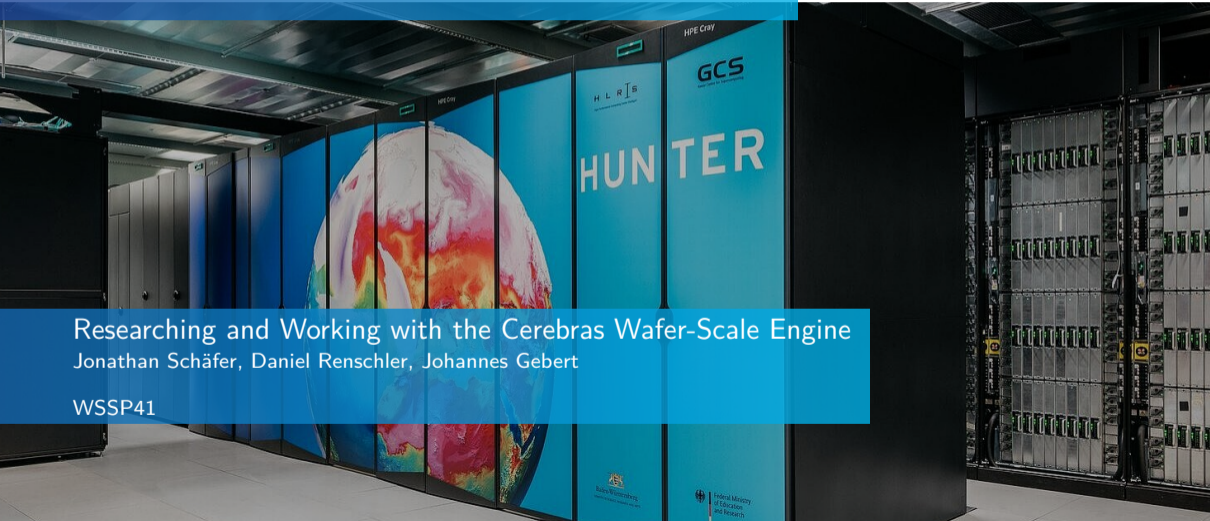


## Researching and Working with the Cerebras Wafer-Scale Engine

Jonathan Schäfer, Daniel Renschler, Johannes Gebert

WSSP41



Motivation

The Cerebras Wafer-Scale Engine in Numbers

The Cerebras Wafer-Scale Engine in Structure

What Have We (FCG@HLRS) Done Lately?

Further Work



# Technological Limits. For Example: Memory Wall

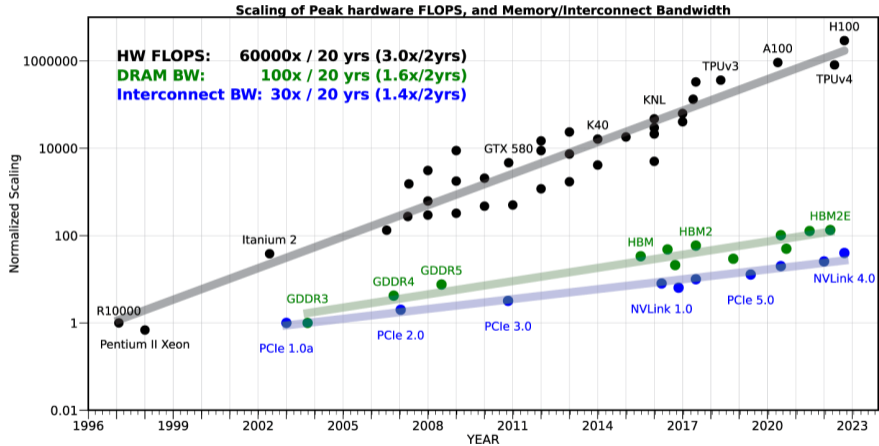
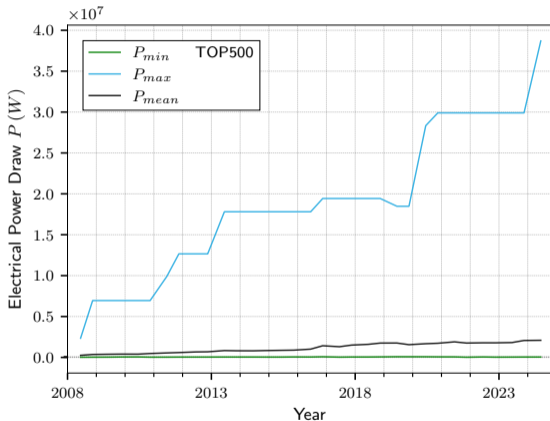


