



Throughput Optimization for Streaming Applications on CPU-FPGA Heterogeneous Systems

Xuechao Wei, Yun Liang, Tao Wang, Songwu Lu and Jason Cong

Center for Energy-efficient and Applications (CECA)

School of EECS, Peking University, China



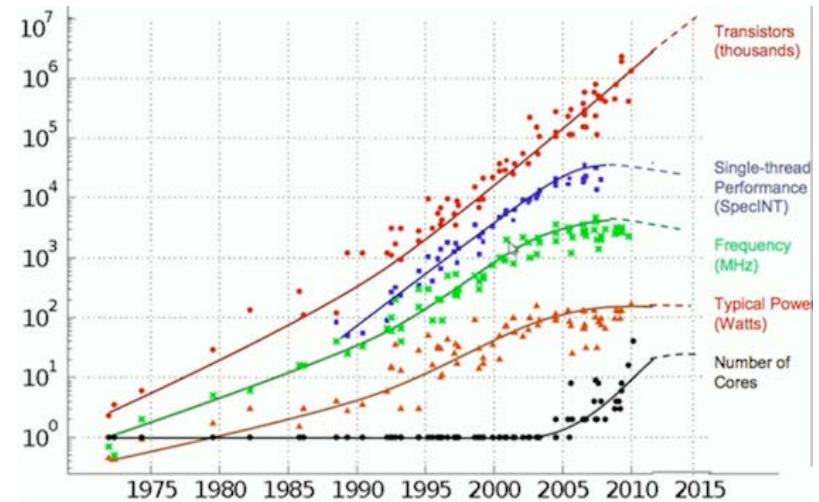
Why We Need Heterogeneous Systems

◆ End of Denard scaling

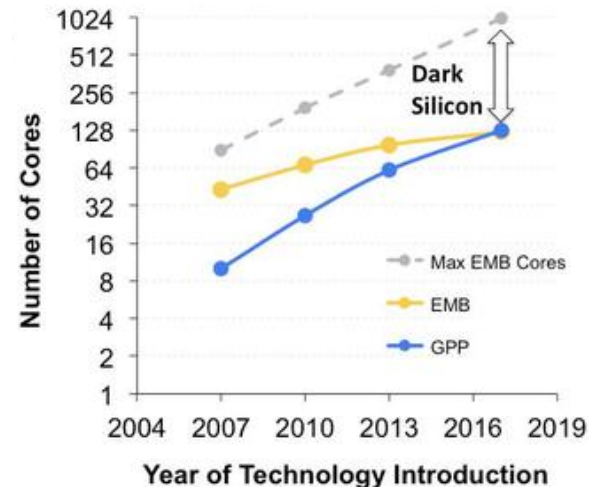
- Power constraint
- Dark silicon

◆ New application fields

- Big volume, time critical, multiple types (e.g., text, images, videos)



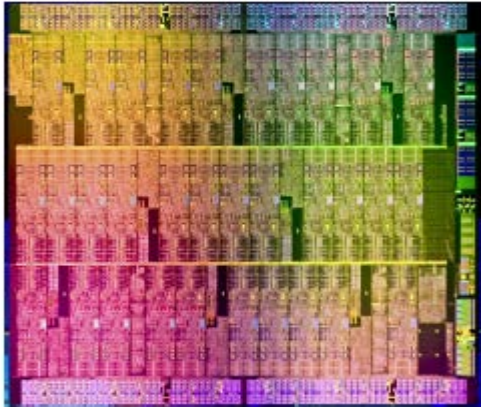
Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten.



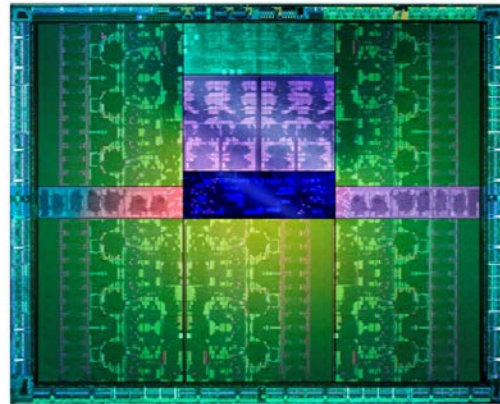
Hardavellas, et al "Towards Dark Silicon in Servers", IEEE Micro, 2011.



