

Enabling Technology for the Cloud and AI – One Size Fits All?



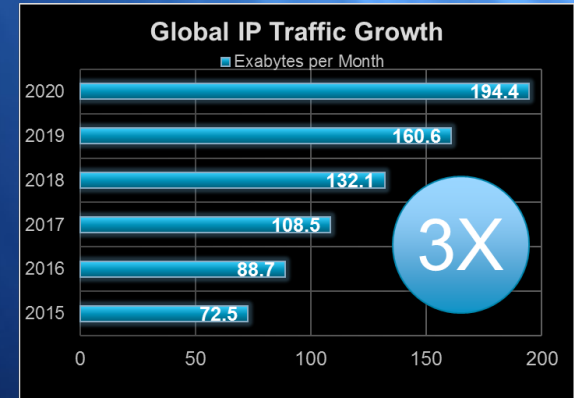
eSilicon®

Collaborate. Differentiate. Win.

Tim Horel

DIRECTOR, FIELD APPLICATIONS

The Growing Cloud – Global IP Traffic Growth



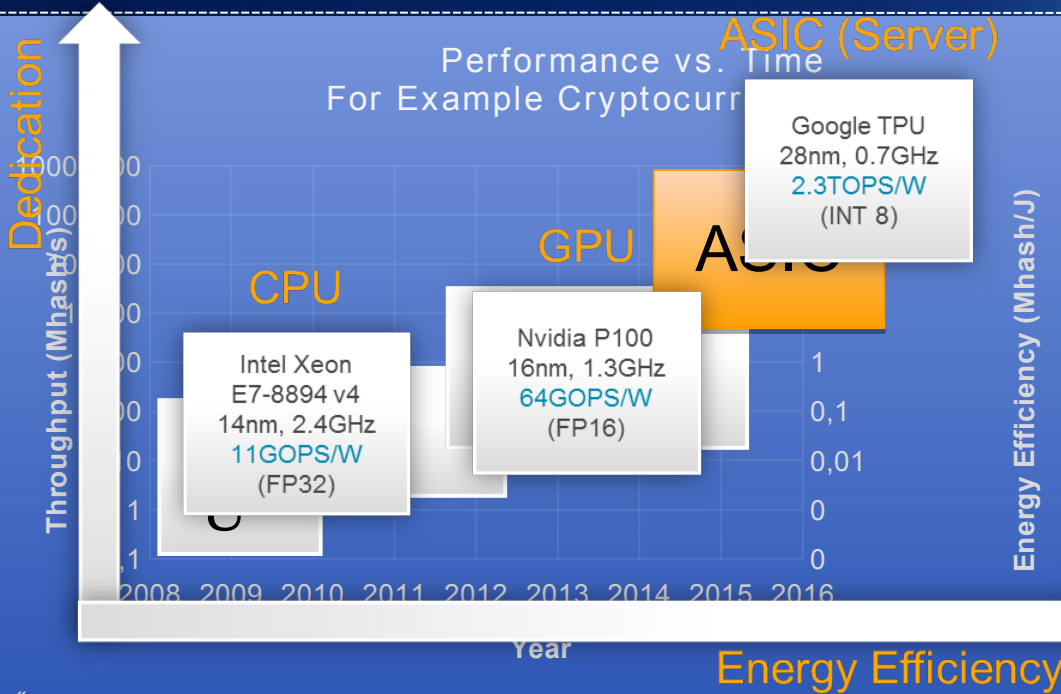
Reference: Cisco VNI Global IP Traffic Forecast, 2015 - 2020

40B+ devices with intelligence will ship in 2020

1B+ additional consumers online, 30B+ connected devices

~200 Exabytes of data traffic per month

ASIC Offers Performance Improvement



- Performance measurement:
 - Computational capacity (or throughput)
 - Energy-efficiency (computations per Joule)
 - Cost-efficiency (throughput per dollar)
- Computational units:
 - Single-core CPU: 1
 - Multi-core CPU: 10
 - GPU: 100
 - FPGA: 1,000
 - **ASIC: 10,000 to 1,000,000**

“The Future of Machine Learning Hardware”, Phillip Jama, Sept. 2016,
<https://hackernoon.com/the-future-of-machine-learning-hardware-c872a0448be8>

