

EUROSERVER: Share-Anything Scale-Out Micro-Server Design

Manolis Marazakis
Institute of Computer Science, FORTH,
Heraklion, Greece
Email: maraz@ics.forth.gr

John Goodacre
ARM Ltd
Cambridge, UK
Email: john.goodacre@arm.com

Didier Fuin
STM SA
Grenoble, France
Email: didier.fuin@st.com

Paul Carpenter
BSC
Barcelona, Spain
Email: paul.carpenter@bsc.es

John Thomson
ONAPP Ltd
Cambridge, UK
Email: john.thomson@onapp.com

Emil Matus
Technical University of Dresden
Dresden, Germany
Email: emil.matus@ifn.et.tu-dresden.de

Antimo Bruno
NEAT Srl
Rome, Italy
Email: bruno@neat.it

Per Stenstrom
Chalmers University of Technology,
Gothenburg, Sweden
Email: per.stenstrom@chalmers.se

Jerome Martin, Yves Durand, Isabelle Dor
CEA
Grenoble, France
Email: {jerome.martin,yves.durand,isabelle.dor}@cea.fr

Abstract—This paper provides a snapshot summary of the trends in the area of micro-server development and their application in the broader enterprise and cloud markets. Focusing on the technology aspects, we provide an understanding of these trends and specifically the differentiation and uniqueness of the approach being adopted by the EUROSERVER FP7 project. The unique technical contributions of EUROSERVER range from the fundamental system *compute unit* design architecture, through to the implementation approach both at the *chiplet* nanotechnological integration, and the *everything-close* physical form factor. Furthermore, we offer optimizations at the virtualisation layer to exploit the unique hardware features, and other framework optimizations, including exploiting the hardware capabilities at the run-time system and application layers.

I. INTRODUCTION

There is growing business interest in *microservers*, i.e. clusters of high-density, low-power servers, which are suited to the growing number of hyperscale workloads found in modern data-centres [1]. Although still in their infancy and not yet widely used in production, microservers show promise of allowing the total compute, network, and storage resource capacity of a data-centre to be utilized with high flexibility and efficiency for a wide range of diverse workloads.

The EUROSERVER FP7 Project Consortium develops a micro-server solution tailored for the workloads running on today's cloud infrastructure [2]. The EUROSERVER platform combines several architectural key features, such as highly efficient ARMv8 processors, an innovative scalable memory scheme called *Unimem*, a lightweight hypervisor platform suitable for micro-servers called *MicroVisor* and the use of Hybrid Memory Cube (HMC) technology to maximize memory storage density and bandwidth. EUROSERVER adopts an innovative structure based on interconnected compute “coherence islands” for an optimal balance between data locality and transfer efficiency. The compute SoC internal structure is organized around several independent *chipllets* implementing the

coherence islands. The coupling between chiplets is realized via high-speed serial links. Physically, the system is integrated onto a cost-effective organic interposer solution.

The emerging key differentiator for EUROSERVER is improved resource utilisation [3]. Just as cloud computing and virtualisation enables companies to consolidate workloads from many distributed and under-utilised hardware platforms into smaller numbers of servers, EUROSERVER proposes to more efficiently exploit micro-server and low-power hardware in order to pave the way towards the next generation of more power-efficient servers.

The rest of this paper is organized as follows. Section II presents a review of industry trends towards the development of commercially viable microservers. In Section III we present the system approach of the EUROSERVER consortium to address both technical and economic challenges so as to enable commercially viable devices that can be also optimized for market specific requirements. We also outline our software stack that supports cloud workloads on top of a microserver platform. Finally, Section IV summarizes our conclusions.

II. REVIEW OF INDUSTRY TRENDS

In this section we present a snapshot summary of the trends in the area of micro-server development and their applications in the broader enterprise and cloud markets.

