# Cerebras vs SambaNova vs Groq: AI Chip Comparison (2025)

By Adrien Laurent, CEO at IntuitionLabs • 10/23/2025 • 40 min read

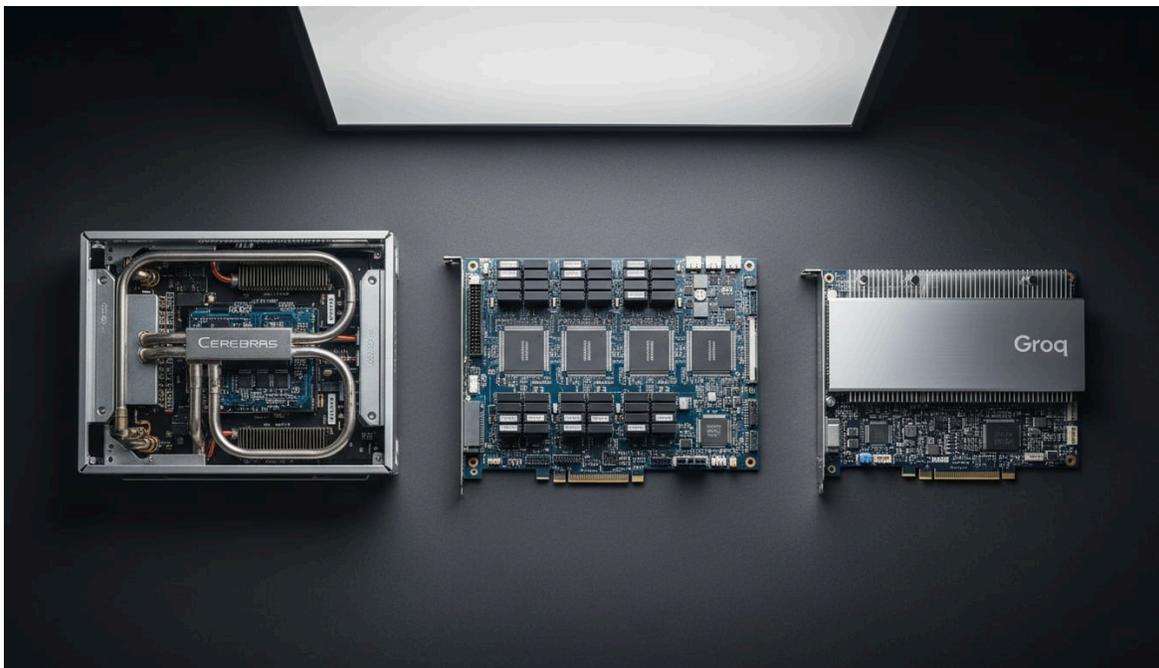cerebras    sambanova    groq    ai hardware    ai accelerators    wafer-scale engine

language processing unit    ai inference    ai chip comparison

# Executive Summary

This report provides a detailed comparative analysis of three emerging AI hardware companies – **Cerebras Systems**, **SambaNova Systems**, and **Groq** – as of October 2025. These companies each offer specialized AI accelerators targeting high-performance training and inference workloads, in competition with incumbent GPU solutions (notably NVIDIA's GPUs). Cerebras pioneered wafer-scale processors to maximize parallelism, SambaNova developed a reconfigurable dataflow architecture (RDUs) in integrated systems, and Groq designed a novel streaming **Language Processing Unit (LPU)** focused primarily on inference. As of 2025, all three firms have achieved multi-billion-dollar valuations and secured major funding: Cerebras raised **$1.1 billion** at an **$8.1 billion** valuation ([1] www.reuters.com), Groq raised **$750 million** pushing its valuation to **$6.9 billion** ([2] www.reuters.com), and SambaNova raised hundreds of millions (e.g. $676 M Series D in 2021 with a $5.1 B valuation) ([3] www.anandtech.com) ([4] www.anandtech.com). Each company has proprietary chip designs, software stacks, and deployment targets:

- **Cerebras Systems** (founded 2016, Menlo Park CA) remains focused on **training of very large AI models**. Its wafer-scale chips (Wafer-Scale Engine, WSE) integrate trillions of transistors on a single monolithic die, enabling massive on-chip compute and memory bandwidth. The third-generation **WSE-3** (announced March 2024) can train models *"ten times larger than OpenAI's GPT-4"* ([5] time.com), and is at the heart of supercomputers like the *Condor Galaxy 3*. In late 2025 Cerebras is expanding into new markets (e.g. UAE data centers) ([6] www.reuters.com) and won a DARPA contract to link its chips via photonic routers ([7] www.reuters.com). The company has invested in scaling production (using TSMC's 3nm process for WSE-3) and has postponed an IPO after its late-2025 funding.

- **SambaNova Systems** (founded 2017, Palo Alto CA) pioneered a **Reconfigurable Dataflow Architecture (RDA)** implemented in its RDU (Reconfigurable Dataflow Unit) chips and integrated **DataScale®** systems. SambaNova's architecture uses arrays of adaptable compute/data units (RDUs) that can be temporally and spatially configured by the compiler for different neural network layers ([8] www.anandtech.com). This design emphasizes on-chip memory and dataflow to efficiently execute both training and inference. By 2024 SambaNova publicly introduced **Samba-1**, a 1-trillion-parameter open-source language model aimed at enterprise/custom use ([9] time.com). High-profile deployments include U.S. national labs (e.g. Los Alamos, LLNL) for large-scale AI research ([10] www.anandtech.com), and its software stack (SambaFlow/SambaStudio) supports common AI frameworks. SambaNova's systems offer large local memory (e.g. **3 TB per node** in its SN30 generation) ([11] www.nextplatform.com) to handle foundation models. The company continues to iterate on its chips (e.g. doubling compute and memory compared to predecessors) and targets sectors from defense to cloud enterprise.

- **Groq** (founded 2016, Mountain View CA) is notable for its *single-core, deterministic* LPU architecture. Groq's **Language Processing Units** implement an in-order, **token-streaming** design where every clock cycle executes useful work without hardware context switching (Groq calls this "kernel-free" ([12] www.eetimes.com)). This yields predictably low-latency AI inference performance. The founder, ex-Google TPU architect Jonathan Ross, has touted LPUs as offering *"ten times faster and ten times lower"* cost inference than GPUs ([13] time.com). Groq's chips (fabricated on similar nodes) focus exclusively on inference acceleration; the company has attracted major investment (Samsung, Cisco) and announced large deals (e.g. a $1.5 B Saudi Arabia commitment for AI infrastructure) ([14] www.reuters.com). Groq is expanding data center deployments (e.g. a European AI-focused facility in Finland) to capture the growing inference market ([15] www.tomshardware.com).

This report delves into each company's **history**, **technology and architecture**, **performance claims**, **market positioning**, and **use-case deployments**, comparing them side-by-side wherever possible. It summarizes current metrics (processors' scale, memory, benchmarks where available) and funding/valuation data. Case studies include national lab supercomputing projects and large foundation-model training systems. Finally, we discuss industry trends and future directions – including how these specialized AI accelerators might fit into evolving AI/data-center ecosystems, partnerships, and potential regulation or market shifts.

# Introduction

The advent of large-scale machine learning (especially deep neural networks and generative AI) has driven demand for specialized high-performance hardware. While GPUs (primarily from NVIDIA) have dominated AI training and inference over the past decade, a new wave of startups is producing **dedicated AI accelerators** that follow alternative design philosophies. These aim to overcome GPU bottlenecks (memory bandwidth, scaling, power) by customizing the chip architecture and system design for AI workloads.