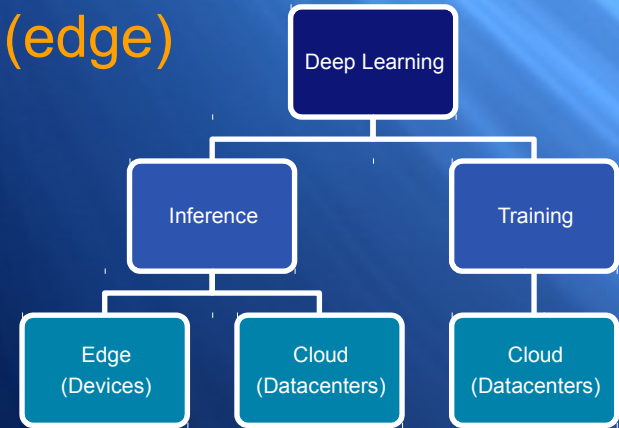
These numbers are based on the performance factors (such as throughput, efficiency) observed through the cryptocurrency-mining evolution.

- Non-specialized hardware comparison https://en.bitcoin.it/wiki/Non-specialized_hardware_comparison
- Mining hardware comparison https://en.bitcoin.it/wiki/Mining_hardware_comparison

Analytics in the Cloud with Artificial Intelligence

Deep learning is deployed today in all major datacenters (cloud) and in many devices (edge)

- Flood of data – images, video, text, transactions, speech
- Hardware acceleration enables and scales deep learning computation

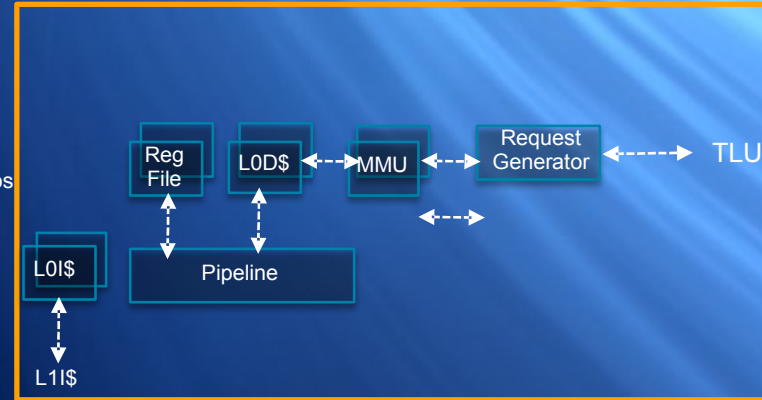
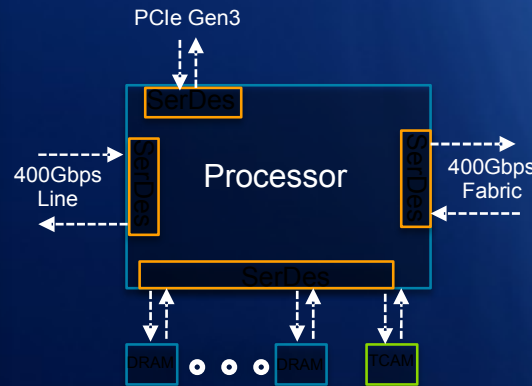
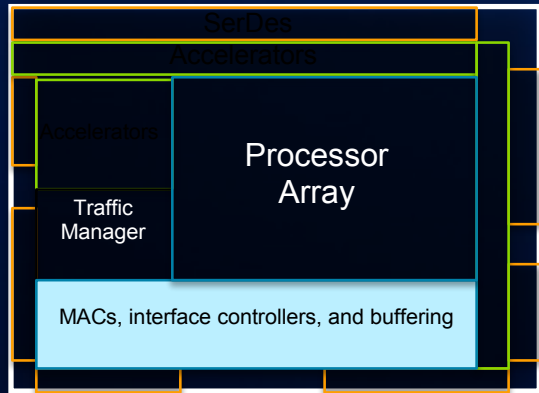


