



RISC-V for Intelligent Edge Application Processing

Charlie Su, Ph.D.
CTO and President, Andes

RISC-V Summit China, 2024/08/21-23



RISC-V Market Research by SHD Group

■ Total RISC-V SoC Market

- \$6.1B in 2023: 276.8% growth over 2022
- \$92.7B by 2030: a CAGR of 47.4%

■ Andes has >30% of RISC-V IP Market

- Applications: Endpoint, Edge, Cloud

Mobile

Performance, code size

NOVATEK **MEDIATEK**

N25F, N45

MPU/MCU/AIoT

Renesas
HPMicro
Kneron
Telink
Internet Company

D25F, D45, AX25MP, AX45MP

Storage

PHISON

Performance, bandwidth, real-time

D23, N25F, N45, AX45MP

5G Networks

EDGE
picocom
Empowering Wireless

N25F, A25, A45MP, AX45MP

CERTIFICATE
N25F-SE

Display & Touch
ILITEK

Auto MCU

Auto DVR Cam

In-Cabin Radar

CMOS Sensor
MetaSicon

N25F-SE, AX45MPV

Large-Scale AI/ML

SRAM CIM
CNN
Photonics
Mamba
Transformer

LIGHTelligence

TetraMem
Accelerate the World®

SK telecom

AXELERA
ARTIFICIAL INTELLIGENCE

RAIN

STREAM COMPUTING

Acceleration, bandwidth, extensibility

NX27V, AX25, AX45MPV, AX65

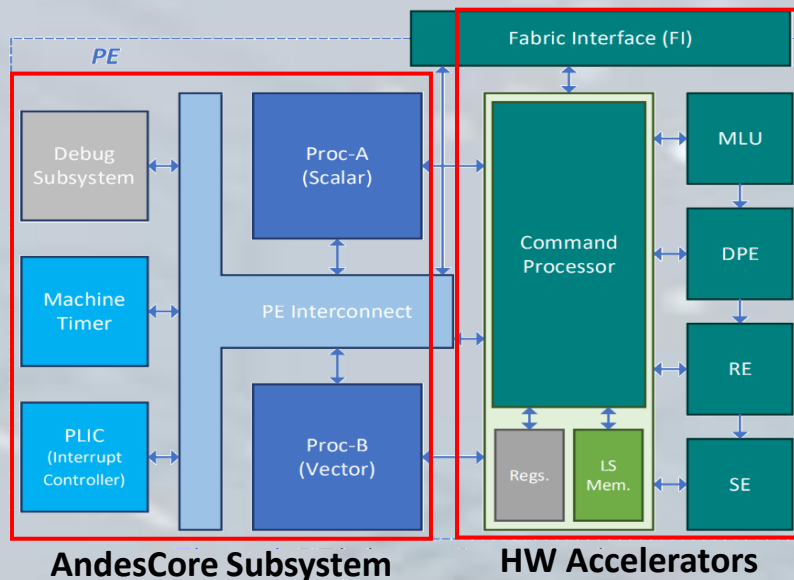
MTIA: Meta Training/Inference Accelerator



- ISCA 2023 paper, “MTIA: First Generation Silicon Targeting Meta’s Recommendation Systems”
- Proc-A/B: AX25-V100, an early version of the popular **NX27V vector processor**
- **Custom extensions (ACE)**: create new interfaces/instructions/registers
- [Next generation MTIA](#) (Meta blog): *Serving models in production*



All photos: courtesy of ACM



RISC-V Enables Innovations for Large-Scale AI/ML



Based on NX27V/AX45MPV (RVV), AX65 (OOO), and AX25

With Andes Automated Custom Extensions™ (ACE)

■ AI Accelerators Using SRAM-based Compute-In-Memory:



■ AI Accelerator Using Photonics



■ AI Accelerator for Cloud Service



■ AI Accelerator for ADAS

Intelligent Edge Application Processing (IEAP)



- **Market projection: Edge AI will be bigger than Cloud AI**
- **Definition of IEAP:**
 - Not Cloud AI/ML, Not AI PC
 - Including application handling and AI/ML accelerations
- **AI/ML for Edge Applications: On-device Intelligence**
 - Intelligent Edge Servers with domain knowledge (e.g. legal, medical, or enterprise)
 - Smart Factory/Healthcare/Transportation: monitoring and predictive maintenance
 - Network Systems: performance optimization, security enhancement, traffic management
 - Automotive: ADAS, Self Driving Vehicles
 - Robots, Surveillance, and more

