# AMD STRATEGY IN EXASCALE SUPERCOMPUTING AND MACHINE INTELLIGENCE

TIMOUR PALTASHEV, D.SC.
SEPTEMBER 20, 2017

▰ Exascale Goals and Challenges

▰ AMD's Vision and Technologies for Exascale Computing

▰ HPC Progress Towards Machine Intelligence

▰ Radeon Instinct and Radeon Open Compute (ROC) Initiatives

▰ AMD Radeon Instinct Accelerators and Naples server SoC for HPC and Machine Intelligence

# DEPARTMENT OF ENERGY'S GOALS FOR EXASCALE COMPUTING SYSTEMS

**AMD**

- The Department of Energy (DOE) plans to deliver exascale supercomputers that provide a 50x improvement in application performance over their current highest-performance supercomputers by 2023

- System should provide a 50x performance improvement over today's fastest supercomputes with 20 MWatts of power while not requiring human intervention due to hardware or system faults more than once a week on average

- Important goals for exascale computing include
  - Enabling new engineering capabilities and scientific discoveries
  - Continuing U.S. leadership in science and engineering

https://asc.llnl.gov/pathforward/

http://science.energy.gov/~/media/ascr/ascac/pdf/meetings/20140210/Top10reportFEB14.pdf
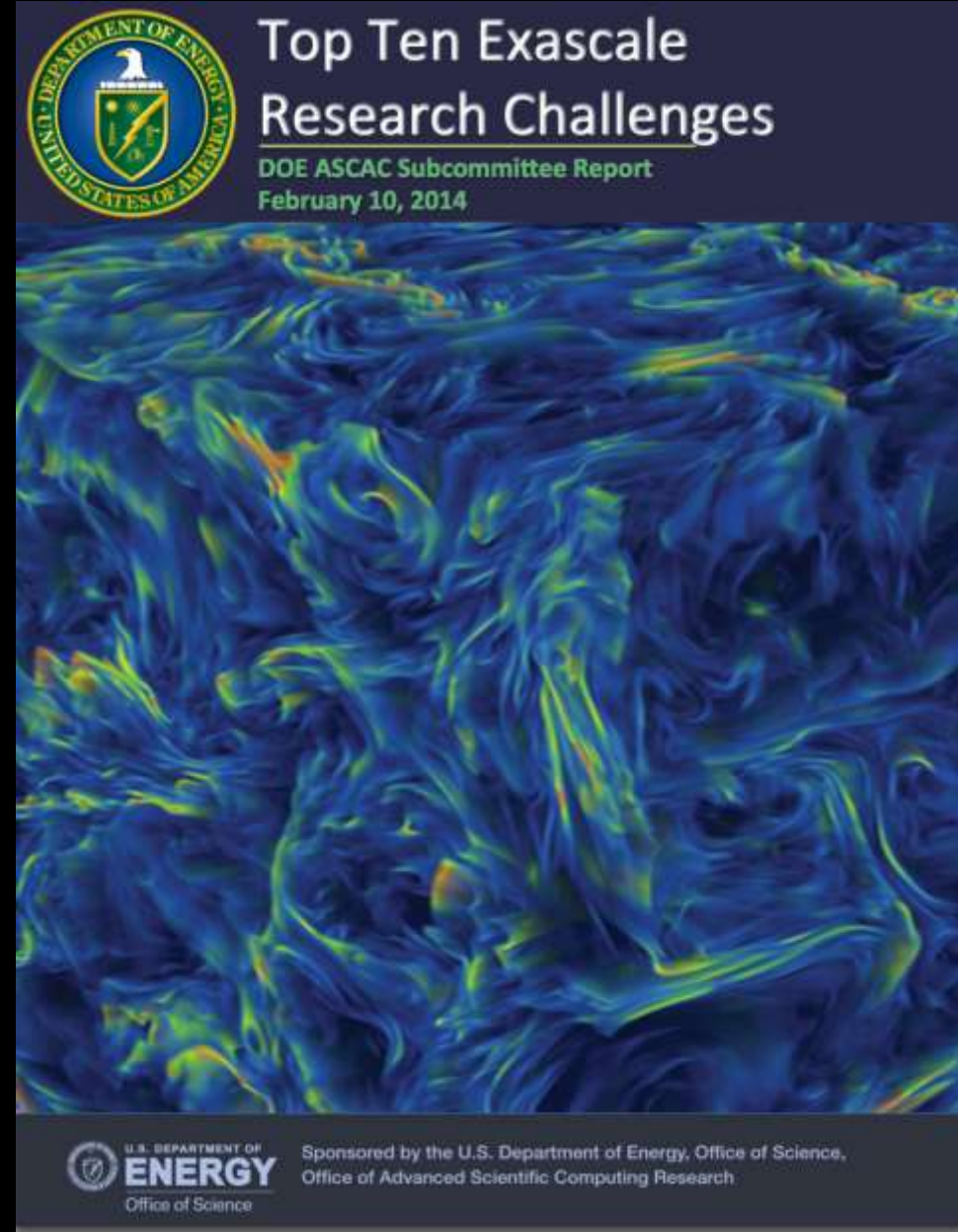
**RADEON** TECHNOLOGIES GROUP

# EXASCALE CHALLENGES

The Top Ten Exascale Research Challenges

1) Energy efficiency

2) Interconnect technology

3) Memory technology

4) Scalable system software

5) Programming systems

6) Data management

7) Exascale algorithms

8) Algorithms for discovery, design, and decision

9) Resilience and correctness

10) Scientific productivity

http://science.energy.gov/~/media/ascr/ascac/pdf/
meetings/20140210/Top10reportFEB14.pdf

Requires significant advances in processors, memory, software, and system design

Top Ten Exascale
Research Challenges
DOE ASCAC Subcommittee Report
February 10, 2014

U.S. DEPARTMENT OF
ENERGY
Office of Science

Sponsored by the U.S. Department of Energy, Office of Science,
Office of Advanced Scientific Computing Research

# DOE EXASCALE TARGET REQUIREMENTS

**AMD**

- ◢ The DOE has aggressive goals and target requirements for exascale systems
  - – Requires research and innovation in a variety of areas

- ◢ One of the most important goals is providing supercomputers that can be effectively utilized for important scientific discoveries

- ◢ Technologies explored for exascale can be applied to a wide variety of computing systems

| Target Requirements | Target Value |
| --- | --- |
| System-Level Power Efficiency | 50 GFLOPS/Watt |
| Compute Performance (per node) | 10 TFLOPS |
| Memory Capacity (per node) | 5TB |
| Memory Data Rate (per node) | 4 TB/sec |
| Message per Second (per node) | 500 million (MPI), 2 billion (PGAS) |
| Mean Time to Application Failure | 7 days |

▶ http://science.energy.gov/~/media/ascr/ascac/pdf/meetings/20140210/Top10reportFEB14.pdf

