

# EVALUATION, OPTIMIZATION, AND APPLICATION OF EXECUTION MODELS FOR EXASCALE COMPUTING

**Exascale computing for DOE mission-critical applications will, by the end of this decade, be obstructed by severe challenges unlikely to be satisfied by incremental extensions to conventional practices. Strategic challenges include sustained performance, energy and power efficiency, dependability, and programmability combined with generality. Ultimately, a paradigm shift in how computing is defined and conducted will be essential to enable practical exascale performance by 2020. Such a paradigm is manifested as an execution model, providing the framework for deriving a consistent set of system elements and forming the basis for hardware/software co-design.**

*Adolfy Hoisie, Pacific Northwest National Laboratory (lead for the top-down approach)*

*Curtis Janssen, Sandia National Laboratories (lead for the bottom-up approach)*

*John Shalf, Lawrence Berkeley National Laboratory*

*Thomas Sterling, Indiana University*

The goals of our research are to define the salient terms for an effective co-design process, devise a methodology of co-design for DOE exascale projects, and develop an experimental execution model as a proof-of-concept foundation for exascale computing systems and applications. We will examine potential exascale execution models and determine their impact on exascale system performance and energy for DOE mission-critical applications. The proposed research will provide DOE with the essential techniques to derive, implement, deploy, and apply one or more revolutionary execution models as needed to enable exascale computing before the end of this decade. The strategy for the research integrates two concurrent and complementary approaches

within a single project focusing on both a bottom-up (EMBU) and a top-down (EMTD) approach. EMTD is related to the development, validation, modeling, analysis (including metrics), and quantitative comparison of advanced execution models, their relationship to a machine abstract model that will serve as the interface to enabling technologies and architecture, and to the physical properties of specific point designs. EMBU will start from concrete examples of execution models and hardware. In support of this the execution model toolkit (EXEMT) will be developed that will be a collection of both fine- and coarse-grained models for execution model components and will be rich enough to construct a variety of execution models of interest for exascale computing. Compact and

skeleton applications will be developed utilizing EXEMT to implement candidate execution models and to compare their estimated time and energy to solution on exascale machine models. Custom hardware features will be explored as a mechanism to accelerate the execution models. A quantitative, predictive co-design methodology for comparing execution models will be derived through both EMBU and EMTD and will be applied in the context of full applications in production use on current leading ASCR supercomputers, and targeted to evolve to exascale through the work of the Co-Design Centers and other DOE SC Exascale projects. These findings will be used in part to engage processor vendors and the ASCR Exascale Co-design Centers.

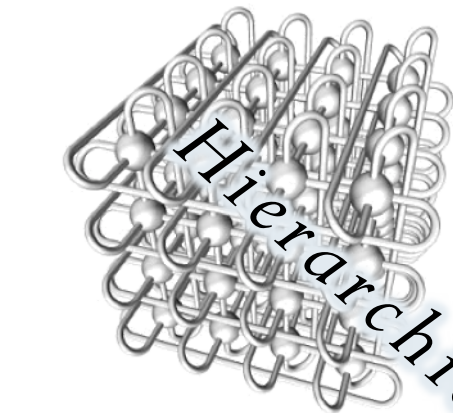


**Sandia  
National  
Laboratories**



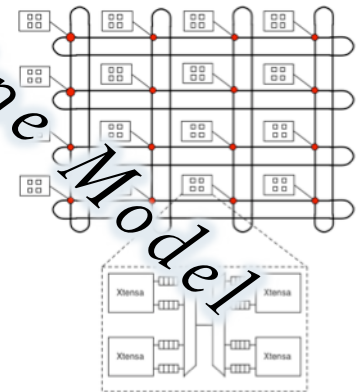
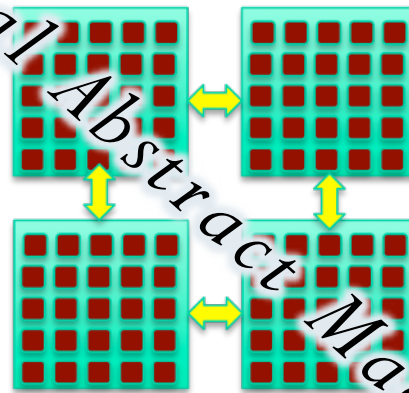
## TOP-DOWN APPROACH

The Top-Down approach aims at modeling Execution Models using accurate *modeling* methodologies. A distinguishing feature of TD is that the models will be based on full applications, as to preserve the workload characteristics of the representative DOE apps selected. The Execution Models will be parameterized for a comprehensive set of attributes, and these parameters will serve as input to the models. Using this approach we will compare novel Execution Model concepts for the full applications optimally mapped onto Exascale architectures being considered by ASCR.



**SST/macro: simulation of different interconnect architectures**

- Driven by traces collected from full application or skeletonized code (either manually or via ROSE)
- Determines degree of congestion on the network and the impact on application performance



SURROGATE	DESCRIPTION
<b>COMPACT APP</b>	Small app. Having fewer features and simplified boundary conditions relative to full app
<b>MINI-APP</b>	Small, self-contained program that embodies essential performance characteristics of keys apps
<b>SKELETON APP</b>	Captures the control flow and communication pattern of an app. Can only be run in a simulator
<b>PROXY APP</b>	General term for all the above
<b>MINI-DRIVER</b>	Small programs that act as drivers of performance-impacting library packages
<b>KERNEL</b>	Captures node-level aspects of an algorithm

Application surrogates for co-design: Use a simplified version of an application for exploration in a restricted manageable design space (inexpensive, fast turnaround), before trying out new approaches in a full application (expensive, multi-year effort).

**RAMP/GreenFlash: Chip-level Simulation**

- Extend GreenFlash/RAMP simulation for more general proxy model
- Create parameterized NoC and memory hierarchy
- Provides model-checking for energy models offered by software simulators

## BOTTOM-UP APPROACH