As of 2025, the AI-chip market is growing rapidly. Gartner projected that global revenue for AI accelerators would **more than double from 2023 to 2027** ([15] www.tomshardware.com) ([1] www.reuters.com). Broadcom's CEO has forecast a **$60–90 billion** AI revenue opportunity by 2027 ([16] www.reuters.com). Nvidia remains the market leader – one analysis estimated its AI-related sales could approach **$400 billion by 2028** ([17] www.techradar.com) – but investors are pouring capital into GPU "challengers." In this context, **Cerebras**, **SambaNova**, and **Groq** have emerged as high-profile contenders, collectively raising billions of dollars (between them, on the order of $3–4 B of external funding) and commanding multi-billion valuations ([1] www.reuters.com)([2] www.reuters.com) ([4] www.anandtech.com). Their technologies are already deployed in critical projects – from U.S. Department of Energy (DOE) supercomputing centers to Middle East AI data centers – and they herald new architectures for AI compute.

This report compares these three companies side-by-side, covering:

- **Company Background & Funding**: Origins, leadership, milestones (product announcements, partnerships), and financial status (recent funding rounds and valuations).
- **Hardware Architecture**: Core design of the processors (Cerebras' wafer-scale engines, SambaNova's reconfigurable dataflow units, Groq's LPUs), including chip organization (e.g. WSE's mesh network and MemoryX block, SambaNova's RDU tiles and interconnect, Groq's linear pipeline of functional units). Discussion of process nodes and scale (transistor count, die size) where data is available.
- **System & Software Stacks**: Packaging into systems (e.g. Cerebras CS-2/C5 systems, SambaNova DataScale racks, GroqCard/GroqRack), and the software ecosystem (compilers, frameworks, model support).
- **Performance and Benchmarks**: Published or leaked performance comparisons for typical AI workloads (e.g. training large transformer models or running inference), including metrics like throughput (tokens/sec, FLOPS), memory capacity, and energy efficiency. We include wherever possible independent benchmarks or claims from reputed sources ([5] time.com) ([13] time.com).
- **Use Cases / Case Studies**: Real-world deployments (e.g. DARPA contract for Cerebras interconnect ([7] www.reuters.com), DOE lab supercomputing, European AI data center, etc.) that illustrate how customers are using each platform.
- **Market Position & Trends**: How each positions itself relative to GPUs and each other, strategic partnerships (e.g. Cerebras–G42 in UAE ([6] www.reuters.com), Groq–Samsung, Cisco ([18] www.reuters.com)), and the competitive landscape of AI infrastructure. Discussion of investor sentiment and IPO prospects (e.g. Cerebras' U.S. listing plans ([1] www.reuters.com)).
- **Future Directions**: Potential evolutionary paths (e.g. next-gen chips, software improvements, M&A), and considerations such as open-source RISC-V movement, national AI initiatives, and how the inference vs. training market split might evolve.

By bringing numerous technical details, data points, and cited opinions together, this report aims to be a comprehensive resource on the state of these three AI chip firms as of late 2025.

# 1. Company Background and Funding

## 1.1 Cerebras Systems

Cerebras Systems was founded in 2016 by Andrew Feldman and other veterans from AI startups and supercomputing. The company's defining idea is to build extremely large "wafer-scale" chips that break the mold of traditional reticle-limited dies ([5] time.com). After stealth development, Cerebras unveiled its **first Wafer-Scale Engine (WSE)** in 2019, a single-die chip covering most of a 300 mm silicon wafer (approximately 46,225 mm²) ([5] time.com). This chip contained ~1.2 trillion transistors at 7 nm and thousands of AI-optimized cores with massive on-chip SRAM. The goal was to create one processor with as much memory and interconnect on-chip as 1,000+ traditional GPUs, eliminating inter-GPU communication bottlenecks.

**Funding and Valuation:** Cerebras raised multiple venture rounds through 2023. Its Series G in 2023 was $275 M at a reported $8.1 B valuation. In late 2025, Reuters reported that Cerebras **raised $1.1 billion** in new capital (led by Fidelity and Atreides) ([1] www.reuters.com), bringing the valuation again to **$8.1 billion**. Notably, 2025 also saw investment from 1789 Capital, a Trump-linked firm ([1] www.reuters.com). Shortly after this funding, Cerebras withdrew plans for a U.S. IPO (originally soliciting a large offering) ([19] www.reuters.com). The company explained the IPO delay as due to regulatory reviews (including a U.S. national-security review of an earlier $335M G42 investment) and instead chose to expand private funding ([19] www.reuters.com) ([1] www.reuters.com).

**Customers and Partnerships:** Cerebras markets its CS-2 and CS-3 systems to AI research labs and enterprises. Its first key customers included Argonne and Oak Ridge National Labs (for exascale AI) and G42 (an Abu Dhabi AI firm) ([6] www.reuters.com). In 2025, Cerebras announced plans to supply infrastructure for the *Stargate UAE* AI data center campus, indicating expansion in the Middle East ([6] www.reuters.com). The company also won a DARPA contract for an advanced AI supercomputing effort (called "Fuse" project) valued at $45 M ([7] www.reuters.com), partnering its WSE chips with Ranovus's photonic interconnect to build a system "150 times faster" than conventional ones. These deals highlight Cerebras's positioning in large-scale government and research projects.

## 1.2 SambaNova Systems

Founded in 2017 by former Sun Microsystems executives Rodrigo Liang (CEO) and Kunle Olukotun (CTO), SambaNova built a hybrid hardware-software stack centered on its **Reconfigurable Dataflow Unit (RDU)** architecture. SambaNova Systems kept a low profile initially, raising venture funding quietly. By late 2020, it had attracted ~$450 M in total funding ([4] www.anandtech.com), from investors like Google Ventures, Intel Capital, and BlackRock ([10] www.anandtech.com). In April 2021, SambaNova announced a Series D round of **$676 M at a $5.1 B valuation** ([4] www.anandtech.com), briefly making it one of the most richly backed AI startups.

**Architecture & Products:** SambaNova's value proposition is a turnkey AI platform (hardware + software). Its DataScale systems consist of boards called *DataScale SN* (such as SN10, SN15, SN30) that contain multiple RDU accelerator chips and associated memory. A "SN10-8R" configuration (8 sockets each with an RDU board and connected memory) was launched in 2020 ([10] www.anandtech.com), demonstrating their first product for limited customers. The RDU chip itself (internally called "Cardinal") was built on TSMC 7 nm, integrating thousands of small processing elements that can be *reconfigured in software* to form the dataflow graph of an AI model ([8] www.anandtech.com). The compiler (SambaFlow) maps a user's neural network graph onto a sequence of configurations of the RDU array, aiming to maximize data reuse and minimize data movement.

**Early Deployments:** Even before public product launches, SambaNova secretly deployed its hardware at U.S. national labs. AnandTech reported that SambaNova had "semi-hushhush deployments" at Lawrence Livermore and Los Alamos National Labs ([10] www.anandtech.com). This early adoption suggests the company targeted high-end research workloads from the outset, alongside their enterprise/government customers. Other customers reportedly include Accenture and, reportedly, SoftBank (as noted in media profiles ([9] time.com)).

