# Proxy and Proto miniApps for Exascale Co-design

**James A. Ang, and Richard F. Barrett**
**Scalable Computer Architectures Department**
**Sandia National Laboratories**
**Albuquerque, NM   USA**

**International Supercomputing Conference**
**Leipzig, Germany**
**June 17-20, 2013**

# Our Exascale Efforts Focus on Co-design

Friday, June 14, 13

# Our Exascale Efforts Focus on Co-design

- Key Co-Design capabilities
  - *SST* Architectural Simulation Framework
  - Pre-production, First-of-a-kind Testbeds
  - Scalable R&D System Software
  - *Mantevo* miniApplications

Friday, June 14, 13

# Our Exascale Efforts Focus on Co-design

- Key Co-Design capabilities
  - *SST* Architectural Simulation Framework
  - Pre-production, First-of-a-kind Testbeds
  - Scalable R&D System Software
  - *Mantevo* miniApplications

- Exascale Implications
  - Sustained commitment of significant funding
  - Requires the Long View >5 years out
    - —> time to *influence* Hardware Architectures

# Which is Harder to change: Hardware or Software?

Friday, June 14, 13

# Which is Harder to change: Hardware or Software?

- Conventional Wisdom
  - Hardware is difficult to change
  - Software is easy to change
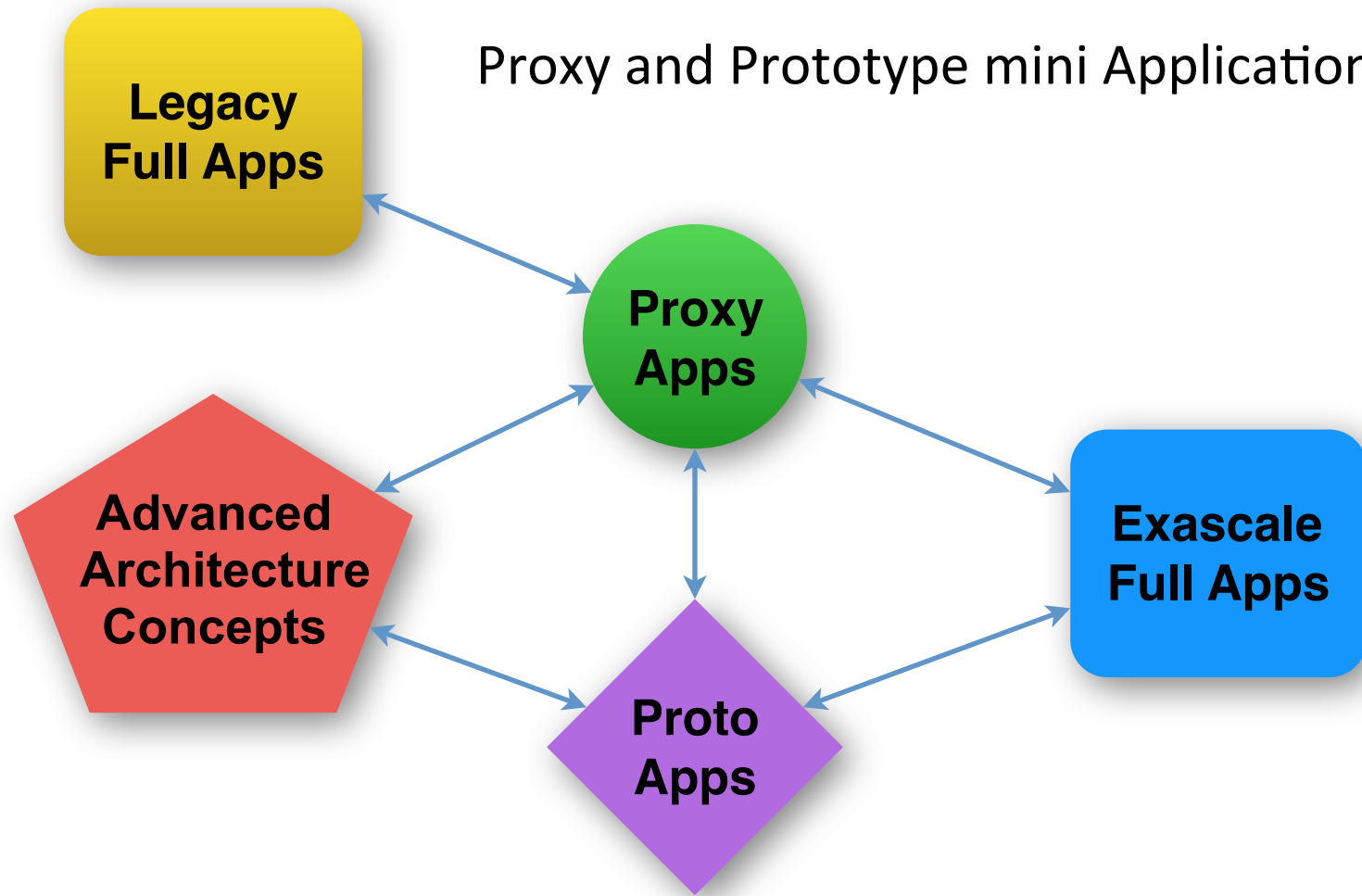
Friday, June 14, 13