**RADEON** TECHNOLOGIES GROUP

# AMD'S VISION FOR SUPERCOMPUTING

**EMBRACING HETEROGENEITY**

**CHAMPIONING OPEN SOLUTIONS**

**ENABLING LEADERSHIP SYSTEMS**

| AMD STRATEGY IN EXASCALE  SUPERCOMPUTING AND MACHINE INTELLIGENCE    |    SEPTEMBER 20, 2017    NANO AND GIGA CHALLENGES , TOMSK, RUSSIAN FEDERATION
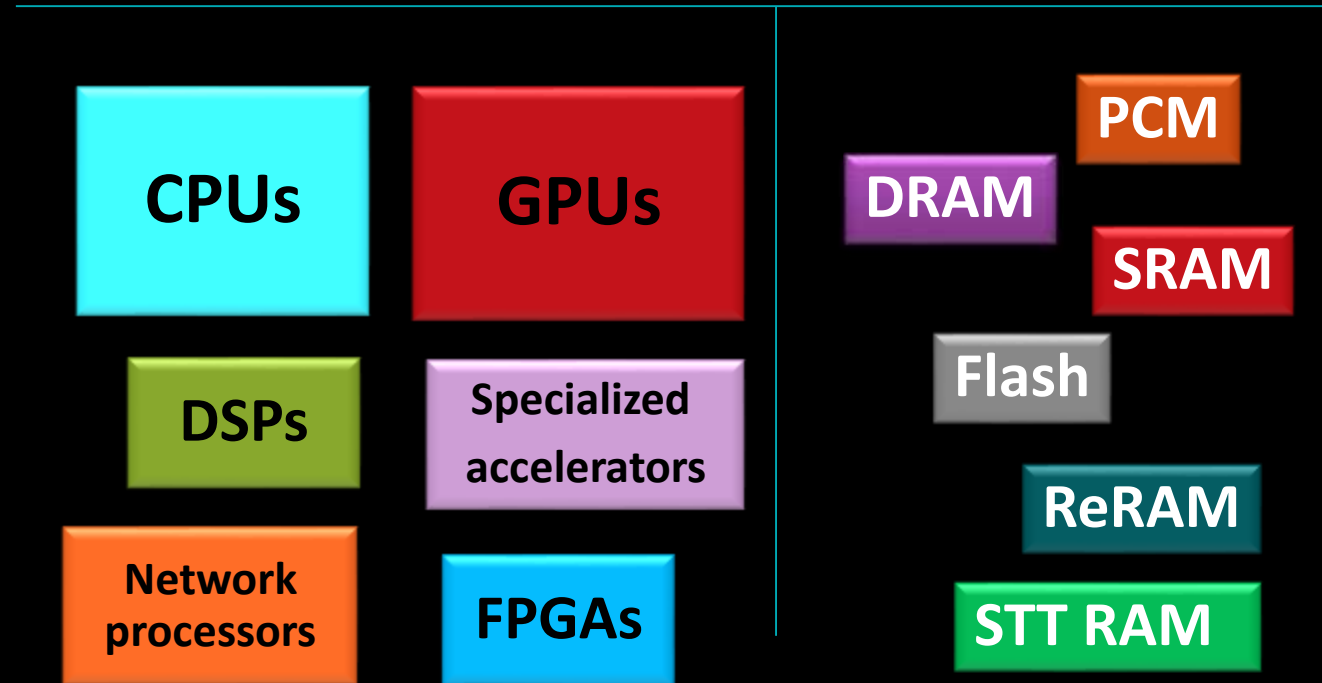
# EMBRACING HETEROGENEITY

**AMD**

- Customers must be free to choose the technologies that suit their problems
- Specialization is key to high performance and energy efficiency
- Heterogeneity should be managed by programming environments and runtimes
- The Heterogeneous System Architecture (HSA) provides:
  - A framework for heterogeneous computing
  - A platform for diverse programming languages

