# Making HPC Systems Resilient with Parallel Objects

Esteban Meneses[12], Laxmikant V. Kalé[3]

[1]Advanced Computing Laboratory, Costa Rica High Technology Center

[2]School of Computing, Costa Rica Institute of Technology

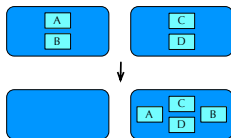[3]Department of Computer Science, University of Illinois at Urbana-Champaign

# Resilience Techniques

A taxonomy

Resilient Parallel Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects

Prevention
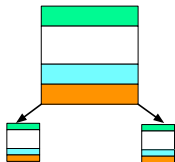
Recovery
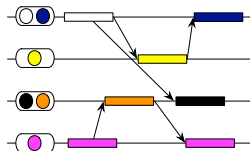
Detection

Containment

**Prevention**
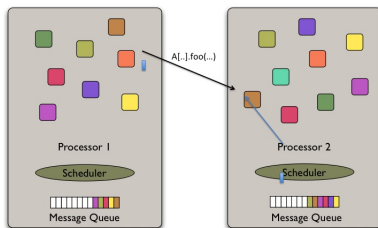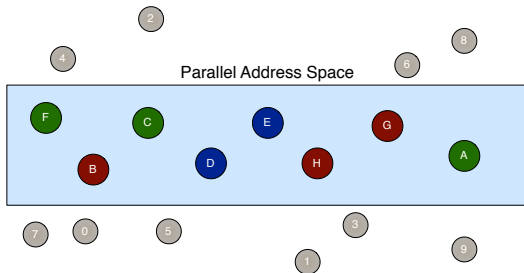
**Detection**

**Containment**

**Recovery**

Esteban Meneses, Xiang Ni, Gengbin Zheng, Celso Mendes, Laxmikant Kalé. **Using Migratable Objects to Enhance Fault Tolerance Schemes in Supercomputers**. IEEE Transactions on Parallel and Distributed Systems (TPDS), 2015.

CeNAT

# Parallel Objects

## The Charm++ programming model
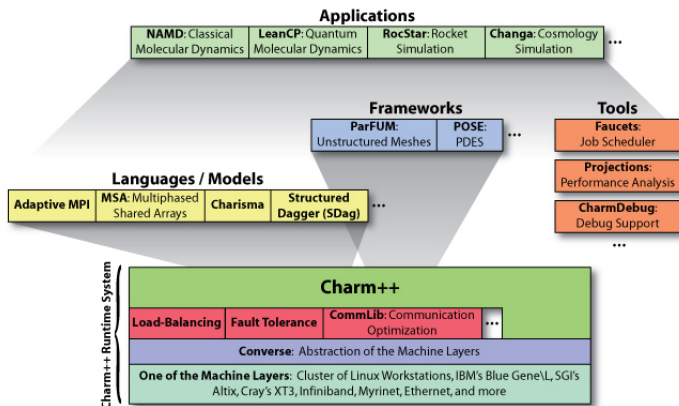
# Introspective Runtime System
## The Charm++ RTS

Resilient Parallel
Objects

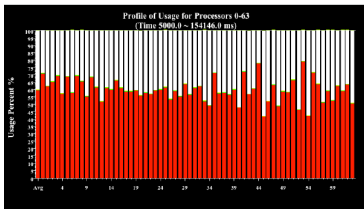Esteban Meneses,
Laxmikant V. Kalé

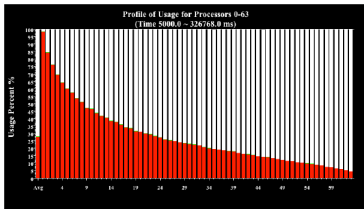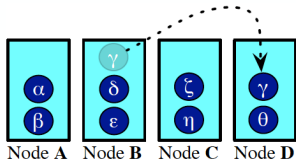Parallel Objects

Prevention

Recovery

Detection

Containment

CeNAT

# Object Migratability

Load balancing framework

Resilient Parallel
Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects

Prevention

Recovery

Detection

Containment

5

# Proactive Fault Tolerance

Evacuation of a faulty node

Resilient Parallel Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects
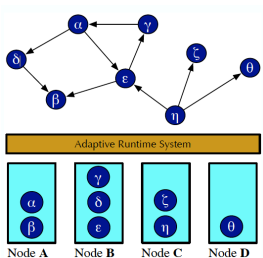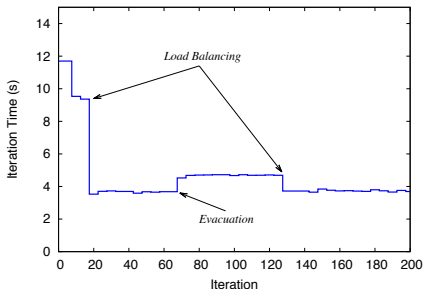
Prevention

Recovery

Detection

Containment

Sayantan Chakravorty, Celso Mendes, and Laxmikant Kalé. **Proactive Fault Tolerance in MPI Applications via Task Migration**. IEEE International Conference on High Performance Computing (HiPC), 2006.

