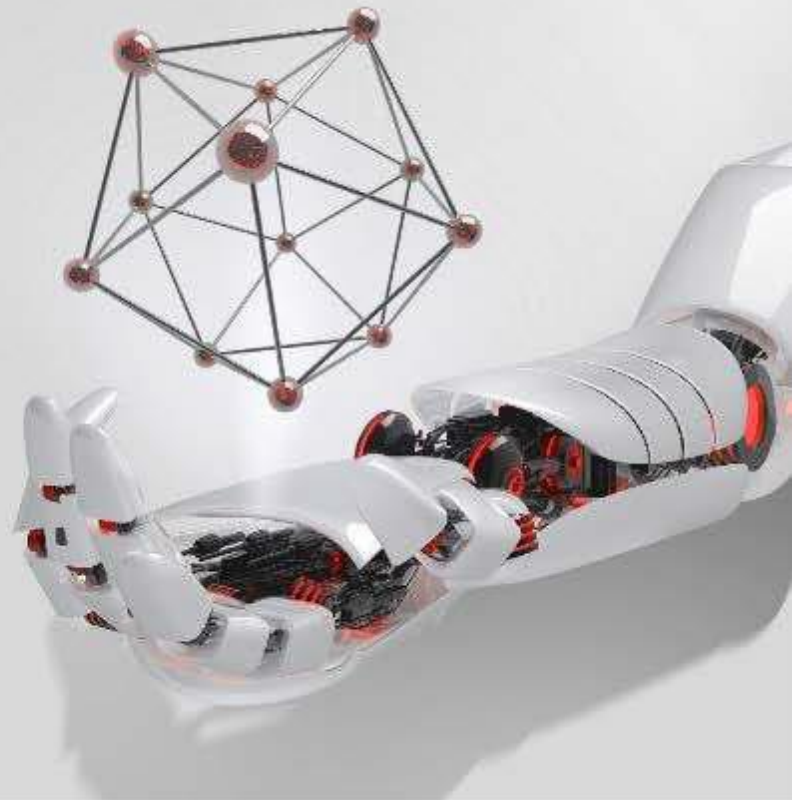
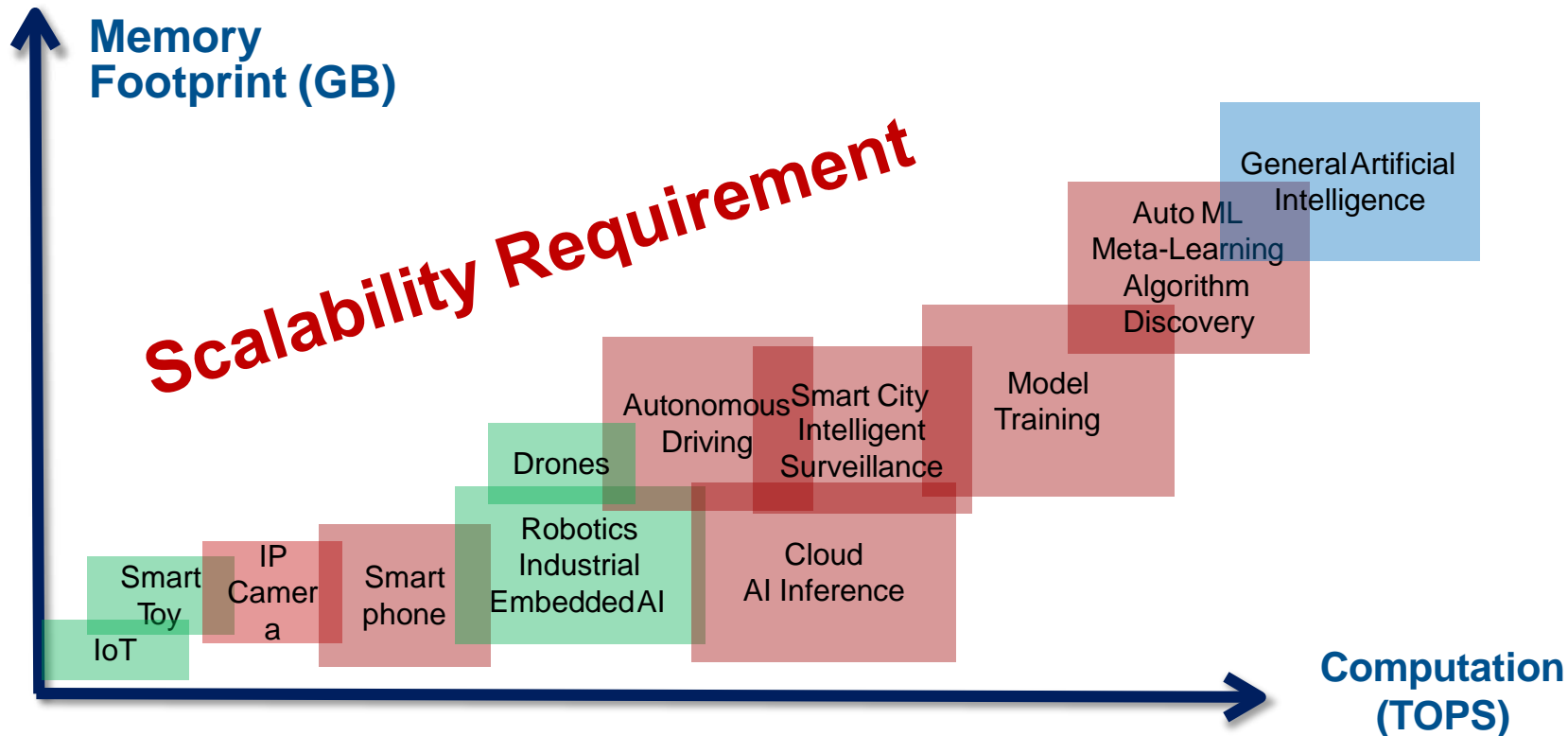


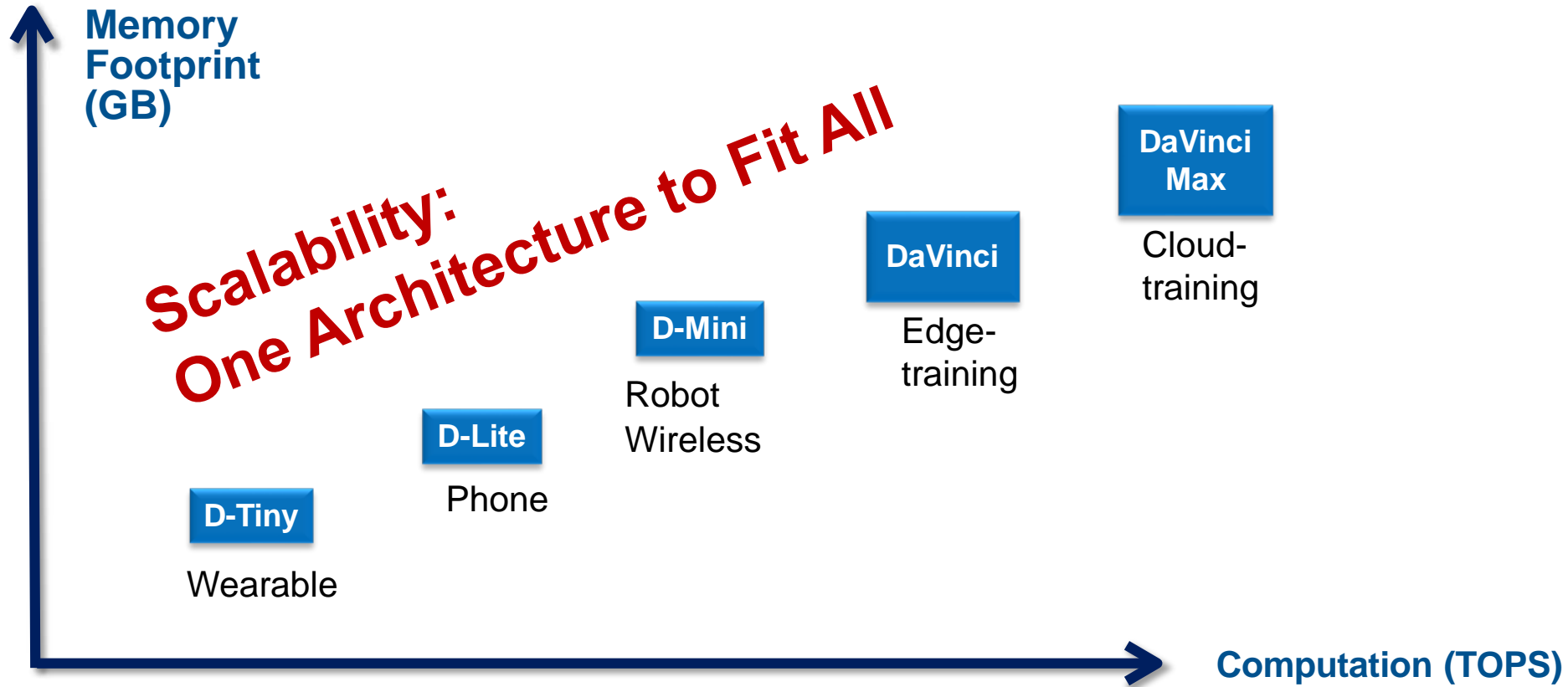
# DaVinci: A Scalable Architecture for Neural Network Computing



