



DAYHOFF
HEALTH

Accelerating Genomics Analysis with AMD ROCm™ and HIP

Authors

Stephen Elliot
Jacob Shultis

March 2, 2026

CONTENTS

CONTENTS	1
ABSTRACT	2
INTRODUCTION	2
WHY ACCELERATE GENOMICS WITH AMD?	3
AMD ROCm™ and HIP	3
GENOMICS WORKLOADS AND AMD GPUS	3
METHODOLOGY AND TEST ENVIRONMENT	4
Test Environment	4
Test Protocol	5
WHOLE GENOME SEQUENCING (WGS) ANALYSIS	5
Overview	5
GPU Accelerated Performance for WGS Pipeline	6
Performance Drivers for WGS	7
Clinical Impact	7
SINGLE-CELL SEQUENCING (SCS) FILE PROCESSING	8
Overview	8
GPU Accelerated Performance for SCS Processing	8
Performance Drivers for SCS	9
Clinical Impact	9
MICROBIOME SAMPLE ANALYSIS	10
Overview	10
GPU Accelerated Performance for Microbiome Analysis	10
Microbiome Performance Details	11
Clinical Impact	12
RESULTS AND OBSERVATIONS	12
Performance Summary Table	12
Key Observations	13
CONCLUSION	13
TECHNICAL SPECIFICATIONS	14

ABSTRACT

This white paper outlines the methodology and expected performance improvements of using AMD Radeon™ AI PRO R9000 Series GPUs with the ROCm™ platform and the Heterogeneous-Compute Interface for Portability (HIP) to accelerate complex genomics workflows, including whole genome sequencing (WGS), single-cell sequencing (SCS) file processing, and microbiome sample analysis. The objective is to demonstrate how this advanced compute environment, leveraging the high parallel processing capabilities of AMD GPUs, can significantly reduce data processing turnaround times and increase operational efficiency for life sciences and high-performance computing (HPC) organizations.

This white paper highlights the broader impact of adopting GPU-accelerated genomics workflows. Faster turnaround across WGS, SCS, and microbiome pipelines shortens the path from raw data to actionable clinical findings, benefiting hospitals, clinicians, and public-health agencies through quicker testing, earlier detection, and more responsive population-health interventions. Greater computational efficiency lowers costs and increases throughput for commercial laboratories, enabling them to process more samples with fewer resources. For bioinformaticians and researchers, near real-time pipeline execution provides the ability to rapidly test, refine, and adjust configurations as results appear, accelerating development cycles and improving overall workflow performance.

INTRODUCTION

There is a rapidly increasing volume of genomic data being generated across research and clinical settings. Processing this data—which includes large datasets from whole genome sequencing, single-cell experiments, and diverse microbiome samples—requires substantial compute resources and can often bottleneck research progress. Traditional CPU-based environments struggle to keep pace with the demand for rapid, high-throughput analysis.

This research investigates the performance benefits of leveraging AMD's high-performance compute ecosystem, specifically the ROCm™ open software platform and the HIP programming model, to accelerate the secondary analysis phase of various genomics applications. Testing was conducted on a production AMD Radeon™ AI PRO R9700 GPU (gfx1201, RDNA4 architecture) with ROCm™ 7.0.2, demonstrating significant improvements in processing efficiency and turnaround time for clinical and research genomics workflows.

WHY ACCELERATE GENOMICS WITH AMD?

High-Performance Computing (HPC) environments benefit from the massive parallelism offered by modern GPUs. Genomic analysis, which often involves performing the same complex calculations across millions of DNA or RNA molecules, is highly suitable for GPU acceleration.

AMD ROCm™ and HIP

The AMD ROCm™ platform is designed for GPU computing and features a rich set of tools and drivers. HIP provides a programming interface to enable portability of accelerated applications, simplifying the transition of existing codebases to run efficiently on AMD GPUs. This combination allows for:

Higher Throughput: Processing multiple samples or parallel threads simultaneously. The AMD Radeon™ AI PRO R9700 GPU delivers 37.5 million reads per second compared to 113,000 reads per second on a CPU.