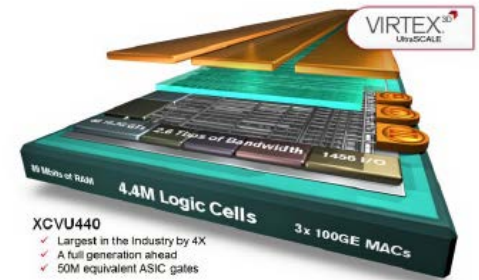
How Promising of Heterogeneous Computing



Intel Xeon Phi
~5B transistors



NVIDIA GK110 Kepler
~7.1B transistors



Xilinx Virtex-UltraScale
XCVU440
~20B transistors

Platform	Normalized Speedup	Normalized Performance/Watt
FPGA	545:1	1090:1
GPU	50:1	21:1
GPP	1:1	1:1

Source: "Reconfigurable Computing in the Multi-Core Era," Khaled Benkrid, HEART'2010



Background of Streaming Applications

- ◆ Streaming applications are centered around streams, and process input data frames iteratively
- ◆ Widespread parallelism and regular communication patterns
- ◆ High performance digital signal processing, time critical
- ◆ Widely applied on embedded platforms that have rigid power supply



CPU is Idle

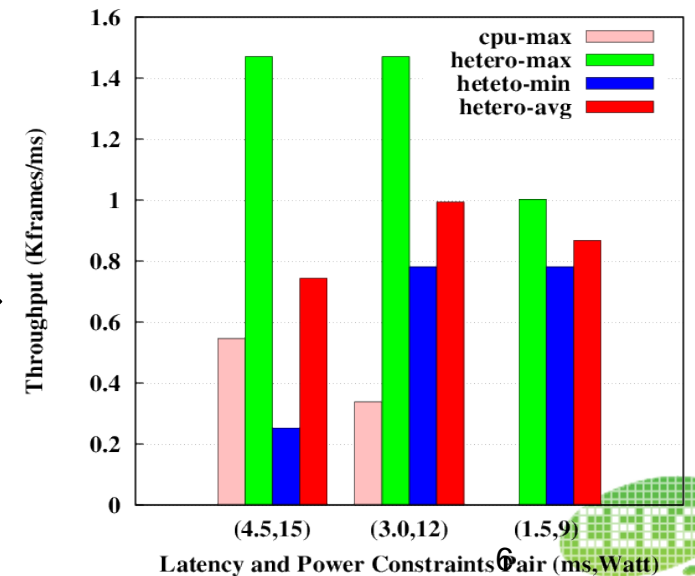
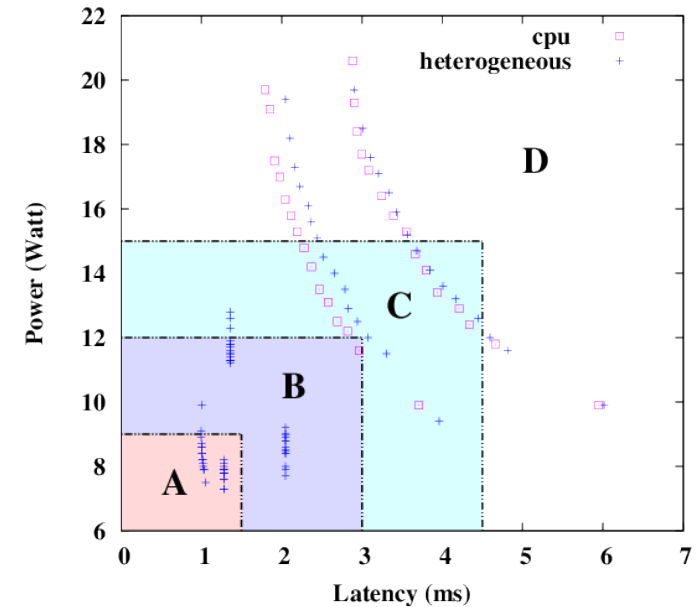
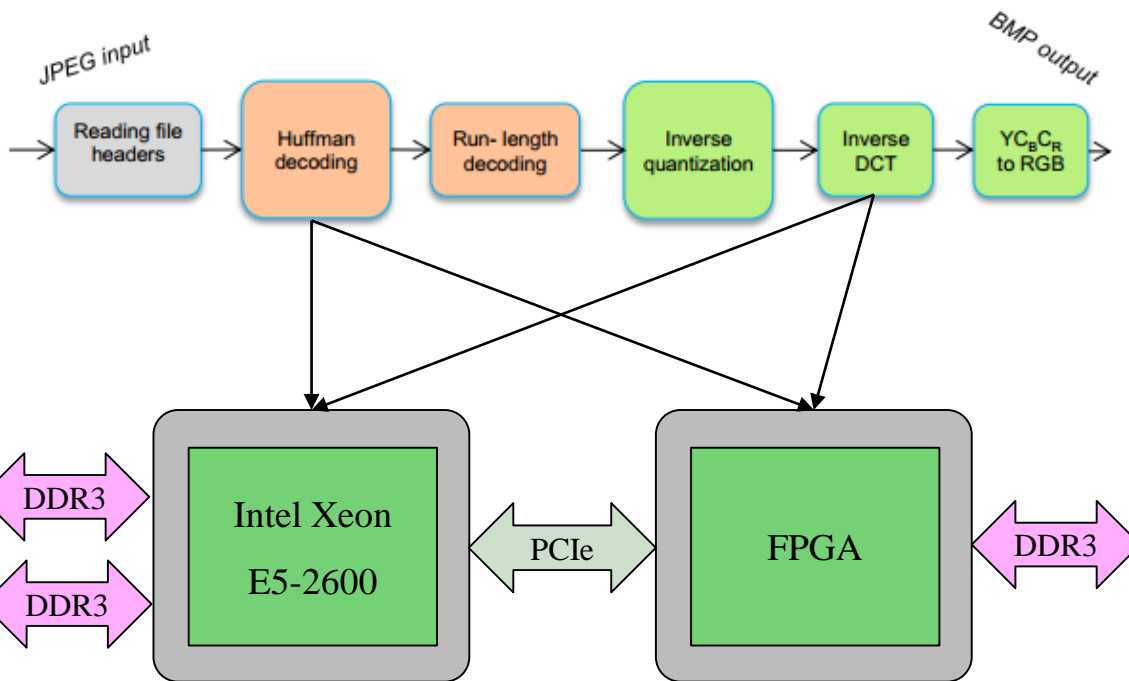
- ◆ Mapping streaming graphs onto accelerators
 - DATE'12, FPGA'14, CASES'08, DAC'09, PLDI'08, PPOPP'12, CGO'16
- ◆ Power or energy aware execution
 - DAC'13, FPL'15



