

ON THE AUTOMATIC INTEGRATION OF HARDWARE ACCELERATORS INTO FPGA-BASED EMBEDDED SYSTEMS

*Christian Pilato[†], Andrea Cazzaniga[†], Gianluca Durelli[†],
Andres Otero^{*}, Donatella Sciuto[†], Marco D. Santambrogio[†]*

[†]Politecnico di Milano – Dipartimento di Elettronica ed Informazione – Milano (Italy)
{pilato,acazzaniga,sciuto,santambr}@elet.polimi.it, gianluca.durelli@mail.polimi.it

^{*}Universidad Politecnica de Madrid – Centro de Electrónica Industrial – Madrid (Spain)
joseandres.otero@upm.es

ABSTRACT

This paper proposes an automatic framework for the seamless integration of hardware accelerators, starting from an OpenMP-based application and an XML file describing the HW/SW partitioning. It extends a fully software architecture by generating and integrating the cores, along with the proper interfaces, and the code for scheduling and synchronization. Experimental results show that it is possible to validate different solutions only by varying the input code.

1. INTRODUCTION

Heterogenous embedded systems [1] are becoming very popular to perform one or a few dedicated functions, often with strong requirements concerning reliability, performance, low power, real-time and cost. They allow to accelerate different parts of the applications with several general purpose, digital signal and hardware accelerators (realized using Field Programmable Gate Arrays - FPGAs - or dedicated Application Specific Integrated Circuits), interconnected through various communication mechanisms.

Nowadays, the optimization of applications for such complex embedded systems is a difficult problem. In fact, according to the traditional design flow, the hardware and software parts of an application are developed separately, often by hand, and their final integration is based on ad-hoc methods, usually requiring a long time and a high expertise by the designer. Research in hardware/software co-design aims at providing a systematic and concurrent development of the hardware and software components, providing interaction and exploration of trade-offs along the whole design process. We strongly believe that a great improvement in system-level design can be obtained when the designer will mainly focus on the development of the application and the high-level specification of the target architecture. Then, the low-level implementation (e.g., generation of cores, drivers

and synchronization directives) should be assisted by automatic design flows.

In this paper, we propose an automated design flow for the integration of hardware cores that aims at bridging the gap between high-level synthesis and system-level design. The proposed framework receives as input the application to be implemented and its decomposition in tasks, with their hardware/software partitioning. It also receives the description of the desired target platform and its components, such as the number of software processors, the interconnection topology and the information about the memory (e.g., addressing spaces). As output, it generates the hardware specification of the entire target architecture and the corresponding software code. As a proof of concept, we developed a backend for targeting Xilinx Virtex5 FPGA devices with Xilinx EDK 12.1 and generating different architectures for a representative test case. The resulting architectures can be directly synthesized, deployed and executed onto the target boards without any additional effort for the designer.

The rest of the paper is organized as follows. Section 2 discusses the related work in the area. Section 3 introduces and details the proposed automatic framework that is then validated in Section 4. Finally we draw some concluding remarks and outline future directions of work in Section 5.

2. RELATED WORK

The automatic generation of hardware platforms is a hot topic in system-level design. For example, Impulse CoDeveloper [2] proposes a design framework that can be configured via XML file, but it leaves to the designer the generation of the synchronization directives. On the other hand, in our opinion, one of the most challenging tasks to create heterogeneous platforms is the generation and the integration hardware cores, along with the development of the corresponding modified application. This work is usually performed manually and it requires a high expertise by the

designer. For example, in [3], the authors propose a framework to automatically generate FPGA-based architectures, but focusing on the prototyping of different interconnection topologies and building cost models. However, the initial application has to be enriched with the driver code to manage the execution of the different cores. Finally, if the application is decomposed in parallel tasks, the execution has to be properly synchronized. On the other hand, in recent years, different algorithms, tools and frameworks for High-Level Synthesis (HLS) have been proposed to leverage the designer from the manual generation of hardware cores (e.g., [4, 5, 6]), but without considering the complete system generation. The interested reader can find a detailed analysis of the different approaches in [7].

In our work, we adopt the structure of the cores proposed in [8] since it allows a seamless integration with software code, and we generate a wrapper for interfacing with the memory as, for example, in [5]. Then, we also generate the complete hardware and software specifications, including software directives for scheduling and synchronization. The framework is built in a modular way, such as it is possible to target different architectures or target devices simply by implementing the related *platform-dependent* part.