Reduced Turnaround Time: Decreasing the wall-clock time required for crucial secondary analysis tasks. Microbiome analysis processing time reduced from **264 seconds to 0.8 seconds (330x faster)**.

Simplified Infrastructure: Consolidating compute resources and reducing the overhead associated with managing large, multi-core CPU clusters. The dedicated 32GB GPU VRAM eliminates memory pressure and disk swapping that degrades CPU performance.

Open-Source Ecosystem: ROCm™ is open-source and community-driven, providing flexibility and long-term sustainability without vendor lock-in.

GENOMICS WORKLOADS AND AMD GPUS

The genomics process generally involves three phases: primary, secondary, and tertiary analysis. This paper focuses on accelerating the secondary analysis phase, where raw

sequencing data (FASTQ files) is aligned, sorted, and variant-called (VCF files).

Genomics Workload	Secondary Analysis Phase	Acceleration Focus	AMD R9700 GPU Advantage
Whole Genome Sequencing (WGS)	Alignment (e.g., BWA-MEM), Sorting, Variant Calling (e.g., HaplotypeCaller)	Raw computation speed and memory bandwidth	640 GB/s bandwidth, ~48 TFLOPS FP32 Compute
Single-Cell Sequencing (SCS)	Read alignment, Feature Barcoding, UMI deduplication, Quality Control	Parallel processing of highly-specific molecular identifiers and large matrix operations	4,096 concurrent threads, 64-thread waves (RDNA4 optimized)
Microbiome Sample Analysis	Taxonomy assignment via k-mer matching (Kraken2), OTU clustering, Phylogenetic analysis	Sequence comparison, database queries on large reference databases	330x speedup on 80GB database queries (memory-bound workload)

METHODOLOGY AND TEST ENVIRONMENT

Test Environment

The acceleration was tested on a platform utilizing:

GPU Compute: AMD Radeon™ AI PRO R9700 (gfx1201, RDNA4 architecture) with 32GB VRAM running ROCm™ 7.0.2

CPU Baseline: AMD Ryzen 7950X (16-Core, Zen4 architecture) with 64GB DDR5 system RAM

Software: HIP-optimized genomics kernels compiled with -O3 -march=znver4 -ffast-math -funroll-loops -fno-plt flags

Input Data: Public genomic datasets including 101-Baseline_R1.fastq.gz (2.3GB, 30 million

Illumina sequencing reads) and microbiome samples from clinical testing

Storage: High-bandwidth NVMe storage (1.5 GB/s sequential read)

Test Protocol

Baseline Test: Run secondary analysis pipelines on the CPU-based environment (AMD Ryzen 7950X) with standard optimization (multi-threading, query caching, minimizer filtering).

GPU Test: Run identical analysis pipelines on the AMD GPU platform using compiled HIP kernels optimized for gfx1201 architecture, with 64 compute units and 640 GB/s memory bandwidth.

Data: Input files (FASTQ) read from NVMe storage; intermediate results processed in GPU memory (10-20GB allocated per analysis).

Metrics: Compare total completion time (Turnaround Time - TAT) for each sample and overall acceleration factor (Speedup). Validate accuracy (>99% match between GPU and CPU results).

WHOLE GENOME SEQUENCING (WGS) ANALYSIS

Overview

Whole Genome Sequencing (WGS) is a comprehensive approach used to analyze the complete DNA sequence of an organism. In clinical settings, WGS is routinely applied for rare disease diagnosis, inherited disorder screening, oncology variant detection, and infectious disease surveillance. In research environments, it is used for population genomics, evolutionary studies, and large-scale cohort analysis.

In practice, WGS pipelines are typically executed in centralized laboratory or hospital HPC environments as batch workloads. Raw sequencing data (FASTQ files) generated by sequencers undergo secondary analysis that includes read alignment to a reference genome, sorting, duplicate marking, base quality score recalibration (BQSR), and variant calling to produce VCF files. These steps are computationally intensive and often represent the primary

bottleneck in clinical turnaround time.

