



A Purpose-Built

Gen AI

Accelerator

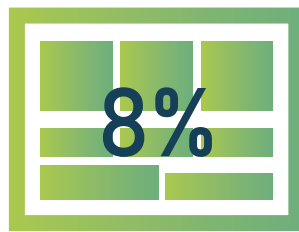


Gen AI N3000 Accelerator Designed for AI Inference

The N3000 AI Accelerator leads the industry in performance and efficiency. Purpose-built, the Gen AI Accelerator improves inferring accuracy, optimizes bandwidth, and increases resource utilization - all while delivering a significantly lower TCO than existing solutions.



Work Group



Host CPU



Offload



N3000 Accelerator

CPU + ASIC Combo

The Perfect Blend of Power & Efficiency

Combining CPU with our ASIC technology provides a more balanced and efficient solution compared to existing solution . N3000 PCIe, supports up to 256B parameters, is the one of the few that has using 7nm ASIC with LPDDR5 and PCIe Gen5 spec at max. 55W TDP.



Optimized for Max. Returns

Maximize your returns with a product that combines low power consumption and high accuracy, offering an attractive ROI. Our ASIC technology is incredibly energy-efficient, slashing operational costs and reducing your carbon footprint.



Precision Meets Accuracy

Achieve unparalleled accuracy with minimal perplexity while enjoying rapid inferring speeds, while keeping power consumption in check, ensuring your business reaps high profits.



Zero Resource Wastage

With ASICs, every bit of processing power is dedicated to your tasks, making your operations more efficient. ASICs are purpose-built for specific tasks, ensuring optimal performance tailored to your needs.

| Usage | Format | #bit | S | Exponent (8bit) | Mantissa (23bit) |
|-------------|---------------|------|---|-----------------|------------------|
| Training | FP32 | 32 | | | |
| | TF32 | 19 | | | |
| | FP16 | 16 | | | |
| | BF16 | 16 | | | |
| Inferencing | Neuchips FFP8 | FP8 | 8 | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

- Hardware convert to BF16 for computing in N3000.
- Flexible 8-bit Floating-Point Format (FFP8) reduce memory transactions & delivers highest inference

Proprietary Patented FFP8

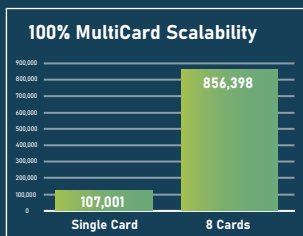
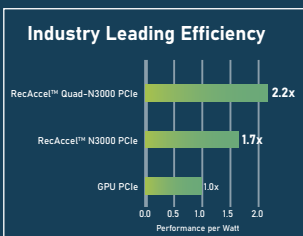
The Ultimate Memory Fix

Solve LLM Memory Bound Challenge

N3000 PCIe Feature Highlights

- Easy-to-use Software Stack
 - Frictionless transition from existing development kit
- Patented Flexible Float Point 8
 - Improves inferencing accuracy
 - Quantizes data for memory capacity and bandwidth optimization in LLM workload
- 100% Linear Scaling to Multi-cards
- 4x 64bit LPDDR5 (6.4GHz) with ECC
- Up to 128GB on Card LPDDR5
- PCIe Gen5

Industry Leading Performance



Product Specification

| | |
|----------------------------|---|
| BFLOAT16 | 32 TFLOPS |
| INT8 | 206 TOPS |
| Memory | 32GB LPDDR5 |
| Memory Bandwidth | 200 GB/s |
| Thermal Design Power (TDP) | 55W |
| Form Factor | Full-height, Full-length (FHFL) 10.5" Dual-slot (266mm/10.5 inch) |
| PCI Express Interface | PCI Express 5.0 x 16 Lane and Polarity Reversal Supported |
| Interconnect | PCIe Gen5: 64 GB/s |
| Server Options | Partner and NEUCHIPS-certified Systems with 1-8 Accelerator |

Environmental Specification

| | |
|-------------------------------|--------------------------|
| Ambient Operating Temperature | 0°C~50°C |
| Storage Temperature | -40°C~75°C |
| Operating Humidity | 5%~85% Relative Humidity |
| Storage Humidity | 5%~95% Relative Humidity |



MLPerf™ Proven AI Solution

Applications

Recommender, Gen AI...

AI Models

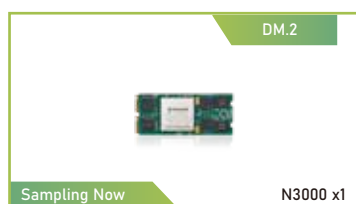
DLRM | LLaMA | GPT...

Frameworks

ONNX PyTorch C++

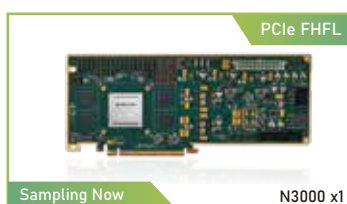
SDK

RecAccel™ SDK



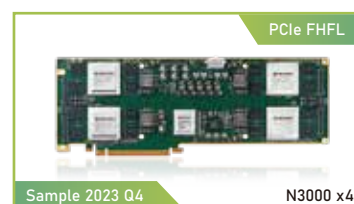
Sampling Now N3000 x1

N3000 DM2



Sampling Now N3000 x1

N3000 PCIe



Sample 2023 Q4 N3000 x4

Quad-N3000 PCIe

| | | | |
|--|-----------------------------------|----------------------------------|------------------------------------|
| DLRM Performance | 20M inf/sec | 30M inf/sec | 120M inf/sec |
| Llama-2 7B ¹ Output Performance (Single Card) | 40 tokens/sec | 100 tokens/sec | 400 tokens/sec |
| Llama-2 7B ¹ Output Performance (Single Card) | 480 tokens/sec (YV 2.5: 12 Cards) | 800 tokens/sec (Server: 8 Cards) | 3,200 tokens/sec (Server: 8 Cards) |
| TDP | 25 W | 55 W | 300 W |
| SRAM | 160 MB | 160 MB | 640 MB |
| LPDDR 5 Capacity | 32 GB | 32 GB | 256 GB |
| LPDDR Bandwidth | 200 GBps | 200 GBps | 800 GBps |

¹ Llama-2 7B performance is based on batch size 16, content length 2k.