Source: NTT Group



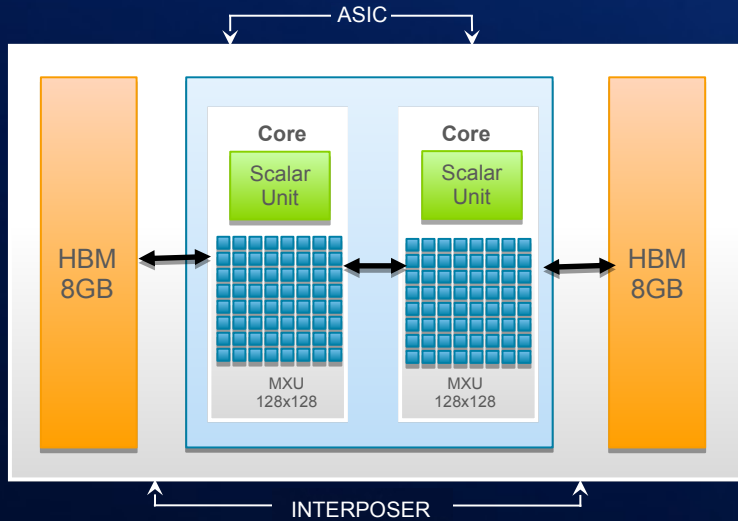
Source: Datacenter Frontier

Example - 400Gbps Multi-Core Network Processor



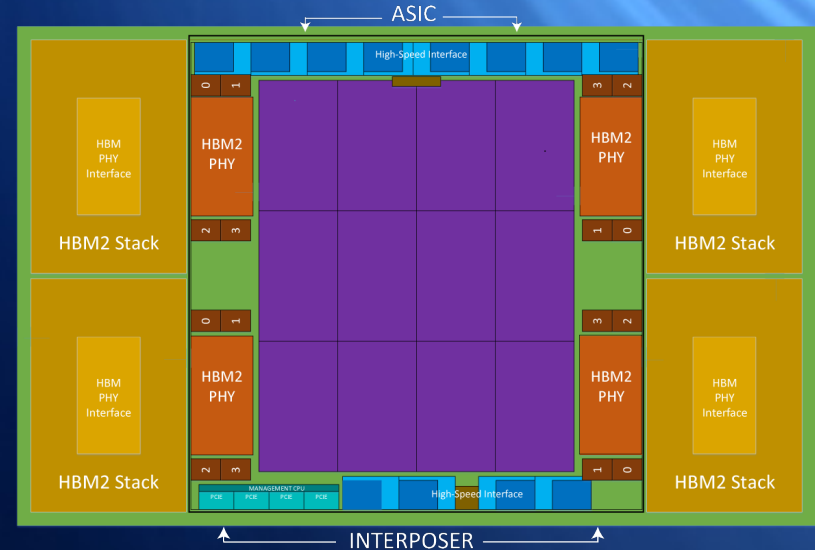
- 9.2 billion transistors; 643mm² die
- 672 total processors: 42 processor clusters
 - Link out to a four-way set associative L2 instruction
 - On-chip interconnect that links the clusters and caches to each
- 276 SerDes
- Interconnect runs at 1 GHz
- 9 Tb/sec of aggregate bandwidth
- 353 Mb of SRAM per core
 - L0 data and instruction cache for each thread
 - L1 instruction and data caches (for a cluster of sixteen cores)
- External DRAM for large data structures and packet buffering
- External TCAM for large data structures
- Integrated Ethernet MACs from 10GE to 100GE

Example AI Processors In The Market Today



- MXU: 32b float accumulation but reduced precision for multipliers
- 45 TFLOPS
- 16GB HBM
- 600GB/s memory bandwidth
- Scalar unit: 32b float

“Case Study on the Google TPU and GDDR5 from Hot Chips 29”,
Patrick Kennedy, Aug. 2017



- Tensor-based architecture
- Flexpoint® – high parallelism, low power
- 4X HBM2 memory stacks

“<https://www.anandtech.com/show/11942/intel-shipping-nervana-neural-network-processor-first-silicon-before-year-end>”, Nate Oh, October 2017