C/C++    FORTRAN    Java

UPC/UPC++    python    MPI

Kokkos/RAJA    OpenMP    OpenACC

CPUs    GPUs

DSPs    Specialized accelerators

Network processors    FPGAs

PCM
DRAM
SRAM
Flash
ReRAM
STT RAM

Heterogeneity Options

**RADEON** TECHNOLOGIES GROUP

# CHAMPIONING OPEN SOLUTIONS

**AMD**

◢ Harness the creativity and productivity of the entire industry

◢ Partner with best-in-class suppliers to enable leading solutions

◢ Multiple paths to open solutions
  – **Open standards**
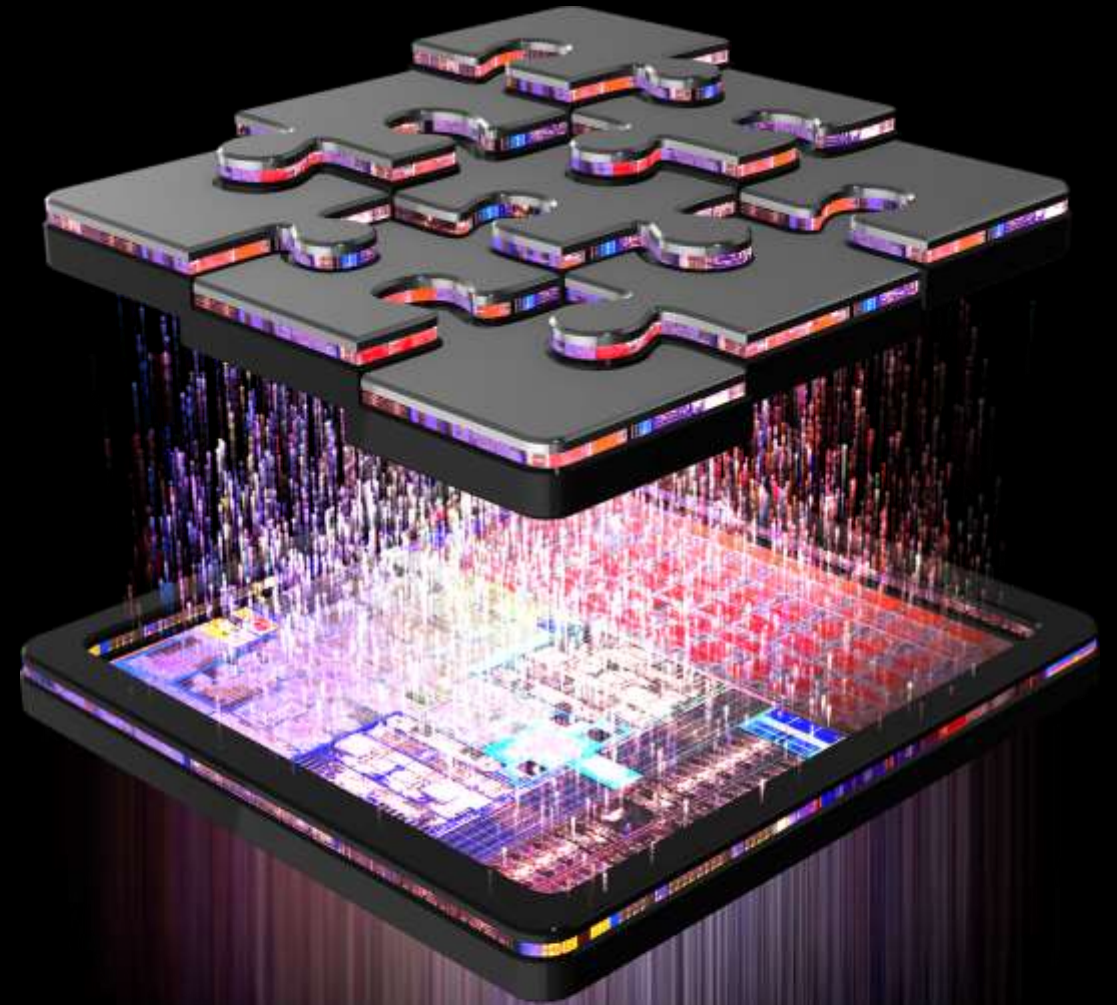  – **Open-source software**
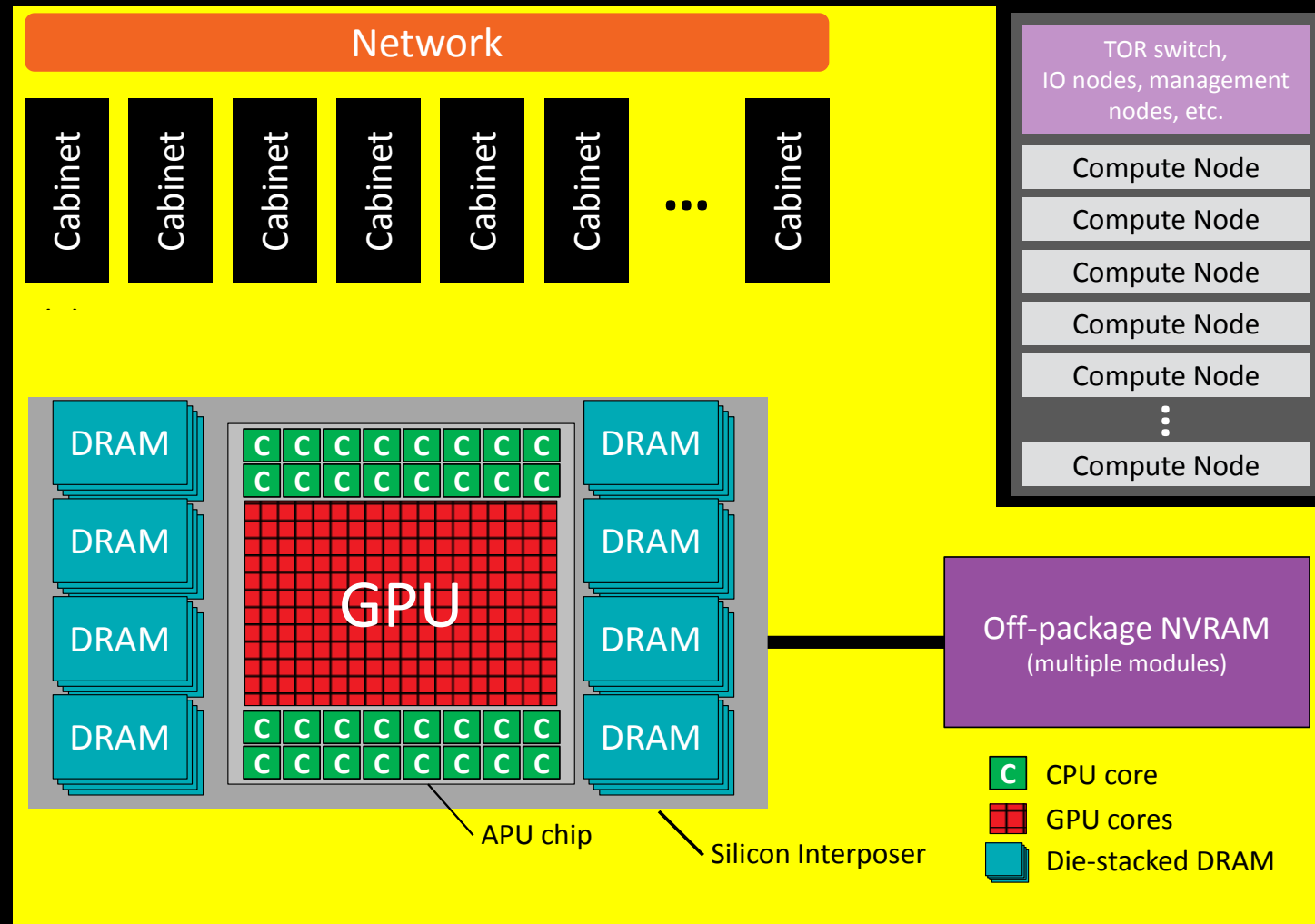  – **Open collaborations**

# ENABLING LEADERSHIP SYSTEMS

**AMD**

- **Re-usable, high-performance technology building blocks**

- **High-performance network on chip**

- **Modular engineering methodology and tools**

- **Software tools and programming environments**

**RADEON** TECHNOLOGIES GROUP

# FUTURE HIGH DENSITY COMPUTE CONFIGURATIONS

- Exascale systems require enhanced performance, power-efficiency, reliability, and programmer productivity
  - Significant advances are needed in multiple areas and technologies

- Exascale systems will be heterogeneous
  - Programming environments and runtimes should manage this heterogeneity

- New computing technologies provide a path to productive, power-efficient exascale systems



Network

Cabinet | Cabinet | Cabinet | Cabinet | Cabinet | Cabinet | ... | Cabinet

TOR switch, IO nodes, management nodes, etc.

Compute Node
Compute Node
Compute Node
Compute Node
Compute Node
Compute Node

DRAM | GPU | DRAM

Off-package NVRAM (multiple modules)

APU chip

Silicon Interposer

C — CPU core
GPU cores
Die-stacked DRAM

For further details see: "Achieving Exascale Capabilities through Heterogeneous Computing," IEEE Micro, July/August 2015.