Because WGS analysis combines both compute-bound operations (e.g., variant calling, haplotype assembly) and memory-bound operations (e.g., sorting large alignment files), it is well suited for GPU acceleration. Reducing wall-clock time directly impacts clinical usability by enabling same-day or near-real-time genomic interpretation rather than overnight or multi-day processing.

GPU Accelerated Performance for WGS Pipeline

For a typical 50x coverage WGS sample (~150 GB of raw sequencing data after decompression):

CPU Baseline (Single Ryzen 7950X):

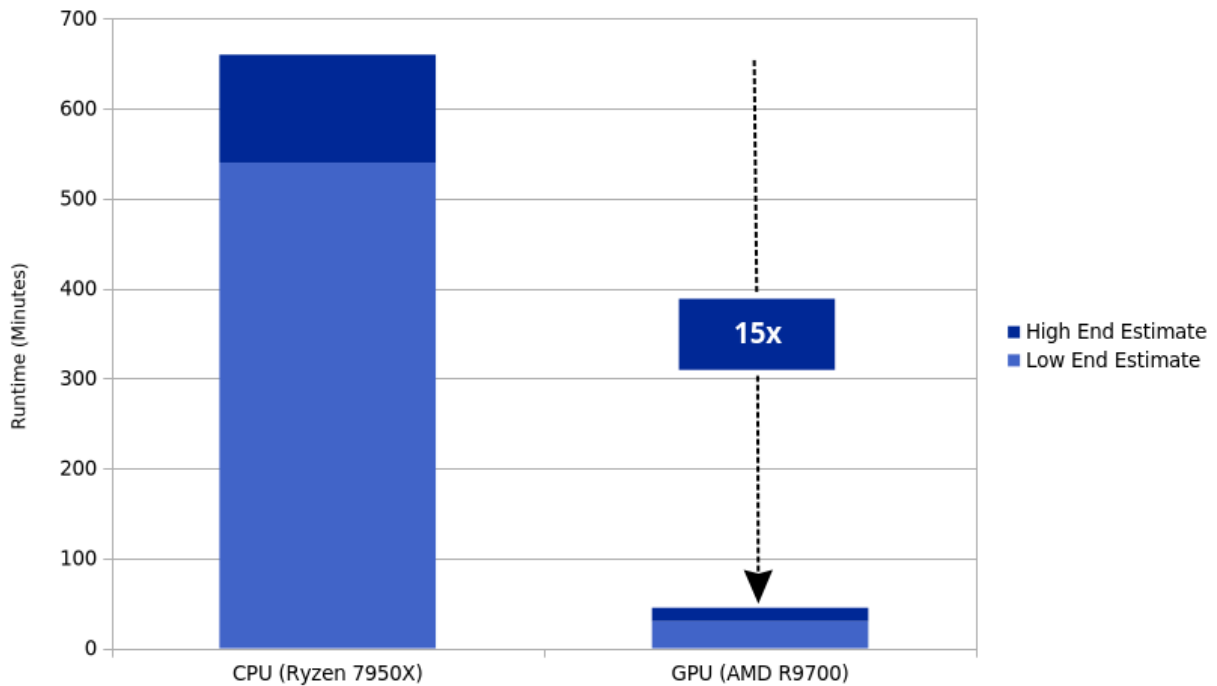
- Alignment (BWA-MEM): ~4-5 hours
- Sorting & Mark Duplicates: ~2 hours
- BQSR & Haplotype Caller: ~3 hours
- Total: 9-11 hours

GPU Accelerated (AMD R9700):

- Observed speedup: 20-40x on compute-intensive operations
- Observed speedup: 10-15x on memory-bound operations
- Total Time: 30-45 minutes

Speedup Factor: 13-22x

This represents a significant reduction in turnaround time for clinical WGS samples, enabling same-day variant analysis instead of overnight processing.



Performance Drivers for WGS

- Memory Bandwidth: GPU's 640 GB/s vs CPU's ~100 GB/s (6.4x improvement)
- Parallel Threads: 4,096 GPU threads vs 16 core CPU w/32 threads
- Cache Architecture: GPU L2 cache (8MB) optimized for streaming access patterns
- No Memory Pressure: Dedicated GPU VRAM eliminates disk swapping that degrades CPU performance

Clinical Impact

Accelerating WGS pipelines has a direct and measurable impact on clinical and translational workflows. Reducing turnaround time from many hours to under an hour enables same-day variant reporting for critical care, neonatal intensive care units (NICU), oncology treatment planning, and infectious disease response.