Similar Yet Different

Immature

General purpose chips offer a costly rapid-to-deployment option

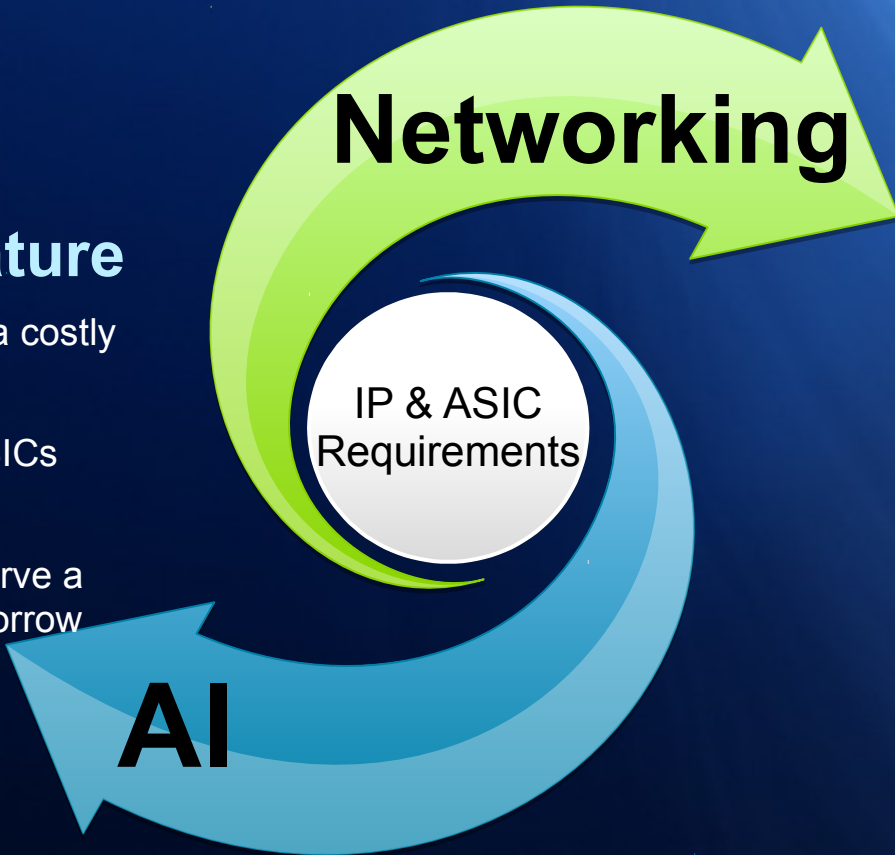
Evolving algorithms make ASICs harder to define

ASICs designed today will serve a more narrow application tomorrow

Networking

Mature

Well established architectures with products that have migrated for multiple technology generations



IP & ASIC Requirements

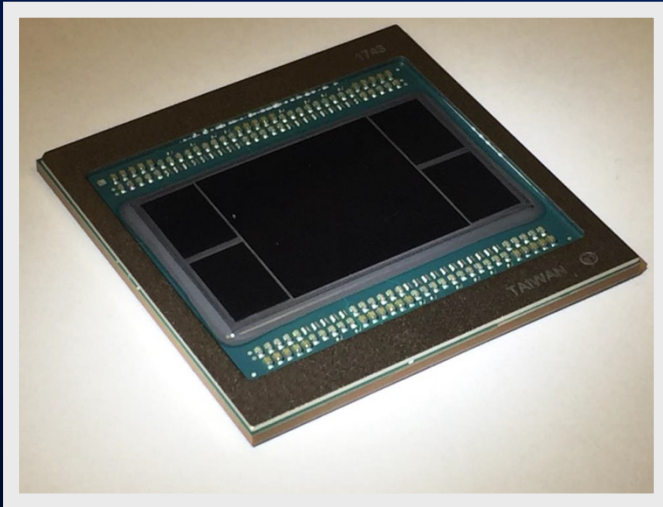
AI

Networking ASIC with High-Bandwidth Memory

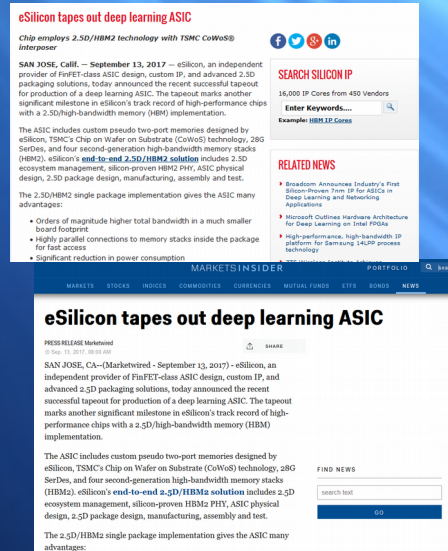


- FinFET technology
- 150M gates
- 900Mb embedded SRAM
- One 4-high HBM2 stack
- 60 lanes multi-protocol SerDes, 1 to 28Gb
- HBM2 PHY
- Memory compilers
- Specialty I/O
- Custom memory instances optimized for power/performance/area
- TCAM (supports 1.2G search/second [GSPS])