# AI Computation for All Scenarios

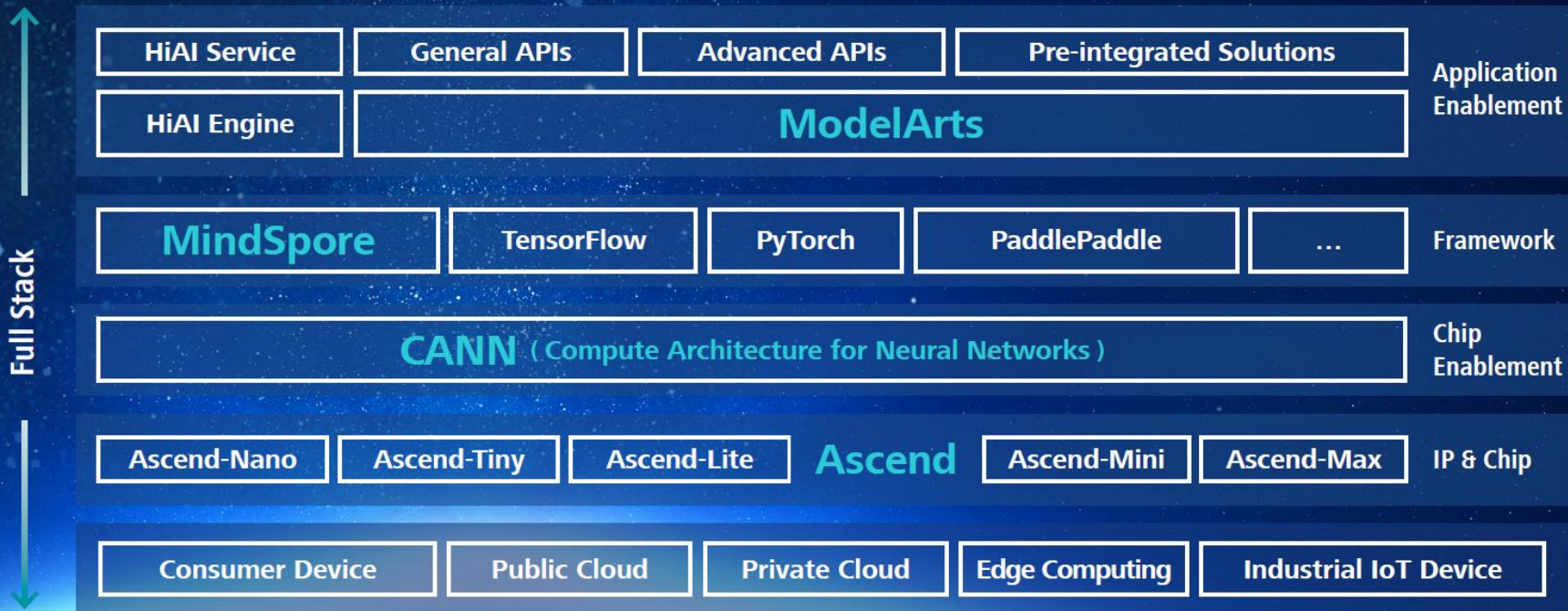


# Huawei AI Processor for All Scenarios



# Huawei's AI Portfolio

## AI Application



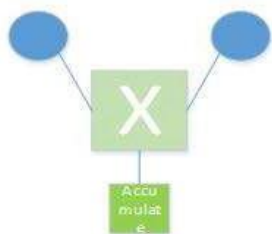
All-scenario

# Architecture Overview of DaVinci

# Building Blocks and their Computation Intensity

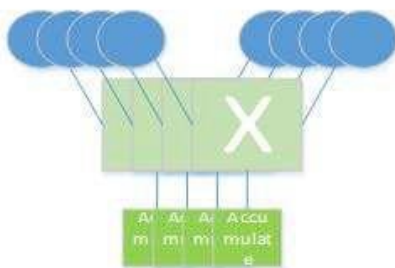
## 1D Scalar Unit

Full flexibility



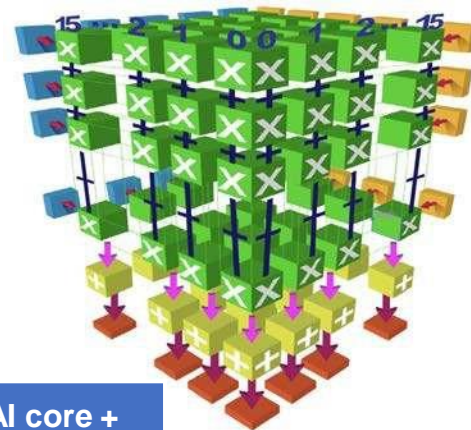
## + 2D Vector Unit

Rich & efficient operations



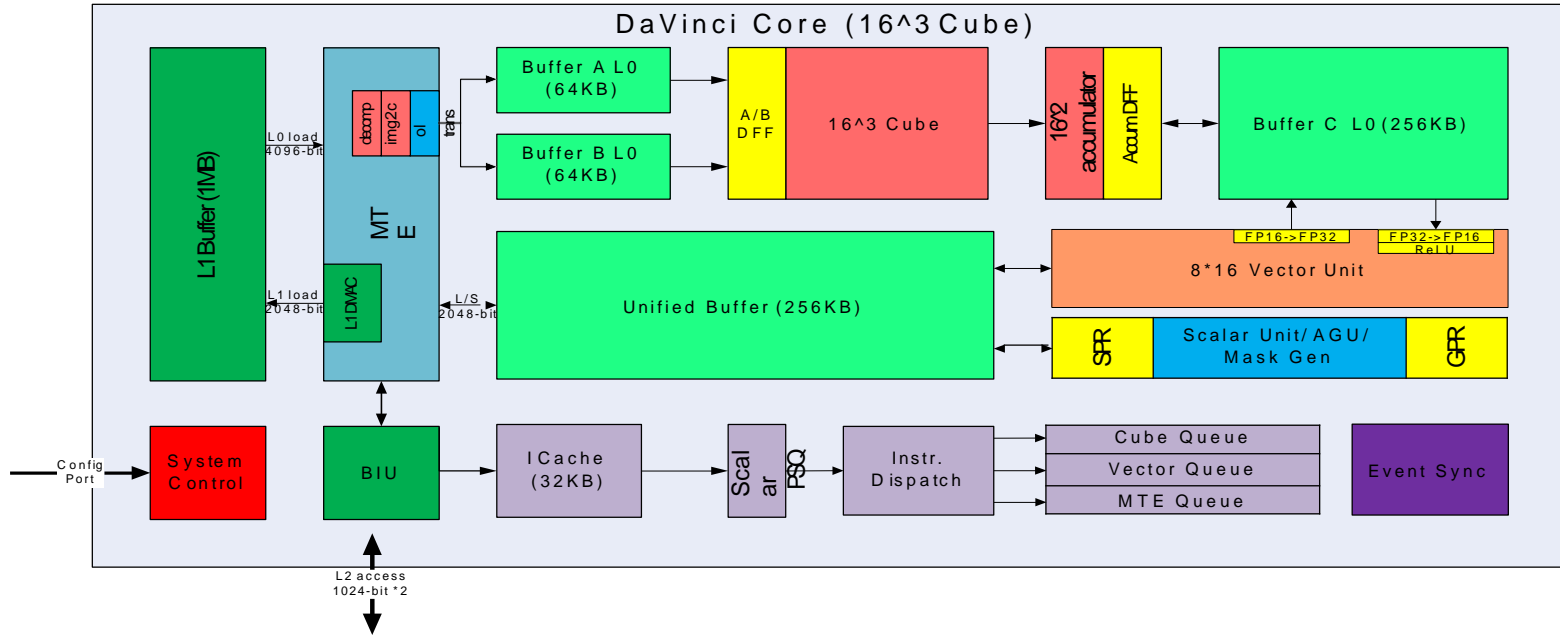
## + 3D Matrix Unit

High intensity



N	N <sup>2</sup>	N <sup>3</sup>		GPU + Tensor core	AI core + SRAM
1	1	1			
2	4	8			
4	16	128			
8	64	512			
<b>16</b>	<b>256</b>	<b>4096</b>	Area (normalized to 12 nm)	5.2mm <sup>2</sup>	13.2mm <sup>2</sup>
32	1024	32768	Compute power	1.7Tops fp16	8Tops fp16
64	4096	262144			

