# Breaking Performance Barriers

RISC-V Summit 2025. Paris, France

Graham Wilson – Product Group

**Processor IP company (RISC-V)**

**Proven Performance Legacy**

Team comprised of veterans from ThunderX2 ARM server chip development, Oracle, MIPS, Intel, Google

KLEINER PERKINS

Mayfield

Fidelity
INVESTMENTS

Akeana's drive is to enable our SoC customers to success with leading edge performance processor system IP

# Akeana Processor IP Product Lines

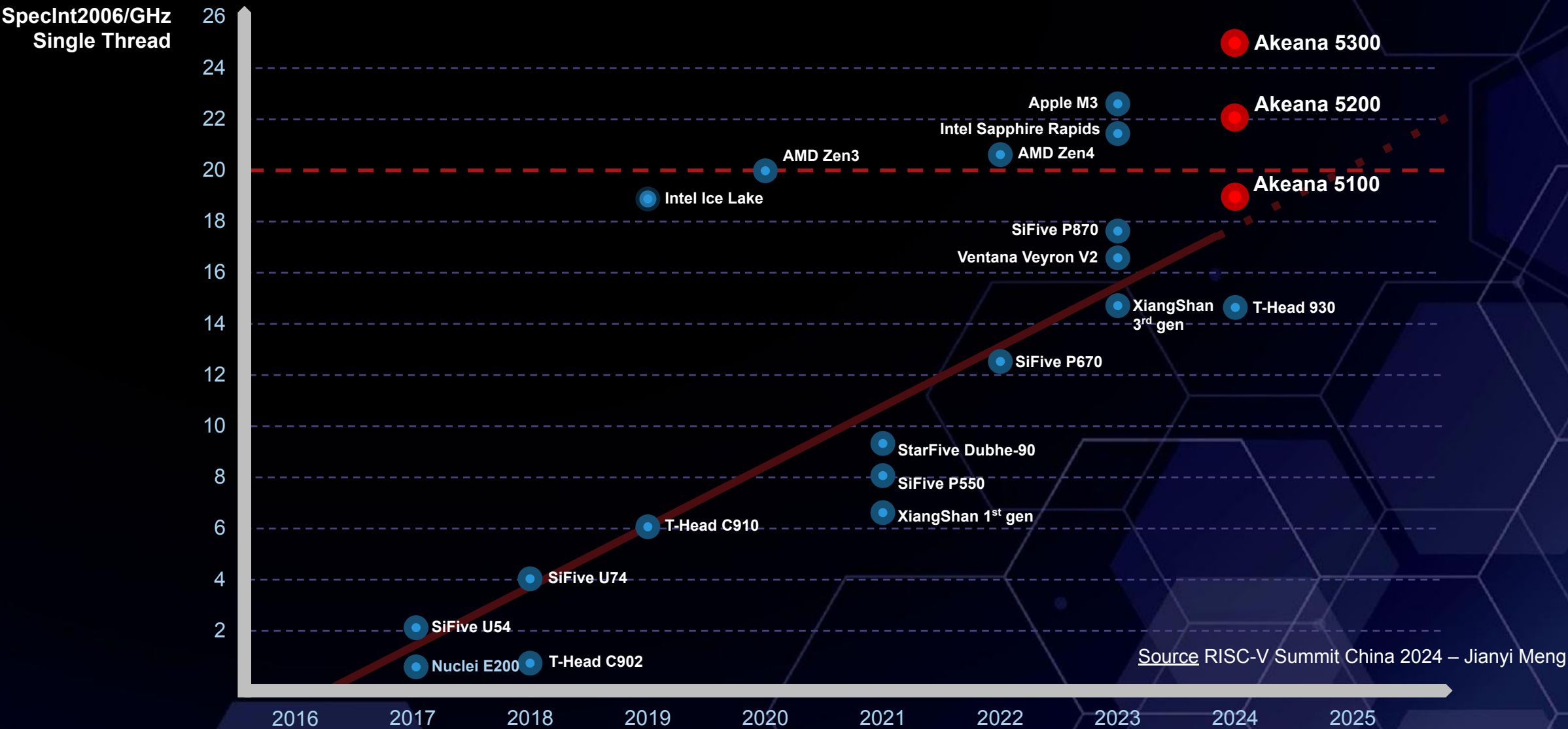| **Akeana 100 Series**<br>**Micro-controller, Embedded** | 32-bit, PMP. 32-bit physical addressing<br>Single to Dual issue In-Order architecture<br>4-stage to 9-stage pipeline<br>Private L1. Shared Coherent L2 Cache<br>Local ICCM, DCCM | **Equivalent to ARM Cortex-M and Cortex-R** |
|---|---|---|
| **Akeana 1000 Series**<br>**Consumer, Automotive** | 64-bit, MMU. Up to 57-bit virtual addressing<br>2- to 4-wide issue In-Order architecture,<br>  with 5-stage to 9-stage pipeline<br>12-stage Out-of-Order architecture (2,3,4-issue)<br>Private L1, L2 Caching. Shared Coherent L3 Cache<br>Vector Extension (up to 2048 bit). AI Acceleration<br>Multi-Threaded support (up to 4 threads) | **Equivalent to ARM Cortex-A** |
| **Akeana 5000 Series**<br>**Mobile / Data Center, Ultra Performance** | 64-bit, MMU. Up to 57-bit virtual addressing<br>6-wide to 10-wide issue Out-of-Order architecture<br>12-stage pipeline<br>Private L1, L2 Caching. Shared Coherent L3 Cache<br>Vector Extension (up to 512-bit). AI Acceleration<br>Multi-Threaded support (up to 4 threads) | **Equivalent to ARM Neoverse N2, N3** |