# Which is Harder to change: Hardware or Software?

- Conventional Wisdom
  - Hardware is difficult to change
  - Software is easy to change

- For the Long View,

  Conventional Wisdom is <u>wrong</u>!

Friday, June 14, 13

# MiniApps: Proxy and Proto

- Relationship among Full Applications

Proxy and Prototype mini Applications

Friday, June 14, 13

# Representative Legacy app: CTH

- Eulerian multi-material modeling application
- 3D, finite volume stencil computation
- BSP with message aggregation (BSPMA)
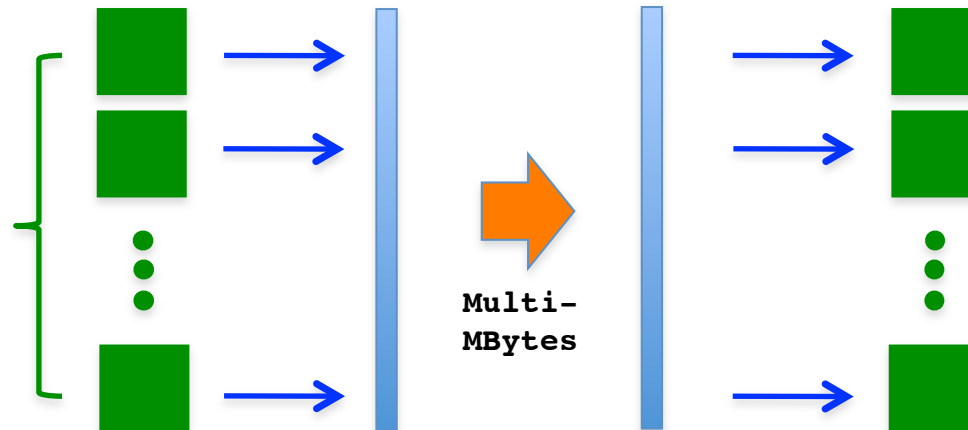


```
do i = 1, num_tsteps
```