Required Architecture Support



■ AI/ML speedup:

● Matrix Multiplication:

- Hardwired solution (APU/NPU/TPU): most efficient and powerful
- Matrix Instructions: most flexible
 - IME (Integrated Matrix Extension) Task Group: leveraging VRF/MAC's in VPU

● Non-linear functions: softmax, sigmoid, sin/cos, GeLu/SiLu

● General compute instruction extensions: Vector Extension (RVV), SIMD/DSP (RVP)

■ Rich application enablement:

● Profiles: RVA20/22/23



AndesAIRE™ End-to-End AI/ML Solutions

Andes AI Runs Everywhere



PyTorch ONNX TensorFlow Lite TensorFlow

AndeSight™ IDE

- GCC/LLVM Toolchains
- **Vector/DSP Library**
- Build, debug, deploy, profile
- Analysis and tuning
- RTOS & Linux
- Device drivers
- Sample codes
- Simulator
- Documentation

AndesClarity™

Pipeline Optimizer

ACE/COPILOT

NN models

AndesAIRE™ Software

AndesAIRE™ NNPIlot™

- Graph-level optimization (Pruning/Quantization)
- Backend-aware optimization (Fusion/Tensor Allocation)

Generated TFL Models

TensorFlow
Lite



Generated C Template

- NN Library API
- AnDLA driver and runtime
- AnDLA command image

AndesAIRE™
XNNPACK

AndesAIRE™
NN Library

AI Compilers



SoC

Linux/RTOS Host

AX45MP, AX65

Compute Acceleration

Vector: 27V, 45V

DSP/SIMD: D23, D25F, D45

**Domain
Extensions**

ACE

Accelerator

AnDLA™ I350

Andes Vector Processor Families



More features for AX47MPV

- . CHI-based multi-cluster
- . Latest RVA profile
- . More FP data formats

AX25-V100
(Meta MTIA-1)

NX27V

AX45MPV

AX46MPV*

AX47MPV*

more features

Integrated Matrix Ext. (IME)

8-core cluster

16-core cluster with private L1/L2

HVM (High-speed Vector Memory)

ACE for RVV

ACE for RVV (advanced)

ACE (Andes Automated Custom Extensions™)

int4~64, fp16~64; bf16 (conv.)

+bf16 (full arithmetic)

+fp8

VLEN 128~512

+1024

+2048

RVV 0.8

RVV 1.0

5-stage single-issue

8-stage dual-issue

* Future products subject to change

AX45MPV and ACE: RAIN AI's Adoption*

Leveraging Andes Custom Extensions (ACE) **automated**

In-memory Compute Integration

ACE Streaming Port

Introduce 2nd vector load/store unit; direct access to VRF

Leverage custom addressing and address control registers to control a plurality of in-memory compute blocks with different data sizes

Non-linear Operator Integration

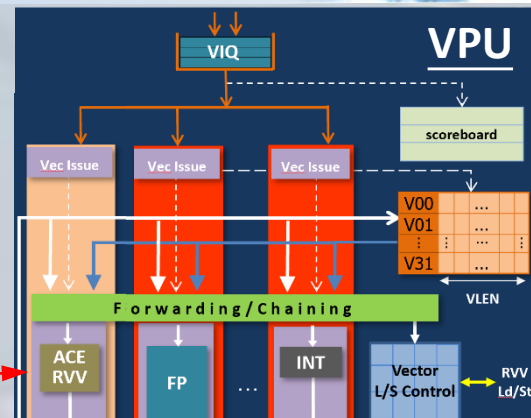
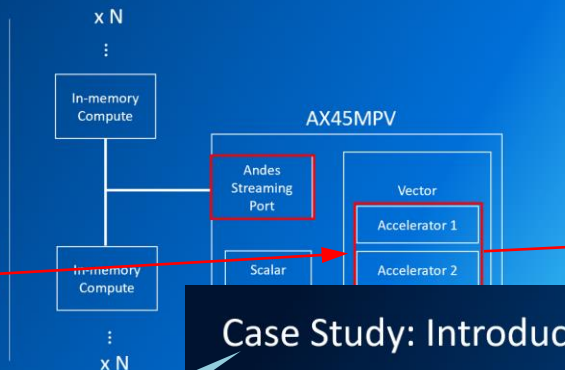
ACE RVV

