# Exascale Co-Design Center for Materials in Extreme Environments
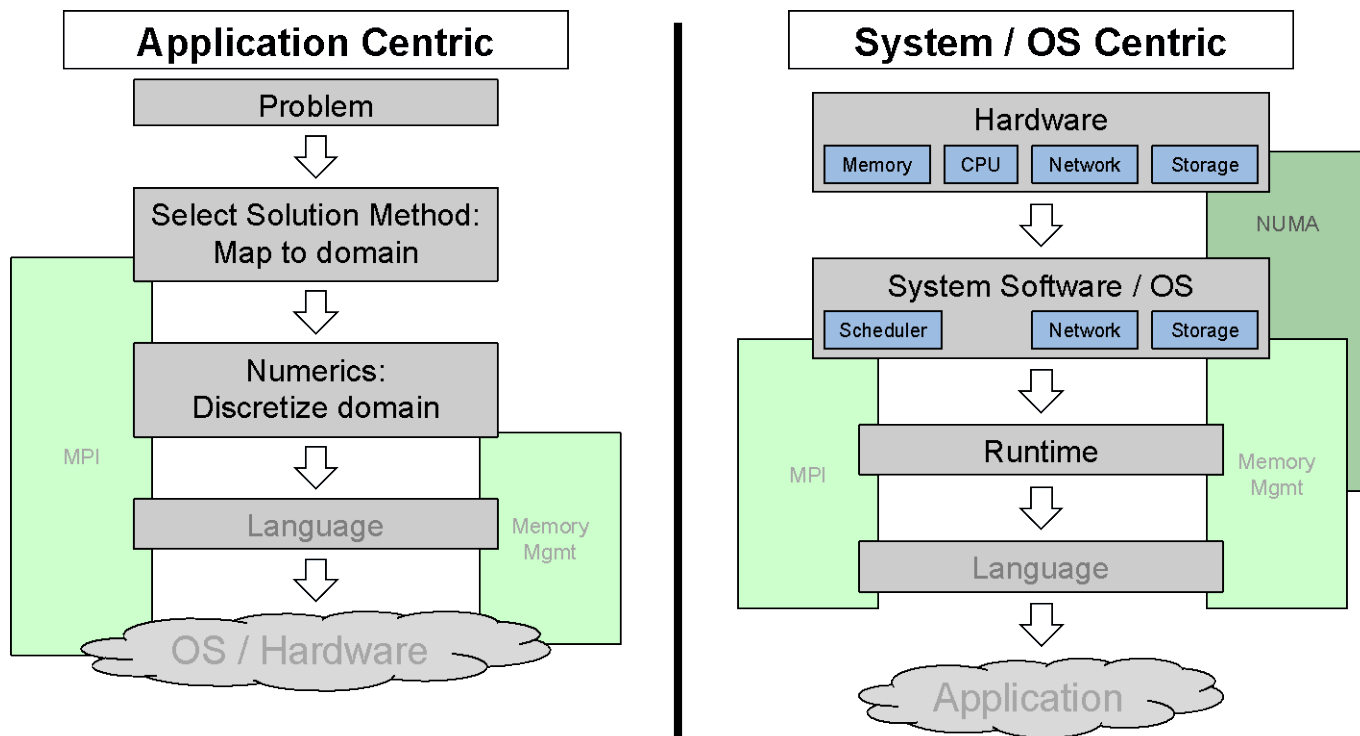


# Timothy C. Germann

# 23 August 2011

# Current science application development strategies are difficult to sustain

An air gap has been encouraged between application developers & system / OS developers and the hardware
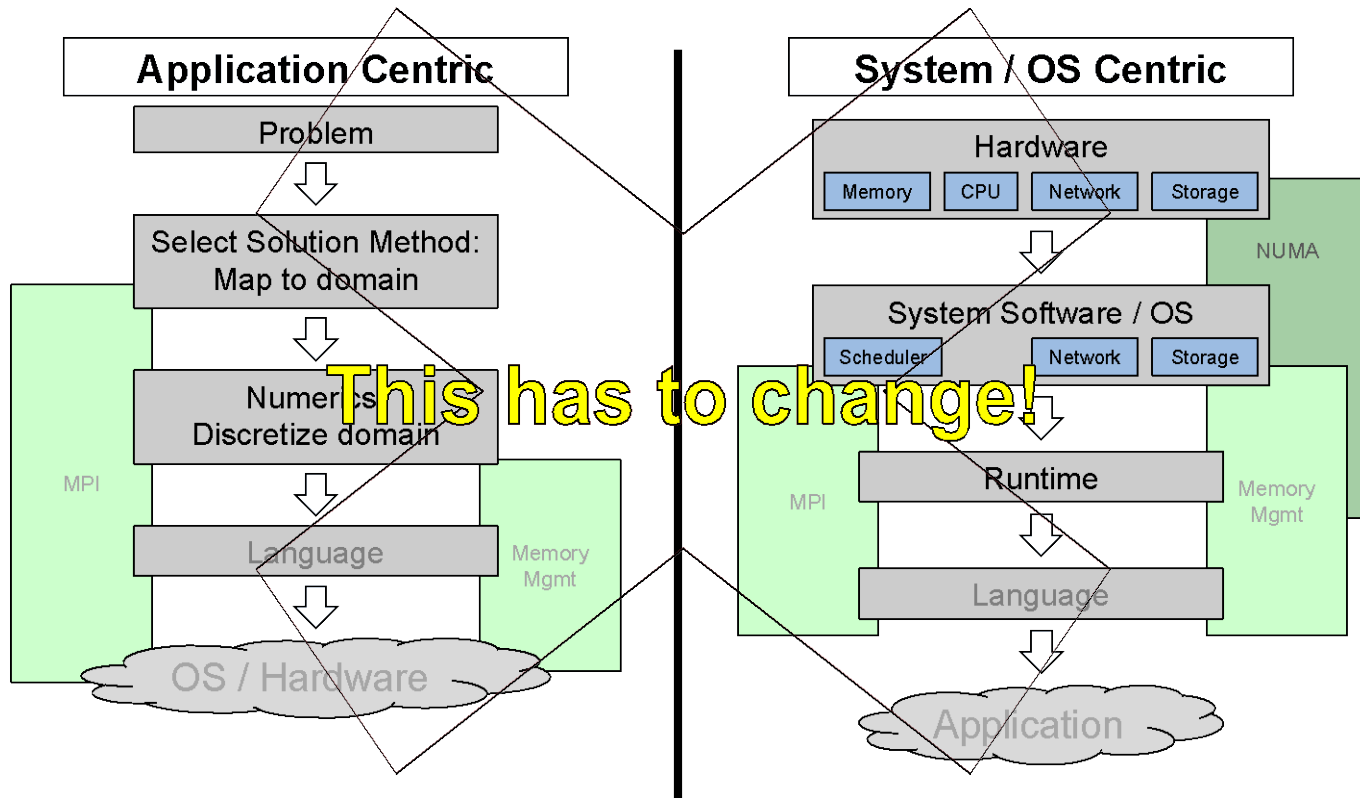
# Current science application development strategies are difficult to sustain