The most dominant trend across the server market continues to be based on adaptation of technology around the historic PC architecture based platform. This platform has evolved over the decades since it was first introduced in support of the early personal computer, but fundamentally, even the most modern server still follows the same model:

- The platform is complete, in the sense that it can only exist in a complete form. The balance between the processing/memory and I/O are fundamentally defined

TABLE I
EUROSERVER COMPETITIVE LANDSCAPE: CHARACTERISTICS OF CURRENT ARM-BASED (64-BIT) DEVICES.

Company	Product	Location	Silicon Tech	# Cores	Core Type	Chip power	Memory support	I/O support	Chip-to-chip support	Status	Ref
ARM	JUNO SoC	UK	TSMC 28HPM	6	Dual-core A57 + quad-core A53 + quad-core Mali-T624	N/A	2x channel DDR3L 32-bit 1600	USB 2.0, Other I/O in IOFPGA: 10/100 Eth, PCIe not functional (r1)	Yes, TLX400, total BW 0.7 Gbps	Available	[4]
EURO-SERVER	HW Prototype	Europe	FDSOI 28nm	32	A53	33W max.	4x HMC links: 4x60GB/s	FPGA I/F, 4x10 GB/s towards FPGA (I/O)	Yes, Unimem, RDMA	Jan. 2016	[5]
Freescale	LS1043A	U.S.A	TSMC 28nm HPM	4	A53	8W	32-bit DDR3L/4	6x 10G Eth, PCIe Gen2, USB3, SATA3	No	Sampling Q1 2015	[6]
Freescale	LS2085A	U.S.A	TSMC 28nm HPM	8	A57	30W TDP	2x 64-bit DDR4	8x 10G Eth, PCIe Gen3, USB3, SATA3	No	Q3 2015	[7]
Broadcom	Vulcan	U.S.A	TSMC 16nm FinFET	32	ARMv8 custom (Vulcan)		DDR4	PCIe	Yes, coherent	1H 2016	[8]
Cavium	ThunderX CN88xx	U.S.A	GF 28nm HKMG	24 to 48	ARMv8 custom	<100W TDP	4x 72-bit DDR3/4	10/40/100G Eth, SATA3, PCIe Gen3	CCPI Dual Socket NUMA	Available	[9]
Applied Micro	X-Gene1	U.S.A	TSMC 40nm	8	ARMv8 custom	40W	DDR3	4x 10G Eth, 6x PCIe, SATA	No	Available	[10]
Applied Micro	X-Gene2	U.S.A	TSMC 28nm	16	ARMv8 custom	25W?	4x DDR3L/4	PCIe Gen3, 1/10Gbe, SATA3	RoCE	Q3 2015	[11]
Applied Micro	X-Gene3	U.S.A	TSMC 16nm FinFET	32	ARMv8 custom		8x DDR4	42x PCI G3		2H 2016	
AMD	Opteron A1100	U.S.A	GF 28nm HKMG	8	A57	<45W	DDR3/DDR4	2x 10G Eth, 8x PCIe Gen3, 8x SATA3	CCN-504?	March 2015	[12]
HiSilicon	Phosphor-V660 Hip05	China	TSMC 16nm FinFET	2x 32	A57		DDR3	3D integration with IO chip	No	Engineering samples	[13]
EZCHIP (Tilera)	TILE-Mx	Israel (U.S.A)	28nm HPM	100	A53	70W	DDR4	1Gbe to 100GbE, PCIe Gen3	Yes (Interlaken) on-chip coherency	Sample H2 2016	[14]

by the general processor capabilities with only limited configuration through cache capacity and core frequency.

- The use of a general purpose processor acting as the master of the platform. As such, the processor is designed as a general solution to address many markets, where the economics define its capabilities, being driven still today by the consumer market, not the server market. This leads to an imbalance between, for example, the cost of delivering the SIMD media capability for consumer markets and the cost of delivering DRAM capacity for server markets.
- Scalability of the solution is through duplication of the entire platform. This leads to duplication of costly support infrastructure. To enable communication between CPUs, the software must make a communication path through high latency, generalized-abstraction interfaces.

