Comprehensive Review of TOTHOST's High-Performance GPU Cloud VM

A Deep Dive into System Provisioning, Environment Setup, and Performance Benchmarking

Do Phuc Hao, Ph.D. Candidate Dinh Truong Duy, Ph.D.

November 13, 2025

Abstract

This report provides a detailed, hands-on evaluation of an upgraded high-performance cloud Virtual Machine (VM) provided by TOTHOST. The review chronicles the end-to-end user journey, from initial system verification and provisioning to the complex process of setting up a stable AI development environment and conducting rigorous performance benchmarks. We focus on the core infrastructure components critical for AI/ML workloads: the NVIDIA RTX 5880 Ada Generation GPU, CPU, RAM, storage I/O, and network connectivity. The findings document not only the platform's robust performance but also the real-world technical challenges encountered and their solutions, offering valuable insights for both prospective users and the TOTHOST engineering team.

Contents

| 1 Introduction | | roduction | 2 | | |
|--|-----|---|---|--|--|
| 2 | Par | t 1: System Provisioning and Environment Setup | 2 | | |
| | 2.1 | Storage Volume Expansion: A SysAdmin Deep Dive | 2 | | |
| | 2.2 | AI Environment Setup: A Tale of Two Package Managers | 2 | | |
| 3 | Par | Part 2: Infrastructure Performance Benchmarks | | | |
| | 3.1 | GPU Compute Performance | 3 | | |
| | | 3.1.1 Synthetic Benchmark: Matrix Multiplication | 3 | | |
| | | 3.1.2 Real-World Benchmark: Stable Diffusion Image Generation | 3 | | |
| | 3.2 | | 4 | | |
| | | 3.2.1 Storage I/O Speed | 4 | | |
| | | 3.2.2 Internet Connectivity | 4 | | |
| 4 Part 3: Conclusion and Recommendations | | t 3: Conclusion and Recommendations | 5 | | |
| | 4.1 | Summary of Findings | 5 | | |
| | 4.2 | Strengths | 5 | | |
| | 4.3 | Constructive Feedback and Recommendations | Ę | | |
| | 4 4 | Final Verdict | 6 | | |

1 Introduction

The objective of this review is to conduct a comprehensive assessment of a TOTHOST cloud VM following a significant resource upgrade. This report details the practical steps and findings from a researcher's perspective, aiming to simulate a real-world project setup. All tests were performed on a single VM with the final, verified specifications:

• GPU: 1x NVIDIA RTX 5880 Ada Generation (48GB GDDR6 VRAM)

• CPU: 8 vCPUs (Intel Xeon Gold 6330N @ 2.20GHz)

• **RAM**: 16 GB

• Storage: 200 GB (Expanded from 48GB)

• Virtualization: VMware

This report is structured into three main parts: an analysis of the system provisioning and environment setup journey, a quantitative evaluation of the core infrastructure's performance, and a concluding summary with constructive recommendations.

2 Part 1: System Provisioning and Environment Setup

This section details the user experience from receiving the upgraded VM to achieving a stable, functional AI development environment.

2.1 Storage Volume Expansion: A SysAdmin Deep Dive

The initial state of the upgraded VM presented a common yet critical system administration task. While the underlying virtual disk was expanded to 200GB by TOTHOST, the operating system's root partition was still only utilizing 48GB.

Challenge: The additional $\sim 150 \mathrm{GB}$ of provisioned space was unallocated and unusable by the OS.

Solution: We successfully reclaimed the unallocated space by performing a standard Linux LVM (Logical Volume Manager) expansion procedure. This multi-step process involved using fdisk, pvcreate, vgextend, lvextend, and resize2fs.

Outcome: The root filesystem was successfully and safely expanded to its full 200GB capacity. This demonstrates the flexibility of the VM environment, which grants users the root-level control necessary for advanced system configuration.

2.2 AI Environment Setup: A Tale of Two Package Managers

With the infrastructure correctly provisioned, the next step was to set up a robust and isolated Python environment using Conda.

Challenge: Initial PyTorch Installation Failure

Our first attempt to install PyTorch using the standard Conda command resulted in a persistent runtime error:

1 ImportError: ... undefined symbol: iJIT_NotifyEvent

Listing 1: Recurring ImportError with Conda installation.

This error, related to the Intel Math Kernel Library (MKL), pointed to a deep-seated dependency conflict between the official Conda channels, a non-trivial issue in AI environment setup.

Solution: Pivoting to a Hybrid Pip-based Installation

To resolve this, we adopted a more robust strategy: using Conda to manage the environment but leveraging **Pip** to install PyTorch directly from its official wheel repository.

Outcome: This approach was 100% successful. PyTorch was installed correctly and functioned without any errors, highlighting a key best practice for setting up complex AI environments.

3 Part 2: Infrastructure Performance Benchmarks

This section presents the quantitative results from our benchmarks.

3.1 GPU Compute Performance

We conducted two tests to measure GPU performance: a synthetic benchmark for raw power and a real-world benchmark for practical application speed.