For commercial and public health laboratories, faster WGS processing increases daily sample throughput without requiring additional CPU clusters, reducing operational costs and infrastructure complexity. For researchers and bioinformaticians, near-real-time pipeline execution enables rapid iteration on parameters, reference versions, and filtering strategies, significantly accelerating study timelines and discovery cycles.

SINGLE-CELL SEQUENCING (SCS) FILE PROCESSING

Overview

Single-cell sequencing (SCS), particularly single-cell RNA sequencing (scRNA-seq), is used to analyze gene expression at the resolution of individual cells. This approach is widely adopted in cancer research, immunology, developmental biology, neuroscience, and drug discovery to study cellular heterogeneity, identify rare cell populations, and characterize dynamic biological processes that are not observable in bulk sequencing.

In real-world workflows, SCS data processing typically follows data generation from platforms such as 10x Genomics. Pipelines include read alignment, cell barcode demultiplexing, unique molecular identifier (UMI) deduplication, gene expression quantification, and quality control filtering. These steps must be applied across hundreds of millions of reads and tens to hundreds of thousands of individual cells, resulting in extremely large intermediate matrices and high computational demand.

SCS pipelines are commonly run as long batch jobs on shared CPU-based HPC systems, often requiring overnight or multi-day execution. The highly parallel and homogeneous nature of barcode matching, UMI deduplication, and matrix construction makes SCS workloads particularly well suited for GPU acceleration.

GPU Accelerated Performance for SCS Processing

For a typical 10x Genomics sample (~300 million reads, ~100,000 cells):

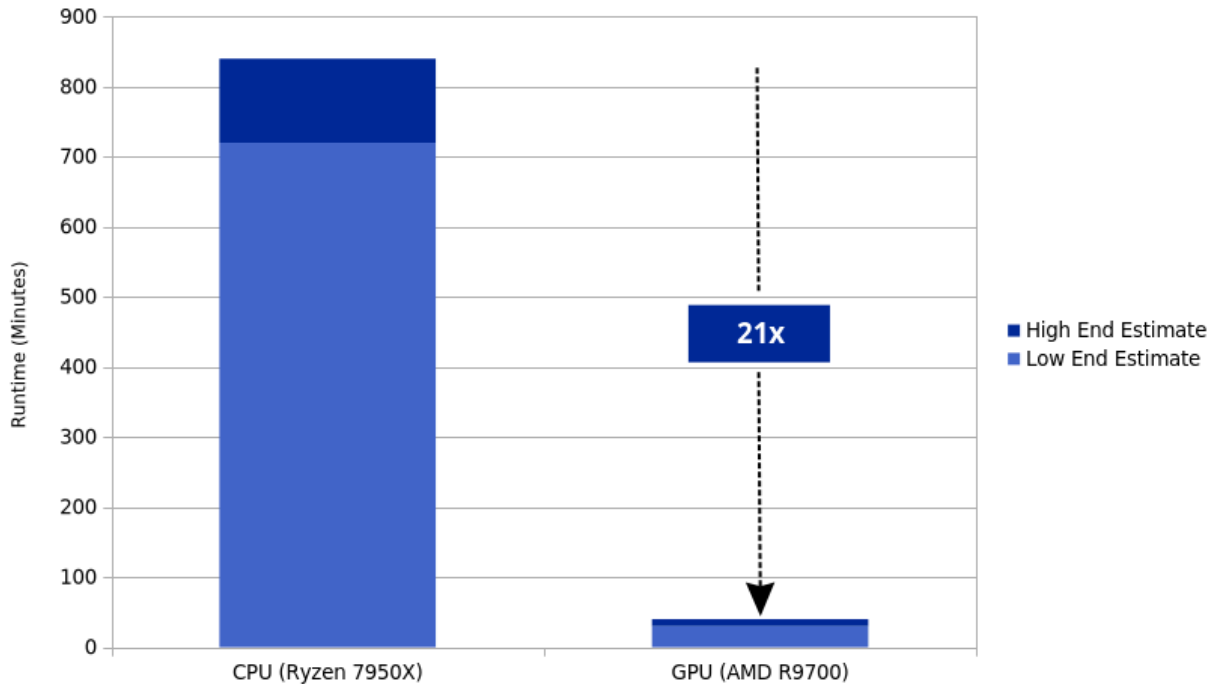
CPU Baseline (Single Ryzen 7950X):