Heterogeneity in Streaming Applications

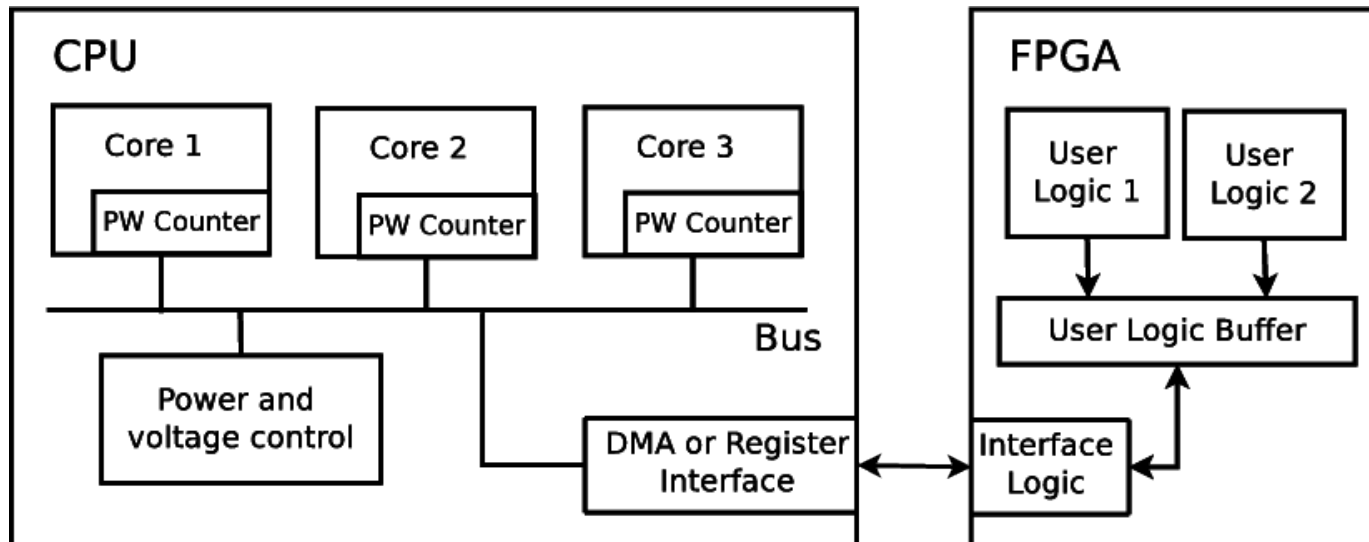
◆ Streaming applications

- Task graph
- Heterogeneity (Perf./Power)



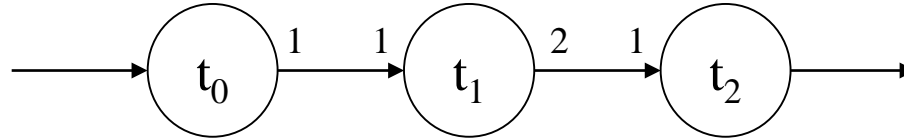
System Architecture

- ◆ FPGA coupled with multicore CPU
 - Voltage and frequency can be changed on CPU core
 - HLS tool converts high level descriptions of tasks into HDL
- ◆ High speed transfer integration

