# **Exploring Emerging Technologies in the Extreme Scale HPC Co-Design Space with** Aspen

Jeffrey S. Vetter

Jeremy Meredith

**MODSIM Workshop 2015** Seattle 13 Aug 2015



MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

ORNL is managed by UT-Battelle for the US Department of Energy





http://ft.ornl.gov + vetter@computer.org





|                              | Speed | Ease | Flexibility | Accuracy | Scalability |           |
|------------------------------|-------|------|-------------|----------|-------------|-----------|
| Ad-hoc Analytical Models     | 1     | 3    | 2           | 4        | 1           |           |
| Structured Analytical Models | 1     | 2    | 1           | 4        | 1           |           |
| Simulation – Functional      | . 3   | 2    | 2           | 3        | 3           |           |
| Simulation – Cycle Accurate  | 4     | 2    | 2           | 2        | 4           |           |
| Hardware Emulation (FPGA)    | 3     | 3    | 3           | 2        | 3           |           |
| Similar hardware measurement | 2     | 1    | 4           | 2        | 2           |           |
| Node Prototype               | 2     | 1    | 4           | 1        | 4           |           |
| Prototype at Scale           | 2     | 1    | 4           | 1        | 2           |           |
| Final System                 | -     | -    | -           | -        | -           |           |
|                              |       |      |             |          |             | CAK RIDGE |

## **Prediction Techniques Ranked**

|                              | Speed | Ease | Flexibility | Accuracy | Scalability |           |
|------------------------------|-------|------|-------------|----------|-------------|-----------|
| Ad-hoc Analytical Models     | 1     | 3    | 2           | 4        | 1           |           |
| Structured Analytical Models | 1     | 2    | 1           | 4        | 1           |           |
| Aspen                        | 1     | 1    | 1           | 4        | 1           |           |
| Simulation – Functional      | 3     | 2    | 2           | 3        | 3           |           |
| Simulation – Cycle Accurate  | 4     | 2    | 2           | 2        | 4           |           |
| Hardware Emulation (FPGA)    | 3     | 3    | 3           | 2        | 3           |           |
| Similar hardware measurement | 2     | 1    | 4           | 2        | 2           |           |
| Node Prototype               | 2     | 1    | 4           | 1        | 4           |           |
| Prototype at Scale           | 2     | 1    | 4           | 1        | 2           |           |
| Final System                 | -     | -    | -           | -        | -           |           |
|                              |       |      |             |          |             | CAK RIDGE |

## **Prediction Techniques Ranked**

## **ASPEN:** Abstract Scalable Performance Engineering Notation



National Laboratory

K. Spafford and J.S. Vetter, "Aspen: A Domain Specific Language for Performance Modeling," in SC12: ACM/IEEE International Conference for High Performance Computing, Networking, Storage, and Analysis, 2012