Figure 1: The memory subsystems can not keep up with the development of the floating-point performance. <sup>1</sup>

<sup>1</sup><https://doi.ieeecomputersociety.org/10.1109/MM.2024.3373763/>

## Physical Limits. For Example: Power Draw



- TOP500 likely won't decrease the maximum power usage
- Especially, smaller centers can't keep up
- We can not reach zettascale without one or a few paradigm shifts

Figure 2: The power draw of the TOP500 increased significantly.



## Manufacturing Limits. For Example: Suppliers

High-NA-EUV-Lithography:

Taiwan Semiconductor Manufacturing Company (**TSMC**) relies on  
Advanced Semiconductor Materials Lithography (**ASML**) relies on  
**Zeiss optics** and  
**Trumpf** lasers

## Post-Moore and the Crystal Ball



- CPUs and GPUs will be there for many years.
- Certain algorithms will benefit from specialized accelerators.



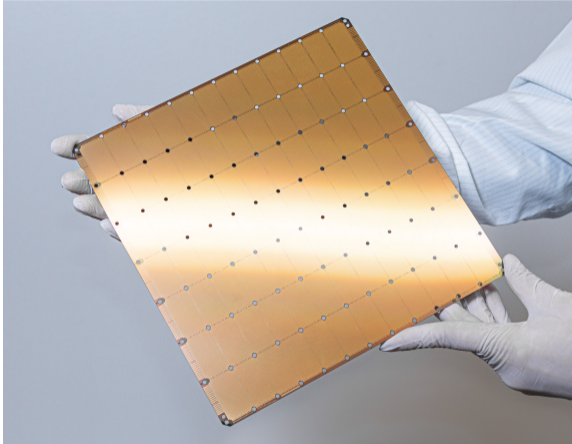


Figure 3: Cerebras Wafer [1]

Important numbers [2]:

- 900 000 processing elements (PEs)
- 125 PetaFLOPs (AI Performance)
- 44 GB on-chip SRAM
- 1,2 Tbit/s system I/O
- 46 250 mm<sup>2</sup> silicon area

Table 2: Performance Comparison for CS-3, H100, and B200. The FP8 and FP16 numbers are peak FLOPS in PetaFLOPS [3].

System	FP8	FP16	Power (kW)	FP8/W	FP8/W/\$	FP16/W	FP16/W/\$
CS-3	250	250	46	5.43	1.09	5.43	1.09
H100	64	32	41.6	1.54	1.10	0.77	0.55
B200	216	108	42.9	5.03	3.35	2.52	1.68
<b>CS-3 Advantage</b>							
vs. H100	3.9x	7.8x	-	3.52x	1.00x	7.05x	2.00x
vs. B200	1.16x	2.31x	-	1.08x	0.32x	2.15x	0.65x



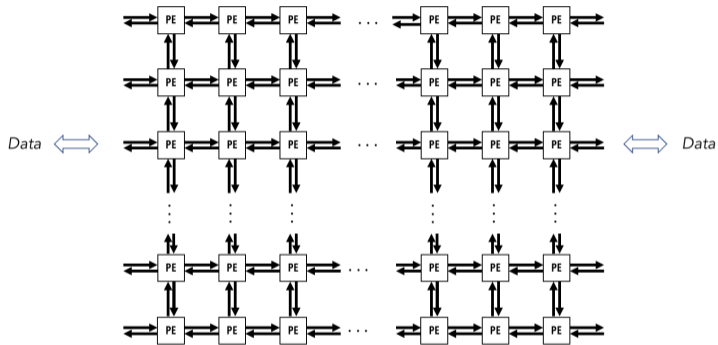


Figure 4: Schematic Wafer Layout (Image by Cerebras)

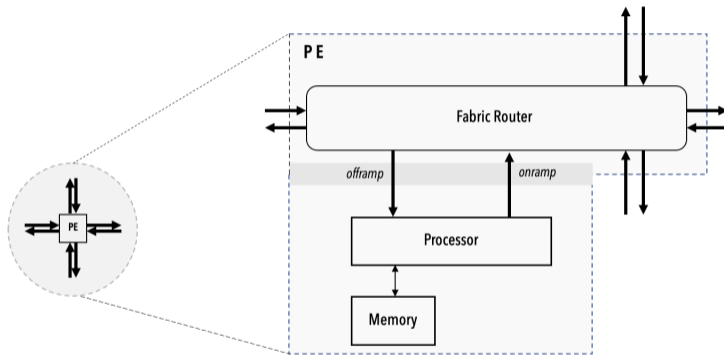


Figure 5: Single PE (Image by Cerebras)