An air gap has been encouraged between application developers & system / OS developers and the hardware

# Co-design is a process by which computer science, applied math, and domain science experts work together to enable scientific discovery
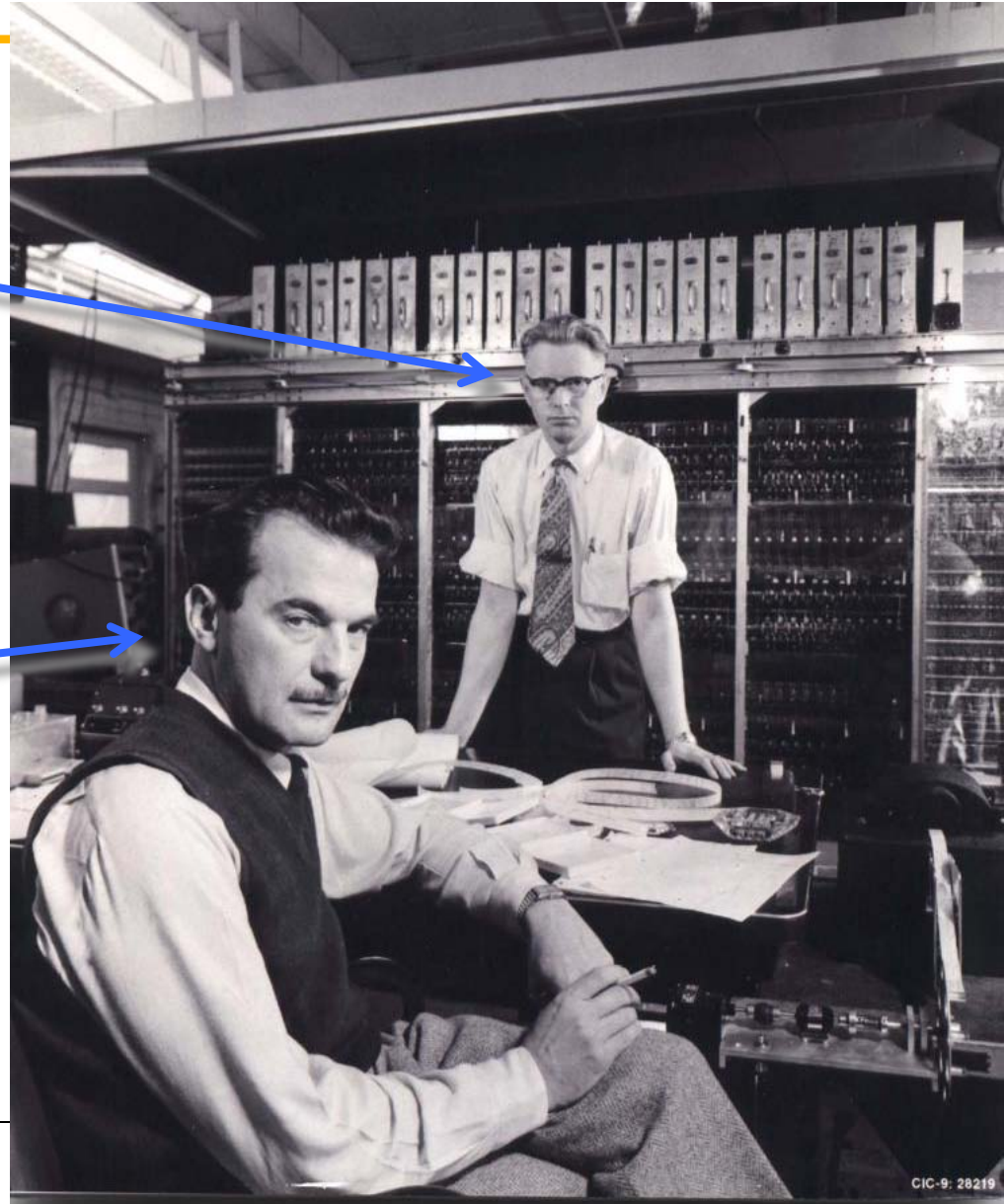
- Hardware is changing dramatically
  - *Increased concurrency*
  - *Increased heterogeneity*
  - *Reduced memory per core*
  - *"Business as usual" is not going to work*

- Algorithms and methods will have to be rethought / revisited
  - *Flops are (almost always) free*
  - *Memory is at a premium*
  - *Power is a constraint for large scale systems*
  - *Resiliency is a challenge*

- Few domain scientists have the extended expertise "from hardware to application" to enable applications to run at exascale

- Success on the next generation of machines will require extensive collaboration between domain scientists, applied mathematicians, computer scientists, and hardware manufacturers