With the advancement in the capabilities of ARM-based processors, there are a number of ARM technology licensees targeting server solutions, and in particular the micro-server based markets. Although ARM is European IP, none of the other ARM-based server chips are being built by European companies (see Table I). EUROSERVER is the first and only European processor for micro-servers.

Table I summarises essential characteristics of current ARM-based devices. The primary differentiation of these

ARM-based devices is the use of low-power ARM processor implementations. To further enhance their low-power capabilities, these devices also apply direct integration of application-specific interfaces for their target market, such as one or more Gigabit or 10 Gigabit Ethernet interfaces and/or SATA interfaces. Although such integration may make these devices more interesting to a specific market that can be addressed by its specific Processor + Memory + I/O combination, if successful, these devices will tend towards a general-purpose solution due to economic pressures.

It is well documented that the non-recurring engineering cost (NRE) to develop a current competitive semiconductor device is in the range of \$50–100M USD [15]. The correction of a design fault or a re-spin to address a yield issue will cost many additional \$10Ms USD. With such a significant NRE cost, the total available market (TAM) for which the device is suitable must extend to many millions of units, to enable the potential for the specific device to capture a substantial share of that market. With the annual TAM of devices for servers between 10 and 20 million, it is clear to see why today the consumer market of around 200-300 million devices dominates the characteristics of today's server processor. Our integration approach does not enforce a fixed compute-I/O-DRAM configuration for devices, but rather encourages a business model of market-specific configurations.

Table II compares currently available low-power processors, comparing not only ARM-based processor designs, but also recent low-power designs by Intel. As ARM-based implementations have moved towards the performance levels necessary for server hardware, Intel in recent months have released a variety of low-power Xeon variants that look to challenge ARM on the energy efficiency front. For completeness, we have also surveyed the energy efficiency of Intel Xeon Phi co-processors. Although not directly targeted to mainstream servers, Intel’s Knights Corner and Knights Landing designs offer a large number of cores and energy-efficiency levels that indicate possibilities for future servers.

TABLE II
ENERGY EFFICIENT COMPUTE PLATFORMS

Processor	Frequency	Cores	Threads	Fab.	Cost	TDP	Watts per Thread	Status
A15	1.9GHz	1-4	1-4	28nm	<\$15	1.5-3.7W	0.38-0.93W	2H 2012
A53	1.6GHz	4-8	4-8	28nm	\$17-25	1.7-6.2W	0.43-1.55W	1H 2015
A57	2.1GHz	4-8	4-8	20nm	\$30-60	2.0-7.1W	0.5-1.78W	1H 2015
A72	2.3GHz	4-8	4-8	16nm	N/A	1.5-5.8W	0.4-1.45W	4Q 2015
X-Gene1	2.4GHz	8	8	40nm	\$316	45W TDP	5.63W	2H 2015
X-Gene2	2.8GHz	8	8	28nm	N/A	45W TDP	5.63W	2H 2016
Pentium N3700	1.6GHz	4	4	14nm	\$161	6W	1.5W	Q1 2015
Atom C2758	2.4GHz	8	8	22nm	\$208	20W	2.5W	Q3 2013
Xeon D-1540	2GHz	8	16	14nm	\$581	45W TDP	2.81W	Q1 2015
Xeon E-1240Lv3	2GHz	4	8	22nm	\$278	25W TDP	3.13W	Q2 2014
ThunderX	2.5GHz	24-48	24-48	28nm	N/A	80W TDP	1.67W	1H 2016
Knights Corner (Phi)	1.1-1.238 GHz	57-61	244	22nm	\$1695 to \$4235	300W TDP	1.23W	Q2 2013
Knights Landing (Phi)	1.2-1.3GHz est.	72	288	14nm	N/A	160-215W TDP	0.75W	1H 2016

III. THE EUROSERVER SYSTEM ARCHITECTURE

The EUROSERVER consortium is developing a system-level approach to address both technical and economic challenges so to enable commercially viable devices that can be also optimized for market specific requirements. In this section we describe the key design properties of the EUROSERVER approach: compute unit scalability, share-anything scale-out, scalable shared memory architecture, chiplet-based device design, virtualisation enhancements and memory optimization.

