

# Versioned Distributed Arrays for Resilience in Scientific Applications: Global View Resilience

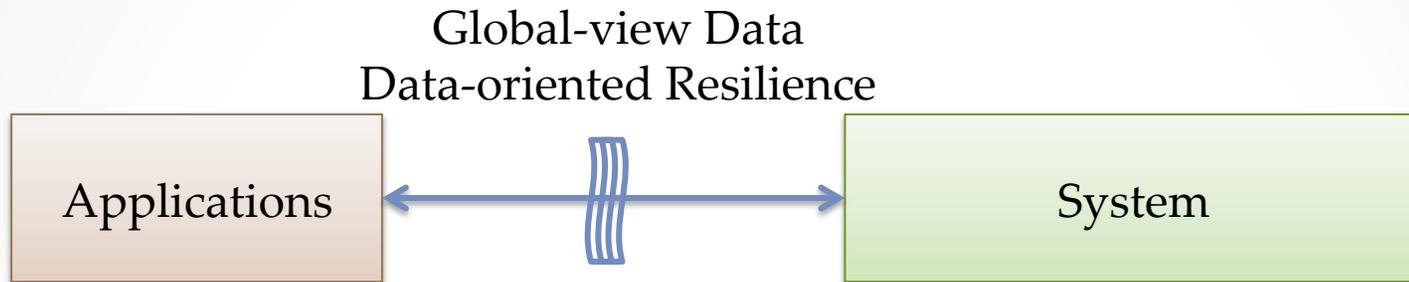
Andrew A. Chien, The University of Chicago  
and Argonne National Laboratory  
ICCS 2015, June 1  
Reykjavik, Iceland



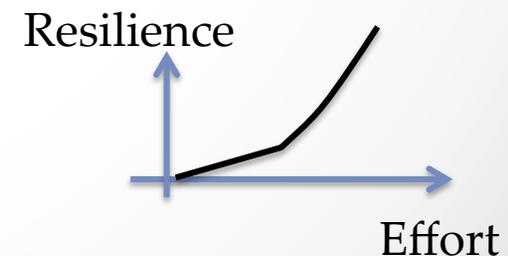
# Outline

- GVR Approach and Flexible Recovery
- GVR in Applications Programming Effort
- GVR Versioning and Recovery Performance
- Summary
- ...More Opportunities with Versioning

# GVR Approach

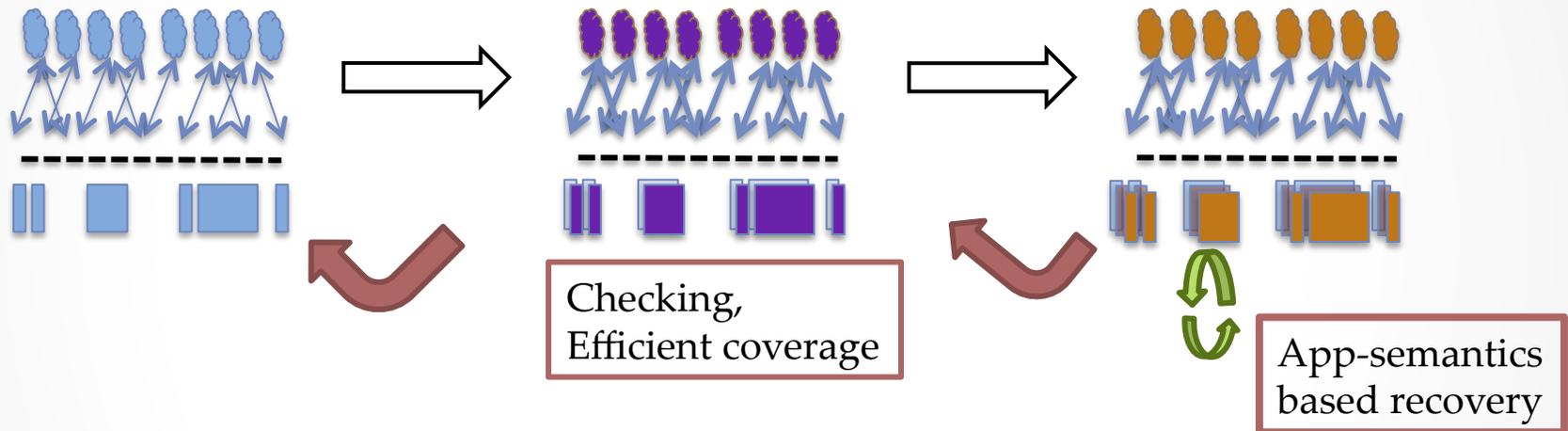


- Application-System Partnership: System Architecture
  - Exploit algorithm and application domain knowledge
  - Enable “End to end” resilience model (outside-in), Levis’ Talk
- Portable, Flexible Application control (performance)
  - Direct Application use or higher level models (task-parallel, PGAS, etc.)
  - GVR Manages storage hierarchy (memory, NVRAM, disk)
  - GVR ensures data storage reliability, covers error types
- Incremental “Resilience Engineering”
  - Gentle slope, Pay-more/Get-more, Anshu’s talk



# Data-oriented Resilience based on Multi-versions

Phases create new  
logical versions



- Global-view data – flexible recovery from data, node, other errors
- Versioning/Redundancy customized as needed (per structure)
- Error checking & recovery framed in high-level semantics (portable)

# GVR Concepts and API

- Create Global view structures

- New, federation interfaces
- `GDS_alloc(...)`, `GDS_create(...)`

- Global view Data access

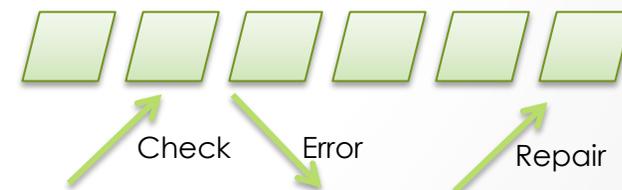
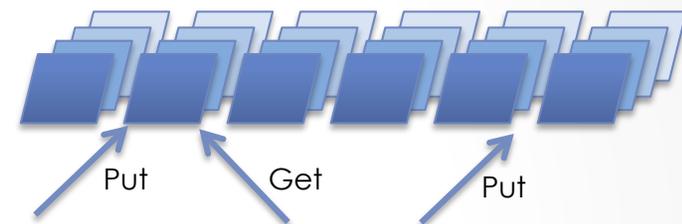
- Data: `GDS_put()`, `GDS_get()`
- Consistency: `GDS_fence()`, `GDS_wait()`, ...
- Accumulate: `GDS_acc()`, `GDS_get_acc()`, `GDS_compare_and_swap()`

- Versioning

- Create: `GDS_version_inc()`, Navigate: `GDS_get_version_number()`, `GDS_move_to_newest()`, ...