# DaVinci Core



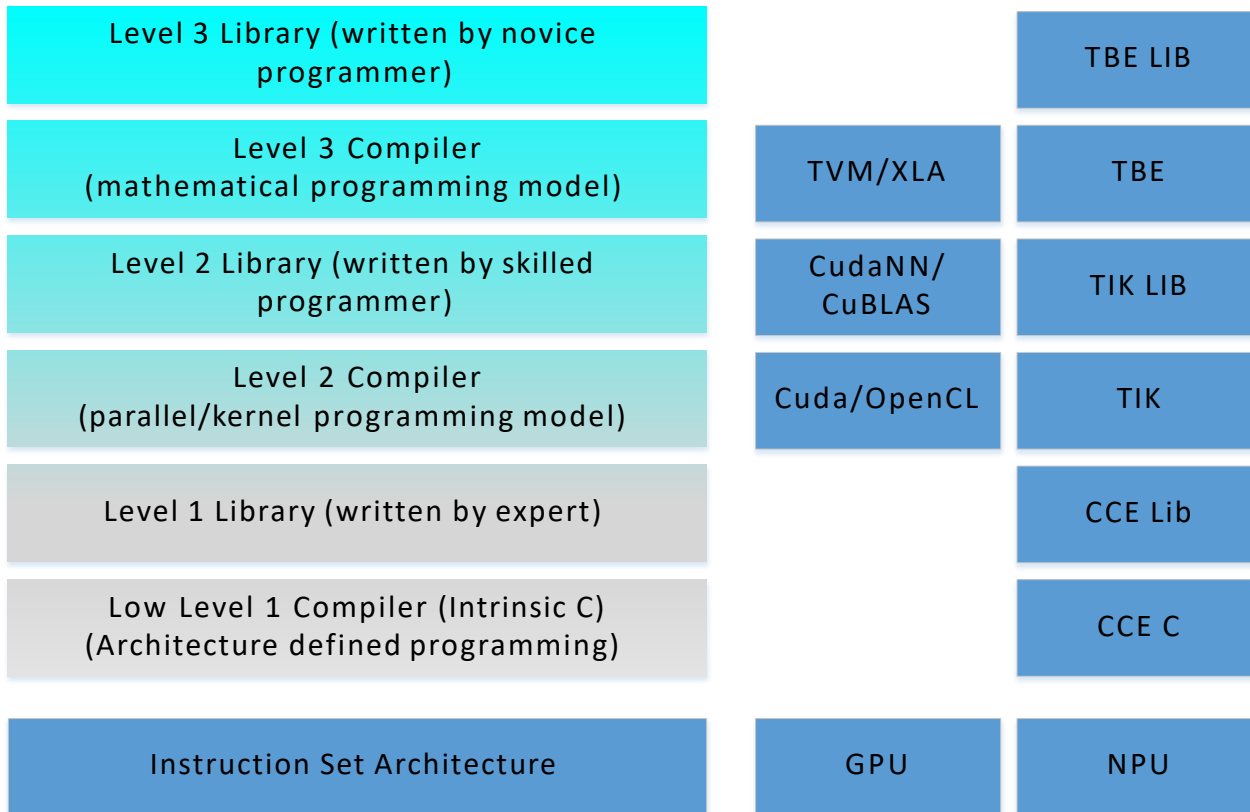
- **Cube:** 4096( $16^3$ ) FP16 MACs + 8192 INT8 MACs
- **Vector:** 2048bit INT8/FP16/FP32 vector with special functions (activation functions, NMS- Non Minimum Suppression, ROI, SORT)
- Explicit memory hierarchy design, managed by MTE

# Micro Architecture Configurations

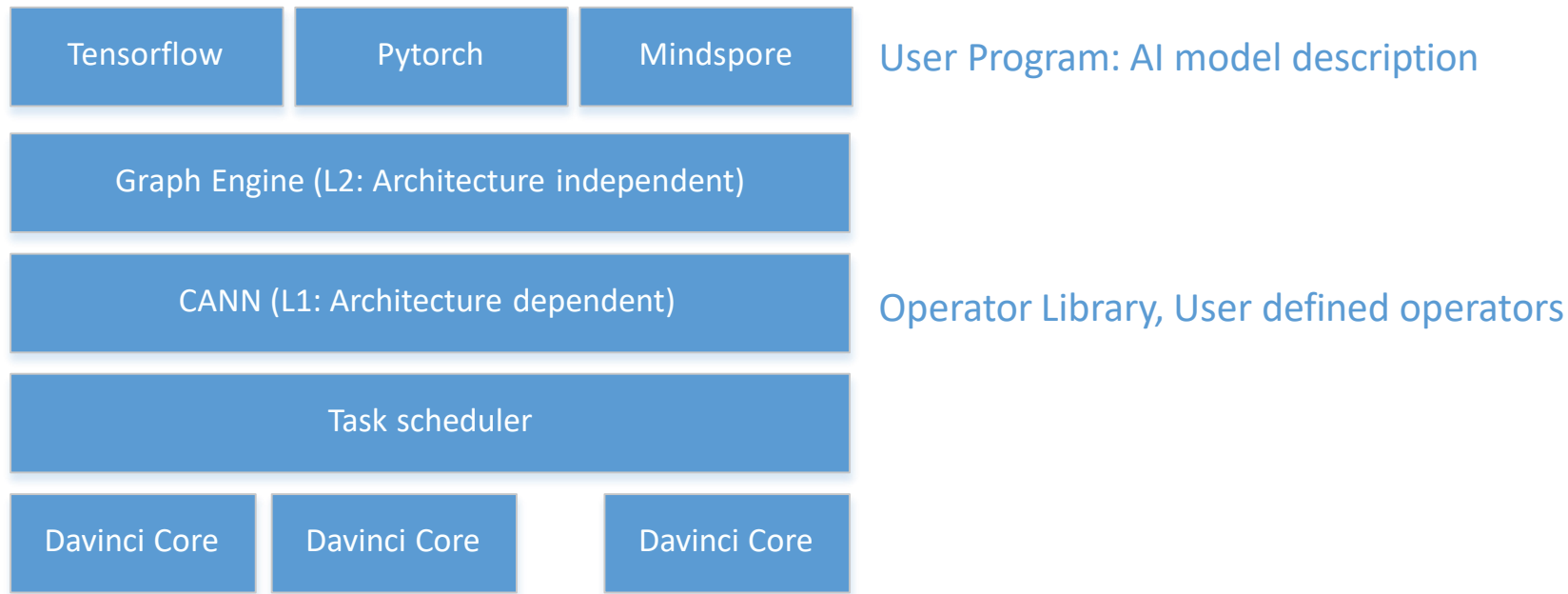
Core Version	Cube Ops/cycle	Vector Ops/Cycle	L0 Bus width	L1 Bus Width	L2 Bandwidth
Davinci Max	8192	256	Match Execution Units  Not bottleneck	A:8192 B:2048	910: 3TB/s ÷32 610: 2TB/s ÷8 310: 192GB/s÷2
Davinci Lite	4096	128		A:8192 B:2048	38.4GB/s
Davinci Tiny	512	32		A:2048 B:512	None
	Set the performance baseline	Minimize vector bound		Ensure this is not a bound	Scarce, limited by NoC, avoid bound where possible



# Overview of the Software Stack



# Putting All This Together



- User program AI model using familiar frameworks
- Extends operator library when necessary
- The tasks are executed in a single node, or over a network cluster