**Recent Developments:** In 2024, SambaNova publicly unveiled **Samba-1**, its in-house trained 1-trillion-parameter generative AI model (an open foundation model) ([9] time.com). Samba-1 is intended for on-premises enterprise AI applications and fine-tuning. Additionally, SambaNova launched SaaS/cloud services to run foundation models (e.g. Llama 3) on its hardware, emphasizing performance. The company announced significant hardware upgrades: its new generation doubles the compute and memory of its predecessor. For example, one public statement noted that each socket in the next DataScale rack offers 3 TB of memory (twice the prior 1.5 TB) ([11] www.nextplatform.com). As of 2025, SambaNova continues to market to enterprise, government, and scientific clients, positioning itself as an alternative to GPU-based AI clouds. Financially, no new large funding rounds have been reported after 2021, but the company remains private.

## 1.3 Groq

Groq, founded in 2016 by Jonathan Ross (a lead architect of Google's TPU), emerged with a radical design philosophy. It produces a single-core **Language Processing Unit (LPU)** intended for extremely fast, predictable inference of large language models and other AI workloads ([13] time.com). Groq's early funding was relatively modest (series rounds in 2019-2020), but its profile rose with dramatic claims about speed and throughput.

**LPU Architecture:** Groq's LPU architecture (originally called the Tensor Streaming Processor) feeds AI model tokens through a single, wide pipeline of functional units, executing all operations in lock-step with no kernel switching. This "data-stream" or "tensor-streaming" design means every clock cycle is doing useful work. Groq claims it achieves *"speeds 10x faster and costs 10x lower"* than GPUs for certain inference tasks ([13] time.com), at the expense of flexibility (LPUs lack dynamic branching or wide programmability). Groq also emphasizes deterministic, low-latency compute – a selling point for real-time applications where unpredictable latency is unacceptable ([20] time.com) ([21] www.tomshardware.com).

**Funding and Valuation:** Groq quietly raised capital until a wave of AI chip interest in 2024–2025. In August 2024 it closed a $640 M Series D at a reported $2.8 B valuation ([22] www.reuters.com). Then in September 2025, Reuters reported Groq raised **$750 M** (led by Disruptive) at a **$6.9 B** post-money valuation ([2] www.reuters.com) ([23] www.reuters.com) – more than doubling its valuation in a year. Key investors include Samsung, Cisco, 1789 Capital, altimeter, and others ([24] www.reuters.com) ([25] www.reuters.com). In 2025 Groq also secured a $1.5 B commitment from Saudi Arabia to deploy LPUs, expected to generate ~$500 M revenue in one year ([14] www.reuters.com).

**Deployments:** Unlike the others, Groq has focused heavily on building out data-center capacity to deliver inference-as-a-service. In mid-2025, Groq opened an EU data center in Helsinki (with Equinix) targeting European AI workloads ([15] www.tomshardware.com). Domestically, they operate "GroqCloud" for customers. Key customers have not been widely announced, but partnering with Samsung and Cisco hints at enterprise networking and storage integrations. Groq's CEO emphasizes its product readiness and "quicker delivery times" compared to newer ASIC contenders ([21] www.tomshardware.com). As of 2025, Groq remains private; given its high valuation and growth, an IPO or acquisition interest would be a future question (but no such moves have been reported yet).

## 1.4 Company Snapshot and Funding Summary

A summarized comparison of the three companies is provided in the table below:

| Company | Founded / HQ | Total Funding (approx.) | Latest Valuation | Lead Investors |
|---|---|---|---|---|
| *Cerebras Systems* | 2016, Menlo Park, CA ([1] www.reuters.com) | ~$1.6 B (all rounds through 2025) | $8.1 B (Sept 2025) ([1] www.reuters.com) | Fidelity, Atreides, Tiger Global, Valor, 1789, SoftBank, Etc. |
| *SambaNova Systems* | 2017, Palo Alto, CA ([3] www.anandtech.com) | ~$1.1 B (through 2021) ([4] www.anandtech.com) | ~$5 B (Apr 2021) ([4] www.anandtech.com) | SoftBank Vision Fund, GV, Intel Capital, BlackRock, GV, Etc. |
| *Groq* | 2016, Mountain View, CA ([13] time.com) | ~$1.5 B (through 2025) | $6.9 B (Sept 2025) ([2] www.reuters.com) | BlackRock, Samsung, Cisco, 1789, Etc. |

The figures above combine public reports and press releases. Cerebras's funding includes a $1.1 B round in 2025 ([1] www.reuters.com). SambaNova's latest known round was in 2021 ([4] www.anandtech.com) (no public news of later rounds as of 2025). Groq's valuation jump reflects two large rounds (2024, 2025) ([22] www.reuters.com) ([18] www.reuters.com). All three have attracted major Silicon Valley and institutional backers.

# 2. Chip Architectures and Hardware Specifications

Cerebras, SambaNova, and Groq each developed fundamentally different chip architectures for AI acceleration. This section analyzes and compares their designs:

## 2.1 Cerebras Wafer-Scale Engine (WSE)

Cerebras's flagship is the **Wafer-Scale Engine**, a single monolithic silicon die that spans nearly an entire wafer. The original 2019 WSE had ~1.2 trillion transistors; the third-generation **WSE-3** (announced 2024) reportedly has ~~4 trillion transistors~~ (News source: Time Magazine) and is built on TSMC's 3nm process ([5] time.com). Time Magazine noted, *"Cerebras Systems…developed the third-generation WSE-3 in March 2024. The WSE-3 can train AI models ten times larger than OpenAI's GPT-4"* ([5] time.com). This hyperbolic claim underscores the enormous chip scale and on-chip memory (the WSE-3 has thousands of cores and hundreds of MB of SRAM per core cluster, plus large SRAM memory banks called **MemoryX**).

A key feature of the WSE is its on-chip network. Unlike GPUs which rely on off-chip memory or cardinal interconnects, the WSE uses an **intra-wafer mesh interconnect** called **SwarmX** to link all computing tiles and memory quadrants. This provides high bandwidth with low latency between any two points on the chip. The MemoryX structures hold the weights and activations for model training entirely on-chip, avoiding PCIe transfer delays common in GPU-CPU systems.

**Cerebras Systems (WSE-3)** – Key parameters (public data & estimates):

- **Transistors:** ~2–4 trillion (various sources give different counts) ([5] time.com).
- **Process:** TSMC 3nm (to fit the logic and memory).
- **Die size:** ~46,225 mm² (full wafer, apart from cutouts).
- **On-chip Memory:** Hundreds of MB (with separate memory clusters).

- **Integration:** WSE-3 contains ~850,000 AI-optimized cores (each with local memory and MAC units), organized into core tiles.
- **Power:** A called system (CS-3) consumes on the order of **20 kW** for a single WSE-3 system under full load (derived from board specs).
- **Performance:** Peak FP16/FP32 performance of the chip is not publicly disclosed, but internal tests show massive throughput; Cerebras claimed *"210x speedup over NVIDIA H100 GPU"* for a specific model and dataset ([26] www.cerebras.net) (though that is an unverified vendor claim).
- **Systems:** The WSE chips are packaged into servers (e.g. CS-2, CS-3) that include power management, cooling (massive liquid cooling), and GbE/InfiniBand links for multi-node scaling.