# Akeana, Leader in Core Performance



SpecInt2006/GHz Single Thread

- Akeana 5300
- Akeana 5200
- Akeana 5100
- Apple M3
- Intel Sapphire Rapids
- AMD Zen4
- AMD Zen3
- Intel Ice Lake
- SiFive P870
- Ventana Veyron V2
- XiangShan 3rd gen
- T-Head 930
- SiFive P670
- StarFive Dubhe-90
- SiFive P550
- XiangShan 1st gen
- T-Head C910
- SiFive U74
- SiFive U54
- Nuclei E200
- T-Head C902

Source RISC-V Summit China 2024 – Jianyi Meng

2016  2017  2018  2019  2020  2021  2022  2023  2024  2025

# AI CPU Compute Evolving

- NVIDIA shifted towards a customizable CPU implementation to achieve required computation performance for AI / HPC

- To achieve required performance in CPU compute array system, Simultaneous Multi-Threading (SMT) has been used
  - NVIDIA has planted the SMT flag for AI CPU compute

- AI SoC developers are recognizing this push to higher performance in the CPU Compute array

- Akeana is the leading provider for Multi-core, Multi-Threaded, Coherent Infrastructure IP
  - Supported in Akeana 1000 series(In-Order), 5000 series(Out-of-Order), up to 4 threads