A. Compute Unit Scalability

One of the foundational new innovations adopted by the EUROSERVER approach is to consider the processor not as the master core of a platform, but as a unit of system scalability for compute and memory. Unlike processors in today’s systems, the compute unit defines a compute plus memory structure in which the processor has ownership of its entire local memory address space, however it also permits other compute units in the system to access its memory through a

new global memory address space, independent of its local memory through a system-based memory management unit. The compute unit can therefore both address its entire local physical memory space, and also map a portion of this space into the globally shared address space, and as such map reads and writes to the local memory space of other compute units directly from within the virtual address space of an application or OS runtime.

Since each compute unit owns its own local memory, the total memory of a compute unit based system is limited only by the number of compute units and the physical addressability of the compute unit. Therefore, unlike other NUMA-based systems, the total shared system memory is not limited by the physical address bus of a single processor.

Since the compute unit also allows remote access to its local memory through a system-based memory management unit, the hardware also supports direct memory access between the memories of multiple compute units, with a level of coherence that removes the need for a costly coherence protocol to be supported between units. This significantly increases the scalability of the shared memory system paradigm. At the same time, remote access is provided with direct processor interconnect-level latencies, opening new opportunities and capabilities for shared memory optimizations in application and operating system frameworks.

Work is also being carried out on a novel light-weight hypervisor platform called the *MicroVisor*. This hypervisor platform takes the notion of compute unit scalability and distributed compute units to the next level by removing the overhead of a single control domain and making the hypervisor more suitable for micro-server type systems that will share all the fundamental resources available in the platform, leveraging any hardware assistance that might be available.

The high scalability of the EUROSERVER architecture enables the use of “little” ARM Cortex-A53 cores, which have a much lower pJ/operation ratio than the “big” Cortex-A57. As shown in Table I, EUROSERVER is one of just three approaches based on Cortex-A53 level performance. The other platforms are the Freescale LS1043A, which is targeted at embedded networking, and the Cavium ThunderX. The EUROSERVER use of 28 nm fully-depleted silicon-on-insulator (FDSOI) also permits the Cortex-A53 to be operated at a higher frequency, allowing it to also provide high single thread performance. To further quantify this energy efficiency, the Cortex-A53 uses approximately half the power and half the area of the Cortex A57, with only a 50% reduction in instructions per cycle of the Cortex-A57, leading to an efficiency gain for a given workload of around three times [16].

B. Share-Anything Scale-Out

With the compute and memory now defined within a scalable unit of compute, the third dimension of system configuration is the interface and I/O, which is traditionally associated along with the CPU and memory, in what is, for others, a fixed platform’s specific configuration. Unlike the traditional approach, EUROSERVER allows the compute unit

and the associated level of memory required to be scaled independently of any I/O or interfaces: our system architecture adopts a ‘share-anything’ approach for such peripheral components. This share-anything approach allows the compute unit to be scaled independently of the level of I/O and interface technology required by the target market. Using the direct memory access capability of each compute unit, the shared I/O devices appear as directly integrated into the local memory of the processor. This is fundamentally different from any of today’s virtualised I/O sharing schemes, in that the I/O payload does not need to be buffered or transferred from the memory owned by the hypervisor, and as such is moved in hardware directly from the wire to the local memory of the compute unit(s) hosting the interface. This separation of scalability between compute and I/O is also a key characteristic that reduces the cost and energy consumption relative to the traditional share-nothing cluster used in today’s servers. For example, a EUROSERVER based solution can scale the compute and I/O independently to match the market requirements without having to duplicate unnecessary resources.