**Los Alamos**
NATIONAL LABORATORY
EST.1943

**NNSA**

# Los Alamos computational co-design, circa 1950

**Hardware architect
(Richardson)**

**Application scientist
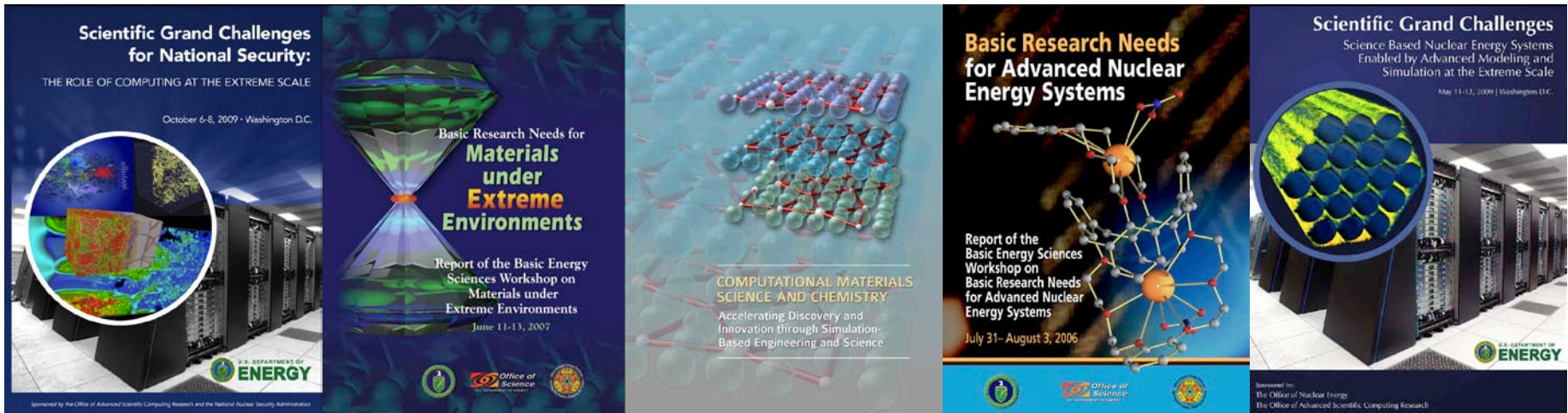(Metropolis)**

H. L. Anderson,
"Metropolis, Monte Carlo, and the MANIAC,"
*Los Alamos Science* **14**, 96-107 (1986).

CIC-9: 28219

# Los Alamos computational co-design, circa 2008

- Roadrunner was a leap into the future
  - *First computer to reach a petaflop*
  - *First heterogeneous supercomputer*
  - *First accelerated supercomputer*
  - *Demonstrated that accelerated supercomputing was possible*
  - *96% of compute power concentrated in accelerators*
- Success required domain scientists, applied mathematicians, and computer scientists working together to identify the correct abstractions for domain science, applied mathematics, programming models, and hardware
- Many successful applications, including
  - Large Scale MD
  - Long time MD
  - Roadrunner Universe
  - DNS of turbulence
  - VPIC laser backscatter
  - VPIC magnetic reconnection
  - Supernova simulations
  - HIV phylogenetics

# A predictive understanding of the response of materials to extreme conditions (mechanical and/or irradiation) underpins many DOE missions.

# Traditional computational materials science: a hierarchy of length/time scales

| Ab-initio Methods | Molecular Dynamics | Phase-Field Modeling | Continuum Methods |
|---|---|---|---|
| **Inter-atomic force model, equation of state** | **Defect and interface mobility, nucleation** | **Direct numerical simulation of multi-phase evolution** | **Multi-phase material response, experimental observables** |
|  |  |  |  |
| Length/time: nm, ps | Length/time: µm, ns | Length/time: 100 µm, µs | Length/time: cm, ms |
| Codes: Qbox/LATTE | Codes: SPaSM/ddcMD | Codes: AMPE/GL | Codes: VP-FFT/ALE3d |
| Motif: Particles and wavefunctions, plane wave DFT with nonlocal norm-conserving, ScaLAPACK, BLACS, and custom parallel 3D FFTs | Motif: Particles, domain decomposition, explicit time integration, neighbor and linked lists, dynamic load balancing, parity error recovery, and *in situ* visualization | Motif: Regular and adaptive grids, implicit time integration, real-space and spectral methods, complex order parameter (phase, crystal, species) | Motif: Regular and irregular grids, implicit time integration, 3D FFTs, polycrystal and single crystal plasticity, |
| Prog. Model: MPI | Prog. Model: MPI + Threads | Prog. Model: MPI | Prog. Model: MPI |

# Current trends will increase the *length*, but not *time*, scales accessible by molecular dynamics simulation



| System attributes | 2010 | "2015" | "2018" |
|---|---|---|---|
| System peak | 2 Peta | 200 Peta | 1 Exa |
| Power | 6 MW | ~15 MW | ~20 MW |
| System memory | 0.3 PB | 5 PB | 32-64 PB |
| Node performance | 125 GF | 0.5 TF or 7 TF | 1 TF or 10x |
| Node memory BW | 25 GB/s | 0.1 TB/s or 10x | 0.4 TB/s or 10x |
| Node concurrency | 12 | O(100) | O(1k) or 10x |
| Total Node Interconnect BW | 1.5 GB/s | 20 GB/s or 10x | 200 GB/s or 10x |
| System size (nodes) | 18,700 | 50,000 or 1/10x | O(100,000) or 1/10 x |
| MTTI | days | O(1day) | O(1 day) |

Source: DOE Exascale Initiative Technical Roadmap

**Clock speeds and bandwidths will not increase substantially, so the *timescale* challenge is going to become increasingly critical.**