# COMPUTING PROGRESS: CLIENT-SERVER

'00,000s    Units

Client-server

Main Frames

0

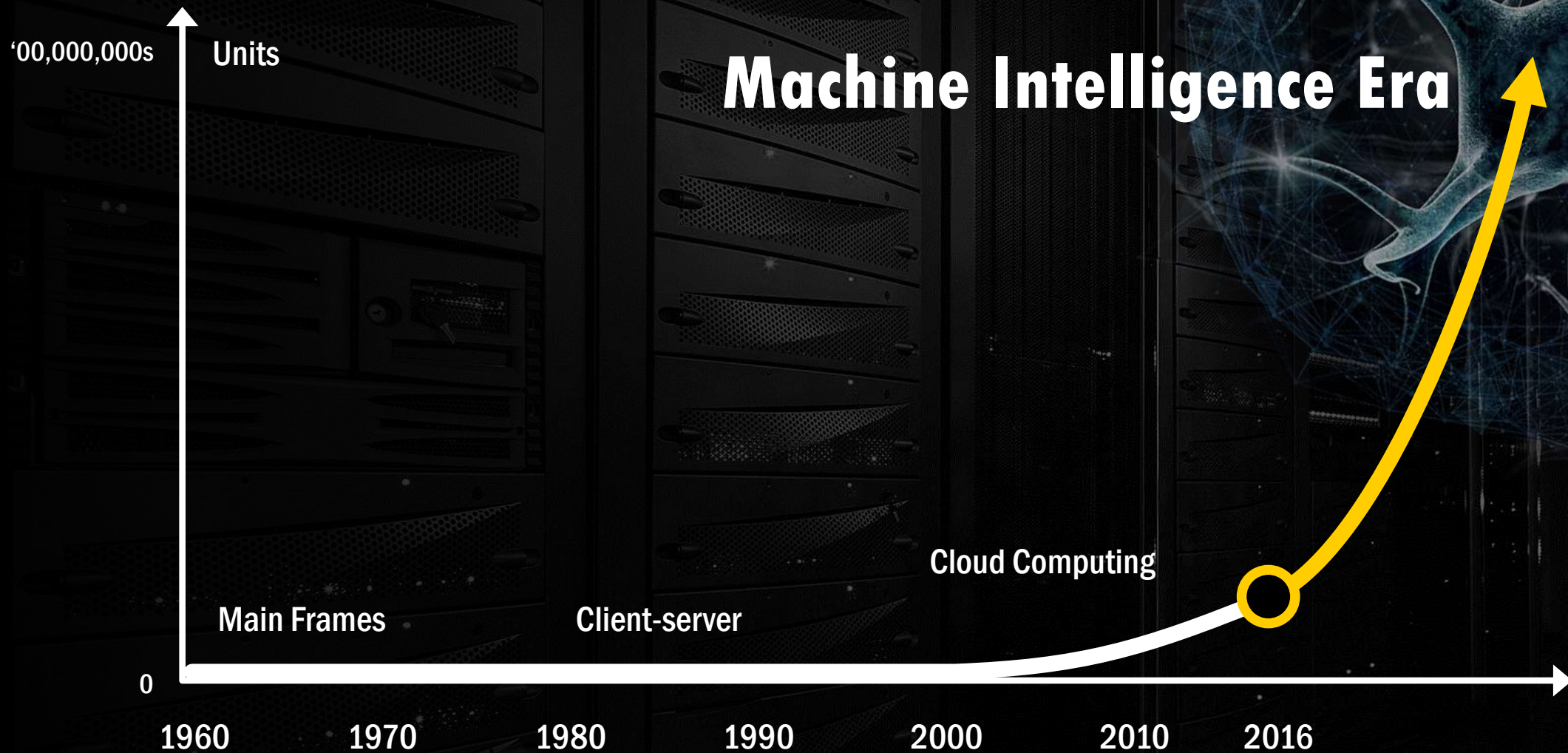1960    1970    1980    1990    2000    2010    2016

Chart for illustrative purposes

AMD

# 2.5 Quintillion Bytes
# of Data is Generated Every Day

**500 million** Tweets

**4 million** hours of content on Youtube

**4.3 billion** Facebook entry

**3.6 billion** Instagram

**6 billion** Google searches

**205 billion** emails

**and many more...**

AMD

# Human Brain in your Hand

# Radeon Instinct Initiative

Cloud / Hyperscale

Financial Services

Energy

Life Sciences

Automotive

Optimized Machine Learning / Deep Learning Frameworks and Applications

ROCm  Software Platform

Radeon™ Instinct Hardware Platform

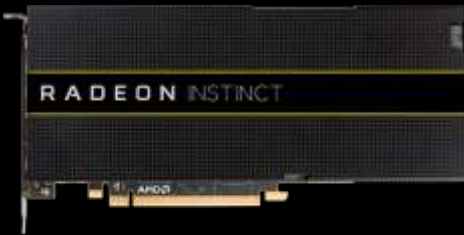## Address market verticals that use a common infrastructure to leverage the investments and  scale fast across multiple industries

18
| AMD STRATEGY IN EXASCALE  SUPERCOMPUTING AND MACHINE INTELLIGENCE   |   SEPTEMBER 20, 2017    NANO AND GIGA CHALLENGES , TOMSK, RUSSIAN FEDERATION

AMD | RADEON

# Accelerators | RADEON INSTINCT

## MI6

Passively Cooled Inference Accelerator

5.70 TFLOPS

224 GB/s Memory Bandwidth

<150W

## MI8

Small Form Factor Accelerator

8.2 TFLOPS

512 GB/s Memory Bandwidth

<175W

## MI25 Vega with NCU

Passively cooled Training Accelerator

2X Packed Math

High Bandwidth Cache and Controller

<300W

# ROCm PROGRAMMING MODEL OPTIONS

**AMD**

## HIP

*Convert CUDA to portable C++*

- Single-source Host+Kernel
- C++ Kernel Language
- C Runtime
- Platforms: AMD GPU, NVIDIA (Designed to have the same or better perf as native CUDA)

When to use it?
- Port existing CUDA code
- Developers familiar with CUDA
- New project that needs portability to AMD and NVIDIA

## HCC

*True single-source C++ accelerator language*

- Single-source Host+Kernel
- C++ Kernel Language
- C++ Runtime
- Platforms: AMD GPU

When to use it?
- New projects where true C++ language preferred
- Use features from latest ISO C++ standards

## OpenCL

*Khronos Industry Standard accelerator language*

- Split Host/Kernel
- C99-based Kernel Language
- C Runtime
- Platforms: CPU, GPU, FPGA

When to use it?
- Port existing OpenCL code
- New project that needs portability to CPU,GPU,FPGA

# INTRODUCING ROCm SOFTWARE PLATFORM

## A new, fully "Open Source" foundation for Hyper Scale and HPC-class GPU computing

**AMD**

### Graphics Core Next Headless Linux® 64-bit Driver

- Large memory single allocation
- Peer-to-Peer Multi-GPU
- Peer-to-Peer with RDMA
- Systems management API and tools

### HSA Drives Rich Capabilities Into the ROCm Hardware and Software

- User mode queues
- Architected queuing language
- Flat memory addressing
- Atomic memory transactions
- Process concurrency & preemption

**HSA™ FOUNDATION**

### Rich Compiler Foundation For HPC Developer

- LLVM native GCN ISA code generation
- Offline compilation support
- Standardized loader and code object format
- GCN ISA assembler and disassembler
- Full documentation to GCN ISA

### "Open Source" Tools and Libraries

- Rich Set of "Open Source" math libraries
- Tuned "Deep Learning" frameworks
- Optimized parallel programing frameworks
- CodeXL profiler and GDB debugging