C. Unimem – Scalable Shared Memory Architecture

The EUROSERVER approach provides a new memory model across the server system, known as *Unimem* and introduced by our consortium. This model breaks the traditional dual memory types available in today’s systems: cached memory versus device memory. To scale the cached memory type, today’s hardware system must implement a coherence protocol that both limits performance and consumes energy. This places a hard limit on the ability of applications and OS to use the shared memory paradigm, even though this is often the best solution to access shared state across a wider system. The software must then move to a device memory type and implement a communications paradigm in order to share any state between various processor systems.

The key observation is that applications in data-centres tend to partition their datasets across servers, with the expectation that server processors and their caches will be placed near the dataset of a particular application task. Rather than moving datasets around at great energy expense in terms of moving data and imposing cache coherency requirements, it is more energy efficient to keep the dataset still and move the task to a processor that is near the required dataset. Thus a software-centric architectural approach is encouraged. A related recent effort has been demonstrated by the TERAFLUX [17] FP7 project, which investigated the use of a unified address space and means to achieve scalable memory without global coherence, and proposed a programming and execution model based on data-flow instead of traditional control-flow.

The Unimem approach not only maintains a consistent and coherent access from each compute node to its local DRAM, but also manages access to the system-wide memory resources. Most importantly, it can be implemented using available ARM technology with little additional hardware overhead. In addition to supporting the device communication schemes, the Unimem memory architecture also supports

shared memory accesses over the same scale as the I/O based communication, either through direct CPU read/write, or, for larger quantities of data, through remote DMA. Since this Unimem sharing path is directly from within the hardware memory system of the compute units, the EUROSERVER consortium has also demonstrated that moving the traditional TCP/IP based I/O communication stack directly on top of the Unimem RDMA stack can accelerate existing TCP/IP based application communication, such as in Hadoop, by half to one order of magnitude. The majority of data-centre applications are I/O [18] and memory-bound [19], and should therefore benefit from remote DRAM borrowing via Unimem. Ensuring backwards compatibility of the EUROSERVER approach is key in supporting the existing application models and software investments, however, Unimem’s raw capability also opens up significant opportunity for further optimizations in future runtimes and application frameworks.

D. Everything Physically Close

Since the EUROSERVER project also considers the deployment of the new scalable hardware devices into market-specific form factors, it also has the opportunity to consider the costs and implications of the physical structure of a server solution. It is becoming well known that, as silicon geometries have shrunk, so has the energy cost associated with processing the data within the silicon device. The energy consumed by moving data between silicon devices and the system I/O has increased in relative importance.

With today’s processors being able to manipulate and move data within the device at around 10 pJ/operation, while, for example, PCIe consumes 100 pJ/bit of data, there is a significant challenge to the overall energy consumption of any solution whose I/O or storage utilize a shared PCIe bus or other “backplane” level bus. The fundamental cause for energy consumption when moving data is the physical Resistance and Capacitance (RC) that increases the energy required to transmit data at a given bandwidth when the distance increases. This is one area of low-power design that today is already leading to the integration of I/O interfaces natively inside the silicon device. Such integrations will continue, but under constraints due to NRE cost concerns and device generalization.

The EUROSERVER approach also brings the I/O into the silicon device, in order to remove the energy overheads of a backplane bus, but the nano-technological device integration in EUROSERVER extends this further, also enabling cost effective connectivity between compute units and I/O interface. Although the supporting components for mass storage today still use shared backplane bus technologies, and hence are high energy consumers, the modular approach to hardware and software of EUROSERVER will enable storage devices to be positioned a few millimetres away from the compute, and will thus further reduce the power consumption of such systems. The economic silo between the computing and storage in today’s server vendors makes such a step difficult and the consortium is not aware of any other solution using an “everything close” physical design paradigm.