# Preparing for exascale: issues to confront

- Computer architectures are becoming increasingly **heterogeneous** and **hierarchical**, with greatly increased flop/byte ratios.
- The algorithms, programming models, and tools that will thrive in this environment must mirror these characteristics.
- SPMD bulk synchronous ($10^9$-way) parallelism will no longer be viable.
- Power, energy, and heat dissipation are increasingly important.
- Traditional global checkpoint/restart is becoming impractical.
  - *Local flash memory?*
- Fault tolerance and resilience
  - *Recovering from soft and hard errors, and anticipating faults*
  - *MPI/application ability to drop or replace nodes*
  - *The curse of silent errors*
- Analysis and visualization
  - *In situ, e.g. "active storage" using I/O nodes?*

Los Alamos
NATIONAL LABORATORY
EST.1943

NNSA

# ExMatEx Co-Design Project Goals

- Our **goal** is to establish the interrelationship between hardware, middleware (software stack), programming models, and algorithms required to enable *a productive exascale environment* for multiphysics simulations of materials in extreme mechanical and radiation environments.

- We will exploit, rather than avoid, the greatly increased levels of concurrency, heterogeneity, and flop/byte ratios on the upcoming exascale platforms.



  - This *task-based* approach leverages the extensive concurrency and heterogeneity expected at exascale while enabling fault tolerance within applications.

  - The programming models and approaches developed to achieve this will be broadly applicable to a variety of multiscale, multiphysics applications, including astrophysics, climate and weather prediction, structural engineering, plasma physics, and radiation hydrodynamics.

# ExMatEx Co-Design Project Objectives

- **Inter-communication of requirements and capabilities between the materials science and the exascale hardware and software communities**
  - Proxy apps communicate the application workload to the hardware architects and system software developers, and are used in models/simulators/emulators to assess performance, power, and resiliency.
  - Exascale capabilities and limitations will be continuously incorporated into the proxy applications through an agile development loop.
  - Single-scale SPMD proxy apps (e.g. molecular dynamics) will be used to assess node-level data structures, performance, memory and power management strategies.
  - System-level data movement, fault management, and load balancing techniques will be evaluated via the asynchronous task-based MPMD scale-bridging proxy apps.

- **Perform trade-off analysis between competing requirements and capabilities in a tightly coupled optimization loop**
  - A three-pronged approach combining:
    - Node- to system-level models and simulators
    - Exascale emulation layer (GREMLIN) to introduce perturbations similar to those expected on future architectures
    - Performance analysis on leadership-class machines
  - Co-optimization of algorithms and architectures for price, performance, power (chiefly memory and data movement), and resilience (P$^3$R)

**Los Alamos**
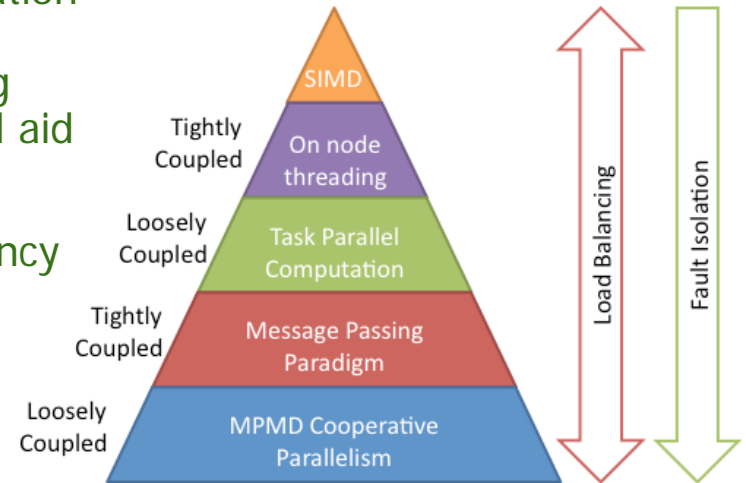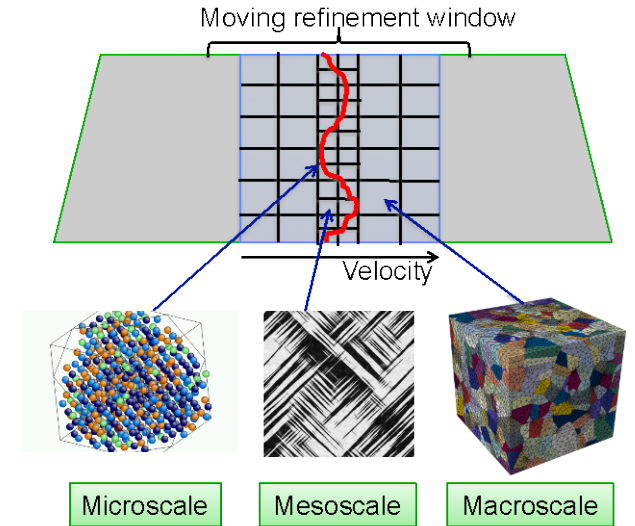NATIONAL LABORATORY
EST.1943

NNSA

# ExMatEx Co-Design Project Objectives

- **Full utilization of exascale concurrency and locality**

  - Heterogeneous, hierarchical MPMD algorithms map naturally to anticipated heterogeneous, hierarchical architectures.

  - Escape the traditional bulk synchronous SPMD paradigm, improve data locality and reduce I/O burden.

- **Application friendly programming models**

  - Must expose hardware capabilities to the application programmer while at the same time hiding the continuous flux and complexity of the underlying hardware through a layer of abstraction that will aid portability.

  - Task-based MPMD approach leverages concurrency and heterogeneity at exascale while enabling novel data models, power management, and fault tolerance strategies.