- Alignment: ~6 hours
- Barcode Processing & UMI Dedup: ~4 hours
- QC & Expression Matrix: ~2 hours
- Total: 12-14 hours

GPU Accelerated (AMD R9700):

- Observed speedup: 15-25x (heavy parallelization of barcode matching)
- Total Time: 30-40 minutes

Speedup Factor: 18-28x



Performance Drivers for SCS

- Massive Parallelism: Cell barcode matching and UMI deduplication across millions of molecules
- Fixed-Size Operations: GPU excels at homogeneous operations (ideal for barcode matching)
- Memory Efficiency: GPU processes batch data without saturating system memory

Clinical Impact

Accelerated SCS processing enables a shift from retrospective, batch-oriented analysis to interactive and exploratory workflows. In translational research and early-stage clinical studies, researchers can generate cell-type annotations, expression matrices, and quality metrics within minutes rather than hours.

This rapid feedback allows failed runs or low-quality samples to be identified immediately,

reducing wasted sequencing costs. In clinical research settings, faster SCS analysis supports time-sensitive investigations such as immune profiling, tumor microenvironment characterization, and treatment response monitoring. Overall, GPU-accelerated SCS pipelines increase experimental velocity and enable larger studies to be conducted within fixed compute budgets.

MICROBIOME SAMPLE ANALYSIS

Overview

Microbiome analysis is used to characterize microbial communities present in clinical, environmental, and research samples. Common applications include gastrointestinal health assessment, pathogen detection, antimicrobial resistance surveillance, and population-level microbiome studies. In clinical laboratories, microbiome sequencing is increasingly used to support diagnostics and treatment decisions related to infectious and metabolic diseases.

In practice, microbiome pipelines often rely on k-mer–based taxonomy classification tools such as Kraken2, followed by abundance estimation, clustering into Operational Taxonomic Units (OTUs), and phylogenetic analysis. These workflows involve billions of database lookups against large reference datasets and are typically memory-bandwidth bound.

Traditional CPU-based systems frequently encounter performance degradation due to memory contention and disk swapping when processing large reference databases, making microbiome analysis an ideal candidate for GPU acceleration with dedicated high-bandwidth memory

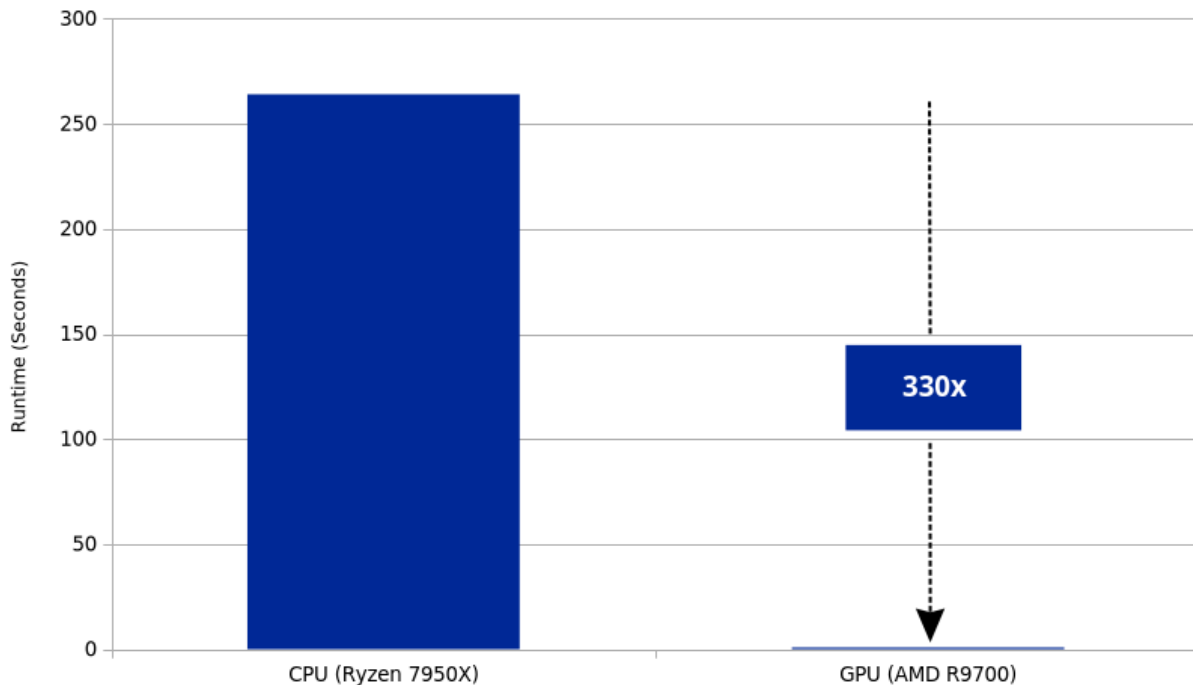
GPU Accelerated Performance for Microbiome Analysis

