

Adaptive Parallelism: Integrated Performance, Power, and Resilience (PPR) Modeling

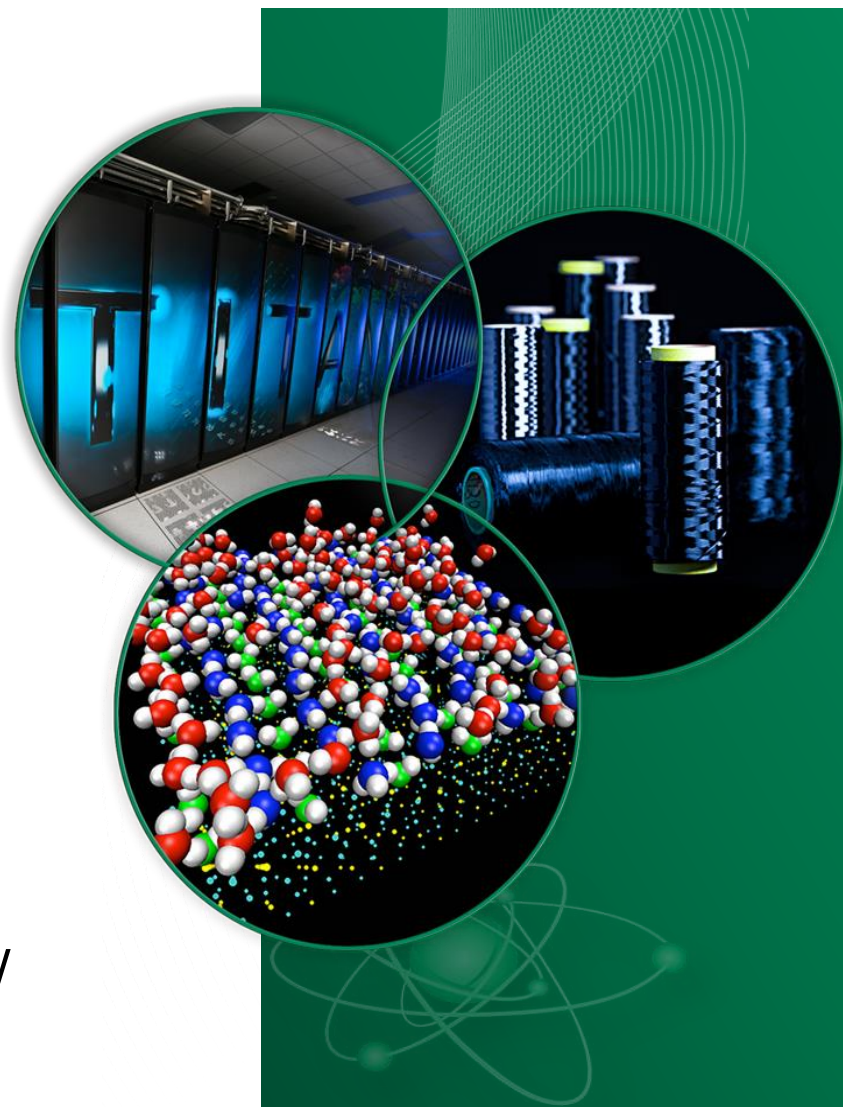
Dong Li⁺

Edgar A. Leon^{*}

Bronis R. de Supinski^{*}

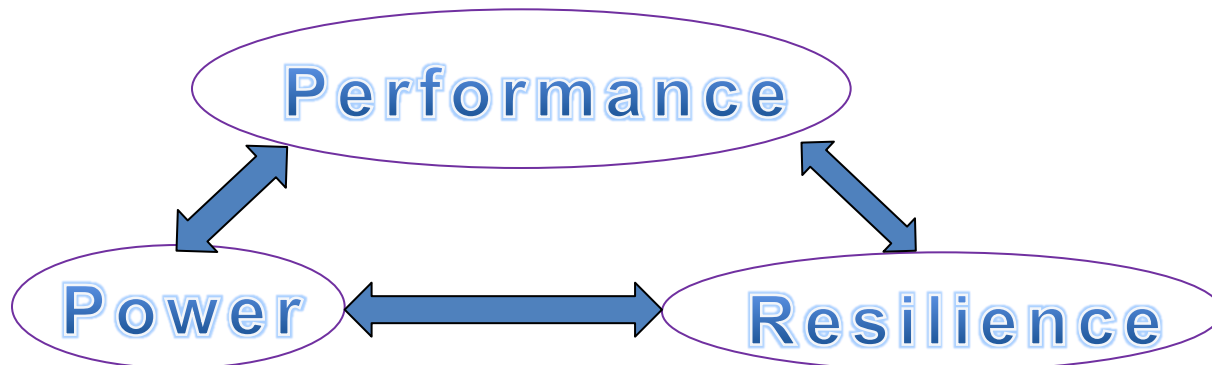
⁺Oak Ridge National Laboratory

^{*}Lawrence Livermore National Laboratory



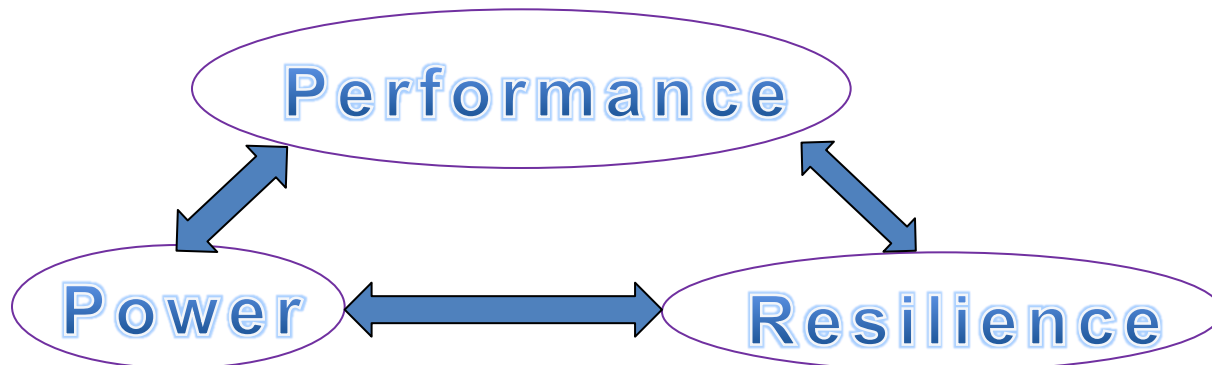
Performance, power and resilience are the critical challenges for large scale systems

- Increased parallelism creates unprecedented challenges to achieve the expected levels of application **performance**
- Unmanageable system **power consumption** prevents system scaling and could negatively impact performance
- A variety of faults and errors with increasing frequency necessitates system **resilience** to enable applications to run efficiently to completion and achieve correct results in a timely manner



The need for integrated modeling

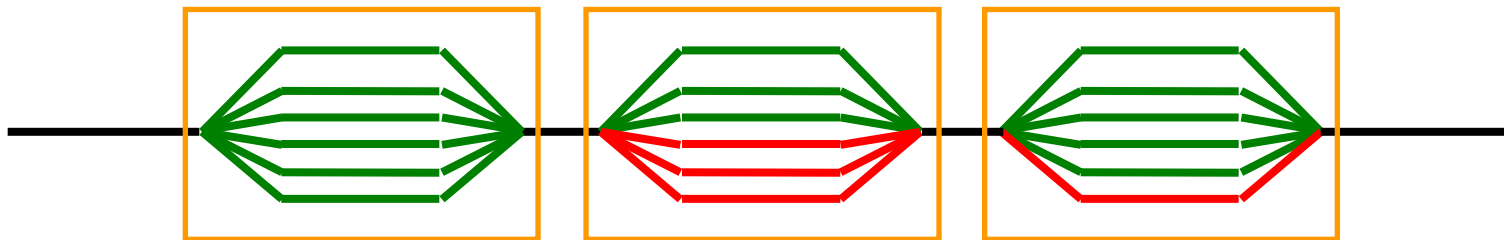
- The cost and benefit trade-offs across the three factors are not well understood
 - Reason 1: No rapid and accurate solution
 - Reason 2: No good resilience metric
- Performance, power, and resilience are often considered in isolation
 - E.g., when optimizing performance and energy (e.g., DVFS), what is the impact on resilience?
 - E.g., when optimizing resilience (e.g., checkpoint and ABFT), what is the impact on energy?



Our research theme: adaptive parallelism

- Motivations
 - The inherent parallelism of scientific applications varies across execution phases
 - Matching the degree of parallelism (parallel configuration) for an application has complex PPR implications
 - What is the appropriate parallelism in terms of PPR?
- We propose an infrastructure for modeling parallelism and its combined effects on performance, power and resilience

Parallel regions



General methodology

- We use a model-driven approach to guide the selection and adaptation of (thread-level) parallel configurations based on two observations
 - Performance, power, and resilience have first-order or second-order correlations with hardware component utilization
 - Given a parallel region, PPR and thread-level parallelism are strongly correlated statistically
- We introduce a new resilience metric, the vulnerability factor (VF)

Adaptive Parallelism

First observation: PPR have correlation with hardware component utilization

- Performance perspective
 - The number of accesses to the memory hierarchy and the number of executed instructions serve as strong indicators of performance with various levels of parallelism
- Power perspective
 - Power consumption is related to hardware usage intensity
- Resilience perspective
 - Resilience is related to both application execution time and number of hardware accesses



Vulnerability factor (VF): a new resilience metric

- The definition of VF is straightforward to measure and quantify resilience

$$VF = FIT * T * N_{ha}$$

Hardware failure rate Execution time Number of accesses

N_{ha} : Number of accesses to hardware component

FIT : Hardware failure rate

T : Application execution time

- Capturing the effects of both hardware and application
- VF must be tied to a specific hardware component