Moving refinement window

Velocity

Microscale | Mesoscale | Macroscale



SIMD

Tightly Coupled — On node threading

Loosely Coupled — Task Parallel Computation

Tightly Coupled — Message Passing Paradigm

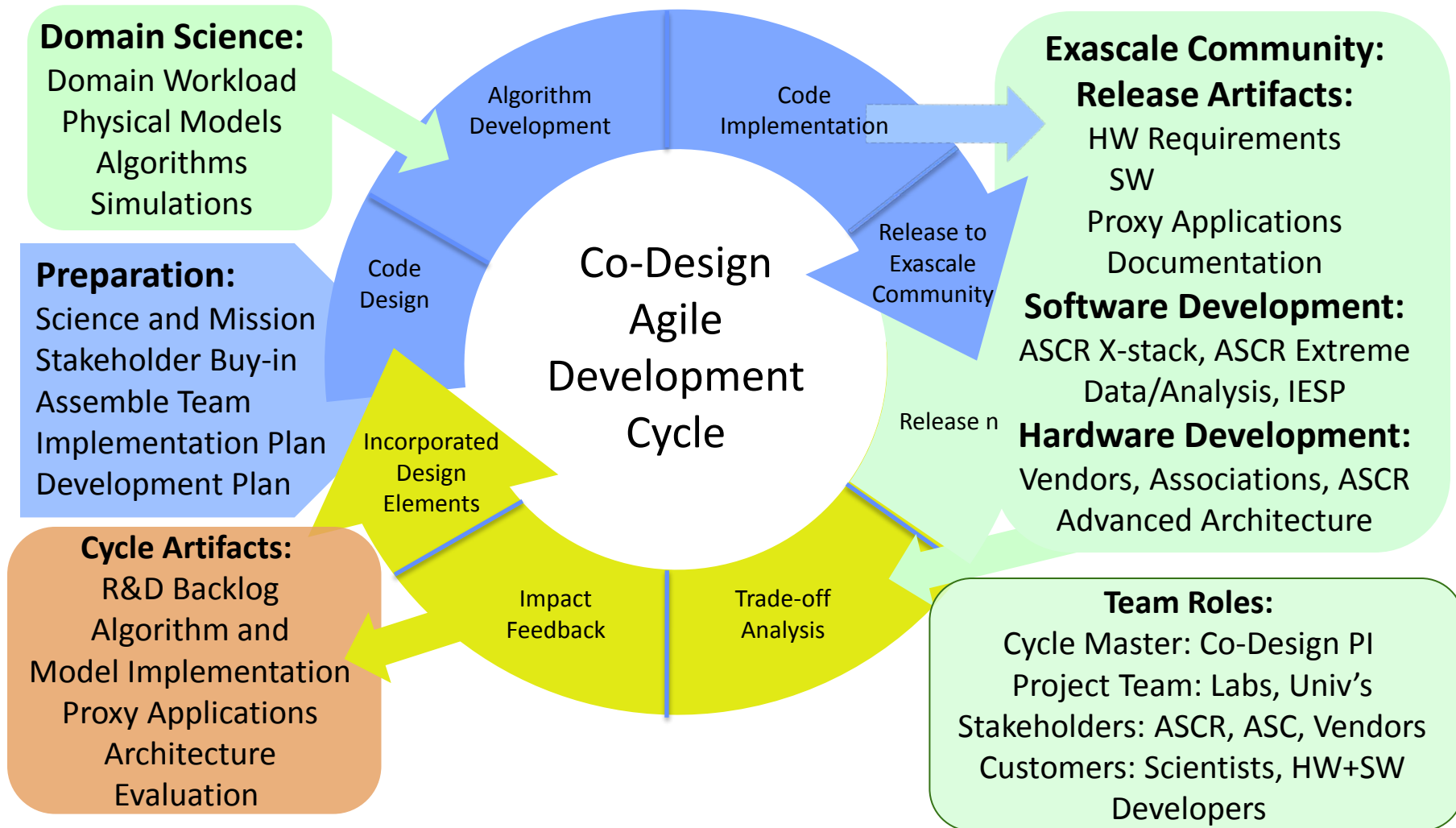Loosely Coupled — MPMD Cooperative Parallelism

Load Balancing

Fault Isolation

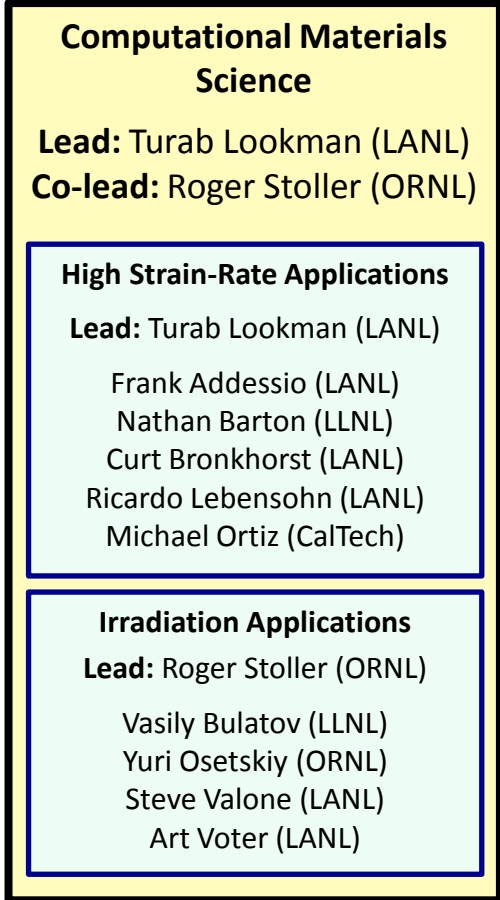# We will manage by adaptive rather than predictive planning.

