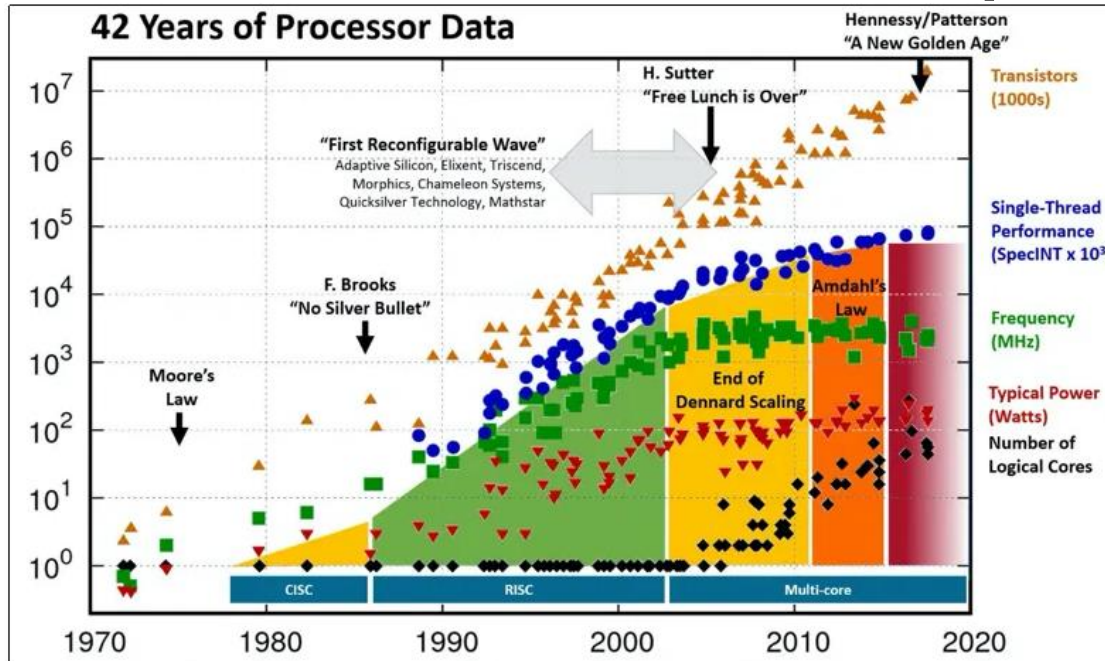


Xilinx Versal
ACAP
Adaptive
Compute
Acceleration
Platform

<https://www.xilinx.com/products/silicon-devices/acap/versal.html>

A blue arrow points from the URL to the hardware accelerator chip image.

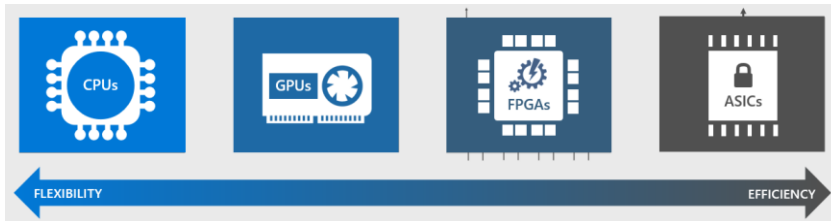
State of Computing



<https://iscaconf.org/isca2018/docs/HennessyPattersonTuringLectureISCA4June2018.pdf>

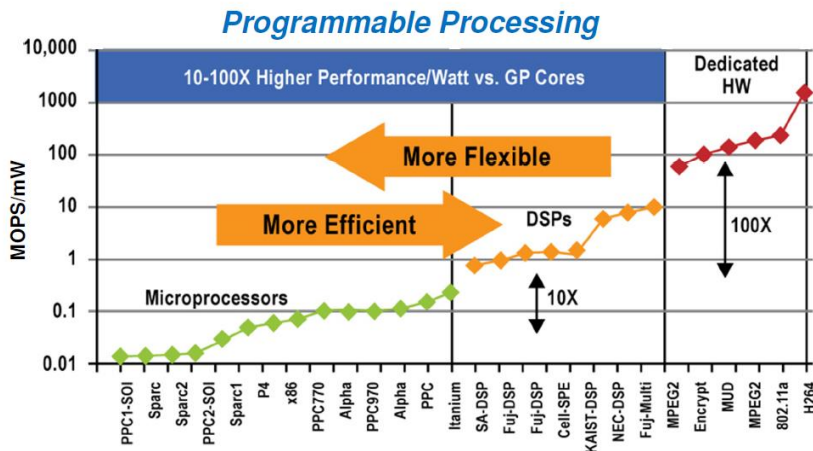
- Moore's law ends → Thermal constraints
- Dennard's scaling ends → Gains from multiprocessor architectures slow down
- Hennessy & Patterson 2018 Turing Lecture's solution: Domain Specific Architectures → Hardware Accelerators

