**SE Seminar #9 Report**

**Latest AMD AI Technology**

**01286391 Seminar in Software Engineering**

**Software Engineering Program**

**Faculty of Engineering, KMITL**

By

65011277 Chanasorn Howattanakulphong

**Introduction**

In this seminar, Mr. Ponchai Prueksarattanon and Mr.Somyod Raksasat, a special speaker, discusses the topic "Latest AMD AI Technology"

**The AMD Instinct MI300X Accelerator**

- The AMD Instinct MI300X Accelerator is a specialized computing accelerator designed by AMD for high-performance computing tasks. The "AX IOD 256" likely refers to the I/O Die (IOD) associated with the AMD CDNA architecture, with "256" potentially indicating a specific aspect of its design.

- The accelerator features AMD Infinity Cache Technology, an on-die cache architecture aimed at improving bandwidth and reducing latency by providing high-speed cache memory closer to the compute units. It also incorporates Next Gen I/O, signifying the utilization of next-generation input/output technology for high-speed data transfer capabilities.

- With a 304 AMD CDNA 3 Compute Unit configuration, the accelerator utilizes AMD's Compute-Dense Next-Gen Architecture optimized for compute-intensive workloads. The "8XCD" designation suggests a compute density that is eight times higher compared to previous generations or similar devices.

- The presence of "8X HBM3" indicates that the accelerator incorporates eight stacks of High-Bandwidth Memory (HBM) version 3, a type of memory known for its high bandwidth and commonly used in high-performance computing applications.

- Regarding memory specifications, the accelerator boasts a capacity of 192GB and a peak memory bandwidth of 5.3 terabytes per second (TB/s), highlighting its high memory performance.

AMD, which stands for Advanced Micro Devices, is a prominent semiconductor company known for its development of CPUs (Central Processing Units), GPUs (Graphics Processing Units), and various other computing technologies. Here are key aspects of AMD:

**History:**

Established in 1969 by Jerry Sanders, AMD initially focused on manufacturing semiconductor memory chips. In the 1980s and 1990s, the company shifted its emphasis to microprocessors, emerging as a significant competitor to Intel in the CPU market.

**Product Lines:**

- CPUs: AMD offers a diverse range of CPUs for desktops, laptops, servers, and workstations. The lineup includes Ryzen for consumer desktops, Ryzen Thread ripper for high-end desktops and workstations, EPYC for servers, and Athlon for entry-level systems.

- GPUs: AMD produces GPUs under the Radeon brand, utilized in gaming PCs, workstations, and servers. Known for competitive performance and value, AMD's GPUs have made a notable impact, especially in the gaming market.
- Accelerated Processing Units (APUs): Combining CPU and GPU cores on a single chip, AMD's APUs are commonly used in laptops and low-power desktops.
- Semi-Custom Chips: AMD designs semi-custom chips for gaming consoles, including the PlayStation and Xbox.

**Technological Innovations:**

- Zen Architecture: Introduced in 2017, AMD's Zen microarchitecture brought a substantial improvement in CPU performance and efficiency. Subsequent iterations have continued to refine these aspects.
- RDNA Architecture: RDNA serves as AMD's graphics architecture designed for gaming and high-performance computing, offering enhanced performance-per-watt compared to previous generations.
- Infinity Fabric: A high-speed interconnect technology connecting CPU cores, GPU cores, and other components within AMD processors, Infinity Fabric enables efficient communication between different parts of the chip.

**Competition:**

AMD is a major competitor to Intel and NVIDIA in the CPU and GPU markets. Notably, AMD has gained significant market share in recent years, particularly with its Ryzen processors in the CPU segment.

**Industry Impact:**

AMD's competitive products have spurred innovation and pricing competition in the CPU and GPU markets. This has resulted in a broader range of choices for consumers and improved performance across various price points.

**NVIDIA and AMD**

Determining the superior performance between AMD and NVIDIA depends on various factors, including specific product lines, use cases, and individual preferences. Both companies manufacture high-quality GPUs that excel in different scenarios.

- Gaming Performance: Traditionally, NVIDIA has dominated the gaming market with its GeForce GPUs. Particularly in high-end graphics cards, such as the GeForce RTX series, NVIDIA often leads in raw gaming performance, offering features like ray tracing capabilities and optimized drivers for popular gaming titles.

- Price-to-Performance Ratio: AMD GPUs, especially those in the Radeon RX series, typically present competitive price-to-performance ratios. In the mid-range and budget segments, AMD GPUs frequently deliver commendable performance at more affordable price points compared to their NVIDIA counterparts. This makes AMD GPUs an appealing choice for gamers seeking solid performance without a significant financial investment.

- Compute Performance: In specific compute-intensive workloads like content creation, deep learning, and scientific simulations, both AMD and NVIDIA GPUs demonstrate robust performance. AMD's Radeon Pro and Instinct series, along with NVIDIA's Quadro and Tesla series, target professional and data center markets, offering specialized features and optimizations.

- Software and Ecosystem: Historically, NVIDIA has maintained a more mature ecosystem with features like CUDA for general-purpose GPU computing and strong developer support. Nevertheless, AMD has been actively enhancing its software stack, with initiatives like ROCm (Radeon Open Compute) and endeavors to optimize for popular software frameworks such as TensorFlow and PyTorch. This reflects AMD's commitment to improving its software and ecosystem offerings.

**What I Have Learned:**

The seminar presentation on "Recent Advances in AMD AI Technology" explores the specifications of the AMD Instinct MI300X Accelerator, a dedicated computing accelerator tailored for high-performance computing tasks. Noteworthy features of the accelerator include AMD Infinity Cache Technology, Next Gen I/O capabilities, enhanced compute density (8XCD), a significant number of compute units (304 AMD CDNA 3 Compute Unit), and the incorporation of High-Bandwidth Memory (HBM3) to achieve an impressive peak memory bandwidth of 5.3 terabytes per second (TB/s). Beyond the detailed examination of the MI300X, the report offers an overview of AMD as a semiconductor company, elucidating its history, diverse product lines (encompassing CPUs, GPUs, and APUs), technological innovations (such as the Zen and RDNA architectures), its competitive stance against Intel and NVIDIA, and the broader industry impact stemming from its competitive products. In essence, the seminar provides valuable insights into AMD's recent strides in AI technology and its broader contributions to the computing landscape.