- Agile development is an adaptive cycle in which

  - *Initial requirements are gathered from the hardware, software, and domain application communities (e.g. Gordon Bell Prize-winning applications).*

  - *Application requirements for hardware and software are continuously released to the exascale community in the form of proxy applications and documentation (release artifacts).*

  - *Application, software, and hardware communities analyze and respond to trade-offs with new requirements and capabilities, both from and to the application.*

  - *Changes in hardware and software designs are rapidly adapted into proxy applications (cycle artifacts).*

  - *Repeated iterations converge to the optimal design for the exascale simulation environment for real science applications.*

**Co-Design Requires Adaptive Methodologies.**

# Management Plan

**Domain Science:**
Domain Workload
Physical Models
Algorithms
Simulations

**Preparation:**
Science and Mission
Stakeholder Buy-in
Assemble Team
Implementation Plan
Development Plan

**Cycle Artifacts:**
R&D Backlog
Algorithm and
Model Implementation
Proxy Applications
Architecture
Evaluation

Algorithm
Development

Code
Implementation

Code
Design

Release to
Exascale
Community

**Co-Design Agile Development Cycle**

Incorporated
Design
Elements

Impact
Feedback

Trade-off
Analysis

Release n

**Exascale Community:**
**Release Artifacts:**
HW Requirements
SW
Proxy Applications
Documentation
**Software Development:**
ASCR X-stack, ASCR Extreme
Data/Analysis, IESP
**Hardware Development:**
Vendors, Associations, ASCR
Advanced Architecture

**Team Roles:**
Cycle Master: Co-Design PI
Project Team: Labs, Univ's
Stakeholders: ASCR, ASC, Vendors
Customers: Scientists, HW+SW
Developers

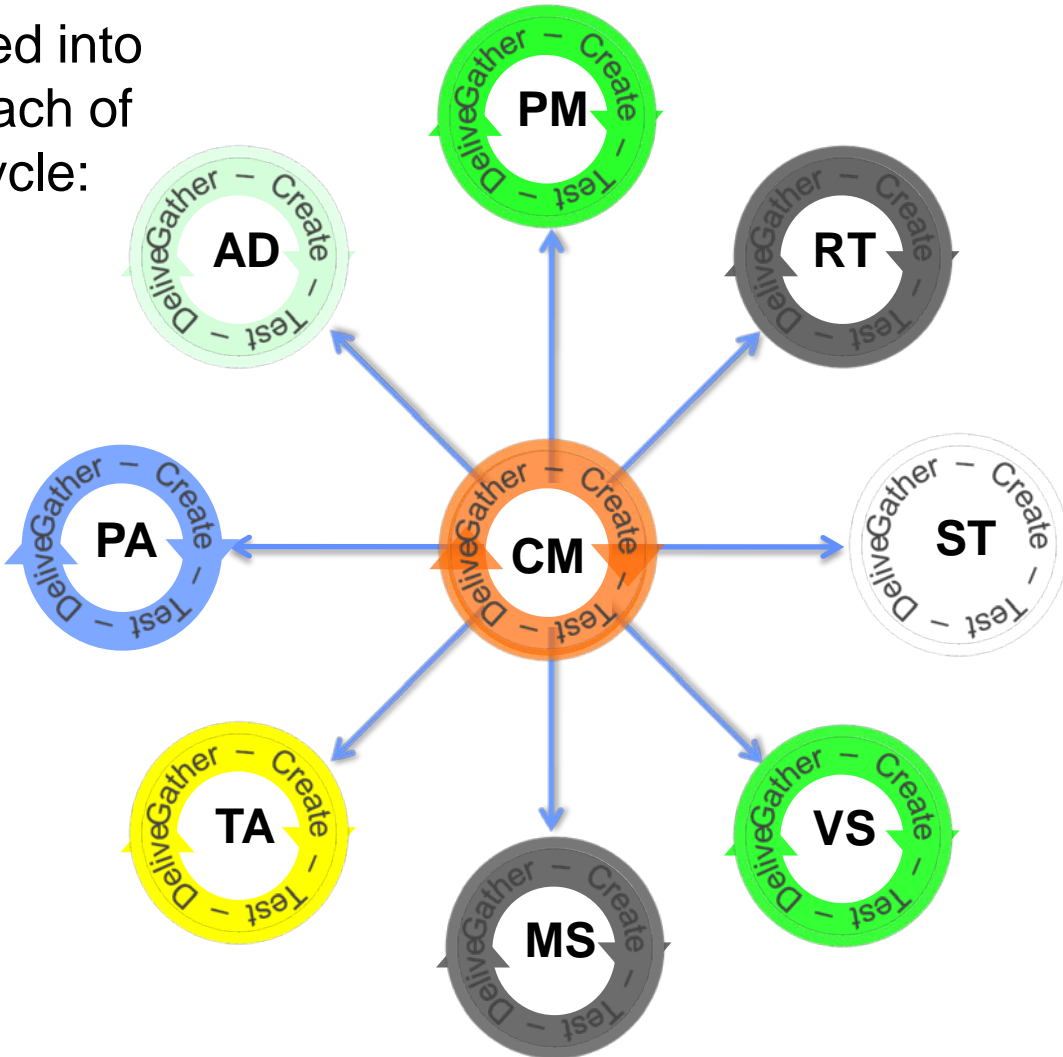**To successfully define this exascale simulation environment, our co-design process must be *adaptive, iterative,* and *lightweight* – i.e. <u>agile</u>.**

**Executive Advisory Board**

Alan Bishop (LANL)
Tomás Díaz de la Rubia (LLNL)
Rick Stevens (ANL)
Kathy Yelick (LBNL)
Steve Zinkle (ORNL)

**Exascale Co-Design Center for Materials in Extreme Environments**

**Center Director:** Tim Germann (LANL)
**Deputy Director:** Jim Belak (LLNL)

**SC/ASCR**

**Exascale Co-Design Consortium**

**Advanced Algorithms & Co-design "Code-Team"**
**Lead:** David Richards (LLNL)   Erik Draeger (LLNL), Tim Kelley (LANL), Bryan Lally (LANL), Danny Perez (LANL)