**GPUOpen**

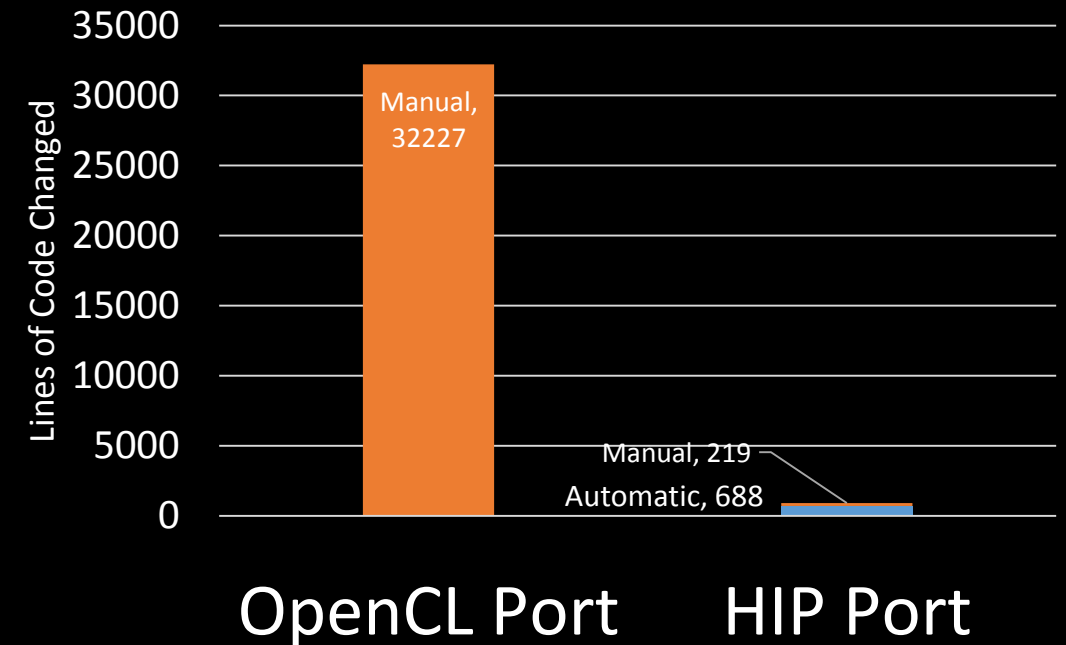# ROCm : DEEP LEARNING GETS HIP

Bringing a faster path to bring deep learning application to AMD GPUs

- The Challenge: CAFFE
  - Popular machine-learning framework
  - Tip version on GitHub has 55000+ lines-of-code
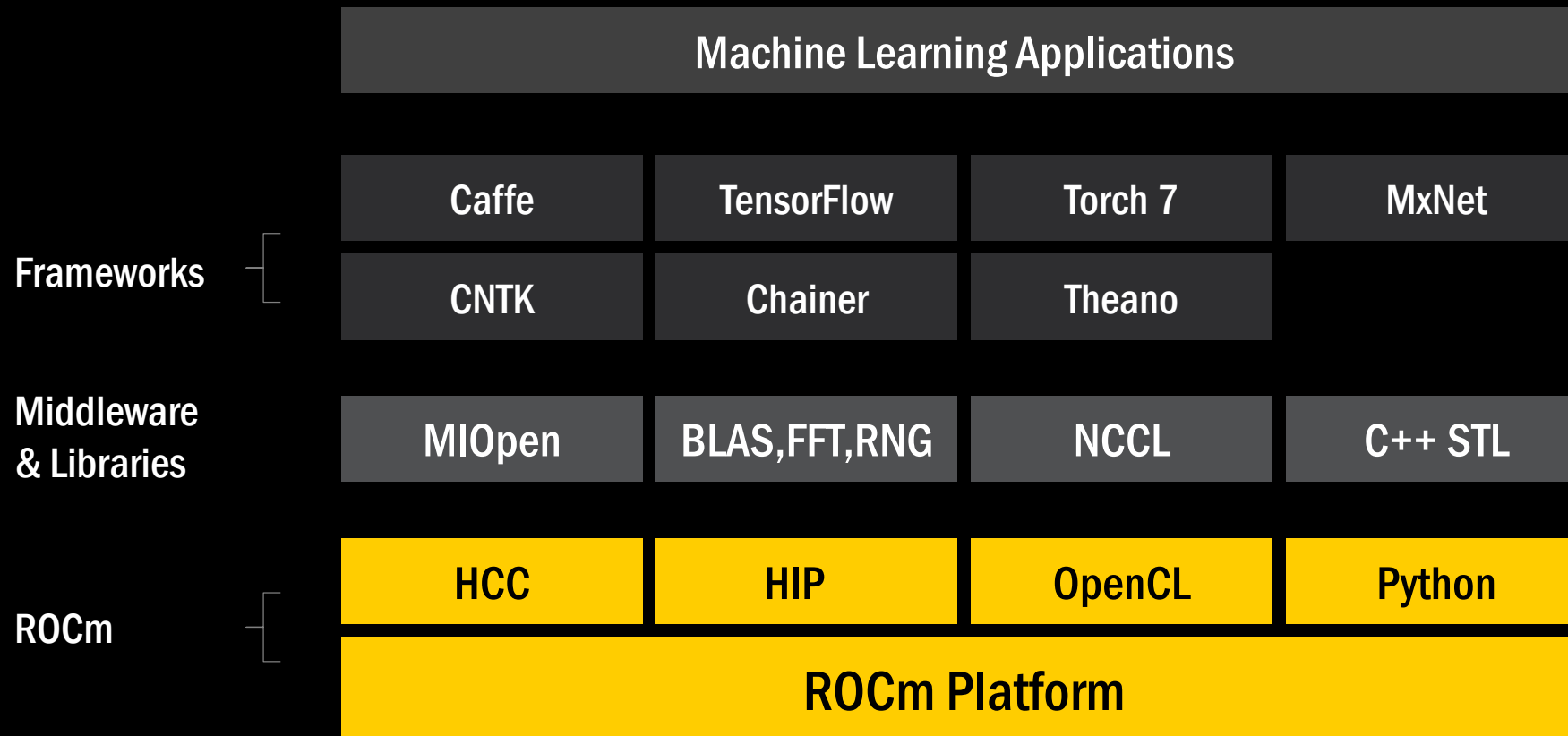  - GPU-accelerated with CUDA

