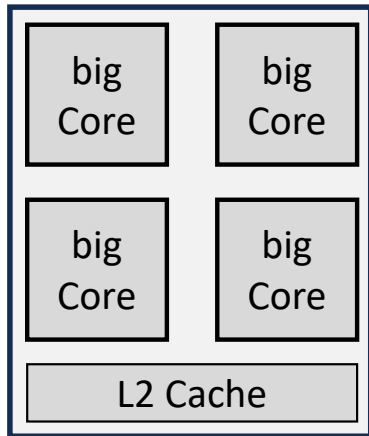
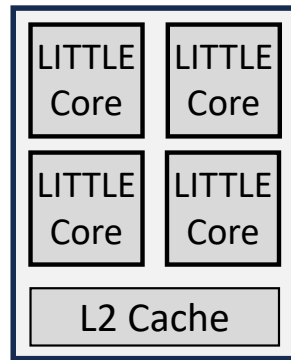


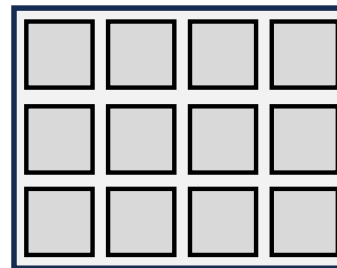
**“big” cores cluster
(Cortex-A15, 2.8GHz)**



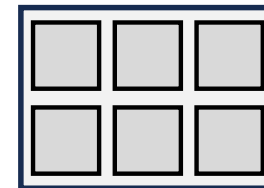
**“LITTLE” cores cluster
(Cortex-A7, 1.8GHz)**



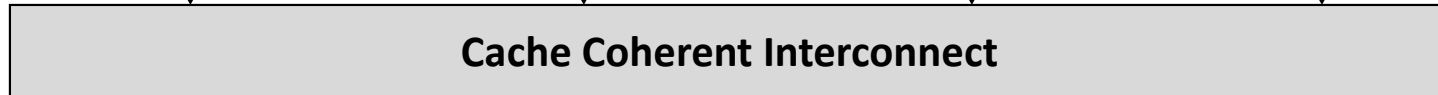
GPU

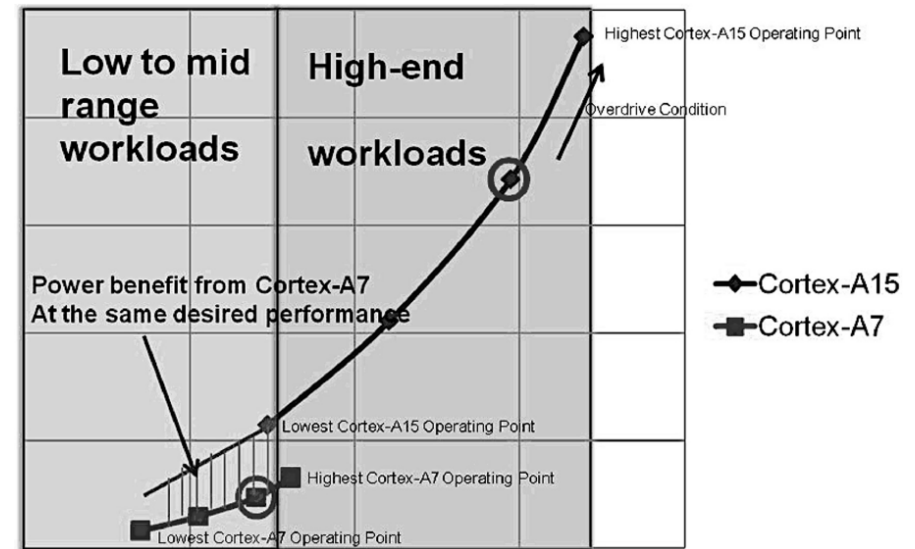
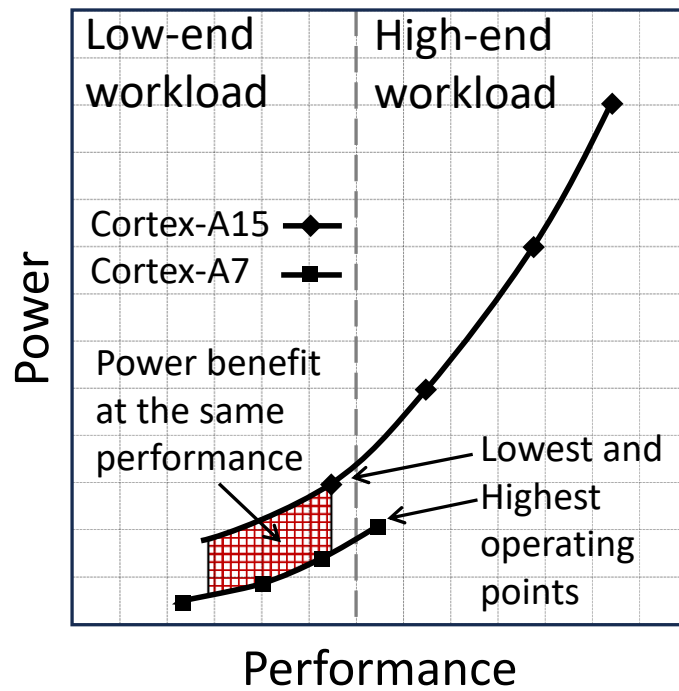


