



NVIDIA ADA GPU ARCHITECTURE

Designed to deliver outstanding gaming and creating, professional graphics, AI, and compute performance.

*Updated to include information on NVIDIA L40 and L4 Data Center GPUs

Table of Contents

Introduction	4
Ada GPU Architecture In-Depth	6
Ada AD102 GPU	7
Memory Subsystem	12
4N Manufacturing Process and Power Efficiency	12
Ray Tracing	15
2x Faster Ray-Triangle Intersection Testing	16
2x Faster Alpha Traversal Performance with Opacity Micromap Engine	16
10x Faster BVH Build in 20X Less BVH Space with Ada’s Displaced Micro-Mesh Engine	18
Shader Execution Reordering (SER)	21
DLSS 3 and Optical Flow Acceleration	23
Ada Fourth-Generation Tensor Core	24
NVIDIA Broadcast/Video	24
Conclusion	26
Appendix A - GeForce RTX 4090 GPU Full Specifications	29
Appendix B - Ada AD103 GPU Full Specifications	32
Appendix C - NVIDIA L40 GPU Full Specifications	36
Appendix D - NVIDIA L4 GPU Full Specifications	39

List of Figures

- Figure 1. AD102 GPU Full-Chip Block Diagram..... 7
- Figure 2. Ada GPC Block Diagram 8
- Figure 3. RT Core Second Generation Block Diagram (Ampere architecture)..... 9
- Figure 4. RT Core Third Generation Block Diagram (Ada architecture)..... 9
- Figure 5. Ada Streaming Multiprocessor (SM) 11
- Figure 6. Complex Shapes Use Texel's Alpha Channel..... 16
- Figure 7. Opacity Mask Applied to Leaf 17
- Figure 8. Ada Opacity Micromap Engine Compared to Ampere Arch..... 18
- Figure 9. Displacement Micro-Mesh - Base Mesh and Micro-Meshes 19
- Figure 10. DMM Simplified BVH, Base Triangle, and Displacement Map.....20
- Figure 11. DMMs Reduce BVH Build Time and Storage Requirements.....21
- Figure 12. Shader Execution Reordering Pipeline22
- Figure 13. DLSS 3 Motion Vectors + Optical Flow = Accurate Motion Estimation23

List of Tables

- Table 1. GeForce RTX 4090 vs GeForce RTX 3090 Ti / 2080 Ti Specifications..... 13
- Table 2. GeForce RTX 4090 vs RTX 3090 Ti vs 2080 Ti29
- Table 3. GeForce RTX 4080 16 GB vs 3080 Ti vs 2080 Super.....32
- Table 4. NVIDIA L40 vs NVIDIA A40.....36
- Table 5. NVIDIA L4 GPU vs NVIDIA T439

Introduction

Launched in 2018, NVIDIA's® Turing™ GPU Architecture ushered in the future of 3D graphics and GPU-accelerated computing. Turing provided major advances in efficiency and performance for PC gaming, professional graphics applications, and deep learning inferencing. Using new hardware-based accelerators, Turing fused rasterization, real-time ray tracing, AI, and simulation to enable incredible realism in PC games, and cinematic-quality interactive experiences. Two years later in 2020, the NVIDIA® Ampere architecture incorporated more powerful RT Cores and Tensor Cores, along with a novel SM structure that offered 2x FP32 performance, clock-for-clock, compared to Turing GPUs. These innovations allowed the Ampere architecture to run up to 1.7x faster than Turing in traditional raster graphics, and up to 2x faster in ray tracing.

The new NVIDIA Ada Lovelace GPU architecture, named after mathematician Ada Lovelace, who is often regarded as the world's first computer programmer¹, raises the bar far above Turing and Ampere GPUs. While improvements in the silicon manufacturing process have slowed, modern computer graphics have seen an exponential rise in complexity. Increases in geometric complexity and innovations in lighting have resulted in graphics that look more lifelike than ever before. *Battlefield V* was the first title to take advantage of NVIDIA's hybrid rendering ray tracing approach, requiring 39 ray tracing operations per pixel to calculate the lighting effects in a typical scene. Four years later, *Cyberpunk 2077* running with its new RT: Overdrive Mode pushes over 600 ray tracing calculations per pixel. To generate environments with this level of complexity at high frame rates, Ada is up to 2x faster in rasterized games, and up to 4x faster in ray-traced games than the prior NVIDIA Ampere GPU architecture.

Ada provides the largest generational performance upgrade in the history of NVIDIA. This is made possible by three key innovations:

- **Revolutionary New Architecture:** NVIDIA Ada architecture GPUs deliver outstanding performance for graphics, AI, and compute workloads with exceptional architectural and power efficiency. After the baseline design for the Ada SM was established, the chip was scaled up to shatter records. Manufacturing innovations and materials research enabled NVIDIA engineers to craft a GPU with 76.3 billion transistors and 18,432 CUDA Cores capable of running at clocks over 2.5 GHz, while maintaining the same 450W TGP as the prior generation flagship GeForce® RTX™ 3090 Ti GPU. The result is the world's fastest GPU with the power, acoustics, and temperature characteristics expected of a high-end graphics card.
- **New Ada RT Core for Faster Ray Tracing:** For decades, rendering ray-traced scenes with physically correct lighting in real time has been considered the holy grail of graphics. At the same time, geometric complexity of environments and objects continues to increase as 3D games and graphics continually strive to provide the most accurate representations of the real world. The Ada RT Core has been enhanced to deliver 2x faster ray-triangle intersection testing and includes two important new hardware units. An Opacity Micromap Engine speeds up ray tracing of alpha-tested geometry by a factor of 2x, and a Displaced Micro-Mesh Engine generates Displaced Micro-Triangles on-the-fly to create additional geometry. The Micro-Mesh Engine provides the benefit of increased geometric complexity without the traditional performance and storage costs of complex geometries.
- **Shader Execution Reordering:** NVIDIA Ada GPUs support Shader Execution Reordering which dynamically organizes and reorders shading workloads to improve RT shading

efficiency. This improves performance by up to 44% in *Cyberpunk 2077 with Ray Tracing: Overdrive Mode*.

- **NVIDIA DLSS 3:** The Ada architecture features an all-new Optical Flow Accelerator and AI frame generation that boosts DLSS 3's frame rates up to 2x over the previous DLSS 2.0 while maintaining or exceeding native image quality. Compared to traditional brute-force graphics rendering, DLSS 3 is ultimately up to 4x faster while providing low system latency.

The GeForce RTX 4090 is the first GeForce graphics card based on the new Ada architecture. At the heart of the GeForce RTX 4090 is the AD102 GPU, which is the most powerful GPU based on the NVIDIA Ada architecture. AD102 has been designed to deliver revolutionary performance for gamers and creators, and enables the RTX 4090 to consistently deliver frame rates over 100 frames per second at 4K resolution in many games.

For the datacenter, the new NVIDIA L40 GPU based on the Ada architecture delivers unprecedented visual computing performance. Compared to the previous generation NVIDIA A40 GPU, NVIDIA L40 delivers 2X the raw FP32 compute performance, almost 3X the rendering performance, and up to 724 TFLOPs² of Tensor operation performance at the same 300W power envelope. NVIDIA L40 is the ideal GPU for servers running applications such as NVIDIA Omniverse, Generative AI, autonomous vehicle drive simulations, FP32 high performance computing (HPC), virtual workstations, cloud gaming, and single GPU AI training and inferencing. NVIDIA Certified Systems with L40 are optimized to deliver Omniverse at scale, such as the reference NVIDIA OVX system that include eight L40 GPUs that can be scaled up to deliver RTX-enabled digital twin renderings. The new NVIDIA L40 delivers this impressive performance speedup at the same 300W TDP and physical profile of the previous generation NVIDIA A40 GPU.

Finally, the Ada based NVIDIA L4 is designed to be the best low power universal GPU for AI, Graphics and Video workloads in the datacenter. Compact and versatile, the low-profile, single-slot, 72W L4 GPU fits in any server, making it ideal for global deployments all the way from regional datacenters out to the edge, including outdoor locations. NVIDIA L4 is the perfect choice for wide variety of applications such AI powered video services, Speech AI (ASR+NLP+TTS), small model Generative AI, search & recommenders, cloud gaming, and virtual Workstations, among many others.

Please refer to Appendices C and D for more details on NVIDIA L40 and L4, respectively.

1 – https://en.wikipedia.org/wiki/Ada_Lovelace

2 – FP8 performance with structured Sparsity enabled

Ada GPU Architecture In-Depth

NVIDIA engineers set clear design goals for every new GPU architecture.

With its groundbreaking RT and Tensor Cores, the Turing architecture laid the foundation for a new era in graphics, which includes ray tracing and AI-based neural graphics. Ampere's revamped SM, enhanced RT and Tensor Cores, and innovative GDDR6X memory subsystem established the bridge between traditional raster-based and ray traced graphics, accelerating both, and providing tremendous performance gains at the highest screen resolutions.

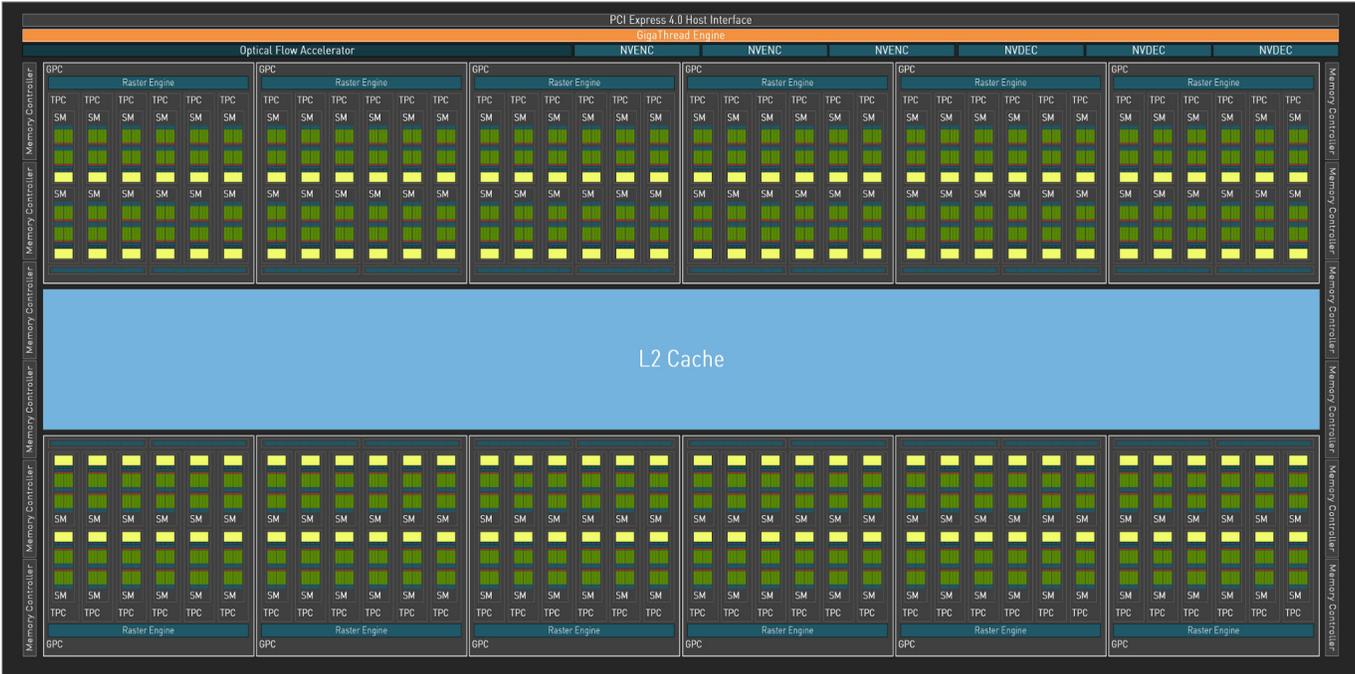
The launch of the new Ada GPU architecture is a breakthrough moment for 3D graphics: the Ada GPU has been designed to provide revolutionary performance for ray tracing and AI-based neural graphics. Performance improvements of 2-4x (up to 4x with the use of DLSS 3) over prior generation Ampere GPUs are possible. The Ada architecture provides a higher level of baseline GPU performance, and marks the tipping point where ray tracing and neural graphics become mainstream.