- Results:
  - 99.6% of code unmodified or automatically converted
  - Port required less than 1 week developer time
  - Supports all CAFFE features (multi-gpu, P2P, FFT filters)
  - HIPCAFFE is the fastest CAFFE on AMD hardware – 1.8X faster than CAFFE/OpenCL
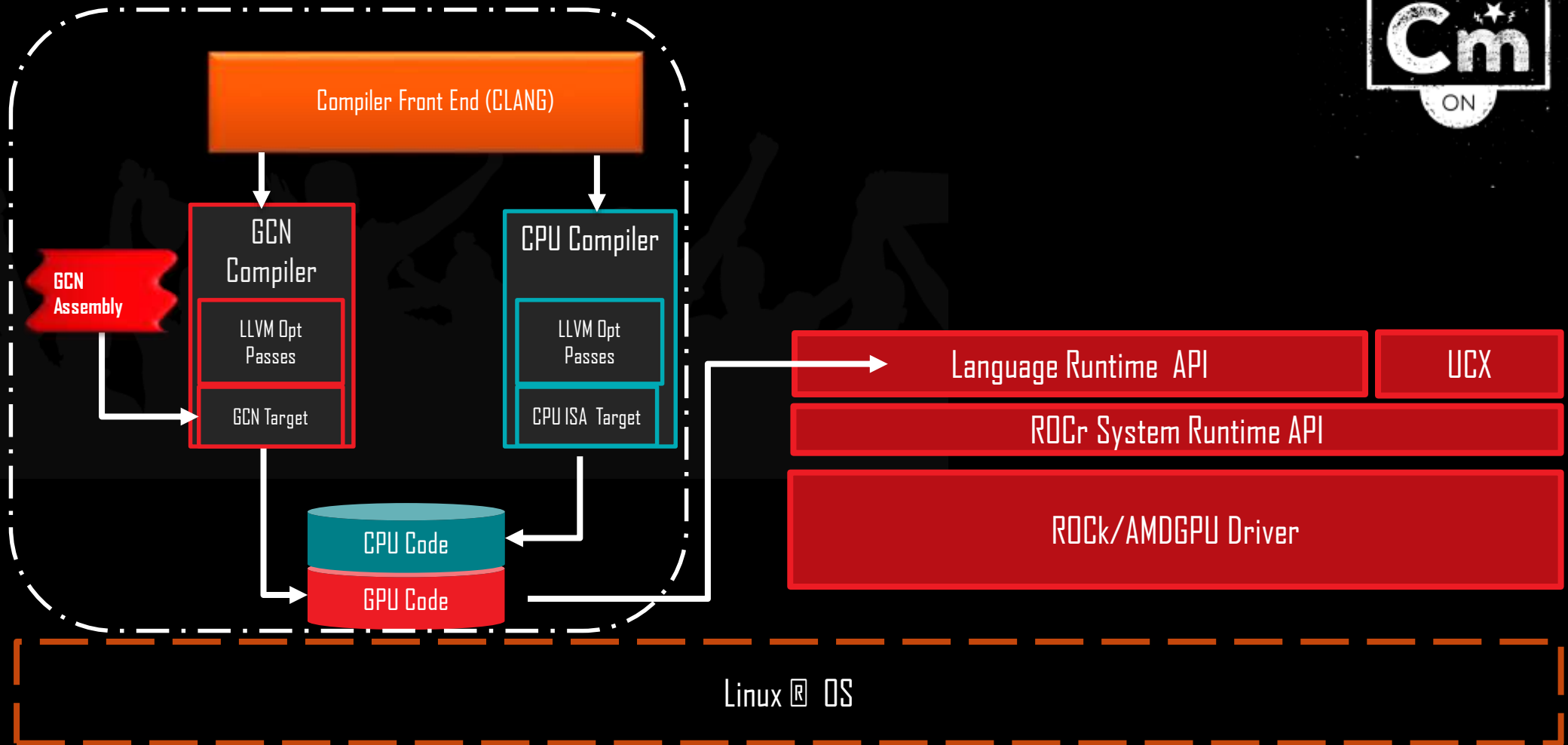
## Complexity of Application Porting: CAFFE

Lines of Code Changed

35000
30000  Manual, 32227
25000
20000
15000
10000
5000   Manual, 219
       Automatic, 688
0

OpenCL Port    HIP Port

AMD Internal Data

# ROCm SOFTWARE

**AMD**

| Machine Learning Applications |
|:---:|

| | | | |
|:---:|:---:|:---:|:---:|
| Caffe | TensorFlow | Torch 7 | MxNet |
| CNTK | Chainer | Theano | |

**Frameworks**

| MIOpen | BLAS,FFT,RNG | NCCL | C++ STL |
|:---:|:---:|:---:|:---:|

**Middleware & Libraries**

| HCC | HIP | OpenCL | Python |
|:---:|:---:|:---:|:---:|

| ROCm Platform |
|:---:|

**ROCm**

ROCm ON

**RADEON**

# DELIVERING AN OPEN PLATFORM FOR GPU COMPUTING

Language neutral solution to match developer needs as heterogeneous programing models evolve

**GCN Compiler**
- Direct-to-ISA
- GCN Docs
- CLANG/LLVM
- GCN Assembler
- *Open-source*

Compiler Front End (CLANG)

GCN Assembly

**GCN Compiler**
- LLVM Opt Passes
- GCN Target

**CPU Compiler**
- LLVM Opt Passes
- CPU ISA Target

CPU Code

GPU Code

Language Runtime API

UCX

ROCr System Runtime API

ROCk/AMDGPU Driver

Linux ® OS

# EXTENDING SUPPORT TO A BROADER HARDWARE ECOSYSTEM

ROCm "Open Source" foundation brings a rich foundation to these new ecosystems

AMD64 Support
AMD ZEN

Intel Xeon E5 v3 v4

ARM AArch64 Support
Cavium Thunder X

CAVIUM
THUNDERX

IBM OpenPower Support
– IBM Power 8

POWER8

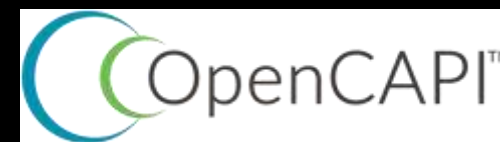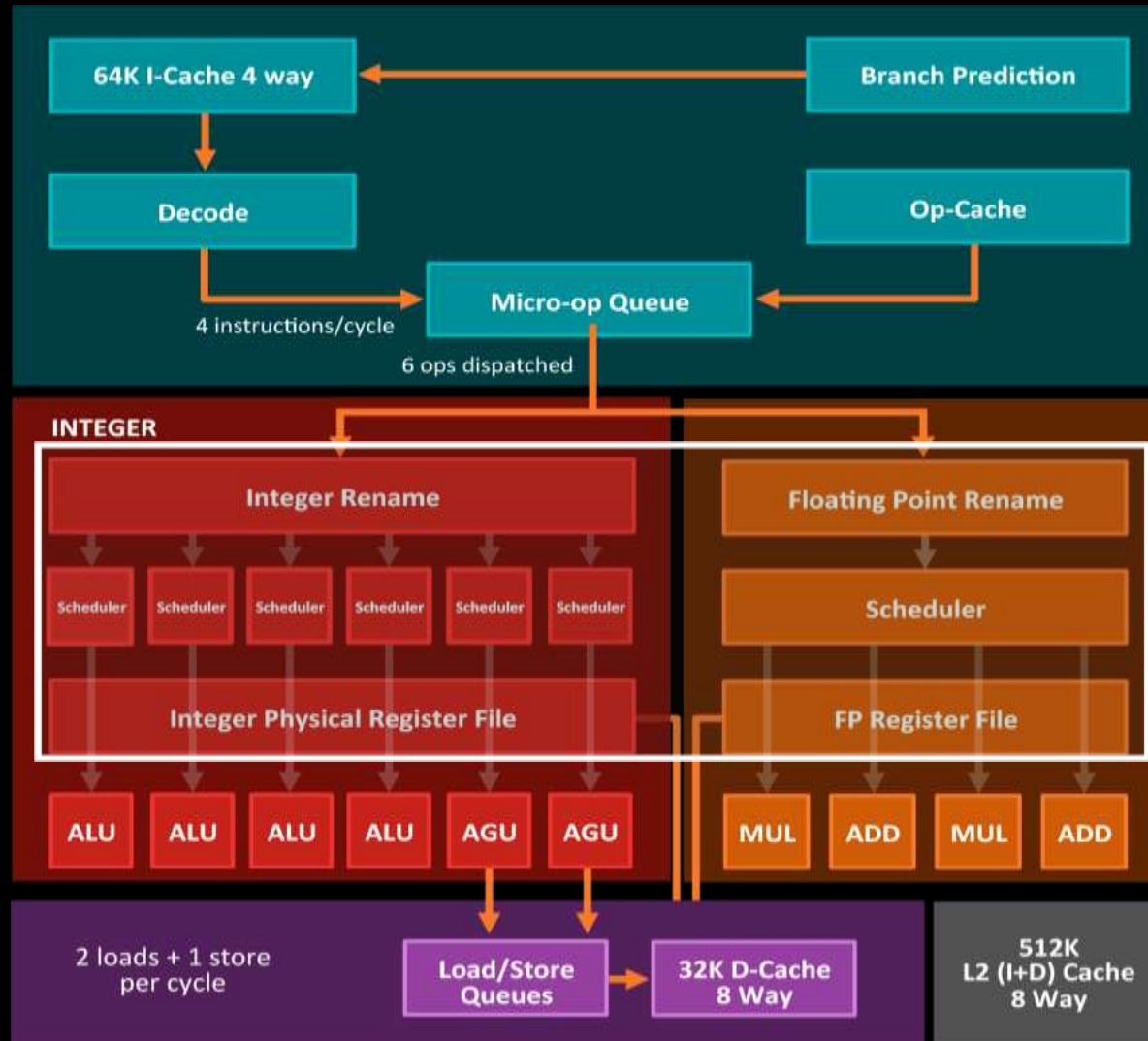ROCm is being built to support next generation I/O Interfaces