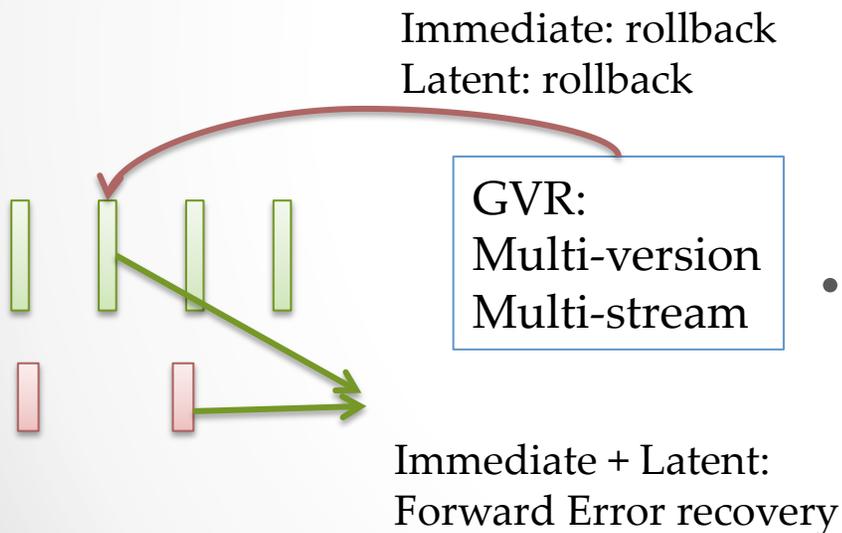
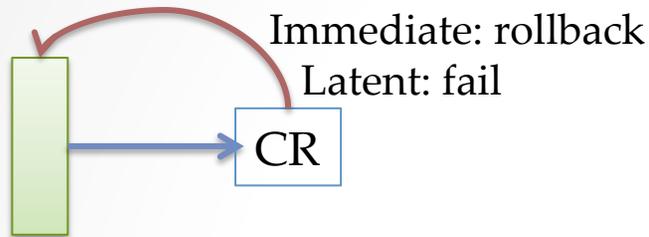
- Error handling

- Application checking, signaling, correction: `GDS_raise_error()`,



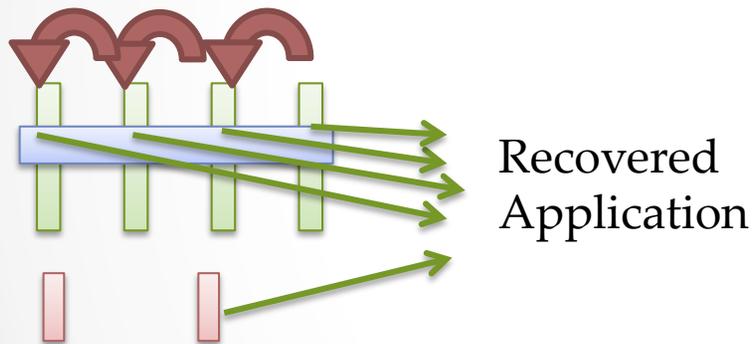
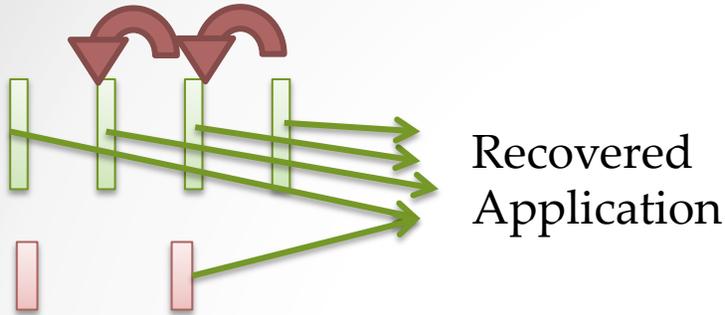
Applications have portable control over coverage and overhead of resilience.

# GVR Flexible Recovery I



- Immediate errors: Rollback
- Latent/Silent errors: multi-version
  - Application recovery using multiple streams
- Immediate + Latent: novel forward error recovery
  - System or application recovery using approximation, compensation, recomputation, or other techniques
- Tune version frequency, data structure coverage, increased ABFT and forward error recovery for rising error rates

# GVR Flexible Recovery II



- Complex errors, Rollback-diagnosis-forward
  - Flexible, Application-based recovery
  - Walk multiple versions
  - Diagnose
  - Compute corrections/approximations, execute forward
- Complex errors, Forward from multiple versions
  - Flexible, Application-based recovery
  - Partial materialization of multiple versions
  - Compute approximations, execute forward
- Tune version frequency, data structure coverage, increased ABFT and forward error recovery for rising error rates

GVR flexibility enables scalability across a wide range of error types and rates.

