HUAWEI | MUNICH RESEARCH CENTER

Al-enabled science

— Challenges and opportunities arising from the convergence of artificial intelligence, high-performance computing and big data

Heiko Joerg Schick Chief Architect | Advanced Computing

Version 1 - DRAFT

HUAWEI

Presenting the work of many people at Huawei

HUAWEI TECHNOLOGIES DÜSSELDORF GmbH

Agenda

Use cases

Challenges for system architecture (hardware, software, workflow and algorithms)

Artificial intelligence

Special computer architectures? How come, for what reason and why?

Examples

In ten years



Heiko Joerg Schick Chief Architect | Advanced Computing Munich Research Center

HUAWEI TECHNOLOGIES Duesseldorf GmbH Riesstrasse 25, C3

Riesstrasse 25, C3 80992 Munich

Mobile +49-151-54682218 E-mail <u>heiko.schick@huawei.com</u>

LinkedIn https://www.linkedin.com/in/heikojoergschick/



In symbols one observes an advantage in discovery which is greatest when they expressed the exact nature of a thing briefly and, as it were, picture it; then indeed the labor of thought is wonderfully diminished.

In Symbolen beobachtet man einen Entdeckungsvorteil, der am größten ist, wenn sie die genaue Natur einer Sache kurz zum Ausdruck bringen und gleichsam abbilden; dann wird in der Tat die Arbeit des Denkens wunderbar verringert.

— Gottfried Wilhelm Leibniz





In models one observes an advantage in discovery which is greatest when they expressed the exact nature of a thing briefly and, as it were, picture it; then indeed the labor of thought is wonderfully diminished.

In Symbolen beobachtet man einen Entdeckungsvorteil, der am größten ist, wenn sie die genaue Natur einer Sache kurz zum Ausdruck bringen und gleichsam abbilden; dann wird in der Tat die Arbeit des Denkens wunderbar verringert.

— Gottfried Wilhelm Leibniz





Integration of artificial intelligence and high performance computing

Use cases:

- 1. Using <u>HPC technologies</u> to execute and enhance <u>AI performance</u>.
- 2. Using <u>HPC simulations</u> to train <u>AI algorithms</u>.
- 3. Using <u>AI algorithms</u> to configure and autotune <u>AI workloads</u> or <u>HPC simulations</u>.
- 4. Using <u>AI algorithms</u> to analyse results of <u>HPC simulations</u>.
- 5. Using <u>AI algorithms</u> to learn from <u>HPC simulations</u> and produce learned surrogates.

Workload-optimized Systems Requirements

Real-world workflows test a broad range of system design points:

- Challenge is to determine the best design region boundaries
- Benchmarks are necessary but not sufficient ...
- Need to examine workflows to get the full picture of requirements
- Need to understand how multiple workflows run simultaneously across system(s)



Powerful Information: Al-enabled Science — a new Area of Computing



Dimensions of data growth



8

On the Roads to Software Defined Environments (SDE)

- Workload types are growing and becoming more flexible and diverse.
- Cloud infrastructure is becoming programmable to meet the requirements in efficiency and resiliency.
- Heterogeneity is increasingly present and important.



- workloads
- Fixed system hardware, manual scaling
- Hard-wired workload, minimal configuration



- Diverse workloads, limited patterns
- Homogenous resource pooling
- Expert configuration and mapping of workloads



- Rapidly changing workloads, dynamic patterns
- Dynamic automatic composition of heterogeneous systems
- Autonomic and proactive management

Four Types of Parallel Architectures for <u>Advanced Computing</u>

	Compute Intensive	Data Intensive: Data at Rest	Data Intensive: Streaming Data	Al-enabled Science	
Programming Language	C/C++, Fortran, MPI, OpenMP, CUDA	Java, JAQL, Python	SPL, C, Java	C/C++, Python, UPC, SHMEM, MPI, OpenMP	
Characteristics	 Data is Generated Long Running Small Input Massive Output 	 Data at Rest Low Velocity Mixed Variety High Volume 	 Data in Motion High Velocity Mixed Variety High Volume 	 Data is Moving Long Running All Data View Small Messages 	
Category	 Network Dependent Structured Communication 	 Embarrassingly Parallel Structured Communication 	Embarrassingly ParallelRandom Communication	Network DependentRandom Communication	
Applications	LINPACK	Apache Hadoop / Key- Value Databases	Apache Storm		
Network Topology		Input Data (on disk) Mappers Mappers Reducers Output Data	$\rightarrow \square \checkmark \square \land \square \land \square \rightarrow \square \land \square \rightarrow \square \land \square \rightarrow \square \land \square \rightarrow \square \rightarrow$		
Scaling	Strong	Weak		Both	
Cluster Stack Management	Custom	Diverse control systems	N/A		
Schedulers	Slurm	YARN, Mesos	N/A		
File System	Cluster File-System	Distributed across nodes (e.g. HDF	POSIX or HDFS		
Operating system	RedHat / CentOS	RedHat / CentOS	Ubuntu		



HUAWEI | MUNICH RESEARCH CENTER

HUAWEI | MUNICH RESEARCH CENTER



Computational Methods

Hardware	_							A	
i la analo	Processors							App	
	Co-processors								Very High
	GPGPUs								High
	FPGAs								Medium
	ASICs								Very Low
	DSPs								Ignorable
									N. A.

