# Few hints for the post exascale architectures

**Marc Duranton and Denis Dutoit**
Commissariat à l'énergie atomique et aux énergies alternatives

May 23th, 2019
**Aristote seminar: "En route vers l'exascale!"**

# This presentation is based on the work done in ETP4HPC and in HiPEAC

Common document between
**ETP4HPC, BDVA & HiPEAC**

A blueprint for the new Strategic Research Agenda for High Performance Computing

ETP4HPC  exdci

April 2019

https://www.etp4hpc.eu/pujades/files/Blueprint%20document_20190429.pdf

## Contributors

Gabriel Antoniu, INRIA (**BDVA**)

Marc Asch, U-PICARDIE (**BDEC-2**)

Peter Bauer, ECMWF

Costas Bekas, IBM

Pascale Bernier-Bruna, **ETP4HPC**

Francois Bodin, IRISA

Laurent Cargemel, Atos

Paul Carpenter, BSC

Marc Duranton, CEA (**HiPEAC**)

Maike Gilliot, ETP4HPC

Hans-Christian Hoppe, INTEL

Jens Krueger, ITWM-FRAUNHOFER

Julian Kunkel, Univ. of Reading

Erwin Laure, KTH

Jean-Francois Lavignon, TECHNOLOGY-STRATEGY

Guy Lonsdale, SCAPOS

Michael Malms, **ETP4HPC**

Fabio Martinelli, CNR (**ECSO**)

Sai Narasimhamurthy, SEGATE

Marcin Ostasz, BSC

Maria Perez, UPM (**BDVA**)

Dirk Pleiter, JSC

Andrea Reale, IBM (**BDVA**)

Pascale Rosse-Laurent, Atos

# This presentation is based on the work done in ETP4HPC and in HiPEAC



**A blueprint for the new Strategic Research Agenda for High Performance Computing**

ETP 4 HPC    exdci

April 2019



HiPEAC

HiPEAC Vision 2019
HIGH PERFORMANCE AND EMBEDDED ARCHITECTURE AND COMPILATION

Editorial board:
Marc Duranton, Koen De Bosschere, Bart Coppens, Christian Gamrat, Madeleine Gray, Harm Munk, Emre Ozer, Tullio Vardanega, Olivier Zendra

https://www.etp4hpc.eu/pujades/files/Blueprint%20document_20190429.pdf

https://www.hipeac.net/roadmap

# Outline

1) Evolution of application scope: the continuum
2) Hardware heterogeneity and orchestration
3) Software?

# Outline

1) Evolution of application scope: the continuum
2) Hardware heterogeneity and orchestration
3) Software?

# Outline

1) Evolution of application scope: the continuum
  1) From smart sensors to HPC
  2) Artificial Intelligence (Deep Learning) loads
  3) Implications for the architecture
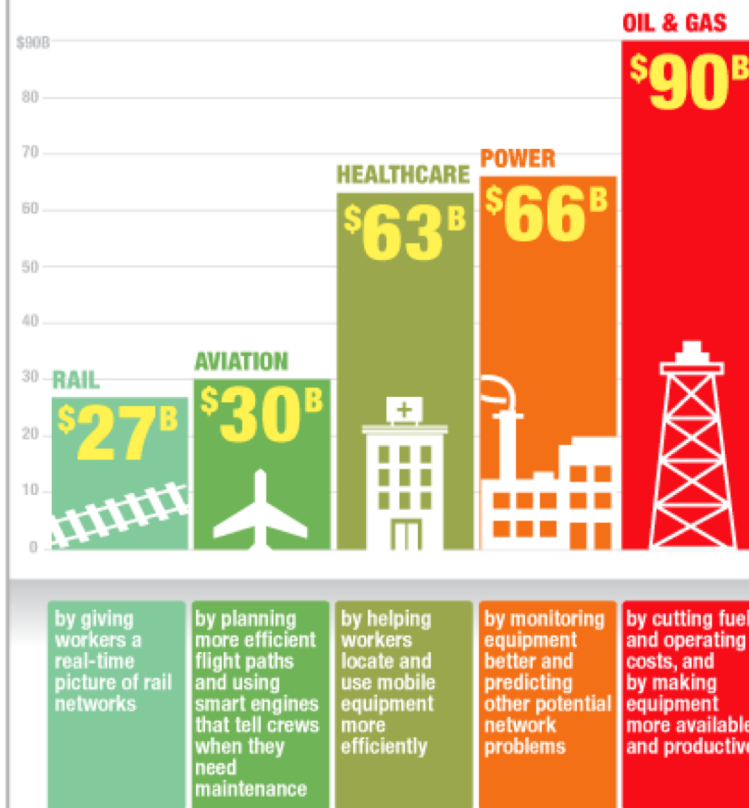2) Hardware heterogeneity and orchestration
3) Software?