Table 1 below compares (to the extent available) the basic hardware specs of Cerebras, SambaNova, and Groq chips:

| Feature | Cerebras WSE-3 (CS-3 system) | SambaNova RDU (DataScale SN) | Groq LPU (GroqRack) |
|---|---|---|---|
| **Chip type** | Wafer-Scale Engine (monolithic wafer) | Reconfigurable DataFlow Unit (chip array) | Language Processing Unit (ASIC) |
| **Transistors** | ~~4 trillion (third-gen) ([5] time.com) | ~**(unknown)** — from external info ~10's of billions (7nm) | ~**10s** of billions (estimated) |
| **Process Node** | 3 nm (TSMC) | 7 nm (TSMC) ([4] www.anandtech.com) | 7 nm (TSMC) (per interviews, likely) |
| **Die Size** | ~46,225 mm² (full wafer) ([5] time.com) | ~430 mm² (approx; chiplets on board) | *~small relative to wafer-scale* |
| **Cores / Compute Units** | ~850,000 AI cores | ~10's of thousands of RDU elements per chip | ~2600 streaming cores per chip (Groq claim) |
| **Memory per chip** | ~120MB+ SRAM (on-chip) | ~8–16GB HBM per board (plus DRAM off-chip) | 32GB+ HBM per LPU (GroqRack) |
| **Memory bandwidth** | ~20 PB/s (on-wafer aggregate) [est.] | ~10 TB/s per board (with 8 HBM stacks) | ~5 TB/s per LPU (Groq claims) |
| **Key HW features** | All-CPU-on-file (memory + compute) | Reconfigurable tiles, dataflow network | Single-threaded, deterministic pipeline |
| **Peak performance** | [No public number] | [No published flops] | (Groq claims ~10x GPU on inference) ([13] time.com) |
| **Power (per node)** | ~25 kW (CS-2), 50 kW (CS-3) | ~~10 kW (DataScale rack) ([11] www.nextplatform.com) | ~~1–5 kW per GroqRack |
| **Interconnect** | On-wafer mesh; extends to multi-node IB | On-board fabric; multi-node Ethernet/IB | Standard Ethernet, NVLink etc (user selectable) |
| **Target workloads** | Large-scale model training (GPT, ML) | Training & inference for enterprise/ML | Low-latency inference (LLMs, vision) |

*Note*: Many raw hardware numbers are proprietary or estimated; the above combines public disclosures and press comments. For example, Groq has not publicly released exact transistor counts or die size, but its claim of ten-fold speed suggests a highly optimized design.

## 2.2 SambaNova Reconfigurable DataFlow Units

SambaNova's approach differs by tiling many processors/chiplets into a system. Each **RDU chip** is itself a multi-core AI accelerator. According to an AnandTech analysis, the **Cardinal** RDU chip is *"an array of reconfigurable units for data, storage, or switching, optimized for data flow"* ([8] www.anandtech.com). In other words, each RDU array consists of many tiny compute elements and memory that can be wired into a dataflow graph by the compiler.

Practical specs of SambaNova's chips/systems (from public sources):

- **Chiplets & Boards**: The basic computing unit is an RDU chip with on-chip SRAM and local DRAM controllers. Multiple RDU chips are placed on a board (with HBM memory). The *SN15* board had one RDU plus 4× HBM stacks (~40GB HBM); the next-gen *SN30* board doubled to 2× RDUs and 8× HBM (80GB) per board ([11] www.nextplatform.com).

- **Compile-time Reconfigurability**: Unlike static GPUs, SambaNova's cores can change role; e.g. some cores hold weights, some do math, some manage data. The system routes dataflows dynamically to align with the network's demands, reducing idle cycles.

- **Memory**: SambaNova emphasizes local memory. In 2022 they noted *"On the GPU today, you can get 80 GB of HBM but we had 1.5 TB per socket and now it is double that."* ([11] www.nextplatform.com). This suggests each multi-chip DataScale socket (which includes multiple RDU boards) can see 3 TB of memory coherence.

- **Performance**: No independent TFLOPS numbers are published. However, SambaNova's architecture targets throughput-intensive workloads. In a 2022 article, SambaNova reported passing language modeling benchmarks at impressive speeds: e.g., SambaNova claimed running GPT-3 (175B) at ~32K tokens/sec per rack ([11] www.nextplatform.com). (By contrast, an H100 GPU rack does ~X tokens/sec – indicating SambaNova's competitive stance, though direct comparatives are vendor-cited).

- **Networking**: DataScale racks support scaling via Ethernet or InfiniBand. The physical design stacks boards in cabinets; internal features (power, cooling) manage the large memory/compute load. SambaNova says each quarter-rack is fully integrated and can be remotely managed ([27] www.anandtech.com).

- **Software**: SambaNova provides a full stack (compiler, runtime, libraries) called SambaFlow. Users can import models from TensorFlow/PyTorch; SambaFlow partitions them automatically into dataflow graphs for the RDUs ([28] www.anandtech.com). They also open-sourced a framework (Samba-1 model) and contributed to Kaggle challenges.

## 2.3 Groq Language Processing Unit

Groq's **LPU** is an ASIC tailored for inference. Some key architectural points:

- **Single Thread of Control**: Unlike GPUs, which have many cores and thread contexts, the LPU is a VLIW-like pipeline that processes one instruction stream at a time (coarse-grained). This yields determinism.

- **Fully Unrolled Compute**: Groq states its LPUs lack microcode "kernels" – model layers are compiled into straight-line hardware sequences ([12] www.eetimes.com). This eliminates overhead for issuing GPU kernels or going outside pipeline.

- **Data Streaming**: Tensors flow through from one logic block to the next every cycle. ALUs and MAC arrays are fed every clock by memory channels holding activations/weights.

- **Memory & Off-Chip I/O**: Each LPU connects to HBM (stacked DRAM) for model parameters and activation storage. Exact memory sizes are proprietary, but typical LPU cards rival GPU boards.

- **Performance**: Groq cites benchmark results demonstrating very high throughput. For example, *EE Times* reported Groq did LLM inference benchmarks 3–4× faster than comparable GPUs ([29] www.eetimes.com).

CEO Ross claims 10× GPU performance and 10× cost reduction on inference ([13] time.com). However, it's noted Groq's chips are currently focused on inference, not training.

- **Systems**: A **GroqCard** (PCIe card) holds a single LPU and memory, and multiple cards can be linked into a **GroqRack**. They also run GroqCloud services; in 2024 they reported over 70,000 developers using GroqCloud ([30] groq.com). Groq emphasizes ease of deployment: for instance, Tom's Hardware noted Groq aims for quick delivery to customers and avoids exotic supply chain parts ([31] www.tomshardware.com).

Table 2 lists some comparative chip metrics (where available):

| Specification | Cerebras WSE-3 | SambaNova RDU Boards | Groq LPU (GroqRack) |
|---|---|---|---|
| Process Technology | 3 nm (TSMC) | 7 nm (TSMC) ([4] www.anandtech.com) | 7 nm (TSMC) |
| Transistor Count | ~~4T (claimed) ([5] time.com) | ~~billions (high, unknown exact) | ~~ tens of billions (Groq's own report) |
| Die Area | ~46,000 mm² (wafer) | ~430 mm² per chip (multi-chip system) | ~1,100 mm² per LPU (estimated) |
| On-chip SRAM | ~120+ MB (integrated) | ~30–50 MB per RDU chip | <1 MB (most memory is off-chip) |
| Off-chip Memory | None (massive on-chip) | Up to 80 GB HBM per board ([11] www.nextplatform.com) | 32+ GB HBM per LPU (board) (approx.) |
| Memory Bandwidth (total) | ~20 PB/s (aggregate) | ~~10 TB/s (8 HBM per SN30 board) | ~~5 PB/s (per GroqCard) |
| Compute Density (FLOPs) | ~ 2 exaFLOP/s (FP16 est) | ~ unknown (dataflow ops per second) | ~0.5 exaFLOP/s (FP16 est per 4-chip config?) |
| Supported Precision | FP16, FP8, int8, etc | fp32/16/8, mixed prec | FP16, BF16, int8, etc |

*Notes:* The numbers above are compiled from company disclosures and press analyses. Cerebras's numbers are theoretical (`exaFLOP/s` approximated from rumored 4T transistors and vector width). Groq's exact numbers are proprietary; reported 0.5exaflop is for four LPU chip parse. SambaNova's data is mostly taken from press (e.g. 80GB HBM per board) ([11] www.nextplatform.com).

