

MuMMI_R: Analyzing and Modeling Power and Time under Different Resilience Strategies

Xingfu Wu and Valerie Taylor
Dept. of Comp. Sci. and Eng., Texas A&M University

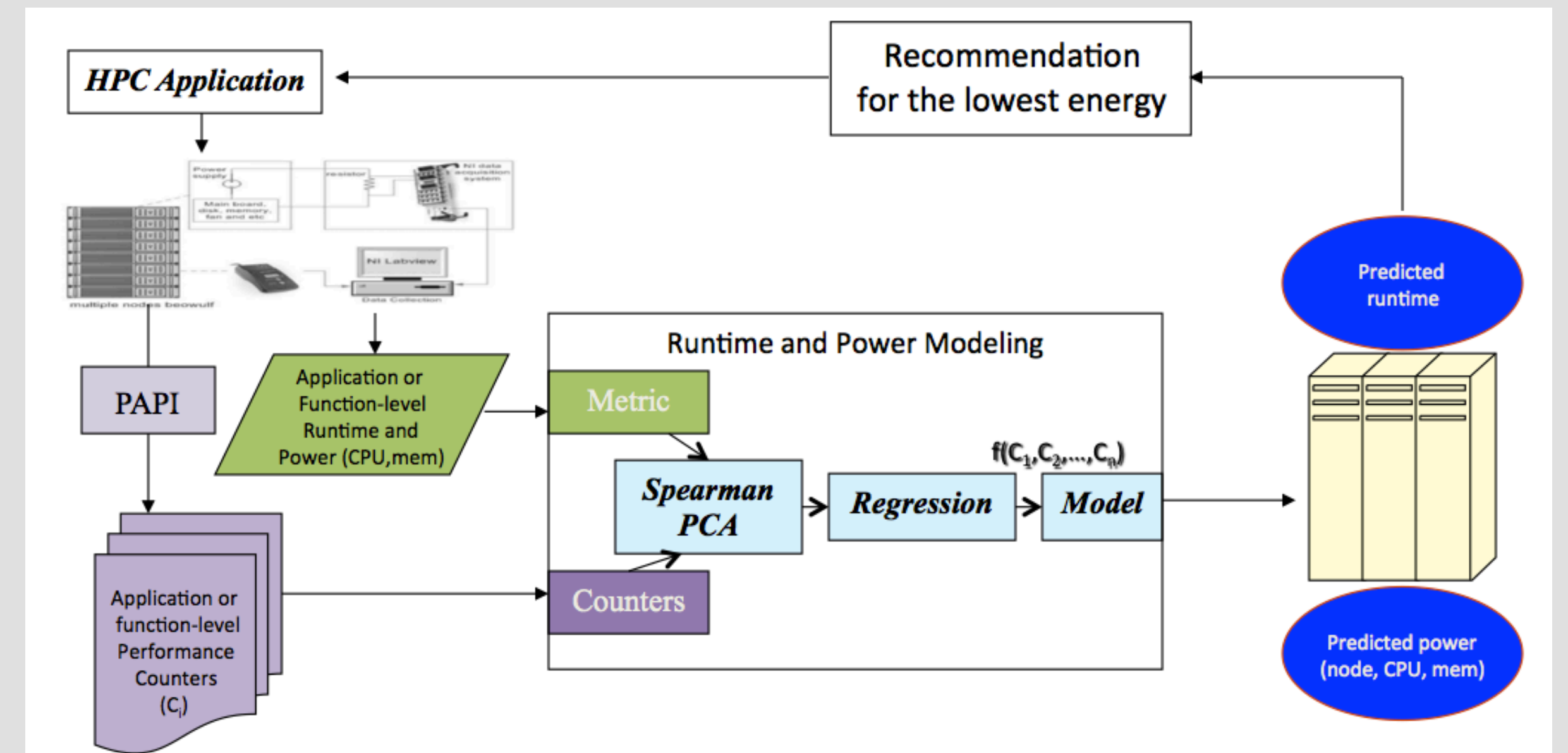
Zhiling Lan
Dept. of Comp. Sci., Illinois Institute of Technology



MuMMI_R Description

- **Goal:** To develop energy-efficient fault-tolerant HPC applications, the MuMMI (Multiple Metrics Modeling Infrastructure) is extended to analyze and model their energy and performance under different resilience strategies
- In this work, we extend the MuMMI to examine the tradeoffs among resilience, execution time, power and energy of the STREAM benchmark on three different architectures
- In the future, we will extend the MuMMI to model the tradeoffs between execution time, power, energy and resilience for various application-system configurations