Integrate non-linear approximations as a functional unit in vector pipeline with direct access to VRF

Introduce vector instructions that look and feel familiar

Support for LMUL, pipelined multi cycle operations

Also added **Softmax, SiLu,**
and other instructions
using ACE



Case Study: Introducing a Sigmoid Instruction

```
rvv_insn ace_sigmoid {  
  op = {out vr:vr result, in vr:vr x};  
  vector_unit = clen;  
  csim = %c{  
  
  }  
  %};  
  latency = 6;  
  cycle_per_result = 1; // pipelined  
};
```

.ace

.v

COPILLOT

Updated core RTL with new vector unit

AndeSysC model with support for new instruction

C Wrappers (Intrinsics)

Compiler collateral

```
static inline vfloat16m4_t __attribute__((always_inline))  
ace_sigmoid_f16m4(vfloat16m4_t x, size_t vl){  
  vfloat16m4_t result_RetVar;  
  vsetvl_e16m4(vl);  
  __asm__ __volatile__ {  
    "ace_sigmoid %[result], %[x], %[vm]"  
    : (result) "=dvr" (result_RetVar)  
    : [x] "v" (x), [vm] "i" (1)  
    : "memory"  
  };  
  return result_RetVar;  
}  
  
static inline vfloat16m4_t __attribute__((always_inline))  
ace_sigmoid_f16m4(vfloat16m4_t x, size_t vl){  
  vfloat16m4_t result_RetVar;  
  vsetvl_e16m4(vl);  
  __asm__ __volatile__ {  
    "ace_sigmoid %[result], %[x], %[vm]"  
    : (result) "=dvr" (result_RetVar)  
    : [x] "v" (x), [vm] "i" (1)  
    : "memory"  
  };  
  return result_RetVar;  
}
```

Sample wrappers for our sigmoid instruction

* Presented in Andes RISC-V CON San Jose, June 2024

https://www.bilibili.com/video/BV1eDWpe6EcZ/?vd_source=b34de578c60c056a381185399b7b80ad

https://youtube.com/live/TpGFCTu_OFw

Andes Application Processor Families



AX66 VPU:

- has 2 pipes shared by FP
- execute 2 ALU/load/store

Android
Base

AX65

Best-Balanced

8.25 specint2k6/GHz

RVA22+

AX66*
Advanced

>10 specint2k6/GHz

RVA23

V/VK (VLEN=128)

Hypervisor + AIA + (IOMMU + IOPMP)

Private L1/L2 Caches, CHI Multi-Cluster Coherency

AX67*

Most performant

>11 specint2k6/GHz

RVA24

further perf boost

V/VK (VLEN up to 512)

AX63 customer-driven
Power-optimized

>7.0 specint2k6/GHz

13-Stage, Out-of-Order Execution, Multicore Coherency, Linux-Capable, Up to 8 Cores/Cluster

Roadmap for AX60 Series

** Roadmap/features: subject to change. Future performance: estimated.*

Andes Application Processor Families



				Cuzco Scalable performance 15~20 specint2k6/GHz
AX63 <small>customer-driven</small> Power-optimized >7.0 specint2k6/GHz	AX65 Best-Balanced 8.25 specint2k6/GHz	AX66 Advanced >10 specint2k6/GHz	AX67 Most performant >11 specint2k6/GHz	
		RVA23	RVA24	RVA24
		V/VK (VLEN=128)	further perf boost	Private L1/L2, Shared L3
		Hypervisor + AIA + (IOMMU + IOPMP)	V/VK (VLEN up to 512)	Vector/Vector Crypto
	RVA22+	Private L1/L2, CHI Multi-Cluster Coherency	8-Core Cluster with CHI	14-stage 8-way/6-way OOO with Patented Time-Based Scheduling
13-Stage, Out-of-Order Execution, Multicore Coherency, Linux-Capable, Up to 8 Cores/Cluster				
Roadmap for AX60 Series				Cuzco Series

