

# Impact of Manufacturing Variability in Power Constrained Supercomputing

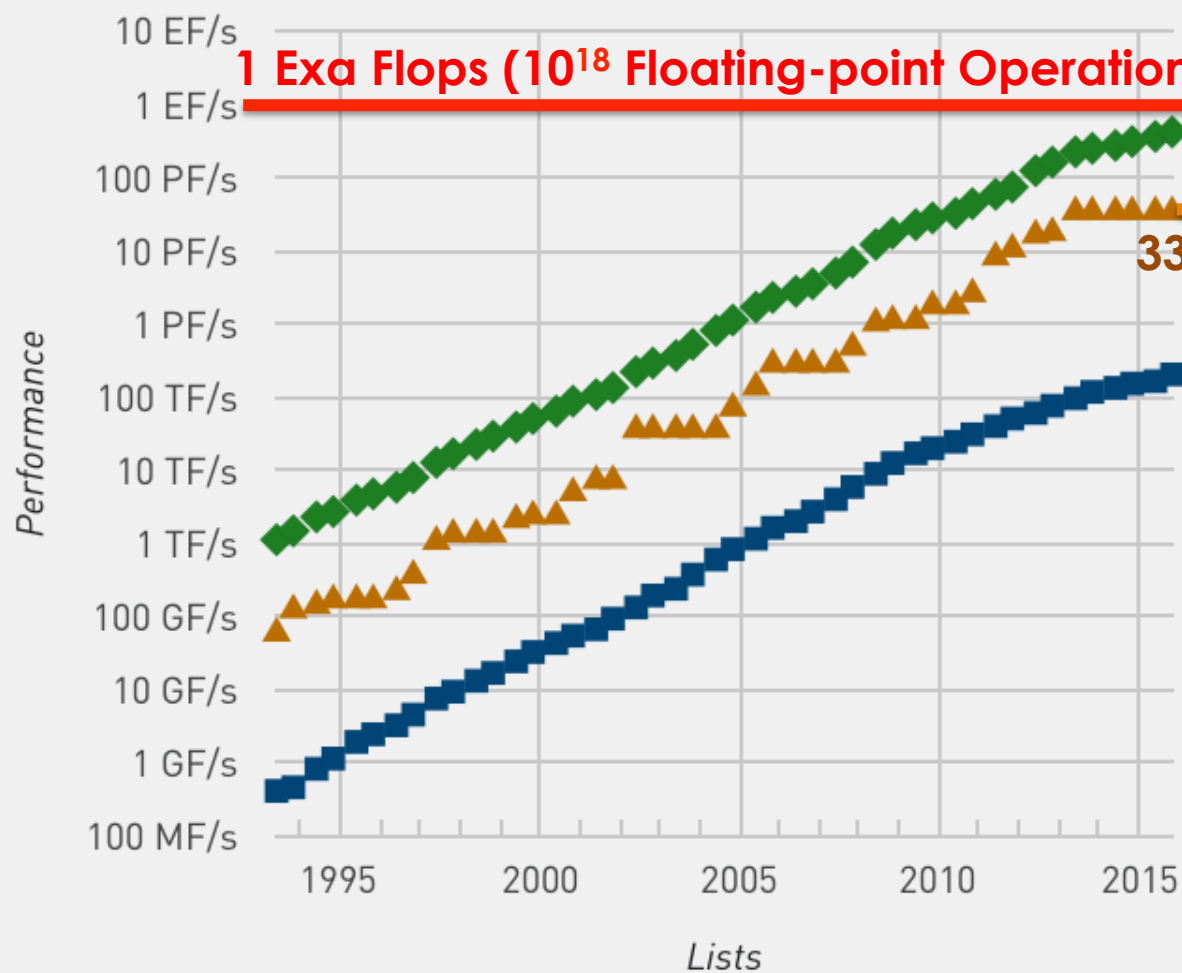
Koji Inoue

Kyushu University



# Trends of Supercomputing

## Performance Development



**1 Exa Flops (10<sup>18</sup> Floating-point Operations Per Second)**

World-Wide Next Target  
1 Exa Flopas @ 20 – 30 MW

30X

33.9 Peta Flops  
@17.8 MW

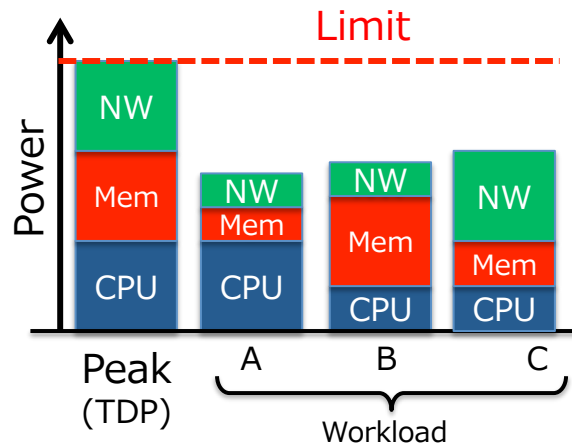
1.7X

Need to improve  
power efficiency!