## Algorithmic Design Choices?

- Exploit the core features of the WSE:
  - Massive parallelism
  - Extremely fast on-chip nearest-neighbor communication (non-uniform access time)

⇒ What problem (sizes) can be effectively computed with these constraints?  
⇒ We need to rethink algorithms with these features in mind!



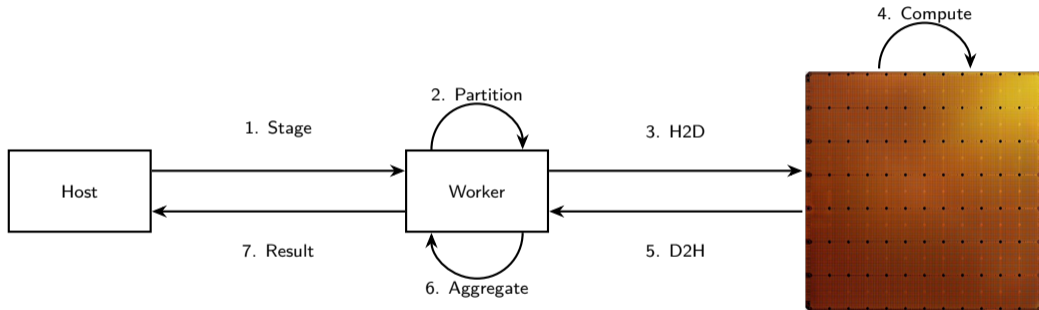


Figure 6: SpMV execution pipeline with the WSE.

## Our Research: Results

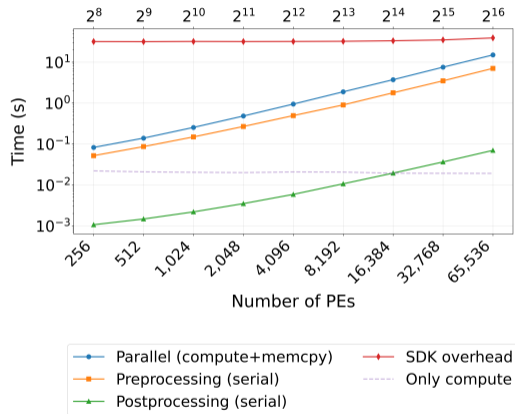


Figure 7: Time breakdown

### Problems:

- Host communication is extremely expensive compared to on-chip communication
  - Data locality on the chip is not exploited
- ⇒ Constraints on the algorithm and problem

# Our Research: Results

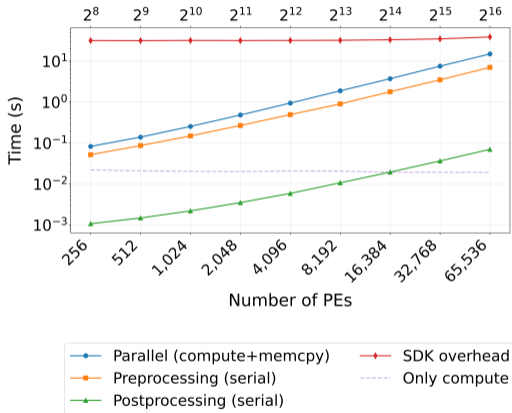


Figure 7: Time breakdown

## Problems:

- Host communication is extremely expensive compared to on-chip communication
  - Data locality on the chip is not exploited
- ⇒ Constraints on the algorithm and problem

HLRIS  
Universität Stuttgart

### SpMV for the Cerebras Wafer-Scale Engine

Daniel Renschler<sup>1,3</sup> Jonathan Schäfer  
Johannes Gebert<sup>1</sup> Mark Parsons<sup>2</sup>

<sup>1</sup>High-Performance Computing Center Stuttgart (HLRIS), University of Stuttgart  
<sup>2</sup>Edinburgh Parallel Computing Centre (EPCC), The University of Edinburgh  
<sup>3</sup>Waltlab, Institute of Technology

---

**Introduction**

Sparse Basic Linear Algebra Subprograms (BLAS) routines are a cornerstone of high-performance scientific computing, providing highly optimized, standardized operations for manipulating sparse matrices and vectors. One particularly important BLAS routine is the Sparse Matrix-Vector Product (SpMV) for a given sparse matrix and a vector.