AD102 is the flagship of the Ada GPU lineup and launches first with the GeForce RTX 4090 graphics card. NVIDIA will also soon be providing follow-on Ada GPUs including AD103 and AD104, utilizing the same basic architecture as AD102.

This section will be focused on the AD102 GPU. For further information on AD103 and AD104, please consult *Appendix B, The Ada AD103 GPU* and *Appendix C, The Ada AD104 GPU*.

Ada AD102 GPU

The full AD102 GPU includes 12 Graphics Processing Clusters (GPCs), 72 Texture Processing Clusters (TPCs), 144 Streaming Multiprocessors (SMs), and a 384-bit memory interface with 12 32-bit memory controllers.

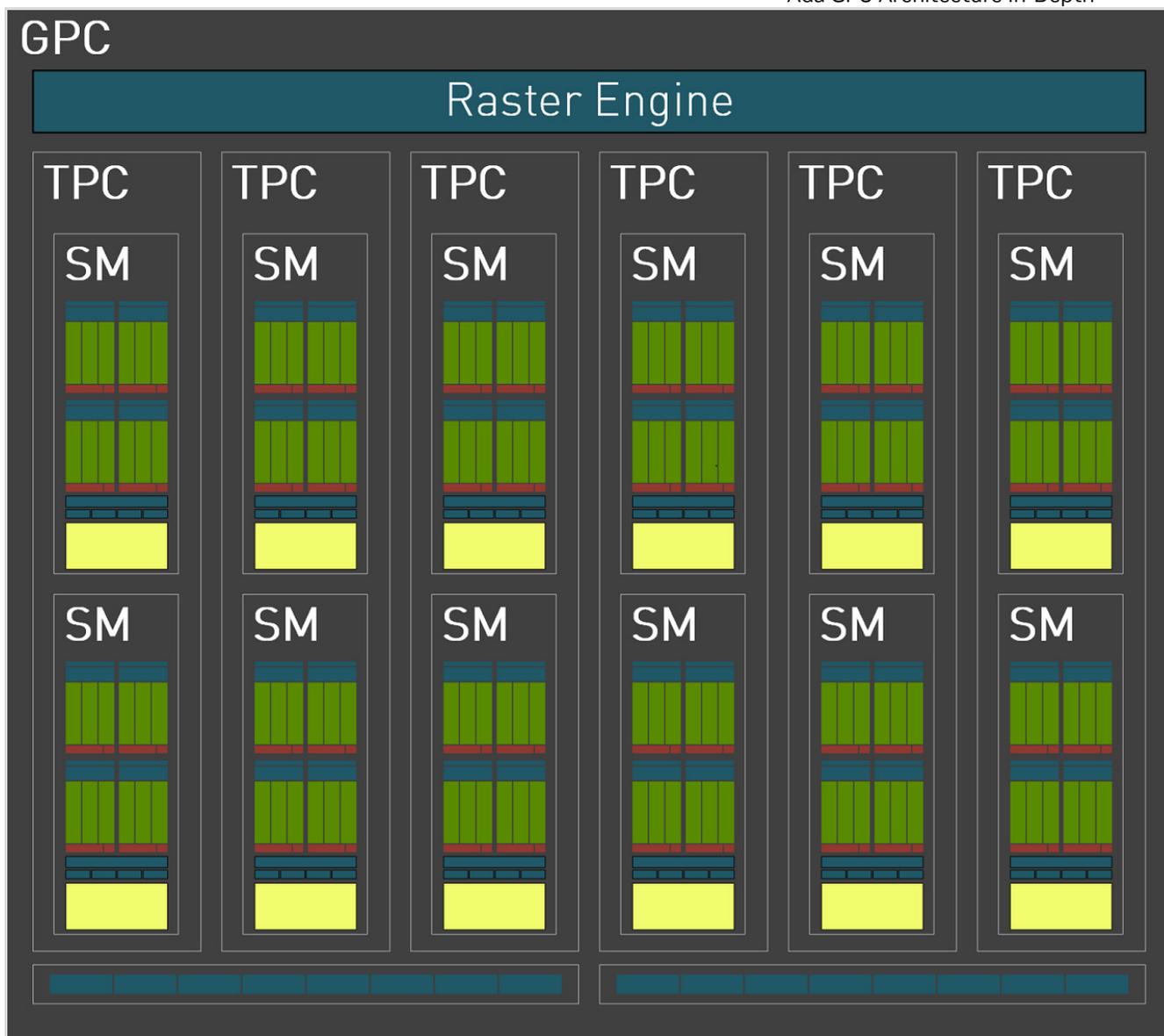


Note: The AD102 GPU also includes 288 FP64 Cores (2 per SM) which are not depicted in the above diagram. The FP64 TFLOP rate is 1/64th the TFLOP rate of FP32 operations. The small number of FP64 Cores are included to ensure any programs with FP64 code operate correctly, including FP64 Tensor Core code.

Figure 1. AD102 GPU Full-Chip Block Diagram.

The full AD102 GPU includes:

- 18432 CUDA Cores
- 144 RT Cores
- 576 Tensor Cores
- 576 Texture Units



Ada GPC with Raster Engine, 6 TPCs, 12 SMs, and 16 ROPs (8 per ROP partition).

Figure 2. Ada GPC Block Diagram

The GPC is the dominant high-level hardware block within all AD10x Ada family GPUs, with all of the key graphics processing units residing within a GPC. Each GPC includes a dedicated Raster Engine, two Raster Operations (ROPs) partitions, with each partition containing eight individual ROP units, and six TPCs. Each TPC includes one PolyMorph Engine and two SMs.

Each SM in AD10x GPUs contain 128 CUDA Cores, one Ada Third-Generation RT Core, four Ada Fourth-Generation Tensor Cores, four Texture Units, a 256 KB Register File, and 128 KB of L1/Shared Memory, which can be configured for different memory sizes depending on the needs of the graphics or compute workload.

The RT Core in Turing and Ampere GPUs includes dedicated hardware units for accelerating Bounding Volume Hierarchy (BVH) data structure traversal, and performing the ray-triangle and ray-bounding box intersection testing calculations that are critical for ray tracing. In the Ampere

RT Core diagram below, BVH traversal is accelerated by the **Box Intersection Engine** (represented by the group of bounding boxes on the left) and ray-triangle intersection testing is accelerated by the **Triangle Intersection Engine** (triangle on the right). By providing dedicated resources for these highly important ray tracing functions, work is offloaded from the SM, freeing it up to perform other pixel, vertex, and compute shading tasks. In testing with synthetic benchmarks and real-world games and applications, the RT Core found in Turing and Ampere GPUs has proven to be the highest performing engine for processing RT workloads to date.

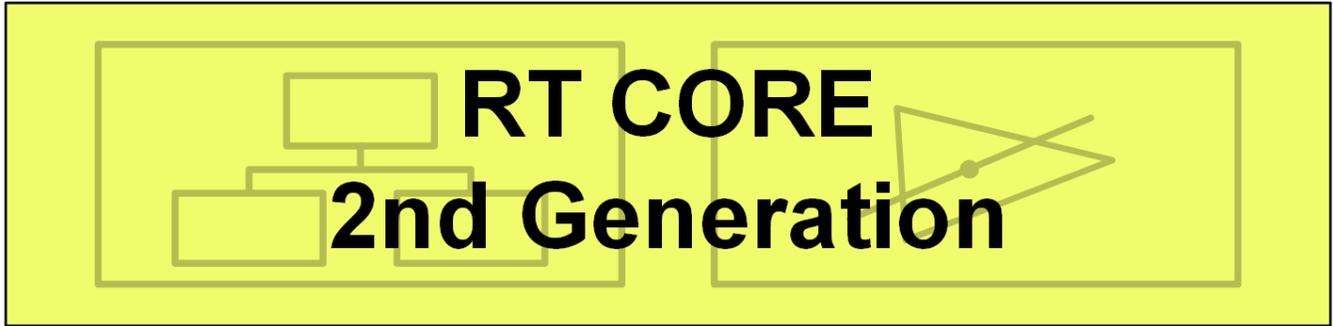


Figure 3. RT Core Second Generation Block Diagram (Ampere architecture)

In addition to these two functions, the Third-Generation RT Core found in Ada GPUs includes dedicated units known as the **Opacity Micromap Engine** and the Displaced Micro-Mesh Engine. The Opacity Micromap Engine evaluates Opacity Micromaps (represented by the triangle with foliage on the bottom left), which are used to accelerate alpha traversal. The **Displaced Micro-Mesh Engine** generates meshes of micro-triangles that are known as Displaced Micro-Meshes (represented by the triangle on the bottom right in the diagram below). Displaced Micro-Meshes allows the Ada RT Core to ray trace geometrically complex objects and environments with significantly less BVH build time and storage costs. Finally, ray-triangle intersection testing is 2x faster in Ada’s Third-Generation RT Core compared to the Ampere GPU generation.

Box Intersection Engine

Triangle Intersection Engine



NEW Opacity Micromap Engine

NEW Displaced Micro-Mesh Engine

Figure 4. RT Core Third Generation Block Diagram (Ada architecture)

Altogether these enhancements make the Ada Third-Generation RT Core the most powerful RT Core NVIDIA has ever built.

Like prior GPUs, the AD10x SM is divided into four processing blocks (or partitions), with each partition containing a 64 KB register file, an L0 instruction cache, one warp scheduler, one dispatch unit, 16 CUDA Cores that are dedicated for processing FP32 operations (up to 16 FP32 operations per clock), 16 CUDA Cores that can process FP32 or INT32 operations (16 FP32 operations per clock OR 16 INT32 operations per clock), one Ada Fourth-Generation Tensor Core, four Load/Store units, and a Special Function Unit (SFU) which executes transcendental and graphics interpolation instructions.

