Towards Heterogeneous Multi-core
Systems-on-Chip for Edge Machine Learning

Vikram Jain • Marian Verhelst

# Towards Heterogeneous Multi-core Systems-on-Chip for Edge Machine Learning

Journey from Single-core Acceleration
to Multi-core Heterogeneous Systems

Vikram Jain
ESAT-MICAS
KU Leuven
Leuven, Belgium

Marian Verhelst
ESAT-MICAS
KU Leuven
Leuven, Belgium

# Preface

Machine Learning (ML), specifically deep learning (DL), has become the workhorse for emerging applications in vision, audio, sensing, and data analytics. State-of-the-art DL models are incredibly costly regarding model size, computational resources required, and the energy consumption of running the models. Owing to their size and complexity, they can only be deployed on large devices like GPUs typically used in cloud servers or data centers. The cloud-computing paradigm, however, comes with several drawbacks. As the large amount of data collected at the end-user devices needs transmission to the resourceful cloud servers, the energy consumption of data communication increases. This is further exacerbated by the fact that there is an exponential growth in the number of end-user devices resulting in copious amounts of data that needs to be transmitted. Moreover, the security of the collected data is not guaranteed, and recent privacy concerns have also come to the forefront. Finally, issues like latency and reliability have additionally become concerns that degrade a seamless real-time experience. Given the drawbacks of cloud computing, the logical consequence is to process data closer to the end-user devices and only sparsely transmit the data necessary for further processing or making decisions.

The paradigm of near-end-user or near-sensor processing has been coined the term (extreme-)edge-computing. (Extreme-)edge-computing can alleviate the drawbacks (energy, privacy, latency, and reliability) of cloud computing by using (extreme) edge devices for data processing as close to the sensor as possible. However, it comes with new challenges, as these devices are battery-operated and severely resource-constrained. Furthermore, as the Internet of Things becomes more pervasive, the number of sensors connected to each edge device proliferates, thereby increasing the data to be processed. Computing at the (extreme) edge requires highly energy-efficient and flexible hardware to map diverse ML and DL workloads, enabling various applications on a single platform. Moreover, it needs algorithms and models specifically designed for resource-constrained devices, thereby requiring a careful co-optimization of hardware and software. This book focuses on the first aspect of the abovementioned challenges of (extreme-)edge-computing, i.e., the design of energy-efficient and flexible hardware architectures

and hardware-software co-optimization strategies to enable early design space exploration of hardware architectures.

The book first focuses on the design of the highly specialized single hardware accelerator optimized for the application of object detection in drones. As the application and model to be accelerated are fixed, the hardware is optimized for mapping only convolutional and dense layers of a DL model in an object detection pipeline. Emerging DL applications deployed on the (extreme) edge devices, however, require multi-modal support, which demands, on the one hand, the need for much more flexible hardware accelerators and, on the other hand, complete standalone systems with the always-on and duty-cycled operation. Heterogeneity in hardware acceleration can enhance a system's flexibility and energy efficiency by utilizing various energy-efficient hardware accelerators supporting multiple DL workloads on a single platform. With this motivation, the book presents a versatile all-digital heterogeneous multi-core system-on-chip with a highly flexible ML accelerator, a RISC-V core, non-volatile memory, and a power management unit. A highly energy-efficient heterogeneous multi-core system-on-chip is presented next by combining a digital and analog in-memory computing core controlled by a single RISC-V core.

Increasing the core count further can benefit the performance of a system. However, data communication in multi-core platforms can quickly become a bottleneck if the design is not optimized. Multi-core CPUs have extensively used classical network-on-chips (NoCs) to address the data communication bottleneck. However, these NoCs use serial packet-based protocols suffering from significant protocol translation overheads toward the endpoints. In the book's final part, an open-source, fully AXI-compliant NoC fabric is proposed to better address the specific needs of multi-core DL computing platforms requiring significant burst-based communication. The NoC enables scaling DNN platforms to multi-accelerator systems, thus allowing the journey toward high-performance heterogeneous multi-core systems.

Accumulating the learning from the results obtained throughout the book, the aim is to enable flexible, high-performance, and energy-efficient (extreme) edge hardware platforms through cross-domain optimization to cater to the needs of the evolving field of machine learning. This would help win the *hardware lottery*—the problem of hardware implementation lagging behind the pace of algorithmic evolution. It will also democratize the impressive power of artificial intelligence to build a more efficient world.

Leuven, Belgium                                                                    Vikram Jain
                                                                                 Marian Verhelst

# Acknowledgments

# Contents

# List of Abbreviations

# List of Figures

# List of Tables