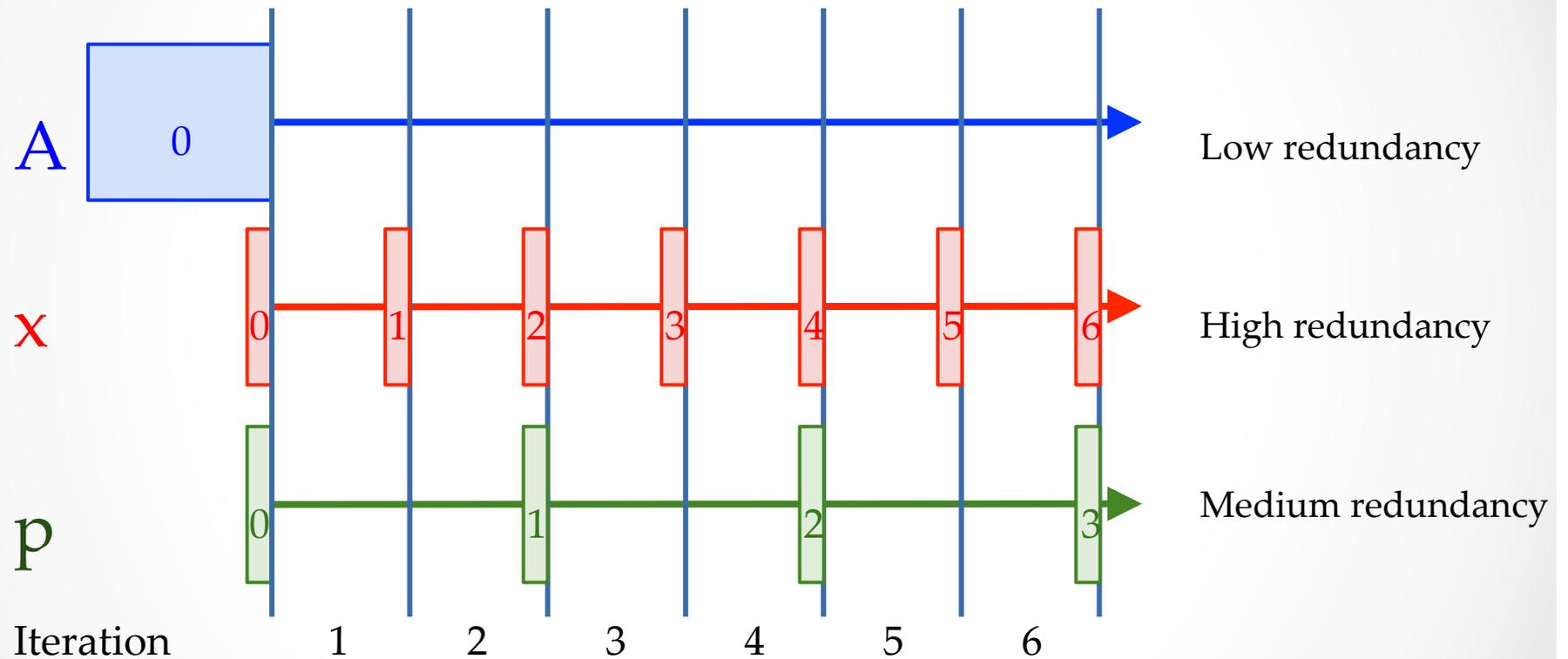
# Simple Version Recovery: Preconditioned Conjugate Gradient

$A = \dots$

```
1:  $r = b - Ax$ 
2:  $iter = 0$ 
3: while ( $iter < max\_iter$ ) and  $\|r\| > tolerance$  do
4:    $iter = iter + 1$ 
5:    $z = M^{-1}r$ 
6:    $\rho_{old} = \rho$ 
7:    $\rho = (r, z)$ 
8:    $\beta = \rho / \rho_{old}$ 
9:    $p = z + \beta p$ 
10:   $q = Ap$ 
11:   $\alpha = \rho / (p, q)$ 
12:   $x = x + \alpha p$ 
13:   $r = r - \alpha q$ 
14: end while
```

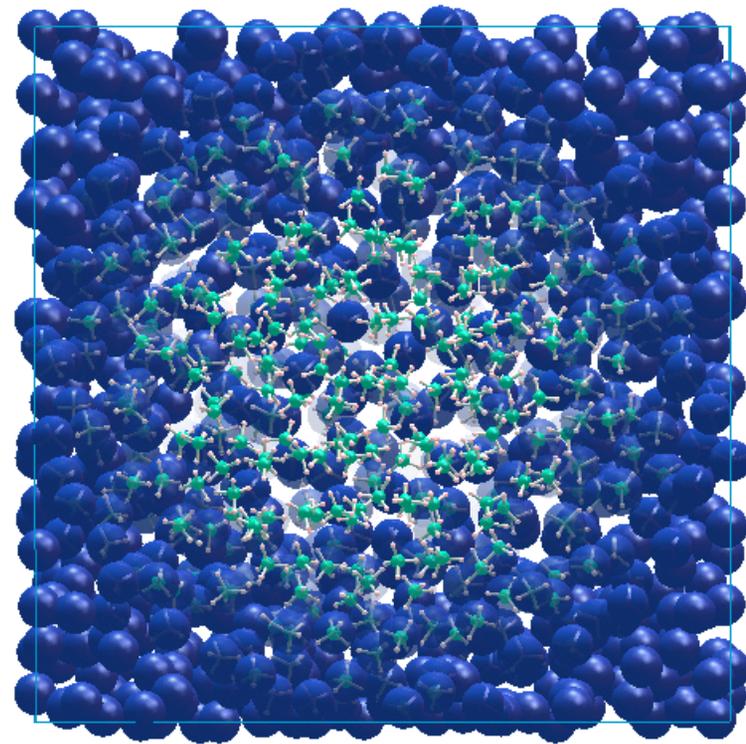
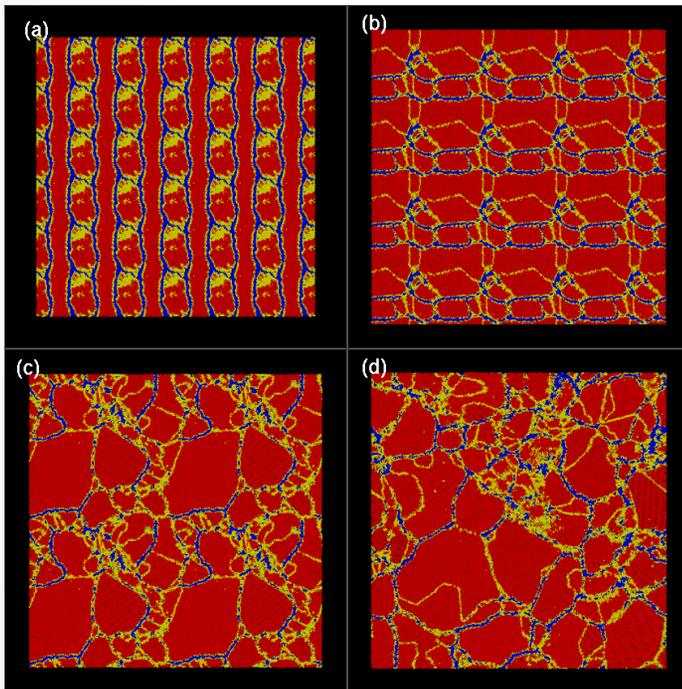
- Version x “solution vector”
  - Restore x on error
- Version p “direction vector”
  - Restore on error
- Version A “linear system”
  - Restore on error
- Restore from which version?
  - Most recent (immediately detected errors)
  - Older version (latent or “silent” errors)

# Multi-stream in PCG: Matching redundancy to need



# Molecular Dynamics: miniMD, ddcMD

- miniMD: a SNL mini-app, a version of LAMMPS
- ddcMD is the atomistic simulation developed by LLNL -- scalable and efficient.