# DaVinci AI Processors and Products

# Mobile AP SoC: Kirin 990 contain D-lite version NPU

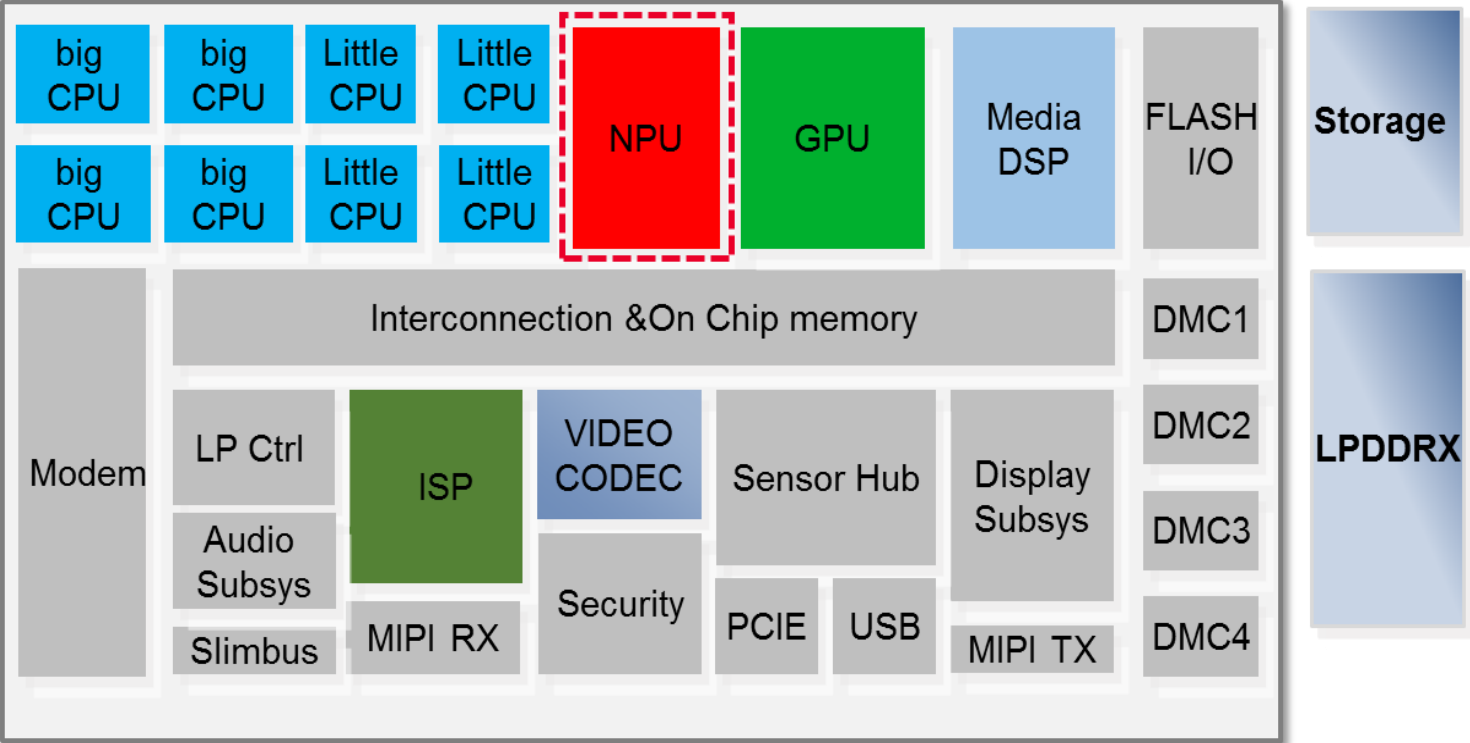


CPU 8 - Core, NPU 2+1 Core, GPU 16-Core, 2G/3G/4G/5G Modem,  
ISP 5.0, LPDDR 4X, UFS 3.0 / 2.1 HiFi Audio, 4K HDR Video, Security Engine

World's st 5G SoC Poweed by 7nm+ EUV  
World's 1st 5G NSA & SA Flagship SoC

Wolrd's 1st 16-Core Mali-G76 GPU  
World's 1st Big-Tiny Core Architechture NPU

# Mobile AP SoC: Kirin 990 contain D-lite version NPU



# Ascend AI Processor: 310 and 910



**Ascend 310**

**High Power Efficiency**

- Ascend-Mini
- Architecture: DaVinci
- FP16: 8 TeraFLOPS
- INT8 : 16 TeraOPS
- 16 Channel Video Decode – H.264/265
- 1 Channel Video Encode – H.264/265
- Power: 8W
- Process: 12nm



**Ascend 910**

**High Computing Density**

- Ascend-Max
- Architecture: DaVinci
- FP16: 256 TeraFLOPS
- INT8: 512 TeraOPS
- 128 Channel Video Decode – H.264/265
- Power: 300W
- Process: 7+ nm EUV

# Huawei AI Solutions For Inference



The advertisement features a dark blue background with a starry space pattern. At the top center, the word "Atlas" is written in a large, white, sans-serif font, with "AI Computing Platform" in a smaller font below it. Below this, four AI hardware products are displayed on glowing blue cylindrical pedestals. From left to right: 1. Atlas 200 AI Accelerator Module, a small black rectangular device. 2. Atlas 200 DK AI Developer Kit, a larger black rectangular device with a red stripe. 3. Atlas 300 AI Accelerator Card, a long, thin black circuit board. 4. Atlas 500 AI Edge Station, a black rectangular device with two antennas. At the bottom center, the text "Now meet the market" is written in a light blue font. In the bottom right corner, the Huawei logo (a red flower) and the word "HUAWEI" are displayed.

**Atlas**  
AI Computing Platform

**Atlas 200**  
AI Accelerator Module

**Atlas 200 DK**  
AI Developer Kit

**Atlas 300**  
AI Accelerator Card

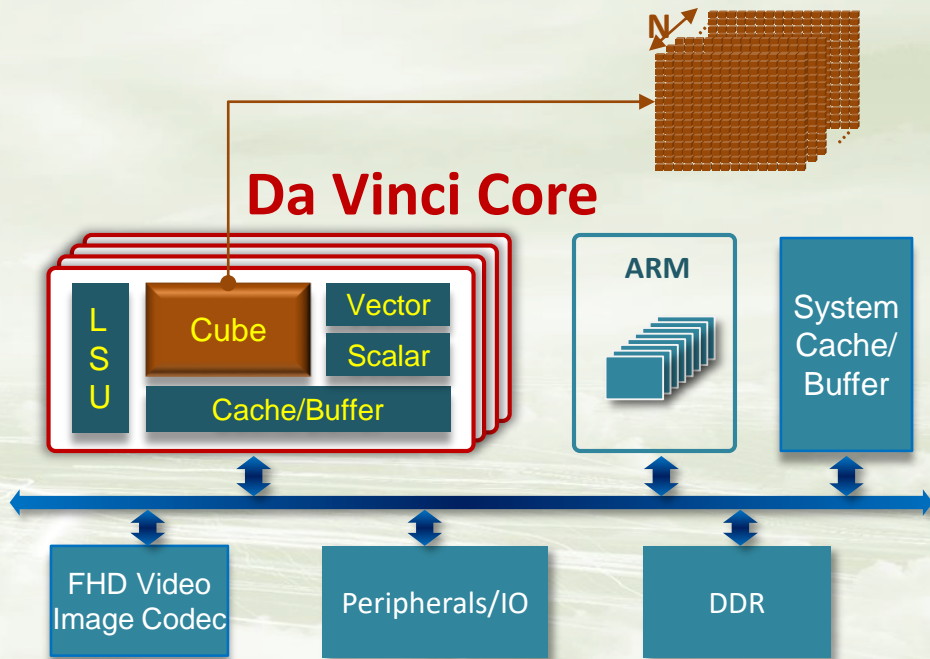
**Atlas 500**  
AI Edge Station

Now meet the market

 HUAWEI

# Ascend 310 Specification: D-mini version

SPECIFICATIONS	Description
Architecture	AI co-processor
Performance	Up to 8T @FP16
	Up to 16T@INT8
Codec	16 Channel Decoder – H.264/265 1080P30 1 Channel Encoder
Memory Controller	LPDDR4X
Memory Bandwidth	2*64bit @3733MT/S
System Interface	PCIe3.0 /USB 3.0/GE
Package	15mm*15mm
Max Power	8Tops@4W, 16Tops@8W
Process	12nm FFC