NPU

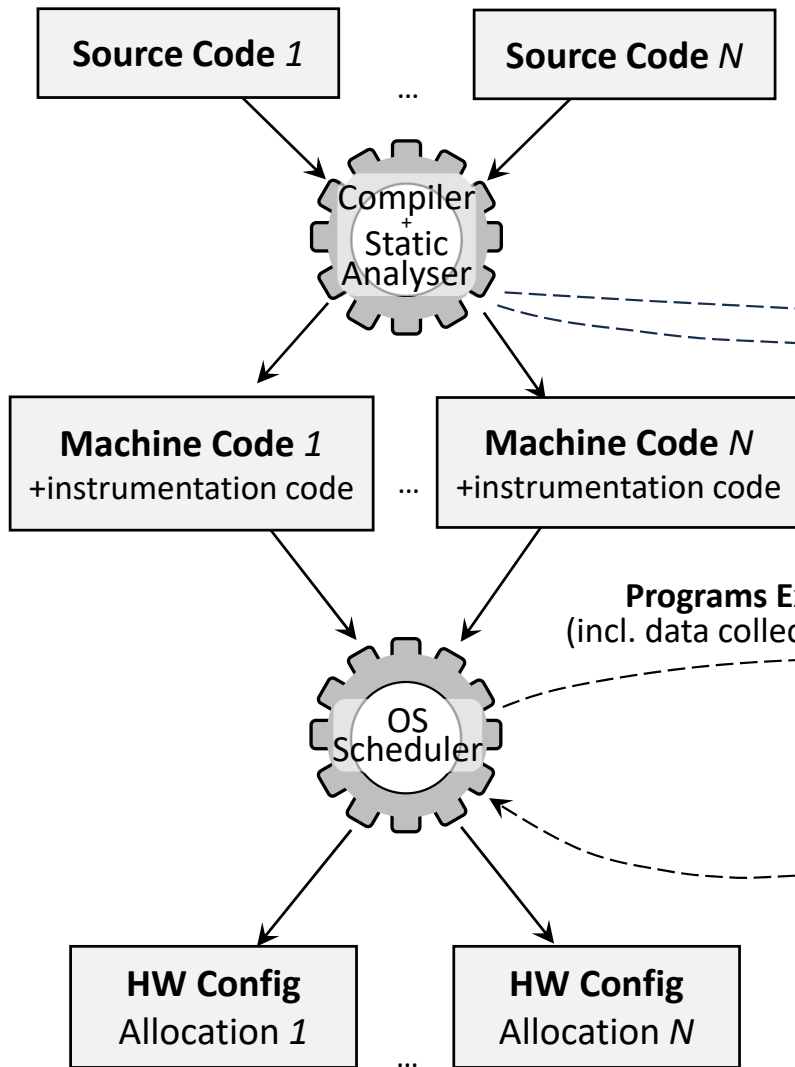


Cache Coherent Interconnect

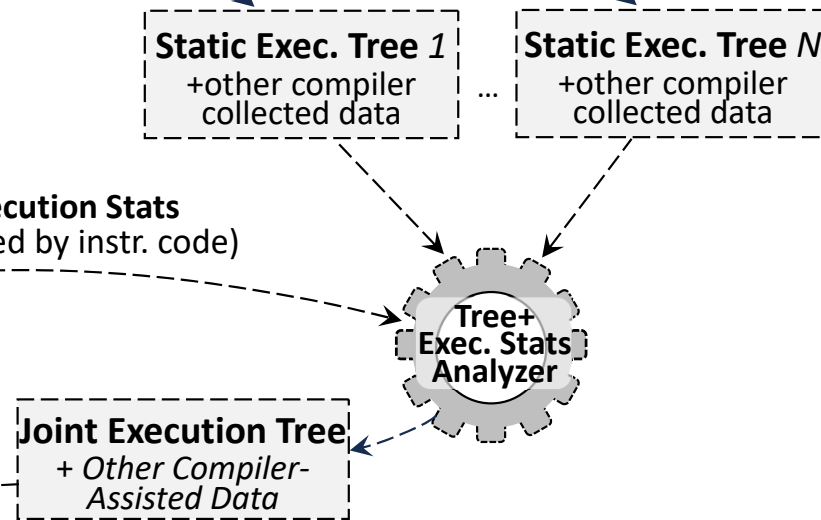


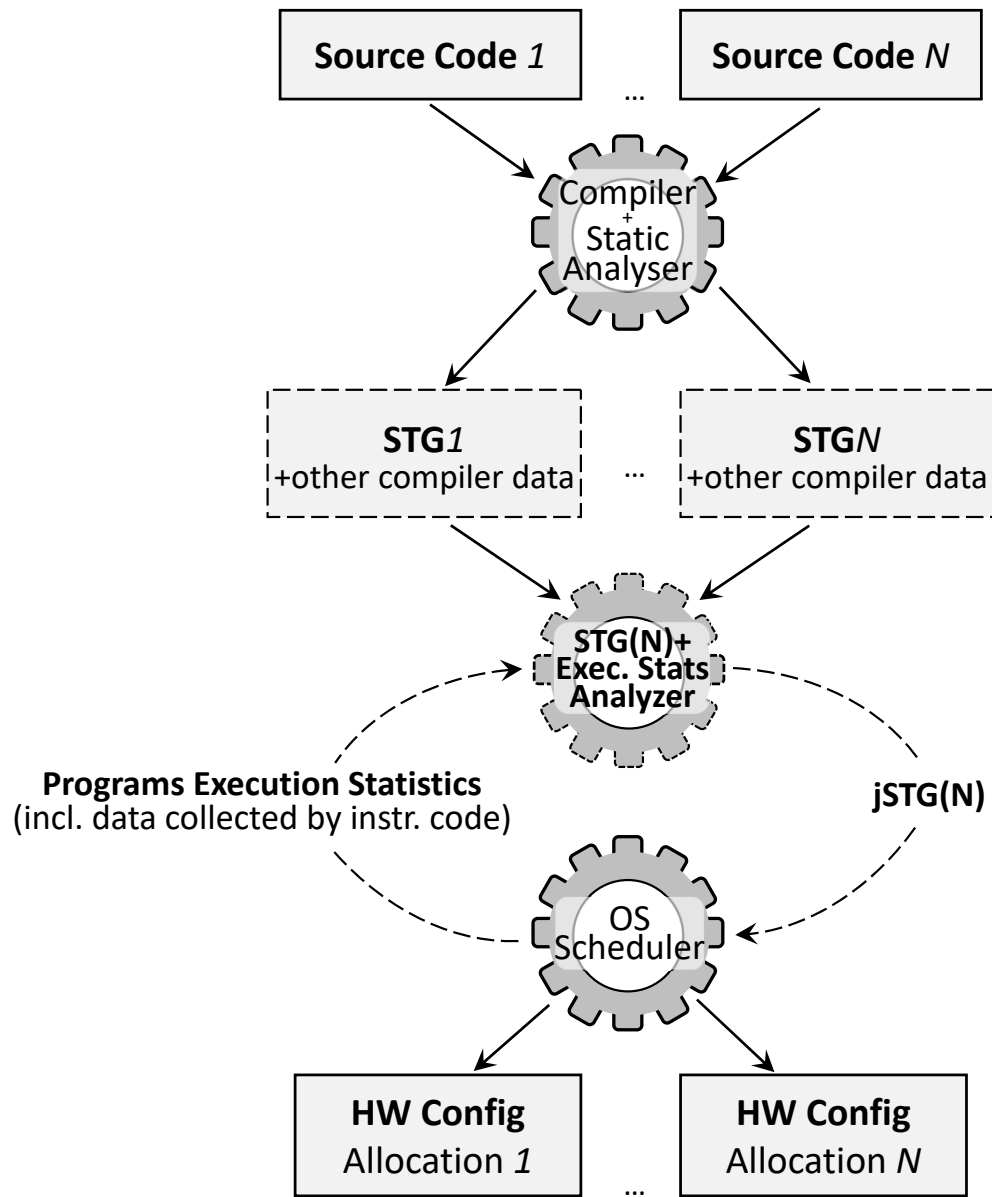


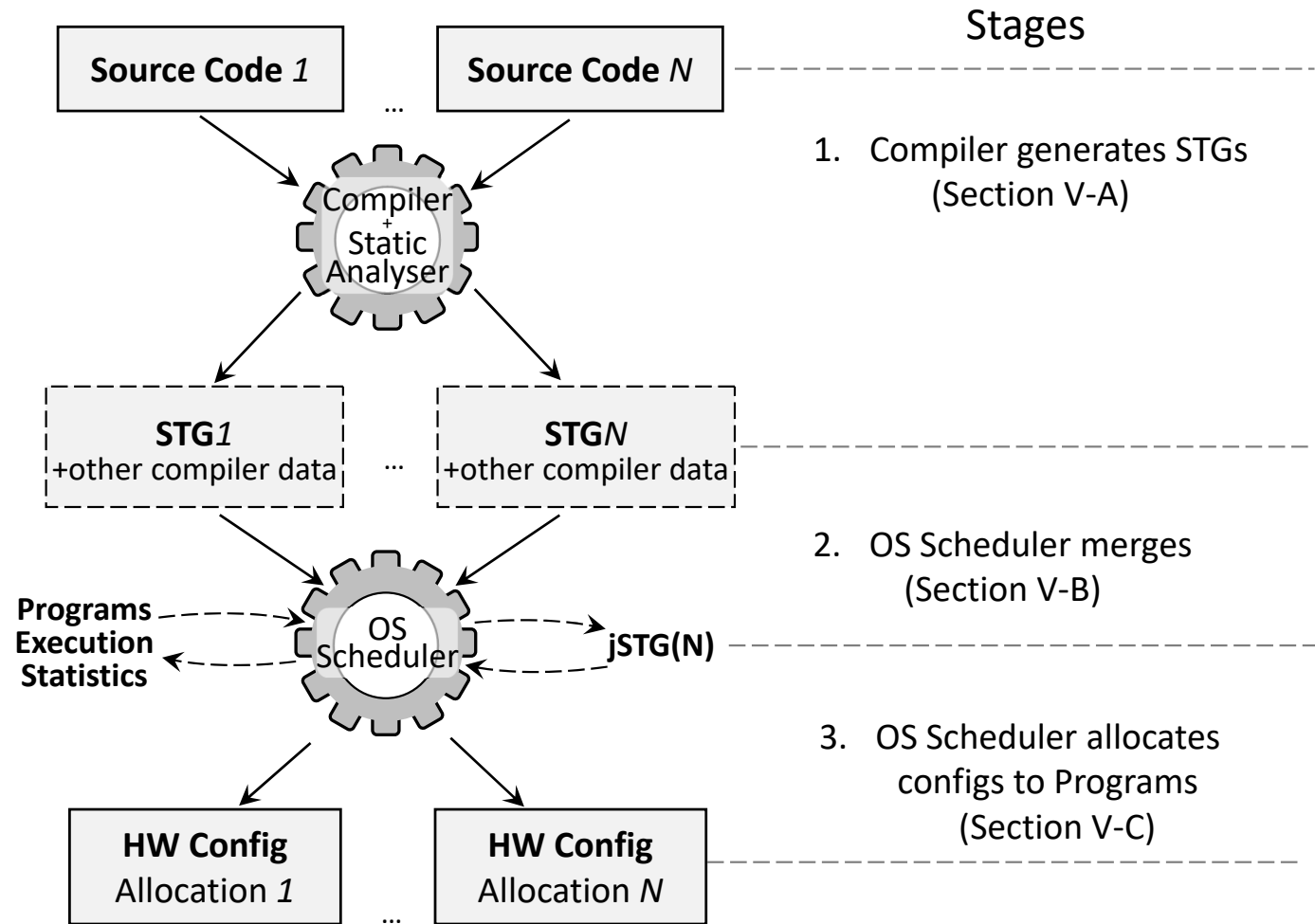