Implementation of end-to-end lifecycle in AI projects [Alake, 2020], [Sato et al., 2019]







The hype roller coaster of artificial intelligence [Villain, 2019]



IM & GENET

Neural networks beat human performance /1 [Giró-i-Nieto, 2016], [Gershgorn, 2017]

— <u>Example:</u> Image classification on ImageNet

15 million images in dataset, 22,000 object classes (categories) and 1 million images with bounding boxes.

IM . GENET

Neural networks beat human performance /1 [Giró-i-Nieto, 2016], [Gershgorn, 2017]

- Example: Image classification on ImageNet

15 million images in dataset, 22,000 object classes (categories) and 1 million images with bounding boxes.



Neural networks beat human performance /2 [Russakovsky et al., 2015], [Papers With Code, 2020]

- Example: Image classification on ImageNet



19

Two distinct eras of compute usage in training AI systems [McCandlish et al., 2018], [Amodei et al., 2019]











Rich variety of computing architectures



- Wide range of options to optimise for performance and efficiency:
 - Central processing unit (CPU) executes general purpose applications (e.g.
 N-body methods, computational logic, map reduce, dynamic programming)
 - General-purpose computing on graphics processing units (GPGPU)
 accelerates compute intensive and time consuming applications for the CPU
 (e.g. dense linear algebra and sparse linear algebra)
 - Digital signal processor (DSP) accelerates signal processing for post camera operations (e.g. spectral methods)
 - Image signal processor (ISP) executes processing for camera sensor pipeline
 - Vision processing unit (VPU) accelerates machine vision tasks
 - Network processor (NP) accelerates packet processing
 - Neural processing unit (NPU) accelerates artificial intelligence applications
 (e.g. matrix-matrix multiplication, dot-products, scalar *a* times *x* plus *y*)

Each of these options represents different power, performance, and area trade-offs, which should be considered for specific application scenarios.

Rich variety of computing architectures



- Wide range of options to optimise for performance and efficiency:
 - Central processing unit (CPU) executes general purpose applications (e.g.
 N-body methods, computational logic, map reduce, dynamic programming)
 - General-purpose computing on graphics processing units (GPGPU)
 accelerates compute intensive and time consuming applications for the CPU
 (e.g. dense linear algebra and sparse linear algebra)
 - Digital signal processor (DSP) accelerates signal processing for post camera operations (e.g. spectral methods)
 - Image signal processor (ISP) executes processing for camera sensor pipeline
 - Vision processing unit (VPU) accelerates machine vision tasks
 - Network processor (NP) accelerates packet processing
 - Neural processing unit (NPU) accelerates artificial intelligence applications
 (e.g. matrix-matrix multiplication, dot-products, scalar *a* times *x* plus *y*)

Each of these options represents different power, performance, and area trade-offs, which should be considered for specific application scenarios.

Ubiquitous and future AI computation requirements



Focus on innovation, continuous dedication and backward compatibility



"Once a technology becomes digital—that is, once it can be programmed in the ones and zeros of computer code—it hops on the back of Moore's law and begins accelerating exponentially."

– Peter H. Diamandis & Steven Kotler, The Future Is Faster Than You Think

Focus on innovation, continuous dedication and backward compatibility



Al Accelerator Module 16 TOPS of INT8

Atlas 200

Atlas 300

Al Accelerator Card

- 64 TOPS of INT8 @ 67 W
- 32 GB memory
- 64-channel HD video real-time analytics
- Standard half-height half-length PCIe card form factor, applicable to general-purpose servers

Ascend 310

16-channel full-HD video decoder: H.264/265

1-channel full-HD video encoder: H.264/265

AI SoC with ultimate efficiency

Max. power consumption: 8 W

Half precision (FP16): 8 TFLOPS

Integer precision (INT8): 16 TOPS



reddot award 2019 winner

AI Edge Stations

- 16 TOPS of INT8
- 25–40 W

Atlas 500

- Wi-Fi & LTE
- 16-channel HD video real-time analytics
- Fanless design, -40°C to +70°C environments



Atlas 800



Plug-and-play installation Ultimate Performance





Atlas 900 AI Cluster

The pinnacle of computing power

- Thousands of Ascend 910 AI processors
- High-speed interconnection
- . Delivers up to 256 to 1024 PetaFLOPS at FP16
- Can complete model training based on • ResNet-50 within 59.8 seconds
- 15% faster than the second-ranking product
- · Faster AI model training with images and speech