### **Representation in Aspen**

## **Manual Example of LULESH**

| ဖို br | ranch: master ▼ aspen / models / lulesh / lulesh.aspen 🗄 🔁                                 |  |  |  |
|--------|--|--|--|--|
| 3      | jsmeredith on Sep 20, 2013 adding models   |  |  |  |
| 1 cor  | ntributor  |  |  |  |
|        |  |  |  |  |
| 336    | lines (288 sloc) 9.213 kb 🛛 🖉 🌶 📋  |  |  |  |
| 1      |  |  |  |  |
| 2      | // lulesh.aspen  |  |  |  |
| 3      |  |  |  |  |
| 4      | // an Astew application model for the LOLESS 1.01 challenge problem. Based                 |  |  |  |
| 5      | // bit cle cook version of the source code found at.                                       |  |  |  |
| 7      | // // // // // // // // //////////////   |  |  |  |
| 8      | param nTimeSteps = 1495  |  |  |  |
| 9      |  |  |  |  |
| 10     | // Information about domain  |  |  |  |
| 11     | param edgeElems = 45   |  |  |  |
| 12     | param edgeNodes = edgeElems + 1  |  |  |  |
| 13     |  |  |  |  |
| 14     | param numElems = edgeElems^3   |  |  |  |
| 15     | param numNodes = edgeNodes^3   |  |  |  |
| 16     |  |  |  |  |
| 17     | // Double precision  |  |  |  |
| 18     | param wordSize = 8   |  |  |  |
| 19     |  |  |  |  |
| 20     | // Element data  |  |  |  |
| 21     | data mNodeList as Array(numElems, wordSize)  |  |  |  |
| 22     | data mMatElemList as Array(numElems, wordSize)   |  |  |  |
| 23     | data mNodeList as Array(8 * numElems, wordSize) // 8 nodes per element                     |  |  |  |
| 24     | data mixim as Array(numilems, wordsize)  |  |  |  |
| 25     | data mixip as Array(numilems, WordSize)  |  |  |  |
| 26     | data mietam as Array(numeiems, wordsize)   |  |  |  |
| 27     | data mirtap as Array(numcitems, WORDSIZE)  |  |  |  |
| 28     | data miztam as Array(numiliems, WORDI22)<br>data miztan as Array(numiliems, WORDI22)       |  |  |  |
| 29     | data melamBr as Array(numFlems, Wordsize)  |  |  |  |
| 31     | data meteneo da Array (numeteneo, mondatza)<br>data meteneo da Array (numeteneo, mondatza) |  |  |  |
| 32     | data me us a track/numElems, wordSize)   |  |  |  |

| 147 | Kernel CalcMonotonicQuradients {         |
|-----|--|
| 148 | execute [numElems]                       |
| 149 | {  |
| 150 | loads [8 * indexWordSize] from nodelist  |
| 151 | // Load and cache position and velocity. |
| 152 | loads/caching [8 * wordSize] from x      |
| 153 | loads/caching [8 * wordSize] from y      |
| 154 | loads/caching [8 * wordSize] from z      |
| 155 |  |
| 156 | loads/caching [8 * wordSize] from xvel   |
| 157 | loads/caching [8 * wordSize] from yvel   |
| 158 | loads/caching [8 * wordSize] from zvel   |
| 159 |  |
| 160 | loads [wordSize] from volo               |
| 161 | loads [wordSize] from vnew               |
| 162 | // dx, dy, etc.                          |
| 163 | flops [90] as dp, simd                   |
| 164 | // delvk delxk                           |
| 165 | flops [9 + 8 + 3 + 30 + 5] as dp, simd   |
| 166 | stores [wordSize] to delv_xeta           |
| 167 | // delxi delvi                           |
| 168 | flops [9 + 8 + 3 + 30 + 5] as dp, simd   |
| 169 | stores [wordSize] to delx_xi             |
| 170 | // delxj and delvj                       |
| 171 | flops [9 + 8 + 3 + 30 + 5] as dp, simd   |
| 172 | stores [wordSize] to delv_eta            |
| 173 | }  |
| 174 | }  |



g

## **Example Uses: Resource Exploration**

| Benchmark     | Runtime Order            |   |
|---------------|--------------------------|---|
| BACKPROP      | H * O + H * I            |   |
| BFS           | nodes + edges            |   |
| CFD           | nelr*ndim                |   |
| CG            | nrow + ncol              |   |
| HOTSPOT       | $sim_time * rows * cols$ |   |
| JACOBI        | $m\_size * m\_size$      |   |
| <b>KMEANS</b> | nAttr*nClusters          |   |
| LAPLACE2D     | $n^2$                    |   |
| LUD           | $matrix_dim^3$           |   |
| MATMUL        | N * M * P                | C |
| NW            | $max\_cols^2$            |   |
| SPMUL         | size + nonzero           |   |
| SRAD          | niter*rows*cols          |   |





Fig. 8: GPU Memory Usage of each Function in LULESH, where the memory usage of a function is inclusive; value for a parent function includes data accessed by its child functions in the call graph.