# Mainstream "business" model



- Smart sensors
- Internet of Things
- Big Data
- Cloud / HPC
- Data Analytics / Cognitive computing
- New services

**HOW MUCH COULD WE SAVE WITH CONNECTED MACHINES?**

A **1%** improvement in efficiency in these five industries could add up to **$276 Billion** over 15 years

RAIL — **$27**B — by giving workers a real-time picture of rail networks

AVIATION — **$30**B — by planning more efficient flight paths and using smart engines that tell crews when they need maintenance

HEALTHCARE — **$63**B — by helping workers locate and use mobile equipment more efficiently

POWER — **$66**B — by monitoring equipment better and predicting other potential network problems

OIL & GAS — **$90**B — by cutting fuel and operating costs, and by making equipment more available and productive

# ECONOMICAL DRIVE OF CONNECTED THINGS: BETTER EFFICIENCY IN RESOURCES AND ENERGY

# Mainstream "business" model



- Smart sensors
- Internet of Things
- Big Data
- Cloud / HPC
- Data Analytics / Cognitive computing
- New services

# HPC in the loop



- Smart sensors
- Internet of Things
- Big Data
- Cloud / HPC
- Data Analytics / Cognitive computing
- New services
- Cyber Physical Entanglement

# HIGH PERFORMANCE SYSTEMS IN THE LOOP

aka Cognitive CPS
aka *Intelligent Embedded Systems*
aka *Autonomous CPS (ACPS)*

**Reduce latency**
**Safety requirements**
**Reduce bandwidth**
**Privacy constraints**
**Reduce energy**

Smart sensors

*Cyber Physical Entanglement*

Internet of Things

**New services**

**CPS + AI: Processing, Understanding *as soon as possible***

Data Analytics Cognitive computing

Big Data

Cloud / HPC

*Enabling **Intelligent** data processing at the **edge**:*
***Fog computing***
***Edge computing***
***Stream analytics***
***Fast data…***

**True** collaboration between edge devices and the HPC/cloud ⇨ creating a

**continuum**

**ENABLING EDGE INTELLIGENCE**
Transforming **data** into *information as early as possible*

# SRA-4: THE INCREASING INTERPLAY OF SIMULATION, AI, IOT AND ANALYTICS



## Societal challenges / user demand

**Science**
- energy
- life science
- weather and climate
- future materials
- fundamental sciences
- ....

**MFF2021-2027: Single Market, Digital and Innovation(*)**

Digital Europe programme
funding digital transformation
beyond 2020

Horizon Europe
Thematical clusters,
R&I missions

**Industrial users**
- energy
- aviation
- automotive
- manufacturing
- pharmaceuticals
- ....

**co-design**

Application and use scenarios    (Simulation - Analytics - AI - IOT)

Deployment in:    HPC centre    Cloud    Fog    Edge

Applications development: design, algorithms, methods, workflows

Technology infrastructure: architectures, hardware, software, I/O, storage, algorithms, programming env., tools...

Upstream technologies
- new memory/storage techno
- nanoelectronics
- photo-electronics
- ....

(*) http://www.europarl.europa.eu/RegData/etudes/
BRIE/2018/628231/EPRS_BRI(2018)628231_EN.pdf

From ETP4HPC

**Simulation**

**Data**

**Machine Learning**

Interacting with the world:
Intertwined with
*CPS* requirements

# 3 PILLARDS OF FUTURE HPC



Simulation

Data

Machine Learning

Interacting with the world: Intertwined with *CPS* requirements

# The Hype cycle - 2018

- **Deep Learning**
- **Virtual assistants**
- **DNN Asics**
- **Autonomous Driving**

# *ONE ASPECT OF AI: PERSONAL ASSISTANTS....*



Google Assistant
(1 billon devices)

Apple Siri
(+500 millions devices)

Amazon Alexa
(+100 millions devices)

Baidu's DuerOS
(+100 millions devices)

# DEEP LEARNING AND VOICE RECOGNITION

# DEEP LEARNING AND VOICE RECOGNITION

" The need for TPUs really emerged about six years ago, when we started using computationally expensive deep learning models in more and more places throughout our products. The computational expense of using these models had us worried. If we considered a scenario where **people use Google voice search for just three minutes a day** and we ran deep neural nets for our speech recognition system on the processing units we were using, we would have had to *double the number of Google data centers*!"