3. PROPOSED FRAMEWORK

The overview of the proposed design flow is shown in Figure 1. It receives in input the C code of the application, annotated with OpenMP [9] for describing the parallelism, along with XML files containing the decisions about the HW/SW partitioning and the specification of the target platform. Then, it produces the complete hardware and software specifications for the entire system, ready for the synthesis with commercial tools. The proposed framework allows the user to test applications on the top of different master-slave architectural models, where both the number of master elements (i.e. the processors) and slave elements (i.e. the hardware accelerators) can be customized. All the considered models allow to map the application tasks both in software, by means of general purpose soft processors (i.e., Xilinx Microblaze), or in hardware, by generating dedicated accelerators for specific tasks.

The design flow is mainly composed of three different parts: (1) the platform-independent analysis, (2) the platform-dependent synthesis and (3) the system synthesis. In the follows, we will describe the first two phases that are in our framework. In fact, in the last part (represented by the gray box in Figure 1), we use the generated specification to interface with third-party commercial tools for the generation of the bitstream to configure the target device. As a result, the proposed design flow results complementary to existing methods for system synthesis, allowing to raise the specification of the architecture to a higher level of abstraction.

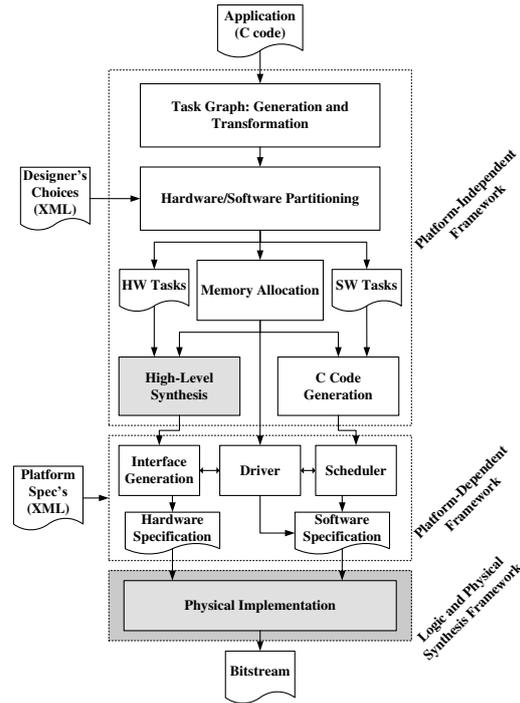


Fig. 1. Overview of the proposed design flow, where gray boxes can represent external tools.

3.1. Platform-Independent Analysis

In this phase, we analyze the input application to derive information about the task graph and the tasks to be executed on the different processing elements. Then we also identify the variables to be allocated in the shared memory and we generate both the hardware and software cores based on this information.

Task Graph Generation: As a first step, we extract the Hierarchical Task Graph [10] (HTG) (i.e., the representation we adopted for the partitioned application) for the original source code, annotated with OpenMP pragmas, and we represent each task with a function. Then, additional transformations may be performed on the HTG. For example, a task needs to be split up if it contains a sequence of function calls that, accordingly to the mapping, have to be run on different processing elements. At this point, we integrate the decisions about the hardware/software partitioning and the tasks are classified with respect to the processing elements assigned for their execution. The hardware tasks require the generation of the corresponding modules, while the code for the software ones can be directly generated. However, in both the cases, it is necessary to define the allocation of shared variables in advance.

Allocation: After determining the tasks to be executed and integrating the information about hardware/software partitioning, the subsequent step consists in the identification of the variables that require an allocation in the shared mem-

Table 1. Area occupation and execution time of the kernel computation for the different experiments. S_p represents the speed-up with respect to the sequential software execution.

	# Cores		Area		Time	
	SW	HW	(# slices)	(%)	(# cycles)	S_p
1	1	0	4,102	23.74	6,593,760	1.00
	0	1	5,088	29.44	2,663,764	2.48
2	1	1	4,452	25.76	3,348,185	1.97
	0	2	5,709	33.04	1,362,660	4.84
4	1	3	6,977	40.38	1,777,853	3.71
	0	4	7,908	45.76	1,244,516	5.30
6	1	5	8,408	48.66	1,396,466	4.72
	0	6	9,432	54.58	1,331,082	4.95

this application by splitting up the computations in different cores, from 1 up to 6, differently assigned to the software processors or the dedicated hardware cores, with an ASAP dynamic scheduling algorithm.