NVIDIA'S UPCOMING AI CHIP FAMILY TO REVOLUTIONIZE DATA CENTER PERFORMANCE

Vera Rubin NVL144
Second Half 2026

3.6 EF FP4 Inference
1.2 EF FP8 Training
3.3X GB300 NVL72

13 TB/s HBM4
75 TB Fast Memory
1.6X

260 TB/s NVLink6
2X

28.8 TB/s CX9
2X

Vera
88 Custom Arm Cores
176 Threads
1.8 TB/s NVLink-C2C

Rubin

2   Sized GPUs
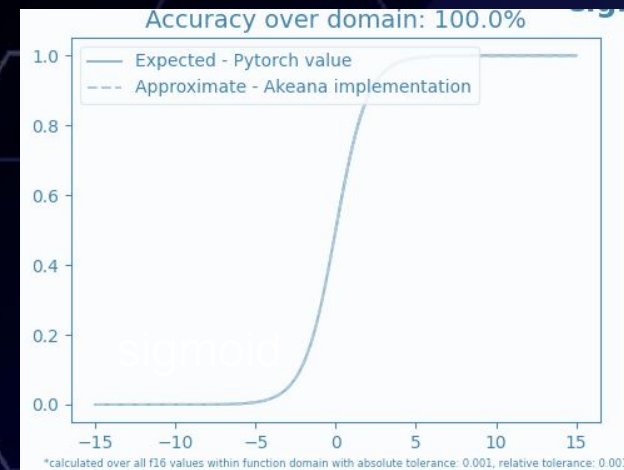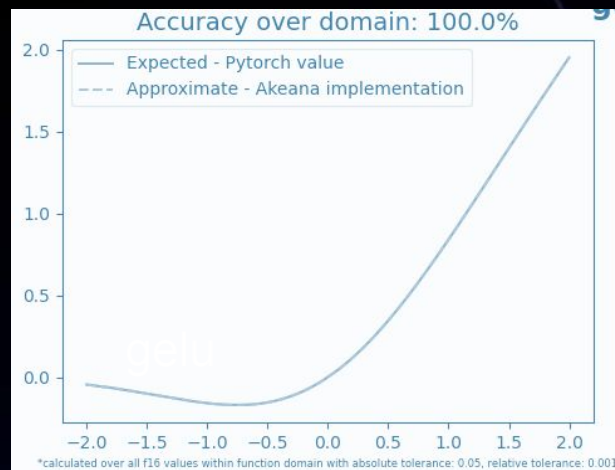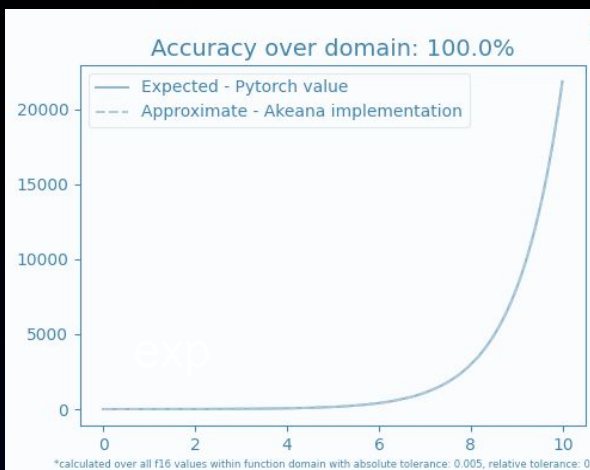50P   2   GB HBM4

# Performance Boost with SMT

| DataBase Processing | | 1 thread | 2 threads | 4 threads |
|---|---|---|---|---|
| | Performance Increase * | 1.00 | 1.79x | 2.25x |

| SpecINT Industry Standard Testbench | | 1 thread | 2 threads | 4 threads |
|---|---|---|---|---|
| | Performance Increase * | 1.00 | 1.18x | 1.28x |