## Computer Science

**Lead:** Sriram Swaminarayan (LANL)
**Co-lead:** Scott Futral (LLNL)

**Programming Models**
**Lead:** Allen McPherson (LANL)
Pat Hanrahan (Stanford)
David Jefferson (LLNL)

**Data/Resource Sharing**
**Lead:** Jim Ahrens (LANL)

**Analysis Tools At Scale**
**Lead:** Martin Schulz (LLNL)

**Performance Modeling**
**Lead:** Jeff Vetter (ORNL)
Jim Ang (SNL)
Arun Rodrigues (SNL)

**Software Stack Engagement**
Jim Ahrens (LANL)
Martin Schulz (LLNL)

**Vendor Engagement**
Matt Leininger (LLNL)
Pat McCormick (LANL)

## Applied Math

**Lead:** Milo Dorr (LLNL)
**Co-lead:** Dana Knoll (LANL)

**Scale-Bridging Algorithms**

**Lead:** Dana Knoll (LANL)

Frank Alexander (LANL)
Milo Dorr (LLNL)
Jean-Luc Fattebert (LLNL)
Ed Kober (LANL)

**V&V+UQ**

**Lead:** Houman Owhadi (CalTech)

Richard Klein (LLNL)
Earl Lawrence (LANL)
Clint Scovel (LANL)

## Computational Materials Science

**Lead:** Turab Lookman (LANL)
**Co-lead:** Roger Stoller (ORNL)

**High Strain-Rate Applications**

**Lead:** Turab Lookman (LANL)

Frank Addessio (LANL)
Nathan Barton (LLNL)
Curt Bronkhorst (LANL)
Ricardo Lebensohn (LANL)
Michael Ortiz (CalTech)

**Irradiation Applications**

**Lead:** Roger Stoller (ORNL)

Vasily Bulatov (LLNL)
Yuri Osetskiy (ORNL)
Steve Valone (LANL)
Art Voter (LANL)

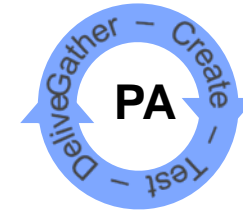# Management Plan: Interconnected Task Areas

- Milestones have been organized into 9 interconnected task areas, each of which operates an agile sub-cycle:

  – **CM:** Center management
  – **PA:** Proxy applications
  – **AD:** Algorithm development and uncertainty quantification
  – **PM:** Programming models
  – **RT:** Resource/task management
  – **ST:** Scalable tool development
  – **MS:** Performance models and simulators
  – **TA:** Tradeoff analysis and simulation
  – **VS:** Vendor and software (ecosystem) engagement

# Management Plan: Internal Communication Plan

- Each task area operates as its own agile development process, with a monthly to quarterly development cycle

- Example: Proxy Applications (PA) task area

  - *Cycle master: David Richards (Code team lead)*
  - *Stakeholders: Computational materials science, applied math pillars; Center management*
  - *Customers: Tradeoff analysis task area, computer science pillar, vendor and software partners*
  - *Project Team: Co-design code team, programming models task area, computational materials science pillar*
  - *Weekly conference call for team to discuss status, challenges*
  - *Monthly releases drive interactions with customers and stakeholders*

- Quarterly synchronization among all task areas, with in-person Integration Management Team meeting to evaluate progress, prioritize effort (backlog) for next quarter, and reallocate resources if needed.

# Embedded Scale-Bridging Algorithms

- Our goal is to introduce more detailed physics into computational materials science applications in a way which escapes the traditional synchronous SPMD paradigm and exploits the heterogeneity expected in exascale hardware.

- To achieve this, we are developing a UQ-driven *adaptive physics refinement* approach.

- Coarse-scale simulations dynamically spawn tightly coupled and self-consistent fine-scale simulations as needed.

- This *task-based* approach naturally maps to exascale heterogeneity, concurrency, and resiliency issues.



Moving refinement window

Velocity

Microscale    Mesoscale    Macroscale

# Adaptive sampling techniques have been successfully demonstrated by LLNL

- A coarse-scale model (e.g. FEM) calls a lower length-scale model (e.g. polycrystal plasticity) and stores the response obtained for a given microstructure, each time this model is interrogated

- A microstructure-response database is thus populated

- The fine-scale workload varies dramatically over the coarse-scale spatial and temporal domain

- Dynamic workload balancing in a task parallel context



N. R. Barton, J. Knap, A. Arsenlis, R. Becker, R. D. Hornung, and D. R. Jefferson. Embedded polycrystal plasticity and adaptive sampling. *Int. J. Plast.* **24**, 242-266 (2008)

# "A call to arms for task parallelism"



FS queries

FS evaluations

**464 cores: 51x speedup**

Table I. Machine configuration for 29 compute nodes (464 cores).

| Component | Instances | Processes/instance | Total nodes |
|---|---|---|---|
| CS | 1 | 192 | 12 |
| ServerProxy | 1 | 1 | 1 |
| FS | 8 | 32 | 16 |

**2272 cores: 97x speedup**

Table II. Machine configuration for 142 compute nodes (2272 cores).

| Component | Instances | Processes/instance | Total nodes |
|---|---|---|---|
| CS | 1 | 192 | 12 |
| ServerProxy | 2 | 1 | 2 |
| FS | 64 | 32 | 128 |

## A call to arms for task parallelism in multi-scale materials modeling[‡]

Nathan R. Barton[1,*,†], Joel V. Bernier[1], Jaroslaw Knap[2], Anne J. Sunwoo[1], Ellen K. Cerreta[3] and Todd J. Turner[4]

[1]Lawrence Livermore National Laboratory, Livermore, CA 94550, U.S.A.
[2]U.S. Army Research Laboratory, Aberdeen Proving Ground, MD 21005, U.S.A.
[3]Los Alamos National Laboratory, Los Alamos, NM 87545, U.S.A.
[4]U.S. Air Force Research Laboratory, Wright Patterson AFB, OH 45433, U.S.A.