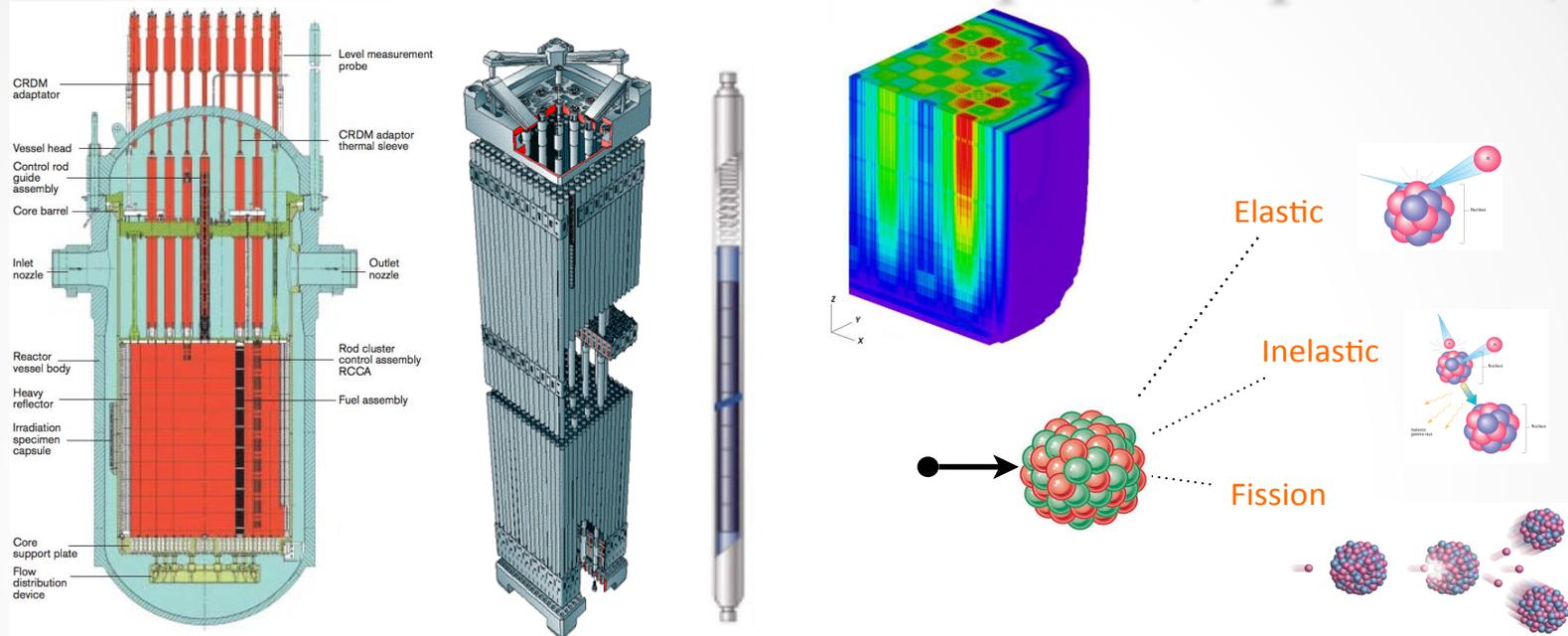
• (c) Andrew A. Chien

LLNL (Dave Richards & Ignacio Laguna)

# ddcMD + GVR

```
main() {
    /* store essential data structures in gds */
    GDS_alloc(&gds);
    /* specify recovery function for gds */
    GDS_register_global_error_handler(gds, recovery_func);
    simulation_loop() {
        computation();
        error = check_func() /* finds the errors */
        if (error) {
            error_descriptor = GDS_create_error_descriptor(GDS_ERROR_MEMORY)
            /* signal error */
            /* trigger the global error handler for gds */
            GDS_raise_global_error(gds, error_descriptor);
        }
        if (snapshot_point) {GDS_version_inc(gds);
            GDS_put(local_data_structure, gds);};
    }
}
/* Simple recovery function, rollback */
recovery_func(gds, error_desc) {
    /* Read the latest snapshot into the core data structure */
    GDS_get(local_data_structure, gds);
    GDS_resume_global(gds, error_desc);
}
```

# Monte Carlo Neutron Transport (OpenMC)



- High fidelity, computation intensive and large memory (100GB~ cross sections and 1TB~ tally data)
- Particle-based parallelization is used with data decomposition
- Partition tally data by global array
- OpenMC: best scaling production code
- DOE CESAR co-design center “co-design application”

# OpenMC + GVR

Initialize initial neutron positions

*GDS\_create(tally & source\_site); // Create global tally array and source sites*