[https://cloudplatform.googleblog.com/2017/04/quantifying-the-performance-of-the-TPU-our-first-machine-learning-chip.html]

# GOOGLE'S CUSTOMIZED HARDWARE…

… required to increase energy efficiency
  with **accuracy adapted to the use (e.g. float 16)**



Google's TPU2 : training and inference in a **180 teraflops$_{16}$ (180 x 10$^{12}$ Flops$_{16}$)** board
(over 200W per TPU2 chip according to the size of the heat sink)

# GOOGLE'S CUSTOMIZED TPU (V2) HARDWARE…

… required to increase energy efficiency
with accuracy adapted to the use (e.g. float 16)



Google's TPU2 : 11.5 petaflops$_{16}$ of machine learning number crunching
(and guessing about **400+ KW**…, 100+ GFlops$_{16}$/W)

From Google                                          Peta = $10^{15}$ = million of milliard

# GOOGLE'S CUSTOMIZED TPU (V3) HARDWARE…



| Chip | TPUv1 | TPUv2 | TPUv3 |
|---|---|---|---|
| Announced | 2016 | May-17 | May-18 |
| Access | Internal-Only | Service Beta | Undisclosed |
| Introduction | 2015 | Feb 2018 | Undisclosed |
| Process | 28nm | 20nm est. | 16/12nm est. |
| Die Size | ~300mm2 | Undisclosed | Undisclosed |
| TOPS | 92 / 23 | 45 | 90 |
| Matrix Input | INT8 / INT16 | bfloat16 | bfloat16 |
| Memory | 8GB DDR3 | 16GB HBM | 32GB HBM |
| CPU Interface | PCIe 3.0 x16 | PCIe 3.0 x8 | PCIe 3.0 x8 est. |
| Power Consumption | 40W | 200-250W est. | 200W est. |

## A Brief Guide to Floating Point Formats

**fp32: Single-precision IEEE Floating Point Format**     Range: ~1e$^{-38}$ to ~3e$^{38}$

Exponent: 8 bits     Mantissa (Significand): 23 bits



**fp16: Half-precision IEEE Floating Point Format**     Range: ~5.96e$^{-8}$ to 65504

Exponent: 5 bits     Mantissa (Significand): 10 bits



**bfloat16: Brain Floating Point Format**     Range: ~1e$^{-38}$ to ~3e$^{38}$

Exponent: 8 bits     Mantissa (Significand): 7 bits



#io18

From https://www.nextplatform.com/2018/05/10/tearing-apart-googles-tpu-3-0-ai-coprocessor/

# ALPHAGO ZERO: SELF-PLAYING TO

# EXPONENTIAL INCREASE OF COMPUTING POWER FOR AI TRAINING

*"Since 2012, the amount of compute used in the largest AI training runs has been increasing exponentially with a 3.5 month-doubling time…*

*(by comparison, Moore's Law had an 18-month doubling period)*\**"*

**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute**



Peta= $10^{15}$

# ALWAYS MORE COMPUTING RESOURCES



HPC: exaflop $(10^{18}$ flops)

From Paul Messina, Argonne National Laboratory

**Traditional Machine Learning Workflow**



**AutoML Workflow**

From Forbes

**Auto-ML uses optimization approaches to select a "good" set of parameters "*automagically*"**

- **It is generally very computing expansive (configuration space search)**
- **Use clever algorithms to avoid exploring all the configuration space**
  **More details for example in http://automl.org**

**MACHINE LEARNING (DEEP LEARNING) LEARNING PHASE**

Specialist

Labelled data set

Specialist

Data Analytics

Big Data

Observations

Environment

- Human defines the learning data set, not the algorithm
- Large set of input data for learning phase
- Low precision floating point
- Large number of operations
- (Stochastic) gradient descent

**MACHINE LEARNING (DEEP LEARNING) INFERENCE PHASE**

- Low precision arithmetic
- Medium to low number of operations
- Co-location computing and storage ("*computing in memory*")
- Should satisfy the application non-functional requirements

Environment

Inference phase

**But for large number of inferences (users) -> more cloud like structure, high throughput**

# REINFORCEMENT LEARNING:
# DYNAMIC PROGRAMMING + DEEP LEARNING

Goal

HAL 9000

Rewards

Observations

Environment

Agent

Actions

Learns to maximize rewards

Respond to action

**REINFORCEMENT LEARNING:
DYNAMIC PROGRAMMING + DEEP LEARNING**

Goal

HAL 9000

Reward

Observations

Environment

Agent

- *Mixed* precision arithmetic
- Very high number of operations
- Large *internal* data manipulation
- Mainly co-location computing and storage ("*computing in memory*")
- High level of parallelism
- Minimization of energy functions

Actions

Learns to maximize rewards

Respond to action

# REINFORCEMENT LEARNING:
## DYNAMIC PROGRAMMING + DEEP LEARNING



Goal

HAL 9000

Rewards

Observations

Agent

Actions

Simulation

Learns to maximize rewards

Respond to action

# COMPLEMENTARITY OF SIMULATION AND IA TECHNIQUES



Simulation for improving IA

**Converged architecture:**
- From float16 to double precision
- Increasing memory per node
- Flexible partitioning
- Increase ratio communication / compute or compute in memory

IA

**Simulation**

IA for improving simulations

# Outline

1) Evolution of application scope: the continuum
2) Hardware heterogeneity and orchestration
3) Software?