Hardware Accelerators

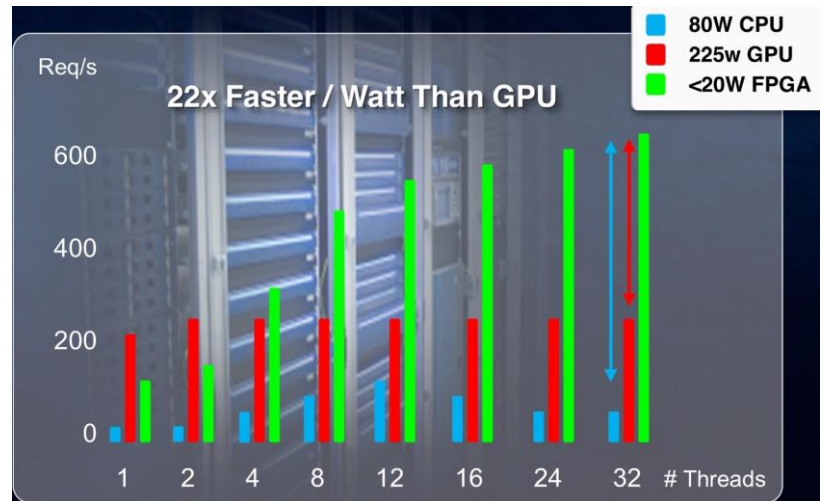


- Specialized hardware instead of general purpose hardware
- Performance and energy-efficiency improvements

The Dilemma: Flexibility vs. Efficiency



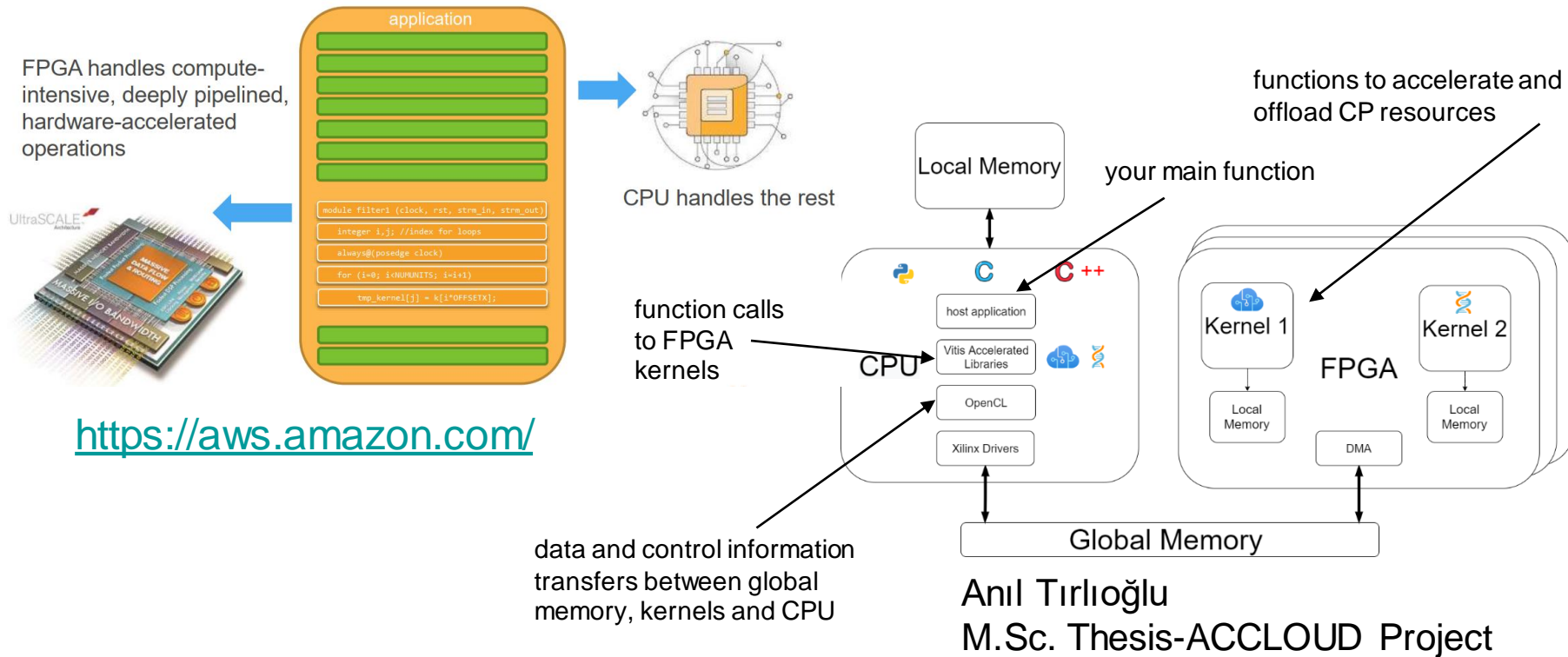
Source: "High-performance Energy-Efficient Reconfigurable Accelerator Circuits for the Sub-45nm Era" July 2011 by Ram K. Krishnamurthy, Circuits Research Labs, Intel Corp.