# Overprovisioned Systems

## Under-provisioned (Conventional)

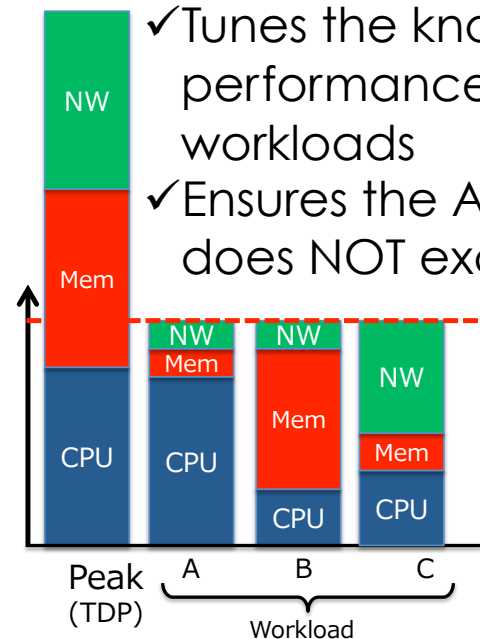
- HW Design
  - ✓Ensures the PEAK system power does NOT exceed the limit
- SW Design
  - ✓Tries to maximize the activity of HW components



TDP: Thermal Design Power

## Over-provisioned

- HW Design:
  - ✓Allows to install HWs w/o considering the power limit
  - ✓Provides power-performance knobs
- SW Design:
  - ✓Tunes the knobs to maximize the performance based on SW workloads
  - ✓Ensures the ACTUAL system power does NOT exceed the limit



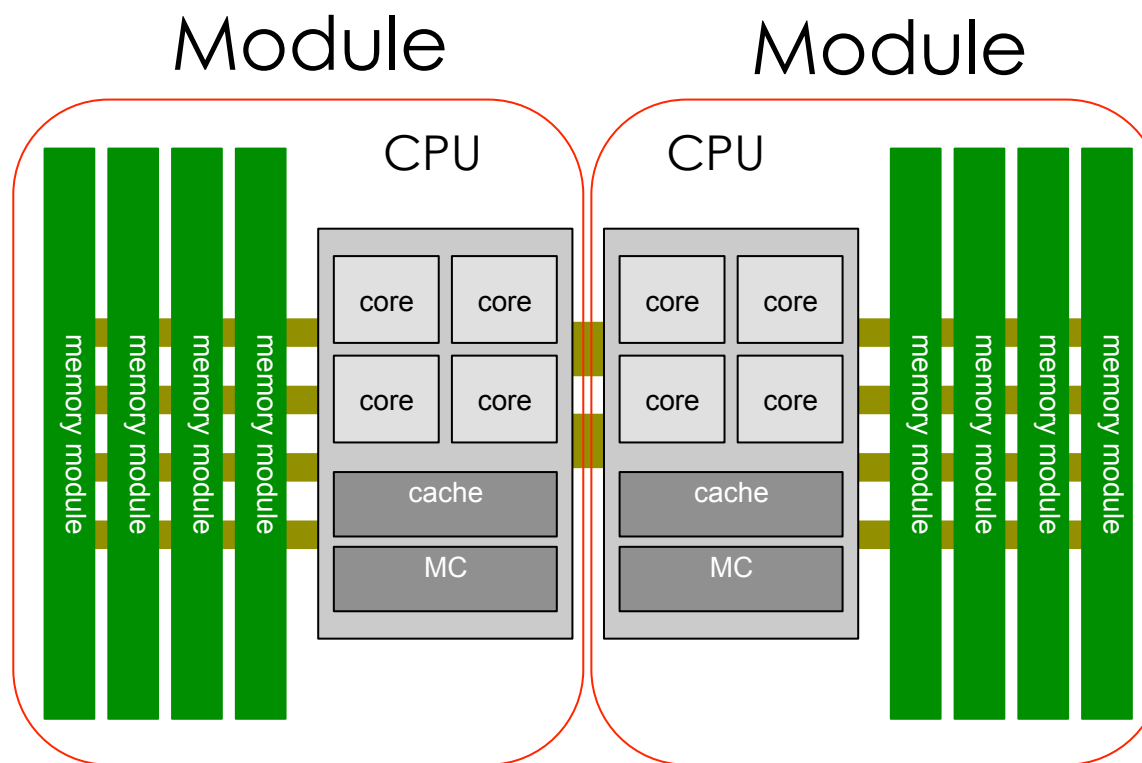
# Experimental Setup

- HPC Challenge: **star DGEMM**, **star STREAM(Triad)**
- NPB: **BT**, **SP**, **EP**
- Magneto Hydro-Dynamics(**MHD**) simulation
  - Typical stencil app. to simulate space plasma
  - Calculations and communications appear in turn
- Fiber benchmark suite: mVMC-mini (**mVMC**)
  - Variational Monte-Carlo simulation for strongly correlated electron system

Blue=EP type  
Red=With Comm. & Sync.