BLAS routines [1, 2], especially the SpMV, are essential in traditional applications such as finite element analysis (FEA) and computational fluid dynamics (CFD), as well as in applications where large, sparse systems of linear equations dominate the computational workload, including an increasing number of specialized and emerging algorithms.

On the hardware side, today's and tomorrow's high-performance computing architectures are increasingly leaning towards artificial intelligence and machine learning workloads, including systems based on domain-specific accelerators such as those developed by Cerebras and other AI-focused vendors. While these custom are often marketed for AI applications, they also offer unique advantages for traditional scientific computing.

**Results and Discussion**

In the following, we present strong and weak scaling analysis, coupled with our interpretation.

**Results and Discussion (continued)**

Figure 5 shows the results of the conducted weak scaling experiment. We fix the load per PE to 3000 mu, PE, and scale the global problem with the number of PEs. The relevant run in table 1 is "Weak 1".

The left graph shows the absolute run-time of different components.

# Our Research: Results

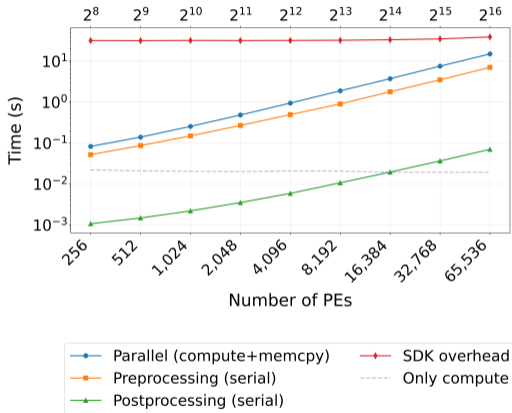


Figure 7: Time breakdown

## Problems:

- Host communication is extremely expensive compared to on-chip communication
  - Data locality on the chip is not exploited
- ⇒ Constraints on the algorithm and problem

Universität Stuttgart

### SpMV for the Cerebras Wafer-Scale Engine

Daniel Renschler<sup>1,3</sup> Jonathan Schäfer  
Johannes Gebert<sup>4</sup> Mark Parsons<sup>2</sup>

<sup>1</sup>High-Performance Computing Center (HPC2), University of Stuttgart  
<sup>2</sup>Edinburgh Parallel Computing Centre (EPCC), The University of Edinburgh  
<sup>3</sup>Waltwhire Institute of Technology

**Introduction**

Sparse Basic Linear Algebra Subprograms (BLAS) routines are a cornerstone of high-performance scientific computing, providing highly optimized, standardized operations for manipulating sparse matrices and vectors. One particularly important BLAS routine is the Sparse Matrix-Vector Product (SpMV) for a given sparse matrix and a vector.

BLAS routines [1, 2], especially the SpMV are essential in traditional applications such as finite element analysis (FEA) and computational fluid dynamics (CFD), as well as in applications where large, sparse systems of linear equations dominate the computational workload, including an increasing number of specialized and emerging algorithms.

On the hardware side, today's and tomorrow's high-performance computing architectures are increasingly leaning towards artificial intelligence and machine learning workloads, including systems based on domain-specific accelerators such as those developed by Cerebras and other AI-focused vendors. While these custom are often marketed for AI applications, they are increasingly being used for traditional scientific workloads.

**Results and Discussion**

In the following, we present strong and weak scaling analysis coupled with our interpretation.

Figure 1: SpMV execution pipeline

**Results and Discussion (continued)**

Figure 5: Weak Scaling: Time Breakdown (MFL Efficiency) right coupled with our interpretation.



Problem constraint: Massive computing power per memory unit  $\Rightarrow$  Algorithms that reuse data multiple times

Algorithm constraint: As much structured data as possible; how to deal with unstructured data?