**for each** batch

**for each** particle in batch

**while** (not absorbed)

      move particle and sample next interaction

**if** fission

*GDS\_acc(score, tally) // tally, add score asynchronously*

        add new source sites

**end**

*GDS\_fence() // Synchronize outstanding operations*

    resample source sites & estimate eigenvalue

*if (take\_version) GDS\_ver\_inc(tally) // Increment version*

*GDS\_ver\_inc(source\_site) // Increment version*

**end**

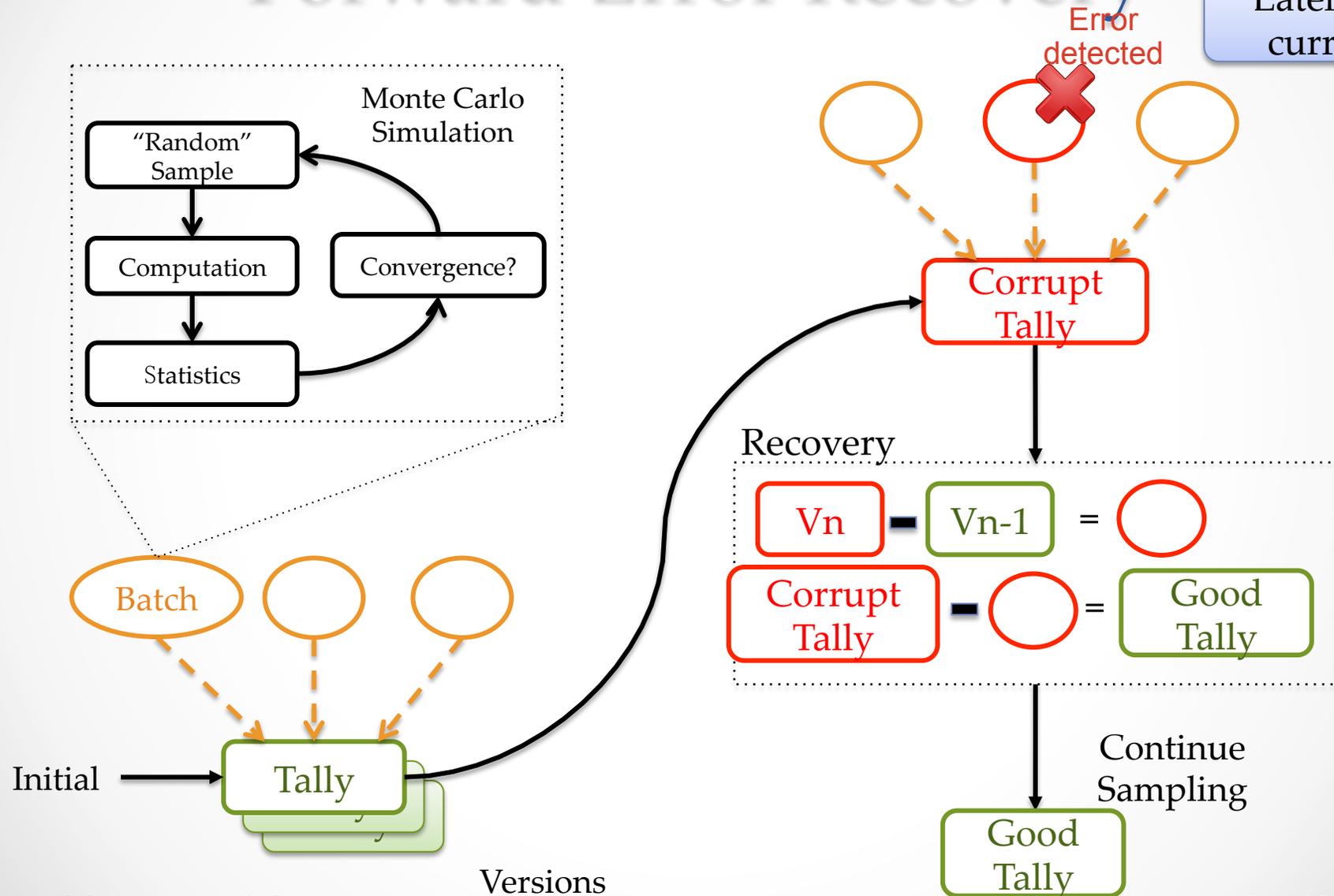
**end**

• (c) Andrew A. Chien

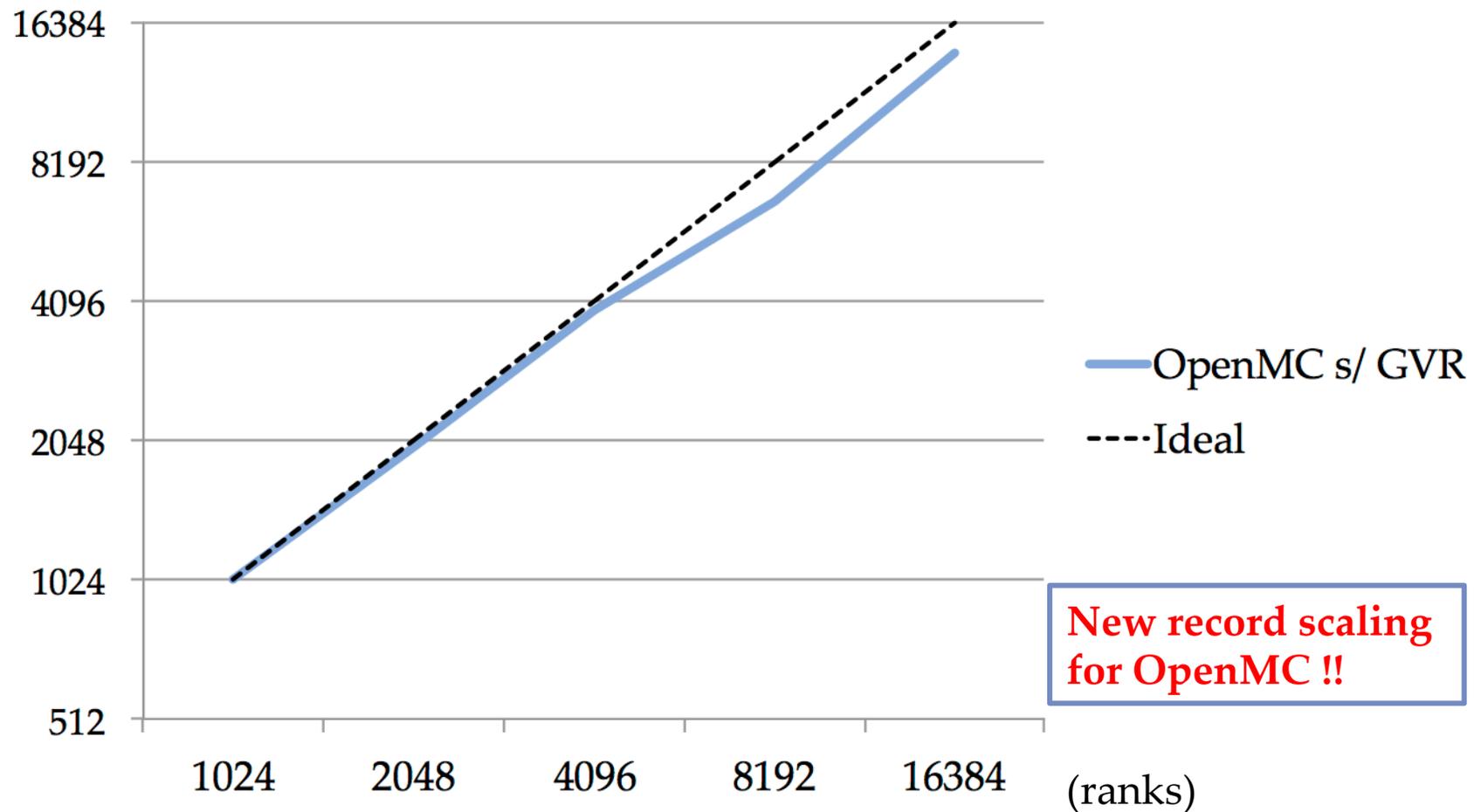
- Create Global view tallies
- Versioning: 259 LOC (<1%)
- Forward recovery: 250 (<1%)
- Overall application: 30 KLOC