Site	Node Micro-Architecture	Total nodes	Procs. Per Node	Cores Per Procs.	Power Msrmt.
Cab(LLNL)	Intel E5-2670 Sandy Bridge	1,296	2	8	RAPL
BG/Q Vulcan (LLNL)	IBM PowerPC A2	24,576	1	16(compute)	EMON
Teller (SNL)	AMD A10-5800K Piledriver	104	1	4	PI
<b>HA8K(Kyushu Univ.)</b>	<b>Intel E5-2697v2 Ivy Bridge</b>	<b>965</b>	<b>2</b>	<b>12</b>	<b>RAPL</b>

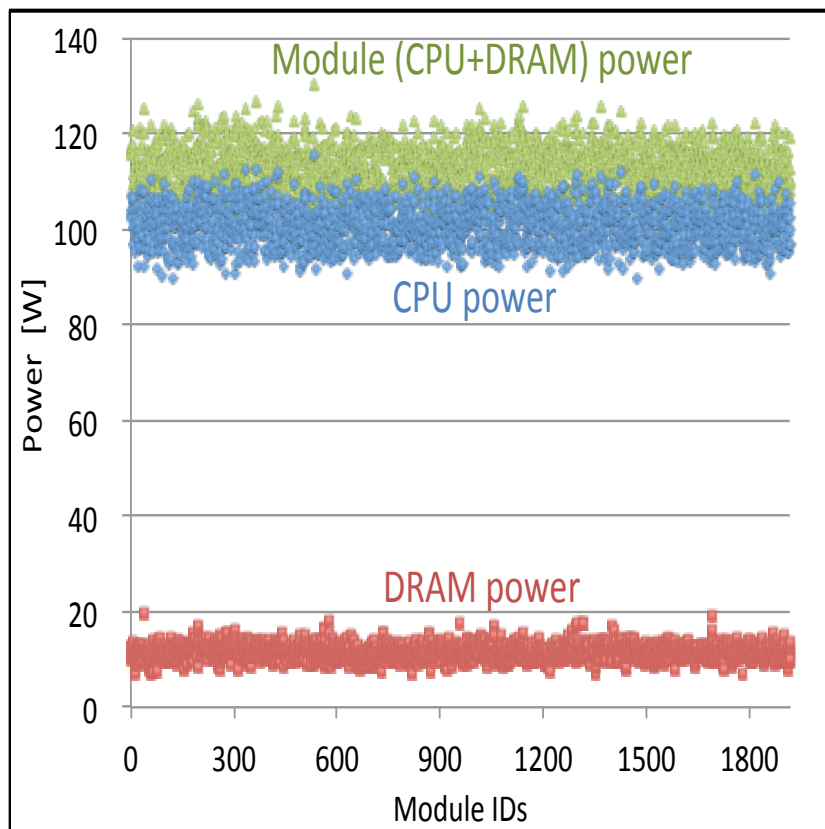
# Terminology



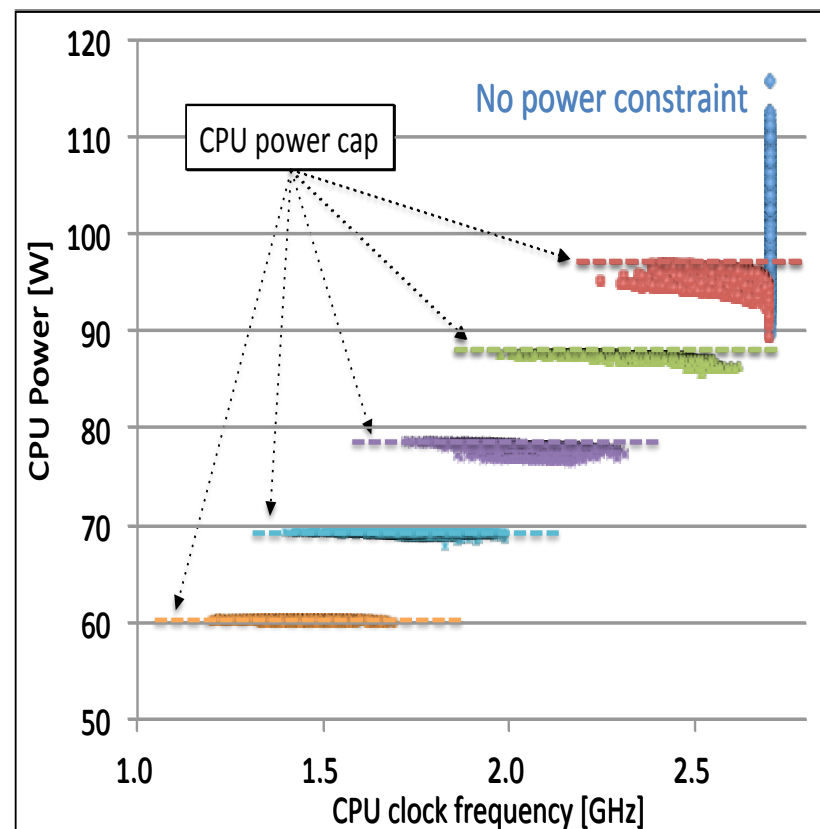
CPU = Processor chip (including cores, cache, MC, etc.)  
 Module = A pair of a CPU and DRAMs directly connected to it