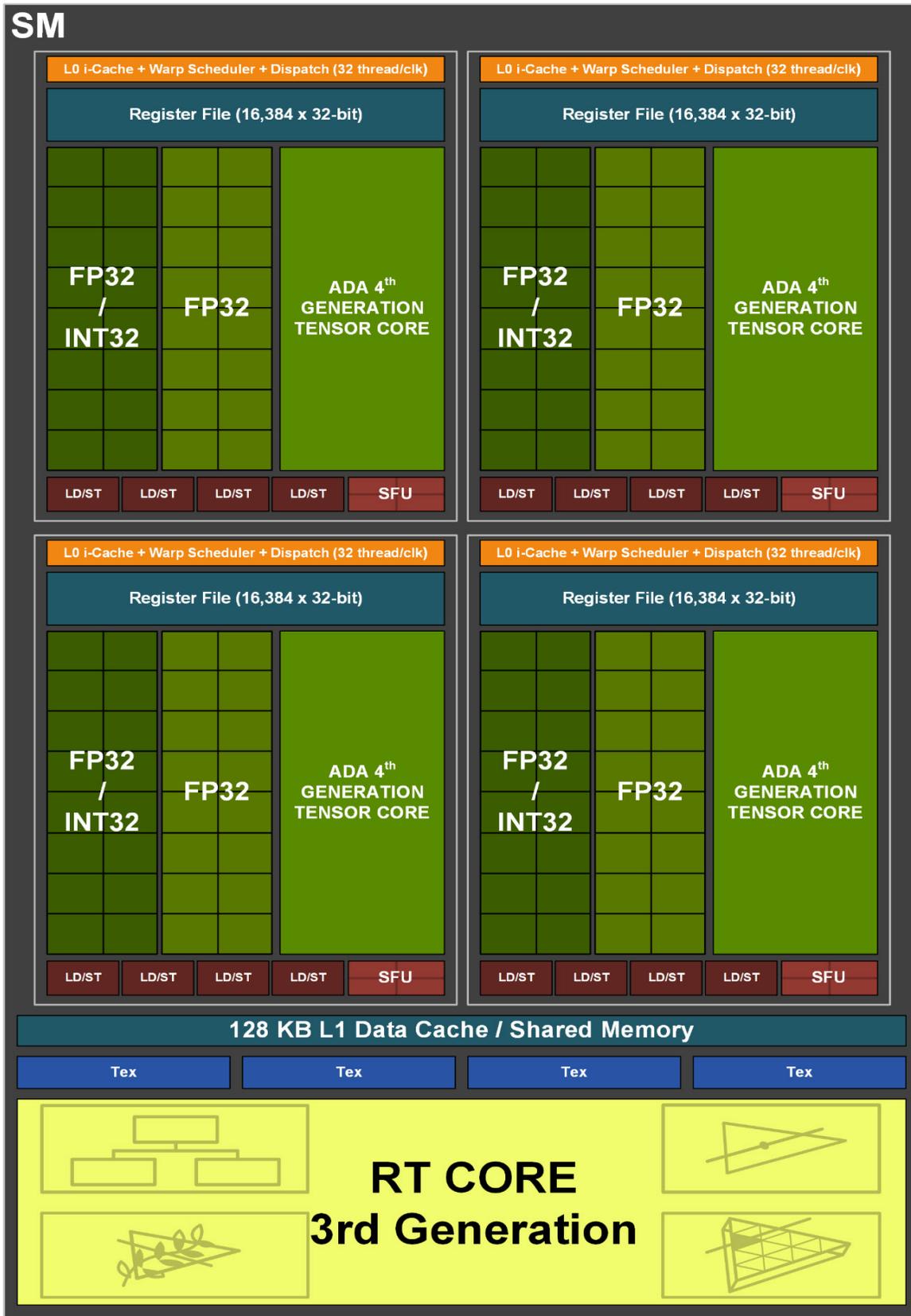


Figure 5. Ada Streaming Multiprocessor (SM)

Memory Subsystem

The Ada SM contains 128 KB of Level 1 cache. This cache features a unified architecture that can be configured to function as an L1 data cache or shared memory depending on the workload. The full AD102 GPU contains 18432 KB of L1 cache (compared to 10752 KB in GA102).

Compared to Ampere, Ada's Level 2 cache has been completely revamped. AD102 has been outfitted with 98304 KB of L2 cache, an improvement of 16x over the 6144 KB that shipped in GA102. All applications will benefit from having such a large pool of fast cache memory available, and complex operations such as ray tracing (particularly path tracing) will yield the greatest benefit.

NVIDIA works closely with the DRAM industry, collaborating on circuit design and signaling to enable the highest GPU memory speeds. With the launch of the GeForce RTX 3090 and 3080 Ampere GPUs, NVIDIA and Micron shipped the first GDDR6X devices, offering speeds up to 19.5 Gbps. Now, two years later, we've worked together to facilitate even higher memory speeds for Ada GPUs: the GeForce RTX 4080 ships with 22.4 Gbps GDDR6X memory; this is the highest speed of any GPU with GDDR-based memory, while the GeForce RTX 4090 offers 1 TB/sec of peak memory bandwidth.

4N Manufacturing Process and Power Efficiency

Ada GPUs are fabricated on TSMC's 4N manufacturing process. NVIDIA engineers worked closely with TSMC to optimize the process for GPU production. Using the 4N process enabled NVIDIA to integrate dramatically more cores: AD102 contains 70% more CUDA Cores than the prior generation GA102 GPU. In total, the AD102 GPU contains 76.3 billion transistors, making it one of the most complex chips ever made.

Ada also operates at high clock frequencies. NVIDIA optimized the design of the GPU, using high speed transistors in critical paths that could otherwise restrict the rest of the chip. Running at a GPU Boost clock of 2.52 GHz, the GeForce RTX 4090 ships with the highest clock frequency of any NVIDIA GPU.

At the same time however, RTX 4090's high clock speeds and core count deliver the highest performance per watt. When running at the same power as the RTX 3090 Ti, the RTX 4090 GPU delivers over 2x more performance.

Table 1. GeForce RTX 4090 vs GeForce RTX 3090 Ti / 2080 Ti Specifications

Graphics Card	GeForce RTX 2080 Ti	GeForce RTX 3090 Ti	GeForce RTX 4090
CUDA Cores	4352	10752	16384
GPCs	6	7	11
TPCs	34	42	64
SMs	68	84	128
GPU Boost Clock (MHz)	1635	1860	2520
FP32 TFLOPS	14.2	40	82.6
Tensor Cores	544 (2nd Gen)	336 (3rd Gen)	512 (4th Gen)
Tensor TFLOPS (FP8)	N/A	N/A	660.6/1321.2 ¹
RT Cores	68 (1st Gen)	84 (2nd Gen)	128 (3rd Gen)
RT TFLOPS	42.9	78.1	191
Texture Units	272	336	512
Texture Fill Rate	444.7	625	1290.2
ROPS	88	112	176
Pixel Fill Rate	143.9	208.3	443.5
Memory Size and Type	11 GB GDDR6	24 GB GDDR6X	24 GB GDDR6X
Memory Clock (Data Rate)	14 Gbps	21 Gbps	21 Gbps
Memory Bandwidth	616 GB/sec	1008 GB/sec	1008 GB/sec

L1 Cache/Shared Memory	6528 KB	10752 KB	16384 KB
L2 Cache	5632 KB	6144 KB	73728 KB
TGP	260 W	450 W	450 W
Transistor Count	18.6 Billion	28.3 Billion	76.3 Billion
Die Size	754 mm ²	628.4 mm ²	608.5 mm ²
Manufacturing Process	TSMC 12 nm FFN (FinFET NVIDIA)	Samsung 8 nm 8N NVIDIA Custom Process	TSMC 4N NVIDIA Custom Process

1- Using Sparsity feature

For the full list of GeForce RTX 4090 specifications, please see Appendix A at the back of this document.

Ray Tracing

The advent of real-time ray tracing has elevated the visual quality of games by delivering realistic lighting effects, physically accurate shadows, and better reflections, creating a final rendered image that approaches photorealism; all while running in real-time. Simultaneously, developers are rapidly increasing the geometric richness of scenes to keep pace with these leaps in shading quality.

In the pursuit of producing more realistic graphics, there is an exploding demand for deeply detailed environments. Vast libraries of objects are available as ingredients. Developers draw on the talents of their artists to craft compellingly intricate custom models.

Generally, developers mine geometric content from two major veins: scans of physical objects and artistically and/or algorithmically synthesized models. The former technique, photogrammetry, captures every minute geometric detail as well as material properties critical for accurate shading. Game studio artists also create fantastically detailed models. The result is objects often composed of millions of triangles, and environments composed of billions.

When working with environments like these, developers face two major challenges: storage and rendering performance. In a given frame, level of detail (LOD) techniques can mitigate some of the performance impact of scene complexity, but it's limited, since there is little control over where the camera/player may wander, and what scattering rays may hit (e.g., behind the camera).

NVIDIA engineers have developed three new features in the Ada RT Core to enable high-performance ray tracing of highly complex geometry:

- First, Ada's Third-Generation RT Core features **2x Faster Ray-Triangle Intersection Throughput** relative to Ampere; this enables developers to add more detail into their virtual worlds.
- Second, Ada's RT Core has **2x Faster Alpha Traversal**; the RT Core features a new **Opacity Micromap Engine** to directly alpha-test geometry and significantly reduce shader-based alpha computations. With this new functionality, developers can very compactly describe irregularly shaped or translucent objects, like ferns or fences, and directly and more efficiently ray trace them with the Ada RT Core.
- Third, the new Ada RT Core supports **10x Faster BVH Build in 20X Less BVH Space** when using its new **Displaced Micro-Mesh Engine** to generate micro-triangles from micro-meshes on-demand. The micro-mesh is a new primitive that represents a structured mesh of micro-triangles that the Ada RT Core processes natively, saving the storage and processing compared to what is normally required when describing complex geometries using only basic triangles.

Taken together these three advances incorporated into the Ada RT Core enable order-of-magnitude increases in richness without commensurate increases in processing time or memory consumption.

As we continue to approach photorealistic rendering with real-time ray tracing, increasing the accuracy with which we model the movement of light through extremely detailed, diverse environments means the raw processing workload becomes less and less coherent. Secondary rays used for reflections, indirect lighting, and translucency effects, for example, tend to shoot in

different directions and hit different materials, resulting in secondary hit shaders being less ordered and less efficient.

Left unaddressed, a loss in execution regularity can lead to inefficient use of the GPU's processing units, the SMs.

To address this issue, the Ada architecture introduces **Shader Execution Reordering**. This feature intelligently schedules shading work on-the-fly, so that complex materials like brushed metal can be processed more effectively.

For more information on the fundamentals of ray tracing and the GeForce RT Core, please refer to the [NVIDIA Turing GPU Architecture Whitepaper](#) and the [NVIDIA Ampere GA102 GPU Whitepaper](#).

2x Faster Ray-Triangle Intersection Testing

Ray-triangle intersection testing is a computationally expensive operation that is commonly performed when rendering a ray-traced scene. Recognizing the importance of this function, with each new RTX GPU NVIDIA engineers have strived to improve intersection testing performance and efficiency. The Third-Generation RT Core in the Ada architecture provides double the throughput for ray-triangle intersection testing over Ampere (and 4x faster than the first-generation RT Core used in Turing GPUs).