Deep Learning ASIC with High Bandwidth Memory and Custom Memory



- 8 Terabits/s of memory access speed
- ~400Mb embedded SRAM
- 4 HBM2 stack - 1TB/s HBM bandwidth
- 2.5D Packaging
- 28Gb SerDes
- HBM Gen2 PHY
- Memory compilers
 - 2PRF, DP SRAM, Fast Cache
- Deep Learning specialty SRAM
 - Custom Pseudo 2-Port SRAM



eSilicon tapes out deep learning ASIC

Chip employs 2.5D/HBM2 technology with TSMC CoWoS® interposer

SAN JOSE, Calif. — September 13, 2017 — eSilicon, an independent provider of FinFET-class ASIC design, custom IP, and advanced 2.5D packaging solutions, today announced the recent successful tapeout for production of a deep learning ASIC. The tapeout marks another significant milestone in eSilicon's track record of high-performance chips with a 2.5D/high-bandwidth memory (HBM) implementation.

The ASIC includes custom pseudo two-port memories designed by eSilicon, TSMC's Chip on Wafer on Substrate (CoWoS) technology, 28G SerDes, and four second-generation high-bandwidth memory stacks (HBM2). eSilicon's end-to-end 2.5D/HBM2 solution includes 2.5D ecosystem management, silicon-proven HBM2 PHY, ASIC physical design, 2.5D package design, manufacturing, assembly and test.

The 2.5D/HBM2 single package implementation gives the ASIC many advantages:

- Orders of magnitude higher total bandwidth in a much smaller board footprint.
- Highly parallel connections to memory stacks inside the package for fast access.
- Significant reduction in power consumption.



News Feed Item

eSilicon tapes out deep learning ASIC

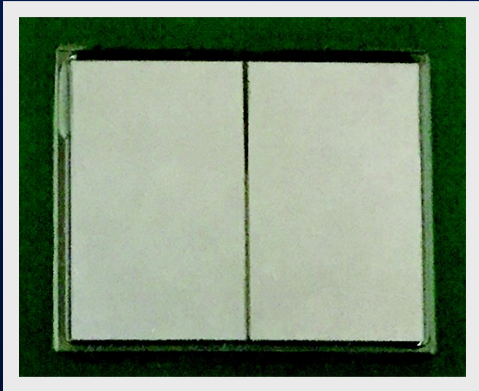
Chip employs 2.5D/HBM2 technology with TSMC CoWoS® interposer

SAN JOSE, Calif. — September 13, 2017 — eSilicon, an independent provider of FinFET-class ASIC design, custom IP, and advanced 2.5D packaging solutions, today announced the recent successful tapeout for production of a deep learning ASIC. The tapeout marks another significant milestone in eSilicon's track record of high-performance chips with a 2.5D/high-bandwidth memory (HBM) implementation.

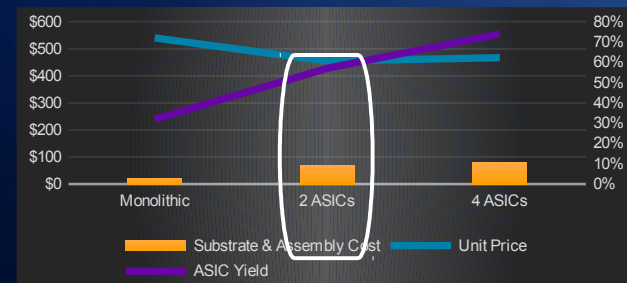
The ASIC includes custom pseudo two-port memories designed by eSilicon, TSMC's Chip on Wafer on Substrate (CoWoS) technology, 28G SerDes, and four second-generation high-bandwidth memory stacks (HBM2). eSilicon's end-to-end 2.5D/HBM2 solution includes 2.5D ecosystem management, silicon-proven HBM2 PHY, ASIC physical design, 2.5D package design, manufacturing, assembly and test.