# Impact on CPU Frequency

star DGEMM



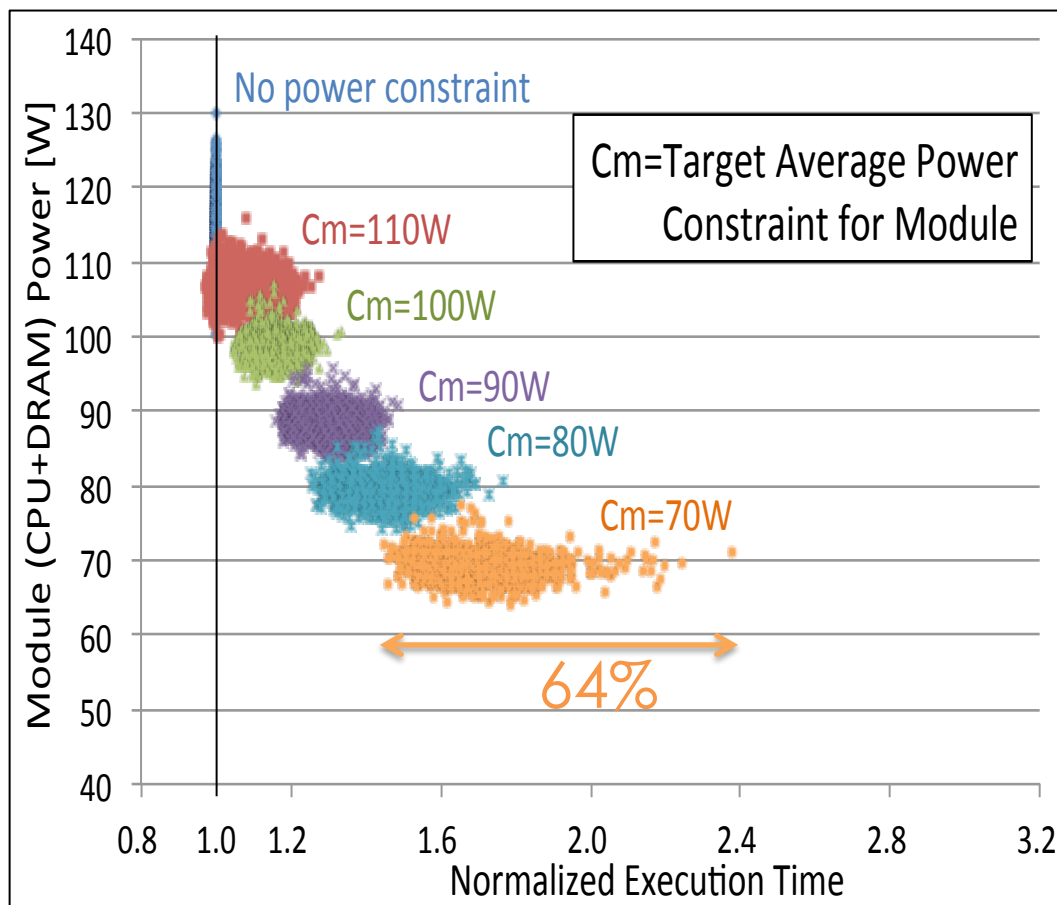
w/ a uniform power constraint



Power variation is translated into CPU frequency variation applying uniform power constraint!