E. Silicon Chiplet-based Devices

Since the EUROSERVER system architecture has defined an independent scalability model between memory capacity, compute and I/O, it is possible to design and build an application specific device at this level of abstraction rather than at the level of a standalone PC architectural platform. Such application-specific design is not unique and has been a defining feature of ARM-based SoCs for generations. The EUROSERVER approach is investigating new approaches to System-in-Package designs that enable application-specific optimization, but with an economic model that reduces the cost of such device by an order of magnitude while potentially increasing the overall capabilities of a silicon device. The Chiplet approach is one where, rather than building a large costly monolithic silicon die, the capabilities of the device are split across multiple, possibly duplicated, dies. The Compute Unit is well suited for such a duplicated die approach, and Figure 1 shows the implication of moving from a single die of four compute units plus I/O, to a die containing the I/O, and supporting four chiplet instances of the compute unit, with technology nodes and integration technology (active interposer) in line with the next-generation EUROSERVER architecture. The chiplet approach can enhance the single-die yield compared with building a single large SoC. The chiplet approach can enhance the single-die yield compared with building a single large SoC. As a consequence, the example in Figure 1 shows a $20\times$ manufacturing cost reduction for the same compute and I/O capabilities dependent on specific yielding and actual wafer cost.

F. Virtualisation Enhancements

The virtualisation of processor, memory, and platform resources has matured around the structure of the PC architecture platform and as a solution to address the misbalance of the platform. With the EUROSERVER promise of solutions, in which the platform is more balanced to the application requirements, the consortium has undertaken to enhance the virtualisation manager capabilities and prepare for such platforms. This work includes lowering the cost of device abstraction and sharing, which is key to enabling efficient use of the EUROSERVER share-anything approach to resources. Another unique aspect of this work is the integration of efficient migration of task data in the extended shared-memory approach and RDMA available in EUROSERVER.

G. Memory Optimization

Another foreseen benefit in systems adopting the EUROSERVER approach is the capability to move less data over any given distance, as a research goal beyond the everything-close design paradigm of EUROSERVER. With the continued scaling of silicon, the cost of data movement to memory will become the next major power bottleneck after I/O is brought close. The EUROSERVER consortium therefore sees great value in creating a solution for memory compression, a scheme by which the larger energy costs of data movement are addressed through additional but smaller costs associated

with compressing the data before its transfer to memory. It is only when a solution, such as EUROSERVER, has addressed the everything-close and the shareanything approaches, that the benefits of memory compression will be most evident. As such, these developments are unique in these regards.

H. Software Stack

Figure 2 outlines the software stack under development and evaluation for the EUROSERVER microserver prototypes. Several areas are being addressed: (1) virtualization via a distributed hypervisor (*MicroVisor*), (2) I/O virtualization for memory, network, and storage resources, (3) efficient resource sharing, (4) optimizations for in-memory workloads, (5) power-aware scheduling, and (6) memory space optimization.

IV. CONCLUSIONS

We have described a number of unique developments and strategies for the evolution of future micro-servers. As EUROSERVER is the first and only processor for micro-servers being developed in Europe, the project is important to the development of a globally competitive solution within Europe.

The unique technical contributions of EUROSERVER range from the fundamental system *compute unit* design architecture, through to the implementation approach both at the *chiplet* nanotechnological integration, and the *everything-close* physical form factor. This paper has described how these unique hardware approaches to system design have enabled a new *Unimem* application scalability model and *share-anything* efficiencies well beyond the known state of the art. Together with optimizations at the virtualisation layer to exploit the unique hardware features, and other framework optimizations, including exploiting the hardware capabilities at the run-time system and application layers, the EUROSERVER project is well placed to enable significant technological and economic exploitation of its findings and deliverables.

ACKNOWLEDGMENTS

We thankfully acknowledge the support of the European Commission under the 7th Framework Programs through the EUROSERVER project (FP7-ICT-610456).