[Source: Xilinx, 2016]

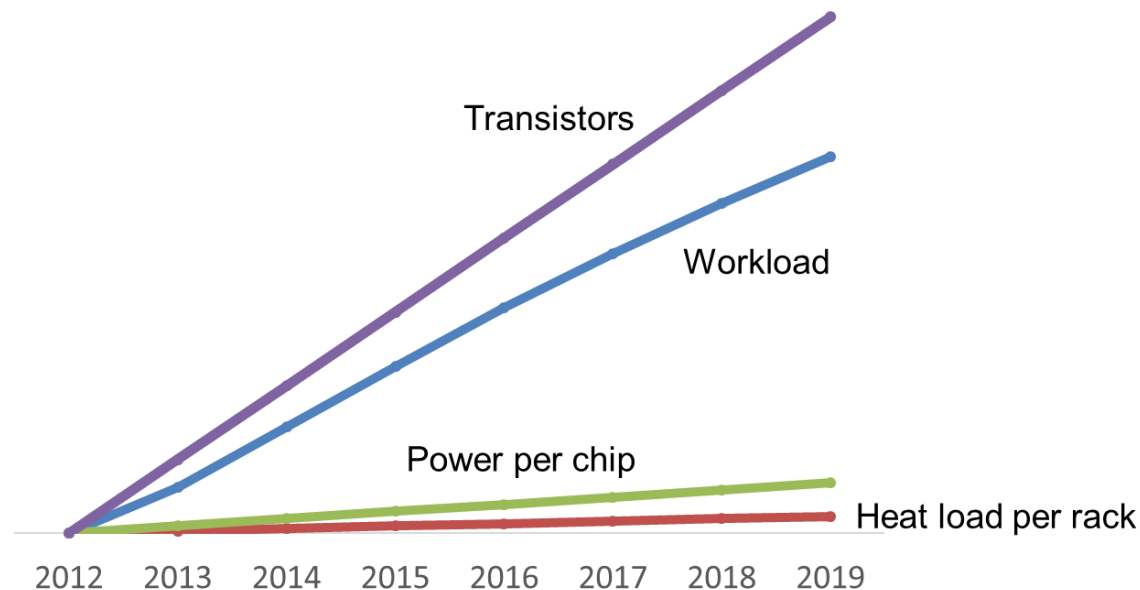
Kachris, Christoforos, and Dimitrios Soudris. "A survey on reconfigurable accelerators for cloud computing." *2016 26th International conference on field programmable logic and applications (FPL)*. IEEE, 2016.

Hardware Accelerators: Development and Use



State of Cloud Data Center

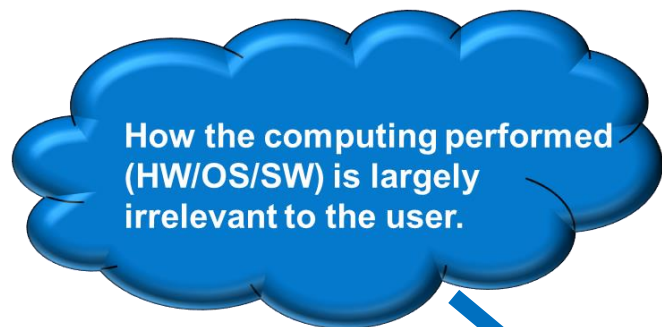
- Workload and Transistors increase fast
- Power and heat budget stay the same



Kachris, Christoforos, and Dimitrios Soudris. "A survey on reconfigurable accelerators for cloud computing." *2016 26th International conference on field programmable logic and applications (FPL)*. IEEE, 2016.