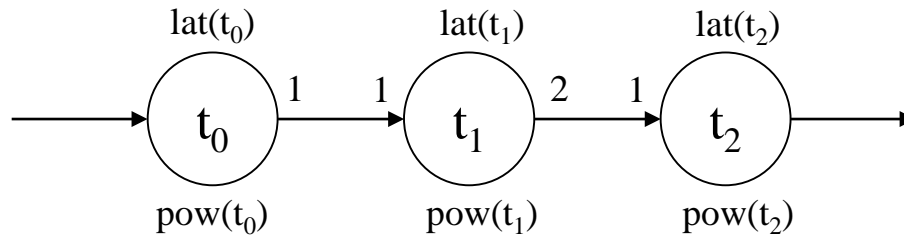


Application Model

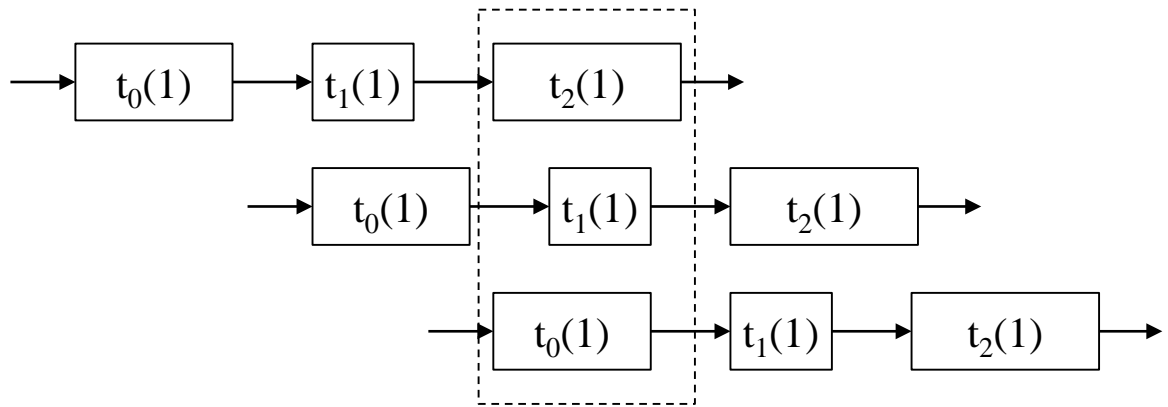
Synchronous Data Flow
Graph (SDFG)