52 mm x 38 mm x 10 mm



Atlas 200 DK

Quickly build development environments in 30 minutes

16-channel HD video real-time analytics, JPEG decoding

4 GB/8 GB memory, PCIe 3.0 x4 interface Operating temperature: -25°C to +80°C

- 16 TOPS of INT8 @ 24 W
- 1 USB type-C, 2 camera interfaces, 1 GE port, 1 SD card slot

• 12nm

4 GB/8 GB memory

Storage-intensive



5280 4U 40-drive storage model

2280 2U 2S balanced model

1280 1U 2S high-density model























• 7nm

Kunpeng 920

The industry's highest-performance ARM-based server CPU

- ARM v8.2-architecture
- up to 64 cores, 2.6 GHz



Ascend 910

Highest compute density on a single chip

- Half precision (FP16): 256 TFLOPS
- Integer precision (INT8): 512 TOPS
 - 128-channel full-HD video decoder: H.264/265 Max. power consumption: 350 W



- 8 DDR4 memory channels
- Integrated 100GE LOM and encryption and compression engines
- Supports 2- or 4-socket interconnects







- - PCIe 4.0 and CCIX





Focus on innovation, continuous dedication and backward compatibility



4 x 4 Data

Building blocks and compute intensity

Scalar Unit

Full flexibility in computation

Cube Unit High intensity computation

Vector Unit 4 x 4 Data Rich and efficient operations

16 x 4 Multiply units

4 x 4 Add Units

Advantages of special compute units



Advantages of special compute units





Number of parameters and floating point operations per second (FLOPS) for each layer of the AlexNet artificial intelligence model.

	99% of the computations are	Typical CNN networks						
	matrix-matrix multiplications	AlexNet	VGG16	Inception-v3				
	Model memory (MB)		> 500	90-100				
Cucles - 1	Parameter count (Million)	60	138	23.2				
Data per cycle = $Rd 2 * 16 * 16; Wr 16*16$	Computation amount (Million)	720	15300	5000				

Upscaling and colourisation of video footage



Despina Manaki is the earliest-born person on film. In 1905, when she was 114 years old (born 1791), she was filmed by her grandsons, Yanaki and Milton Manaki, cinema pioneers in the Balkans and the Ottoman Empire. This video shows the full force of artificial intelligence restoring old footages.

Source: https://www.reddit.com/r/interestingasfuck/comments/idbtrg/i upscaled and colorized the footage about the/

Retinal blood vessel segmentation in the eyeground



- The fundus retinal blood vessel segmentation application was developed for the Atlas 200 DK inference system, in partnership with the Nankai University, led by Professor Li Tao of Intelligent Computing System Research Office .
- This project makes full use of the neural network computing power of the Atlas 200 DK system to segment the fundus vessels in real-time.
- The total inference time of **20 pictures is 761.8 milliseconds**, and the average inference time of one image is 38 milliseconds.



An overview of the vascular segmentation model







Things we can do in science with AI technologies now /1

	Underhood air flow	External aerodynamics	HVAC			
Standard case	3 stationary points - vmax, 250 km/h - Idle - Mountain 30 km/h	 140 km/h With and without underhood flow 	 Defrost Pull-down Heater mode 			
Target values	Cooling air mass flowp distribution	 Drag Lift cp distribution 	 Air distribution Velocity distribution Δp 			
Turnaround times	5-8 h/over-night	2-3 days	Few hours/over-night			

Things we can do in science with AI technologies now /2

DeepSUM: Deep neural network for Super-resolution of Unregistered Multitemporal images

Andrea Bordone Molini, Diego Valsesia, Giulia Fracastoro, and Enrico Magli

Abstract—Recently, convolutional neural networks (CNN) have been successfully applied to many remote sensing problems. However, deep learning techniques for multi-image super-resolution from multitemporal unregistered imagery have received little versions of the same scene to be combined by means of MISR techniques, where the reconstruction of high spatial-frequency details takes full advantage of the complementary information coming from different observations of the same scene.



In ten years ...

- 1. Learned AI model begin to replace data
- 2. The discovery process of experiments is dramatically refactored
- 3. We will persue many questions semi-autonomously
- 4. HPC simulations and AI technologies will merge
- 5. Al algorithms will contribute to advancing theories
- 6. Al technologies are becoming a common part of scientific laboratories activities

- \rightarrow AI technologies will help to solve new challenges
- \rightarrow AI technologies will become a new partner in simulation and data analytics
- →AI technologies will generate new <u>computing architectures</u>, new <u>software environments</u>, new <u>policies</u>, new <u>user</u> <u>communities</u> and new <u>ways of dissemination</u>.



Stay safe — stay healthy

Copyright © 2020 Huawei Technologies Düsseldorf GmbH. 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.