# Outline

1) Evolution of application scope: the continuum
2) Hardware heterogeneity and orchestration
   1) End of Dennard's scaling
   2) Heterogeneous accelerators
   3) Heterogeneous integration
3) Software?

# END OF DENNARD'S SCALING

# WHAT WILL BE THE NEXT TECHNOLOGY?



**And after CMOS?**

# Exponential increase of performances in 33 years



To infinity and beyond…

Production car of 1985
Lamborghini Countach 5000QV
Max speed 300 Km/h

X 100 000 000
in 33 years

Star Trek Enterprise
Year: about 2290
27 times the speed of light?

# THE END OF ~~MOORE'S LAW~~ DENNARD SCALING

| Parameter (scale factor = a) | Classic Scaling |
|---|---|
| Dimensions | $1/a$ |
| Voltage | $1/a$ |
| Current | $1/a$ |
| Capacitance | $1/a$ |
| Power/Circuit | $1/a^2$ |
| Power Density | **1** |
| Delay/Circuit | $1/a$ |

Everything was easy:
- Wait for the next technology node
- Increase frequency
- Decrease Vdd
  $\Rightarrow$ Similar increase of sequential performance
  $\Rightarrow$ No need to recompile (except if architectural improvements)

Source: Krisztián Flautner "From niche to mainstream: can critical systems make the transition?"

# Technology evolution



**Transistor 2D**

Fully Depleted Silicon
on Insulator (FDSOI)
Transistor

# Technology evolution

**FDSOI**

**22FD** | **12FD** | **Next Gen**

**Non planar / trigate / stacked Nanowires**

**Silicon Quantum bits**

**FinFET**

*28nm* | *10nm* | *2019* | *5nm*

*14nm* | *2017* | *7nm*

**Disruptive scaling**

**Hybrid logic**

**Steep slope devices**

**Mechanical switches**

**Si Quantum bits**

**Alternative to scaling and diversification**

**Monolithic 3D for 3D VLSI**

# The problem:

## Expected case scenario



**Exponential power consumption**

From "Total Consumer Power Consumption Forecast", Anders S.G. Andrae, October 2017

# SPECIALIZATION LEADS TO MORE EFFICIENCY EFFICIENCY

**CPU**
1690 pJ/flop

**GPU**
140 pJ/flop

| Type of device | Energy / Operation |
|----------------|--------------------|
| CPU | 1690 pJ |
| GPU | 140 pJ |
| Fixed function | 10 pJ |

FPGA with HLS
"software programming space and not only time"

Westmere
32 nm

Kepler
28 nm

Source from Bill Dally (nVidia) « Challenges for Future Computing Systems »
HiPEAC conference 2015

## TODAY'S HPC ARE HETEREGENEOUS



**NVIDIA TESLA V100**
World's First Fused HPC and AI Processor

HIGH-SPEED NVIDIA NVLINK™ TO GPUs, IBM POWER
300GB/s NVLink

NVIDA VOLTA™ TENSOR CORE GPU
**640 TENSOR CORES**
125 TFLOPS Tensor Ops

**5120 NVIDIA CUDA® CORES**
15.7 TFLOPS FP32
7.8 TFLOPS FP64

MEMORY
32GB/16GB HBM2

TENSOR CORE GPU | 21 BILLION TRANSISTORS | 125 TFLOPS | REVOLUTIONARY HPC AND AI PERFORMANCE

- 3.3 peak exaops for emerging AI workloads
- 4,608 compute nodes, each containing two 22-core IBM Power9 processors and **six Nvidia Tesla V100 GPUs**
- Interconnected with dual-rail Mellanox EDR 100Gb/s InfiniBand.

# TOP500 #1 & #2: NVIDIA TESLA V100 GPU + IBM POWER9 CPU



Server Block Diagram
Power Systems AC922 with NVIDIA Tesla V100 with Enhanced NVLink GPUs

Source IBM

POWER

EDR/HDR InfiniBand

POWER

NVMe Flash Storage (PCI-E x8 gen 4.0)

IBM POWER9 SMP bus

Direct Attach DDR4 memory (~170GB/s BW per CPU)

PCI-Express x8 (gen 4.0) bus with CAPI for IB (12.8GB/s)
1x PCI-E x8 4.0 from each CPU to IB (multi-socket host direct)

PCI-Express x8 (gen 4.0) bus with CAPI (12.8GB/s)

25GB/s NVIDIA NVLink Interconnect (50GB/s bi-directional)
75GB/s of bandwidth between points (3 links)