The kernel execution time for the different experiments is shown in Table 1, along with their area occupation. We also reported the speed-up S_p with respect to the sequential software execution. As expected, the best execution times are obtained when the tasks are totally executed in hardware, reaching a speed up of 5.3x with four cores. In fact, when one of the tasks is executed in software, this usually limits the performance of the system if the the core iterations are not properly balanced. On the other hand, as shown in Figure 3, increasing the number of the cores results in a contention when accessing the shared memory. This can lead to better explore the memory allocation or to rewrite the application to reduce the memory accesses.

Note that these experiments have been obtained just by varying the input source code and the XML file describing the HW/SW partitioning. Then, each of these architectures required few seconds to be generated by the proposed design flow, ready for the synthesis.

5. CONCLUSIONS AND FUTURE WORK

This paper presented an automated framework for integrating hardware accelerators in FPGA-based embedded systems directly from the C code of the application. It only requires information about the hardware/software partitioning and it can automatically interface commercial tools for the generation of the configuration bitstream.

The proposed framework will serve to evaluate the different techniques that we are developing in system-level design. In particular, the integration of the support for dynamic partial reconfiguration will allow to deploy architectures with a larger part of the application to be executed in hardware.

Acknowledgments

This work was partially funded by the European Commission in the context of the FP7 FASTER project (#287804).

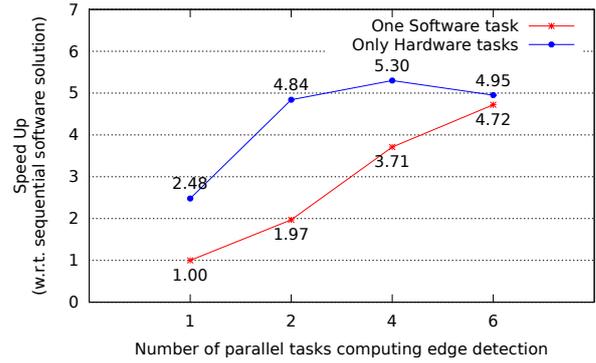


Fig. 3. Speed-up resulting from the hardware acceleration.

6. REFERENCES

- [1] M. Duranton, "The Challenges for High Performance Embedded Systems," in *Proceedings of DSD*, 2006, pp. 3–7.
- [2] R. Bodenner, "Creating Platform Support Packages," *Impulse Accelerated technologies, Inc., Application Note: IATAPP-109*, 2011.
- [3] S. Secchi, P. Meloni, and L. Raffo, "Exploiting FPGAs for technology-aware system-level evaluation of multi-core architectures," in *Proceedings of ISPASS 2010*, march 2010, pp. 194–202.
- [4] P. Coussy, C. Chavet, P. Bomel, D. Heller, E. Senn, and E. Martin, "GAUT: A high-level synthesis tool for DSP applications," in *High-Level Synthesis: From Algorithm to Digital Circuit*, P. Coussy and A. Morawiec, Eds. Springer, 2008.
- [5] J. Cong, Y. Fan, G. Han, W. Jiang, and Z. Zhang, "Platform-Based Behavior-Level and System-Level Synthesis," in *IEEE International SOC Conference*, 2006, pp. 199–202.
- [6] Politecnico di Milano, "Bambu web site." [Online]. Available: <http://panda.dei.polimi.it>
- [7] J. Cong, B. Liu, S. Neuendorffer, J. Noguera, K. Vissers, and Z. Zhang, "High-Level Synthesis for FPGAs: From Prototyping to Deployment," *IEEE Trans. on CAD of Integrated Circuits and Systems*, vol. 30, no. 4, pp. 473–491, April 2011.
- [8] C. Pilato, F. Ferrandi, and D. Sciuto, "A design methodology to implement memory accesses in high-level synthesis," in *Proceedings of CODES+ISSS '11*, 2011, pp. 49–58.
- [9] M. Sato, "OpenMP: parallel programming API for shared memory multiprocessors and on-chip multiprocessors," in *ISSS '02*, 2002, pp. 109–111.
- [10] M. Girkar and C. D. Polychronopoulos, "The hierarchical task graph as a universal intermediate representation," *International Journal of Parallel Program*, vol. 22, no. 5, pp. 519–551, 1994.
- [11] G. Durelli, A. Cazzaniga, C. Pilato, D. Sciuto, and M. D. Santambrogio, "Automatic run-time manager generation for reconfigurable mpsoC architectures," *Proceedings of ReCoSoC '12*, 2012.
- [12] GNU Compiler Collection, "GCC, version 4.5." [Online]. Available: <http://gcc.gnu.org/>