Table 2: Order analysis, showing Big O runtime for each benchmark in terms of its key parameters.

| Method Name                      | FLOPS/byte |
|----------------------------------|------------|
| InitStressTermsForElems          | 0.03       |
| CalcElemShapeFunctionDerivatives | 0.44       |
| SumElemFaceNormal                | 0.50       |
| CalcElemNodeNormals              | 0.15       |
| SumElemStressesToNodeForces      | 0.06       |
| IntegrateStressForElems          | 0.15       |
| CollectDomainNodesToElemNodes    | 0.00       |
| VoluDer                          | 1.50       |
| CalcElemVolumeDerivative         | 0.33       |
| CalcElemFBHourglassForce         | 0.15       |
| CalcFBHourglassForceForElems     | 0.17       |
| CalcHourglassControlForElems     | 0.19       |
| CalcVolumeForceForElems          | 0.18       |
| CalcForceForNodes                | 0.18       |
| CalcAccelerationForNodes         | 0.04       |
| ApplyAccelerationBoundaryCond    | 0.00       |
| CalcVelocityForNodes             | 0.13       |
| CalcPositionForNodes             | 0.13       |
| LagrangeNodal                    | 0.18       |
| AreaFace                         | 10.25      |
| CalcElemCharacteristicLength     | 0.44       |
| CalcElemVelocityGrandient        | 0.13       |
| CalcKinematicsForElems           | 0.24       |
| CalcLagrangeElements             | 0.24       |
| CalcMonotonicOGradientsForElems  | 0.46       |



Fig. 7: Measured and predicted runtime of the entire LULESH program on CPU and GPU, including measured runtimes using the automatically predicted optimal target device at each size.



Figure 1: A plot of idealized concurrency by chronological phase in the digital spotlighting application model.

## **Aspen allows Multiresolution Modeling**



## **Node Scale Modeling with COMPASS**



# **COMPASS System Overview**

## Detailed Workflow of the COMPASS Modeling Framework

Optional feedback for advanced users



S. Lee, J.S. Meredith, and J.S. Vetter, "COMPASS: A Framework for Automated Performance Modeling and Prediction," in ACM International Conference on Supercomputing (ICS). Newport Beach, California: ACM, 2015, 10.1145/2751205.2751220.



## **MM example generated from COMPASS**

int N = 1024;1 void matmul(float \*a, float \*b, float \*c){ int i, j, k ;  $\mathbf{2}$ 3 #pragma acc kernels loop gang copyout(a[0:(N\*N)]) \ copyin(b[0:(N\*N)],c[0:(N\*N)]) 4 for (i=0; i<N; i++) $\mathbf{5}$ #pragma acc loop worker 6 for (j=0; j<N; j++) { float sum = 0.0;  $\overline{7}$ for (k=0; k<N; k++) {sum+=b[i\*N+k]\*c[k\*N+j];} 8  $a[i*N+j] = sum; \}$ 9 10} //end of i loop } //end of matmul() 1112int main() { 13int i; float \*A = (float\*) malloc(N\*N\*sizeof(float));float \*B = (float\*) malloc(N\*N\*sizeof(float));14float \*C = (float\*) malloc(N\*N\*sizeof(float));15for (i = 0; i < N\*N; i++)16 $\{ A[i] = 0.0F; B[i] = (float) i; C[i] = 1.0F; \}$ 17#pragma aspen modelregion label(MM) 18matmul(A,B,C); 19free(A); free(B); free(C); return 0; 20} //end of main() 21