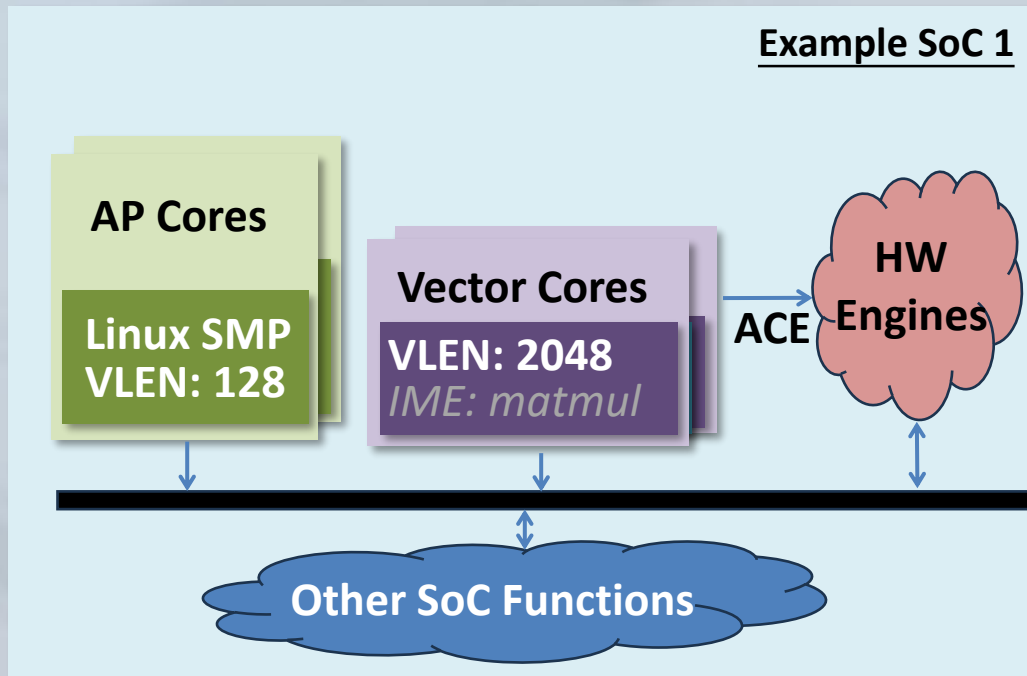
Roadmap/features: subject to change. Future performance: estimated.

Choices of SoC Architecture



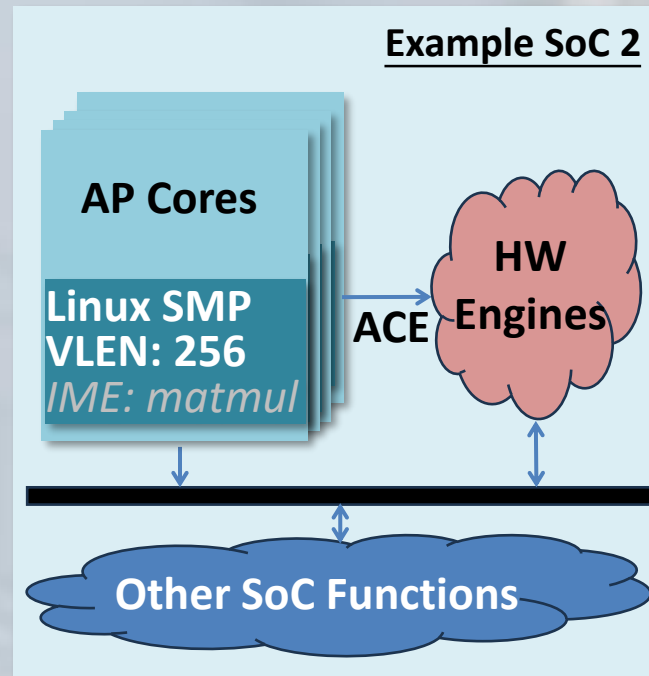
- Comparison: performance, area efficiency (AE), ease of programming/optimization

Example SoC 1



(perf, AE, programming)= (high, mid, low)