Traditional Scheduling

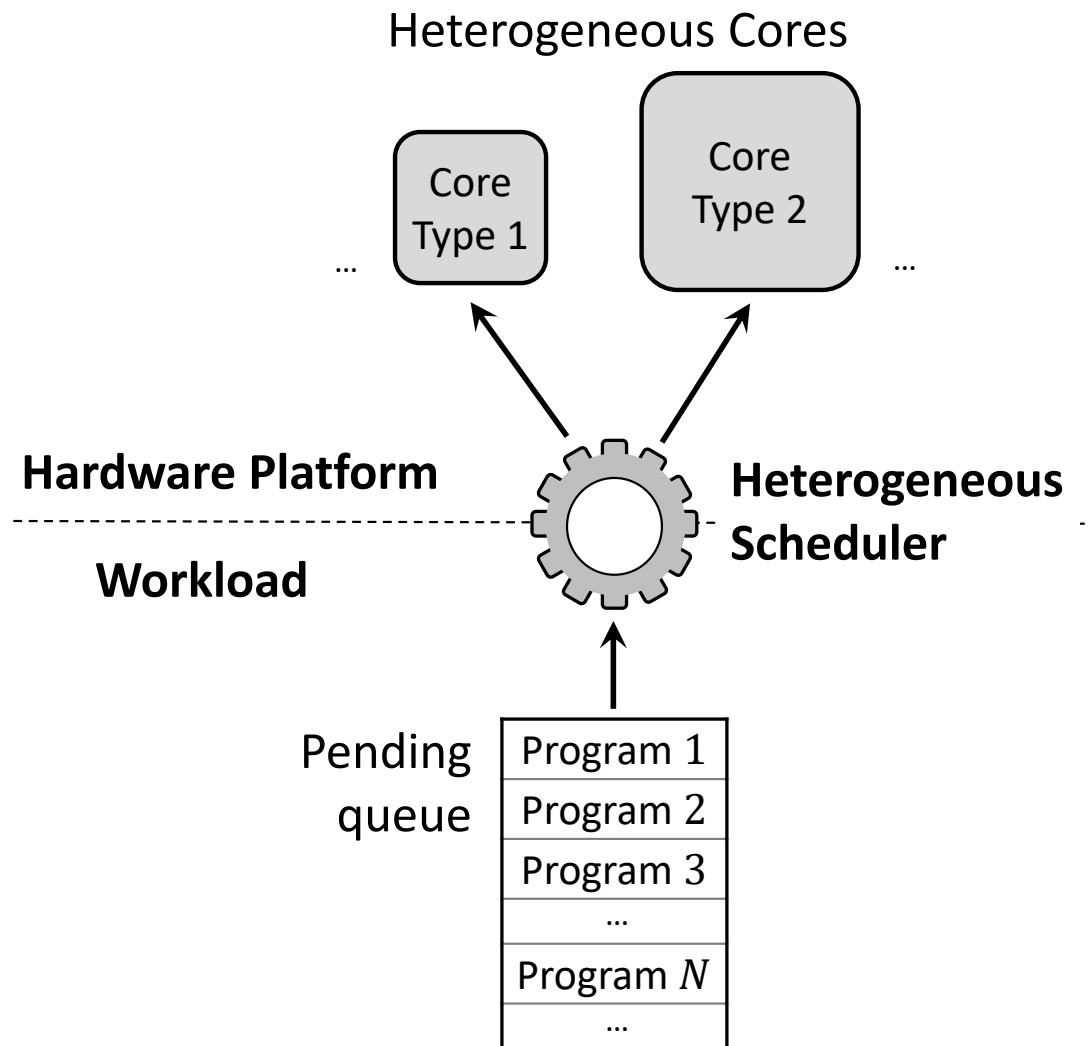


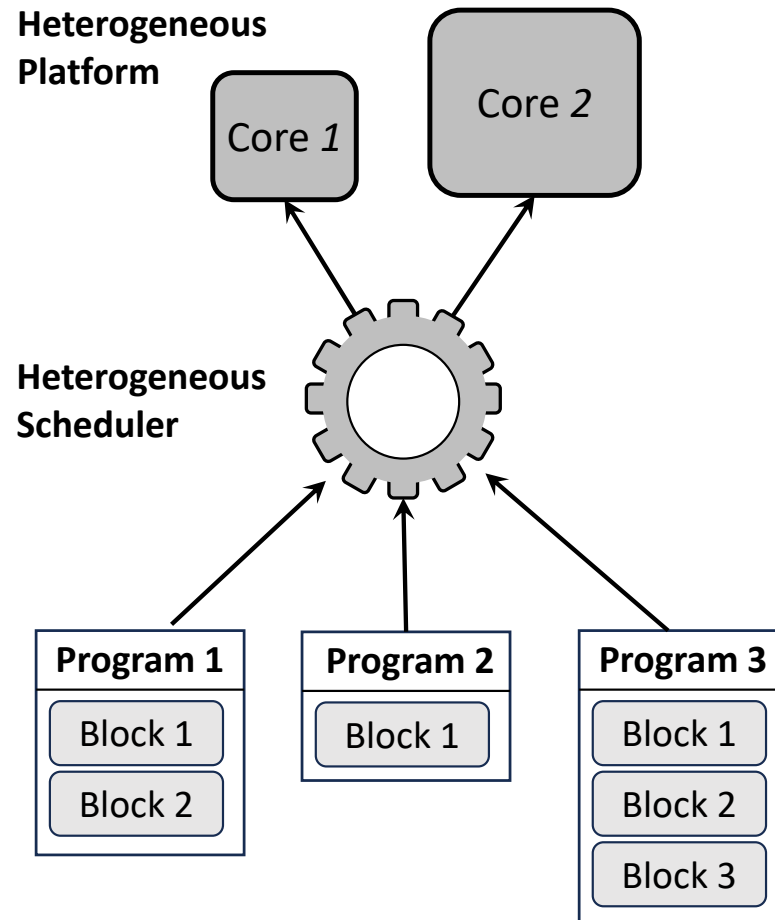
Extension to Proposed Compiler-Assisted Scheduling