SDFG labelled with latency
and power attributes



Overlapped consecutive frames
in different pipeline stages

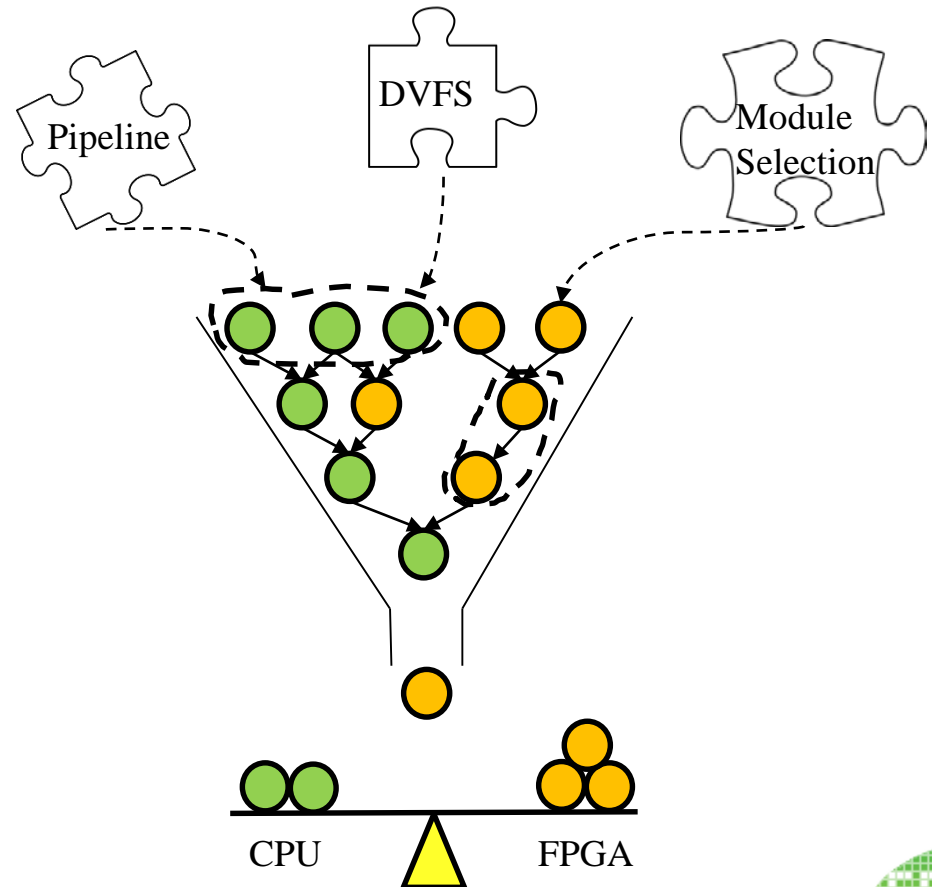


Problem Definition

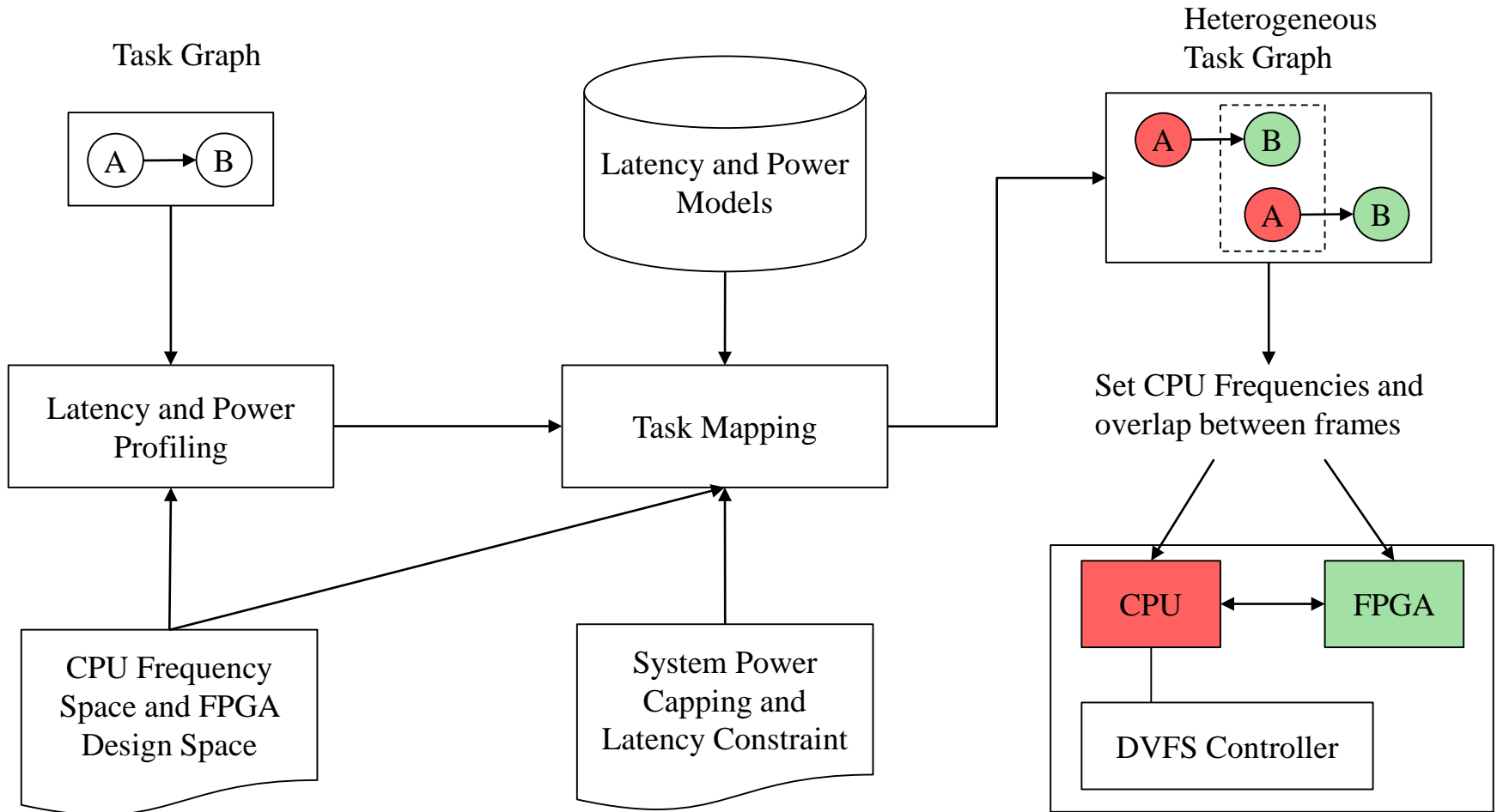
◆ Mapping tasks of streaming applications on CPU-FPGA system to optimize throughput under

- System power capping
- Application latency

Task	CPU	FGPA
task0	(200ms, 30w)	(400ms, 8w)
task1	(100ms, 40w)	(60ms, 15w)
.....		



Optimization Flow



Latency and Power Models

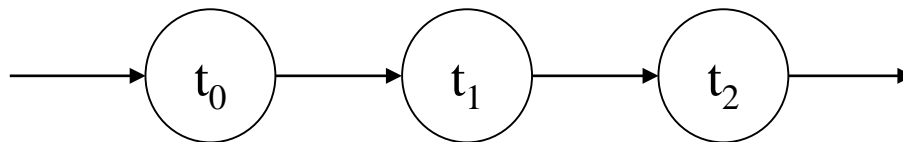
◆ Latency model

$$Lat = m \times \max_{1 \leq s_i \leq s_m} lat_{s_i} \quad Throughput = \frac{1}{\max_{s_1 \leq s_i \leq s_m} lat_{s_i}}$$

