NANDA 2024:

Think Different: The Unique Role of FPGAs in a GPU-Dominated World

Sept 9, 2024

Andrew Putnam

Partner GM, Cloud & Al Hardware Engineering Microsoft Azure



The AI Era

A NEW ERA

The Age of AI has begun

Artificial intelligence is as revolutionary as mobile phones and the Internet.

By Bill Gates | March 21, 2023 • 14 minute read



In my lifetime, I've seen two demonstrations of technology that struck me as revolutionary.

The first time was in 1980, when I was introduced to a graphical user interface...

The second big surprise came just last year.

https://www.gatesnotes.com/The-Age-of-Al-Has-Begun

Microsoft

Eras Tour – Computing Edition



Mainframe



PC

Mobile



Cloud



Al







0





ASIC Prototype

Networking & Telecom

Specialized Compute

Cloud

Al

Success in Prior Eras



ASIC Prototype Ne

Networking & Telecom

Specialized Compute

Cloud

AI

Success in Prior Eras



ASIC Prototype

Networking & Telecom

Specialized Compute

Cloud

AI

Can FPGAs compete when the biggest players in computing are focused on the *same things* as FPGA researchers and developers?

Be First

Be Special

Can FPGAs compete when the biggest players in computing are focused on the *same things* as FPGA researchers and developers?

Be Special

FPGA Target Markets

Cloud / Datacenter



Azure Accelerated Networking / Azure Boost



Project Brainwave











Microsoft previews Azure Boost for offloading

virtualized network functions - Converge Dige...









adidas Yeezy Boost 380 Azure FZ4986 Release Date | SneakerNews.com



adidas Yeezy Boost 700 Faded Azure GZ2002



adidas Yeezy Boost 700 Faded Azure GZ2002





Adidas Yeezy

Boost 700 'Fad...

adidas Yeezy Boost 700 'Faded Azure' For Sale - T ... for developers

Microsoft makes Azure Bot Service generally available

First Look at the adidas Yeezy Boost 380 Azure • KicksOnFi ...



Release Details: Yeezy BOOST 700 'Faded Azure' -Sneaker Freaker



adidas Yeezy Boost 700 Faded adidas Yeezy Boost 380 Azure FZ4986 Azure GZ2002 Release Date Info | Sneak



First Look at the adidas Yeezy Boost 380 Azure • KicksOnFire.com



adidas Yeezy Boost 700 Faded Azure GZ2002

Release Date - Sneaker Ba...



adidas Yeezy Boost 700 Faded Azure Release Date - SBD



Processes



Azure Boost Launches to Offload Virtualization



adidas YEEZY BOOST 380 'Azure' (Azure/Azure/Azure) - FZ4986 - Co ...

Releases Tomorrow | Laptrinh...



Azure

Microsof



adidas Yeezy Boost 700 Faded Azure Release Date Info - SNKRS WORLD





Adidas Yeezy Boost 700 'Faded Azure' GZ2002 R.



The new Microsoft Azure Boost delivers accelerated improvements in performance, s...











The adidas Yeezy Boost 700 "Fade Azure"





adidas Yeezy Boost 380 Azure FZ4986 Release Date - Sneaker Bar Detroit





0

5+ stores

Announcing the general availability of Azure Boost

By Max Uritsky Published Nov 15 2023 08:00 AM @ 8,433 Views

ଲ

We're excited today to announce the general availability of <u>Azure Boost</u>, a system designed by Microsoft that offloads server virtualization processes traditionally performed by the hypervisor and host OS onto purpose-built software and hardware, enabling faster storage and networking performance for Azure VM customers.

Going forward, every new Azure virtual machine (VM) series will benefit from Azure Boost technologies, making the networking and storage components of your virtual workloads run even faster, whether you have deployed general purpose compute workloads or specialized AI clusters.

Azure Boost was engineered by a team of Microsoft hardware and software engineers to enhance the performance, security, and reliability of Microsoft Azure, and is already being used at Microsoft datacenters around the world, delivering benefits to millions of customers VMs in production today.

To learn more, watch the product <u>overview video</u>, read the <u>documentation</u>, or experience Azure Boost benefits today by provisining Azure Boost enabled VMs sizes listed in the documentation. Continue reading to learn more about the benefits of using Azure Boost.



Azure Boost

Azure Boost enables customers to achieve a **remote storage throughput up to 12.5 GBps and 650K IOPS and a local storage throughput of up to 17.3 GBps and 3.8 million IOPS**, enabling the fastest storage workloads available today.

In addition, Azure customers can to achieve a up to **200 Gbps networking throughput**.

https://azure.microsoft.com/en-us/updates/public-preview-new-generation-amd-vms-eav6dav6fav6 https://techcommunity.microsoft.com/t5/azure-compute-blog/announcing-preview-of-new-azure-dlsv6-dsv6-esv6-vms-with-new-cpu/ba-p/4192124



FPGA Target Markets

Cloud / Datacenter



Azure Accelerated Networking / Azure Boost





Microsoft

What is special in BrainWave?