- **Compute performance from GPU**

## NVIDIA TESLA V100 SPECIFICATIONS

Source NVIDIA

| | Tesla V100 for NVLink |
|---|---|
| PERFORMANCE with NVIDIA GPU Boost | DOUBLE-PRECISION 7.8 teraFLOPS |
| | SINGLE-PRECISION 15.7 teraFLOPS |
| | DEEP LEARNING 125 teraFLOPS |
| INTERCONNECT BANDWIDTH Bi-Directional | NVLINK 300 GB/s |
| MEMORY CoWoS Stacked HBM2 | CAPACITY 32/16 GB HBM2 |
| | BANDWIDTH 900 GB/s |

**➡ Heterogeneous integration driven by compute energy efficiency**

From Denis Dutoit

# AMD'S EPYC AND RADEON TO POWER EXASCALE SUPERCOMPUTER

May, 2019

- Performance target: 1.5 exaflops; 40 MW; 37 GFLOPS/W
- Compute node: x1 CPU + x4 GPU + coherent fabric



**HIGH PERFORMANCE CPU CUSTOMIZED FOR HPC**

| Custom AMD EPYC processor optimized for HPC and AI | Utilizes Future "Zen" Core High-Performance Architecture | AI-Optimized for Supercomputing Workloads |

**INFINITY FABRIC**

**HIGH PERFORMANCE GPU OPTIMIZED FOR HPC AND AI**

| HPC-Customized Compute Engines | Extensive Mixed Precision Ops for Optimum Deep Learning Performance | High-Bandwidth Memory (HBM) for Maximum Throughput |

| High-Bandwidth, Low-Latency Connection Between CPU and GPU | Custom Coherent Fabric | Connects 4:1 GPU to CPU Per Node |

➡ Heterogeneous integration requires high-bandwidth, low-latency connection

From Denis Dutoit

| 54

# 52ND EDITION OF THE TOP500 LIST (NOVEMBER 11TH, 2018)

**TOP 500**
The List.

**Heterogeneous integration ?**

| Rank | Site | System | Rmax (TFlop/s) | Rpeak (TFlop/s) | Het. |
|------|------|--------|----------------|-----------------|------|
| 1 | DOE/SC/Oak Ridge National Laboratory United States | **Summit** - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM | 143,500.0 | 200,794.9 | Yes |
| 2 | DOE/NNSA/LLNL United States | **Sierra** - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM / NVIDIA / Mellanox | 94,640.0 | 125,712.0 | Yes |
| 3 | National Supercomputing Center in Wuxi China | **Sunway TaihuLight** - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRCPC | 93,014.6 | 125,435.9 | No |
| 4 | National Super Computer Center in Guangzhou China | **Tianhe-2A** - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000 NUDT | 61,444.5 | 100,678.7 | Yes |
| 5 | Swiss National Supercomputing Centre (CSCS) Switzerland | **Piz Daint** - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100 Cray Inc. | 21,230.0 | 27,154.3 | Yes |

# Deep Learning Chipset Shipments to Increase to 2.9 Billion Units Annually by 2025, According to Tractica

**GPUs and CPUs Currently Lead in Market Share, but ASICs will Capture the Lead by 2022, with Expanded Opportunities for SoC Accelerators and FPGAs**

May 06, 2019 07:20 AM Eastern Daylight Time



Deep Learning Chipset Unit Shipments by Type, World Markets: 2018-2025

Source: Tractica

# GOING NEURO-INSPIRED: "SPIKING" NEURAL NETWORKS

Using another way of coding information…not using bits

| | IBM TrueNorth | Intel Loihi | DynapSEL |
|---|---|---|---|
| Technology | 28nm CMOS | 14 nm CMOS | 28 nm FDSOI |
| Supply Voltage | 0.7-1.05 V | 0.5-1.25 V | 0.73-1 V |
| Design Type | Digital | Digital | Mixed-signal |
| Neurons per core | 256 | Max 1k | 256 |
| Core Area | 0.094 mm$^2$ | 0.4 mm$^2$ | 0.36 mm$^2$ |
| Computation | Time multiplexing | Time multiplexing | **Parallel processing** |
| Fan In/Out | 256/256 | 16/4k | **2k/8k** |
| On-line Learning | No | Programmable | **STDP** |
| Synaptic Operation / Second / Watt | 46 GSOPS/W | | **300 GSOPS/W** |
| Energy per synaptic operation | 26 pJ | 23.6 pJ | **<2 pJ** |

University of Zurich

# FUSING PARADIGMS AT HARDWARE LEVEL

**At the hardware level, the good old Von Neumann/ CMOS partnership can act as a computing substrate, or <span style="color:red">orchestrator</span> of various accelerators/technologies**

- Acting as coordination / communication node
- Allowing Hardware / Software integration



Qubits on Silicon

Maurand et al, *Nature Com.*, Jul. 2016.