num_vars

Multi-
MBytes

```
    do j = 1, num_vars
        compute
    end do
end do
```

# Representative Legacy app: CTH

- Eulerian multi-material modeling application
- 3D, finite volume stencil computation
- BSP with message aggregation (BSPMA)
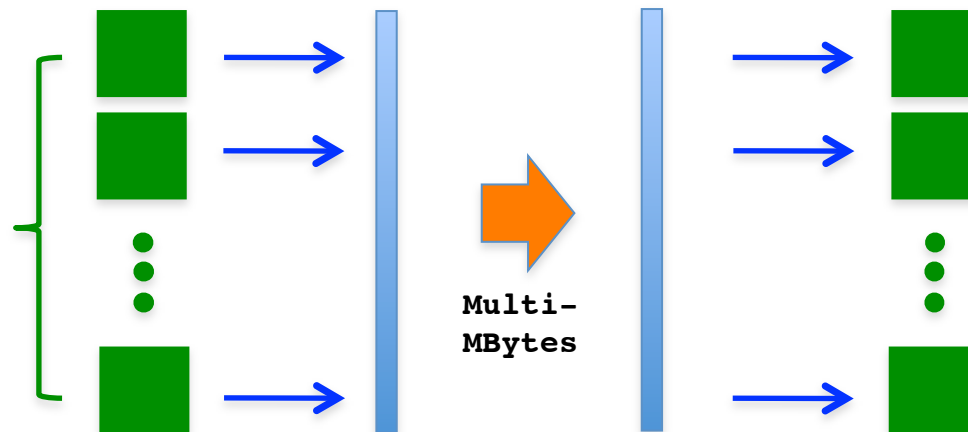
```
do i = 1, num_tsteps
```