# Impact on Application Performance

star DGEMM w/ a uniform power constraint



# Problem and Goal

- Power-Constraint Supercomputing
  - will be applied to future HPC systems
- Manufacturing Variability
  - leads to performance variation under power constraint

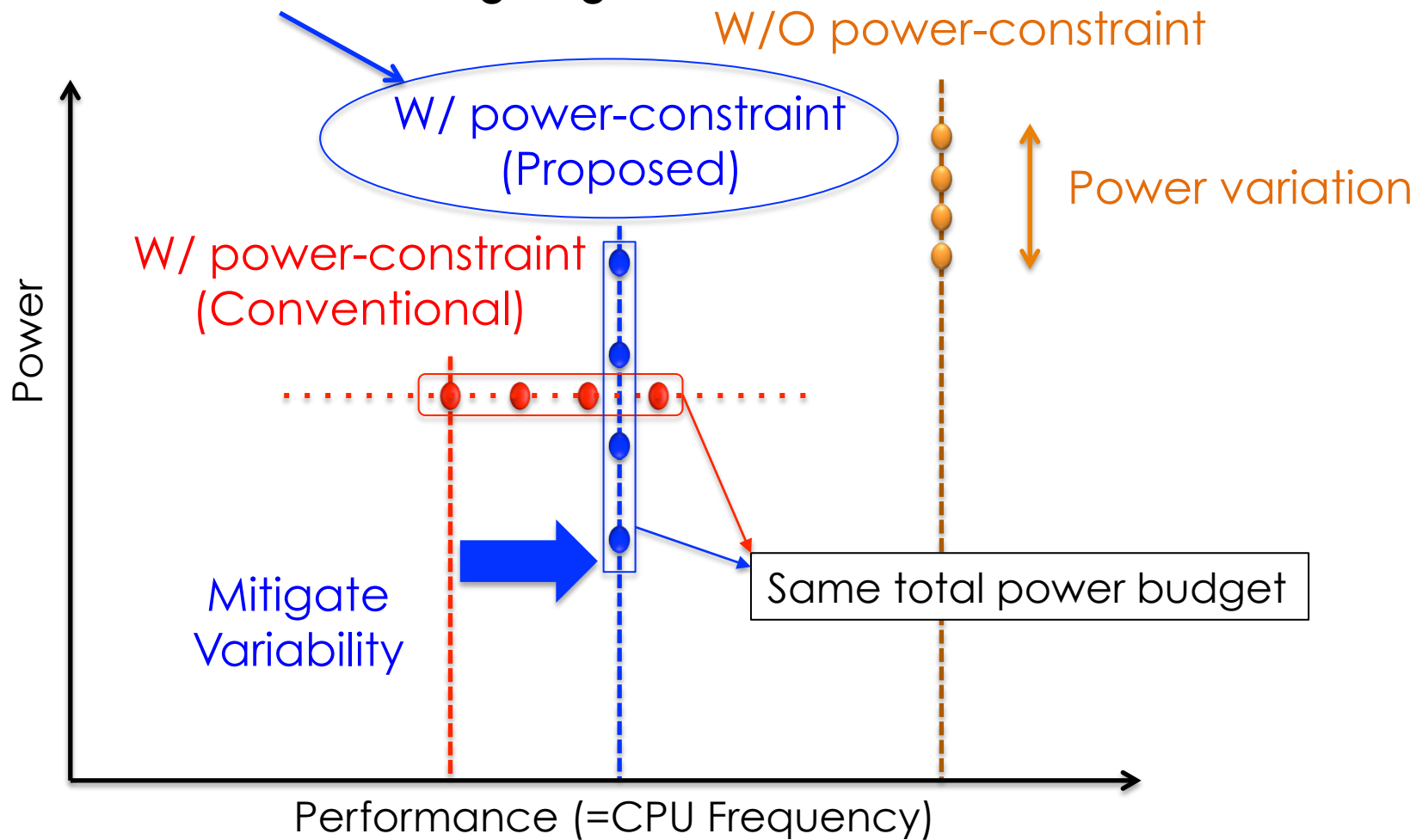
## Our Goal

Mitigate the impact of manufacturing variability on performance of HPC apps. under power constraint !