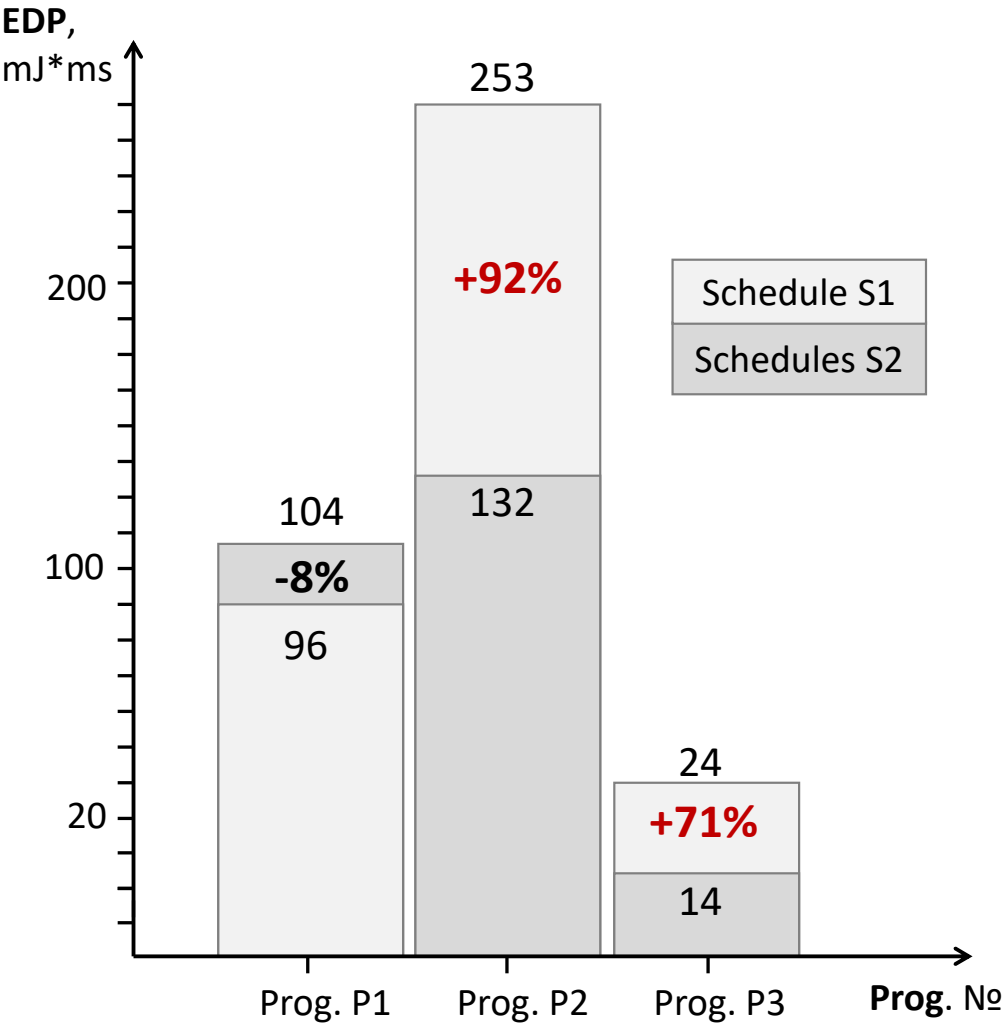


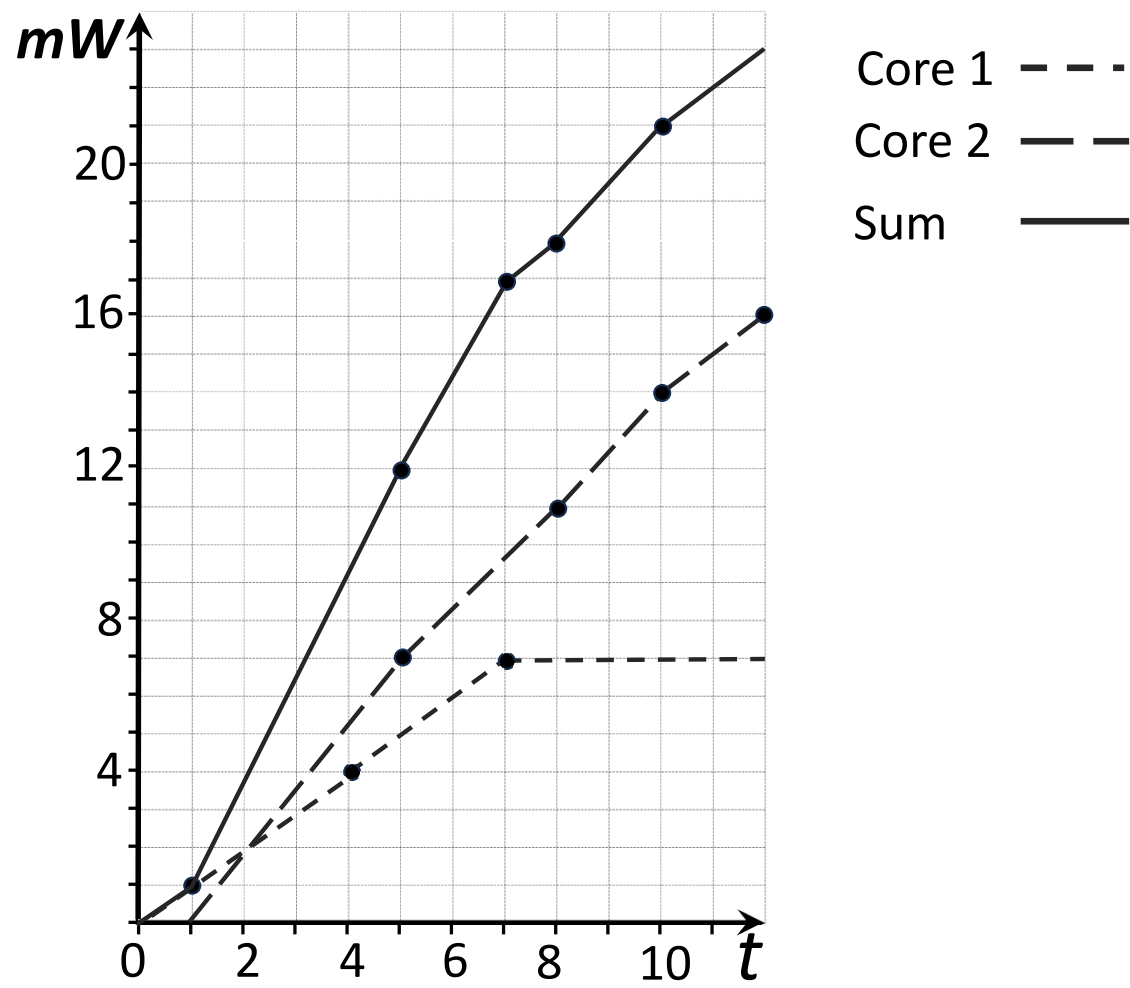
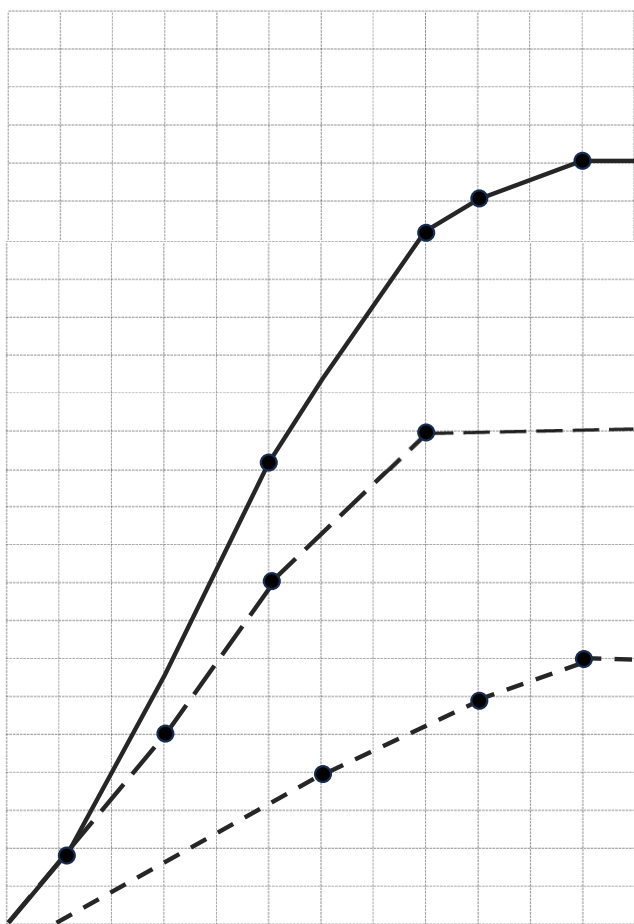
High-level scheduling of 3 programs of several blocks to be executed over a heterogeneous platform of 2 cores.

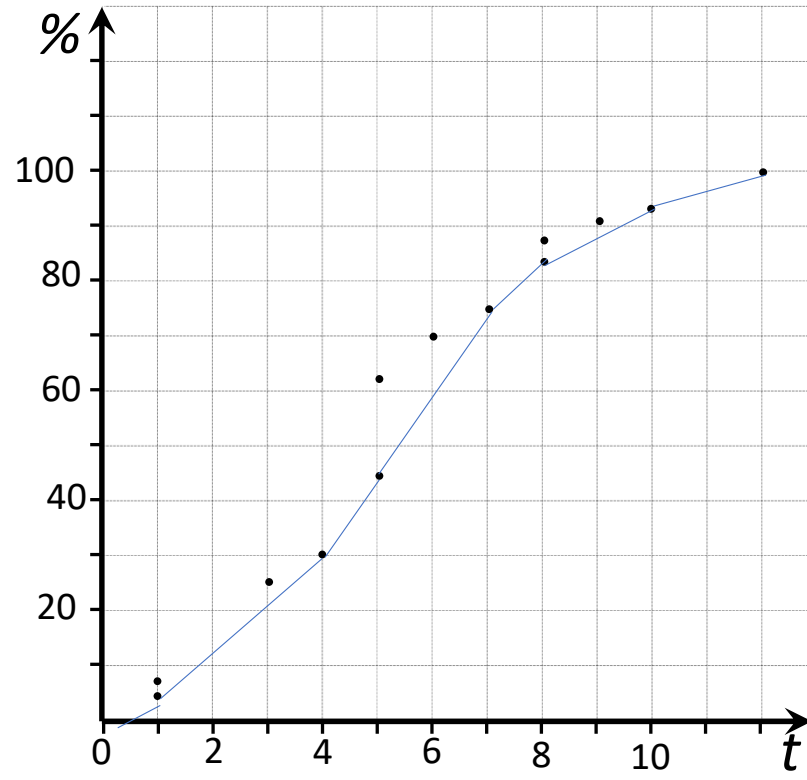
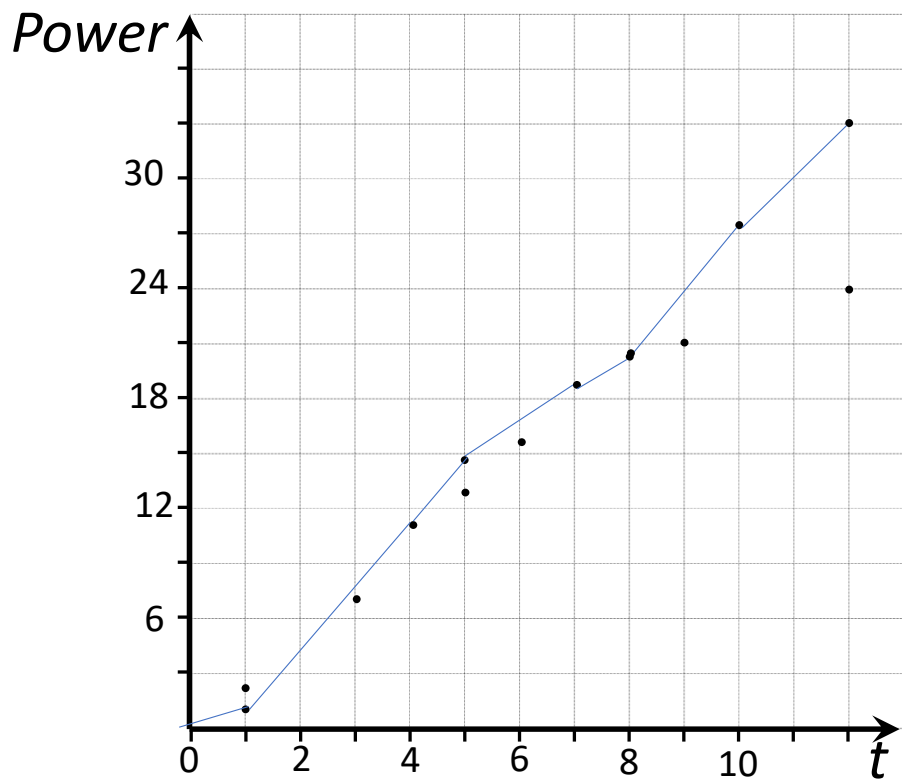
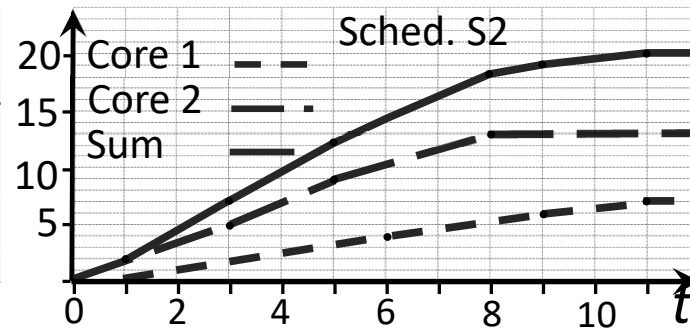
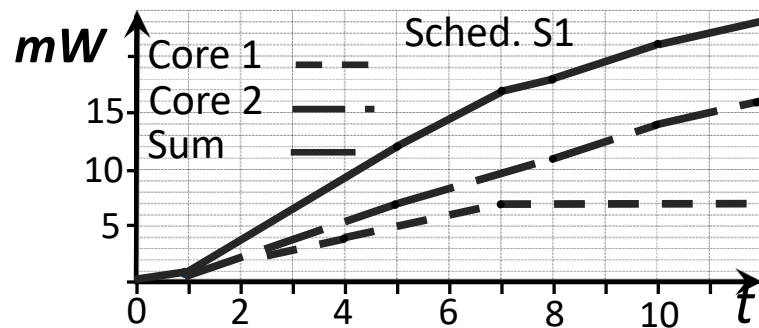
Sample execution requirements of program blocks of 3 programs for 2 heterogeneous cores.