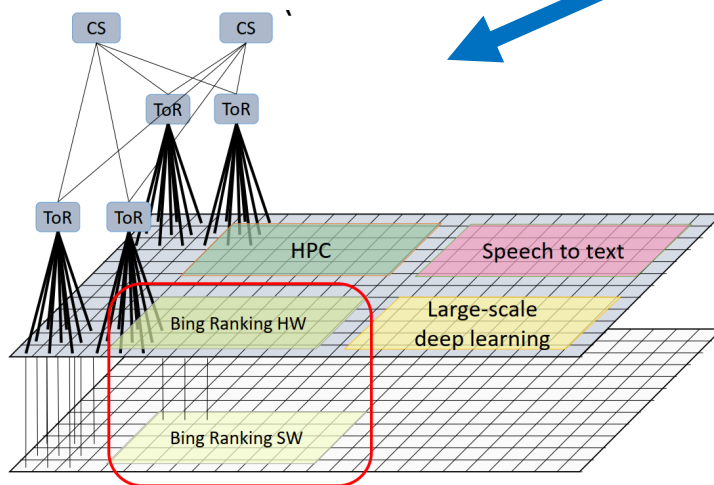


Part III: Hardware Accelerated Cloud Data Centers



Xilinx Versal
ACAP
Adaptive
Compute
Acceleration
Platform

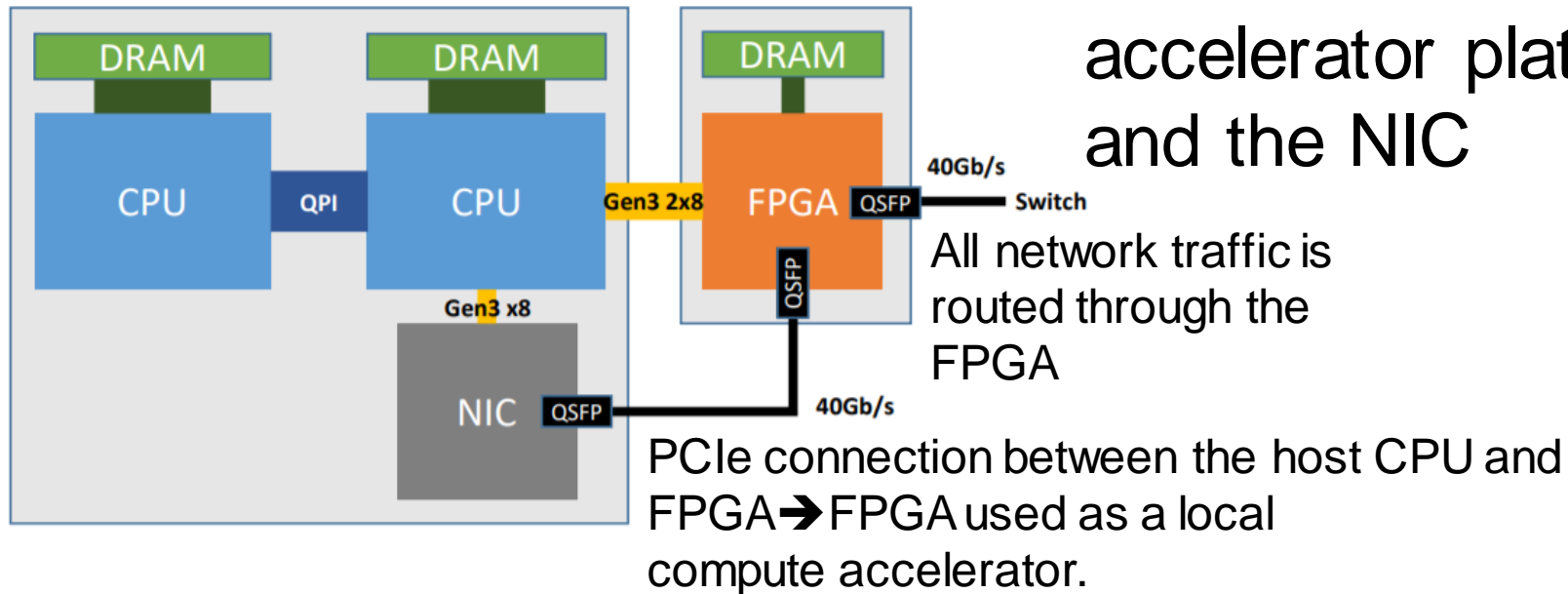
<https://www.xilinx.com/products/silicon-devices/acap/versal.html>