2x Faster Alpha Traversal Performance with Opacity Micromap Engine

Developers frequently use a texture's alpha channel to economically cut out complex shapes or more generally to represent translucency. A leaf might be described using a couple of triangles, employing a texture's alpha channel to economically capture the complex shape. A flame's complex shape and translucency can also be approximated by alpha.



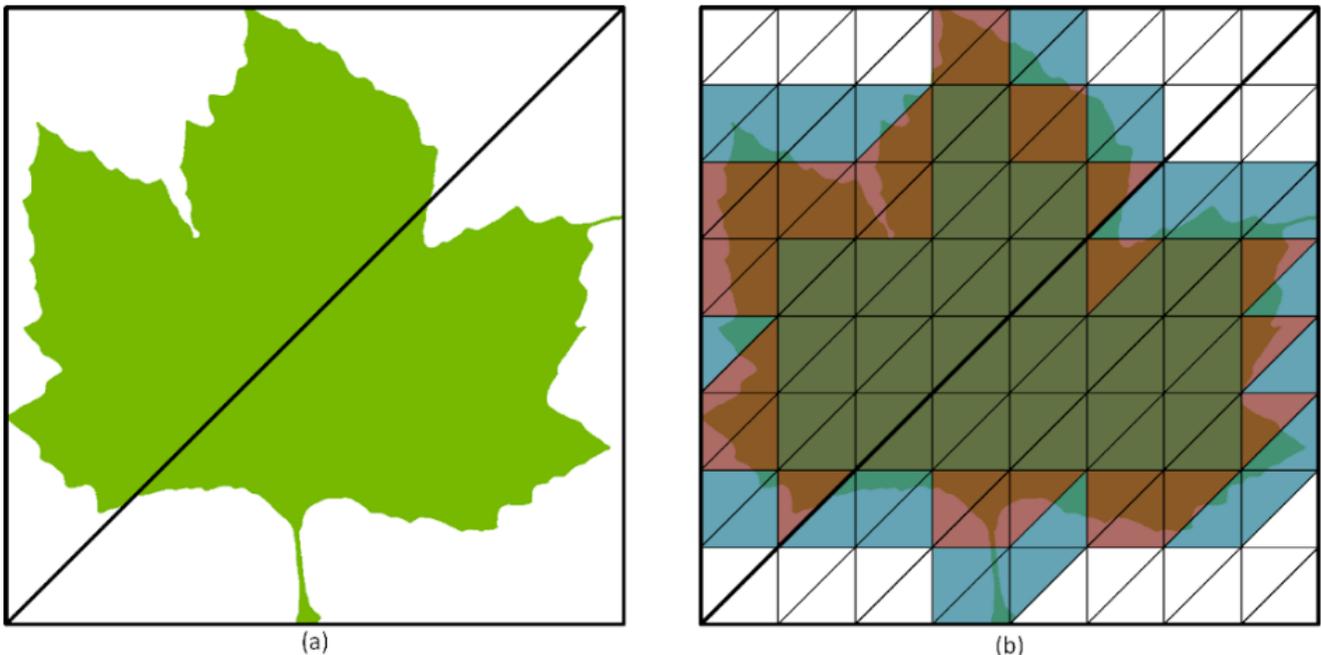
Complex shapes such as a leaf or a flame often use a texel's alpha channel to represent levels of transparency and opaqueness.

Figure 6. Complex Shapes Use Texel's Alpha Channel

Prior to Ada’s RT Core, a developer could incorporate these kinds of content into a ray traced scene by tagging them as not opaque. When a leaf is hit by a ray, a shader is invoked to determine how to treat the intersection, even if the ray is simply characterized as a hit or a miss. This incurs noticeable cost. Specifically, when a warp of rays is cast towards non-opaque objects, individual ray queries may require multiple shader invocations to resolve, while other rays terminate immediately. The result is lingering live threads and commensurate inefficiency.

To efficiently handle these kinds of content, NVIDIA engineers have added an Opacity Micromap Engine to Ada’s RT Core. An opacity micromap is a *virtual* mesh of micro-triangles, each with an opacity state that the RT Core uses to directly resolve ray intersections with non-opaque triangles. Specifically, the barycentric coordinates of an intersection are used to address the corresponding micro-triangle’s opacity state. The opacity state may be opaque, transparent, or unknown. If opaque, then a hit is recorded and returned. If transparent, the intersection is ignored and the search for an intersection continues. If unknown, then control is returned to the SM, invoking a shader (“anyhit”) to programmatically resolve the intersection.

The new Opacity Micromap Engine evaluates the opacity mask, which is a regular triangular mesh defined using the barycentric coordinate system used for reporting ray/triangle intersections. These meshes may be sized from one to *sixteen million* micro-triangles, with one or two bits associated with each micro-triangle. As a simple illustrative example, consider a detailed maple leaf described using two triangles and an alpha texture (see sub-Figure (a) in Figure 7).



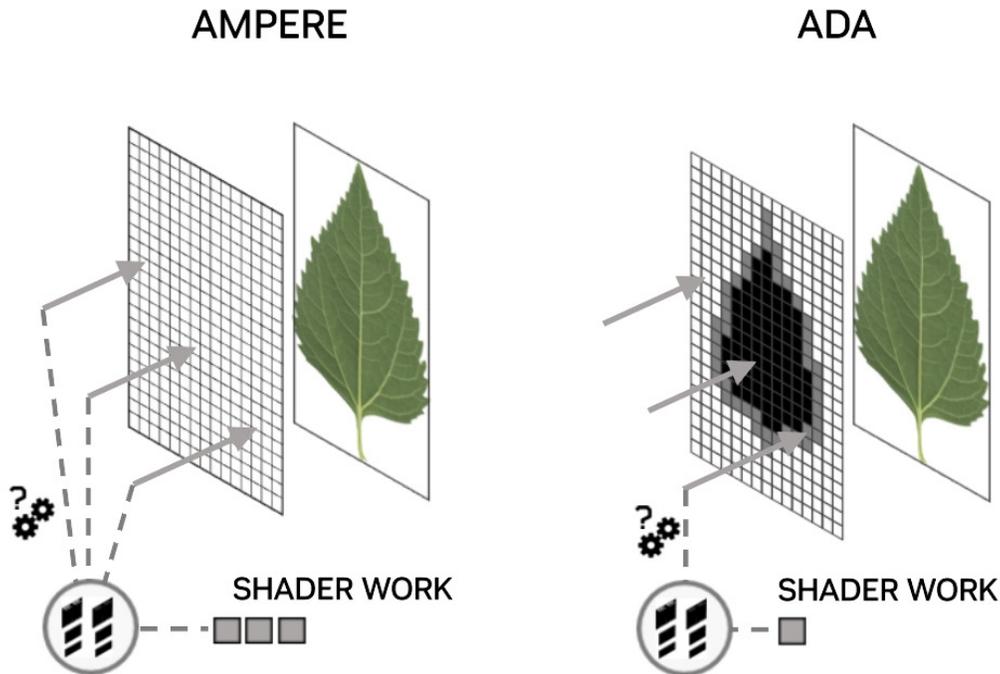
Opacity mask is applied to maple leaf which is composed of 2 triangles. The opacity engine evaluates the leaf and determines which sections are opaque, transparent, or unknown (in which case it must be sent back to the SM).

Figure 7. Opacity Mask Applied to Leaf

Sub-figure in Figure 7 above shows a pair of opacity masks, one per triangle. In the figure, transparent regions are white, they contain no leaf whatsoever. Dark green micro-triangles

correspond to opaque areas of the leaf, lastly red and blue correspond to regions of mixed opacity (*unknown*). In the example above, the Opacity Micromap Engine tags 30 of the micro-triangles as transparent, 41 as opaque, and 57 as unknown. This means that over half of the leaf is fully characterized, and that more than half of the rays intersecting these triangles either miss the leaf, or unambiguously intersect the leaf's interior. The result is that the Ada RT Core can fully characterize these rays without invoking any shader code, while preserving the full resolution and fidelity of the original alpha texture. When an *unknown* state is encountered, control is returned to a shader for resolution.

Figure 8 below shows an example of how alpha-tested content would be handled in prior GPUs.



Ada's Opacity Micromap Engine with opacity mask reduces shader work compared to Ampere

Figure 8. Ada Opacity Micromap Engine Compared to Ampere Arch

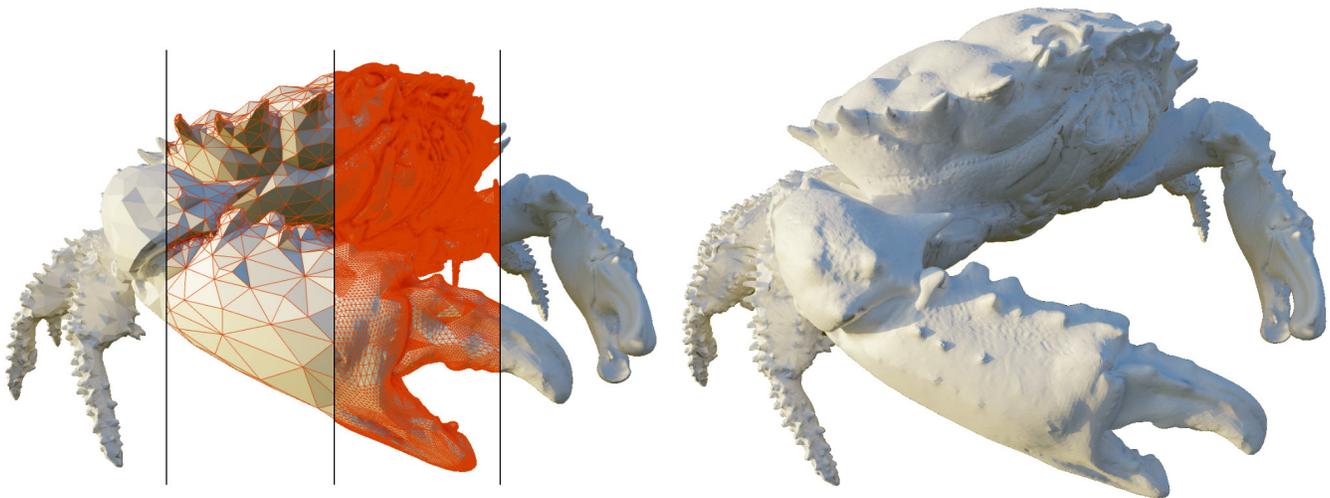
With the addition of the Opacity Micromap Engine to Ada's RT Core we have measured a doubling of scene traversal performance in applications with alpha-tested geometry. Performance gains heavily depend on usage, typically shadow rays cast against alpha-tested geometry see the largest gains. Ada's opacity mask support can significantly increase the amount and fidelity of detailed geometry within scenes, raising the realism bar.

