

# EuroCloud: Energy-conscious 3D Server-on-Chip for Green Cloud Services

Emre Özer<sup>α</sup>, Krisztián Flautner<sup>α</sup>, Sachin Idgunji<sup>α</sup>, Ali Saidi<sup>α</sup>, Yiannakis Sazeides<sup>β</sup>,  
Bushra Ahsan<sup>β</sup>, Nikolas Ladas<sup>β</sup>, Chrysostomos Nicopoulos<sup>β</sup>, Isidoros Sideris<sup>β</sup>, Babak  
Falsafi<sup>γ</sup>, Almutaz Adileh<sup>γ</sup>, Michael Ferdman<sup>γ</sup>, Pejman Lotfi-Kamran<sup>γ</sup>, Mika  
Kuulusa<sup>δ</sup>, Pol Marchal<sup>ε</sup>, Nikolas Minas<sup>ε</sup>

<sup>α</sup> ARM, <sup>β</sup> University of Cyprus, <sup>γ</sup> EPFL, <sup>δ</sup> Nokia and <sup>ε</sup> IMEC

Cloud Computing (CC) is radically changing the way information technology (IT) is provided and used. Applications are delivered as services over the Internet on-demand based on Data Centers that cost-effectively consolidate and manage the Cloud servers. Already CC services are provided to hundreds of millions of users by millions of servers in the Data Centers. As increasingly more IT services are moved to the Cloud, the demand for new Data Centers and, consequently, for servers increases at phenomenal rates.

Data Center construction and operating costs, currently often in excess of tens of millions of dollars, have alarmed the public about the negative environmental impact of Data Centres and created an impetus for the design of low-cost, compact, energy-efficient Data Centers with much less carbon footprint. Server equipment costs and unprecedented electricity bills due to power consumption [1] are the main contributing factor in the total cost of a Data Center. Furthermore, microprocessors and the memory systems are both the most costly and power-consuming components in a server [2,3].

In the past decades, microprocessor designers have harnessed power growth by reducing transistor supply voltages. Unfortunately, voltage scaling is hitting the diminishing returns because further reductions in voltage will eventually cause increases in power due to larger transistor leakage currents. In the absence of voltage scaling, future chips may either need to include many low-power processor cores or a new class of "embedded-based" server architectures to sustain the scalability of future Data Centers.

The EuroCloud project addresses this need by proposing a 3D "Server-on-Chip" that provides a very dense low power server using many ARM cores, hardware accelerators and integrated 3D DRAM. EuroCloud represents a pioneering approach since for the first time an ARM based 3D stacked chip will be built grounds up aiming to address the needs of future Cloud servers.

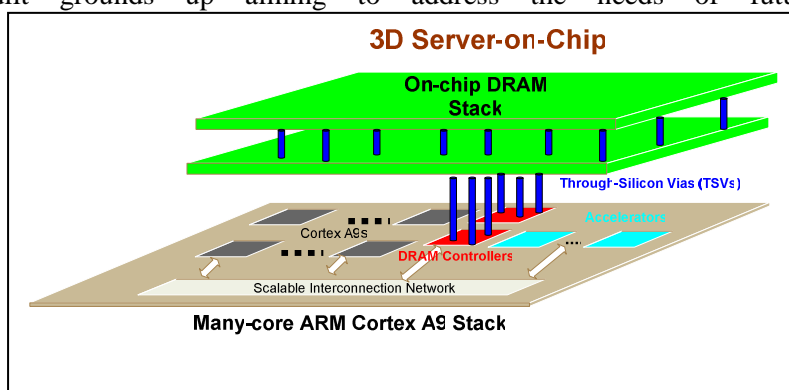


Figure 1 EuroCloud 3D Server-on-Chip concept

The rationale behind the EuroCloud concept is that a number of existing (e.g. web and database) and emerging Cloud computing applications (e.g., video/music streaming, song recognition and data mining) favour an increase in thread- and memory-level parallelism and benefit little from instruction-level parallelism. Each client request is serviced by a single or multiple independent threads each running on a dedicated processor core. These applications need many simple processor cores with high-bandwidth/low-latency access to very large memories. Using standard off-chip DRAM is

bandwidth-constrained due to limited pin count, slow due to chip crossing, and power-hungry due to I/O pads and driving circuitry. To eliminate these inefficiencies and address these issues, we propose to 3D stack DRAM chip on top of the ARM Cortex A9 [4] processor cores and hardware accelerators (e.g. GPU, video and crypto engines) as shown in **Figure 1**.

Our vision and current developments point to a major shift to novel server models using embedded processors around 2011 timeframe, and this shift will slow down the growth of PC-based commodity servers. These embedded-based servers will be equipped with off-chip DRAM, and emerge in the server market to meet the demand for energy-efficiency and green Cloud services. Our forecast is that 3D Server-on-Chips, e.g. EuroCloud server chip, based on many-core embedded processors will emerge in the server market around 2016 providing many low-cost low-power server processor cores with high-bandwidth on-chip memory subsystem and smaller form factor.

EuroCloud is realized in the context of a three year (2010-2012) project funded by European Union Framework Program 7 with the following five clear objectives:

- 1) **Workload characterization of Mobile Cloud applications:** Characterizing the *ovi.com* [5] Cloud applications in traditional and Cortex A9-based server systems, and providing feedback on the memory, processor and hardware acceleration requirements on the future 3D Server-on-Chip architectures
- 2) **Scalable 3D Server-on-Chip Architecture Specifications and Server-on-Chip Power Management:** Specifying the non-coherent and coherent 3D Server-on-Chip architectures, and hardware accelerators, implementing on-chip server power management.
- 3) **On-chip Hierarchies and Interconnects:** Optimizing the memory/cache hierarchies in the presence of on-chip DRAM, and exploring the scalable interconnection networks for performing scalable core-core and core-memory communications.
- 4) **Reliability, Availability and Fault Tolerance:** Investigating reliability/fault-tolerance/availability issues from core-level to system level including software and hardware mitigation techniques.
- 5) **Virtual and physical prototype specifications of 3D Server-on-Chip:** Prototype specification and implementation of a 3D Server-on-Chip with Cortex A9-based logic SoC and DRAM integrated in 3D packaging.

At project completion we anticipate a wide adoption of the EuroCloud concept and approach in the server design for future Data Centers.

## References

- [1] K. G. Brill, "The Invisible Crisis in the Data Center: The Economic Meltdown of Moore's Law," White Paper, Uptime Institute, 2007.
- [2] X. Fan, W. Weber, L. Barraso, "Power Provisioning for a Warehouse-sized Computer", Proceedings of the 34th Annual International Symposium on Computer Architecture (ISCA), Pages 13-23, 2007.
- [3] K. Lim, P. Ranganathan, J. Chang, C. Patel, T. Mudge, S. Reinhardt, "Understanding and Designing New Server Architectures for Emerging Warehouse-Computing Environments", 35th International Symposium on Computer Architecture (ISCA), pp.315-326, June 2008.
- [4] ARM Cortex A9, Technical Reference Manual, 2008-09, [http://infocenter.arm.com/help/topic/com.arm.doc.ddi0388e/DDI0388E\\_cortex\\_a9\\_r2p0\\_trm.pdf](http://infocenter.arm.com/help/topic/com.arm.doc.ddi0388e/DDI0388E_cortex_a9_r2p0_trm.pdf)
- [5] Nokia ovi store, <http://www.ovi.com>.