Example SoC 2



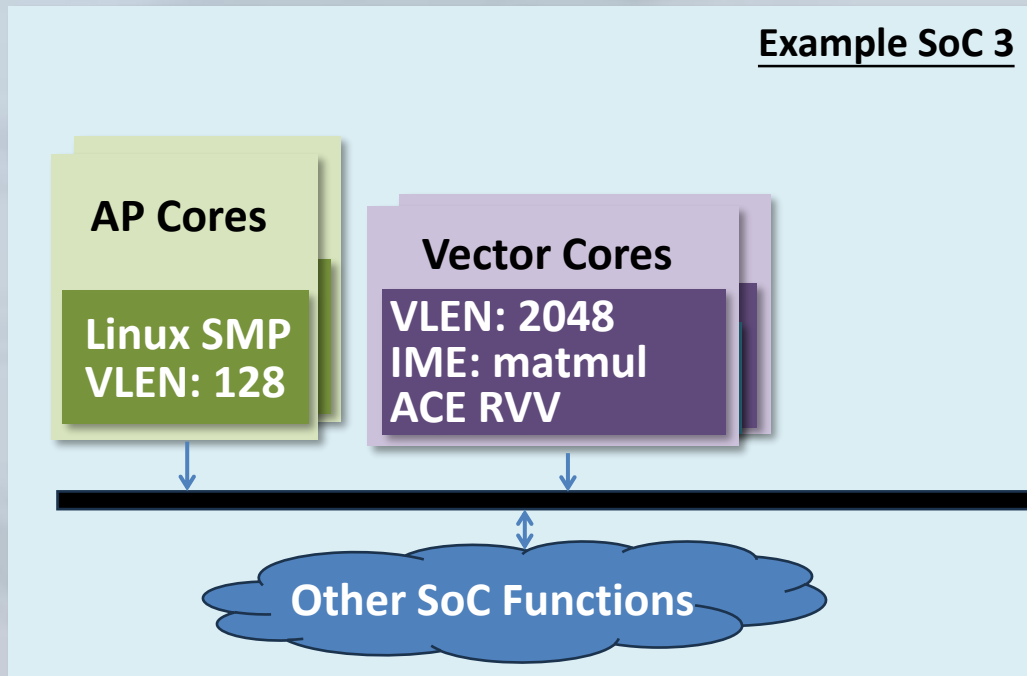
(perf, AE, programming)= (high, low, mid)

Choices of SoC Architecture



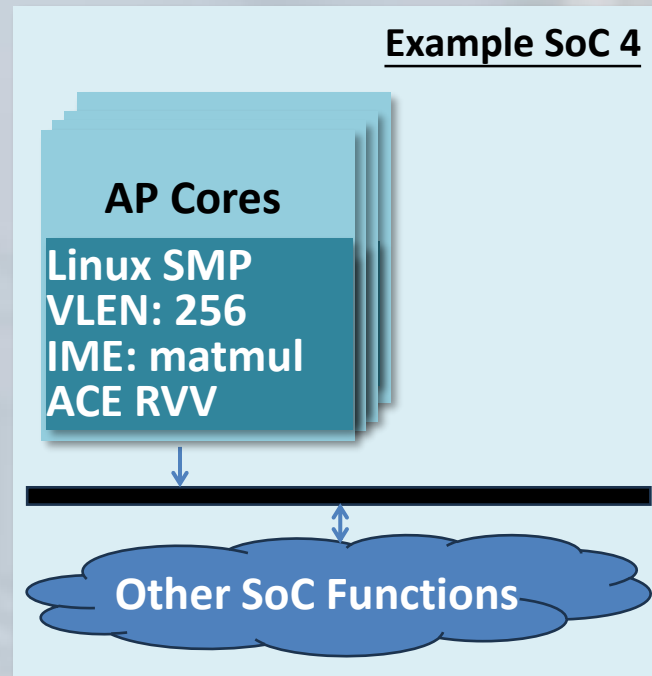
- Comparison: performance, area efficiency (AE), ease of programming/optimization

Example SoC 3



(perf, AE, programming) = (mid, high, mid)

Example SoC 4



(perf, AE, programming) = (low, mid, high)

Concluding Remarks



■ Intelligent Edge Applications

■ RISC-V Architecture Advantages

- Vector/Matrix processing: RVV, IME
- Custom Extensions
- RVA20/22/23: Linux/Android
- *Capacity/Bandwidth Management*