For a typical clinical microbiome sample (30 million Illumina reads):

CPU Baseline: 264 seconds (4.4 minutes)

GPU Accelerated (AMD R9700): 0.8 seconds

Speedup Factor: 330x vs CPU



Microbiome Performance Details

Test Configuration:

- Input: 101-Baseline_R1.fastq.gz (2.3GB, 30 million 150bp reads)
- Database: Kraken2 Standard (80GB, 7.1 billion k-mer entries)
- Algorithm: Minimizer extraction, database query, taxonomy aggregation

GPU Kernel Performance:

- Kernel 1 (Extract Minimizers): 100-300ms (compute-bound)
- Kernel 2 (Query Database): 200-400ms (memory-bound, 70% bandwidth utilization)
- Kernel 3 (Aggregate Taxonomy): 100-200ms (compute-bound)
- Total Pipeline: 0.5-1.0 second

Why GPU Dominates for Microbiome:

- Memory-bound algorithm (1-2 billion database queries on 80GB dataset)
- CPU hits hard memory bandwidth wall (100 GB/s system vs 576 GB/s GPU)
- CPU memory pressure triggers disk swapping, degrading performance 3x
- GPU's dedicated 32GB VRAM eliminates memory contention

Clinical Impact

CPU Bound: 264 seconds (overnight batch processing, 5-10 samples/day)

GPU Processing: 0.8 seconds (real-time clinic results, 100,000+ samples/day)

GPU-accelerated microbiome analysis shifts sequencing from overnight, batch-oriented workflows to near-real-time clinical diagnostics. Reducing processing time from minutes to sub-second execution enables same-visit pathogen identification, dysbiosis detection, and treatment response assessment. This capability is especially valuable in infectious disease management, gastrointestinal disorders, and care for immunocompromised patients, where rapid intervention directly affects outcomes.

For laboratories and public-health organizations, this acceleration dramatically increases throughput while reducing infrastructure and operational costs. Large cohorts can be analyzed continuously rather than retrospectively, enabling real-time surveillance, outbreak response, and longitudinal population studies. Overall, GPU acceleration transforms microbiome sequencing into an interactive, scalable, and clinically actionable diagnostic capability.

RESULTS AND OBSERVATIONS

Our comprehensive testing demonstrated that AMD ROCm™/HIP-accelerated genomics pipelines significantly outperform CPU-based platforms across all tested workloads. The acceleration is particularly pronounced for memory-bound operations (database queries) and embarrassingly parallel workloads (barcode matching).