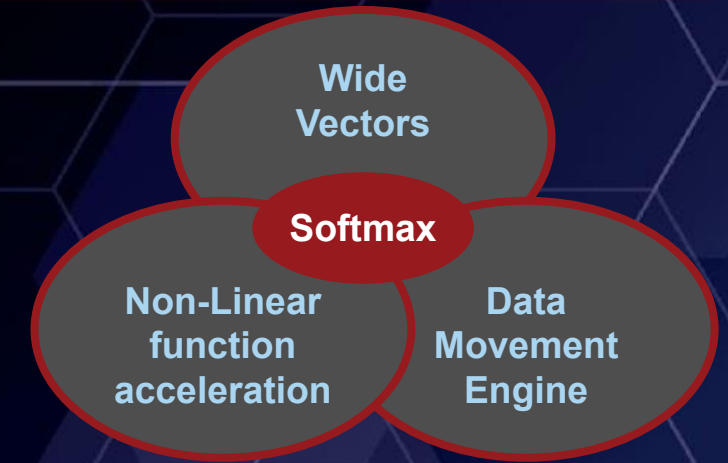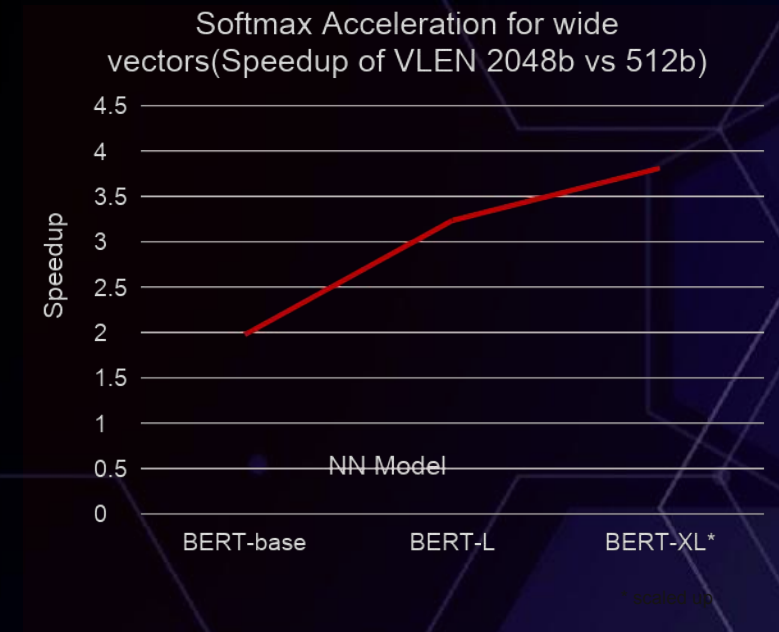* Numbers based upon current results, subject to change

# Performance Data Compute Enabled

- Akeana AI Nonlinear acceleration instructions provide >$10x$ performance increase (example FP16 datatype with sigmoid, gelu, tanh, exp functions)

  - Benefit of lower core count needed, reducing area and power consumed

  - Supported within Akeana AI Performance library functions

- Implemented with RISC-V Vector Extensions up to 2048-bit VLEN for further data compute acceleration
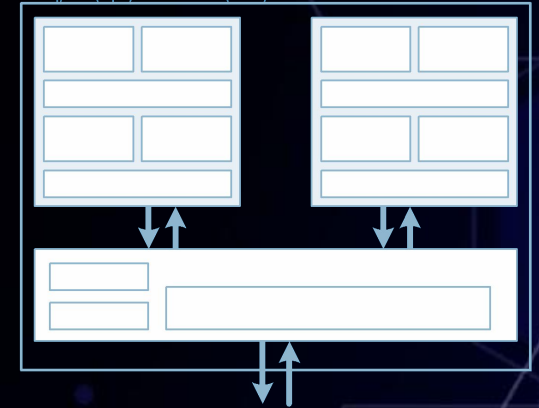
# Accelerating Softmax

- Softmax needs acceleration in 3 domains; vectors, Nonlinear functions and data movement

    - Akeana data movement engine IP available to further accelerate Softmax implementations

- Example Softmax implementation running through various vector length computation

    - Auto-vectorization compiler able to efficiently map vectorizable code to Akeana vector cores

    - Ability to map over range of Transformer based models



Softmax Acceleration for wide vectors(Speedup of VLEN 2048b vs 512b)