Low batch size optimizations

Microsoft

RDMA transport for AI Networks

Custom Data Types for Efficient Inference

(MSFP / MXFP)



Largest deployment of AI FPGAs

2018

A configurable cloud-scale DNN processor for real-time AI

by Jeremy Fowers, Kalin Ovtcharov, Michael Papamichael, Todd Massengill, Ming Liu, Daniel Lo, Shlomi Alkalay, Michael Haselman, Logan Adams, Mahdi Ghandi, Stephen Heil,

Prerak Patel, Adam Sapek, Gabriel Weisz, Lisa Woods, Sitaram Lanka, Steven K. Reinhardt, Adr

[Retrospective]

RETROSPECTIVE: A Configurable Cloud-Scale DNN Processor for Real-Time AI

Eric S. Chung Doug Burger Jeremy Fowers Mahdi Ghandi Gabriel Weisz Sitaram Lanka Steven K. Reinhardt

I. CONTEXT running on-chip SRAM-pinned GEMV kernels. It was uniquely As we reflect upon the work we started in 2015 and its optimized for low latency and featured chained instructions that food multiple SIMD operations. These optimizations

Since 2015, Brainwave has undergone significant advancements and evolutions. Over this period, we have developed and deployed at least five generations of soft NPUs, strategically integrating them across a fleet spanning three generations of FPGAs: Stratix V, Arria 10, and Stratix 10. These advancements have enabled us to scale production and effectively power over 100 models across Bing and Office, thereby supporting a wide range of services such as ranking, semantic vectorization, question and answer generation, image processing, and query intent understanding [5].

II. BRAINWAVE ARCHITECTURE

IV. PITFALLS AND SURPRISES

Our project was not devoid of detours and surprises. The on-chip SRAM approach introduced by Project Brainwave gradually lost its relevance as the compute-bound Transformerbased models surpassed the capacity of on-chip memory hierarchies. These models no longer necessitated aggressive, SRAM-based model parallelism to achieve high utilization. Additionally, Nvidia GPUs aggressively scaled their High-Bandwidth Memory (HBM) and hardened tensor units to improve their AI performance, which eventually became predominant in most mainstream AI processors.

This limited the opportunity space for novel architectures to compete, and despite the development of the Stratix 10 NX FPGA [1], which introduced an HBM-enabled FPGA with hard tensor cores, the scope and scale at which Catapult FPGAs could competitively serve flagship AI models were reduced. Ultimately, the hardware lottery [7] effect made mainstream GPUs a much more attractive target for model evolution.



The Brainwave system architecture exploited model paral-

and innovation in hardware IP were embraced. One example of such IP that exists because of Project lelism in a unique way where LSTM weights were carefully Brainwave, MSFP [6], [9], enabled us to fit LSTMs with partitioned and pinned in on-chip SRAM across multiple millions of parameters into FPGAs with merabytes of SRAM

span or weeks to months. This approach instance occurets in our FPGA teams, fostering an environment where risk-taking

The Hardware Lottery

Key Insights

- The term hardware lottery describes a research idea that wins due to its compatibility with available software and hardware, not its superiority over alternative research directions.
- We may be in the midst of a present-day hardware lottery. Hardware design has prioritized delivering on commercial use cases, while built-in flexibility to accommodate the next generation of ideas remains a secondary consideration.



Sara Hooker

 Any attempt to avoid future hardware lotteries must be concerned with making it cheaper and less time consuming to explore different hardware/software/algorithm combinations.

Hardware Lottery Quadrants

CUDA







Custom HW Arch

GPU HW Arch



bric for Accelerating Large-Scale Data













Microsoft





Seattle 🦾 59 🕫 anputnam@microsoft.com 🦂 127530 😲 \$0.00 💲 🚍





Just hanging out

Write a C# program to check if a number is prime

Ask Copilot

Copilot

Yes, Bing still uses FPGAs (field-programmable gate arrays) to enhance the speed and accuracy of its search results. Microsoft initially deployed FPGAs in its data centers to improve the efficiency of Bing's search processing 1. These reconfigurable chips help deliver faster and more relevant search results by accelerating specific tasks that would otherwise be handled by traditional CPUs¹.

Is there anything specific you'd like to know about how FPGAs work or their applications?

Learn more 1 computerworld.com 2 reddit.com 3 computerworld.com

1 of 30 responses Al-generated content may be incorrect.