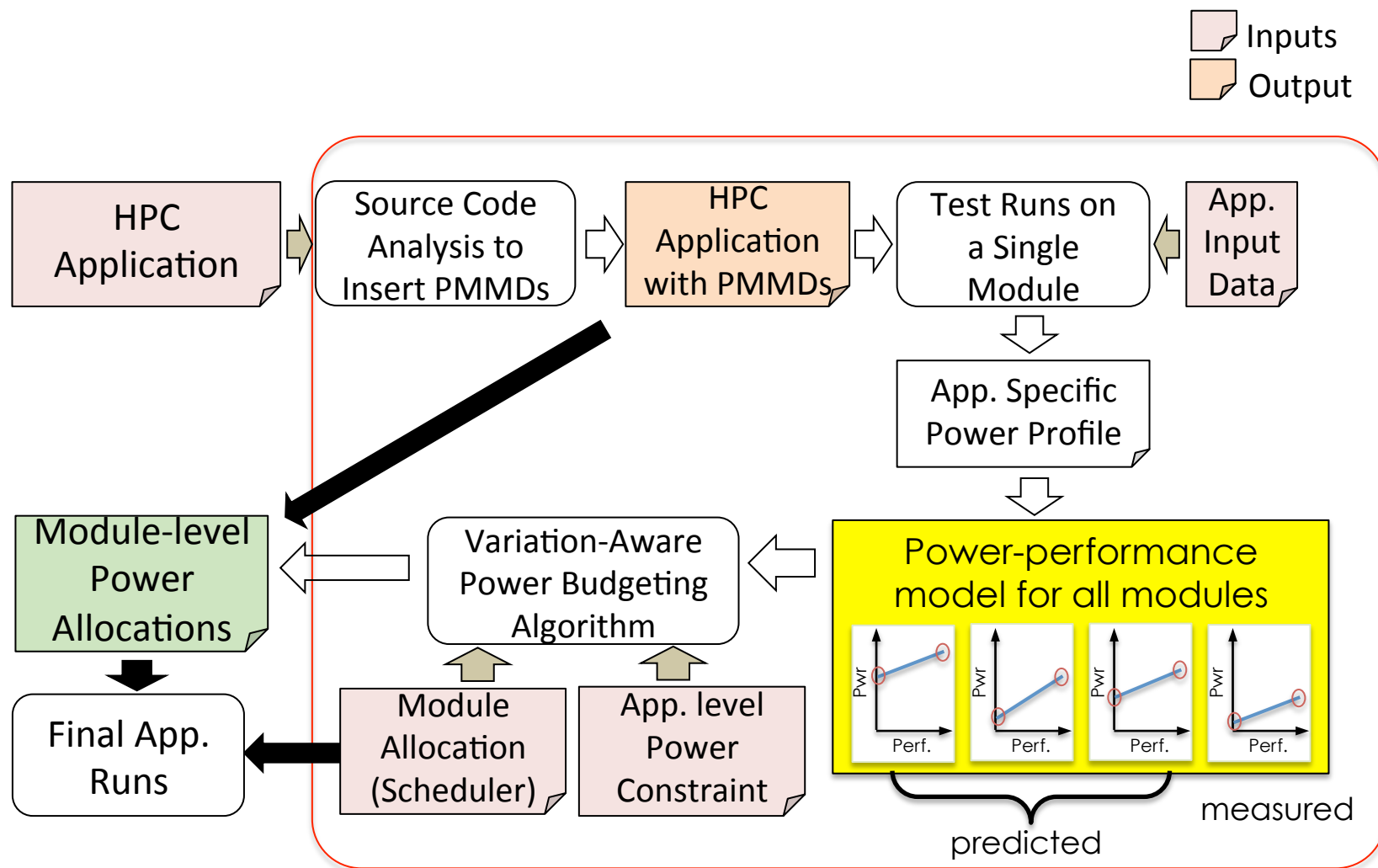


# Concept

## Variation-Aware Power Budgeting



# Variation-Aware Power Budgeting Strategy



# Power Model Calibration

Application-independent Power Variation Table (PVT)

Module ID	Normalized Power
1	1.0
:	:
k	1.2
:	:
N	0.8

Obtained once at system installation

Test run on a module!

Estimated Application Specific Power Consumption

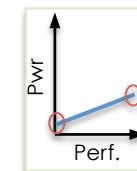
Module ID	Power Consumption
1	
:	
k	<b>120W</b>
:	
N	

Measured power on single module k

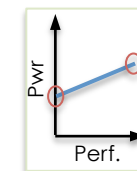
Module ID	Power Consumption
k	<b>120W</b>



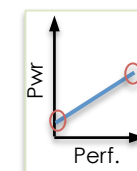
Module 1



Module 2

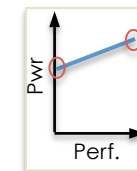


Module 3



⋮

Module N



# Power Model Calibration

Application-independent Power Variation Table (PVT)

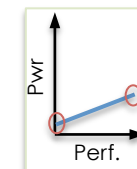
Module ID	Normalized Power
1	1.0
:	:
k	<b>1.2</b>
:	:
N	0.8

Estimated Application Specific Power Consumption

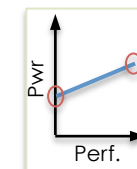
Module ID	Power Consumption
1	
:	
k	<b>120W</b>
:	
N	



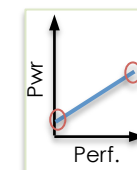
Module 1



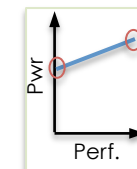
Module 2



Module 3



Module N



120W/1.2

application dependent average power

Measured power on single module k

Module ID	Power Consumption
k	<b>120W</b>

# Power Model Calibration

Application-independent Power Variation Table (PVT)

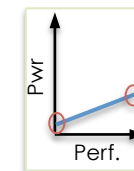
Module ID	Normalized Power
1	1.0
:	:
k	1.2
:	:
N	0.8