# Monte Carlo “Compensating” Forward Error Recovery

Latent or  
current



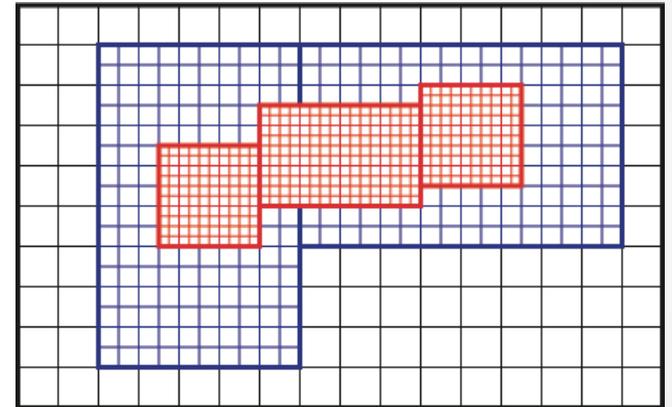
# OpenMC+GVR Performance



N. Dun, H. Fujita, J. Tramm, A. Chien, and A. Siegel. Data Decomposition in Monte Carlo Neutron Transport Simulations using Global View Arrays, IJHPCA, May 2014

# Chombo + GVR

- Resilience for core AMR hierarchy
  - Central to Chombo
  - Lessons applicable to Boxlib (ExaCT co-design app)
- Multiple levels, each with own time-step
- Data corruption and Process Crash Resilience
  - GVR used to version each level separately
  - Exploits application-level snapshot-restart
- GVR as vehicle to explore cost models for “resilience engineering” (Dubey)
  - Future: customize or localize recovery

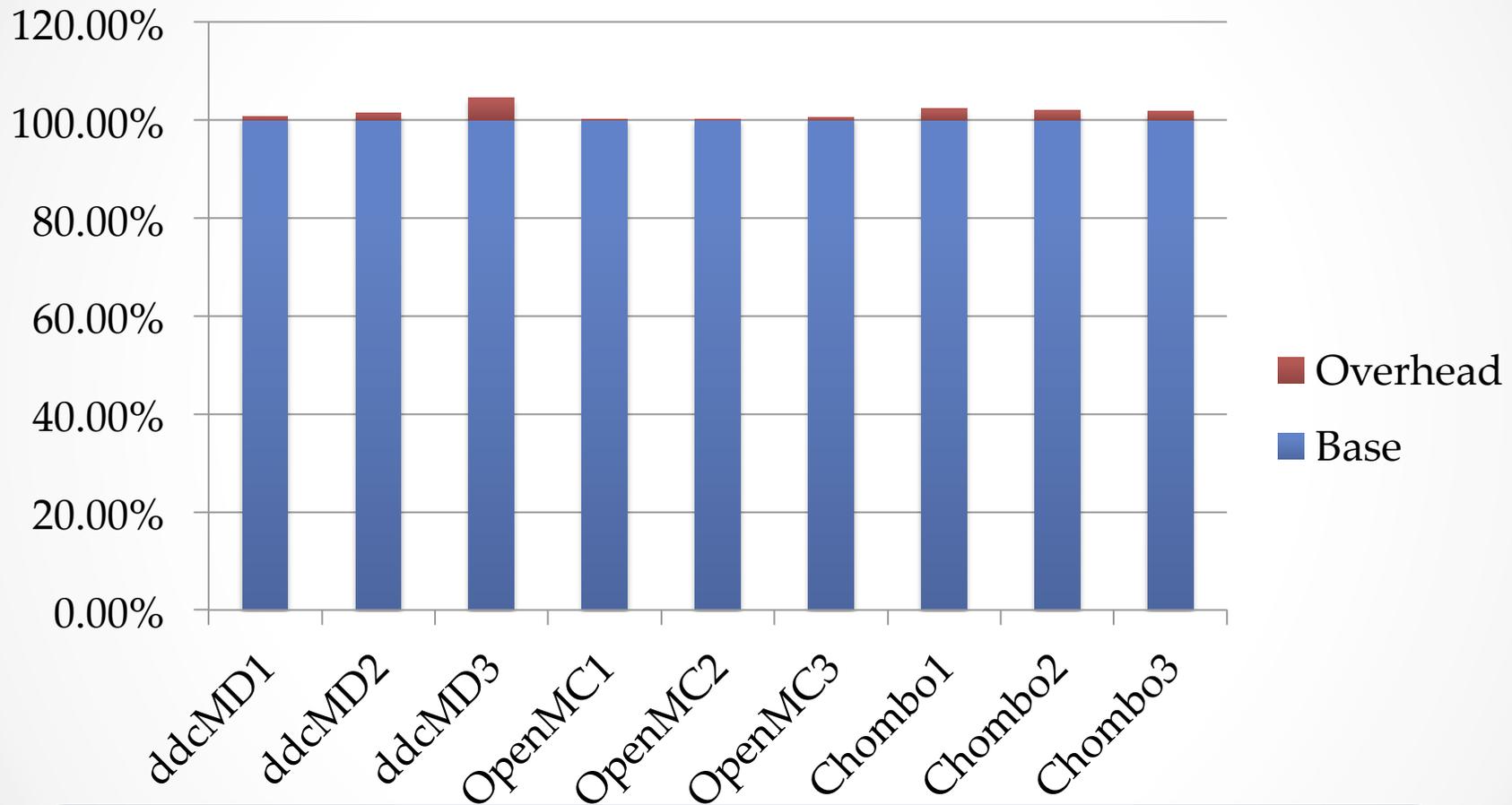


# GVR Gentle Slope

Code/ Application	Size (LOC)	Changed (LOC)	Leverage Global View	Change SW architecture
Trilinos/PCG	300K	<1%	Yes	No
Trilinos/ Flexible GMRES	300K	<1%	Yes	No
OpenMC	30K	<2%	Yes	No
ddcMD	110K	<0.3%	Yes	No
Chombo	500K	<1%	Yes	No

GVR enables a gentle slope to Exascale resilience

# GVR Performance (Overhead)



Varied version frequency, against the native program. All < 2%.