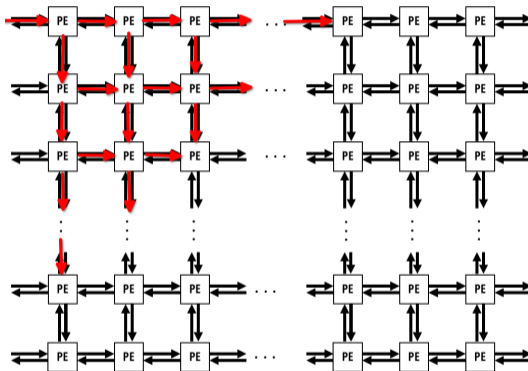


Figure 8: Streaming on the WSE

## More Features to Exploit

H L R I S

Using higher-level languages to ease development!

# More Features to Exploit

Using higher-level languages to ease development!

For example: SPADA [4]

**SPADA: A Spatial Dataflow Architecture Programming Language**

Lukas Gianuzzi Noëda, ETH Zurich Zurich, Switzerland lux@noeda.ai	Tal Ben-Nun Lawrence Livermore National Laboratory Livermore, California, USA	Torsten Hoefler Department of Computer Science ETH Zurich Zurich, Switzerland
--	---	--

**Abstract**—Spatial dataflow architectures like the Cerebras Wafer-Scale Engine achieve exceptional performance in AI and scientific applications by leveraging distributed memory across processing elements (PEs) and localized computation. However, programming these architectures remains challenging due to the need for explicit orchestration of data movement through reconfigurable networks-on-chip and asynchronous computation triggered by data arrival. Existing FPGA and CGRA programming models emphasize loop scheduling but overlook the unique capabilities of spatial dataflow architectures, particularly efficient dataflow over regular grids and intricate routing management. We present SPADA, a programming language that provides precise control over data placement, dataflow patterns, and asynchronous operations while abstracting architecture-specific low-level details. We introduce a rigorous dataflow semantics framework for SPADA that defines routing correctness, data

avoids shared memory contention overhead, achieving memory bandwidth that conventional systems cannot match. However, the absence of shared memory introduces significant programming challenges. Data movement between PEs must be explicitly orchestrated using reconfigurable circuits controlled by routers in the network-on-chip (NoC). Computations are triggered asynchronously as data arrive requiring careful synchronization across the distributed fabric. Furthermore, the circuit-switched network provides only a limited number of concurrent communication channels, necessitating explicit channel assignment to avoid routing conflicts. As a result, programming SDAs for large-scale applications becomes time-intensive and error-prone. Routing conflicts at data rates manifest non-deterministically at runtime across

DCJ 12 Nov 2025

- Simulation of physical systems, e.g., engineering, based in Euclidean space. [5]
- Exploiting data locality, for example, by compiler features. [6]

- [1] Cerebras Systems. *Press Kit*. <https://www.cerebras.ai/company/press-kit>. 2024.
- [2] S Lie. *Cerebras Wafer-Scale AI*. Hot Chips 36, [https://hc2024.hotchips.org/assets/program/conference/day2/72\\_HC2024.Cerebras.Sean.v03.final.pdf](https://hc2024.hotchips.org/assets/program/conference/day2/72_HC2024.Cerebras.Sean.v03.final.pdf). 2024.
- [3] Y Kundu, M Kaur, T Wig, K Kumar, P Kumari, V Puri, and M Arora. “A comparison of the cerebras wafer-scale integration technology with nvidia gpu-based systems for artificial intelligence”. In: *arXiv preprint arXiv:2503.11698* (2025).
- [4] L Gianinazzi, T Ben-Nun, and T Hoefler. “SPADA: A Spatial Dataflow Architecture Programming Language”. In: *arXiv preprint arXiv:2511.09447* (2025).
- [5] T Opielstrup, N Giambianco, DZ Kalchev, I Sharapov, M Taylor, D Van Essendelft, S Rajamanickam, and M James. “Beyond Exascale: Dataflow Domain Translation on a Cerebras Cluster”. In: (Nov. 2025). DOI: [10.48550/ARXIV.2511.11542](https://doi.org/10.48550/ARXIV.2511.11542). arXiv: [2511.11542](https://arxiv.org/abs/2511.11542) [cs.DC].
- [6] N Stawinoga, D Katz, A Lydike, J Zarins, N Brown, G Bisbas, and T Grosser. “An MLIR Lowering Pipeline for Stencils at Wafer-Scale”. In: (Jan. 2026). DOI: [10.1145/3779212.3790124](https://doi.org/10.1145/3779212.3790124). arXiv: [2601.17754](https://arxiv.org/abs/2601.17754) [cs.DC].