



LLMs: Taking the World By Storm

Large Language Models (LLMs) are revolutionizing various sectors, including cybersecurity, government, research, finance, and enterprise communications. These advanced language processing tools can generate, classify, and summarize text with remarkable coherence and accuracy with the ability to predict the next word in a sentence by analyzing vast amounts of data. LLMs are transforming the way we work and communicate and their potential applications are vast, with their impact on various markets only beginning to be realized.

Fully unlocking the potential of LLMs requires:

1. Real-time interactions 2. Differentiated performance from your competition

Groq is an AI-born, software-first platform, made to run real-time AI solutions at scale. Groq invented and is delivering the first Language Processing Unit system, a chip specifically designed to power LLMs for the exploding AI market, and the computers of the future. Our LPU system, built on our unique and pioneering architecture, transforms the pace, predictability, performance, and accuracy of AI solutions for LLMs.

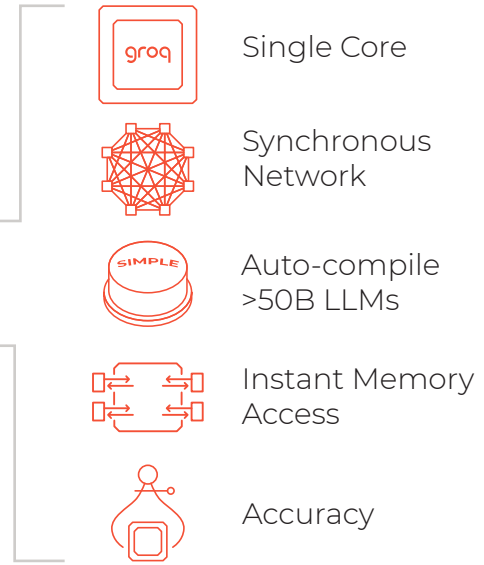
- We offer superior performance for LLMs at scale, delivering real-time outcomes, higher throughput, greater accuracy, and more efficient scalability.
- We accelerate the pace of AI workload development through rapid, kernel-less compilation, improving time to market, reducing resource requirements, and keeping up with the pace of innovation.
- We deliver predictability, providing accurate data on workload performance and costs at compile time so that developers can optimize software design with full understanding of how it will run when deployed.
- We enable higher accuracy, as lower latency allows for software techniques to deliver better predictions and improved business results in real-time.

We envision an AI solutions ecosystem that moves at the pace of software, unconstrained by the slow pace and high costs of big chip makers' hardware development cycles.

KEY ATTRIBUTES OF AN LPU™ SYSTEM



- 10X performance advantage
- Lower cost & lower power
- No supply constraints



Leveraging our software and hardware ecosystem, Groq can get multi-billion parameter LLMs and other AI models up and **running in less than five days**. Today, we're running Llama-2 70B at over 300 tokens per second per user, record-breaking performance that has continued to build over the past three months. Groq offers the fastest inference performance, in tokens per second per user, in the industry. Ask us about our most recently compiled models.

| July 18th | July 24th | July 29th | August 3rd | August 31 | Today |
|----------------------|--|---|---|---------------------------------------|--|
| Model released | Model compiling five days after first download | Performance five days after first compile | Performance 10 days after first compile | Performance 38 days after 1st compile | 30x performance increase since first compile |
| Llama-2 70B released | 10 T/s per user initial performance | 65 T/s per user | 100 T/s per user | 240 T/s per user | >300 T/s per user |

Reach out at contact@groq.com to experience real-time LLM performance.