| 1  | model MM {   |
|----|--|
| 2  | param floatS = 4; param N = $1024$   |
| 3  | data A as Array((N*N), floatS)   |
| 4  | data B as Array((N*N), floatS)   |
| 5  | data C as Array((N*N), floatS)   |
| 6  | kernel matmul {  |
| 7  | execute matmul2_intracommIN  |
| 8  | $\{ \text{ intracomm [floatS}*(N*N) \} \text{ to } C \text{ as copyin} \}$ |
| 9  | intracomm [floatS*(N*N)] to B as copyin $\}$                               |
| 10 | map matmul2 [N] {  |
| 11 | map matmul3 [N] {  |
| 12 | iterate [N] {  |
| 13 | execute matmul5  |
| 14 | $\{ \text{ loads [floatS] from B as stride}(1) \}$                         |
| 15 | loads [floatS] from C; flops [2] as sp, simd }                             |
| 16 | } //end of iterate   |
| 17 | execute matmul6 { stores [floatS] to A as $stride(1)$ }                    |
| 18 | } // end of map matmul3  |
| 19 | } //end of map matmul2   |
| 20 | execute matmul2_intracommOUT   |
| 21 | $\{ \text{ intracomm [floatS*(N*N)] to A as copyout } \}$                  |
| 22 | } //end of kernel matmul   |
| 23 | kernel main { $matmul()$ }   |
| 24 | } //end of model MM  |



## **Annotation Overhead**

| Benchmark Name | Lines of Code | Lines of Annotation | Annotation Overhead<br>(%) |
|----------------|---------------|---------------------|----------------------------|
| JACOBI         | 241           | 2                   | 0.8                        |
| MATMUL         | 128           | 1                   | 0.7                        |
| SPMUL          | 423           | 10                  | 2.3                        |
| LAPLACE2D      | 210           | 7                   | 3.3                        |
| CG             | 1511          | 10                  | 0.6                        |
| EP             | 759           | 9                   | 1.1                        |
| BACKPROP       | 1074          | 4                   | 0.3                        |
| BFS            | 435           | 16                  | 3.6                        |
| CFD            | 752           | 9                   | 1.1                        |
| HOTSPOT        | 525           | 11                  | 2.0                        |
| KMEANS         | 1822          | 11                  | 0.6                        |
| LUD            | 421           | 6                   | 1.4                        |
| NW             | 478           | 8                   | 1.7                        |
| SRAD           | 550           | 12                  | 2.1                        |
| LULESH         | 3743          | 125                 | 3.3                        |

National Laboratory

## **Example: LULESH (10% of 1 kernel)**

kernel IntegrateStressForElems { execute [numElem\_CalcVolumeForceForElems] { loads [((1\*aspen\_param\_int)\*8)] from elemNodes as stride(1) loads [((1\*aspen\_param\_double)\*8)] from m\_x loads [((1\*aspen\_param\_double)\*8)] from m\_z loads [((1\*aspen\_param\_double)\*8)] from determ as stride(1) flops [8] as dp, simd flops [8] as dp, simd flops [8] as dp, simd flops [3] as dp, simd flops [2] as dp, simd stores [(1\*aspen\_param\_double)] as stride(0) flops [2] as dp, simd stores [(1\*aspen\_param\_double)] as stride(0) flops [2] as dp, simd stores [(1\*aspen\_param\_double)] as stride(0) flops [2] as dp, simd stores [(1\*aspen\_param\_double)] as stride(0) flops [2] as dp, simd stores [(1\*aspen\_param\_double)] as stride(0) flops [2] as dp, simd loads [(1\*aspen\_param\_double)] as stride(0) stores [(1\*aspen\_param\_double)] as stride(0) loads [(1\*aspen\_param\_double)]

- Input LULESH program: 3700 lines of C codes
- Output Aspen model: 2300 lines of Aspen codes



## **Model Validation**

|           | FLOPS | LOADS | STORES |
|-----------|-------|-------|--------|
| MATMUL    | 15%   | <1%   | 1%     |
| LAPLACE2D | 7%    | 0%    | <1%    |
| SRAD      | 17%   | 0%    | 0%     |
| JACOBI    | 6%    | <1%   | <1%    |
| KMEANS    | 0%    | 0%    | 8%     |
| LUD       | 5%    | 0%    | 2%     |
| BFS       | <1%   | 11%   | 0%     |
| НОТЅРОТ   | 0%    | 0%    | 0%     |
| LULESH    | 0%    | 0%    | 0%     |

0% means that prediction fell between measurements from optimized and unoptimized runs of the code.