NVM Synapses on Silicon

D. Roclin et al, *IEEE NanoArch*, 2014.

Slide from Christian Gamrat

**NON VOLATILE MEMORIES**

## PCM

**GST
GeTe
GST + HfO$_2$**



*Thermal
effect*

## MR...



- Can change the structure of memory hierarchy?
+ 64/128 addressing scheme
⇒ Do we still need files?
⇒ Direct access of objects

Ag / GeS$_2$



*Electrochemical
effect*

## OXRAM

**TiN/HfO$_2$/Ti/TiN**



*Electronic effect  oxygen vacancies*

# NEW ARCHITECTURE PARADIGMS WITH NVM



From Denis Dutoit

# SOLVING THE ENERGY CHALLENGE: COST OF MOVING DATA



Source: Bill Dally, « To ExaScale and Beyond »
www.nvidia.com/content/PDF/sc_2010/theater/Dally_SC10.pdf

# PROCESSOR ARCHITECTURE EVOLUTION

**End of Dennard's scaling**

**Moore's law slow-down**

**Quantum Computing**

**Mono-core architecture** for single thread performance

**Many-core architecture** for parallelism

**Heterogeneous architecture** for energy efficiency

**Disruptive Architecture** for data management

CPU

Cache

NIC

Bus

Far Mem.

NIC

CPU | CPU
Cache | Cache
CPU | CPU
Cache | Cache
Cache

CPU | CPU
Cache | Cache
CPU | CPU
Cache | Cache
Cache

NoC + LLC

Cache | Cache
Cache | Cache
CPU | CPU
Cache | Cache
CPU | CPU

Cache | Cache
Cache | Cache
CPU | CPU
Cache | Cache
CPU | CPU

Memory

Memory

Memory

Generic processing

Close Mem.

Close Mem

Coherent Link

Close Mem

Close Mem

HW accelerator

Data Centric Interconnect

CPU | CPU | CPU | CPU
Cache | Cache | Cache | Cache

In Memory Computing, Neuromorphic Computing.

Memory invades logic

NIC
(Network InterConnect)

~2006

~2016

~2026…

From Denis Dutoit

| 62

# 2.5 stacking with chiplets and interposers for heterogeneous integration

**Benefits of 3DVLSI (some of..)**

- ➢ Scalability
- ➢ Dedicated die function
- ➢ IP re-use
- ➢ Compacity
- ➢ Performances

  - ➢ …. And cost !

Sensor / logic + Memory

Active Interposer

3D SoC

Photonic Interposer

Optical Switch

SRAMs / Logic

Non-CMOS componant / Logic

And even more to imagine

The future of 3D VLSI @ Leti

# FROM ADVANCED PACKAGING TECHNOLOGIES ….

## … TO CHIPLETS

Advanced Integration

### SiP
**Multi-Chip-Module**

Interconnect density: 100μm x 100μm

*Source: AMD EPYC 7260, 4-chiplet chip*

System-in-Package

### 3D
Die stacking

Interconnect density: 10μm x 10μm

*Source: Micron High-Bandwidth-Memory*

### 2.5D
Interposer based

Interconnect density: 10μm x 10μm

*Source: AMD Fiji GPU*

3D Integrated-Circuit (3D IC)

### CHIPLET partitioning

Custom chiplets   Commercial chiplets

COMM    RADAR EW    SIGINT

*Source: DARPA*

*Source: LETI*

*Source: GeorgiaTech*

From Denis Dutoit

## From AMD ….

## … and INTEL



22 Jun 2018 | 12:46 GMT

**AMD Tackles Coming "Chiplet" Revolution With New Chip Network Scheme**

Active silicon interposers could make for smaller, better computers, but the networks need to mesh

By Samuel K. Moore

https://spectrum.ieee.org/tech-talk/semiconductors/design/amd-tackles-coming-chiplet-revolution-with-new-chip-network-scheme

[J. Yin et al., "Modular Routing Design for Chiplet-based Systems", ISCA'2018]

FOVEROS Technology

Intel unveils a groundbreaking way to make 3D chips

"Foveros" will let Intel stack logic chips on top of each other.

Devindra Hardawar, @devindra
12.12.18 in Gadgetry

39 Comments    870 Shares

As it's getting more difficult to cram transistors next to each other in chips, and we near the end of Moore's Law, the only choice is to go vertical. Literally. That's the essence of 3D chip design, and it's the crux of a major

**2D AND 3D PACKAGING DRIVE NEW DESIGN FLEXIBILITY**

MONOLITHIC    2D INTEGRATION    3D INTEGRATION

https://www.engadget.com/2018/12/12/intel-foverus-3d-chip/?yptr=yahoo&guccounter=2

From Denis Dutoit

| 66

# EUROPEAN PROCESSOR INITIATIVE

**EPI IP's Launch Pad**
**&**
**Pan European Research**
**Platform for HPC and AI**

**Gen3 GPP Family**

2021

2022-2023

2021-2022

2024-...

