# Innovation to meet AI and HPC computing growth

**Mark Papermaster**
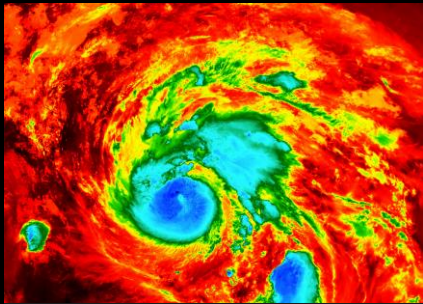**CTO and EVP**

AMD
together we advance_

# Fundamental inflection of AI powered computing
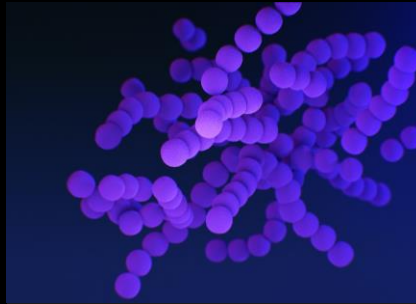
## Benefits are key to solving the world's most pressing problems



**Education and Knowledge base**



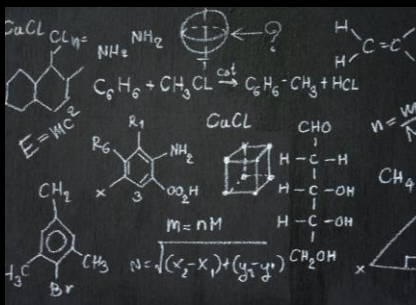**Climate Change**



**Chemical Sciences**
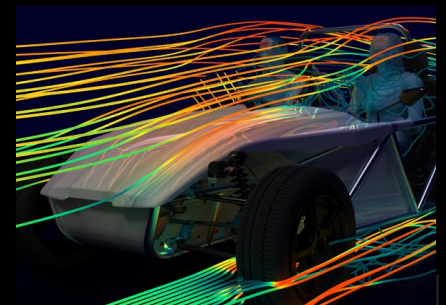


**Energy Solutions**



**Productivity**



**Healthcare**



**Content Creation**



**Research**



**Gaming and Entertainment**



**Real-Time Simulation**

AMD
together we advance_

# Insatiable demand for more compute



AI

Traditional CPU and
GPU Compute

Accelerated
Computing

Compute Demand (log scale)

2005　　2010　　2015　　2020　　2025

AMD
together we advance_

# Constant innovation required in software and hardware



**Open Ecosystem**



**Flexible Chiplet Design**

**AMD** together we advance_
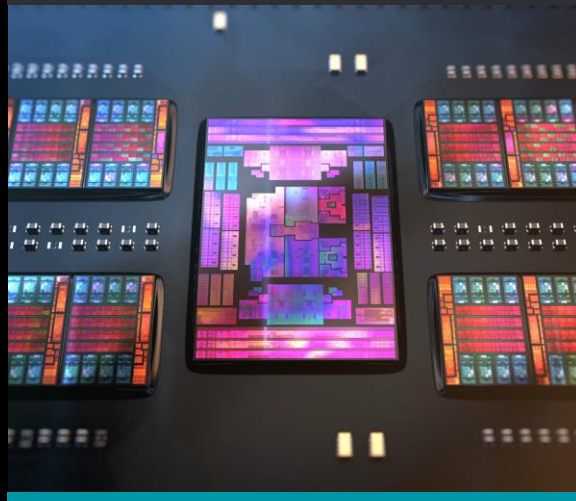
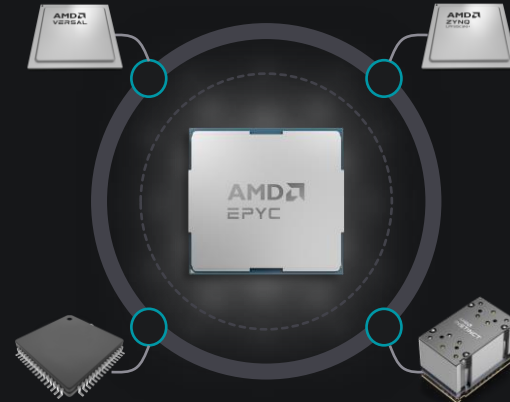# **Modularity** enables efficient and tailored solutions



**IP**

**Chiplets and 3D stacking**

**Component and system**

**Software**

AMD
together we advance_

# Chiplet and packaging evolution



**AMD**

| 2.5D HBM | Multichip Module | Chiplets | AMD 3D V-Cache | AMD Instinct MI300 |
|----------|------------------|----------|----------------|--------------------|
| 2015 | 2017 | 2019 | 2021 | 2023 |

**3D/2.5D Hybrid Compute Future**

**XILINX**

| World's Highest Capacity FPGA in 20nm | World's Highest Capacity FPGA in 16nm | Versal Premium | Versal 2.5D |
|---------------------------------------|---------------------------------------|----------------|-------------|
| 2015 | 2019 | 2021 | 2022 |

# Enabling design innovations

| Multichiplet SoCs | Power delivery | Heterogeneous tailored compute | SW and HW system level integration |
|---|---|---|---|