**GenZ Founding Member**

GENZ

CCIX Founding Member

CCIX Cache Coherent Interconnect
for Accelerators ▶ ▶ ▶

OpenCAPI Founding Member

OpenCAPI™

# ZEN CPU CORE: PERFORMANCE AND THROUGHPUT

**AMD**



QUANTUM LEAP IN CORE EXECUTION CAPABILITY

- ❑ Enhanced branch prediction to select the right instructions
- ❑ Micro-op cache for efficient ops issue
- ❑ 1.75X instruction scheduler window*
- ❑ 1.5X issue width and execution resources*

Result: instruction level parallelism designed for dramatic gains in single-threaded performance

*Compared to predecessor

RADEON
TECHNOLOGIES GROUP

**AMD**

## "RYZEN" aka "SUMMIT RIDGE"

◢ 8 CORES, 16 THREADS

◢ AM4 Platform

- DDR4
- PCI EXPRESS® GEN 3
- NEXT-GEN I/O

▸ https://www.amd.com/en/ryzen?&gclid=CL7W9ZyX-tICFUOXfgodGt8BPg

**RADEON** TECHNOLOGIES GROUP

**32-Core, 64-Thread**

**1st Public Demo, August 2016, San Francisco**

**Enabling Industry Software Support**

2 SOCKET

64 CORE

128 THREAD
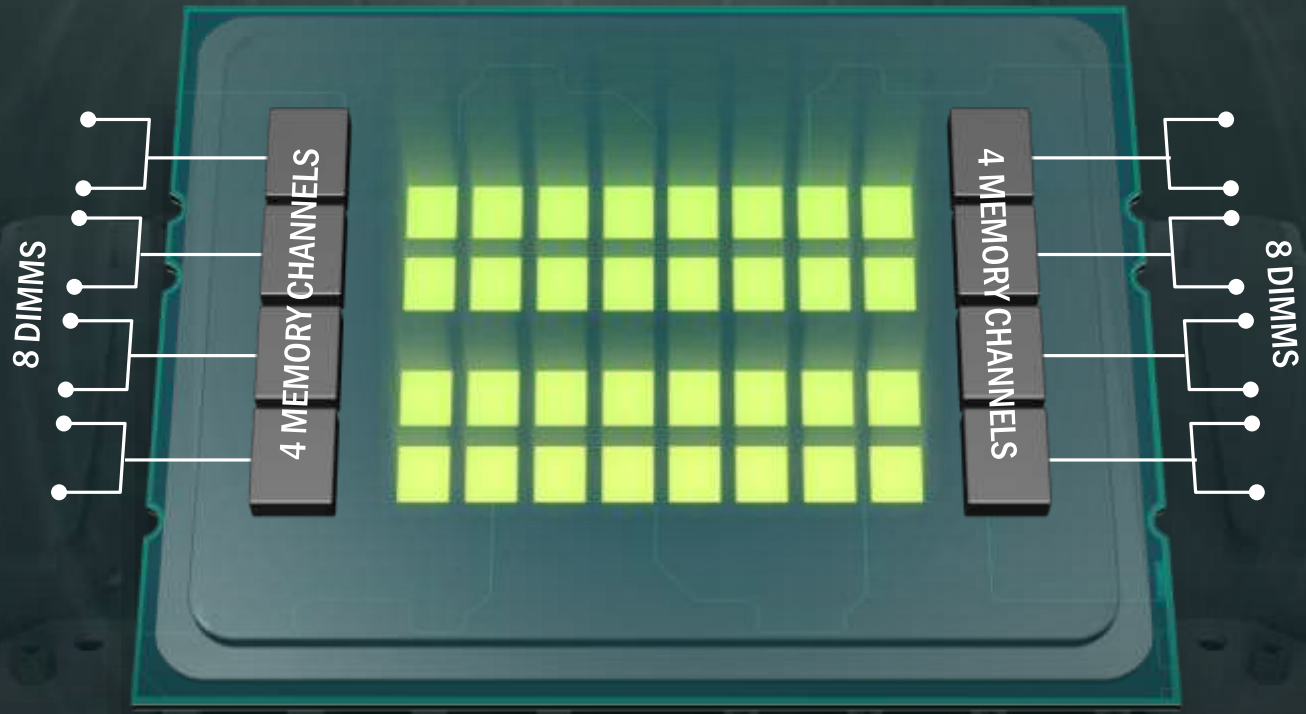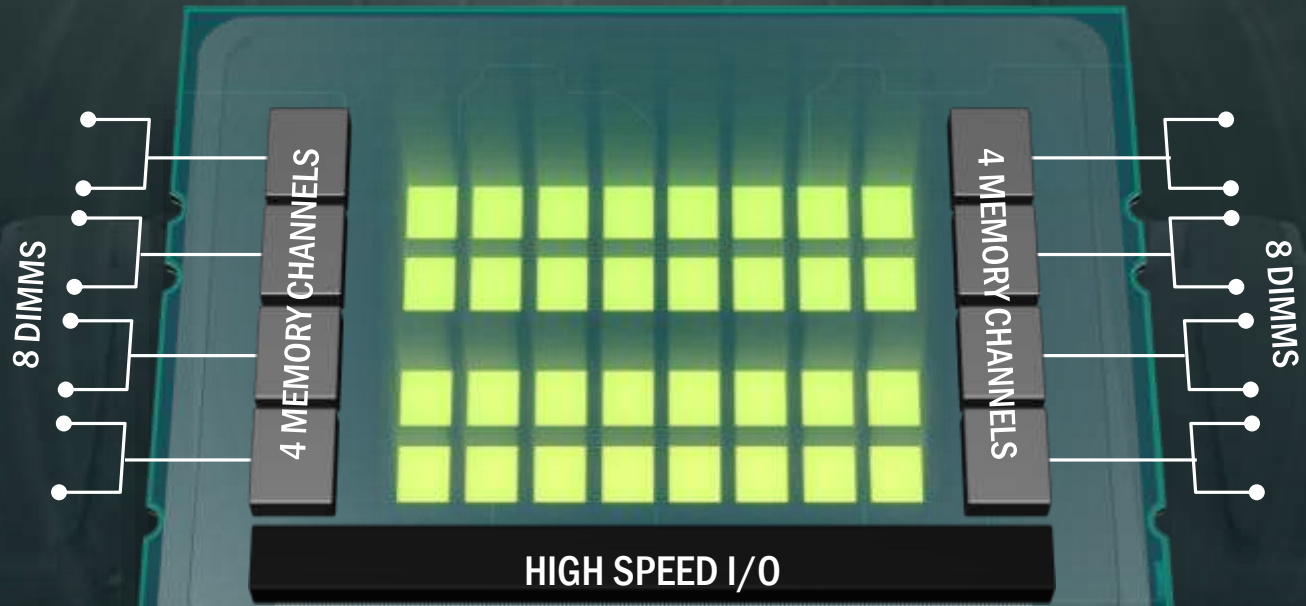
8 DIMMS
4 MEMORY CHANNELS
4 MEMORY CHANNELS
8 DIMMS
HIGH SPEED I/O
64 LANES