- GVR performance scales over versions and partial materialization too!

# GVR Summary

- Easy to add to an application
- Flexible control and coverage
- Flexible recovery (enables variety of forward techniques, approximations, etc.)
- Low overhead
- Efficient version restore (across versions)
- Efficient incremental restore

All Portable!

# Additional GVR Research

...

## Latent or “silent” error model

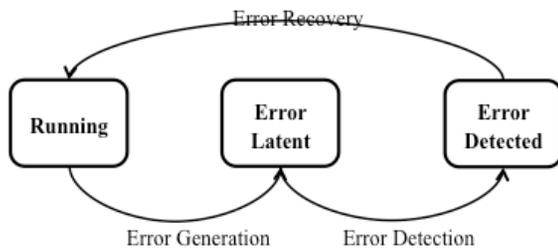
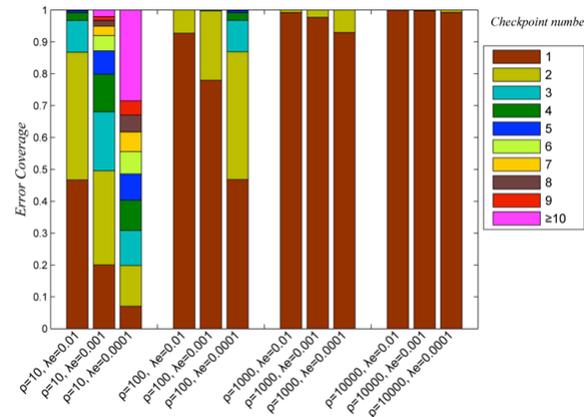
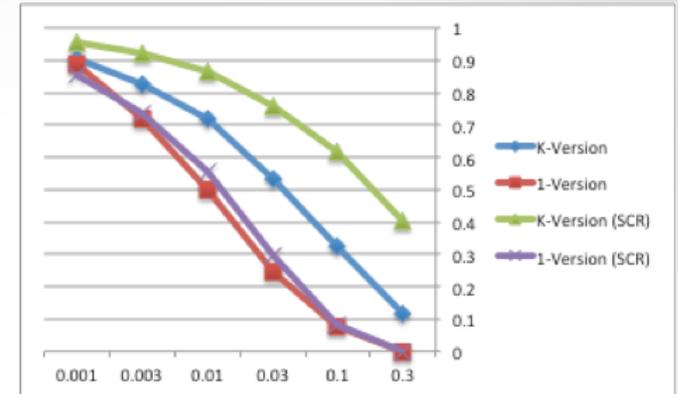


Fig. 1.b Latent Error Model



Multi-version critical for difficult to detect errors

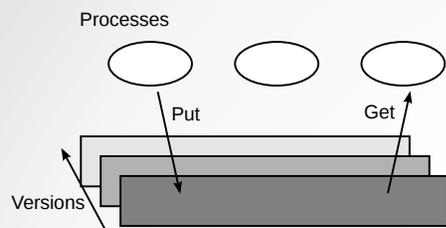


Multi-version increases efficiency at high error rates

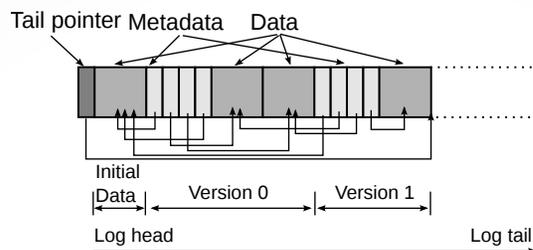
# Latent Error Recovery

• • •

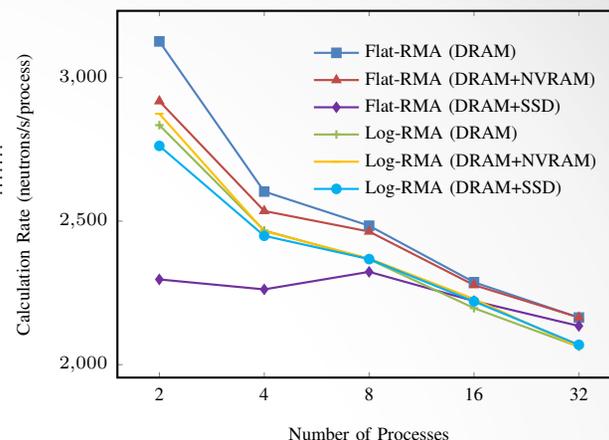
When multiple versions are useful  
 Impact on high-error rate regimes  
 Impact on difficult to detect errors



Flat (Traditional)



Log-structured



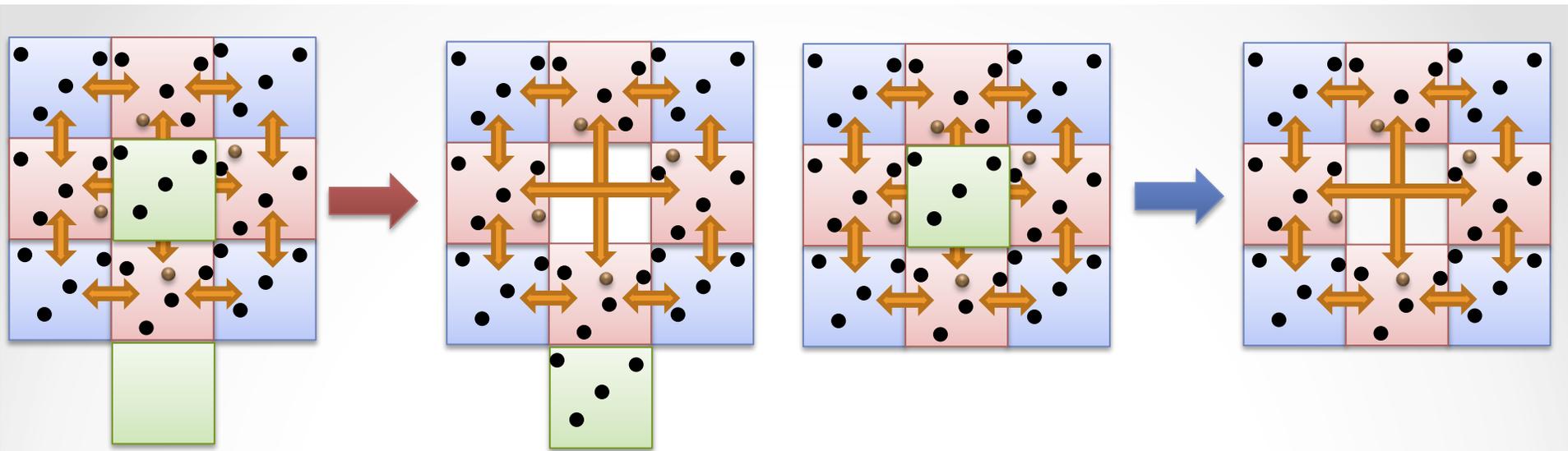
Comparative Studies  
with applications +  
varied memory hierarchies

# Efficient Versioning

- Different implementations (SW, HW, OS, Application)
  - OS page tracking, dirty bits, SW declared
  - Skewed and Multi-version in-memory representations
- Efficient storage and materialization
- Leverages collective view
- Exploit NVRAM, burst buffers, etc.