Second observation: PPR and thread-level parallelism are strongly correlated

- Based on hardware components utilization and PPR information collected from a few samples of parallel configurations (named *seminal configurations*), we can predict PPR for untested parallel configurations



Learn

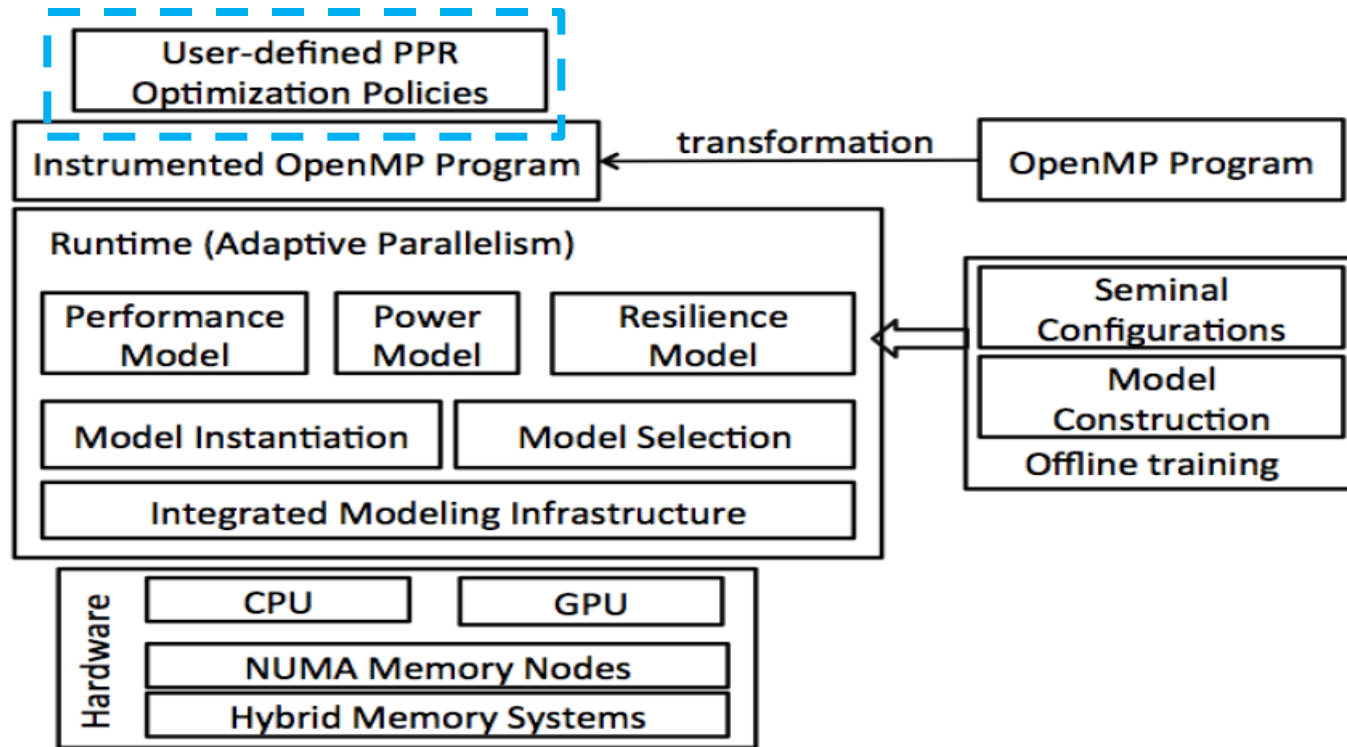


Predict

We construct an integrated PPR modeling infrastructure in two steps

- Offline model training
 - Using correlation analysis to determine the hardware component utilization information that is the most correlated with PPR
 - Measurable with lightweight hardware counters
 - Seminal configurations are determined with the k-means clustering algorithm
 - Building a series of PPR models
- Online model selection
 - Using a few sample iterations of parallel regions to execute with seminal configurations and measuring hardware component utilization
 - We can predict PPR for untested parallel configurations
 - Choosing which configurations result in the best PPR

Enabling runtime management based on PPR modeling



- Ongoing work

- Modeling memory-level parallelism
- Creation of models for hybrid memory architectures

Summary

- What is the major contribution of your research?
 - Designing a Integrated PPR modeling infrastructure
 - Enabling fast exploration of PPR tradeoff at runtime
 - Introducing a new resilience metric
- What is the bigger picture for your research area? What are the gaps?
 - Resilience modeling methodology (especially from the application perspective) is significantly lag behind
 - Sometimes interaction between performance, power, and resilience modeling is difficult
 - There is no single metric to capture the effects of PPR
 - We need to have fast exploration methods
- What major opportunities do you see for cross-pollination?
 - Need cooperation between the fields of application, system software, and hardware using a holistic view

Questions ??