10x Faster BVH Build in 20X Less BVH Space with Ada's Displaced Micro-Mesh Engine

Geometric complexity continues to rise with every new generation. Ray tracing performance scales attractively with increases in scene complexity. When we ray trace complex environments, tracing costs increase slowly, a one-hundred-fold increase in geometry might only double tracing time. However, creating the data structure (BVH) that makes that small increase in time possible requires roughly linear time and memory; 100x more geometry could mean 100x more BVH build

time, and 100x more memory. Ada's Third-Generation RT Core with Displaced Micro-Meshes (DMM) helps significantly with both of the challenges of high geometric complexity - BVH build performance and memory/storage footprint. Asset storage and transmission costs are reduced as well.

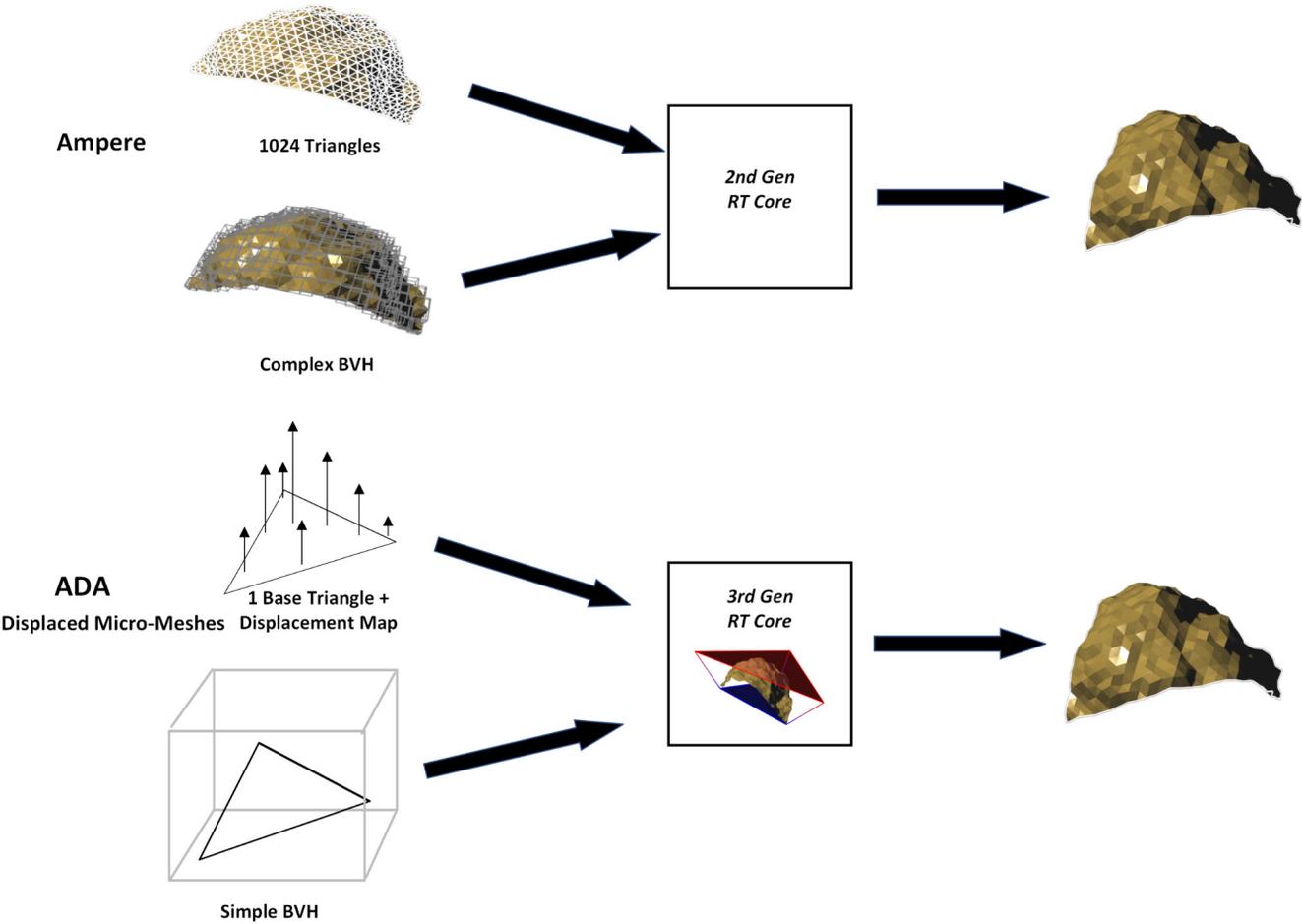
We developed Ada's displaced micro-mesh as a structured representation of geometry that exploits spatial coherence for compactness (compression) and exploits its structure for efficient rendering with an intrinsic level of detail (LOD) and light-weight animation/deformation. When ray tracing we use the displaced micro-mesh structure to avoid a large increase in BVH construction costs (time and space) while preserving efficient BVH traversal. When rasterizing we use the intrinsic micro-mesh LOD to rasterize right-sized primitives with Mesh Shaders or Compute Shaders.



Reef crab broken into base triangles represented in red (on the far left), with higher geometric detail (also in red) represented by the micro-meshes to the right. The final result is represented on the far right.

Figure 9. Displacement Micro-Mesh - Base Mesh and Micro-Meshes

The displaced micro-mesh is a new geometric primitive that was co-designed with the Micro-Mesh Engine in Ada's Third-Generation RT Core. Each micro-mesh is defined by a base triangle and a displacement map. The Micro-Mesh Engine on-demand generates micro-triangles from this definition in order to resolve ray micro-mesh intersections down to the individual micro-triangle. We use a watertight base-mesh of micro-meshes to represent highly detailed objects. We compress displacement magnitude into maps, one map per base triangle. Micro-triangle vertices are on a power-of-two, barycentric grid, and their barycentric coordinates (uv) are used to directly address micro-vertex displacements.



Ada's Third-Generation RT Core with Displaced Micro-Mesh Engine uses a simple BVH, 1 base triangle + displacement map to create a highly detailed geometric mesh with fewer required resources (both triangles and BVH structure) than Ampere's 2nd Generation RT Core.

Figure 10. DMM Simplified BVH, Base Triangle, and Displacement Map

Displaced Micro-Mesh Measured Gains



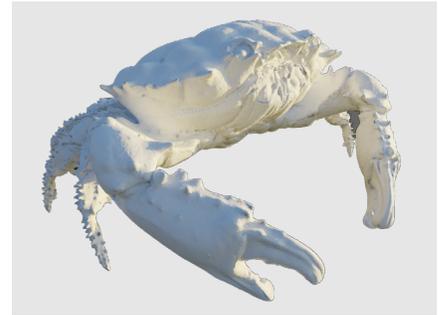
Jewel Box – 11:1

153K micro-meshes,
11M micro-triangles,
13 bits/micro-triangle
BVH build:
8.5x faster, 6.5x smaller



Ewer – 28:1

175K micro-meshes,
57M micro-triangles,
5 bits/micro-triangle
BVH build:
>15x faster, 20x smaller



Reef Crab – 14:1

17K micro-meshes,
1.6M micro-triangles,
10 bits/micro-triangle
BVH build:
7.6x faster, 8.1x smaller

Displaced Micro-Meshes allow Ada's RT Core to generate complex geometry with faster build time and reduced storage requirements.

Figure 11. DMMs Reduce BVH Build Time and Storage Requirements

Targeted launch partners for displaced micro-meshes include Adobe and Simplygon (part of Xbox Game Studios).

Shader Execution Reordering (SER)

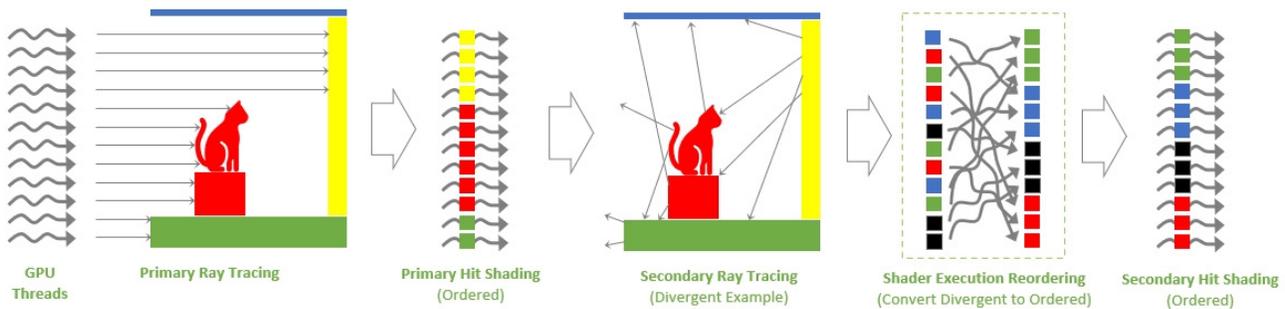
Raw RT Core horsepower is not enough to ensure high frame rates with ray-traced content, as RT workloads can be bottlenecked by a number of factors. In particular, divergent RT shaders are increasingly becoming a limiter, for example when executing multi-bounce, stochastic path tracing algorithms, or when evaluating complex materials.

Divergence takes two forms: execution divergence, where different threads execute different shaders or code paths within a shader, and data divergence, where threads access memory resources that are hard to coalesce or cache. Both types of divergence occur naturally in many ray tracing scenarios. This leaves performance on the table, because GPUs operate most efficiently when the processed work is uniform.

Ada includes a new technology designed to enhance the efficiency of RT shader execution by tackling the divergence problem. Shader Execution Reordering (SER) is a new scheduling system that reorders shading work on-the-fly for better execution and data locality. Years of research and development have been invested in SER in order to minimize overheads and maximize its

effectiveness. The Ada hardware architecture was designed with SER in mind and includes optimizations to the SM and memory system specifically targeted at efficient thread reordering.

SER is fully controlled by the application through a small API, allowing developers to easily apply reordering where their workload benefits most. The API additionally introduces new flexibility around the invocation of ray tracing shaders to the programming model, enabling more streamlined ways to structure renderer implementations while taking advantage of reordering. Furthermore, we are adding new features to the NSight Graphics shader profiler to help developers optimize their applications for SER. Developers can initially use NVIDIA-specific NVAPI extensions to implement SER in their code. We are also working with Microsoft and others to extend the standard graphics APIs with SER.



Shader Execution Reordering. Advanced lighting techniques such as path tracing cause shader divergence as secondary rays bounce off objects in the scene (denoted by various colors). In these scenarios, SER reorders shading work to improve efficiency.

Figure 12. Shader Execution Reordering Pipeline

The above diagram shows a simple ray tracing example. Starting at the top left, a number of GPU threads are shooting primary rays into a scene. Primary rays hitting the same objects can be assumed to be running the same shader program on each of the threads that hit those objects, and they are well-ordered, so the primary hit shading has high execution efficiency and data locality.

Secondary rays are generated at each primary ray hit point in the middle scene. Starting at the primary hit surfaces they shoot off in different directions, hitting different objects. Secondary hit shading tends to be less ordered and less efficient when executing on the GPU, because different shader programs are running on the different threads, and often must serialize execution. Examples of secondary rays that can benefit from SER include those used for path tracing, reflections, indirect lighting, and translucency effects.