*40 vars*
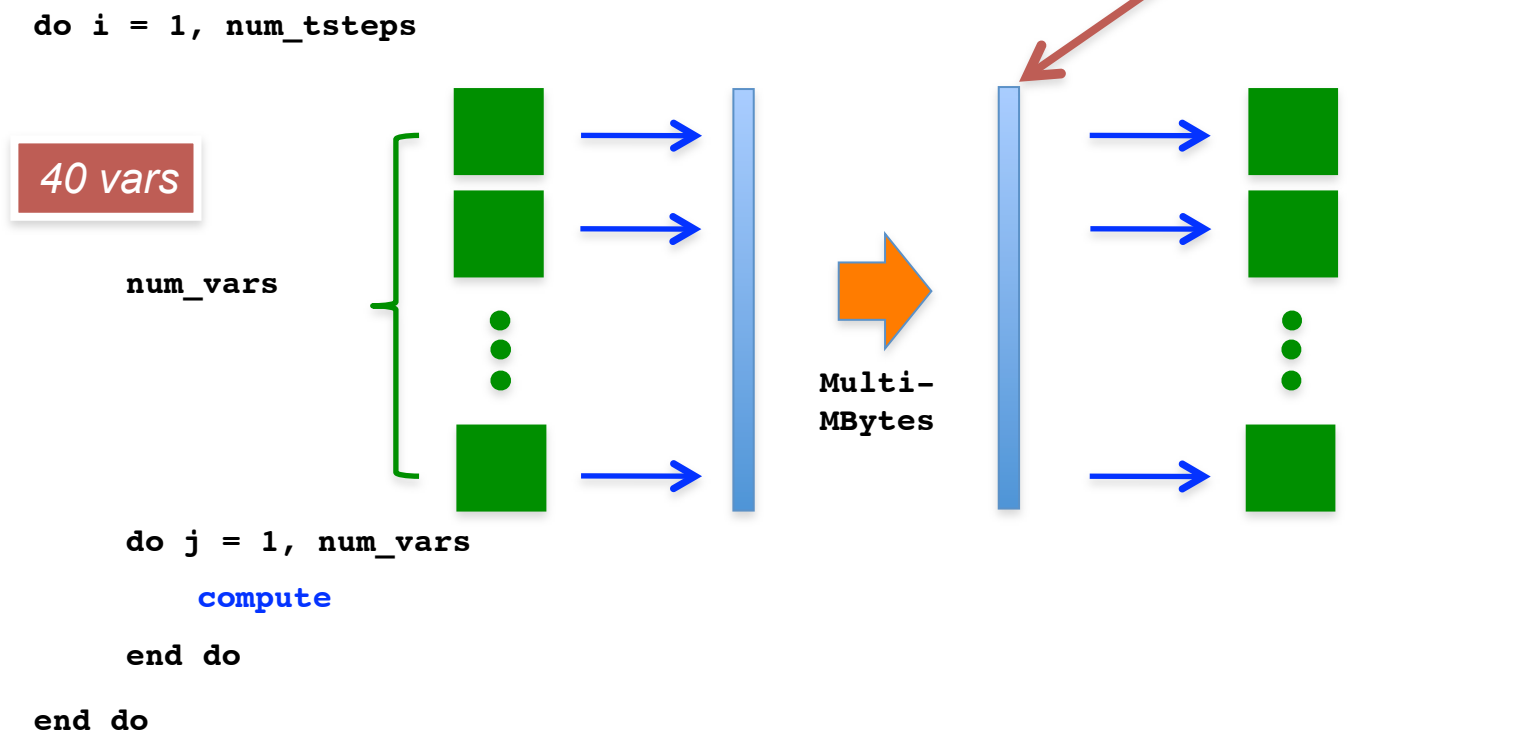
num_vars

Multi-
MBytes

```
do j = 1, num_vars
    compute
end do
end do
```

# Representative Legacy app: CTH

- Eulerian multi-material modeling application

- 3D, finite volume stencil computation

- BSP with message aggregation (BSPMA)

**3.2 MBytes**

```
do i = 1, num_tsteps
```

**40 vars**

**num_vars**

**Multi-MBytes**

```
do j = 1, num_vars
    compute
end do
end do
```

# Alternative inter-node strategy:
# Single Variable Aggregated Faces (SVAF)

```
do i = 1, num_tsteps

  do j = 1, num_vars

      compute
```



```
  end do

end do
```

# Alternative inter-node strategy:
# Single Variable Aggregated Faces (SVAF)

```
do i = 1, num_tsteps

  do j = 1, num_vars

    compute
```