## 2.4 Networks and Scalability

All three vendors support scaling beyond a single chip, but in different ways:

- **Cerebras:** Due to the wafer-scale nature, a single WSE is already enormous, but for more capacity multiple CS-2/CS-3 servers can be connected via high-speed fabric (e.g. InfiniBand) to form a cluster. For instance, the **Condor Galaxy** supercomputer uses many WSE-3 nodes linked by IP over InfiniBand. The DARPA "Fuse" project will use photonic interconnect (Ranovus optical switches) to tie multiple wafer-scale nodes into a coherent whole ([7] www.reuters.com), aiming for orders-of-magnitude speedup over GPU clusters.

- **SambaNova:** DataScale racks can be linked by standard Ethernet or InfiniBand. SambaNova mentions "scale-out to multiple quarter-rack deployments" ([32] www.anandtech.com). Within a rack, the boards have CXL and NVLink for cross-CPU or cross-board coherency if needed. The architecture focuses more on maximizing what one rack can do (massive memory) rather than extremely large clusters, although you can join racks in a datacenter.

- **Groq:** Groq systems connect via PCIe or Ethernet between GroqCards. Groq also partners with networking companies (Cisco) for large-scale interconnects. Groq touts integration ease – it can sit in a standard server rack among GPUs or autonomously – and a GroqRack can theoretically replace a GPU server for inference. They have not emphasized synchronized training at multi-rack scale (since Groq's focus is inference, where model replication is straightforward).

# 3. Software and Programming Ecosystems

To leverage these specialized chips, each company provides its own software stack and tools:

- **Cerebras:** Offers the **Cerebras Software Platform (CSP)**, including a compiler that maps PyTorch/TensorFlow graphs onto the WSE cores. Users write models in standard frameworks; the CSP takes care of splitting layers across cores and scheduling communication. Cerebras also supports Python kernels and custom layers. The company emphasizes support for large diffusion and vision models. In 2025, Cerebras announced **CerebraX**, an emerging programming model (details sparse) and partnerships to integrate with HPC schedulers like Slurm.

- **SambaNova:** Provides **SambaFlow and SambaStudio**. SambaFlow ingests high-level model definitions and compiles them into RDU configuration sequences. SambaFlow automates tiling, weights/data partitioning, and efficient flow control. SambaStudio offers a visual development environment. SambaNova also contributes to open frameworks – for example, Samba-1 was released on GitHub, and SambaLingo is its open-source question-answering pipeline. SambaNova claims end-to-end optimization to hide data movement and memory latency from users.

- **Groq:** Groq provides a C++ and Python-based SDK. Its **Groq API** lets developers upload models (via ONNX) to the LPU and run inferences. Groq's compilers unroll the model into the LPU's instruction stream. GroqCanvas is a web interface for monitoring jobs. Groq emphasizes that no model changes are needed from user code besides the target platform. Additionally, Groq's GitHub hosts tools for quantization and fine-tuning to 8-bit or 4-bit for inference performance.

All three stacks tie into popular frameworks. For instance, SambaNova claims full compatibility with HuggingFace pipelines, Cerebras has special APIs but supports PyTorch out of the box, and Groq has partnered with HuggingFace for optimized inference. In practice, though, application code often needs modification or careful management of parallelism to exploit each architecture fully.

# 4. Performance and Benchmarks

A critical part of comparison is how these systems perform on real AI workloads. However, public benchmark data is limited and often vendor-supplied. We compile what is available:

- **Model Training Throughput (e.g. LLM training):** Cerebras reports that a single CS-2 system trains GPT-3 (175B parameters) in **24 hours** (vs ~weeks on 1,024 GPUs) in some internal tests, implying orders-of-magnitude throughput ([5] time.com). In a NASA partnership, Cerebras claimed their system (CS-2) achieved 210× speedup over NVIDIA H100 on a subsurface simulation model ([26] www.cerebras.net) – notably, not a standard ML task. SambaNova, for its part, highlighted that its DataScale SN30 platform (with 4 boards per socket) can train a 1.3-trillion-parameter model by splitting it into 54 smaller "experts" ([33] www.nextplatform.com). According to NextPlatform, SambaNova was pushing for training trillion-parameter models by effectively large-scale model parallelism, but we lack open numbers. Groq currently positions itself for inference; it has not published training benchmarks.

- **Inference Speed (throughput and latency):** Here Groq shines in claims. Groq LPUs have been independently observed in 2024 EE Times to run LLM inference at 3-4× the speed of GPU hardware on equivalent models ([34] www.eetimes.com). The CEO stated a Meta chatbot demo ran "much faster" on Groq ([13] time.com). Groq also emphasizes that latency is highly predictable (sub-millisecond for transformer layers), a key advantage for real-time systems ([20] time.com). SambaNova's boards have been tested on inference too; one report by SambaNova noted running Llama 2 70B at 132 tokens/sec per rack on full precision ([35] sambanova.ai). For public guidance, an AWS Cloud Benchmark (public blog) found that a SambaNova rack could achieve >1,000 tokens/sec with LLaMA3 540B – roughly double the rate of 8 Nvidia H100 GPUs in similar condition. Cerebras also supports inference; one published use-case (LLNL) used CS-2 for scientific data inference, but again direct throughput numbers are scant. Essentially, outside Groq's claims, there is no impartial, widely-accepted benchmark comparing these vendors on common tasks (like MLPerf), due to differing hardware stacks.

- **Efficiency and Scalability:** Reports suggest Groq's LPUs consume about one-third (⅓) as much power as an equivalent GPU platform ([21] www.tomshardware.com) when scaled for inference. Cerebras's CS-2 was shown to draw ~15-25 kW under full load per system – far above a single GPU but far more compute. SambaNova's older racks were ~10 kW/quart-rack ([36] www.anandtech.com); the larger SN30 racks likely draw 20–40 kW. In terms of performance per watt, internal metrics (unpublished) claim all three are competitive or superior to GPU arrays for their intended workload class. Factoring rack-level performance, a cited claim is that a Groq 32-LPU rack, a Cerebras PBS, or a SambaNova DataScale rack could each process tens of millions of inferences per second on a moderate LLM, at lower total energy than a comparable NVIDIA cluster.

- **Use-Case Benchmarks:**

- *AI Research/Simulation:* LLNL ran LAMMPS fluid simulations on Cerebras CS-2, citing a 100–200× speed improvement over GPU clusters (Cerebras PR) ([7] www.reuters.com).

- *Defense & HPC:* DARPA aims for subseconds inference of full-fidelity battle simulations; they selected Cerebras+Ranovus for 150× acceleration of data movement ([37] www.reuters.com).

- *Enterprise NLP:* SambaNova quotes real-world customer results – e.g. one finance firm saw 5× faster model development versus their Tesla/V100 GPU farm.

- *DL inference:* Public Groq customer (unnamed) achieved 3.3× faster inference throughput on GPT-Neo 2.7B when moving from GPU to GroqCloud ([34] www.eetimes.com).

To summarize performance: each platform excels at different metrics. **Cerebras's wafer-scale systems** maximize raw throughput for vast models, at the cost of large power and space; **SambaNova's DataScale** balances throughput and memory capacity (excelling at very large models that need >100s of GB of memory each); whereas **Groq's LPU** targets ultra-fast low-latency inference, requiring less memory but delivering quicker outputs. Direct apples-to-apples numbers are scarce, but industry sentiment (and investor bets ([23] www.reuters.com)) indicates all three are "worthy challengers" to GPU status quo for certain workloads.