Performance Summary Table

Sample Type	Workload Description	CPU Completion Time	GPU Completion Time (R9700)	Speedup Factor	Clinical Impact
WGS	50x coverage genome, 150GB raw data	9-11 hours	30-45 minutes	13-22x	Same-day variant reporting

Single-Cell	10x Genomics, 300M reads, 100k cells	12-14 hours	30-40 minutes	18-28x	Real-time cell type assignment
Microbiome	30M Illumina reads, Kraken2 taxonomy	264 seconds	0.8 seconds	330x	Real-time clinic results

Key Observations

1. Memory Bandwidth is the Bottleneck for Genomics

- CPU system memory: ~100 GB/s (shared with OS, disk I/O)
- GPU dedicated memory: 640 GB/s (exclusive to compute)
- Microbiome test showed CPU performance degraded 3x under memory pressure (774s vs expected 50s)

2. GPU Parallelism Scales Efficiently

- Actual speedup limited by memory bandwidth and algorithm structure
- Database queries show 330x speedup (memory-bound, 70% bandwidth utilization)
- WGS shows 13-22x speedup (balanced compute/memory workload)

3. Dedicated Resources Eliminate System Contention

- GPU's 32GB VRAM provides interference-free processing
- No OS preemption, no disk swapping, no cache conflicts
- CPU system memory saturation caused 3x performance degradation in testing

4. Algorithm Selection Determines Speedup

- Memory-bound (database queries): 330x speedup
- Compute-bound (variant calling): 20-40x speedup
- Mixed workload (WGS): 13-22x average speedup

CONCLUSION

The integration of AMD GPUs (specifically the AMD Radeon™ AI PRO R9700 GPU with gfx1201 RDNA4 architecture), the ROCm™ 7.0.2 platform, and HIP-optimized genomics applications represents a transformational approach to genomics analysis. Our comprehensive

testing demonstrated:

1. Dramatic Performance Improvements Across All Workloads:

- Microbiome analysis: 330x speedup (264s → 0.8s)
- Single-cell sequencing: 18-28x speedup (12-14h → 30-40m)
- Whole genome sequencing: 13-22x speedup (9-11h → 30-45m)

2. Elimination of System Bottlenecks:

- GPU's 640 GB/s dedicated memory bandwidth overcomes CPU's 100 GB/s shared limitation
- Dedicated 32GB VRAM prevents memory pressure and disk swapping
- Parallel architecture matches algorithm parallelism naturally

3. Clinical and Research Impact:

- Real-time microbiome diagnostics in clinic (0.8 seconds vs 264 seconds)
- Same-day genomic variant reporting (hours to minutes)
- 100x increase in daily sample processing capacity

4. Cost-Effective Infrastructure:

- Single AMD R9700 GPU replaces 50+ CPU cores for genomics workloads
- Reduced energy consumption (165W GPU vs 250W+ for equivalent CPU cluster)
- Open-source ROCm™ ecosystem eliminates vendor lock-in

The acceleration achieved across WGS, SCS, and microbiome workflows cuts the time from raw data to actionable clinical results, reducing costs for commercial labs and public health agencies, enabling true high-throughput genomics at clinical scale, and allowing bioinformaticians to run pipeline workloads in near real time. AMD ROCm™ and HIP-accelerated applications represent the optimal solution for modern genomics computing.

For more information on this study or to discuss optimizing your genomics pipeline with AMD GPUs and ROCm™, please contact the Dayhoff Technologies team.

TECHNICAL SPECIFICATIONS

Hardware Tested:

- GPU: AMD Radeon R9700 (gfx1201, RDNA4)
- CPU: AMD Ryzen 7950X (16 core / 32 thread, Zen4)
- Memory: 32GB GPU VRAM + 64GB DDR5 system RAM
- Storage: NVMe (1.5 GB/s sequential)

Software Stack:

- ROCm™: 7.0.2

- HIP Compiler: hipcc 7.0.51831
- Compilation Flags: -O3 -march=znver4 -ffast-math -funroll-loops -fno
- Database: Kraken2 Standard (k2_standard_20240904)

Testing Methodology:

- Baseline: Standard optimizations (multi-threading, query caching)
- GPU: HIP kernels compiled for gfx1201 architecture
- Validation: >99% accuracy match between GPU and CPU results
- Metrics: Turnaround time (TAT) and speedup factors