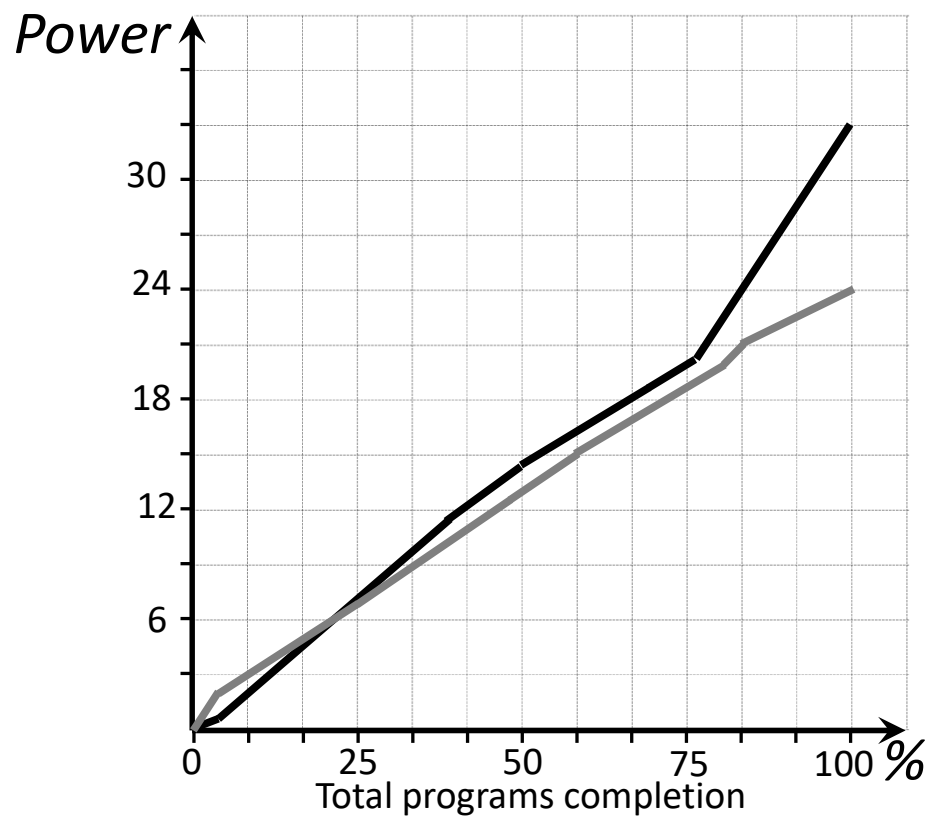
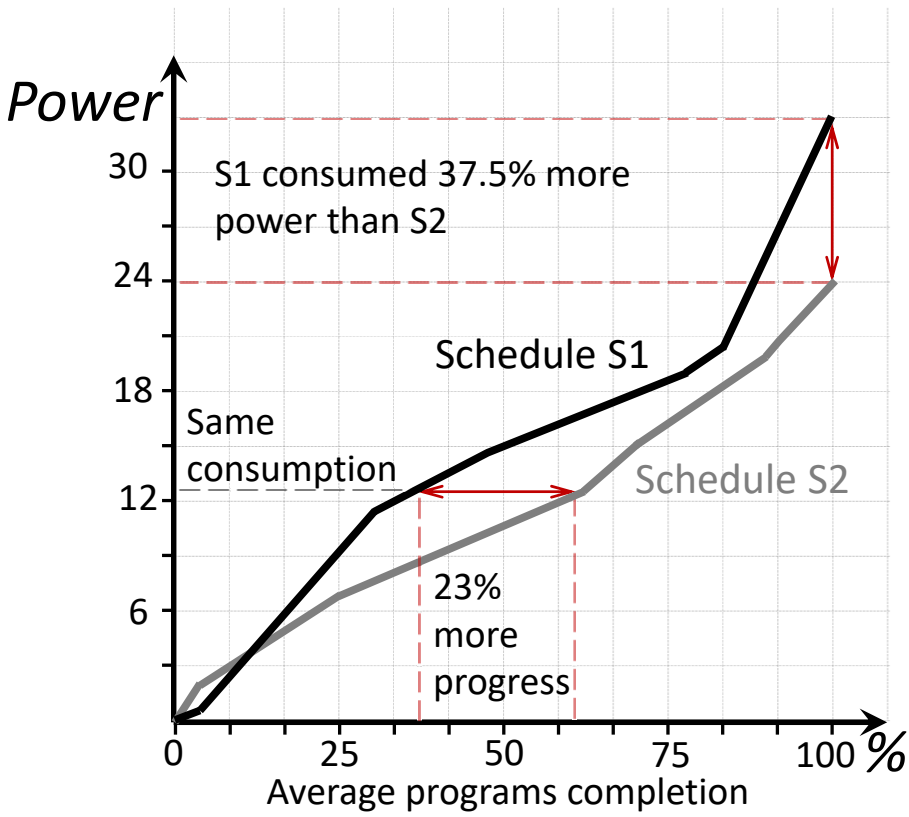
Programs /Blocks		Execution time, ms		Energy use, mJ	
		Core 1	Core 2	Core 1	Core 2
		(slow)	(fast)		
P1	B1	4	3	8	9
	B2	5	3	3	4
P2	B1	5	4	6	10
	B2	3	2	3	7
	B3	3	2	3	6
P3	B1	3	2	6	7

Programs	Execution time, ms		Energy use, mJ	
	Core (no DVFS)	Core (slower)	Core (no DVFS)	Core (slower)
P1	5	7	9	8
P2	4	5	10	6
P3	2	3	7	6