REFERENCES

- [1] N. Heath, "Microservers: What you need to know," ZDnet, <http://www.zdnet.com/article/microservers-what-you-need-to-know>, 2014.
- [2] C. Kachris, "Microservers Brew in Europe's Labs," http://www.eetimes.com/author.asp?section_id=36&doc_id=1324294, 2014.
- [3] M. Coppola, B. Falsafi, J. Goodacre, and G. Kornaros, "From Embedded Multi-core SoCs to Scale-out Processors," in *Proceedings of the Conference on Design, Automation and Test (DATE) in Europe*, 2013.
- [4] "ARM Ltd, V2M-Juno Technical Reference Manual," <http://www.arm.com>.
- [5] "EUROSERVER FP7 project web-site," <http://www.euroserver-project.eu>.
- [6] "Freescale Semiconductor Inc, LS1043A: QorIQ LS1043A and LS1023A Communication Processors," http://www.freescale.com/webapp/sps/site/prod_summary.jsp?code=LS1043A.
- [7] "Freescale Semiconductor Inc, LS2085A: QorIQ LS2045A and LS2085A Multicore Communications Processors," http://www.freescale.com/webapp/sps/site/prod_summary.jsp?code=LS2085A.

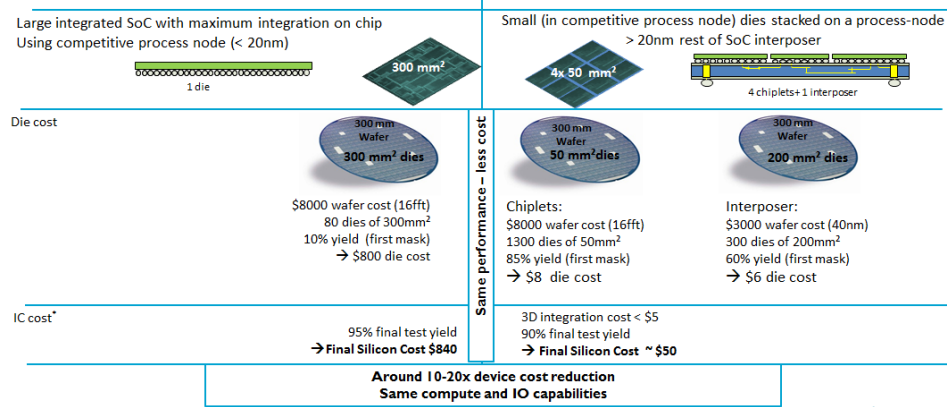


Fig. 1. Large integrated SoC vs Interposer.

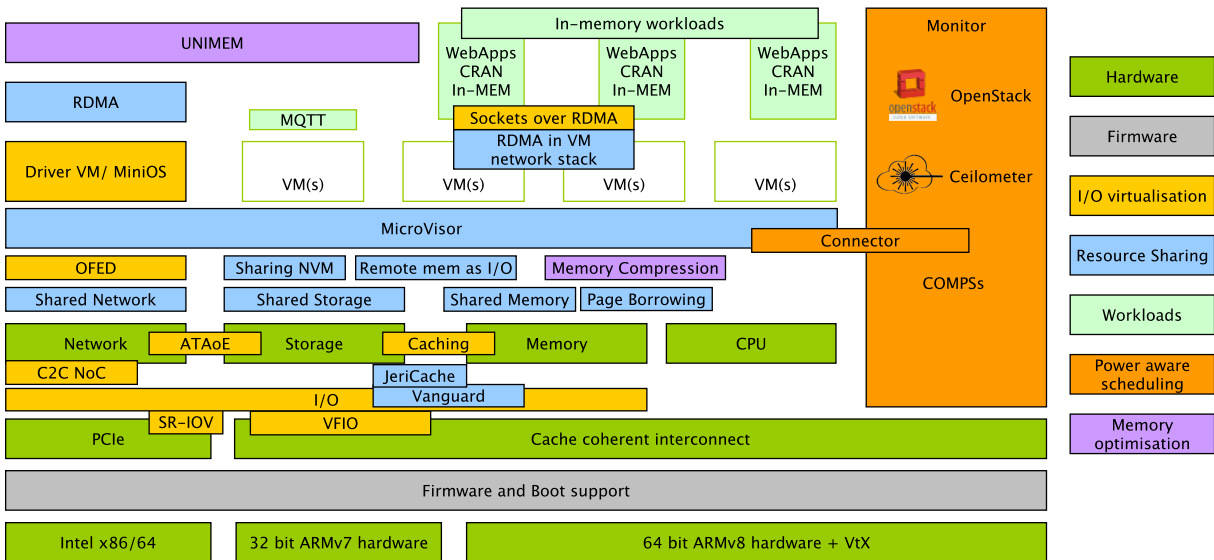


Fig. 2. EUROSERVER software stack.