## **Model Scaling Validation (LULESH)**



## **Performance Modeling for Distributed Scientific Workflows**

# (see our Panorama poster)







## **End-to-end Resiliency Design using Aspen**



## **Resiliency Modeling with Aspen**

- End-to-End system design for Extreme-scale HPC
  - Why pay redundant costs for power, performance, etc?
- We introduce a new metric, the data vulnerability factor (DVF)
  - Quantifying vulnerability of data structures
  - Avoiding the isolation between application and hardware
- We measure DVF based on Aspen, a domain specific language for system modeling
- We categorize memory access patterns of scientific applications from a spectrum of computational domains
  - Dense linear algebra, Sparse linear algebra, N-body method, Structured grids, Spectral methods, and Monte Carlo
- We demonstrate the significance of DVF by two case studies
  - Algorithm optimization
  - Data protection quantification

## End-To-End Arguments in System Design

J. H. SALTZER, D. P. REED, and D. D. CLARK Massachusetts Institute of Technology Laboratory for Computer Science

This paper presents a design principle that helps guide placement of functions among the modules of a distributed computer system. The principle, called the end-to-end argument, suggests that functions placed at low levels of a system may be redundant or of little value when compared with the cost of providing them at that low level. Examples discussed in the paper include bit-error recovery, security using encryption, duplicate message suppression, recovery from system crashes, and delivery acknowl-edgment. Low-level mechanisms to support these functions are justified only as performance enhancements.

CR Categories and Subject Descriptors: C.0 [General] Computer System Organization—system architectures; C.2.2 [Computer-Communication Networks]: Network Protocols—protocol architecture; C.2.4 [Computer-Communication Networks]: Distributed Systems; D.4.7 [Operating Systems]: Organization and Design—distributed systems

General Terms: Design

Additional Key Words and Phrases: Data communication, protocol design, design principles

### 1. INTRODUCTION

Choosing the proper boundaries between functions is perhaps the primary activity of the computer system designer. Design principles that provide guidance in this choice of function placement are among the most important tools of a system

| $DVF_d$     | DVF for a specific data structure                  |
|-------------|--|
| FIT         | Failure rate (i.e., failures per billion hours per |
|             | Mbit)  |
| T           | Application execution time                         |
| $S_d$       | Size of data structure                             |
| $N_{error}$ | Number of errors that could occur to a specific    |
|             | data structure during application execution        |
| $N_{ha}$    | Number of accesses to hardware (the main           |
|             | memory in this work)                               |
| n           | Number of major data structures in an appli-       |
|             | cation   |
| $DVF_a$     | DVF for the application                            |



L. Yu, D. Li et al., "Quantitatively modeling application resilience with the data vulnerability factor (Best Student Paper Finalist)," in SC14: International Conference for High Performance Computing, Networking, Storage and Analysis. New Orleans, Louisiana: IEEE Press, 2014, pp. 695-706, 10.1109/sc.2014.62.

## **Data Vulnerability Factor:** Why a new metric and methodology?

- Analytical model of resiliency that includes important features of architecture and application
  - Fast
  - Flexible
- Balance multiple design dimensions
  - Application requirements
  - Architecture (memory capacity and type)
- Focus on main memory initially
- Prioritize vulnerabilities of application data

L. Yu, D. Li et al., "Quantitatively modeling application resilience with the data vulnerability factor (Best Student Paper Finalist), in **CAK RIDGE** SC14: International Conference for High Performance Computing, Networking, Storage and Analysis. New Orleans, Louisiana: IEEE Press, 2014, pp. 695-706, 10.1109/sc.2014.62.