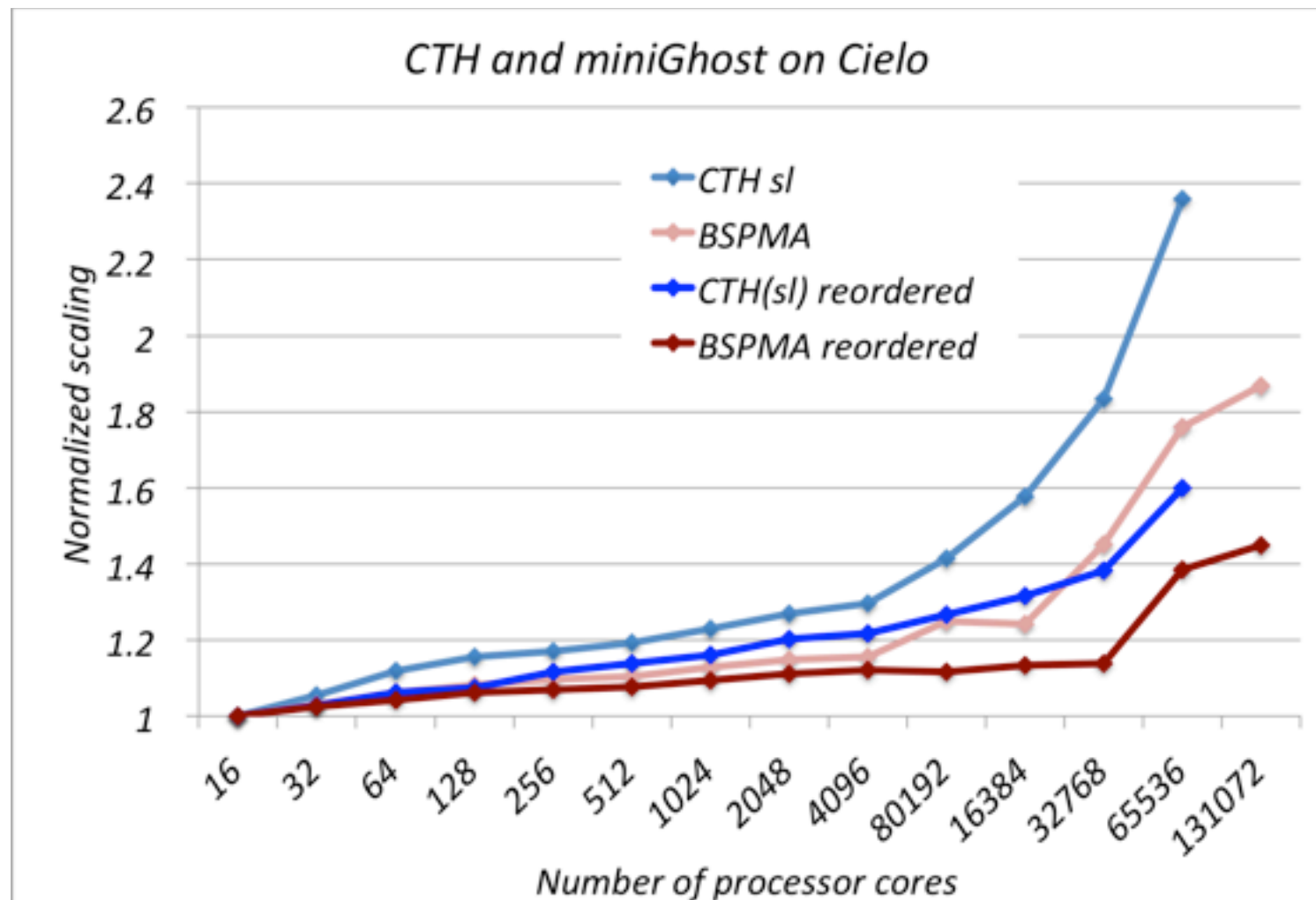
**3.2 Mbytes / 40 = 80kbytes**
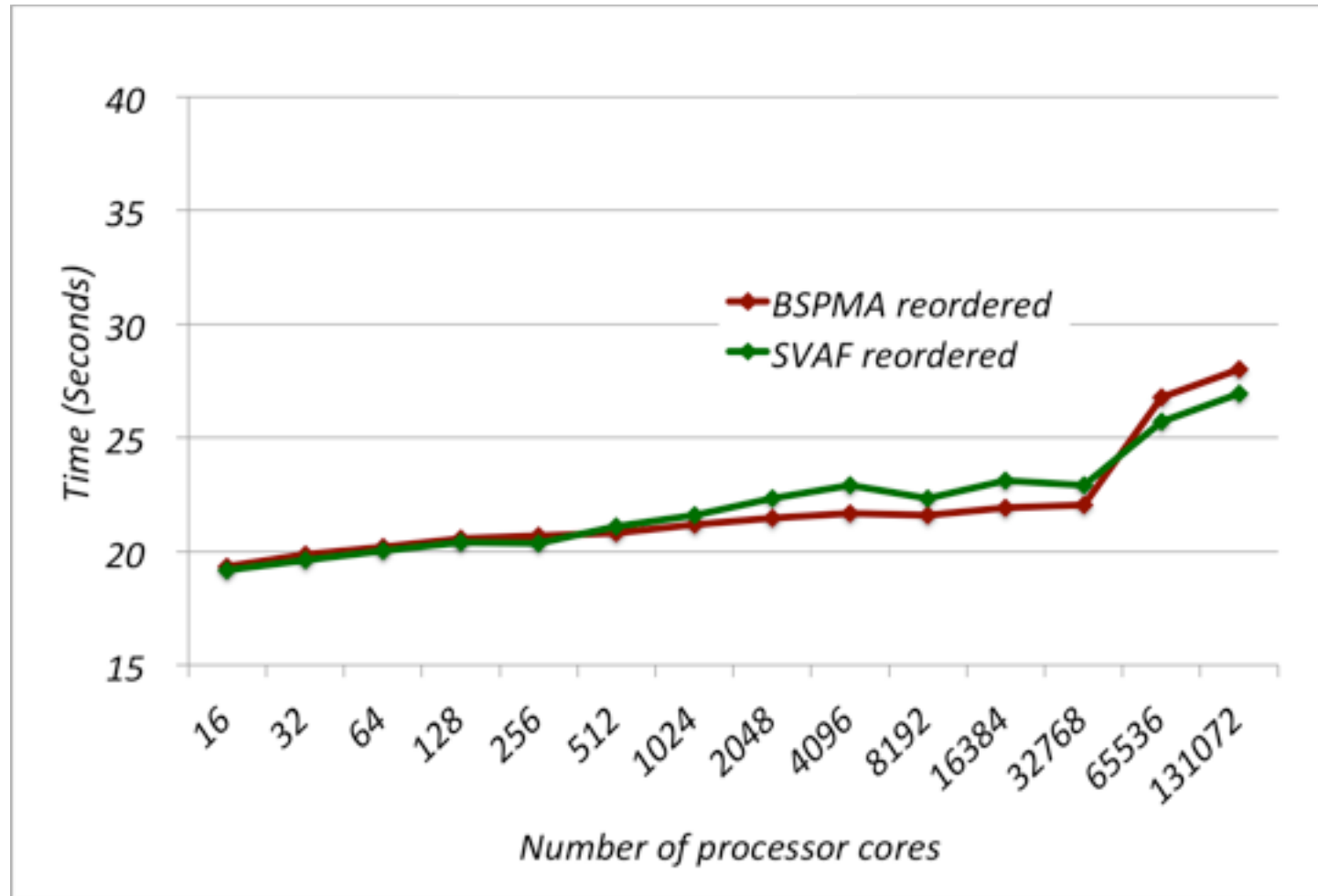
```
  end do

end do
```