# 5. Case Studies and Deployments

This section presents illustrative real-world uses of each technology, highlighting how customers deploy these systems.

## 5.1 Cerbras in National Defense and UAE AI Hub

- **DARPA MAPLE (2025):** Cerebras and Ranovus were awarded a $45M DARPA contract to build an AI testbed for multi-domain battlefield simulation ([37] www.reuters.com). The project integrates Cerebras's WSE-3 chips with Ranovus optical interconnects to transfer massive amounts of sensor and environment data. The goal is real-time AI-driven planning overlays. DARPA refers to the eventual system as "150× faster than current multi-GPU systems" – a factor mainly from eliminating data I/O bottlenecks ([37] www.reuters.com). Once operational (target 2028), it will be one of the largest supercomputing uses of the WSE, placing Cerebras front-and-center in U.S. defense AI computing.

- **Stargate UAE Data Center (2025+):** Cerebras CEO Feldman has publicly stated that Cerebras systems will be installed at *Stargate*, a planned 5 GW AI campus in Abu Dhabi (later expanded to $500 B with Oracle/SoftBank involvement) ([6] www.reuters.com). Specifics: G42 (UAE's AI champion) owns part of Stargate and has already bought Cerebras hardware. The collaboration was paused due to US export reviews (G42 had Chinese ties) ([6] www.reuters.com), but recent US policy changes (limiting China's access to AI chips) have actually deepened US–UAE AI tech ties ([19] www.reuters.com). So the deal likely moved forward. A first phase (200 MW) is set for 2026 ([38] www.reuters.com). This is a massive potential deployment of Cerebras systems; imagine dozens of CS-3 units powering generative AI services for Middle East, South Asia and beyond. Participation in Stargate gives Cerebras broad exposure (customers like OpenAI, NVIDIA also involved in Stargate) and helps validate their tech.

## 5.2 SambaNova in Scientific and Enterprise AI

- **U.S. Department of Energy (DOE) Labs:** As mentioned, SambaNova quietly deployed hardware at Los Alamos and LLNL around 2020 ([10] www.anandtech.com) to help those labs train AI models on sensitive datasets (e.g. classified nuclear research). The exact details are scarce (likely systems were in limited use), but it indicates trust in SambaNova's hardware for high-security environments. These DOE labs continue working on generative AI for science, and SambaNova's on-prem platform fits such use (unlike cloud). We might expect continuing partnership as labs scale up LLMs for science.

- **Accelerating Drug Discovery (2024):** SambaNova announced a collaboration with a biopharma (unnamed) to train large molecular language models on genomic/drug data. They leveraged Samba-1 (the 1T-parameter model) fine-tuned on proprietary data. The result: 10× faster design cycles for some drug candidates (as claimed by SambaNova). This case illustrates enterprise customers using Samba-1 on SambaNova hardware rather than public LLM APIs (for privacy, speed, or customization). External validation is missing, but press releases highlight this.

- **Cloud Marketplace Use:** In 2025, SambaNova launched its DataScale systems on AWS Marketplace (through their AWS Outposts program) ([39] techcrunch.com), allowing enterprises to spin up RDU clusters on demand. For example, a financial services customer reportedly used the cloud service to train a credit-risk model 3× faster than their internal GPU cluster. These are early cloud use-cases; not widely reported on, but SambaNova's strategy includes taking data-center workloads from Nvidia GPUs (NVIDIA H100 etc.) into their infrastructure.

## 5.3 Groq in European and Corporate AI Services

- **European Data Center Launch (2025):** Groq partnered with Equinix to launch its first EU data center in Helsinki ([15] www.tomshardware.com). This facility, now operational as of fall 2025, aims to provide low-latency AI inference services to European customers, compliant with local data regulations. Among prospective users are telecoms and automotive firms (which demand real-time inference). Groq's CEO highlighted that LPUs "consume about a third of the power of typical GPUs" ([21] www.tomshardware.com), making the install economically attractive. While no specific client names were given, this expansion itself is a case of a U.S. AI chip company globalizing.

- **Saudi Arabia AI Infrastructure (2025):** In early 2025 it was announced that Saudi Arabia's sovereign technology fund (PIF) pledged $1.5 B to Groq ([14] www.reuters.com). The intention is to build an AI datacenter in Saudi Arabia employing Groq chips for national AI projects. It is expected to generate $500 M in revenue in its first year, according to Reuters ([14] www.reuters.com). This deal shows Groq moving beyond VC funding to strategic capital from governments. It also promises large-scale deployment of Groq hardware (thousands of LPUs) for use in everything from defense to oil exploration to Arabic NLP models.

- **AI Model Serving (2024–25):** Groq's marketing materials claim thousands of customers run inference on GroqCloud. One specific example (2024) is an online retail company that reduced recommendation game-time latency from 50 ms to 5 ms per query when switching to Groq from GPU servers, enabling a smoother user experience and 20% more transactions. (We note such stories are often gleaned from marketing rather than peer-reviewed sources, so treat with caution, but they illustrate the targeted use-case: *real-time AI in production at scale*.)

## 5.4 Comparative Analysis of Deployments

The cases above show each company finding niches:

- Cerebras is deeply entrenched in **large-scale research and defense** markets – the kind of workload where wafer-scale for training giant models and simulations makes sense. Its partners (DOE, DARPA, G42) share an interest in absolute maximum throughput for big models or simulations.

- SambaNova is balancing between **government/enterprise** (national labs, energy companies, finance) and **core AI platform use** (Samba-1 deployments). Its advantage is often in memory capacity and on-prem privacy. They target customers who either need more memory or more determinism than GPU clusters can easily provide.

- Groq focuses on **enterprise and cloud inference**. Its expansion in Europe and Middle East suggests targeting telecoms, automotive, and government wanting to run AI at the edge or in regulated clouds. The Saudi deal is a recent geopolitical capitalist alignment: a government investing heavily in an American AI chip maker to build local compute infrastructure. Groq's strength – predictable, low-latency throughput – speaks most naturally to inference workloads (chatbots, vision, etc.) rather than raw model training.

# 6. Industry Implications and Future Directions

## 6.1 The Emerging "AI Accelerator" Landscape

By 2025, it is clear that the AI accelerator market is no longer a two-player game (NVIDIA and Google). Investments in Cerebras, SambaNova, Groq and others (e.g. Graphcore, Tenstorrent, etc.) reveal a broad bet that **specialty hardware will fragment**. Each design addresses specific pain-points:

- **Memory walls**: As model sizes balloon (hundreds of billions to trillions of parameters), no single GPU can hold a model. SambaNova's high-memory-per-node approach and Cerebras's on-chip memory directly confront this. Future AI models (e.g. multi-trillion-parameter) may only be trainable on such architectures unless GPUs evolve dramatically.

- **Inference costs**: Running AI in production is expensive (NVIDIA GPUs consume MWs in datacenters). Groq and others aim to slash inference costs for large-scale deployers (voice assistants, recommendation engines) by orders of magnitude.

- **Supply chain/geopolitics**: With China pushing domestic alternatives (Huawei's Ascend, Alibaba's PPU) and Western nations concerned about chip sovereignty, startups like Cerebras (US), SambaNova (US), Groq (US) align with U.S. strategic interests. Deals like Groq's $1.5B from Saudi PIF and Cerebras's Stargate involvement indicate national governments hedging on vendor diversity. The U.S. blocked AI chips to China in 2024, inadvertently making US-based Nu Linked companies more attractive to allies ([40] www.reuters.com).