**Rhea Family - Gen1 GPP**

EPI Common Platform
ARM & RISC-V
External IPs

HPC System PreExascale
Automotive PoC

**Cronos Family - Gen2 GPP**

EPI Common Platform
ARM & RISC-V

HPC System Exascale
Automotive CPU

## PROJECT PILLARS

- Common platform and global architecture stream
- HPC general purpose processor stream
- Accelerator stream
- Automotive platform stream

www.european-processor-initiative.eu

BMW GROUP    BMW    MINI    Rolls-Royce Motor Cars Limited

Atos    Infineon    JÜLICH Forschungszentrum

BSC Barcelona Supercomputing Center Centro Nacional de Supercomputación    FORTH INSTITUTE OF COMPUTER SCIENCE

KALRAY    semidynamics silicon design and verification services

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA    CHALMERS    FER UNIVERSITY OF ZAGREB FACULTY OF ELECTRICAL ENGINEERING AND COMPUTING    UNIVERSITÀ DI PISA

Fraunhofer ITWM    EXTOLL latency matters.

E4 COMPUTER ENGINEERING    CINECA    cea    ST

SURF SARA    ift TÉCNICO LISBOA    EB Elektrobit

ETHzürich    GENCI

PROVE & RUN    KIT Karlsruher Institut für Technologie    menta

# COMMON PLATFORM FOR MULTI-LEVEL HETEROGENEOUS INTEGRATION



- Heterogeneous socket
- Homogenegous/ Heterogeneous chiplet

PCIe gen5 links

HSL links

D2D links to adjacent chiplets

ARM    MPPA

eFPGA    EPAC

HBM memories

DDR memories

- Heterogeneous memories

From Denis Dutoit

- Heterogeneous SoC:
  - General processing core
  - EPAC - EPI Accelerator
  - MPPA - Multi-Purpose Processing Array
  - eFPGA - embedded FPGA

| 69

# ELECTRONS VERSUS PHOTONS

**Electrons:** **Easy to create and interface**

**Attenuation with the distance (Ohm's law)**

**Photons:** **Energy demanding for creation and interfacing**

**Low attenuation with the distance**

# OFF-CHIP PHOTONICS

Photonics: cost in sending information, nearly *nothing in transmission*

# IN-PACKAGE PHOTONICS

# DEMONSTRATION OF A THERMALLY TUNED WDM ELECTRO-OPTICAL LINK



Tx w. tuning

Rx w. tuning

Optical fiber array

CMOS+Si-Photonics
3D stack

Chip-on-board integration

Y. Thonnart & al. ISSCC'2018

➔ **1Tbps/mm²** bandwidth density
➔ Tight technology integration of
E/O ring modulators within a 3D stack
➔ Integrated thermal tuning
robust to compute fabric heating

[LETI: Y. Thonnart, ISSCC2018]



No wobulation          Wob+Thermal tuning deactivated          Wob.+Thermal tuning activated

From Denis Dutoit

| 73

# LETI'S SI-PHOTONICS ROADMAP FOR POST-EXASCALE COMPUTING



- **Target demonstrator → 2021**
- **96-core cache-coherent processor**
- **Generic E/O chiplets**
- **8-node optical NoC**
  - 576 Gbit/s aggregated bandwidth
  - 384 microring resonators
  - ~10 ns electro-optical latency

From Denis Dutoit

# POTENTIAL SOLUTION FOR POST EXASCALE BOARD



Time

SW tools, benchmarks and design methodologies

Photonic

High Density 3D

CoolCube™

New Memory Technologies

Neuromorphic

Heterogeneity & everything close

**SW tools**, benchmarks and design methodologies energy aware

Active silicon interposer, High density 3D

Photonic

New Memories (NVM) close to the logic

Neuro chiplet

Scaling with FF and CoolCube™

From Denis Dutoit

# Outline

1) Evolution of application scope: the continuum
2) Hardware heterogeneity and orchestration
3) Software?

# PARALLELISM AND SPECIALIZATION ARE NOT FOR FREE…



Frequency limit
➔ parallelism
Energy efficiency
➔ heterogeneity

Ease of programming

# PARALLELISM AND SPECIALIZATION ARE NOT FOR FREE…

# MANAGING COMPLEXITY



"And that's why we need a computer."