- [8] "Broadcom Corp, Presentation of Vulcan at the IDC HPC USER FORUM," <https://hpcuserforum.com/presentations/santafe2014/Broadcom Monday night.pdf>, 2014.
- [9] "Cavium Inc, ThunderX ARM Processors," http://www.cavium.com/ThunderX_ARM_Processors.html.
- [10] "Applied Micro-Circuits Corp, X-Gene ARMv8 Server SoC," <https://www.apm.com/products/data-center/x-gene-family/x-gene/>.
- [11] "Applied Micro-Circuits Corp, Presentation of X-Gene2 at HotChips," http://www.hotchips.org/wp-content/uploads/hc_archives/hc26/Hc26-11-day1-epub/Hc26.11-4-ARM-Servers-epub/Hc26.11.430-X-Gene-Singh-AppMicro-HotChips-2014-v5.pdf, 2014.
- [12] "Advanced Micro Devices Inc, Presentation of the AMD Opteron A1100 "Seattle" Processor at HotChips," http://www.hotchips.org/wp-content/uploads/hc_archives/hc26/Hc26-11-day1-epub/Hc26.11-4-ARM-Servers-epub/Hc26.11.410-Opteron-Seattle-White-AMD-HotChipsAMDSeattle_FINAL.pdf, 2014.
- [13] "Taiwan Semiconductor Manufacturing Company Ltd, Announcement of 32-core 16FinFET ARM Cortex-A57 Networking Processor," <http://www.tsmc.com/tsmcdotcom/PRListingNewsArchivesAction.do?action=detail&newsid=THPGGOGO&language=E>, 2014.
- [14] "EZchip Semiconductor Inc, Announcement of the 100-core TILE-Mx100 ARMv8 Processor," <http://www.ezchip.com/News/PressRelease/?ezchip=97>.
- [15] I. B. S. (IBS), "Market Analysis and Key Trends from FD SOI Perspective," http://www.soiconsortium.org/fully-depletedsoi/presentations/september2014fdsoiforum/I9-FDSOIFORUM9.2214_Handel Jones.pdf, 2014.
- [16] J. Goodacre, "The Homogeneity of Architecture in a Heterogeneous World (keynote)," in *Proceedings of the International Conference on Embedded Computer Systems: Architectures, Modeling and Simulation (SAMOS XII)*, 2012.
- [17] R. Giorgi, R. M. Badia, F. Bodin, A. Choen, P. Euripidou, P. Faraboschi, B. Fechner, G. R. Gao, A. Garbade, R. Gayatri, S. Girbal, D. Goodman, B. Khan, S. Koliai, J. Landwehr, M. L. Nhat, F. Li, M. Lujan, A. Mendelson, L. Morin, N. Navarro, T. Patejko, A. Pop, P. Trancoso, T. Ungerer, I. Watson, S. Weis, S. Zuckerman, and M. Valero, "TER-AFLUX: Harnessing dataflow in next generation teradevices," *Journal of Microprocessors and Microsystems: Embedded Hardware Design (MICPRO)*, vol. 38, no. 8, pp. 976-990, 2014.
- [18] M. Awasthi, T. Suri, Z. Guz, A. Shayesteh, M. Ghosh, and V. Balakrishnan, "System-level characterization of datacenter applications," in *Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*. ACM, 2015, pp. 27-38.
- [19] R. Pierce, "Memory solutions for high performance computing," 2015. [Online]. Available: http://memcon.com/pdfs/proceedings2015/TDT104_ALTERA.pdf