

Funded by the Horizon 2020 Framework Programme of the European Union MSCA-ITN-EID

IBM **Research** | Zurich

NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning

<u>Gagandeep Singh</u>, Juan Gomez-Luna, Giovanni Mariani, Geraldo F. Oliveira, Stefano Corda, Sander Stuijk, Onur Mutlu, Henk Corporaal

56th Design Automation Conference (DAC), Las Vegas 4th-June-2019





Executive Summary

- Motivation: A promising paradigm to alleviate data movement bottleneck is nearmemory computing (NMC), which consists of placing compute units close to the memory subsystem
- **Problem:** Simulation times are extremely slow, imposing long run-time especially in the early-stage design space exploration
- **Goal:** A quick high-level performance and energy estimation framework for NMC architectures

Our contribution: NAPEL

- Fast and accurate performance and energy prediction for previously-unseen applications using ensemble learning
- Use intelligent statistical techniques and micro-architecture-independent application features to minimize experimental runs

Evaluation

- NAPEL is, on average, 220x faster than state-of-the-art NMC simulator
- Error rates (average) of 8.5% and 11.5% for performance and energy estimation

We open source Ramulator-PIM: https://github.com/CMU-SAFARI/ramulator-pim/



Michael Wise, ASTRON, "Science data Centre challenges", DOME Symposium, 18 May, 2017



Massive amounts of data



Michael Wise, ASTRON, "Science data Centre challenges", DOME Symposium, 18 May, 2017



* R. Nair et al., "Active memory cube: A processing-in memory architecture for exascale systems", IBM J. Research Develop., vol. 59, no. 2/3, 2015

- **Compute Centric Approach**
- Memory hierarchies take advantage of *locality*
 - Spatial locality

Data movement bottleneck

- Applications are increasingly data hungry
- Data movement energy dominates compute
 - Especially true for off-chip movement



Floating point

Integer core Clock

> Leakage L1D

L1P L1P to L2 bus

L2 cache

2 to DDR hus

Data Movement

4

Processo

Data Access

System-level power break down*

Paradigm Shift - NMC

- Compute-centric to a data-centric approach
- Biggest enabler stacking technology





NMC Simulators

- Simulation for:
 - Design space exploration (DSE)
 - Workload suitability analysis
- NMC Simulators:
 - Sinuca, 2015
 - HMC-SIM, 2016
 - CasHMC, 2016
 - Smart Memory Cube (SMC), 2016
 - CLAPPS, 2017
 - Gem5+HMC, 2017
 - Ramulator-PIM¹, 2019

¹Ramulator-PIM: https://github.com/CMU-SAFARI/ramulator-pim/

NMC Simulators

- Simulation for:
 - Design space exploration (DSE)
 - Morkload cuitability analysis

Simulation of real workloads can be 10000x slower than native-execution!!!

Sindic Memory Cube (SMC), 2010

- CLAPPS, 2017
- Gem5+HMC, 2017
- Ramulator-PIM¹, 2019

¹Ramulator-PIM: https://github.com/CMU-SAFARI/ramulator-pim/

NMC Simulators

- Simulation for:
 - Design space exploration (DSE)
 - Markland quitability analysis

Idea: Leverage ML with statistical techniques for quick NMC performance/energy prediction

Smart Memory Cube (SMC), 2010

- CLAPPS, 2017
- Gem5+HMC, 2017
- Ramulator-PIM¹, 2019

¹Ramulator-PIM: https://github.com/CMU-SAFARI/ramulator-pim/

NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning

NAPEL Model Training



Phase 1: LLVM Analyzer



Phase 2: Microarchitecture Simulation



Phase 3: Ensemble ML Training



NAPEL Framework



NAPEL Prediction



Experimental Setup

- Host System
 - IBM POWER9
 - Power: AMESTER
- NMC Subsystem
 - Ramulator-PIM¹



- Workloads
 - PolyBench and Rodinia
 - Heterogeneous workloads such as image processing, machine learning, graph processing etc.
- Accuracy in terms of mean relative error (MRE)

¹https://github.com/CMU-SAFARI/ramulator-pim/

NAPEL Accuracy: Performance and Energy Estimates



NAPEL Accuracy: Performance and Energy Estimates



MRE of 8.5% and 11.6% for performance and energy



Speed of Evaluation

					0		1200 -		
Application		Prediction Time		lup					
Name	#DoE conf.	DoE run (mins)	Train+Tune (mins)	Pred. (mins)	ec		1000 -	256 DoE	
atax	11	522	34.9	0.49	be	ے		configurations	
bfs	31	1084	34.2	0.48	S	to	800 -	for 12	
bp	31	1073	43.8	0.47	U C	lla	000	TOF 12	
chol	19	741	34.9	0.49	itic	มน	600	evaluated	
gemv	19	741	24.4	0.51	dic.	ar	600 -		
gesu	19	731	36.1	0.51	e G	2		applications	
gram	19	773	36.5	0.52	Рг	ē	400 -		
kme	31	742	36.9	0.55	<u>_</u> ~	8			
lu	19	633	37.9	0.51			200		
mvt	19	955	38.0	0.54	ΔP		200 -		
syrk	19	928	35.7	0.51	Ž				
trmm	19	898	37.6	0.48			0 -		
							1	DoE configura	itions 25

Speed of Evaluation



220x (up to 1039x) faster than NMC simulator

kme	31	742	36.9	0.55	s l	700		
lu	19	633	37.9	0.51	Ц Ц	200		
mvt	19	955	38.0	0.54	P	200 -		
syrk	19	928	35.7	0.51	NZ			
trmm	19	898	37.6	0.48		0 -]
						1	DoE configurations	256

Use Case: NMC Suitability Analysis

- Assess the potential of offloading a workload to NMC
- NAPEL provides accurate prediction of NMC suitability



• MRE between 1.3% to 26.3% (average 14.1%)

Conclusion and Summary

- Motivation: A promising paradigm to alleviate data movement bottleneck is nearmemory computing (NMC), which consists of placing compute units close to the memory subsystem
- **Problem:** Simulation times are extremely slow, imposing long run-time especially in the early-stage design space exploration
- **Goal:** A quick high-level performance and energy estimation framework for NMC architectures

Our contribution: NAPEL

- Fast and accurate performance and energy prediction for previously-unseen applications using ensemble learning
- Use intelligent statistical techniques and micro-architecture-independent application features to minimize experimental runs

Evaluation

- NAPEL is, on average, 220x faster than state-of-the-art NMC simulator
- Error rates (average) of 8.5% and 11.5% for performance and energy estimation

We open source Ramulator-PIM: https://github.com/CMU-SAFARI/ramulator-pim/



Funded by the Horizon 2020 Framework Programme of the European Union MSCA-ITN-EID

IBM **Research** | Zurich

NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning

<u>Gagandeep Singh</u>, Juan Gomez-Luna, Giovanni Mariani, Geraldo F. Oliveira, Stefano Corda, Sander Stuijk, Onur Mutlu, Henk Corporaal

56th Design Automation Conference (DAC), Las Vegas 4th-June-2019