FABRIC

8 DIMMS
4 MEMORY CHANNELS
4 MEMORY CHANNELS
8 DIMMS
HIGH SPEED I/O
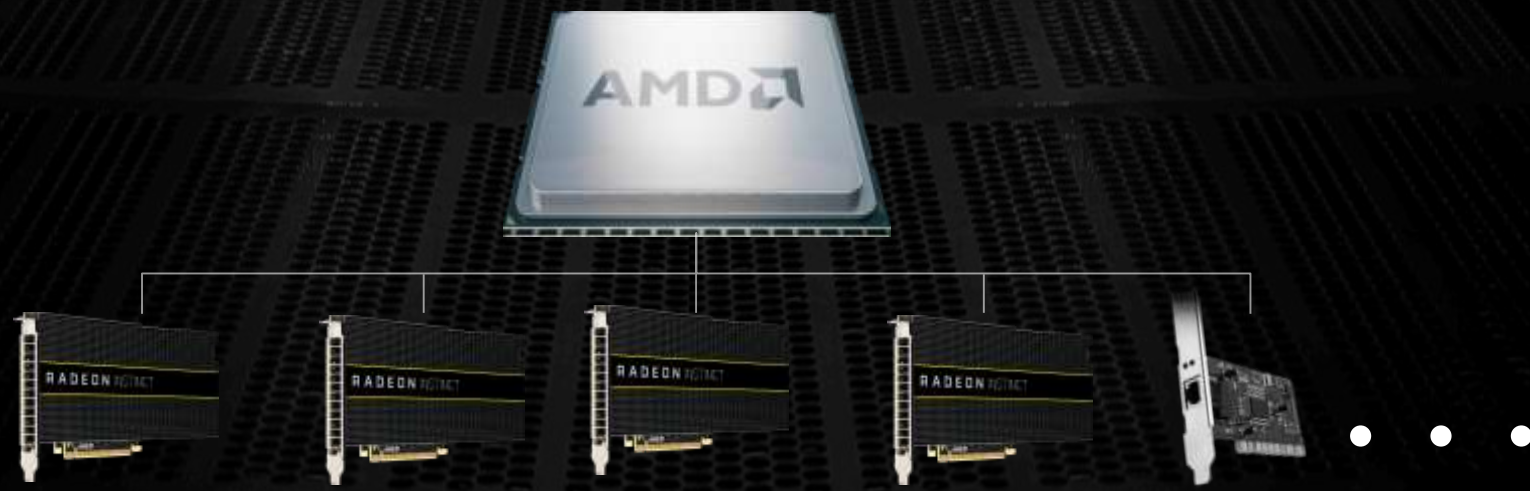64 LANES

# DEMO SETUP: EPYC VS. FASTEST INTEL 2-SOCKET SERVER

Both systems AMD and INTEL have the following features:

| Component | AMD | INTEL |
|---|---|---|
| CPU model | "EPYC" | E5-2699A V4 |
| Total CPUS | 2 | 2 |
| Total cores (SMT/HT on) | 128 | 88 |
| Total memory channels | 16 | 8 |
| Total memory capacity (16 GB DIMMS) | 512 | 384 |
| Memory frequency | 2400 | 1866 |
| Total PCIE gen3 lanes to CPUs | 8x16=128 | 2x40=80 |

o   Intel server is a standard, commercially available server from a major OEM

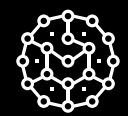# Radeon Instinct with Zen "EPYC" Platform

High-speed Network Fabric
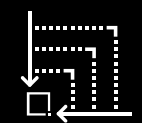
## Optimized for GPU and Accelerator Throughput computing

**Lower System Cost**

**Lower Latency Architecture**

**Peer to Peer Communication**

**High Density Footprint**

AMD | RADEON

# DISCLAIMER & ATTRIBUTION

**AMD**

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions and typographical errors.

The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION.

AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY DIRECT, INDIRECT, SPECIAL OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

**RADEON** TECHNOLOGIES GROUP

# Backup slides

# THE HETEROGENEOUS SYSTEM ARCHITECTURE (HSA)

◢ HSA is a platform architecture and software environment for simplified efficient parallel programming of heterogeneous systems, targeting:
– Single-source language support:
– Mainstream languages: C, C++, Fortran, Python, OpenMP
– Task-based, domain-specific, and PGAS languages
– Extensibility to a variety of accelerators
– GPUs, DSPs, FPGAs,, etc.

◢ The HSA Foundation promotes HSA via:
– Open, royalty-free, multi-vendor specifications
– Open-source software stack and tools
– Runtime stack
– Compilers, debuggers, and profilers

◢ See http://www.hsafoundation.com  and
http://github.com/hsafoundation

| AMD STRATEGY IN EXASCALE  SUPERCOMPUTING AND MACHINE INTELLIGENCE  |  SEPTEMBER 20, 2017    NANO AND GIGA CHALLENGES , TOMSK, RUSSIAN FEDERATION