MuMMI_R Modeling Framework



Resilience Strategies: Multilevel Checkpointing

- FTI is a fault tolerance interface that aims to add a highly reliable layer between the operating system and the application
- It provides five application level subroutines, *FTI_Init()*, *FTI_Protect()*, *FTI_Snapshot()*, *FTI_BitFlip()*, and *FTI_Finalize*.
- It has the following features for the initial configuration: four-level checkpointing (local write (L1), Partner copy (L2), Reed-Solomon coding (L3), and PFS write (L4)), checkpointing frequency, number of bit-flip failure injections, injection bit position, and injection frequency
- *ckp(1,3,5,7)* stands for frequencies for a four-level checkpoint

Three HPC Architectures and Environments

	IBM BG/Q	Intel Haswell	AMD Kaveri
CPU cores per Node	16	32	4
Sockets per node	1	2	1
CPU type and speed	PowerPC A2 1.6GHz	Xeon(R) CPU E5-2698 3.6GHz	AMD A10-7850K 3.7GHz
L1 cache per core	D:16KB/I:16KB	D:32KB/I:32KB	D:16KB/I:96KB
L2 cache per socket	32MB (shared)	256KB (per core)	2MB (shared)
L3 cache per socket	None	40MB (shared)	None
Memory per Node	16GB	128GB	16GB
Network	5D Torus	Mellanox FDR InfiniBand	Mellanox FDR InfiniBand
HW Threads per core	4	2	2
Power tools	EMON/MonEQ	PowerInsight	PowerInsight

Runtime (s), Average Power (W) and Energy (J) under Different Resilience Strategies

Frequency	Runtime	Node Power	CPU Power	Memory Power	Network Power	Energy
Original	216	53.56	32.51	7.76	1.86	11568.96
<i>ckp(1,2,3,4)</i>	268	53.35	32.49	7.59	1.86	14297.80
<i>ckp(1,3,5,7)</i>	254	53.46	32.54	7.64	1.86	13578.84
<i>ckp(2,3,4,5)</i>	256	56.39	34.60	8.52	1.85	14435.84
<i>ckp(2,4,6,8)</i>	227	54.01	32.20	8.70	1.89	12260.27
<i>ckp(3,4,5,6)</i>	229	57.31	35.17	8.66	1.88	13123.99
<i>ckp(3,5,7,11)</i>	229	56.54	34.64	8.63	1.85	12947.66
<i>ckp(3,5,7,9)</i>	231	57.25	35.12	8.65	1.88	13224.75
<i>ckp(4,5,6,7)</i>	217	57.33	35.13	8.71	1.88	12440.61

On BG/Q

Frequency	Runtime	Node Power	CPU Power	Memory Power	Disk Power	Energy
Original	800	311.26	230.34	64.83	2.32	249008.00
<i>ckp(1,2,3,4)</i>	1832	219.24	161.42	41.69	2.31	401647.68
<i>ckp(1,3,5,7)</i>	1848	230.29	170.83	43.32	2.33	425575.92
<i>ckp(2,3,4,5)</i>	1618	243.08	180.58	46.36	2.32	393303.44
<i>ckp(2,4,6,8)</i>	1311	261.27	193.95	51.17	2.33	342524.97
<i>ckp(3,4,5,6)</i>	1460	253.43	187.23	50.04	2.32	370007.80
<i>ckp(3,5,7,11)</i>	1396	261.87	194.46	51.28	2.33	365570.52
<i>ckp(3,5,7,9)</i>	1358	265.54	197.50	51.90	2.32	360603.32
<i>ckp(4,5,6,7)</i>	1450	259.06	191.50	51.39	2.33	375637.00

On Haswell

Frequency	Runtime	Node Power	CPU Power	Memory Power	Disk Power	Energy
Original	486	78.17	44.03	16.45	0.85	37990.62
<i>ckp(1,2,3,4)</i>	535	77.83	44.15	16.01	0.82	41639.05
<i>ckp(1,3,5,7)</i>	530	76.51	43.36	15.59	0.88	40550.30
<i>ckp(2,3,4,5)</i>	527	76.90	43.51	15.69	0.86	40526.30
<i>ckp(2,4,6,8)</i>	502	77.73	43.67	16.39	0.84	39020.46
<i>ckp(3,4,5,6)</i>	519	77.24	43.88	15.76	0.84	40087.56
<i>ckp(3,5,7,11)</i>	509	77.54	43.71	16.13	0.82	39467.86
<i>ckp(3,5,7,9)</i>	509	77.42	43.89	15.90	0.84	39406.78
<i>ckp(4,5,6,7)</i>	518	76.72	43.47	15.70	0.87	39740.96

On Kaveri

Power over Time on Three Architectures