### End-To-End Arguments in System Design

J. H. SALTZER, D. P. REED, and D. D. CLARK Massachusetts Institute of Technology Laboratory for Computer Science

This paper presents a design principle that helps guide placement of functions among the modules of a distributed computer system. The principle, called the end-to-end argument, suggests that functions placed at low levels of a system may be redundant or of little value when compared with the cost of providing them at that low level. Examples discussed in the paper include bit-error recovery, security using encryption, duplicate message suppression, recovery from system crashes, and delivery acknowledgment. Low-level mechanisms to support these functions are justified only as performance enhancements.

CR Categories and Subject Descriptors: C.0 [General] Computer System Organization-system architectures; C.2.2 [Computer-Communication Networks]: Network Protocols-protocol architecture; C.2.4 [Computer-Communication Networks]: Distributed Systems; D.4.7 [Operating Systems]: Organization and Design-distributed systems

General Terms: Design

Additional Key Words and Phrases: Data communication, protocol design, design principles

#### 1. INTRODUCTION

Choosing the proper boundaries between functions is perhaps the primary activity of the computer system designer. Design principles that provide guidance in this choice of function placement are among the most important tools of a system

| $DVF_d$     | DVF for a specific data structure                  |
|-------------|--|
| FIT         | Failure rate (i.e., failures per billion hours per |
|             | Mbit)  |
| T           | Application execution time                         |
| $S_d$       | Size of data structure                             |
| $N_{error}$ | Number of errors that could occur to a specific    |
|             | data structure during application execution        |
| $N_{ha}$    | Number of accesses to hardware (the main           |
|             | memory in this work)                               |
| n           | Number of major data structures in an appli-       |
|             | cation   |
| $DVF_a$     | DVF for the application                            |

National Laboratory

## **Workflow to calculate Data Vulnerability Factor**



Fig. 3. The workflow to calculate DVF.



## **An Example of Aspen Program for DVF**



**DVF Results** 



### Provides insight for balancing interacting factors

## **DVF: next steps**

- Evaluated different architectures
  - How much no-ECC, ECC, NVM?
- Evaluate software and applications
  - ABFT
  - C/R
  - TMR
  - Containment domains
  - Fault tolerant MPI

- End-to-End analysis
  - Where should we bear the cost for resiliency?
    - Not everwhere!



## Summary

- Our community has major challenges in HPC as we move to extreme scale
  - Power, Performance, Resilience, Productivity
  - New technologies emerging to address some of these challenges
    - Heterogeneous computing
    - Nonvolatile memory
  - Not just HPC: Most uncertainty in at least two decades
- We need performance prediction and engineering tools now more than ever!
- Aspen is a tool for structured design and analysis
  - Co-design applications and architectures for performance, power, resiliency

OAK RII

- Automatic model generation
- Scalable to distributed scientific workflows
- DVF a new twist on resiliency modeling

## **Acknowledgements**

- · Contributors and Sponsors
  - Future Technologies Group: http://ft.ornl.gov
  - US Department of Energy Office of Science
    - DOE Vancouver Project: <u>https://ft.ornl.gov/trac/vancouver</u>
    - DOE Blackcomb Project: <u>https://ft.ornl.gov/trac/blackcomb</u>
    - DOE ExMatEx Codesign Center: <a href="http://codesign.lanl.gov">http://codesign.lanl.gov</a>
    - DOE Cesar Codesign Center: <u>http://cesar.mcs.anl.gov/</u>
    - DOE Exascale Efforts: <u>http://science.energy.gov/ascr/research/computer-science/</u>
    - Scalable Heterogeneous Computing Benchmark team: <u>http://bit.ly/shocmarx</u>
  - US National Science Foundation Keeneland Project: <u>http://keeneland.gatech.edu</u>
  - US DARPA
  - NVIDIA CUDA Center of Excellence