AMD
together we advance_

# A **holistic** design approach is required

- System-level optimizations

- Domain-specific heterogeneous architecture

- Tight integration of processors, packaging and interconnect

- Leveraging AI holistically



Silicon and Die Stacking

Application

Packaging

Software Stack

**Holistic Design**

Compute

Interconnect

Memory

Accelerators

AMD
together we advance_

# MI300: Next Gen Chiplet Design

- Leadership **HPC** and **Generative AI accelerator** based on configuration

- 5nm process technology with 3D stacking

- Up to 5.2 TB/s memory bandwidth

- Up to 153 billion transistors

- Frameworks and open models fully supported
  PyTorch, TensorFlow, ONNX, Hugging Face, others

- Easy migration path from CUDA

**Chiplet Configurable**

AMD
together we advance_

# Many more challenges to solve

# Bandwidth demands accelerating power consumption

## Server memory interface power

Based on AMD internal data

AMD
together we advance_

# Even tighter integration of compute and memory



**DRAM layers**

**Compute**

**Silicon interposer**

**Memory layers**

**Compute**

## Integration enables higher bandwidth at lower power

| | DIMMS | 2.5D Micro-bumps (HBM) | 3D Hybrid Bond |
|---|---|---|---|
| **pj/bit** | **~12** | **~3.5** | **~0.2** |

Image source: https://commons.wikimedia.org/wiki/File:SDRAM-Modul.jpg, Creative Commons 4.0.

**AMD**
*together we advance_*

# Optical communication for energy efficient connectivity

Co-packaged optics provide compelling efficiency gains

Single mode, enabling 10m up to 2km reach

Energy efficient at <1pJ/bit receive energy

Tight integration of optical transceivers to compute die is the key to efficiency
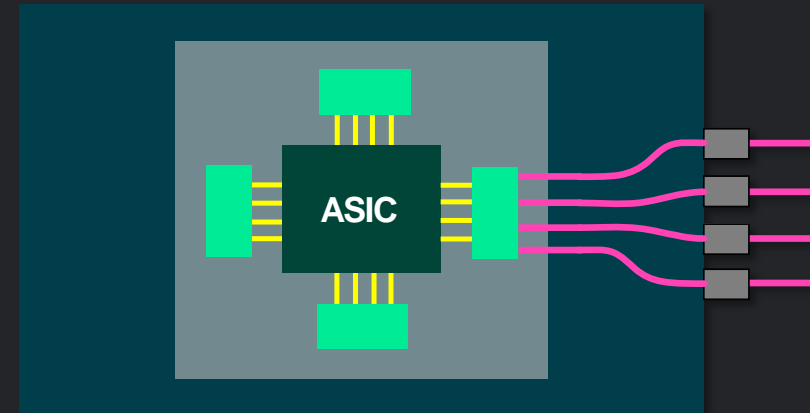
**Top view**

**Side view**

AMD
together we advance_

# Heat and power density

## Managing thermal density remains a priority

| **Thermal interface and lid materials** | **System cooling techniques** | **Power management design** |
| --- | --- | --- |
| Improving thermal conductivity from chip to lid | Improved ambient temperature control in system | More interactive thermal management |



ERI 2.0 Summit | August 23, 2023

AMD
together we advance_

# Improved power efficiency and regulation

## Per core regulation
Up to 19% power reduction (64-cores)



## Heterogenous optimized
Optimized PDN serving multichip HPC SoC



## Package integrated VRs
On-package voltage regulation

Sources: See End Notes.

AMD
together we advance_

# AMD

# Leadership AI and HPC products

| Energy efficient high-performance | Holistic design | Open solutions | Public and private sector partnerships |

AMD
together we advance_

# ENDNOTES

Slide 15.

The source for per core regulation image is [Friedrich14]. J. Friedrich, et. al., "POWER8: A 12-Core Server-Class Processor in 22nm SOI with 7.6T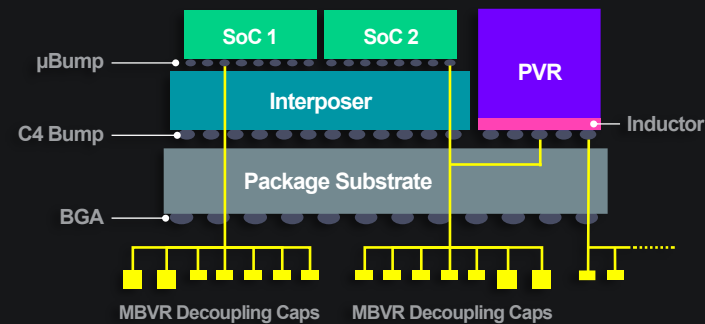b/s Off-Chip Bandwidth," ISSCC, Feb. 2014. The source for package integrated VRs image is [Kurd14] N. Kurd, et. al., "Haswell: A Family of IA 22nm Processors," ISSCC, Feb. 2014.

**AMD**
together we advance_

# Disclaimer

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.

## ATTRIBUTION
© 2023 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, EPYC, Instinct, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions.

AMD
together we advance_