The 2.5D/HBM2 single package implementation gives the ASIC many advantages:

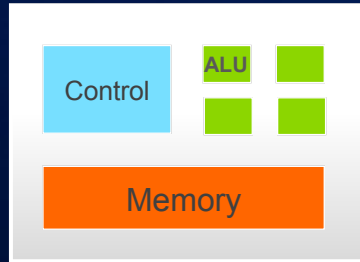
Economics of Billions of Transistors: Divide and Conquer



- FinFET technology
- Die-to-die interconnect on interposer running at >2.5Gbps per pin (~6000 interconnects between two die)
- Leverages HBM PHY technology, adopted for die-to-die
- Pre-requisite: Architect for yield optimization

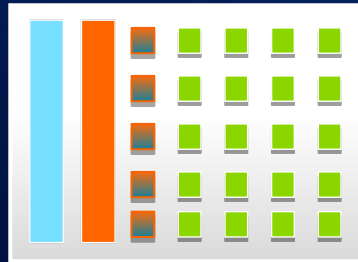


Processing Units – Rapid and Diverse Innovation Area



CPU

- Complex control logic
- Low compute density
- High programmability



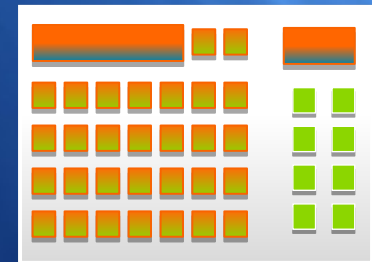
GPU

- Parallel computation
- High compute density
- Floating point datatype



ASIC

- Matrix multiply arch.
- Dedicated memory arch.
- Fixed point datatype



DNPU

- Heterogeneous arch.
- Fixed computation pattern
- Dedicated datatype

➔
Parallel
Processing

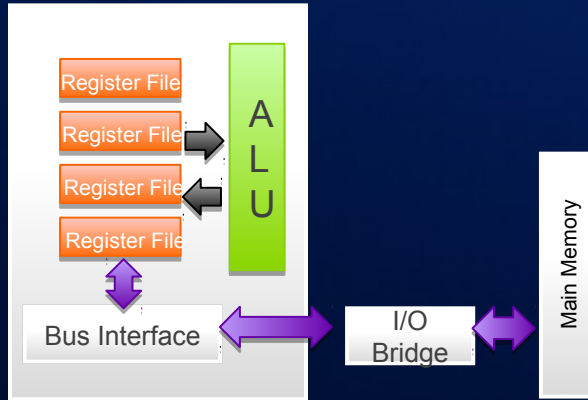
➔
Dedicated
Architecture

➔
Heterogeneous
Architecture

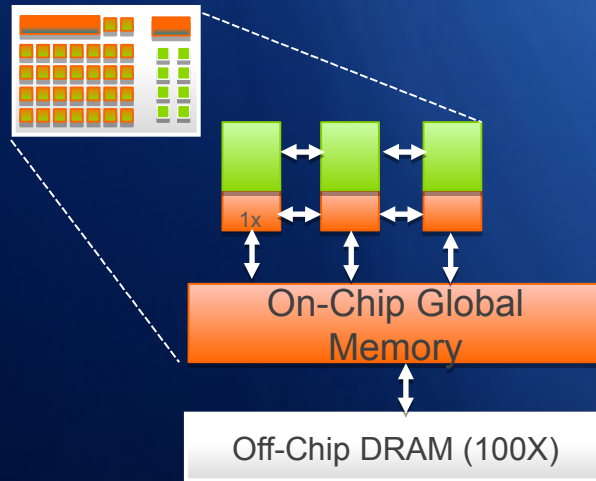
Reference: "DNPU: An Energy-Efficient Deep Neural Network Processor with On-Chip Stereo Matching", Dongjoo Shin, Jinmook Lee, Jinsu Lee, Juhyoung Lee, and Hoi-Jun Yoo, Semiconductor System Laboratory School of EE, KAIST