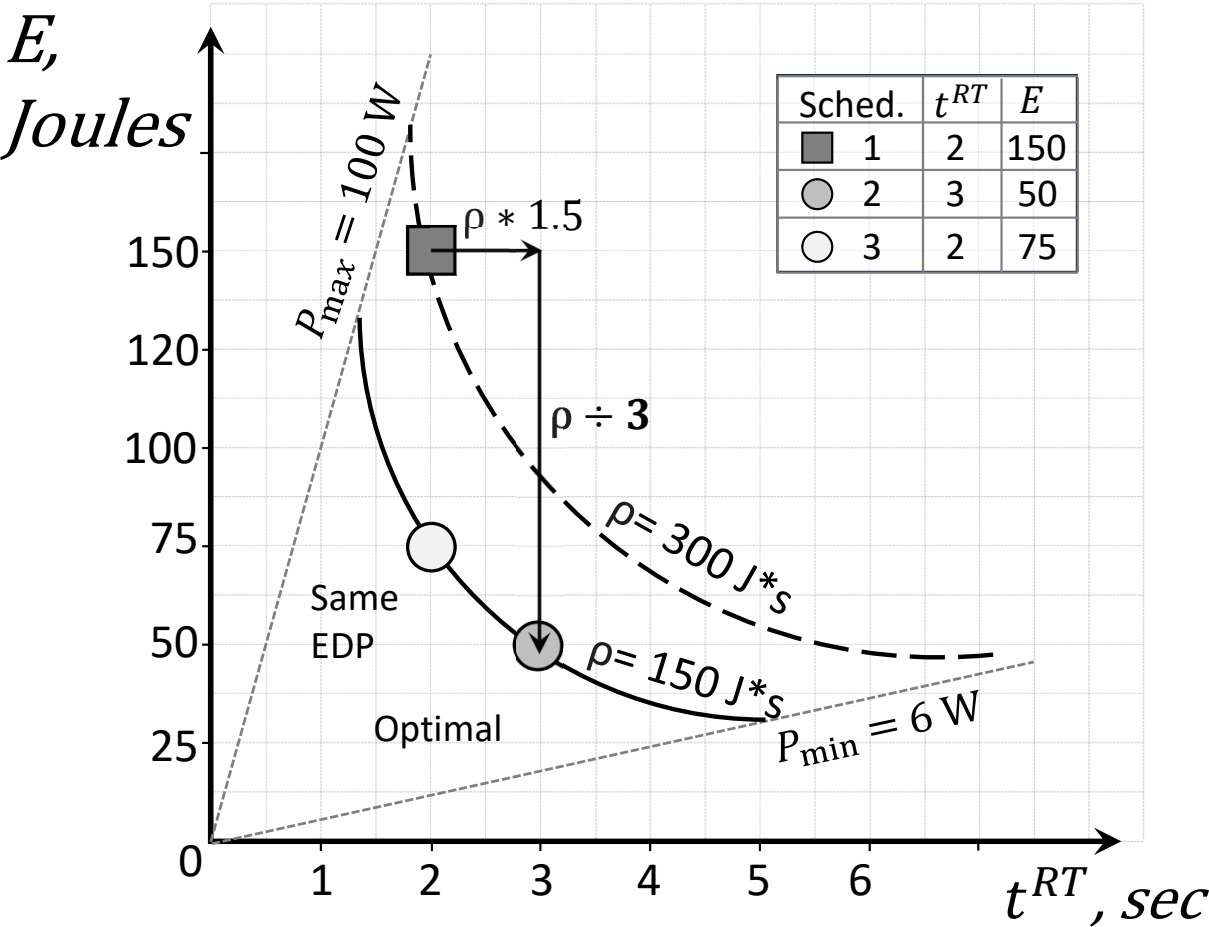




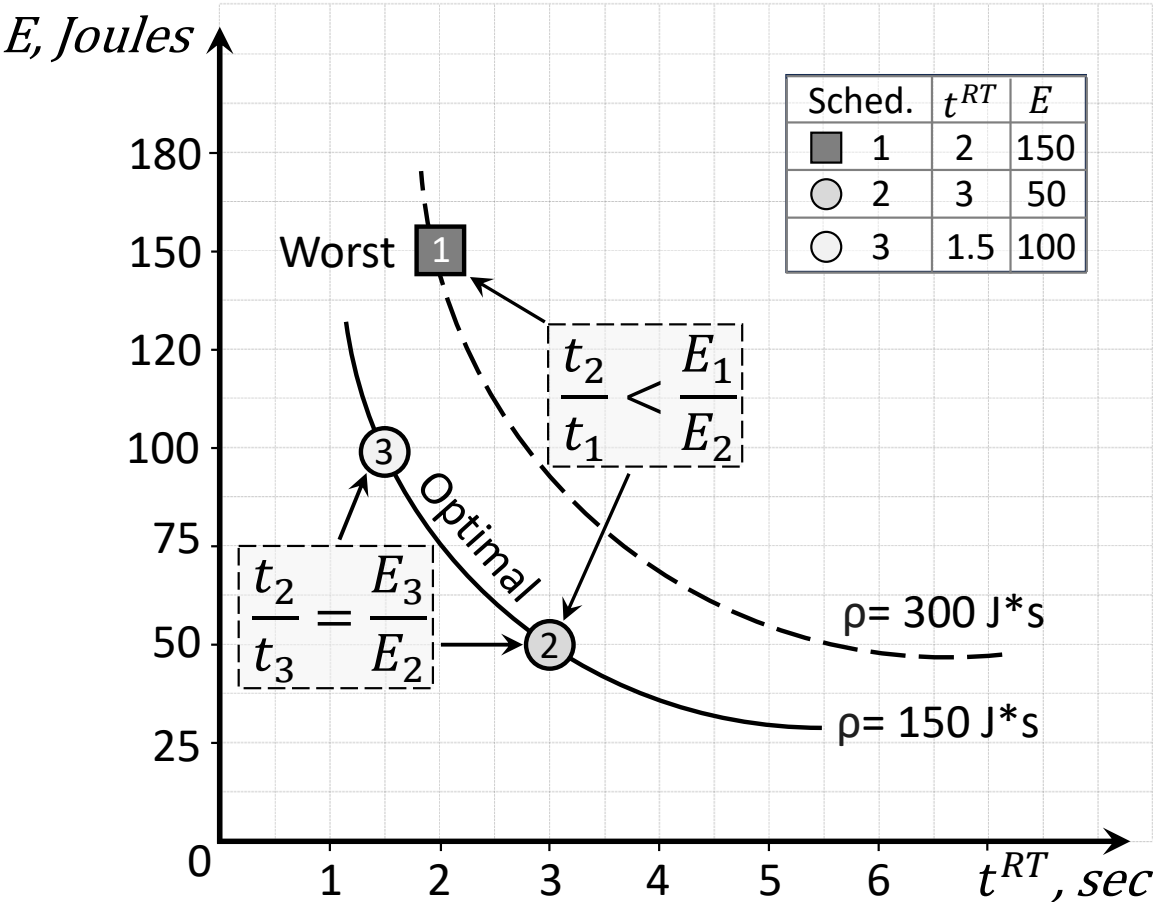




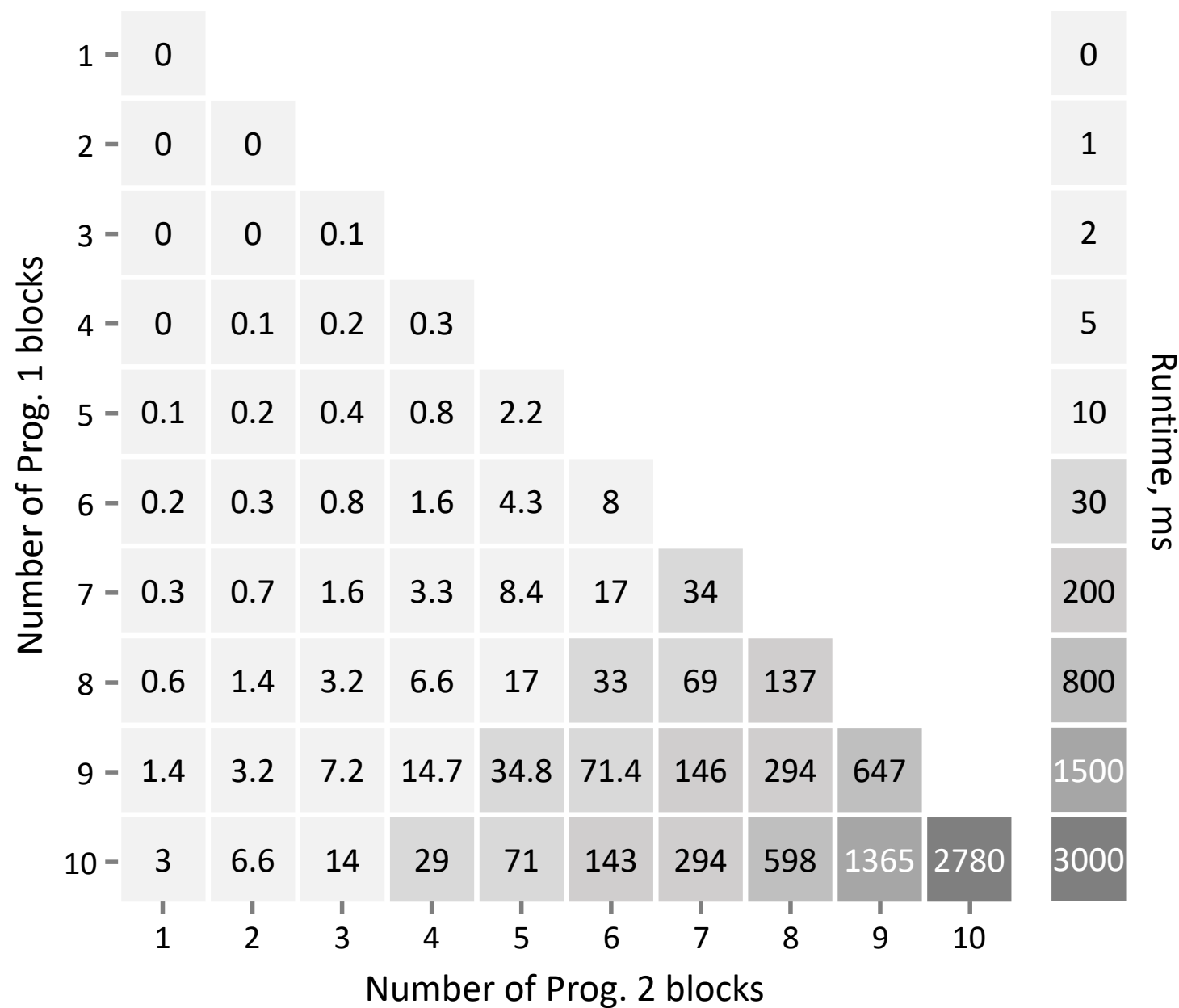
Energy efficiency



Energy efficiency

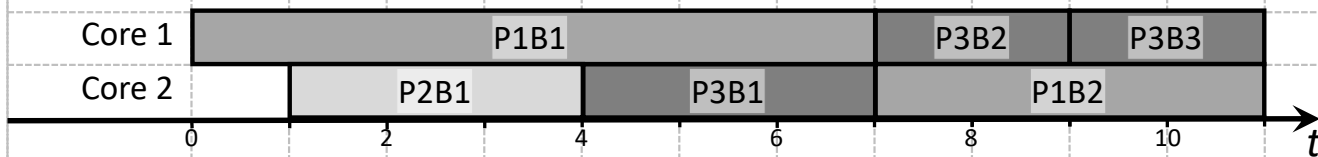


Experiments

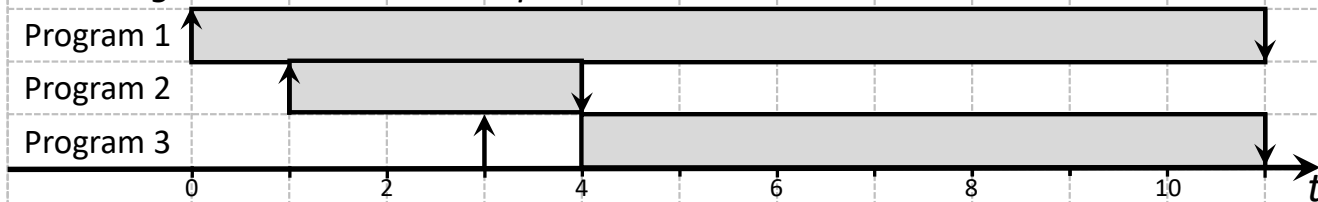


Notation: \uparrow - program release \downarrow - program completion \square - resource allocation

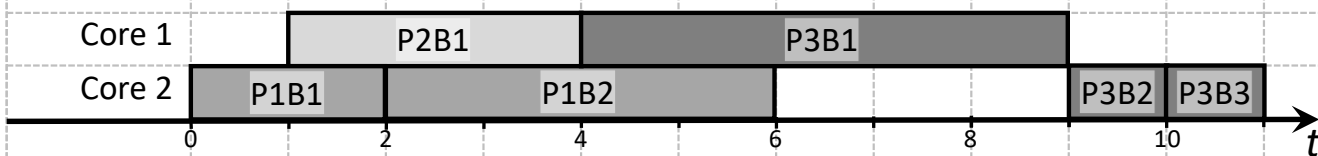
S1 – Cores Load



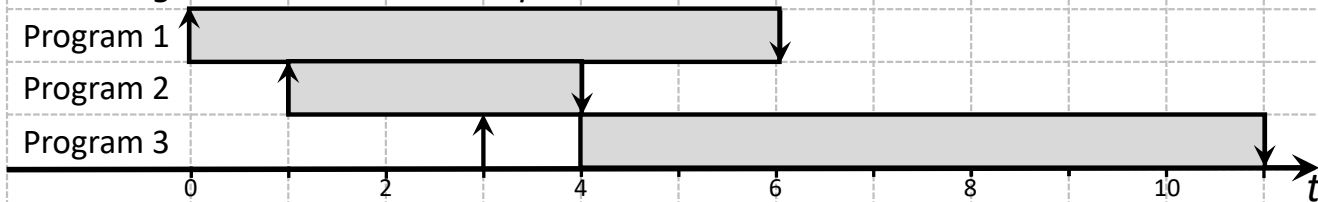
S1 – Programs Arrival and Completion