Shader Execution Reordering adds a new stage in the ray tracing pipeline which reorders and groups the secondary hit shading to have better execution locality, thus much higher overall ray-traced shading efficiency.

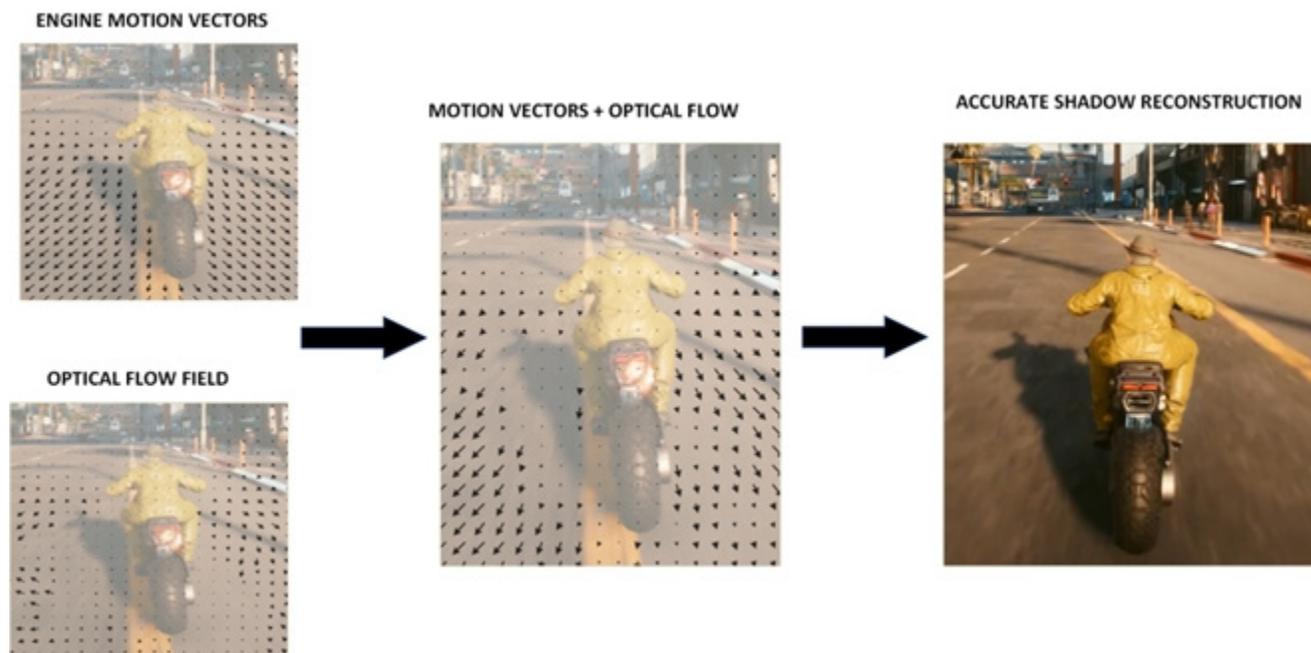
SER can often provide up to 2X performance improvement for RT shaders in cases with a high level of divergence (such as path tracing). In testing with Cyberpunk 2077 running in RT: Overdrive Mode, we've measured overall performance gains of up to 44% from SER.

DLSS 3 and Optical Flow Acceleration

Over the past four years, the NVIDIA Applied Deep Learning Research team has been developing a frame generation technique that combines optical flow estimation with DLSS to improve the gaming experience. Inserting accurate synthesized frames between existing frames improves frame rate and provides a smoother gaming experience.

Optical flow estimation is commonly used in computer vision applications to measure the direction and magnitude of the apparent motion of pixels between consecutively rendered graphics frames or video frames. In the 3D graphics and video fields, typical use cases have included reducing latency in augmented and virtual reality, improving smoothness of video playback, enhancing video compression efficiency, and enabling video camera stabilization. Deep learning uses often include automotive and robotic navigation, and video analysis and understanding.

Optical flow is superficially similar to the motion estimation component of video encoding, but with much more challenging requirements for accuracy and consistency. As a result, different algorithms are used. Starting with the Ampere GPU architecture, NVIDIA's GPUs have had support for a standalone optical flow engine (OFA) that uses state of the art algorithms to ensure high quality results. Ada's OFA unit delivers 300 TeraOPS (TOPS) of optical flow work (over 2x faster than the Ampere generation OFA) and provides critical information to the DLSS 3 network.



To generate more accurate frames without distracting artifacts, DLSS 3 combines game engine motion vectors with the OFA Engine's optical flow field.

Figure 13. DLSS 3 Motion Vectors + Optical Flow = Accurate Motion Estimation

The Ada OFA unit and new motion vector analysis algorithms are fundamental components that enable accurate and performant frame generation capability within the new DLSS 3 technology framework. This new DL-based frame generation method improves frame rates by an additional 2x over DLSS 2. When DLSS 3 is combined with the new RT Core and other Ada architecture enhancements, Ada is up to 4x faster than prior GPUs.

DLSS 3 can also improve performance in cases where the GPU is bottlenecked by the CPU. For instance, *Microsoft Flight Simulator* is a common example of a game that is CPU-limited because of its physics and gigantic draw distances. This limits the performance benefits that traditional super resolution technologies can offer. In this case however, DLSS 3's ability to generate frames still provides up to a 2x performance improvement. For more information on OFA and DLSS 3, please read the **NVIDIA Ada Science Whitepaper**.

Ada Fourth-Generation Tensor Core

Tensor Cores are specialized high performance compute cores that are tailored for the matrix multiply and accumulate math operations that are used in AI and HPC applications. Tensor Cores provide groundbreaking performance for the matrix computations that are critical for deep learning neural network training and inference functions that occur at the edge.

Compared to Ampere, Ada delivers more than double the FP16, BF16, TF32, INT8, and INT4 Tensor TFLOPS, and also includes the Hopper FP8 Transformer Engine, delivering over 1.3 PetaFLOPS of tensor processing in the RTX 4090.

NVIDIA Broadcast/Video

As a result of the global pandemic, the PC broadcast industry has seen explosive growth. Gamers in particular are streaming more than ever using services such as *Twitch*. In the past, streaming on your own was challenging – to improve performance, a dedicated PC was commonly used solely for video capture, and users were often gaming **while** they were also streaming.

NVIDIA innovations have democratized streaming so that more people can easily stream on their PC without the cost and complexity challenges that have traditionally been an issue for users in the past. The introduction of NVIDIA's NVENC encoder and optimizations for OBS (Open Broadcaster Software) eliminated the need for a dedicated PC for video capture, allowing users to play and stream from the same PC with good stream quality and high fps in games. Finally, NVIDIA's Broadcast suite, powered by AI, provides tools for noise and echo removal, virtual backgrounds, video noise removal, and automatic camera tracking so that anyone can start and work their way up to their dream streaming setup without the obligation to purchase professional microphones, cameras, and need a dedicated recording studio to stream.

Ada GPUs take streaming and video content to the next level, incorporating support for AV1 video encoding in the Ada **eighth generation dedicated hardware encoder** (known as NVENC). Prior generation Ampere GPUs supported AV1 decoding, but not encoding. Ada's AV1 encoder is 40% more efficient than the H.264 encoder used in GeForce RTX 30 Series GPUs. AV1 will enable users who are streaming at 1080p today to increase their stream resolution to 1440p while running at the same bitrate and quality, or for users with 1080p displays, streams will look similar to 1440p, providing better quality.

NVIDIA collaborated with OBS Studio to add AV1 — on top of the recently released HEVC and HDR support — within an upcoming software release, expected later this year. OBS is also optimizing encoding pipelines to reduce overhead by 35% for all NVIDIA GPUs. The new release will additionally feature updated NVIDIA Broadcast effects, including noise and room echo removal, as well as improvements to virtual background.

We've also worked with Discord to enable end-to-end livestreams with AV1. In an update releasing later this year, Discord will enable its users to use AV1 to dramatically improve screen sharing, be it for game play, schoolwork, or hangouts with friends.

To further aid encoding performance, GeForce RTX 4090 and RTX 4080 are equipped with **dual NVENC encoders**. This enables video encoding at 8K/60 for professional video editing or four 4K/60. (Game streaming services can also take advantage of this to enable more simultaneous sessions, for instance.) Blackmagic Design's DaVinci Resolve, the popular Vookoder plugin for Adobe Premiere Pro, and Jianying — the top video editing app in China — are all enabling AV1 support, as well as a dual encoder through encode presets. Dual encoder and AV1 availability for these apps will be available in October. NVIDIA is also working with the popular video-effects app Notch to enable AV1, as well as Topaz to enable support for AV1 and the dual encoders.

In addition to NVENC, Ada GPUs also include the **fifth-generation hardware decoder** that was first launched with Ampere (known as NVDEC). NVDEC supports hardware-accelerated video decoding of MPEG-2, VC-1, H.264 (AVCHD), H.265 (HEVC), VP8, VP9, and the AV1 video formats. 8K/60 decoding is also fully supported. In the future, NVIDIA is also working to enable high quality video production using AI. For more information on this topic, please read the **NVIDIA Ada Science Whitepaper**.

Conclusion

With frame rates that are up to 4x faster than the previous generation, the Ada Lovelace GPU architecture provides performance that is beyond fast, delivering NVIDIA's greatest generational upgrade ever. Ada's record-breaking performance is made possible by a number of engineering innovations.

NVIDIA engineers worked closely with TSMC to create the 4N manufacturing process that is tailored for NVIDIA GPUs. The smaller process allows more processing units and memories to be integrated into the chip. NVIDIA's AD102 GPU contains 18,432 CUDA Cores (70% more CUDA Cores than Ampere), 18 MB of L1 cache, 96 MB of L2 cache (16x more than Ampere), and a large 36 MB register file. The entire GPU contains over 76 billion transistors, making it second only to NVIDIA's H100 in terms of GPU complexity. Even though the GeForce RTX 4090 runs at a Boost Clock frequency of 2.5 GHz – 660 MHz higher than the previous GeForce flagship RTX 3090 Ti – it consumes the same TGP of 450W. Ultimately Ada delivers 2x higher power efficiency compared to prior generation Ampere. It is truly a marvel of engineering.

As massive as Ada's core counts, memory, and clocks are, the Ada GPU is about more than just those figures. The Ada SM has been significantly enhanced, especially for ray tracing workloads. Ada's Third-Generation RT Core offers 2x faster ray-triangle intersection throughput over prior generation Ampere GPUs (and 4x faster than Turing). Triangle intersection testing is a computationally expensive operation that is very commonly performed when rendering a ray-traced scene, so providing a 2x improvement is very significant.

The Ada RT Core also includes two new hardware units. The first, an Opacity Micromap Engine, speeds up alpha traversal by 2x. With this new capability, developers can very quickly assign opacity values to irregularly shaped objects (like ferns and fences) or translucent items (like flames or smoke) allowing the Ada RT Core to directly alpha test this geometry instead of relying on the GPU's SM.