H. Fujita, N. Dun, Z. Rubenstein, and A. Chien. Log-structured global array for efficient multi-version snapshots. CCGrid, May 2015.

H. Fujita, K. Iskra, P. Balaji, and A. Chien, "Empirical Characterization of Versioning Architectures", submitted.



# N+1- $\rightarrow$ N and N- $\rightarrow$ N-1 Recovery

- MPI Recovery (ULFM)
- Application Process Recovery
- Load Balancing and Performance
- Post-recovery Efficiency (PRE)

# GVR Software Status

- Open source release, Oct 2014 ([gvr.cs.uchicago.edu](http://gvr.cs.uchicago.edu))
  - Tested with Miniapps – miniMD, miniFE experiments, and Full apps – ddcMD, PCG, OpenMC, Chombo
- Features
  - Versioned distributed arrays with global naming (a portable abstraction)
  - Independent array versioning (each at its own pace)
  - Reliable storage of the versioned arrays in memory, local disk/ssd, or global file system (thanks to Adam and SCR team!)
  
  - Whole version navigation and efficient restoration
  - Partial version efficient restoration (partial “materialization”)
  - C native APIs and Fortran bindings
  - Runs on IBM Blue Gene, Cray XC, and Linux Clusters
- Key: all of the application investment is portable because the abstractions are portable

# More GVR Info I

## Basic API's and Usage

- GVR Team. Gvr documentation, release 0.8.1-rc0. Technical Report 2014-06, University of Chicago, Department of Computer Science, 2014.
- GVR Team. How applications use gvr: Use cases. Technical Report 2014-05, University of Chicago, Department of Computer Science, 2014.

## GVR Architecture and Implementation Research

- Hajime Fujita, Kamil Iskra, Pavan Balaji, and Andrew A. Chien, "Empirical Characterization of Versioning Architectures", in CLUSTER, October 2015.
- A. Fang and A. Chien, "How Much SSD Is Useful for Resilience in Supercomputers", in IEEE Symposium on Fault-tolerance at Extreme-Scale (FTXS), June 2015.
- Hajime Fujita, Nan Dun, Zachary A. Rubenstein, and Andrew A. Chien. Log-structured global array for efficient multi-version snapshots. In CCGrid 2015..
- Guoming Lu, Ziming Zheng, and Andrew A. Chien. When is multi-version checkpointing needed? In Proceedings of the 3rd Workshop on Fault-tolerance for HPC at extreme scale, ACM FTXS '13, July 2013.
- Wesley Bland, Aurelien Bouteiller, Thomas Herault, Joshua Hursey, George Bosilca, and JackJ. Dongarra. An evaluation of User-Level Failure Mitigation support in MPI. Computing, 95(12):1171–1184, 2013.

# More GVR Info II

## Application Studies

- A. Chien, P. Balaji, N. Dun, A. Fang, H. Fujita, K. Iskra, Z. Rubenstein, Z. Zheng, J. Hammond, I. Laguna, D. Richards, A. Dubey, B. van Straalen, M Hoemmen, M. Heroux, K. Teranishi, A. Siegel. Exploring Versioning for Resilience in Scientific Applications: Global-view Resilience, submitted for publication, March 2015. (Best overall project summary)
- A. Chien, P. Balaji, P. Beckman, N. Dun, A. Fang, H. Fujita, K. Iskra, Z. Rubenstein, Z. Zheng, R. Schreiber, J. Hammond, J. Dinan, A. Laguna, D. Richards, A. Dubey, B. van Straalen, M Hoemmen, M. Heroux, K. Teranishi, A. Siegel, and J. Tramm, "Versioned Distributed Arrays for Resilience in Scientific Applications: Global View Resilience", in International Conference on Computational Science (ICCS 2015), Reykjavik, Iceland, June 2015.
- Nan Dun, Hajime Fujita, John R. Tramm, Andrew A. Chien, and Andrew R. Siegel. Data Decomposition in Monte Carlo Neutron Transport Simulations using Global View Arrays. Technical report, Computer Science, University of Chicago, IJHPCA, April 2014.
- Aiman Fang and Andrew A. Chien. Applying gvr to molecular dynamics: Enabling resilience for scientific computations. Technical Report, Computer Science, University of Chicago, April 2014.
- Zachary Rubenstein, Hajime Fujita, Ziming Zheng, and Andrew Chien. Error checking and snapshot-based recovery in a preconditioned conjugate gradient solver. Technical Report, Computer Science, University of Chicago, November 2013.
- Ziming Zheng, Andrew A. Chien, and Keita Teranishi. Fault tolerance in an inner-outer solver: A gvr-enabled case study. In 11th International Meeting High Performance Computing for Computational Science VECPAR 2014, Oregon.

# Acknowledgements

- GVR Team: Hajime Fujita, Zachary Rubenstein, Aiman Fang, Nan Dun, Yan Liu (UChicago), Pavan Balaji, Pete Beckman, Kamil Iskra, (ANL), and application partners Andrew Siegel (Argonne/CESAR), Ziming Zheng (UC/Vertica), James Dinan (Intel), Guoming Lu (UESTC), Robert Schreiber (HP), Jeff Hammond (Argonne/ALCF/NWChem->Intel), Mike Heroux, Mark Hoemmen, Keita Teranishi (Sandia), Dave Richards (LLNL), Anshu Dubey, Brian Van Straalen (LBNL)
- SCR Team – some elements included in GVR system (thanks!)
- Department of Energy, Office of Science, Advanced Scientific Computing Research DE-SC0008603 and DE-AC02-06CH11357
- **For more information:** <http://gvr.cs.uchicago.edu/>