Estimated Application Specific Power Consumption

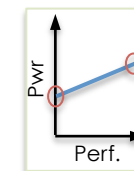
Module ID	Power Consumption
1	
:	:
k	120W
:	:
N	80W



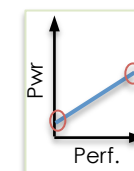
Module 1



Module 2

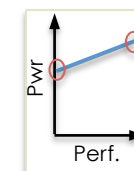


Module 3



⋮

Module N



$$\left( \frac{120W}{1.2} \right) \times 0.8$$

application dependent power on module-N

Measured power on single module k

Module ID	Power Consumption
k	120

# Options for Power Setting

## Two options for power settings

- **Power Capping** (Pc) using RAPL
- **Frequency Selection** (Fs) using CPUFreqlibs

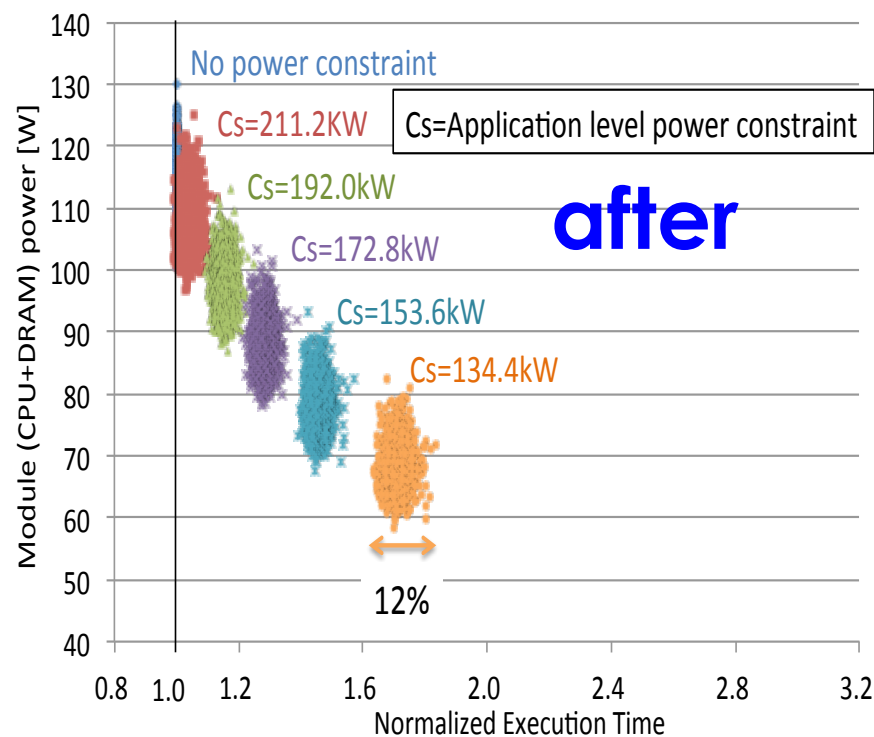
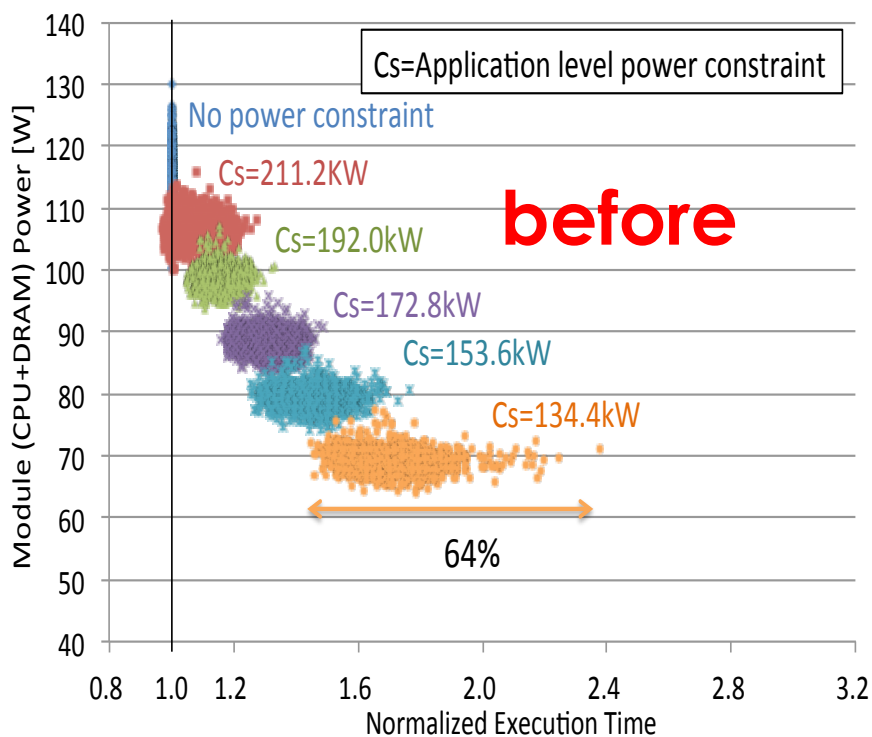
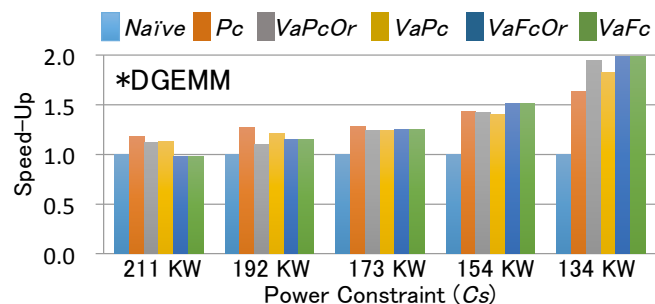
	Power Capping (Pc)	Frequency Selection (Fs)
Power Constraint	⊙ Guaranteed	△ Not guaranteed
Performance Equivalence	△ Not guaranteed	⊙ Guaranteed