Hardware Accelerated
Cloud Data Center

Hardware Accelerated Cloud Data Centers

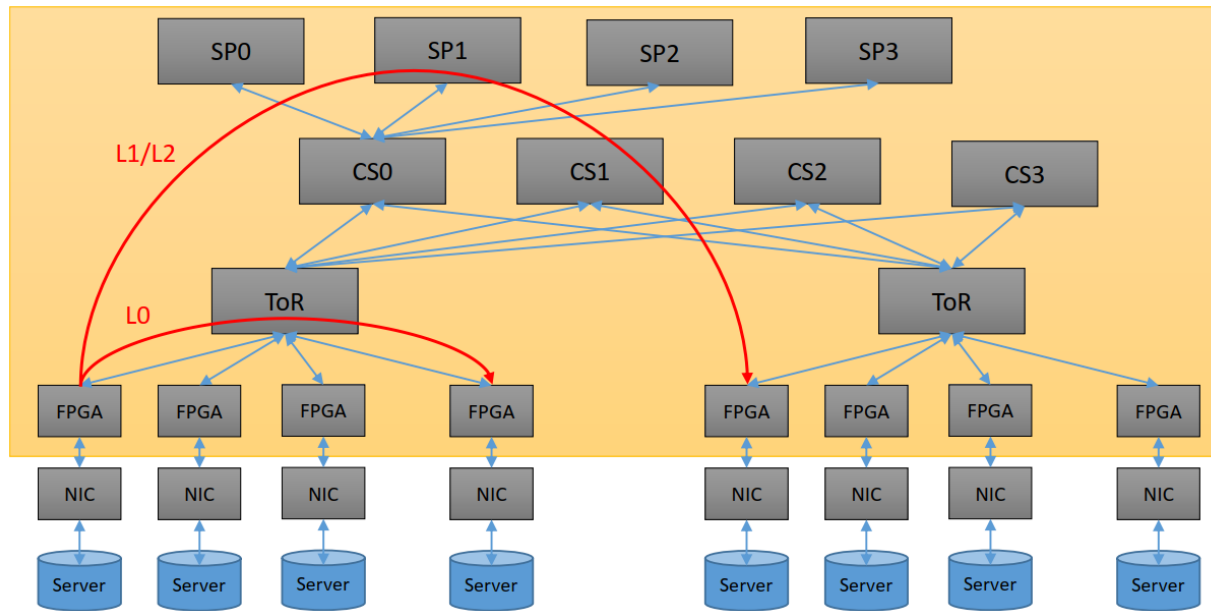
- Microsoft Catapult Project
- FPGA is both the accelerator platform and the NIC



Caulfield, Adrian M., et al. "A cloud-scale acceleration architecture." *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 2016.

Hardware Accelerated Cloud Data Centers

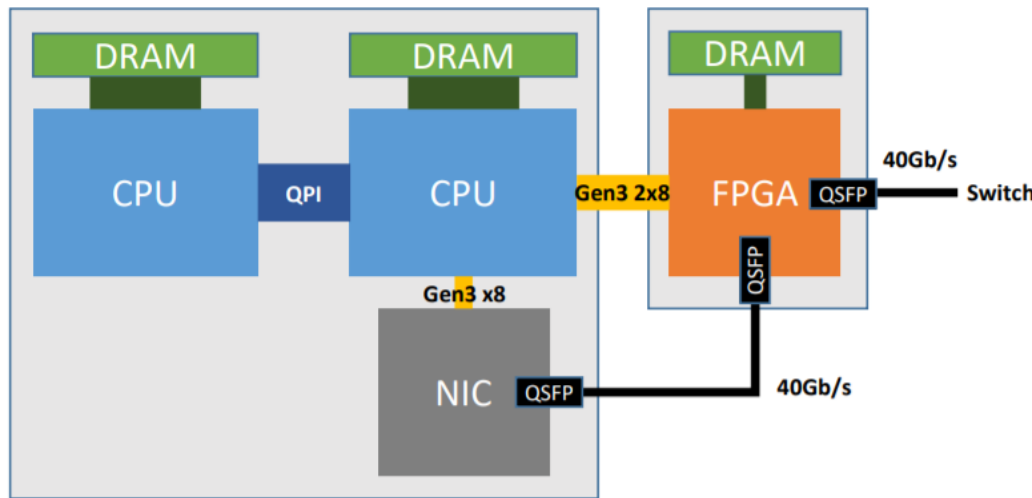
- Microsoft Catapult Project



- Low-latency inter-FPGA communication (Light Transport Layer)