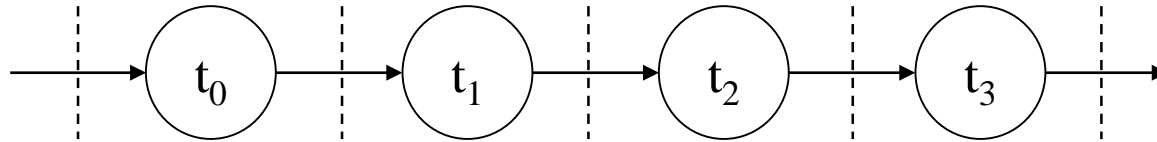
$$lat_{s_i} = \sum_{start(s_i) \leq k \leq end(s_i)} lat(t_k) + e(t_k, t_{k+1})$$

◆ Power model

$$Pow = \sum_{s_1 \leq s_i \leq s_m} P_{s_i}^{active} + P_{cpu}^{idle} + P_{fpga}^{idle} + P_{fpga}^{comm}$$



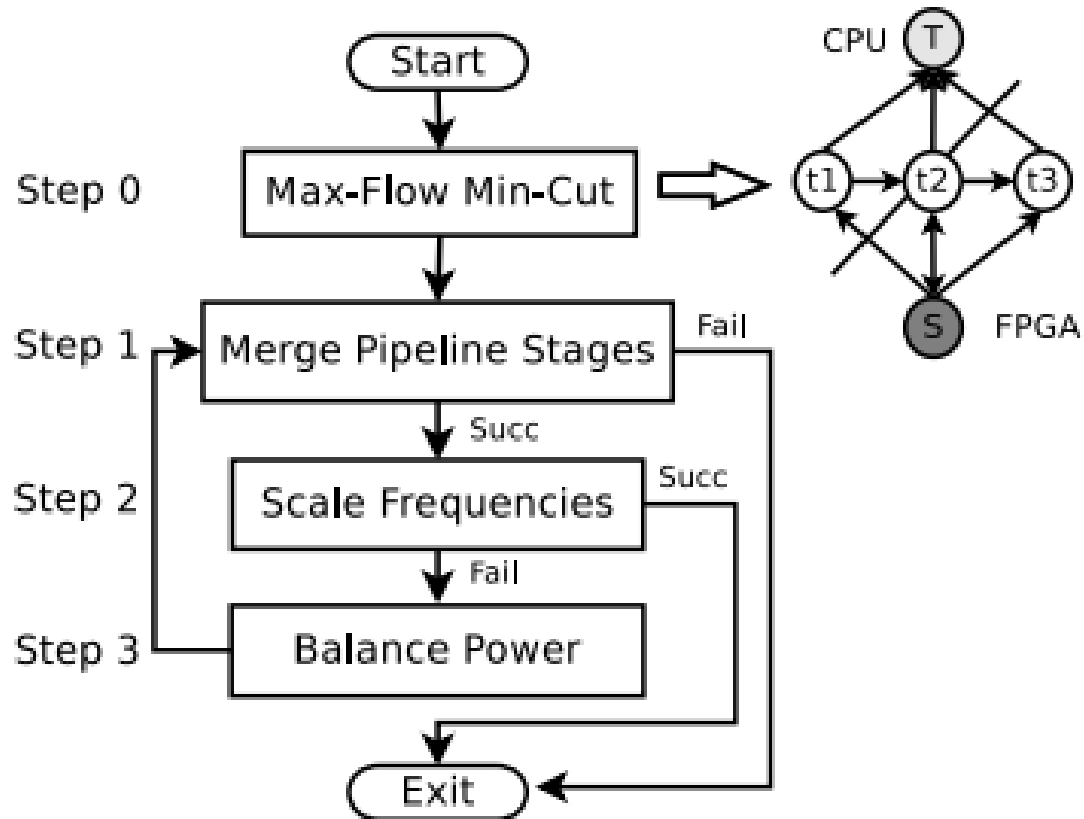
Optimal Algorithm



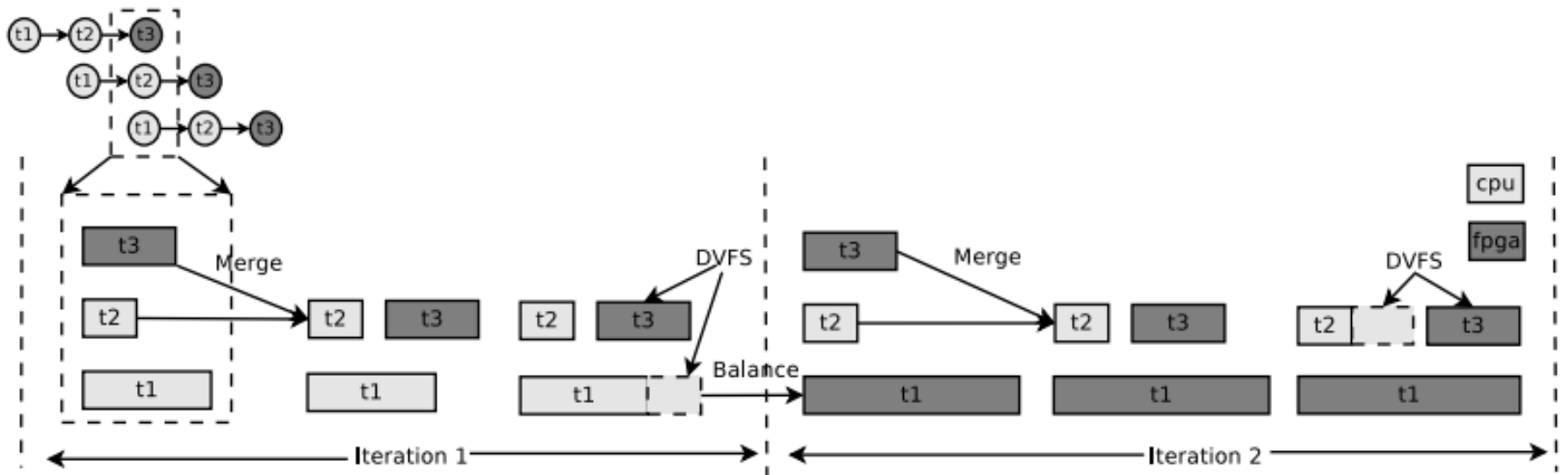
1. How to assign tasks into pipeline stages? (2^n)
 2. Map a task on CPU or FPGA?
 3. Which frequency or implementation?
- } $(f+m)^n$