The second new hardware unit that has been incorporated into the Ada RT Core is the Displaced Micro-Mesh Engine. The new Micro-Mesh Engine has been designed to reduce the BVH build time and storage requirements that are traditionally required when dealing with complex objects with high levels of geometric detail. With this new feature, a new displaced micro-mesh primitive has been developed for ray tracing. The Micro-Mesh Engine evaluates the micro-meshes, and when additional geometric detail is needed, the Micro-Mesh Engine can dynamically generate additional micro-triangles as needed. Compared to traditionally rendering these complex objects, the Micro-Mesh Engine reduces BVH build time by a factor of 10x, while reducing BVH storage requirements by a factor of 20x.

In addition, Ada introduces the Shader Execution Reordering (SER) scheduling system. Shader Execution Reordering organizes and reorders workloads on the fly so they can be processed by the SM and RT Core more efficiently. Shader Execution Reordering is as big of an innovation for GPUs as out-of-order execution was for CPUs back in the 1990s, offering 2-3x speedups for some RT workloads.

When combined, the improvements to SM throughput, higher clocks and core counts, Ada's Third-Generation RT Core, and new features such as Shader Execution Reordering all provide the Ada GPU with a performance uplift of up to 2x over Ampere. So how did we ultimately get to a

generational performance uplift of up to 4x? The remainder comes from the Ada GPU's new Optical Flow Accelerator and DLSS 3.

NVIDIA's DLSS technology pioneered the concept of AI-based neural graphics. To date, 216 titles take advantage of DLSS, and the list continues to grow. Ada's new DLSS 3 technology builds on DLSS 2 Super Resolution, which internally renders using lower resolution pixels and uses AI algorithms to produce beautiful, sharp higher resolution images to dramatically improve performance compared to traditional raster-based graphics rendering.

DLSS 3 harnesses the new Optical Flow Accelerator found in Ada GPUs to calculate the motion flow of every pixel in a given frame. This data is combined with the traditional motion vectors that are already commonly used in games and fed into the AI network, which then generates entire frames rather than just pixels. This combination of DLSS Super Resolution and DLSS Frame Generation also incorporates NVIDIA Reflex to ensure the lowest latency possible for gamers. The result is up to a 2x performance boost over DLSS 2, while preserving the exquisite image quality DLSS is already known for.

DLSS 3 can also improve performance in CPU-bound cases that occur when the GPU is bottlenecked by the CPU and is therefore unable to generate higher frame rates. Because DLSS 3 is able to generate frames independent of the CPU, Ada GPUs with DLSS 3 are still capable of improving performance in these cases. An example of this is *Microsoft Flight Simulator* where performance can be doubled using DLSS 3.

With the introduction of DLSS 3, neural graphics are taken to an entirely new level. DLSS 3 is easy for developers to integrate and is one of the most compelling features that NVIDIA has introduced for a new graphics architecture. 35 titles have already been confirmed to support DLSS 3, with the first games coming in October.

Ada GPUs include a new Fourth-Generation Tensor Core. The GeForce RTX 4090 offers double the throughput for existing FP16, BF16, TF32, and INT8 formats, and its Fourth-Generation Tensor Core introduces support for a new FP8 tensor format. Compared to FP16, FP8 halves the data storage requirements and doubles throughput. With the new FP8 format, the GeForce RTX 4090 delivers 1.3 PetaFLOPS of performance for AI inference workloads.

All Ada GPUs ship with NVIDIA's 8th Generation NVENC encoder, which adds support for AV1 encoding. AV1 is 40% more efficient than the prior H.264 encoder that was commonly used previously. For streamers, this will allow livestreams to look as if they are using 40% higher bitrate. NVIDIA is working with OBS Studio to integrate AV1 support into its next software update, which is expected to be released in October 2022. OBS is also optimizing encoding pipelines to reduce overhead by 35% for all NVIDIA GPUs. Discord is also adding AV1 later this year.

Additionally, GeForce RTX 4090 and RTX 4080 are equipped with dual NVENC encoders, providing support for 8K/60 video encoding or four 4K/60 for video editing. NVIDIA has worked with all of the top software makers to integrate AV1 support, with DaVinci Resolve, Voukoder (a plugin for Adobe Premiere Pro), and Jianying (China's most popular video editing app), providing software updates that are releasing in October. For gamers with GeForce Experience, NVIDIA ShadowPlay will also support 8K/60.

The NVIDIA Ada Lovelace architecture delivers a quantum leap in GPU performance and capabilities, giving GeForce RTX 40 Series users the power to experience the next generation of fully ray-traced games beginning with the introduction of the GeForce RTX 4090 and 4080 GPUs this Fall. They truly are beyond fast.

Appendix A - GeForce RTX 4090 GPU Full Specifications

Table 2. GeForce RTX 4090 vs RTX 3090 Ti vs 2080 Ti

Graphics Card	GeForce RTX 2080 Ti	GeForce RTX 3090 Ti	GeForce RTX 4090
GPU Codename	TU102	GA102	AD102
GPU Architecture	NVIDIA Turing	NVIDIA Ampere	NVIDIA Ada Lovelace
GPCs	6	7	11
TPCs	34	42	64
SMs	68	84	128
CUDA Cores / SM	64	128	128
CUDA Cores / GPU	4352	10752	16384
Tensor Cores / SM	8 (2nd Gen)	4 (3rd Gen)	4 (4th Gen)
Tensor Cores / GPU	544	336 (3rd Gen)	512 (4th Gen)
OFA TOPS ³	N/A	126	305
RT Cores	68 (1st Gen)	84 (2nd Gen)	128 (3rd Gen)
GPU Boost Clock (MHz)	1635	1860	2520
Peak FP32 TFLOPS (non-Tensor) ¹	14.2	40	82.6
Peak FP16 TFLOPS (non-Tensor) ¹	28.5	40	82.6

Appendix A - GeForce RTX 4090 Full Specifications

Peak BF16 TFLOPS (non-Tensor)¹	N/A	40	82.6
Peak INT32 TOPS (non-Tensor)^{1,3}	14.2	20	41.3
RT TFLOPS	42.9	78.1	191
Peak FP8 Tensor TFLOPS with FP16 Accumulate¹	N/A	N/A	660.6/1321.2 ²
Peak FP8 Tensor TFLOPS with FP32 Accumulate¹	N/A	N/A	330.3/660.6 ^{2,4}
Peak FP16 Tensor TFLOPS with FP16 Accumulate¹	113.8	160/320 ²	330.3/660.6 ²
Peak FP16 Tensor TFLOPS with FP32 Accumulate¹	56.9	80/160 ²	165.2/330.4 ²
Peak BF16 Tensor TFLOPS with FP32 Accumulate¹	N/A	80/160 ²	165.2/330.4 ²
Peak TF32 Tensor TFLOPS¹	N/A	40/80 ²	82.6/165.2 ²
Peak INT8 Tensor TOPS¹	227.7	320/640 ²	660.6/1321.2 ²
Peak INT4 Tensor TOPS¹	455.4	640/1280 ²	1321.2/2642.4 ²
Frame Buffer Memory Size and Type	11 GB GDDR6	24 GB GDDR6X	24 GB GDDR6X
Memory Interface	352-bit	384-bit	384-bit
Memory Clock (Data Rate)	14 Gbps	21 Gbps	21 Gbps
Memory Bandwidth	616 GB/sec	1008 GB/sec	1008 GB/sec

Appendix A - GeForce RTX 4090 Full Specifications

ROPs	88	112	176
Pixel Fill-rate (Gigapixels/sec)	143.9	208.3	443.5
Texture Units	272	336	512
Texel Fill-rate (Gigatexels/sec)	444.7	625	1290.2
L1 Data Cache/Shared Memory	6528 KB	10752 KB	16384 KB
L2 Cache Size	5632 KB	6144 KB	73728 KB
Register File Size	17408 KB	21504 KB	32768 KB
Video Engines	1 x NVENC (7th Gen) 1 x NVDEC (4th Gen)	1 x NVENC (7th Gen) 1 x NVDEC (5th Gen)	2 x NVENC (8th Gen) 1 x NVDEC (5th Gen)
TGP (Total Graphics Power)	260 W	450 W	450 W
Transistor Count	18.6 Billion	28.3 Billion	76.3 Billion
Die Size	754 mm ²	628.4 mm ²	608.5 mm ²
Manufacturing Process	TSMC 12 nm FFN (FinFET NVIDIA)	Samsung 8 nm 8N NVIDIA Custom Process	TSMC 4N NVIDIA Custom Process
PCI Express Interface	Gen 3	Gen 4	Gen 4

1. Peak rates are based on GPU Boost Clock.
2. Effective TOPS / TFLOPS using the new Sparsity Feature
3. TOPS = IMAD-based integer math
4. Peak FP8 Tensor TFLOPS with FP32 Accumulate changed to proper number in v.2.02 of this whitepaper

Appendix B - Ada AD103 GPU Full Specifications

The Ada GPU architecture consists of a family of GPUs that are targeted for multiple segments of the graphics market. AD102 is NVIDIA’s flagship GPU offering based on the Ada architecture, delivering revolutionary performance for the ultra-enthusiast graphics segment.

The AD103 GPU is NVIDIA’s product for the high-end graphics segment. AD103 retains all of the key features found in AD102, including all of the innovations introduced with the Ada SM such as Ada’s Third-Generation RT Core and Fourth-Generation Tensor Core as well as DLSS 3 and Ada’s new OFA unit. AD103 brings record-breaking performance to the market and is faster than NVIDIA’s prior generation flagship GA102 GPU that was used in products such as the GeForce RTX 3090 Ti.

The full AD103 chip consists of 45.9 billion transistors and contains 7 GPCs, 40 TPCs, 80 SMs, and eight 32-bit memory controllers (256-bit total). With each SM containing 128 FP32 CUDA Cores, the full chip contains 10,240 CUDA Cores as well as 80 RT Cores, 320 Tensor Cores, 320 Texture Units, and 112 ROPS. The memory subsystem includes 10,240 KB L1 cache, 20,480 KB Register File, and 65,536 KB L2 cache.

The first GeForce RTX 40 series product that will be launching using the AD103 GPU is the GeForce RTX 4080 16 GB.