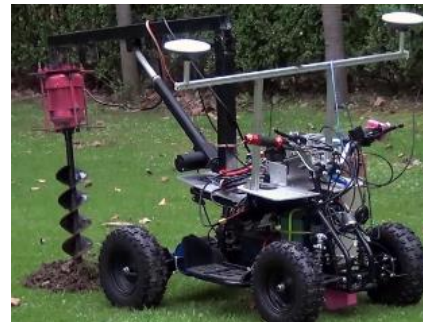


**Note:** This is typical configuration, high performance and low power sku can be offered based on your requirement.

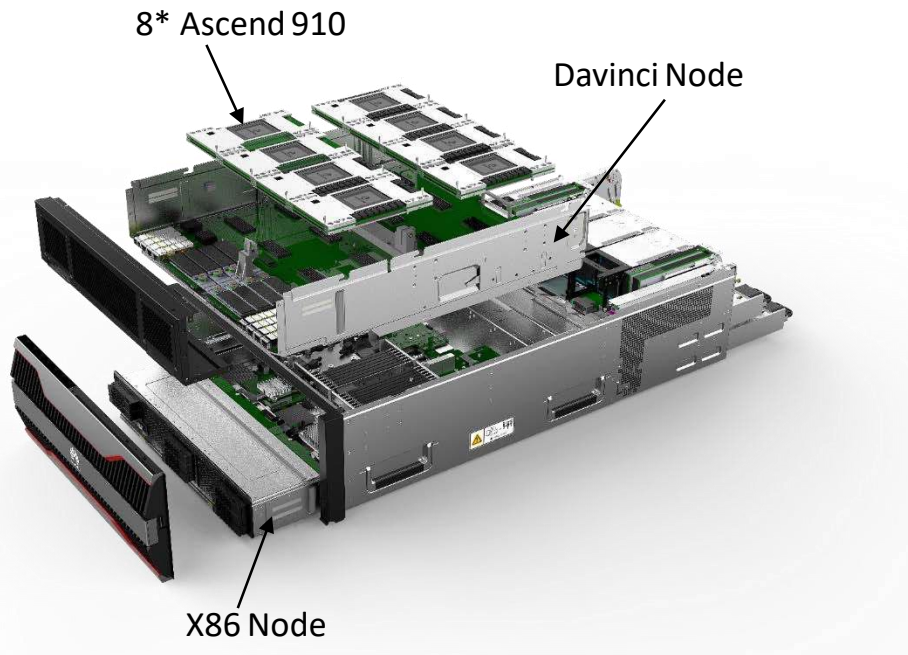


# Applications Built on Atlas 200 DK

Semantic assisted high-precision  
3D reconstruction  
Facial recognition  
Blocked face detection  
Image marking  
Handwritten text recognition  
HDR  
Image processing  
Fundus retinal vascular segmentation  
Classification network  
Cartoon image generation  
Vehicle detection  
Facial expression  
transplant  
**Unmanned vehicle**  
Robots  
AR shadow generation  
Robotic arm  
Age recognition  
NPL  
Protein subcellular position prediction