Hardware Accelerated Cloud Data Centers

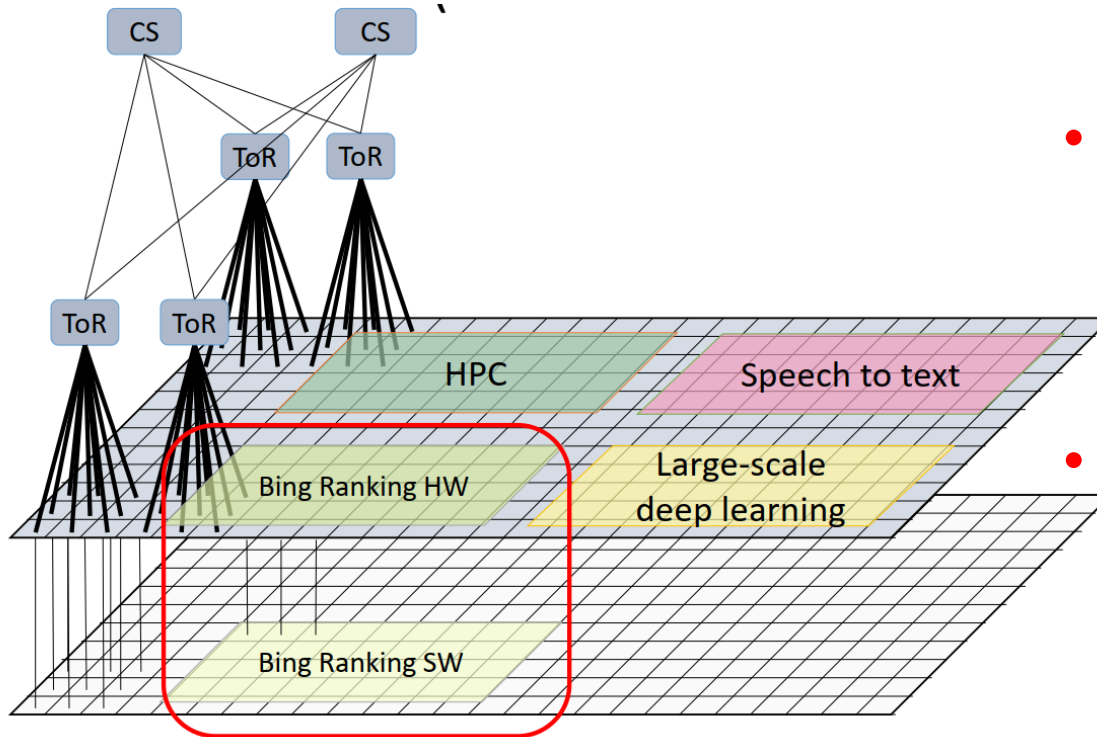
- Microsoft Catapult Project



- Local compute accelerator
- Network/storage accelerator
- Remote compute accelerator

Hardware Accelerated Cloud Data Centers

- Microsoft Catapult Project



- Hardware Acceleration as a Service Across Data Center (or even across Internet)
- FPGA is independent of the server

Part IV: ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

- TUBITAK Funded 1003 Research Project (to finish in April 2021)
- METU and Aselsan are partners
- Participation of many graduate students



<http://accloud.eee.metu.edu.tr/>

ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

Accelerator
as a Service



Optimal, accelerator
aware resource
allocation

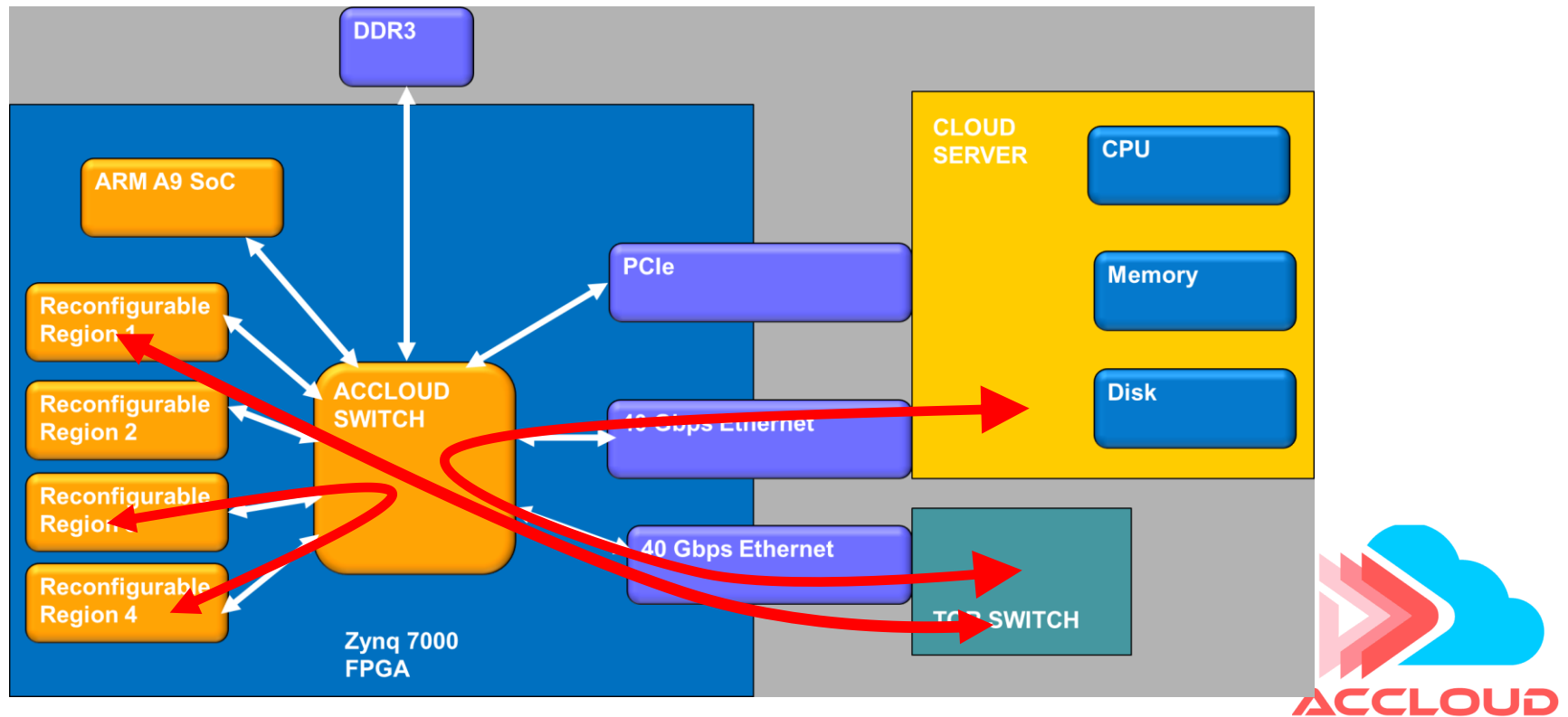
Accelerator implementation
on FPGA reconfigurable
regions