Table 3. GeForce RTX 4080 16 GB vs 3080 Ti vs 2080 Super

Graphics Card	RTX 2080 Super	RTX 3080 Ti	RTX 4080 16 GB
GPU Codename	TU104	GA102	AD103
GPU Architecture	NVIDIA Turing	NVIDIA Ampere	NVIDIA Ada Lovelace
GPCs	6	7	7
TPCs	24	40	38
SMs	48	80	76
CUDA Cores / SM	64	128	128
CUDA Cores / GPU	3072	10240	9728

Tensor Cores / SM	8 (2nd Gen)	4 (3rd Gen)	4 (4th Gen)
Tensor Cores / GPU	384	320 (3rd Gen)	304 (4th Gen)
OFA	N/A	126	305
RT Cores	48 (1st Gen)	80 (2nd Gen)	76 (3rd Gen)
GPU Boost Clock (MHz)	1815	1665	2505
Peak FP32 TFLOPS (non-Tensor)¹	11.2	34.1	48.7
Peak FP16 TFLOPS (non-Tensor)¹	22.3	34.1	48.7
Peak BF16 TFLOPS (non-Tensor)¹	N/A	34.1	48.7
Peak INT32 TOPS (non-Tensor)^{1,3}	11.2	17	24.4
RT TFLOPS	33.6	66.6	112.7
Peak FP8 Tensor TFLOPS with FP16 Accumulate¹	N/A	N/A	389.9/779.8 ²
Peak FP8 Tensor TFLOPS with FP32 Accumulate¹	N/A	N/A	194.9/389.8 ^{2,4}
Peak FP16 Tensor TFLOPS with FP16 Accumulate¹	89.2	136.4/272.8 ²	194.9/389.8 ²
Peak FP16 Tensor TFLOPS with FP32 Accumulate¹	44.6	68.2/136.4 ²	97.5/195 ²

Peak BF16 Tensor TFLOPS with FP32 Accumulate¹	N/A	68.2/136.4 ²	97.5/195 ²
Peak TF32 Tensor TFLOPS¹	N/A	34.1/68.2 ²	48.7/97.4 ²
Peak INT8 Tensor TOPS¹	178.4	272.8/545.6 ²	389.9/779.82 ²
Peak INT4 Tensor TOPS¹	356.8	545.6/1091.2 ²	779.8/1559.6 ²
Frame Buffer Memory Size and Type	8 GB GDDR6	12 GB GDDR6X	16 GB GDDR6X
Memory Interface	256-bit	384-bit	256-bit
Memory Clock (Data Rate)	15.5 Gbps	19 Gbps	22.4 Gbps
Memory Bandwidth	496 GB/sec	912 GB/sec	716.8 GB/sec
ROPs	64	112	112
Pixel Fill-rate (Gigapixels/sec)	116.2	186.5	280.6
Texture Units	192	320	304
Texel Fill-rate (Gigatexels/sec)	348.5	532.8	761.5
L1 Data Cache/Shared Memory	4608 KB	10240 KB	9728 KB
L2 Cache Size	4096 KB	6144 KB	65536 KB
Register File Size	12288 KB	20480 KB	19456 KB

Video Engines	1 x NVENC (7th Gen) 2 x NVDEC (4th Gen)	1 x NVENC (7th Gen) 1 x NVDEC (5th Gen)	2 x NVENC (8th Gen) 1 x NVDEC (5th Gen)
TGP (Total Graphics Power)	250 W	350 W	320 W
Transistor Count	13.6 Billion	28.3 Billion	45.9 Billion
Die Size	545 mm ²	628.4 mm ²	378.6 mm ²
Manufacturing Process	TSMC 12 nm FFN (FinFET NVIDIA)	Samsung 8 nm 8N NVIDIA Custom Process	TSMC 4N NVIDIA Custom Process
PCI Express Interface	Gen 3	Gen 4	Gen 4

1. Peak rates are based on GPU Boost Clock.
2. Effective TOPS / TFLOPS using the new Sparsity Feature
3. TOPS = IMAD-based integer math
4. Peak FP8 Tensor TFLOPS with FP32 Accumulate changed to proper number in v.2.02 of this whitepaper

Appendix C - NVIDIA L40 GPU Full Specifications

The NVIDIA L40 GPU Accelerator is a full height, full-length (FHFL), dual-slot 10.5 inch PCI Express Gen4 graphics solution based on the latest NVIDIA Ada Lovelace Architecture. The card is passively cooled and consuming up to 300 W maximum board power. The NVIDIA L40 supports the latest hardware-accelerated ray tracing, revolutionary AI features, advanced shading, and powerful simulation capabilities for a wide range of graphics and compute use cases in data center and edge server deployments. This includes NVIDIA Omniverse™, cloud gaming, batch rendering, virtual workstations, and deep learning training, as well as inference workloads.

As part of the NVIDIA OVX™ server platform, L40 delivers the highest level of graphics, ray tracing, and simulation performance for NVIDIA Omniverse. With 48 GB of GDDR6 memory, even the most intense graphics applications run with the highest level of performance.

L40 accelerates high-fidelity creative workflows, including real-time, full-fidelity, interactive ray-traced rendering, 3D design, video streaming, and virtual production. Running professional 3D visualization applications with NVIDIA L40 enables creative professionals to iterate more, render faster, and unlock tremendous performance advantages that increase productivity and speed up project completion.

Powerful training and inference performance, combined with enterprise-class stability and reliability, make the NVIDIA L40 the ideal platform for single-GPU AI training and development. The NVIDIA L40 reduces the time to completion for model training and development, and also data science data prep workflows by delivering higher throughput and support for a full range of precisions, including FP8.

The NVIDIA L40 takes streaming and video content workloads to the next level with three video encode and three video decode engines. With the addition of AV1 encoding, the L40 delivers breakthrough performance and improved TCO for broadcast streaming, video production, and transcription workflows.

Table 4. NVIDIA L40 vs NVIDIA A40

Graphics Card	NVIDIA A40	NVIDIA L40
GPU Architecture	NVIDIA Ampere	NVIDIA Ada Lovelace
GPU Code Name	GA102	AD102
GPCs	7	12
TPCs	42	71

Appendix C – NVIDIA L40 GPU Full Specifications

SMs	84	142
CUDA Cores / SM	128	128
CUDA Cores / GPU	10752	18176
Tensor Cores / SM	4 (3rd Gen)	4 (4th Gen)
Tensor Cores/GPU	336	568
RT Cores	84 (2nd Gen)	142 (3rd Gen)
OFA	126	307
GPU Memory	48GB GDDR6 w/ECC	48GB GDDR6 w/ ECC
GPU Memory Bandwidth	696 GB/s	864 GB/s
L2 Cache Size	6144 KB	98304 KB
FP32 Performance	37.4 TFLOPS	90.5 TFLOPS
RT Core Performance	73.1 TFLOPS	209.3 TFLOPS
Tensor Float 32 (TF32) Performance	74.8 149.6 TFLOPS	90.5 181 TFLOPS ¹
BFLOAT16 Tensor Core Performance	149.7 299.4 TFLOPS	181 362 TFLOPS ¹
FP16 Tensor Core Performance	149.7 299.4 TFLOPS	181 362 TFLOPS ¹
FP8 Tensor Core Performance	NA	362 724 TFLOPS ¹
INT8 Tensor Core Performance	299.3 598.6 TOPS	362 724 TOPS ¹
INT4 Tensor Core Performance	598.7 1197.4 TOPS	724 1448 TOPS ¹

Appendix C – NVIDIA L40 GPU Full Specifications

Video Engines	1 Encoder, 2 Decoders	3 Encoders, 3 Decoders, 4 JPEG Decoders
Max Thermal Design Power (TDP)	300 Watts	300 Watts
Form Factor	4.4" H x 10.5" L – Dual Slot	4.4" H x 10.5" L – Dual Slot

1. Effective TOPS / TFLOPS using the new Sparsity Feature

Appendix D - NVIDIA L4 GPU Full Specifications

The NVIDIA L4 datacenter GPU based on the Ada architecture is designed for universal datacenter workloads spanning graphics, AI, and video streaming. It is designed to meet the needs across not only AI inference, but also video, graphics, virtualization, and numerous other applications including cloud gaming, simulation, and data science. It's a true universal GPU in a low-profile form factor with a sub-75W power profile that delivers a cost-effective, energy-efficient solution for high throughput and low latency in servers from the edge, to the data center, to the cloud.

NVIDIA L4 is optimized for 24/7 enterprise data center operations and is designed, built, extensively tested, and supported by NVIDIA and partners for maximum performance, durability, and security. L4 features secure boot with root-of-trust technology, providing an additional layer of security for data centers.

NVIDIA L4 is the best low power universal GPU for a wide variety of applications such AI-powered video services, Speech AI (ASR+NLP+TTS), small model Generative AI, search & recommenders, cloud gaming, and virtual workstations, among many others.

Table 5. NVIDIA L4 GPU vs NVIDIA T4

Graphics Card	NVIDIA T4	NVIDIA L4
GPU Architecture	NVIDIA Turing	NVIDIA Ada Lovelace
GPU Code Name	TU104	AD104
GPCs	5	5
TPCs	20	29
SMs	40	58
CUDA Cores / SM	64	128
CUDA Cores / GPU	2560	7424
Tensor Cores / SM	8 (2nd Gen)	4 (4th Gen)
Tensor Cores/GPU	320	232

Appendix D – NVIDIA L4 GPU Full Specifications

RT Cores	40 (1st Gen)	58 (3rd Gen)
OFA	N/A	281
GPU Memory	16GB GDDR6 w/ECC	24GB GDDR6 w/ ECC
GPU Memory Bandwidth	320 GB/s	300 GB/s
L2 Cache Size	4096 KB	49152 KB
FP32 Performance	8.1 TFLOPS	30.3 TFLOPS
RT Core Performance	28.9 TFLOPS	73.1 TFLOPS
Tensor Float 32 (TF32) Performance	NA	60 120 TFLOPS ¹
BFLOAT16 Tensor Core Performance	NA	121 242 TFLOPS ¹
FP16 Tensor Core Performance	65 TFLOPS	121 242 TFLOPS ¹
FP8 Tensor Core Performance	NA	242 485 TFLOPS ¹
INT8 Tensor Core Performance	130 TOPS	242 485 TOPS ¹
INT4 Tensor Core Performance	260 TOPS	484 969 TOPS ¹
Video Engines	1 Encoder, 2 Decoders	2 Encoders, 4 Decoders, 4 JPEG Decoders
Max Thermal Design Power (TDP)	70 Watts	72 Watts
Form Factor	2.71" H x 6.67" L – Single Slot	2.71" H x 6.67" L – Single Slot

1. Effective TOPS / TFLOPS using the new Sparsity Feature

Notice - The information provided in this specification is believed to be accurate and reliable as of the date provided. However, NVIDIA Corporation (“NVIDIA”) does not give any representations or warranties, expressed or implied, as to the accuracy or completeness of such information. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This publication supersedes and replaces all other specifications for the product that may have been previously supplied.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and other changes to this specification, at any time and/or to discontinue any product or service without notice. Customer should obtain the latest relevant specification before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer. NVIDIA hereby expressly objects to applying any customer general terms and conditions with regard to the purchase of the NVIDIA product referenced in this specification.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer’s own risk.

NVIDIA makes no representation or warranty that products based on these specifications will be suitable for any specified use without further testing or modification. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer’s sole responsibility to ensure the product is suitable and fit for the application planned by customer and to do the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer’s product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this specification. NVIDIA does not accept any liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this specification, or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this specification. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA. Reproduction of information in this specification is permissible only if reproduction is approved by NVIDIA in writing, is reproduced without alteration, and is accompanied by all associated conditions, limitations, and notices.

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, “MATERIALS”) ARE BEING PROVIDED “AS IS.” NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA’s aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the NVIDIA terms and conditions of sale for the product.

Trademarks - NVIDIA, the NVIDIA logo, GeForce, and GeForce RTX are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright - © 2023 NVIDIA Corporation. All rights reserved.