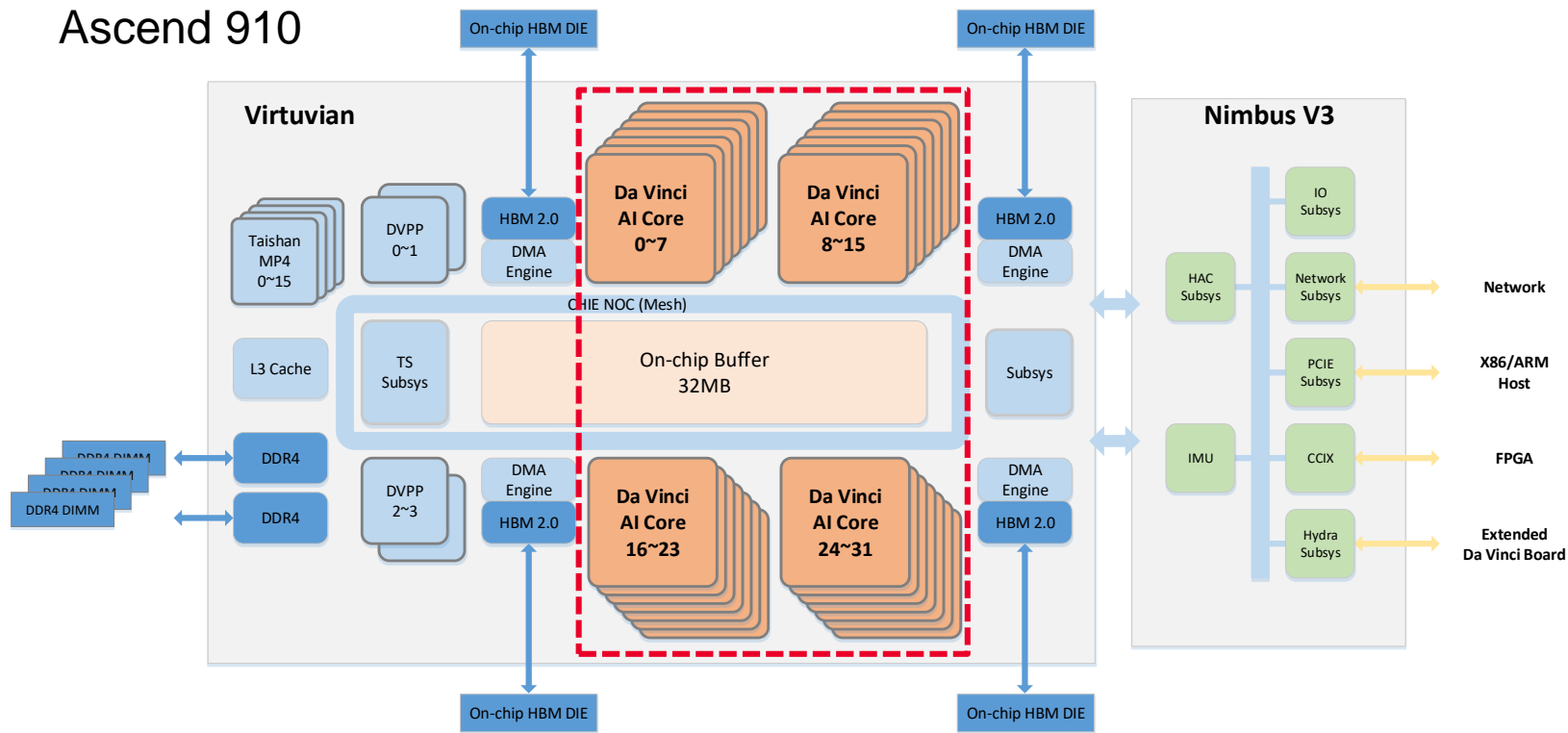


# Huawei AI Solutions For Training: Ascend 910 Server

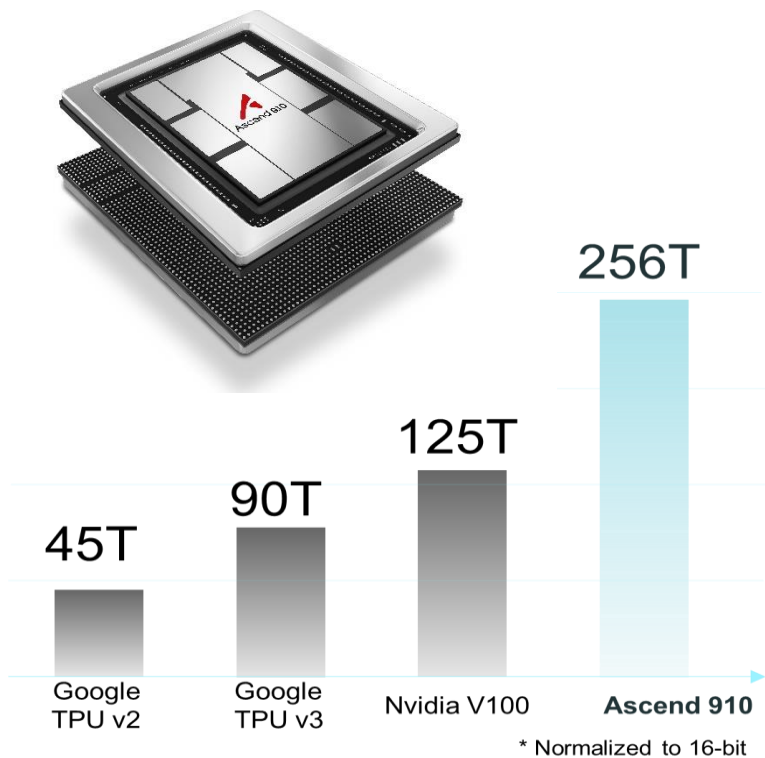


# AI Training SoC: Atlas 910 contains D-max version

## Ascend 910



## Ascend 910 Specification

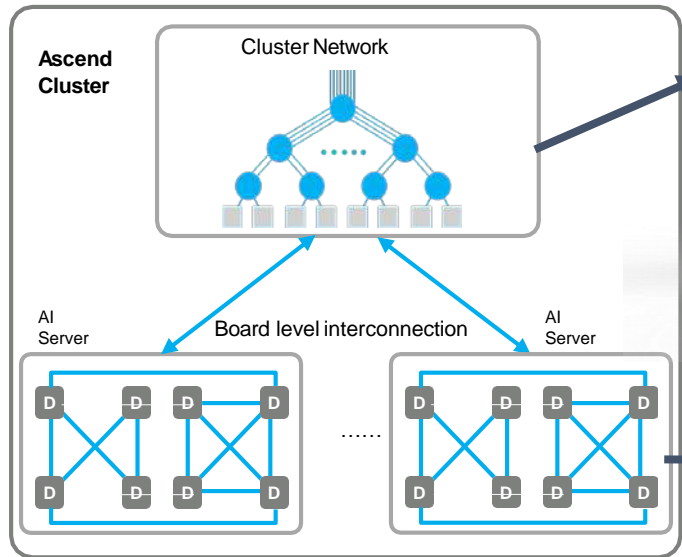


SPECIFICATIONS	Description
Architecture	AI co-processor
Performance	Up to 256T @FP16 Up to 512T@INT8
Decoder Codec	128 Channel FHD Video Decoder – H.264/265
Memory Interface	32GB HBM Gen2
System Interface	PCIe3.0 x16 & HCCS & 100G RoCE
Max Power	Up to 350W
Process	7nm+

# Ascend 910 Cluster

- 2048 Node x 256TFlops = 512 Peta Flops

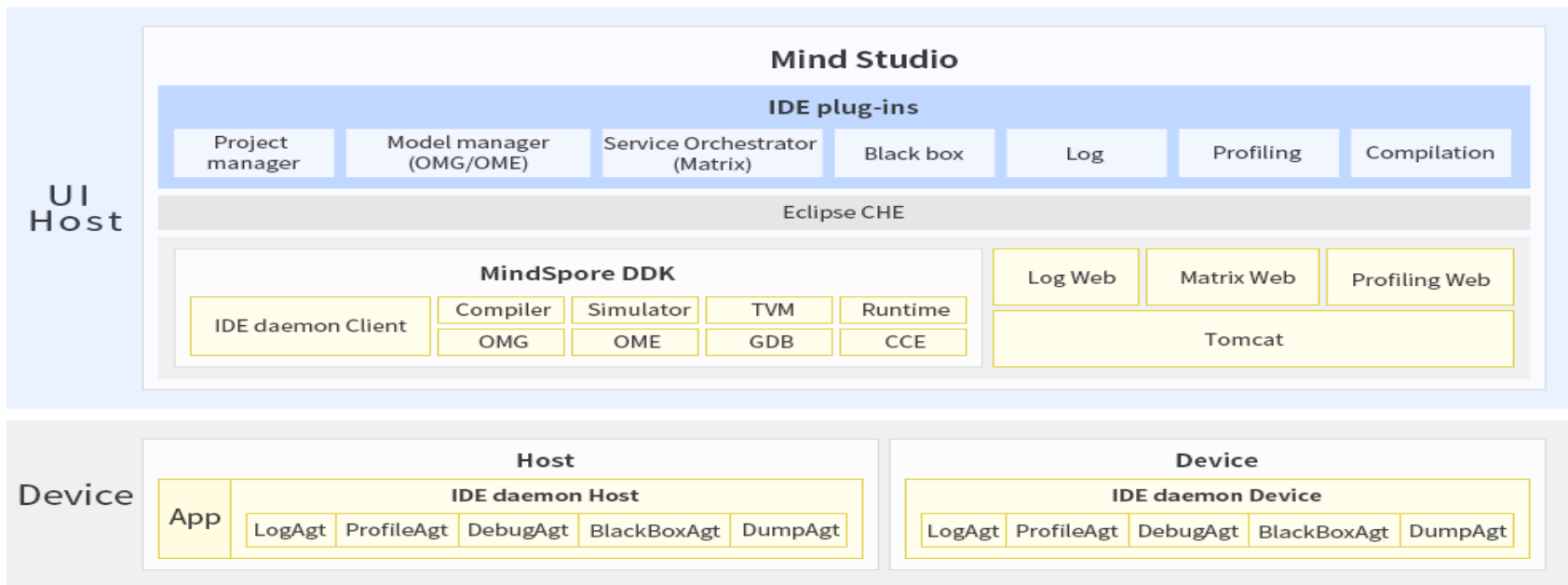
- 1024~2048 Node Cluster



Ascend910 Board

# Development Tool: Mind Studio

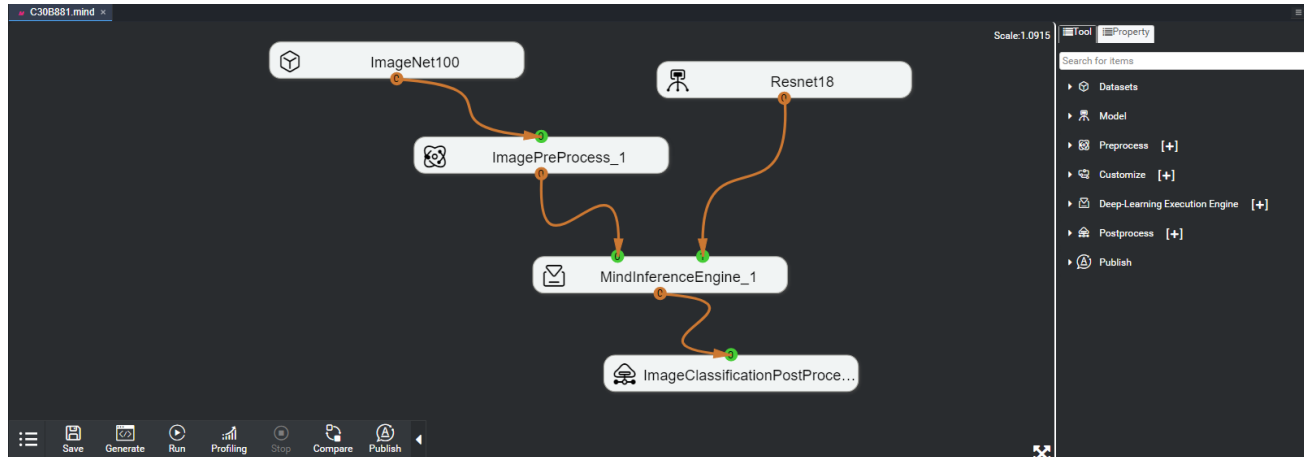
# Mind Studio — Atlas 200 DK Tool Chain



- Operator development
- Offline model tool
- Service orchestration tool
- App Development Tool