Heuristic Algorithm



Heuristic Algorithm



Experiment Setup

◆ CPU-FPGA heterogeneous system

- CPU – Intel Core i5
- FPGA – Xilinx VC707

◆ Benchmark suite

- AES, MD5 (security), DCT, JPEG (image processing), LTE receiver (wireless communication)

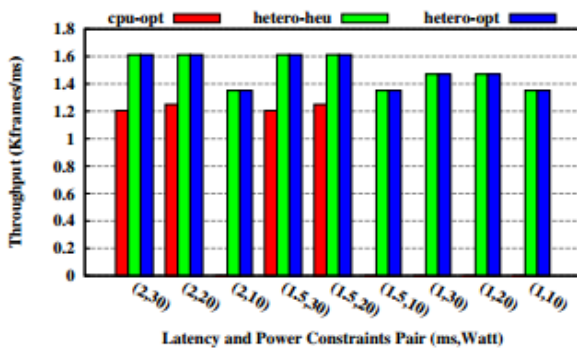
◆ Implementation

- Hardware implementation – Xilinx High Level Synthesis (HLS) tools
- Power – Intel's Running Average Power Limit (RAPL) interface and Xilinx power estimator
- SDF tool – SDF³
- Pipeline – POSIX Threads interfaces

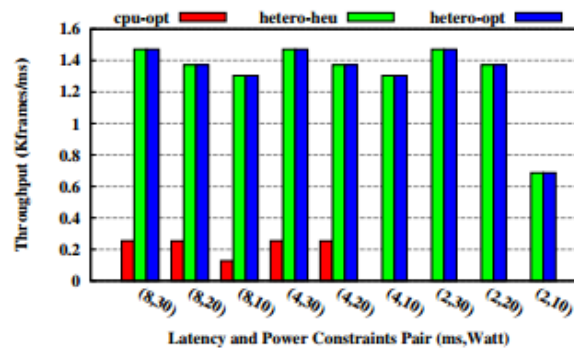


Results

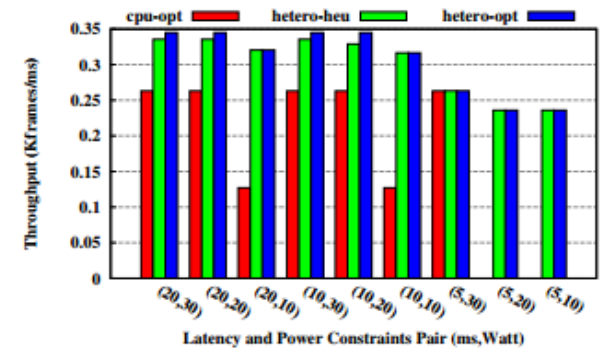