- Brainwave is competitive (but generally inferior) on performance metrics compared to GPUs ٠
- Brainwave is superior on metrics when overall cost is included Microsoft

FPGA Target Markets

Cloud / Datacenter



Azure Accelerated Networking / Azure Boost

Project Brainwave

FPGA Target Markets

Cloud / Datacenter

Azure Accelerated Networking / Azure Boost

AI / ML

Project Brainwave

- Customers don't *actually* want you to differentiate
- Lack of common shell, tools, makes portability difficult
- 3rd party programs largely failed

- Space is evolving too quickly to care about HW innovation
- Not clear how to make money on the mainstream, much less the optimization
- FPGAs didn't win the HW Lottery

Go be special... Compete when the competition of the special... computing are focused on the source trans s as FPGA researcher and developers?

> What is "Special"

Define "Special"

- Special is more than what gets you published
 - Performance Power Perf / W
- **Special** includes anything that has *impact*
 - **Relevant, application-level** Perf, Perf / \$
 - NOT FLOPS, OPS, inferences/s on MNIST
 - What does the customer actually see & experience?

Hardware Companies

Software Companies

Everywhere or Nowhere

- Until an accelerator is in every relevant region, software needs 2 branches
 - no-accelerator
 - with-accelerator
- Code branches... tend to branch
- 65+ Regions, 6 continents worldwide
- And growing...
- 6+ generations of servers active in Azure
- Servers deployed in 2014 are still running
- Hardware for those servers was designed 2009-2013
- Customers won't see a consistent experience
- Even *better-than-expected* performance can still result in complaints
- Special: designs that work with old and new hardware

No Complaints

Customer Complaints

Allow Experimentation

- Bing accelerator (ISCA 2014) began as an accelerator for scoring Decision Tree forests
- Software optimizations reduced 250x speedup to 30x speedup in just a few months
- Had to find new portions of the workload to accelerate
- But software got *better*, and became *ready for accelerators*
- We never shipped DT scoring... but shipped a LOT more
- **Special**: Allow experimentation

	Software	Hardware	Speed
	(us)	(us)	up
Feb	5000	20	250
March	250	8	31
June	125	2	63

Even if Catapult never ships a single FPGA into production, you've had a major impact on the performance and relevance of Bing.

- VP, Bing Platforms

Adaptation for Multiple Program Phases

- Code you accelerate is only part of the full algorithm
 - 100x can become 2x or less easily
- Need tightly-integrated CPU platforms
 - GPU SKUs are ~17X more \$ than CPU-only
- Avoiding work provides the best speedup
- **Special**: Adapting to multiple program phases

Reconstructed as one massive jet with substructure

000

201

And ruine

Microsoft will be carbon negative by 2030

Jan 16, 2020 | Brad Smith - President & Vice Chair

Microsoft President Brad Smith, Chief Financial Officer Amy Hood and CEO Satya Nadella preparing to announce Microsoft's plan to be carbon negative by 2030. (Jan. 15, 2020/Photo by Brian Smale)

Microsoft's pathway to carbon negative by 2030

Annual carbon emissions

2024 Update

Our 2024 Environmental Sustainability Report

May 15, 2024 | Brad Smith – Vice Chair and President; Melanie Nakagawa – Chief Sustainability Officer

https://blogs.microsoft.com/on-the-issues/2024/05/15/microsoft-environmental-sustainability-report-2024/

Carbon reduction continues to be an area of focus, especially as we work to address Scope 3 emissions. In 2023, we saw our Scope 1 and 2 emissions decrease by 6.3% from our 2020 baseline. This area remains on track to meet our goals. But our indirect emissions (Scope 3) increased by 30.9%. In aggregate, across all Scopes 1–3, **Microsoft's emissions are up 29.1% from the 2020 baseline**.

The rise in our Scope 3 emissions **primarily comes from the construction of more datacenters and the associated embodied carbon** in building materials, **as well as hardware components such as semiconductors, servers, and racks**. Our challenges are in part unique to our position as a leading cloud supplier that is expanding its datacenters. But, even more, we reflect the challenges the world must overcome to develop and use greener concrete, steel, fuels, and chips. These are the biggest drivers of our Scope 3 challenges.

What is better for the environment?

Be Special – There's more than just performance and power efficiency

What is better for the environment?

Be Special – There's more than just performance and power efficiency

Conclusions

- FPGAs compete by being **First** and being **Special**... okay... Special.
- Cloud and AI markets are large enough to compete on both metrics faster than FPGA hardware alone can evolve
- Competing head-on isn't the path that historically made FPGAs successful
- There is more to being "special" than perf, power, and cost
- Embrace backwards compatibility and other metrics SW developers & customers care about
- Think about multiple simultaneous generations of hardware, software & applications
- Think about worldwide deployment
- Think about having more than one trick
- "Special" *impactful* metrics come from adaptation as much as perf/cost/power
- Even carbon efficiency comes more from *adaptation* than from lower energy usage

© Copyright Microsoft Corporation. All rights reserved.