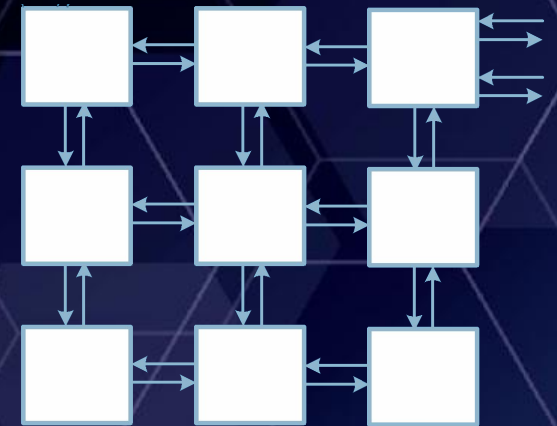# Performance through Scalability

## Single Coherent Cluster

- Utilizes Compute Coherent Block (CCB)
- Shared Coherent Cache, accessed in parallel by all cores
- Up to 8-cores, programmable engines, coherent operation
- Easy scalability to coherent 8-core system
- AI accelerator (GEMM), Customers hardware engines
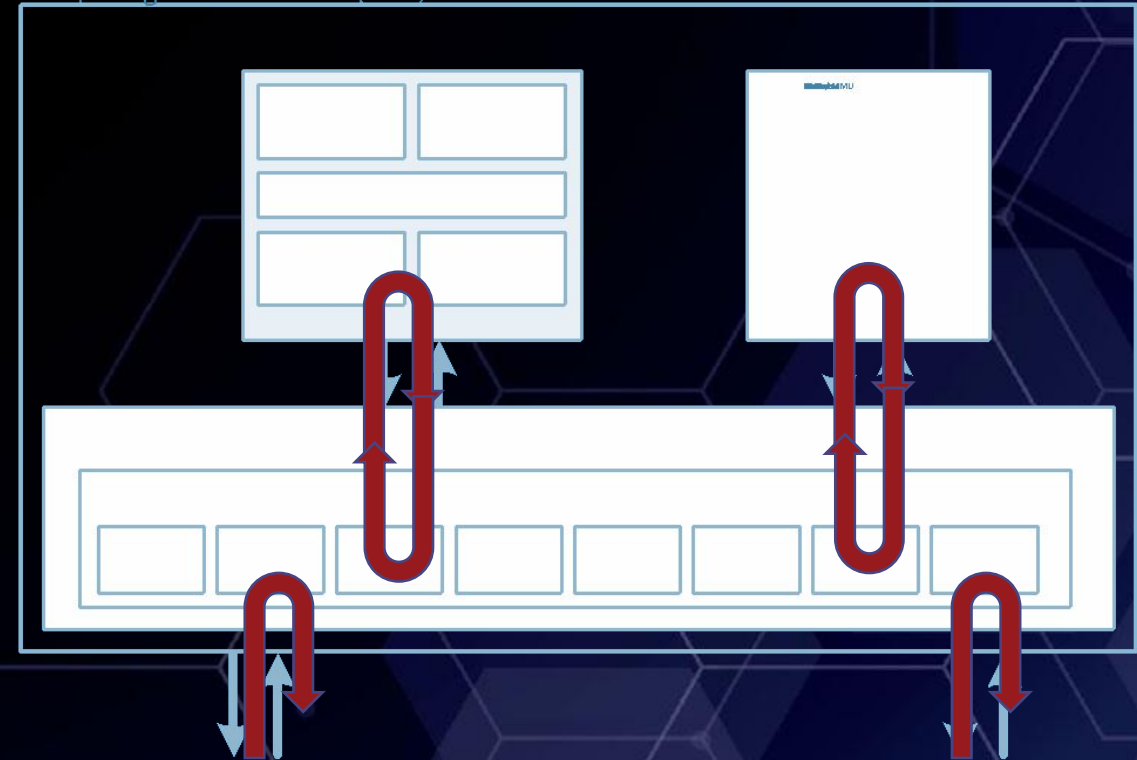
## Scalable Coherent Mesh (Akeana Mesh)

- Utilizes multiple Akeana IP blocks to build up 2D Mesh array
- Single Coherent Cluster (CCB) integrated into 2D Mesh
- AMBA CHI compliant
- Can be built to connect up to 100's of cores, fully coherent
- Akeana provides all the IP blocks, and works with customers to build these larger 2D Mesh coherent interconnect systems

# Optimized Data Movement Performance

- Shared memory can be partitioned into separately accessible banks
  - Allows parallel accesses from multiple processors, Matrix Engines
  - Ping-Pong-ing of Banks when processing of large amounts of data between cores
- 2$^{nd}$ external port to AI CCB block
  - Allows dedicated high bandwidth Matrix Data (Activation and Weights) to be pulled into Shared L2 without blocking from Processor accesses

**Performance** CPU Compute

**Performance** Data Compute

**Performance** Scalable Multi-core

**Performance** Data Movement

# Thank you



AKEANA.COM