# Service Orchestrator

Drag-and-drop mode: auto-generate codes

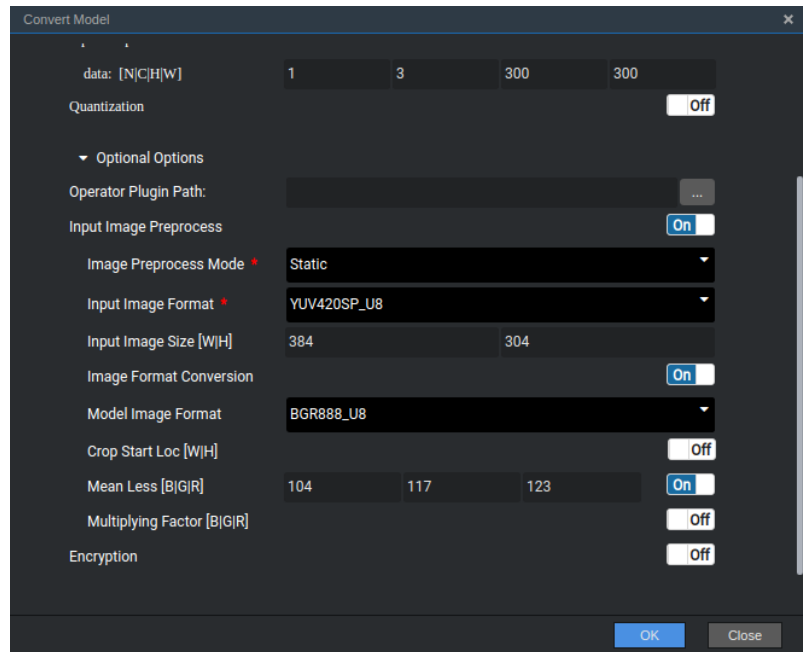
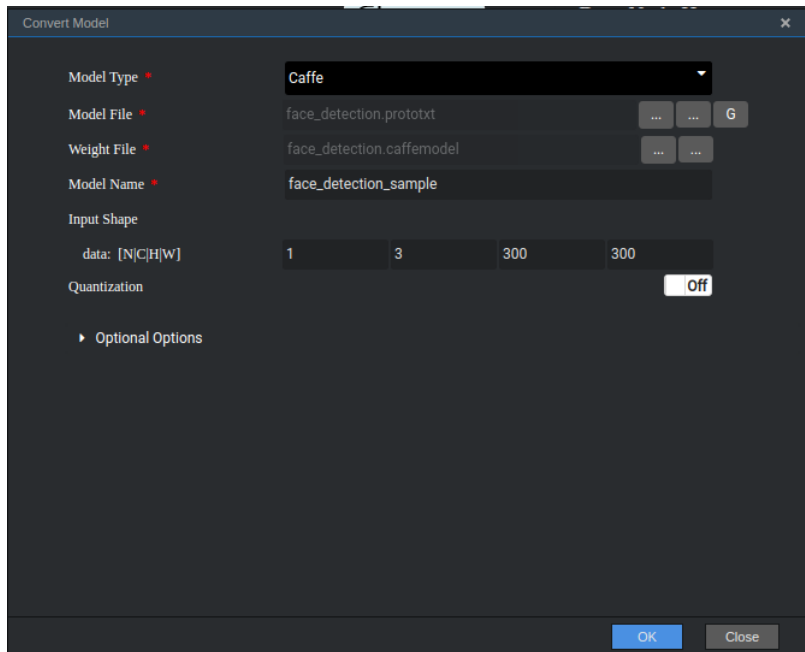




# Model Manager

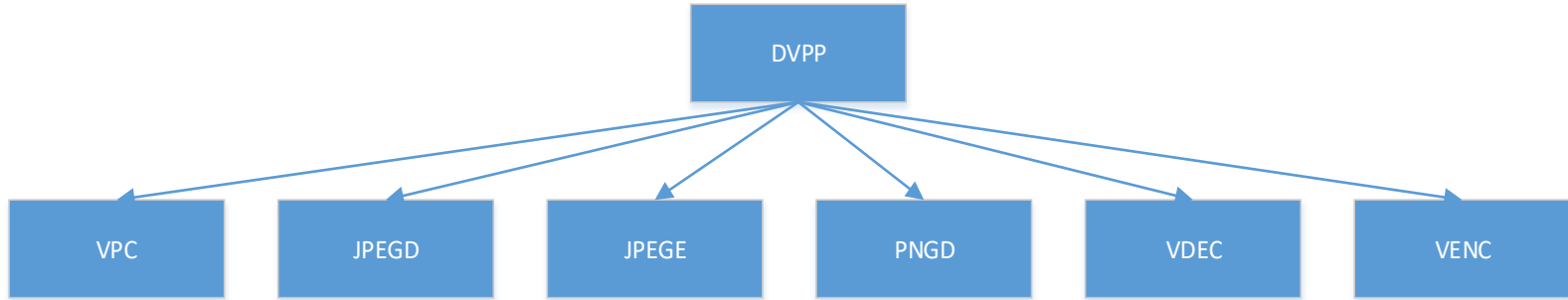
Wizard for OMG, configuring AI preprocessing (AIPP) inside of offline model.

AIPP involves image cropping, color space conversion (CSC), and mean subtraction and multiplication coefficient (pixel changing). All these functions are implemented by the AI core.



# DVPP — Brief Review

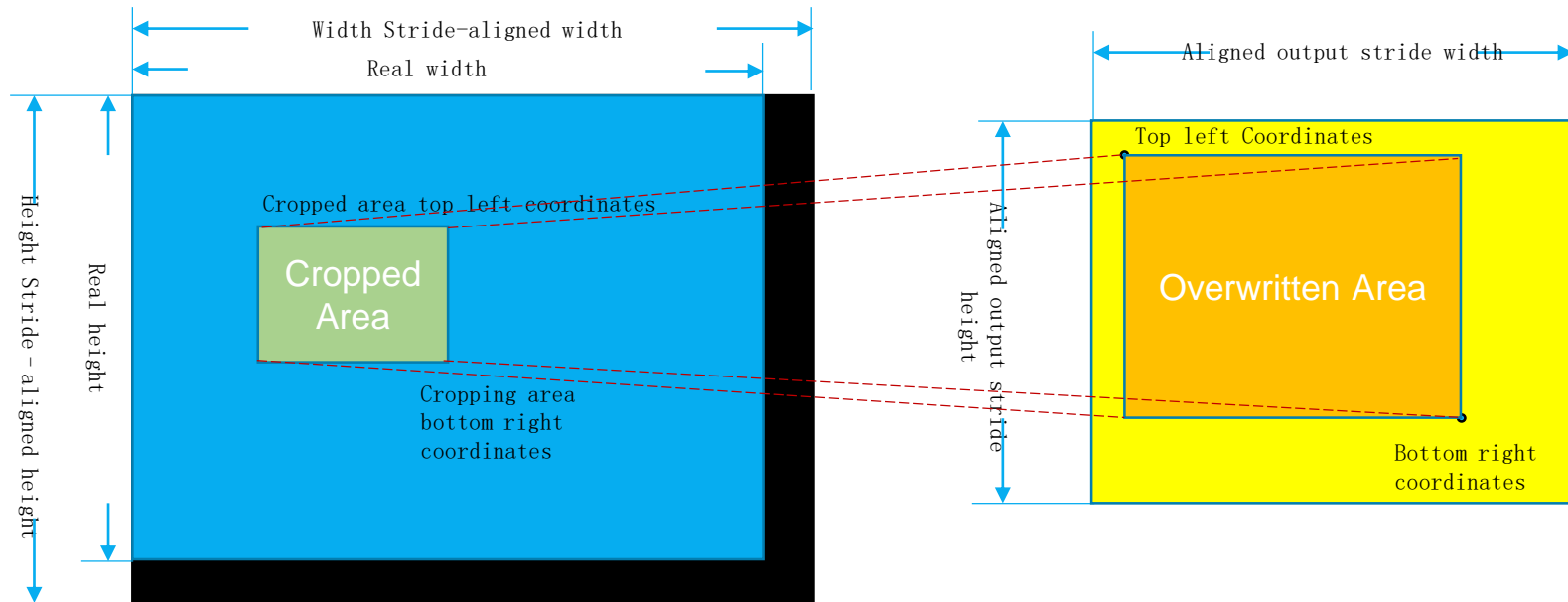
DVPP(Digital Video Pre-Processor) : image/video encoding, decoding, resizing, cropping, format conversion.



DVPP	Digital Vision Pre-Process
VPC	Vision Preprocess Core
JPEGD	JPEG Decode
JPEGE	JPEG Encoder
PNGD	Portable Network Graphics Decoder
VDEC	Video Decoder
VENC	Video Encoder

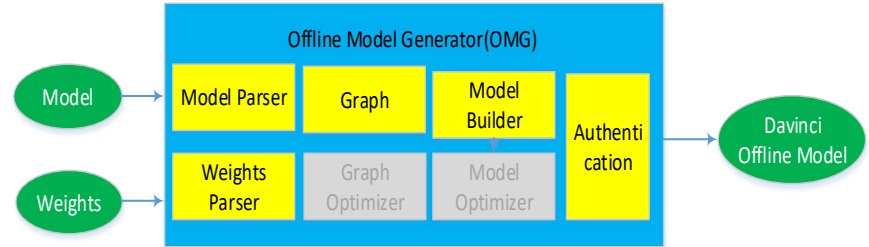
# DVPP — VPC Common Operation

VPC Common Operation: Crop+Resize:



# Conversion Capabilities

Open-Source Framework	Input File	Storage Format
Caffe	*.prototxt, network structure file	Protobuf format
	*.caffemodel, weight file	Binary file
Caffe2	predict_net.pb, network definition file	Protobuf format
	All input data of the init_net.pb network, including all weight data and input data description of the first operator	Protobuf format
Mxnet	xx-symbol.json, network structure file	JSON format
	xx.params, weight file including the data, data type, and data format of the weight node	Binary file
Tensorflow	*.pb: The network model and weight data are in the same file.	Protobuf format



- ❑ The OMG adapts to different network files and weight files in various frameworks and converts them into offline files in DaVinci format for the OME to use.
- ❑ The OMG is independently executed offline on the host (Linux Ubuntu) and is provided to the IDE as an OMG executable program. The OMG can be called by command lines.