◆ Performance improvement



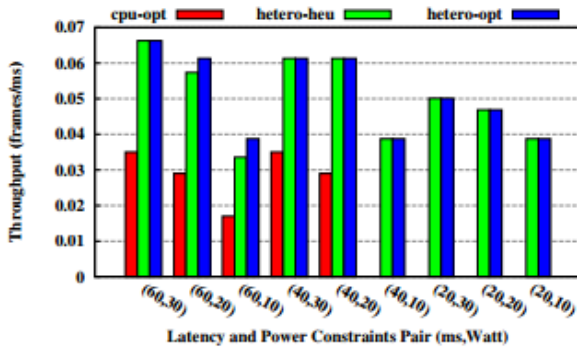
(a) AES



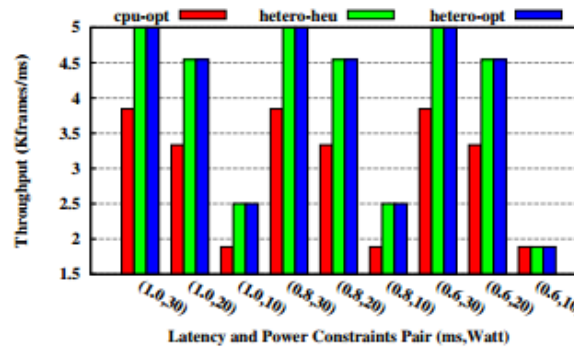
(b) DCT



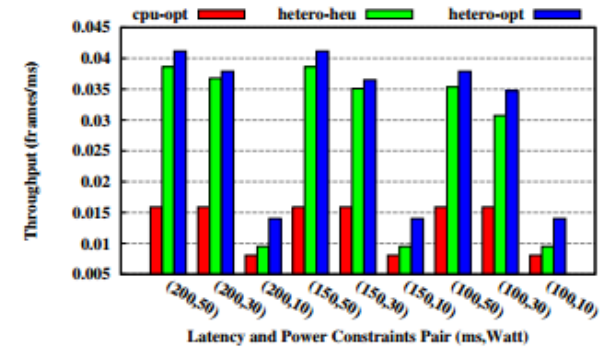
(c) JPEG



(d) LTE-Rx3



(e) MD5



(f) LTE-Rx8



Results

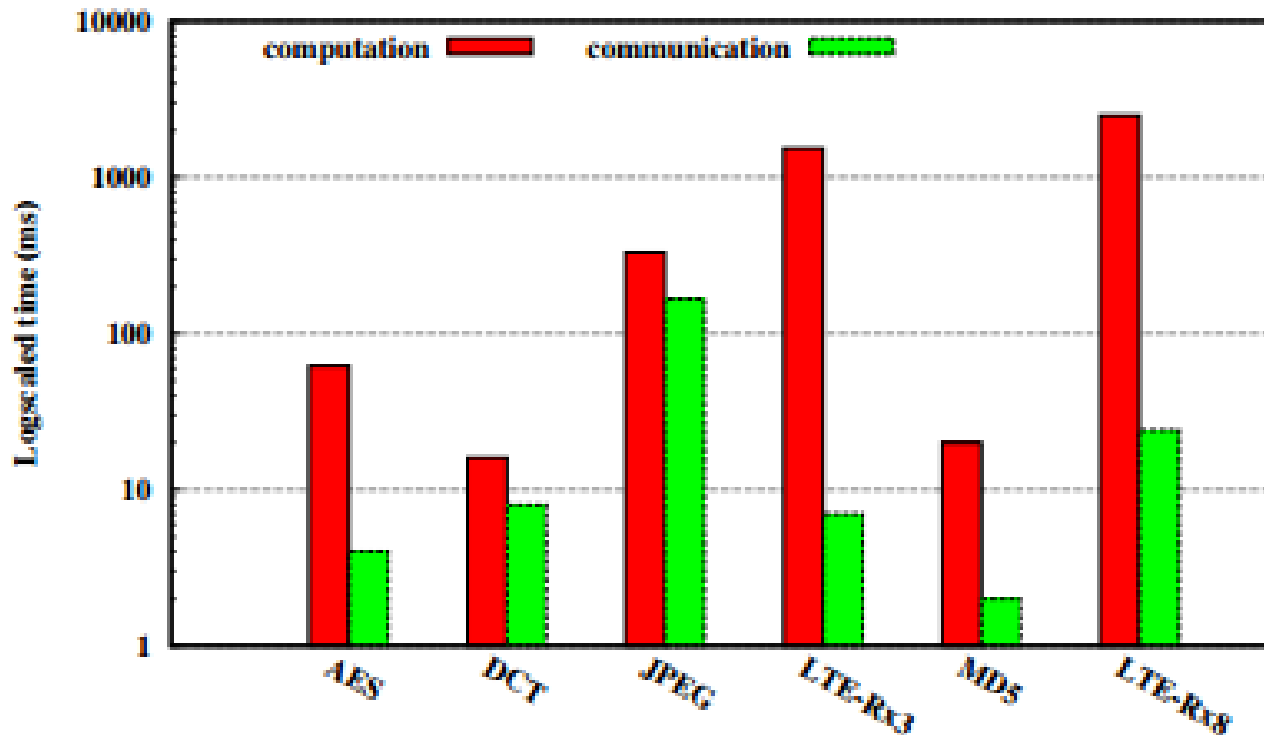
- ◆ Runtime comparison between two algorithms

Benchmark	Hetero-Heuristic	Hetero-Optimal	Speedup
AES	0.29ms	1.20ms	4.14
DCT	1.67ms	40.13ms	24.47
JPEG	6.49ms	105.60s	16271
LTE-Rx3	0.43ms	12.53ms	29.14
LTE-Rx8	20.18ms	>4 hours	N/A
MD5	0.24ms	2.07ms	8.63



Results

◆ Hiding of communication overhead



Thank you!

Q&A

Xuechao Wei

xuechao.wei@pku.edu.cn