# Tested Power Budgeting Methods

Method Name	Application Specific	Variation Aware	Power-Performance Model	Pwr. Set.
Naive	No	No	-----	Power Cap
Pc	Yes	No	Calibration	Power Cap
VaPc	Yes	Yes	Calibration	Power Cap
VaFs	Yes	Yes	Calibration	Freq. Sel.
VaPcOr	Yes	Yes	Oracle	Power Cap
VaFsOr	Yes	Yes	Oracle	Freq. Sel.

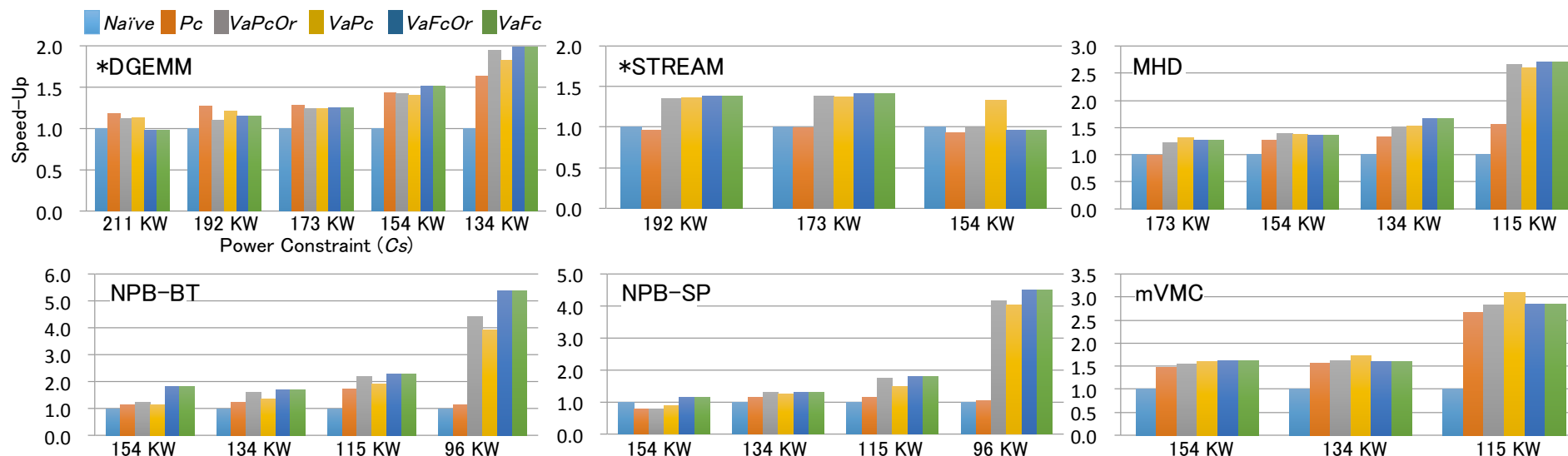
Va=Variation-Aware, Pc=Power Capping, Fs=Frequency Selection  
Or=Observed power data are used

# Speedup Ratios Normalized to Naïve (star DGEMM on 1,920 modules)





# Speedup Ratios Normalized to Naïve (All results on 1,920 modules)



5.4X speedup at maximum (NPB-BT)

1.8X speedup in average

# Conclusions

**Power constrained computing** becomes main-stream!

**Manufacturing variability** causes serious performance issue!

**Optimize** power resource allocation!

# Acknowledgements

- This research was supported by JST CREST.
- Special thanks to Dr. Yuichi Inadomi, Prof. Masaaki Kondo, Dr. Tapasya Patki, Dr. Martin Schulz, and other all members of this project.