# Physical ordering of MPI ranks



CTH and miniGhost on Cielo

Friday, June 14, 13
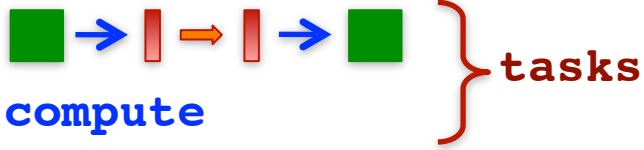
# miniGhost on Cielo: BSPMA and SVAF

# Implications for Co-design at a System Architecture Level

- We know that BSPMA applications create congestion problems for high radix topologies

- Hypothesis: If applications are task parallel, supported with asynchronous runtime software, and multithreaded processors, then congestion issues dissipate

- Focus area for our R&D with *Proto* miniApps
  - Application – Task parallel implementation
  - System Software – Support for asynchronous, adaptive threads
  - Node Architecture – Support for light-weight threading
  - System Architecture – Interconnect Fabric with high radix routers and high injection message rate

Friday, June 14, 13
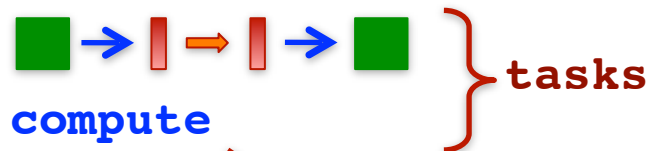
# miniGhost task parallel

```
do i = 1, num_tsteps

  do j = 1, num_vars

    do iblk, jblk, kblk
```



```
    end do

  end do

end do
```

# miniGhost task parallel

```
do i = 1, num_tsteps

  do j = 1, num_vars

    do iblk, jblk, kblk
```



**compute**          **tasks**

```
    end do

  end do

end do
```

Compute *after* comm due to workspace usage model

# miniGhost task parallel

```
do i = 1, num_tsteps

  do j = 1, num_vars

    do iblk, jblk, kblk
```



Needs abstraction (DAG)