# OMG — CMD Introduction

## Sample cmd:

```
./omg --model=./alexnet.prototxt --weight=./alexnet.caffemodel --framework=0 --output=./domi
```

```
./omg --model=$HOME/ResNet18_deploy.prototxt --weight= $HOME/ResNet18_model.caffemodel -  
--framework=0 --output=$HOME/outdata/ResNet18 --input_shape= "data:1 ,3, 224, 224" --  
insert_op_conf= ./aipp_con.cfg
```

```
omg: usage: ./omg <args>  
example:  
./omg --model=./alexnet.prototxt --weight=./alexnet.caffemodel  
--framework=0 --output=./domi  
Arguments explain:  
--model          Model file  
--weight         Weight file. Required when framework is Caffe  
--framework      Framework type(0:Caffe; 3:Tensorflow)  
--output         Output file path&name(needn't suffix, will add .om automatically)  
--encrypt_mode   Encrypt flag, 0: encrypt; -1(default): not encrypt  
--encrypt_key    Encrypt key file  
--certificate    Certificate file  
--hardware_key   ISV file  
--private_key    Private key file  
--input_shape    Shape of input data. E.g.: "input_name1:n1,c1,h1,w1;input_name2:n2,c2,h2,w2"  
--h/help        Show this help message  
--cal_conf       Calibration config file  
--insert_op_conf Config file to insert new op  
--op_name_map    Custom op name mapping file  
--plugin_path    Custom op plugin path. Default value is: "./plugin". E.g.: "path1;path2;path3".  
--om            The model file to be converted to json  
--json          The output json file path&name which is converted from a model  
--mode          Run mode. 0(default): model => davinci; 1: framework/davinci model => json; 3: only pre-check  
--target        Target platform. (mini)  
--out_nodes     Output nodes designated by users. E.g.: "node_name1:0;node_name1:1;node_name2:0"  
--input_format  Format of input data. E.g.: "NCHW"  
--check_report  The pre-checking report file. Default value is: "check_result.json"  
--input_fp16_nodes Input node datatype is fp16 and format is NCHW. E.g.: "node_name1;node_name2"  
--is_output_fp16 Net output node datatype is fp16 and format is NCHW, or not. E.g.: "false,true,false,true"  
--ddk_version   The ddk version. E.g.: "x.y.z.Patch.B350"  
--net_format    Set net prior format. ND: select op's ND format preferentially; 5D: select op's 5D format preferentially  
--output_type   Set net output type. Support FP32 and UINT8  
--fp16_high_prec FP16 high precision. 0(default): not use fp16 high precision; 1: use fp16 high precision
```

## Common Paras:

--model path to original model file

--weight path to weight file (Caffe)

--framework 0: caffe  
3: tensorflow

--output output om file path

--plugin\_path customized Operator path

--insert\_op\_conf aipp config file path

--input\_shape input data shape.

--fp16\_high\_prec generate FP16 model

--mode 0: generate davinci model  
1: om model to json  
3: only pre-check

For other parameters, please refer to [Atlas200DK model conversion guide](#)

# MindSpore: All-scenario AI computing framework

AI applications



## MindSpore

Unified APIs for all scenarios

Automatic differentiation

Automatic parallelization

Automatic tuning

MindSpore intermediate representation (IR) for computation graphs

On-device execution

Pipeline parallelism

Deep graph optimization

Device-edge-cloud cooperative distributed architecture (for deployment, scheduling, communications, etc.)

Easy development: AI Algorithm As Code

Efficient execution: Optimized for Ascend, GPU support

Flexible deployment: On-demand cooperation across all scenarios

Processors: Ascend, GPU, CPU



# ModelArts: Full-pipeline model production services



ModelArts



Supports full pipeline – From data collection and model development to model training and deployment



4,000+ training tasks per day (total of 32,000 training hours)

Visual tasks: 85%; audio tasks: 10%; ML tasks: 5%



30,000+ developers

# Software + hardware co-optimization: Stronger performance

## AI computing challenges

### Complex computing

- Scalar, vector, and tensor computing
- Hybrid precision computing
- Parallelism between data augmentation and minibatch computing
- Parallelism between gradient aggregation and minibatch computing

### Diverse computing power

- CPUs, GPUs, and Ascend processors
- Diverse computing units: scalar, vector, and tensor

## MindSpore

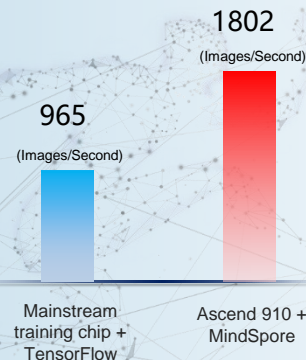
### Framework optimization

- Pipeline parallelism
- Cross-layer memory reuse

### Software + hardware co-optimization

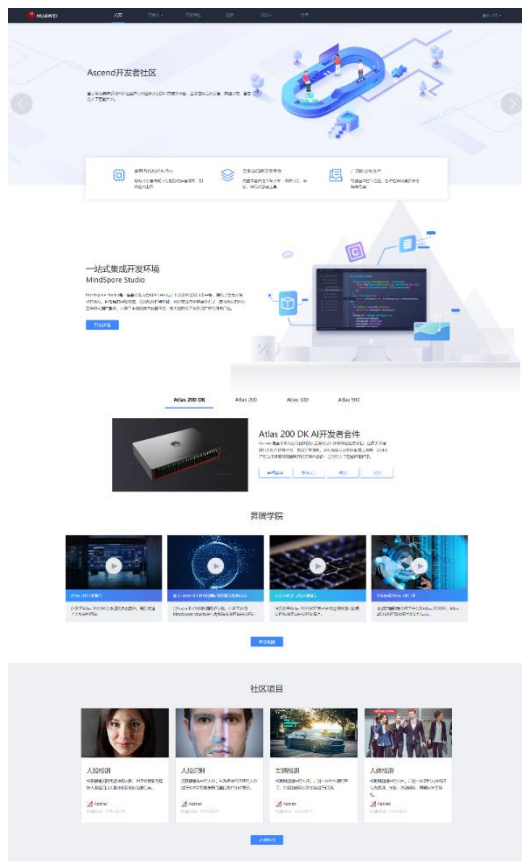
- On-device execution
- Deep graph optimization

- ResNet 50 V1.5
- ImageNet 2012
- Based on optimal batch sizes





# Ascend Developer Community: End-to-End Services for Developers



## Ascend developer portal

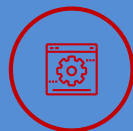
<https://ascend.huawei.com>

Unified channel for open capability release and operation



## One-stop support services

Technical documents, online resources, development tools, video tutorials, support services, interactive communities, and information release



## Developer-centric experience

Provides rich, friendly, easy-to-use, and quick portal experience for developers

## Ascend Eco-Systems in China

### Huawei AI Developer Support Program



2019  
华为Atlas  
人工智能开发者大赛  
AI Developer Challenge

- Public cloud voucher
- Certification course ticket
- Atlas developer kit
- HUAWEI CLOUD MVP
- Prize money



Developer Technical Salon

Developed **50,000** developers and 100+ startups in 2019

### AI Talent Development Program

The first batch of courses co-created with colleges

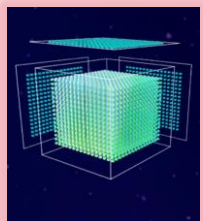
 Shanghai Jiao Tong University	 Zhejiang University	 Tsinghua University	 Peking University	 Tianjin University
 Renmin University of China	 Beijing University of Posts and Telecommunications	 Huazhong University of Science and Technology	 Southern University of Science and Technology	 Xidian University

Co-created courses & labs with **30+** colleges in 2019

### We also already started in Canada:

- SFU: 2020 Spring Term Computer Vision class for Professional Master program;
- MakeUofT Hackathon: Feb 15 -16

# Thank you.



Da Vinci



Ascend



Atlas

Bring digital to every person, home and organization for a fully connected, intelligent world.

Copyright©2020 Huawei Technologies Co., Ltd.  
All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.

Huawei Confidential