S2 – Cores Load

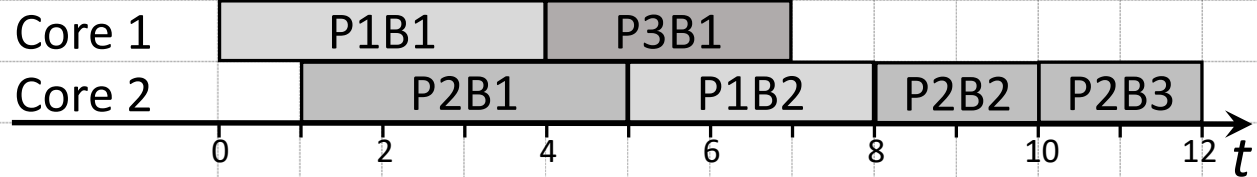


S2 – Programs Arrival and Completion

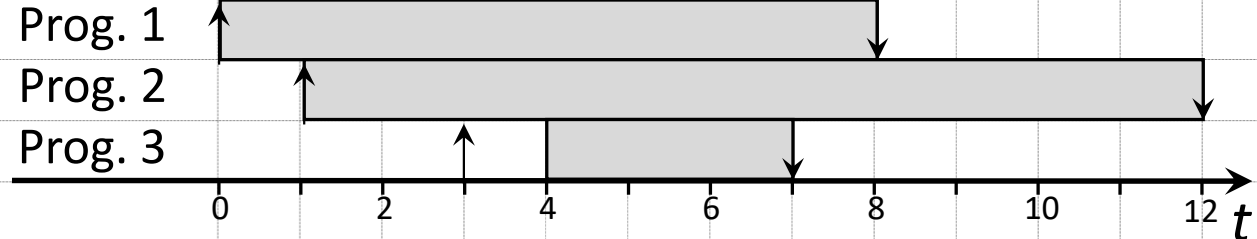


Notation: \uparrow - program release \downarrow - program completion \square - resource allocation

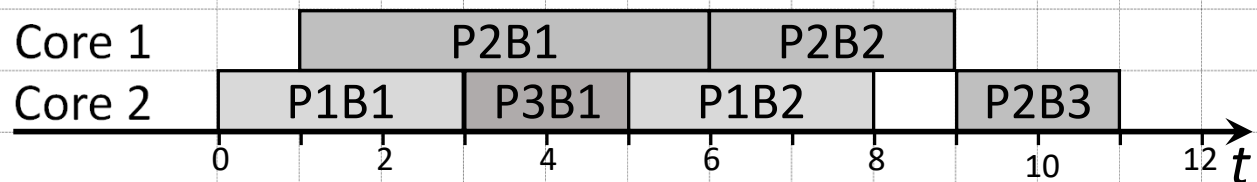
S1 – Cores Load



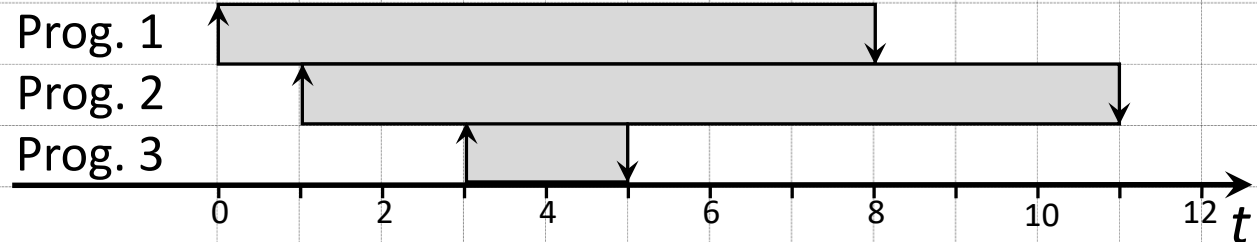
S1 – Programs Arrival

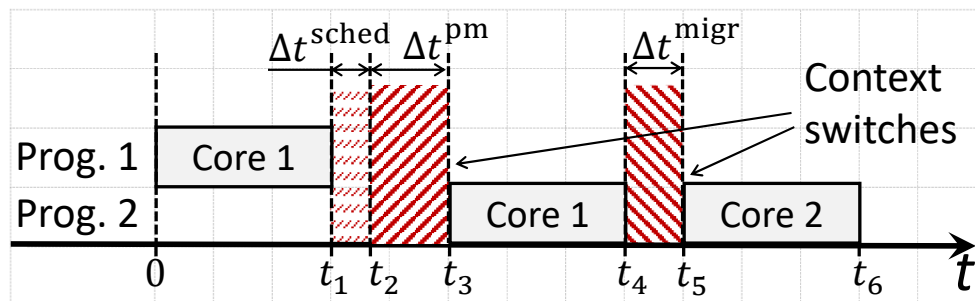
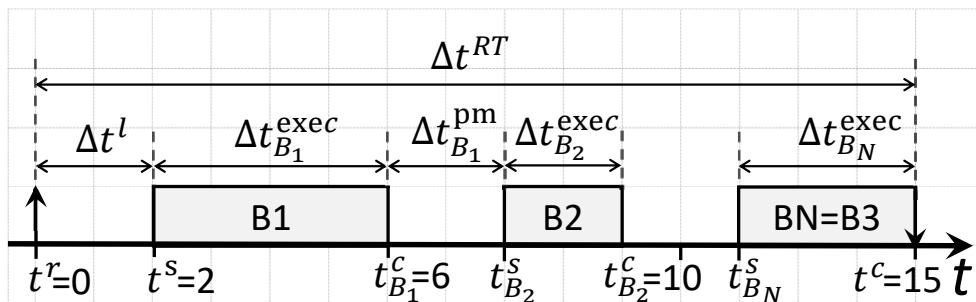
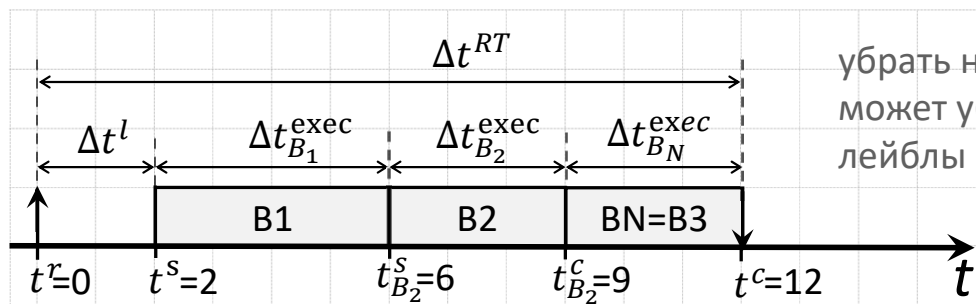