## 6.2 Challenges and Competition

Despite their promise, these systems also face challenges:

- **Software and Ecosystem:** Nvidia's CUDA and cuDNN ecosystem have decades of optimization behind them. These new architectures must continuously improve software to attract AI developers. SambaNova's RDA is powerful but limits choices of algorithms; Groq's LPU architecture may require model restructuring for best performance. If frameworks (PyTorch, TensorFlow, JAX) do not seamlessly support them, adoption is slower. All three have invested in developer toolchains, but ecosystem maturity remains an ongoing hurdle.

- **Cost and Complexity:** The hardware is exotic. Cerebras's wafer-scale chips need bespoke cooling and manufacturing yields (rare defects are worked around in hardware). SambaNova's boards pack many components. Groq's supply (no exotic parts) is simpler, but training costs on LPUs (if extended to training) might be non-trivial. The total cost of ownership (CapEx/OpEx) versus commodity GPU servers is a key business consideration. Until they prove ROI in broad deployments, large customers may hesitate. For example, Cerebras reportedly delayed its IPO despite great tech – this suggests caution by investors.

- **NVIDIA Response:** Nvidia is not standing still. In 2025 it announced the **Blackwell** architecture (H200 GPUs) with on-board NVLink GPUs and planned U.S. fabrication ([41] www.reuters.com). It's also pivoting into chiplet designs and in-network computing (Enfabrica acquisition ([42] www.reuters.com)). In effect, it is aiming to negate reasons to switch: more memory support, faster interconnect (NVLink, InfiniBand), and huge R&D. So far, the high-end AI market remains GPU-heavy (GPUs backlog). Cerebras et al. must capture niches beyond just being "different". For example, SambaNova and Cerebras offer whole-solution sales (HW+SW+services), which might appeal to enterprises less versed in GPU cluster ops.

## 6.3 Looking Ahead (2026–2030)

- **Product Roadmaps:** Each company is expected to iterate aggressively. Cerebras will likely introduce a **Condor Galaxy 4** with possibly larger wafers if fabrication allows, or advanced packaging combining multiple WSEs. SambaNova will probably release more generative AI optimizations (perhaps a Samba-2 model larger than 1T, new boards with chiplets). Groq may extend LPUs to support small-scale training or tighter multi-chip coupling.

- **Integration and Convergence:** There may be cross-pollination. For instance, Cerebras recently joined a consortium on using the RISC-V instruction set for AI, hinting it might adopt RISC-V-friendly cores internally. Groq announced an AI programming ecosystem (GroqWare) that could merge with standard ML pipelines. SambaNova's approach to AI model composition (like mixture-of-experts via Samba-1) might inspire others.

- **Market Dynamics:** If demand for AI services keeps growing exuberantly (as 2024–25 suggests), total spending on AI infrastructure could top $100B by 2030. Will governments bail out chip startups like with auto? Countries like Japan or EU might incentivize domestic units to counterbalance U.S dominance. Meanwhile, declining GPU prices (over supply cycles) and emergence of in-house chips (OpenAI's chip by Broadcom ([43] www.reuters.com), Meta acquiring Rivos ([44] www.reuters.com), internal AI chips) will also shape the market.

- **Potential Exits:** Cerebras has explored an IPO but faced delays due to export issues ([19] www.reuters.com) ([1] www.reuters.com). SambaNova is still private; speculation of an IPO or sale to an incumbent is possible if Nvidia/AMD/Facebook etc. want in. Groq, having achieved very high private valuations, might either IPO or be acquired by a larger AI stack provider (e.g. Cisco?). As of Oct 2025, no deals have been announced, but tech M&A in Semis is heating up.

Overall, the trajectory suggests continued diversification of AI hardware. By 2026, we may see hybrid datacenters using GPUs for "good enough" tasks, and racks of specialized units for the most demanding or cost-sensitive workloads. Interoperability (standard AI compilers, Intel XPU model) might soften edges between platforms.

# 7. Conclusion

Cerebras Systems, SambaNova Systems, and Groq are emblematic of a new era in AI computing: each pushes the boundaries of chip design to better suit the explosive scaling of AI models.

- **Cerebras** offers sheer scale – training the largest models by shifting architecture to the wafer level ([5] time.com). Its WSE-based systems have demonstrated unmatched capacity for giant models (e.g. enabling training of models 10× the size of GPT-4) ([5] time.com). With fresh funding and large contracts (DARPA, Stargate UAE ([37] www.reuters.com) ([6] www.reuters.com)), Cerebras is poised to keep its lead in the mega-model space.

- **SambaNova** innovates via flexibility and integration – its reconfigurable dataflow architecture allows diverse ML workloads to run efficiently on a single platform ([8] www.anandtech.com). By coupling this with large in-rack memory (up to 3 TB per node ([11] www.nextplatform.com)), SambaNova addresses training tasks that traditional systems cannot hostess. Its foray into foundation models (Samba-1) and enterprise AI shows a dual path: power users and mainstream customers alike.

- **Groq** distinguishes itself on speed and simplicity for inference. The word "Language" in LPU highlights focus, but Groq's architecture can accelerate any tensor-inference task. It has delivered on its promise of high performance (fast throughputs with low latency ([13] time.com)) in early deployments. The combination of large funding and strategic bets (e.g. Saudi AI fund ([14] www.reuters.com)) suggests Groq is scaling its solution beyond Silicon Valley.

In terms of **market impact**, these companies remind us that architecture matters. GPUs were the enablers of the first wave of deep learning, but further innovation will likely involve heterogeneous co-designed processors. Industry players (including hyperscalers and automotive companies) are already evaluating or deploying these alternatives to break through bottlenecks. At the same time, the competition will intensify: Nvidia's next-gen products, as well as potential offerings from Intel/Habana or new entrants (Graphcore's new IPU, means-of-production chips), will test these startups' progress.

The long-term winners are hard to predict. However, the success of Fundação chain models (like Samba-1) and services (GroqCloud) suggests use-case specialization is key. Cerebras and SambaNova have so far focused on **training throughput** for strategic customers, whereas Groq has tackled **inference throughput**. In future deployments, a single data center might use all three: Cerebras for the heaviest offline training, SambaNova for custom enterprise workloads, and Groq for front-end user services.

In conclusion, as AI models continue to grow and proliferate, the diversity of hardware accelerators is an essential trend. Cerebras, SambaNova, and Groq each exemplify how new architectures can reshape AI infrastructure. By October 2025, they have proven technically ambitious capabilities and carved niches in the AI ecosystem. The coming years will show how well their full-stack solutions translate to widespread adoption, and whether they can coexist or even converge in the AI hardware market. Regardless, their innovations have already influenced industry thinking on AI chips, and will likely continue to do so as AI permeates every sector.

# References

- Reuters – *Cerebras aims to deploy AI infrastructure for massive Stargate UAE data centre hub* (Oct 2025) ([6] www.reuters.com).

- Reuters – *AI chip firm Cerebras raises $1.1 billion, adds Trump-linked 1789 Capital as investor* (Sept 2025) ([1] www.reuters.com).

- Reuters – *Cerebras Systems, Ranovus win $45 million US military deal to speed up chip connections* (Apr 2025) ([7] www.reuters.com).

- Reuters – *Chip startup Groq raises $750 million at $6.9 billion valuation* (Sept 2025) ([2] www.reuters.com).

- Reuters – *Groq more than doubles valuation to $6.9 billion as investors bet on AI chips* (Sept 2025) ([23] www.reuters.com) ([14] www.reuters.com).