### SUMMARY

Simulations based on multi-scale material models enabled by adaptive sampling have demonstrated speedup factors exceeding an order of magnitude. The use of these methods in parallel computing is hampered by dynamic load imbalance, with load imbalance measurably reducing the achieved speedup. Here we discuss these issues in the context of task parallelism, showing results achieved to date and discussing possibilities for further improvement. In some cases, the task parallelism methods employed to date are able to restore much of the potential wall-clock speedup. The specific application highlighted here focuses on the connection between microstructure and material performance using a polycrystal plasticity-based multi-scale method. However, the parallel load balancing issues are germane to a broad class of multi-scale problems. Copyright © 2011 John Wiley & Sons, Ltd.

# Embedded Scale-Bridging Algorithms

- Scale-bridging algorithms require a consistent two-way algorithmic coupling between temporally evolving distinct spatial levels; they are not "modeling", and not one-way information flow.

- Our focus is on coupling between macro (coarse-scale model) and meso (fine-scale model) scales with all unit physics being deterministic.

- We begin by building off of our adaptive sampling success, but move to the use of temporally evolving mesoscale and spatial adaption.

- Similar concepts apply in the time domain, e.g. using *ab initio* techniques to compute activation energies for a rate theory or kinetic Monte Carlo model ("on-the-fly kMC") applied to radiation damage modeling.

# Agile Proxy Application Development

- Petascale single-scale SPMD and scale-bridging MPMD proxy apps will be used to explore algorithm and programming model design space with domain experts, hardware architects and system software developers.

- These proxy applications will not be "toy models", but will realistically encapsulate the workload, data flow and mathematical algorithms of the full applications.

# Agile Proxy Application Development

- Proxy apps for single-scale SPMD applications (e.g. molecular dynamics) will be used to assess node-level issues including:

  - *Data structures*

  - *Hierarchical memory storage and access*

  - *Power management strategies*

  - *Node-level performance*

- The asynchronous task-based MPMD scale-bridging proxy apps will be used to optimize:

  - *System-level data movement*

  - *Resilience (fault management)*

  - *Load balancing techniques*

  - *Performance scalability*

- These proxy apps are **not** static entities, but the central mechanism for our co-design process.

# Proxy application suite (single-scale)

- First-Principles Molecular Dynamics (MD): Qbox
  - *Dense linear algebra and spectral transform operations*
  - *2006 Gordon Bell Prize (2005 finalist)*
- Tight-Binding MD: LATTE
- Classical MD (Pair-like potentials): SPaSM
  - *Particle-based, spatial (linked-cell) domain decomposition*
  - *In situ visualization demonstrated to 1 trillion atoms on BlueGene/L*
  - *1993, 1998 Gordon Bell Prizes (2005, 2008 finalist)*
- Classical Molecular Dynamics (Many-body potentials): ddcMD
  - *Particle-based, particle domain decomposition*
  - *Soft error recovery demonstrated to CPU-millenium on BlueGene/L*
  - *2005, 2007 Gordon Bell Prizes (2009 finalist)*
- Phase Field Method: AMPE/GL
- Polycrystal plasticity: VP-FFT

# Hierarchical Programming Models

- The challenge for programming models in the context of this project is that they need to expose hardware capabilities to the application programmer while at the same time hiding the continuous flux and complexity of the underlying hardware.

- A hierarchy of programming models exposes and exploits the heterogeneity while providing a transparent layer of abstraction that insulates the application programmer from the flux and complexity of the underlying hardware.

- The programming models and approaches developed to achieve our scale-bridging materials application will be broadly applicable to a variety of multiscale, multiphysics applications:

| | |
|---|---|
| Astrophysics & the structure of the universe | Structural engineering |
| Climate and weather prediction | Plasma physics |
| Nuclear reactor simulation | Radiation hydrodynamics |

# Hierarchical Programming Models

- This hierarchy will replace the traditional bulk synchronous parallel paradigm:

  - *On-node task parallelism* will allow us to couple multiple tightly coupled application components or segments while exploiting on-node resources to their full extent.

  - *Inter-node cooperative parallelism* will provide the necessary capabilities to execute scalable, dynamically structured MPMD applications.

  - *Domain specific languages* aim to encapsulate these levels, enable programmer productivity, and bridge disparate architectures.

SIMD

Tightly Coupled — On node threading

Loosely Coupled — Task Parallel Computation

Tightly Coupled — Message Passing Paradigm

Loosely Coupled — MPMD Cooperative Parallelism

Load Balancing

Fault Isolation

# Holistic Analysis and Optimization

- A hierarchy of performance models, simulators, and emulators are used to explore algorithm, programming model, and hardware design space before the application is fully constructed.

  - ASPEN: Rapid exploration of design space using application skeletons
  - SST: Detailed simulation of data flow, performance and energy/power cost
  - GREMLIN: Emulation layer to mimic exascale complexity by injecting faults, OS jitter, and other noise to "stress test" the application/SW stack

# Summary

- Our objective is to establish the interrelationship between algorithms, system software, and hardware required to develop a multiphysics exascale simulation framework for modeling materials subjected to extreme mechanical and radiation environments.



- This effort is focused in four areas:

  - *Scale-bridging algorithms*
    - » UQ-driven adaptive physics refinement

  - *Programming models*
    - » Task-based MPMD approaches to leverage concurrency and heterogeneity at exascale while enabling fault tolerance

  - *Proxy applications*
    - » Communicate the application workload to the hardware architects and system software developers, and used in performance models/simulators/emulators

  - *Co-design analysis and optimization*
    - » Optimization of algorithms and architectures for performance, memory and data movement, power, and resiliency

# Kickoff meeting: Aug 24-26 @ Santa Fe, NM

- Over 40 participants
  - *Vendors: AMD, Cray, HP, Intel, IBM, Nvidia*
  - *DOE computational materials science community: ASC, CASL, MaRIE*
  - *Program & line management*
- Three sessions:
  - *Stakeholder input*
  - *Task area discussion*
  - *Y1 work plan development*