Transparent allocation
of accelerators as
Virtual Machine
parameters

On-chip switch
architecture for
interconnecting
hardware modules.

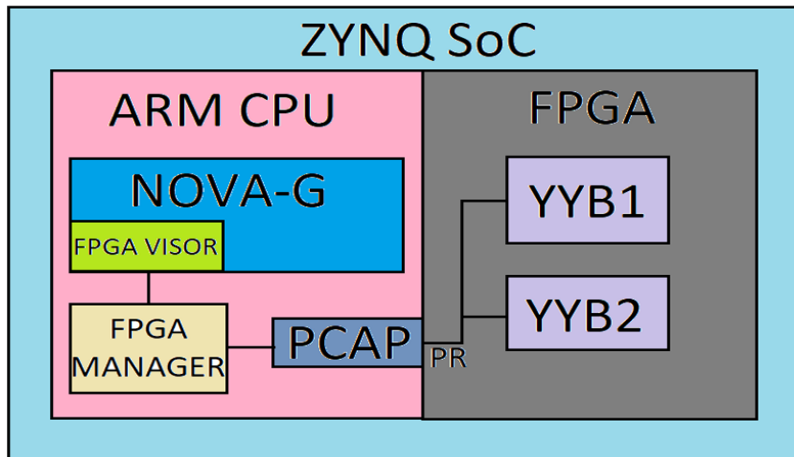
ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

ACCLOUD FPGA Accelerator and Cloud Server Layout



ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

Transparent allocation of accelerators as
Virtual Machine parameters



- Cloud Resource Management Framework for VM Creation
- Modification of Nova Compute component to allocate accelerators Similar to allocating CPU, RAM, Disk

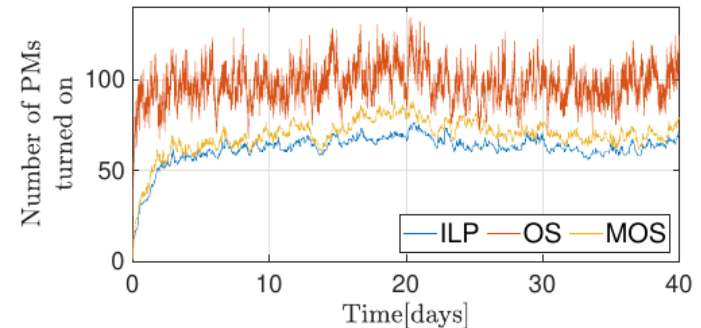
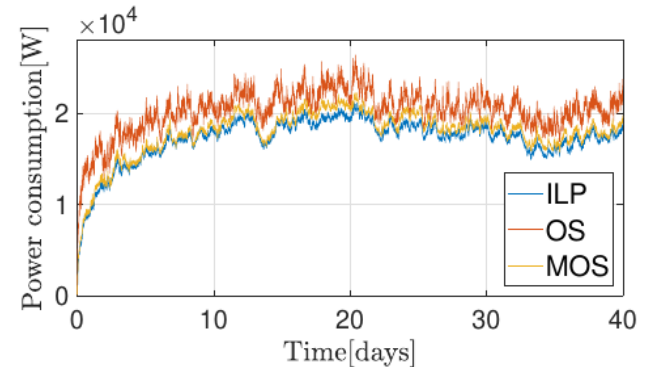
A. Erol, A. Yazar and E. G. Schmidt, "OpenStack Generalization for Hardware Accelerated Clouds," 2019 28th International Conference on Computer Communication and Networks (ICCCN), Valencia, Spain, 2019.



ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

Optimal, accelerator aware resource allocation: ACCLOUD-MAN

- Defining alternatives for SaaS requests
- Example: Video processing
 - 4 CPU cores
 - or 2 FPGA regions
 - or 2 CPU cores and 1 FPGA region
- Resource Allocation with minimum number Physical Machines, minimum power consumption



ILP: ACCLOUD-MAN

OS: Legacy
OpenStack

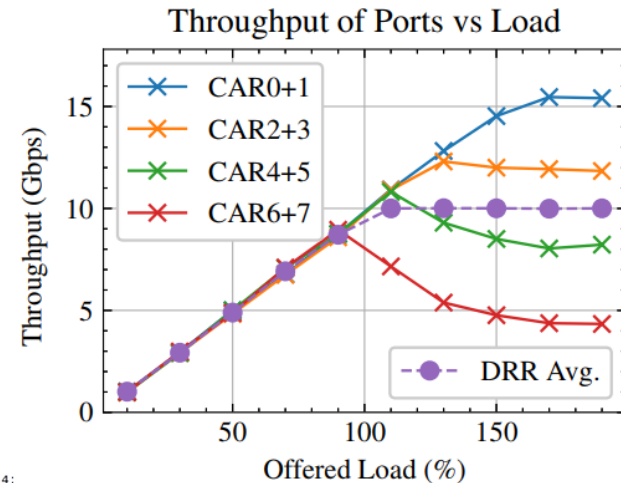
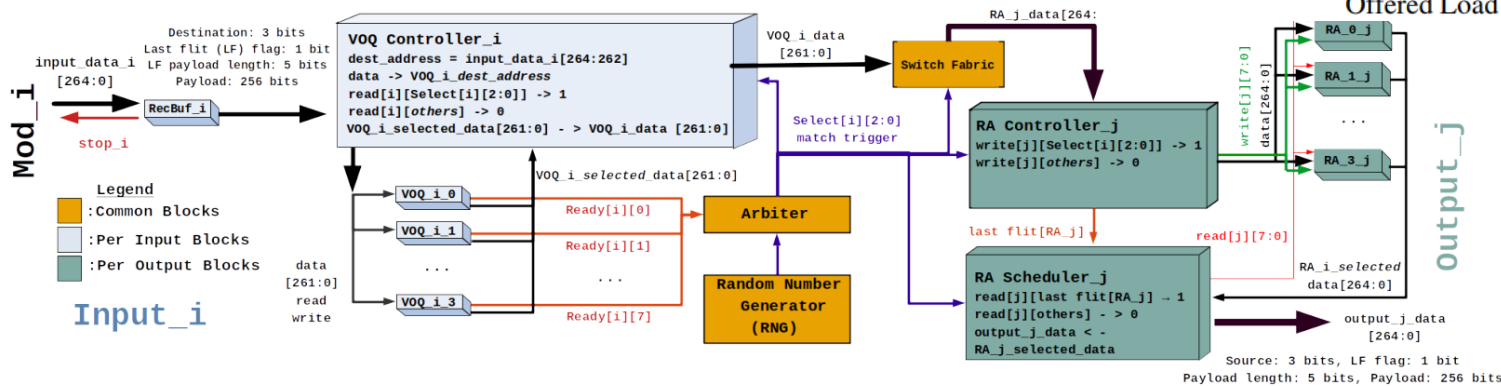
N. U. Ekici, K. W. Schmidt, A. Yazar and E. G. Schmidt, "Resource Allocation for Minimized Power Consumption in Hardware Accelerated Clouds," *2019 28th International Conference on Computer Communication and Networks (ICCCN)*, Valencia, Spain, 2019.



ACCLOUD (Accelerated Cloud): A New Cloud Architecture with FPGA Acceleration

On-chip switch architecture for interconnecting hardware modules: ACCLOUD-SWITCH

- 40 Gbps crossbar fabric
- 256 bit payload/flits
- Novel fabric arbitration with QoS



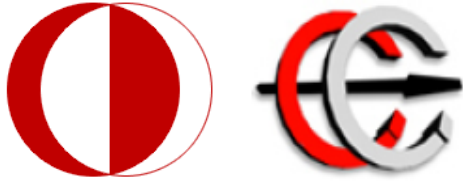
F. Yazıcı, A. S. Yıldız, A. Yazar, E. G. Schmidt, "A Novel Scalable On-chip Switch Architecture with Quality of Service Support for Hardware Accelerated Cloud Data Centers," *IEEE International Conference on Cloud Networking*, 2020.



Concluding Remarks

- Exploiting the *golden age* of hardware acceleration (as put by Henessy and Patterson)
- Seamlessly offering hardware resources to achieve more power efficient and higher performance services
- Wonderful research opportunities with many interesting problems!





Cloud Computing and Hardware Accelerated Clouds

Ece Güran Schmidt

METU Electrical and Electronics
Engineering