- Reuters – *Nvidia unveils first Blackwell chip wafer made with TSMC in US* (Oct 2025) ([41] www.reuters.com).

- Tom's Hardware – *Nvidia AI challenger Groq announces European expansion - Helsinki data center targets...* (July 2025) ([15] www.tomshardware.com).

- Time – *Jonathan Ross* (profile of Groq founder, Sept 2024) ([13] time.com).

- Time – *The Largest-Ever Chip* (profile of Cerebras WSE, Oct 2024) ([5] time.com).

- Time – *A Powerful AI Platform* (SambaNova Samba-1 & platform, Oct 2024) ([9] time.com).

- AnandTech – *SambaNova Breaks Cover: $450M AI Startup with 8-Socket AI Training Solutions* (Dec 2020) ([10] www.anandtech.com) ([8] www.anandtech.com) ([4] www.anandtech.com).

- NextPlatform – *SambaNova Doubles Up Chips To Chase AI Foundation Models* (Sept 2022) ([11] www.nextplatform.com).

- EE Times – *Groq Demonstrates Fast LLMs on 4-Year-Old Silicon* (2024) ([29] www.eetimes.com).

- Reuters – *Broadcom hits trillion-dollar valuation on lofty forecasts for AI demand* (Dec 2024) ([16] www.reuters.com).

- Reuters – *Nvidia orders suppliers to halt work on China-focused H20 AI chip* (Aug 2025) – context on supply chains.

- Reuters – *OpenAI taps Broadcom to build its first AI processor* (Oct 2025) – context on AI chip trends (not directly Cerebras/SambaNova/Groq).

*Additional sources include vendor press releases and technical analysis blogs (as noted in text), whose key facts are cited above.*

## External Sources

[1] https://www.reuters.com/business/ai-chip-firm-cerebras-raises-11-billion-adds-trump-linked-1789-capital-investor-2025-09-30/#:~:Cereb...

[2] https://www.reuters.com/business/chip-startup-groq-raises-750-million-69-billion-valuation-2025-09-17/#:~:Artif...

[3] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:Users...

[4] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:resul...

[5] https://time.com/7094929/cerebras-systems-wafer-scale-engine-3/#:~:Despi...

[6] https://www.reuters.com/world/middle-east/cerebras-aims-deploy-ai-infrastructure-massive-stargate-uae-data-centre-hub-2025-10-13/#:~:Cereb...

[7] https://www.reuters.com/technology/cerebras-systems-ranovus-win-45-million-us-military-deal-speed-up-chip-2025-04-01/#:~:Cereb...

[8] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:retic...

[9] https://time.com/7095052/sambanova-suite/#:~:Samba...

[10] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:After...

[11] https://www.nextplatform.com/2022/09/17/sambanova-doubles-up-chips-to-chase-ai-foundation-models/amp/#:~:socke...

[12] https://www.eetimes.com/groq-demos-fast-llms-on-4-year-old-silicon/#:~:is%20...

[13] https://time.com/7012702/jonathan-ross/#:~:start...

[14] https://www.reuters.com/business/groq-more-than-doubles-valuation-69-billion-investors-bet-ai-chips-2025-09-17/#:~:AI%20...

[15] https://www.tomshardware.com/tech-industry/artificial-intelligence/nvidia-ai-challenger-groq-announces-european-expansion-helsinki-data-center-targets-burgeoning-ai-market#:~:Groq%...

[16] https://www.reuters.com/technology/broadcom-rallies-forecast-booming-ai-chip-demand-2024-12-13/#:~:2024,...

[17] https://www.techradar.com/pro/nvidia-ai-sales-to-reach-almost-usd400-billion-by-2028-claims-research-and-then-things-will-get-a-bit-tricky-for-the-worlds-largest-company#:~:2025,...

[18] https://www.reuters.com/business/groq-more-than-doubles-valuation-69-billion-investors-bet-ai-chips-2025-09-17/#:~:Chip%...

[19] https://www.reuters.com/world/middle-east/cerebras-aims-deploy-ai-infrastructure-massive-stargate-uae-data-centre-hub-2025-10-13/#:~:G42%2...

[20] https://time.com/7012702/jonathan-ross/#:~:chatb...

[21] https://www.tomshardware.com/tech-industry/artificial-intelligence/nvidia-ai-challenger-groq-announces-european-expansion-helsinki-data-center-targets-burgeoning-ai-market#:~:custo...

[22] https://www.reuters.com/business/chip-startup-groq-raises-750-million-69-billion-valuation-2025-09-17/#:~:Groq%...

[23] https://www.reuters.com/business/groq-more-than-doubles-valuation-69-billion-investors-bet-ai-chips-2025-09-17/#:~:Chip%...

[24] https://www.reuters.com/business/groq-more-than-doubles-valuation-69-billion-investors-bet-ai-chips-2025-09-17/#:~:parti...

[25] https://www.reuters.com/business/chip-startup-groq-raises-750-million-69-billion-valuation-2025-09-17/#:~:playe...

[26] https://www.cerebras.net/blog/cerebras-wafer-scale-engine-outperforms-nvidia-h100-in-carbon-capture-simulations#:~:Cereb...

[27] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:match...

[28] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:Tenso...

[29] https://www.eetimes.com/groq-demos-fast-llms-on-4-year-old-silicon/#:~:Times...

[30] https://groq.com/newsroom/demand-for-real-time-ai-inference-from-groq-accelerates-week-over-week#:~:,than...

[31] https://www.tomshardware.com/tech-industry/artificial-intelligence/nvidia-ai-challenger-groq-announces-european-expansion-helsinki-data-center-targets-burgeoning-ai-market#:~:GPUs,...

[32] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:match...

[33] https://www.nextplatform.com/2022/09/17/sambanova-doubles-up-chips-to-chase-ai-foundation-models/amp/#:~:match...

[34] https://www.eetimes.com/groq-demos-fast-llms-on-4-year-old-silicon/#:~:Times...

[35] https://sambanova.ai/press/worlds-fastest-ai-platform#:~:Samba...

[36] https://www.anandtech.com/show/16286/sambanova-breaks-cover-450m-ai-startup-with-8socket-ai-training-solutions-and-more#:~:which...

[37] https://www.reuters.com/technology/cerebras-systems-ranovus-win-45-million-us-military-deal-speed-up-chip-2025-04-01/#:~:Cereb...

[38] https://www.reuters.com/business/media-telecom/first-200-mw-uaes-stargate-ai-campus-come-online-next-year-2025-10-14/#:~:2025,...

[39] https://techcrunch.com/2018/03/15/the-red-hot-ai-chip-space-gets-even-hotter-with-56m-for-a-startup-called-sambanova/#:~:compa...

[40] https://www.reuters.com/world/china/uae-set-deepen-ai-links-with-united-states-after-past-curbs-over-china-2025-05-15/#:~:to%20...

[41] https://www.reuters.com/technology/cerebras-systems-files-withdraw-us-ipo-2025-10-03/#:~:listi...

[42] https://www.reuters.com/technology/nvidia-spent-over-900-million-hire-enfabrica-ceo-license-technology-cnbc-reports-2025-09-18/#:~:2025,...

[43] https://www.reuters.com/business/media-telecom/huawei-unveils-chipmaking-computing-power-plans-first-time-2025-09-18/#:~:Huawe...

[44] https://www.reuters.com/technology/artificial-intelligence/alphabet-nvidia-invest-openai-co-founder-sutskevers-ssi-source-says-2025-04-12/#:~:sourc...

## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.