3.1.1 Synthetic Benchmark: Matrix Multiplication

This test measures the raw computational power using a large matrix multiplication (20000x20000) in FP32 precision.

Table 1: GPU FP32 Matrix Multiplication Benchmark Results.

| Metric | Result |
|---------------------------------|----------------|
| Matrix Dimensions | 20000 x 20000 |
| Average Time per Multiplication | 0.6409 seconds |

An average time of just **0.64 seconds** demonstrates excellent raw compute capability.

3.1.2 Real-World Benchmark: Stable Diffusion Image Generation

This test measures the end-to-end time to generate a 512x512 image from a text prompt using the popular Stable Diffusion v1.5 model.

| Table 2: Stable Diffusion | n v1.5 Benchmark Results. |
|---------------------------|--------------------------------|
| Metric | Result |
| Model | runwayml/stable-diffusion-v1-5 |
| Precision | FP16 |
| Average Time per Image | $1.85 \mathrm{seconds}$ |

Generating a high-quality image in under 2 seconds is an **extremely impressive** result. It confirms the GPU's powerful performance in a practical, real-world generative AI task. During the test, monitoring showed the GPU reaching a high-performance state (P2), drawing up to 243W of power.

3.2 Storage and Network Performance

3.2.1 Storage I/O Speed

The speed of the expanded storage was measured using the dd utility with a 10GB test file.

Table 3: Storage I/O Benchmark Results.

| Operation | Speed |
|-------------|----------|
| Write Speed | 121 MB/s |
| Read Speed | 107 MB/s |

The I/O speeds are consistent and adequate for tasks like saving model checkpoints and handling moderately sized datasets.

3.2.2 Internet Connectivity

Network speed was tested using the speedtest-cli tool against a local server in Hanoi.

Table 4: Network Speed Benchmark Results.

| Operation | Speed |
|----------------|---------------|
| Download Speed | 291.27 Mbit/s |
| Upload Speed | 297.72 Mbit/s |

The network performance is a standout feature. A symmetric speed of nearly 300 Mbps is very strong and highly beneficial for both downloading large datasets and uploading results or models.

4 Part 3: Conclusion and Recommendations

This comprehensive review confirms that the upgraded TOTHOST High-Performance GPU Cloud VM is a powerful and capable platform for demanding AI/ML research and development.

4.1 Summary of Findings

- Infrastructure Performance: The core hardware delivers excellent performance. The NVIDIA RTX 5880 GPU is the star component, providing massive computational power, demonstrated by its sub-2-second image generation time with Stable Diffusion. Network connectivity is also a strong point with symmetric 300 Mbps speeds.
- System Provisioning: The platform grants the necessary root access for essential system administration tasks like LVM expansion, offering great flexibility for advanced users.
- Developer Experience: The initial setup process presents real-world challenges. The need to manually expand storage and navigate complex software dependency issues (Conda vs. Pip for PyTorch) are significant steps in the setup journey. Overcoming these challenges leads to a highly stable and powerful final environment.

4.2 Strengths

- Top-Tier GPU Hardware: Access to the NVIDIA RTX 5880 Ada Generation with 48GB VRAM is a significant competitive advantage, enabling research on large models.
- Excellent Computational and Network Performance: The benchmarks confirm that the hardware's full power is available to the user, complemented by a very strong network connection. The real-world performance in generative AI tasks is particularly impressive.
- Flexibility and Control: Providing a true VM with root access allows researchers to fully customize their environment, a crucial feature for R&D.

4.3 Constructive Feedback and Recommendations

- Streamline Storage Provisioning: Consider offering an option at VM creation to automatically expand the filesystem to the full size of the provisioned disk. This would eliminate the manual LVM steps and significantly improve the out-of-the-box experience.
- Provide Pre-configured AI-Ready Images: To further enhance the developer experience, TOTHOST could offer official VM images that come pre-installed with:
 - 1. An activated Miniconda environment.
 - 2. A tested, stable installation of the latest NVIDIA drivers and PyTorch/TensorFlow (installed via the more reliable Pip method we discovered).

This would reduce setup time from hours to minutes and bypass common dependency issues.

• Publish Technical Guides: The challenges we solved (LVM expansion, Conda/Pip for PyTorch) are excellent topics for technical blog posts or tutorials on the TOTHOST website. This content would be invaluable to other users and would position TOTHOST as a knowledgeable and supportive provider for the AI community.

4.4 Final Verdict

The TOTHOST High-Performance GPU Cloud VM is a formidable platform for AI research. While the initial setup requires a degree of system administration expertise, the end result is a flexible, powerful, and unconstrained development environment. The performance of the core components, especially the NVIDIA RTX 5880 GPU and the network infrastructure, is exceptional.

With some refinements to the initial user onboarding process, such as offering preconfigured and fully expanded images, TOTHOST has the potential to become a leading choice for AI researchers and developers in Vietnam who require maximum performance and control. This platform is a powerful enabler for the local AI ecosystem.