Cognitive solutions for complex computing systems:

- Using **AI and optimization techniques for computing systems**
  - Creating new hardware
  - Generating code
  - Optimizing systems
- Similar to *Generative design* for mechanical engineering

# USING AI FOR MAKING COMPUTING SYSTEMS: "GENERATIVE DESIGN" APPROACH

The user *only states desired goals and constraints*
-> The *complexity wall* might *prevent explaining* the solution



"Autodesk"

Motorcycle swingarm: the piece that hinges the rear wheel to the bike's frame

# EXAMPLE: DESIGN SPACE EXPLORATION FOR DESIGN MULTI-CORE PROCESSORS[1] (2010)

- **Ne-XVP project – Follow-up of the TriMedia VLIW (https://en.wikipedia.org/wiki/Ne-XVP )**

- **1,105,747,200 heterogeneous multicores in the design space**

- **2 millions years to evaluate all design points**

- **"AI inspired" techniques allowed to reduce the induction time to only few days**

**=> *x16 performance increase***



**Calculated Performance versus Area**

[1] M. Duranton et all., "Rapid Technology-Aware Design Space Exploration for Embedded HeterogeneousMultiprocessors" in Processor and System-on-Chip Simulation, Ed. R. Leupers, 2010

# THIS IS ALSO VALID FOR SOFTWARE: AUTOML AND OTHER PROGRAM GENERATORS



Communications of the ACM, 55(2), pp. 70–80, February 2012

www.prog-by-opt.net

## Microsoft's AI is learning to write code by itself, not steal it

Written by Dave Gershgorn

What if instead of searching through menus within programs like Microsoft Excel, our computers could understand the problem we're trying to solve and write the software to solve it? It's a hyper-futuristic idea, but one that has recently seen progress from Microsoft Research and the University of Cambridge.
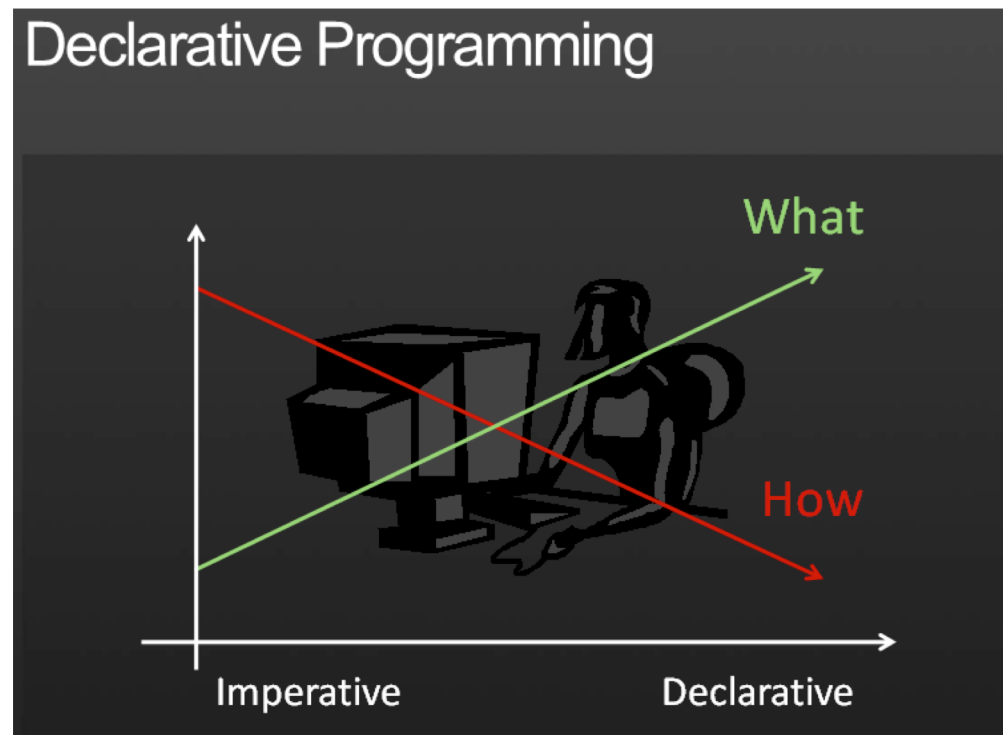
In a November 2016 paper (pdf), which gained notoriety after being accepted into one of the year's largest artificial intelligence conferences, Microsoft and Cambridge built an algorithm capable of writing code that would solve simple math problems. The algorithm, named DeepCoder, would be able to augment its own ability by also looking at potential combinations of code for how a problem could be solved. (It's a bit complicated; we'll break it down later.) However, this doesn't mean it steals code, or copy and pastes it from existing software, or searches the internet for solutions, as some reports have claimed.

# PROGRAMMING 2.0: LET THE COMPUTER DO THE JOB

**Describing *what* the program should accomplish, rather than describing *how* to accomplish it**

- For example, describe the *concurrency* of an application, not how to parallelize the code for it.
- (Good) compilers know better about architecture than humans, they are better at optimizing code…

**CONCLUSION**

# CONCLUSION: WE LIVE AN EXCITING TIME!



*"The best way to predict the future is to invent it."*

*Alan Kay*

# Thank you for your attention

Special thank you to Denis Dutoit, Christian Gamrat,
Carlo Reita for their slides I borrowed.