Latency v. Throughput at the Cost of Power

Traditional Compute Memory Subsystem



Today's Memory Subsystem

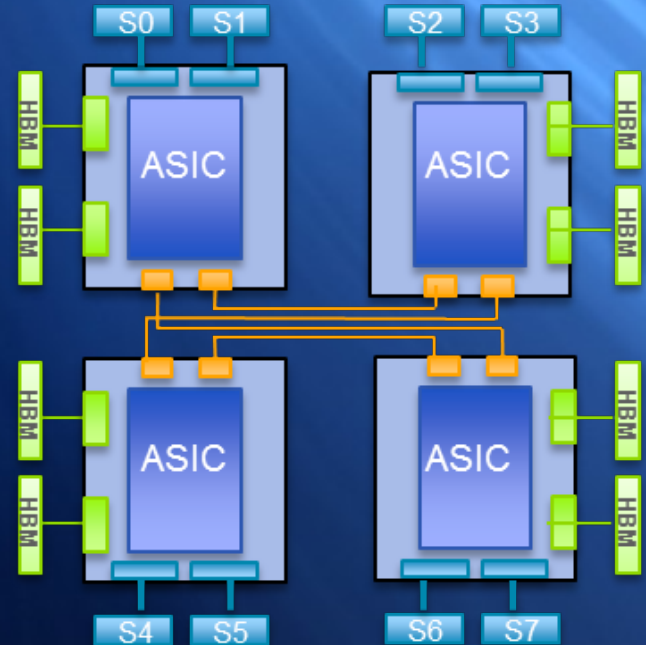


- Data movement has high energy cost
- Processing elements stalled waiting for data
- Processors optimize data flow to reduce energy cost
- *Off-chip memory is 100X less energy efficient than local memory*

- More on-chip heterogeneous SRAM
- Traditional “one size fits all” embedded SRAM is no longer applicable for both Networking and AI
- Operations can be done in parallel, however memory access and power is key bottleneck!

Interconnect – Lacing the Fabric

- Chip-to-chip interconnect
 - Short reach multi-chip modules
 - Packaging of dies with dissimilar semiconductor processes
 - Packaging of smaller dies to increase yield
 - No standard, proprietary solutions
- High throughput data interfaces:
 - Processor or switch to wide, high bandwidth memory
 - Off-board SerDes - 56G/112G SerDes PHY for long-reach backplane requirements for the 400 GB Ethernet
 - Switch to switch links



Rapid Pace of Innovation

Quick Evolution

- Short time to product to hit tight market windows
- Architectural flexibility to allow for incremental refinements
- High degree of re-use to shorten design cycle
- Swappable modules and heterogeneous technologies
- “What if” exploration capabilities



Virtual Acceleration

Unique IP offering based on mathematical functions integrated with memory

Highly parallel, embedded interconnects

Very scalable solutions across multiple dies, chiplets, packages and technologies



Shorten Time To Tape-out for AI

Library of Mega Cells

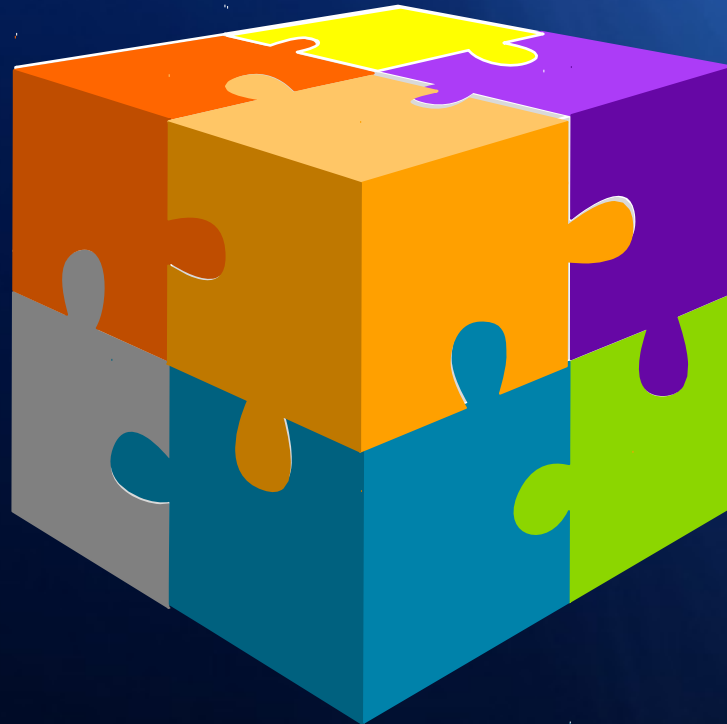
Configurable, hardened cells, referenced from RTL, specialized memory

On-Chip Interconnect scheme

Configurable NOC

High Speed Off-chip Interconnect

Configurable, hardened high-speed interconnect , 56/112G SerDes, HBI



Off-Chip Memory

HBM2 DRAM, PHY/controller hardened and pre-validated

2.5 D Packaging

Interposer design, signal & thermal integrity, package warpage

ASIC Toolbox & Recipe

Database and tools to analyze, assess and determine most optimal technology, logic and memory implementation



eSilicon®

Collaborate. Differentiate. Win.