■ Andes Offerings

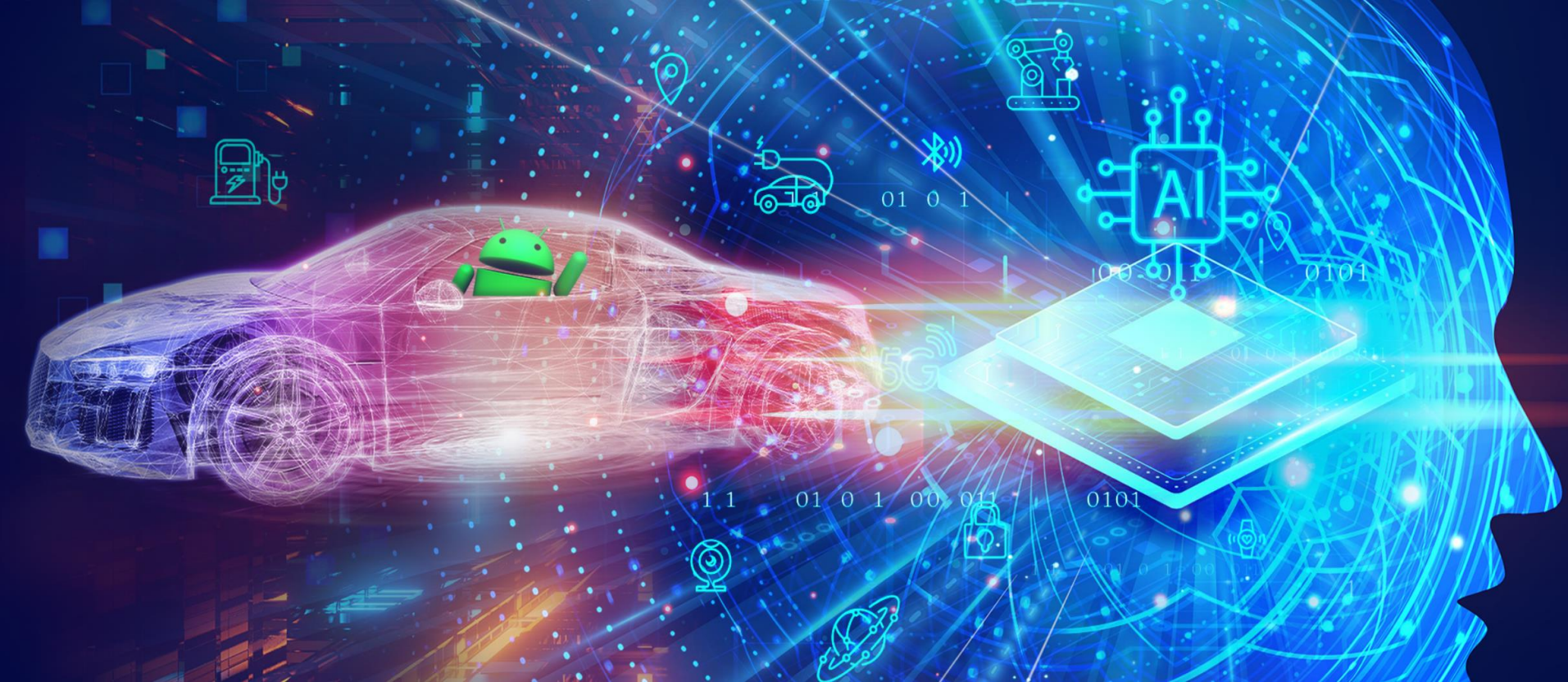
- Vector cores: 40-Series
- Linux cores: 60-Series, and Cuzco
- *AnDLA I350: hardwired engine*

More features for AX47MPV <ul style="list-style-type: none">· CHI-based multi-cluster· RVA23· fp8		AX45MPV	AX46MPV	AX47MPV
		8-core cluster	Integrated Matrix Ext. (IME)	
		HVM (High-speed Vector Memory)		
AX25-V100	NX27V	ACE for RVV	ACE for RVV (advanced)	
ACE (Andes Automated Custom Extension™)				
int4~64, fp16~64; bf16 (conv.)		+bf16 (full arithmetic)		+fp8
VLEN 128~512		+1024	+2048	
RVV 0.8		RVV 1.0		
5-stage single-issue		8-stage dual-issue		

* Future products subject to change

		AX67 Most performant >11 specint2k6/GHz		Cuzco Scalable performance 15~20 specint2k6/GHz
		AX66 Advanced >10 specint2k6/GHz	RVA24	RVA24
		RVA23	further perf boost	Private L1/L2, Shared L3
		V/VK (VLEN=128)	V/VK (VLEN up to 512)	Vector/Vector Crypto
		Hypervisor + AIA + (IOMMU + IOPMP)	8-Core Cluster with CHI	
		Private L1/L2, CHI Multi-Cluster Coherency	14-stage 8-way/6-way OOO with Patented Time-Based Scheduling	
13-Stage, Out-of-Order Execution, Multicore Coherency, Linux-Capable, Up to 8 Cores/Cluster				
Roadmap for AX60 Series				Cuzco Series

Roadmap/features: subject to change. Future performance: estimated.



Thank You !!