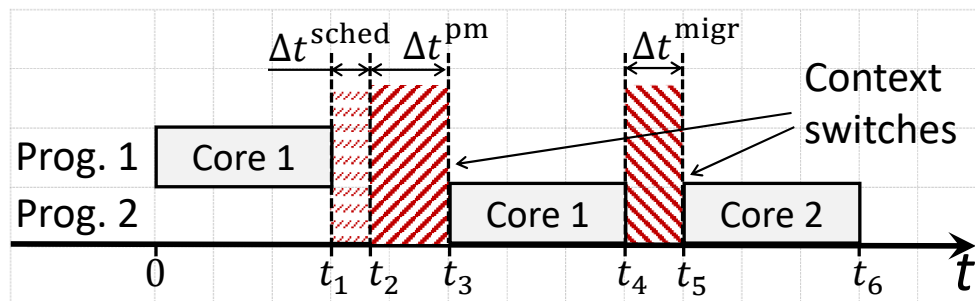
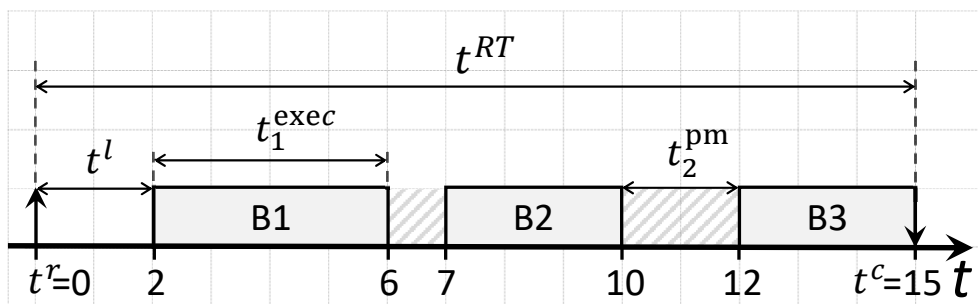
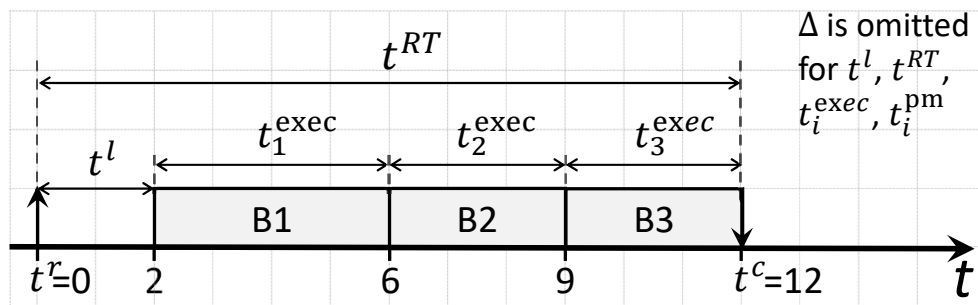
S2 – Cores Load



S2 – Programs Arrival



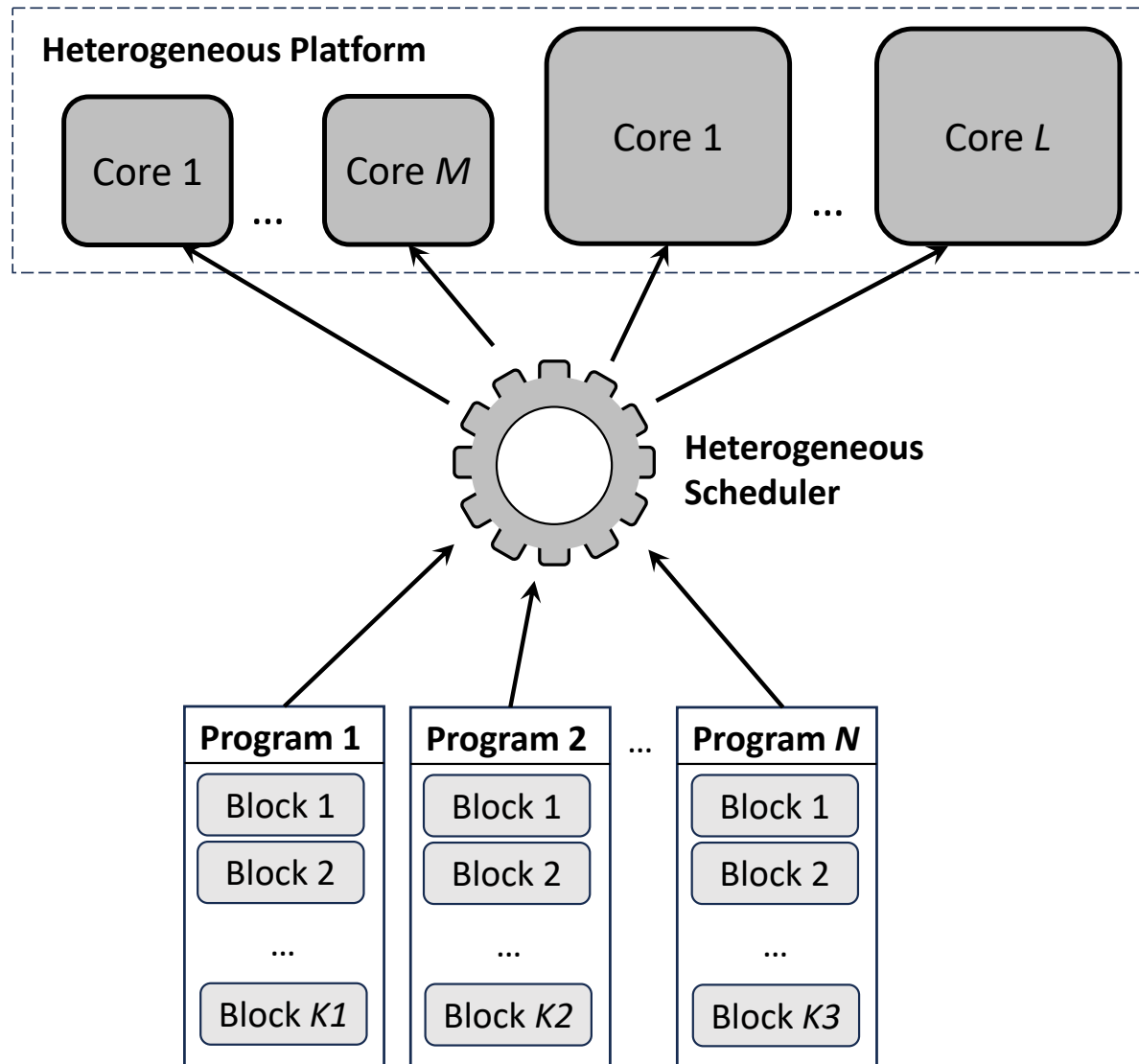




Program	Schedule	Scheduling Objective to Minimize		
		Latency	Execution Time	Response Time
Prog. 1	S1	0	7	8
	S2	0	6	8
Prog. 2	S1	0	8	11
	S2	0	10	10
Prog. 3	S1	1	3	4
	S2	0	2	2
Average	S1	0.33	6	7.67
	S2	0	6	6.67

Schedule	Program	Latency	Execution Time	Response Time	Average Execution Time	Average Response Time
S1	P1	0	7	8	6	7.67
	P2	0	8	11		
	P3	1	3	4		
S2	P1	0	6	8	6.33	7
	P2	0	11	11		
	P3	0	2	2		

Schedule	Program	Latency	Execution Time	Response Time	Average Execution Time	Average Response Time
S2	P1	0	6	8	6.33	7
	P2	0	11	11		
	P3	0	2	2		



Algorithm	Block №	Config 1	Config 2
Winograd	1	71	75
	2	5	5
	3	1813	2240
	4	1717	1985
	5	1773	1736
	6	1773	1802
	7	1781	917
	8	1616	755
	9	1775	757
	10	7	3
LZ4	1	40	35
	2	26	20
	3	50	46
	4	74	52
	5	31	26
	6	19736	9808
	7	965	311
	8	1	1
	9	15	5

Comparison point	Winograd	LZ4
Purpose	Matrix multiplication algorithm (for large matrices)	High-speed lossless data compression
Number of blocks	10	9
Improvement from switching to Config 2	x1.2	X2
Memory usage	High (depends on matrix size)	Low
Cache utilization	Inefficient (matrices are larger than cache)	Efficiently utilizes caches (64KB window)
Main block type	Recursive calls	Loops
Practical usage	Rarely used due to high constant factors	Extensively used in production
Main performance issues	Disk swapping, page faults	Slow memory access