**tasks**

**compute**

```
    end do

  end do

end do
```
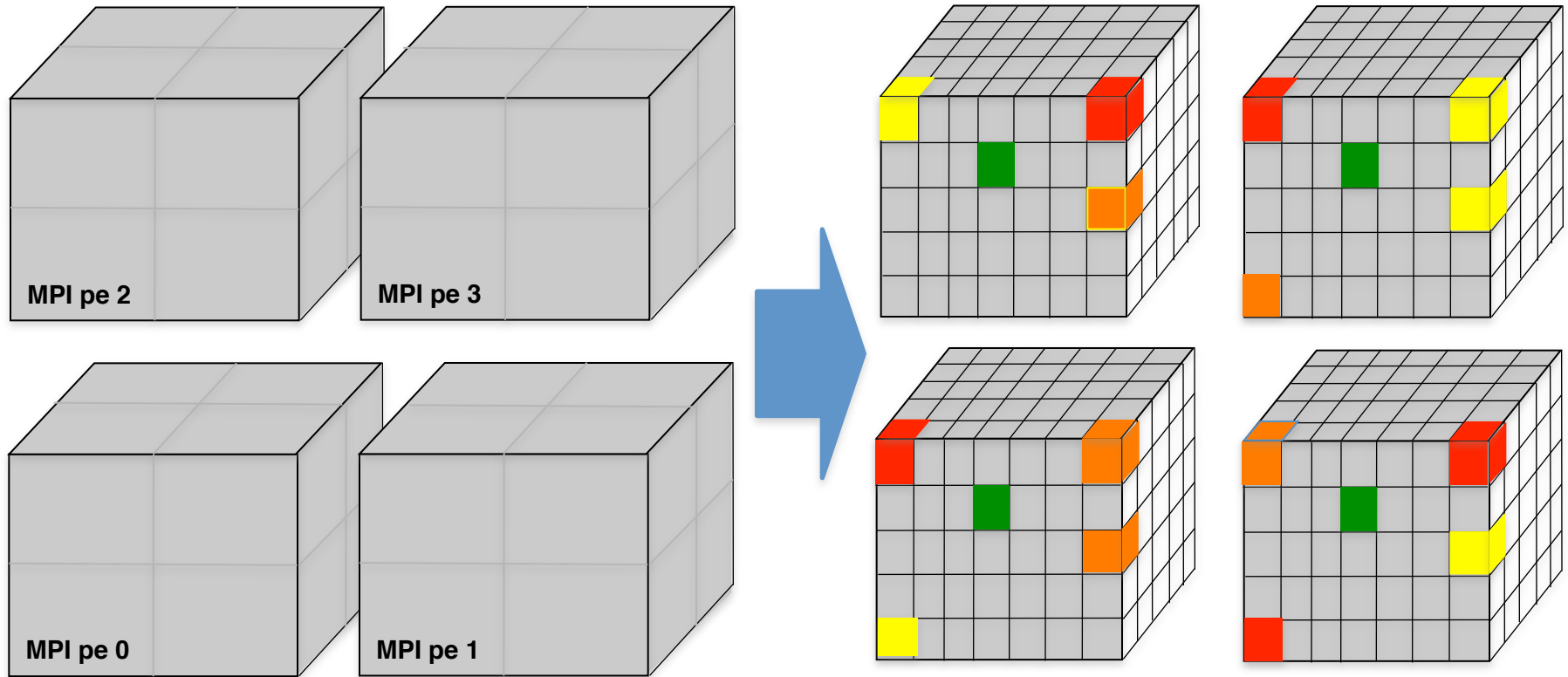
Compute *after* comm due to workspace usage model

# miniGhost: over-decomposition task parallel implementation



Data parallel
☐ thread

Task parallel: some representative task workloads

🟩 Computation          🟧 Computation + MPI
🟨 Computation + BC      🟥 Computation + BC + MPI

MPI pe 2    MPI pe 3
MPI pe 0    MPI pe 1

# miniGhost tp: Adding AMR

- Diffusion over the 3D domain with random initial conditions and reflective boundary conditions
  - Refinement based on the boundary or volume of an object being moved through the mesh and changing size.  So, for example, a shock front can be simulated by refining based on a sphere which starts small and grows in size as the problem advances.

- Refinement within blocks
  - A block is refined into 8 blocks
  - Neighbors must be within one level of refinement
  - Computation is self-contained within a block

- Communication aggregated to BSP model
  - Excellent candidate for task parallelism version

Friday, June 14, 13

# Concluding Thoughts

Friday, June 14, 13

# **Concluding Thoughts**

- Multiple Dimensions of Co-Design
  - HW: Node and System architecture
  - SW: Application and System Software

Friday, June 14, 13

# Concluding Thoughts

- Multiple Dimensions of Co-Design
  - HW: Node and System architecture
  - SW: Application and System Software

- Examine *Conventional Wisdoms* on COTS and system balance
  - Component performance
  - Investment
    - Platform costs
    - R&D investments