CeNAT

# Controllable Resilience

Restraining processor temperature to reduce failure frequency

Resilient Parallel
Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects

Prevention

Recovery

Detection

Containment

Osman Sarood, Esteban Meneses, and Laxmikant Kalé. **A Cool Way of Improving the Reliability of HPC Machines**. International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2013.

CeNAT

# Checkpoint/Restart

Leveraging object migratability

Resilient Parallel Objects

Esteban Meneses, Laxmikant V. Kalé

Parallel Objects

Prevention

Recovery

Detection

Containment
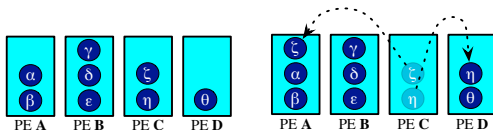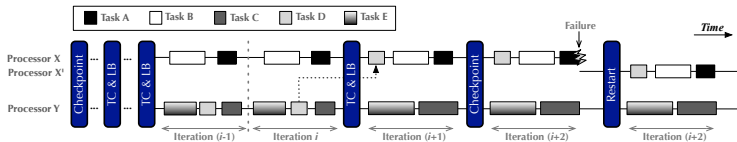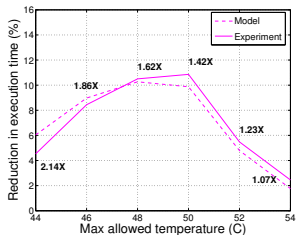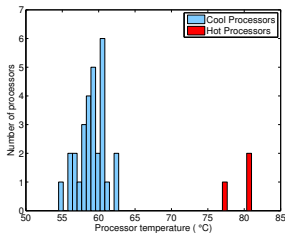


Checkpoint Time – Intrepid(leanMD)
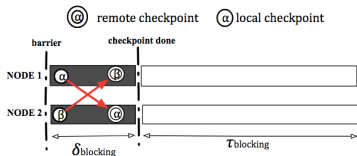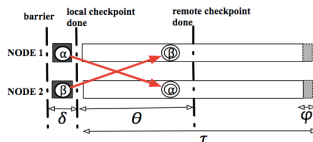


Gengbin Zheng, Xiang Ni, and Laxmikant Kalé. **A Scalable Double In-memory Checkpoint and Restart Scheme towards Exascale**. Workshop on Fault-Tolerance for HPC at Extreme Scale (FTXS), 2012.
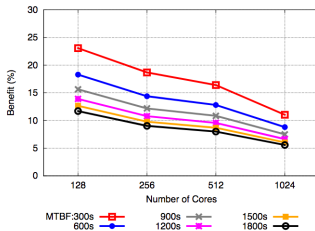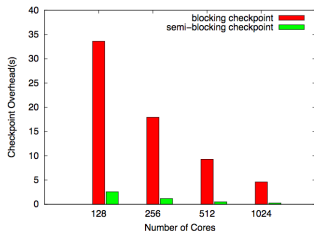
CeNAT

# Semi-blocking Checkpoint/Restart

Overlapping checkpoint and communication transmission
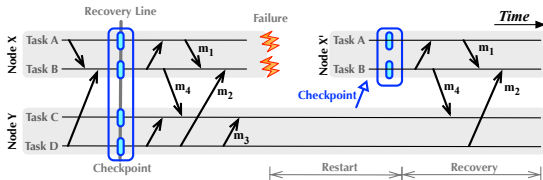
(a) Blocking Checkpoint.
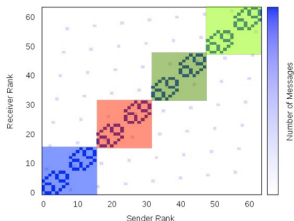


(b) Semi-Blocking Checkpoint.





Xiang Ni, Esteban Meneses, and Laxmikant Kalé. **Hiding Checkpoint Overhead in HPC Applications with a Semi-Blocking Algorithm**. IEEE International Conference on Cluster Computing (Cluster), 2012.

CeNAT

# Message Logging

Communication is stored and replayed after a failure

Resilient Parallel Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects
Prevention
**Recovery**
Detection
Containment

Esteban Meneses, Greg Bronevetsky, and Laxmikant Kalé. **Dynamic Load Balance for Optimized Message Logging in Fault Tolerant HPC Applications**. IEEE International Conference on Cluster Computing (Cluster), 2011.
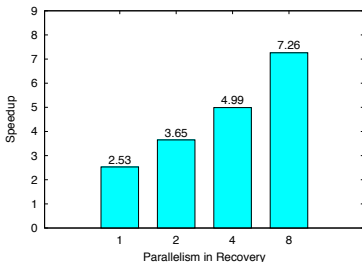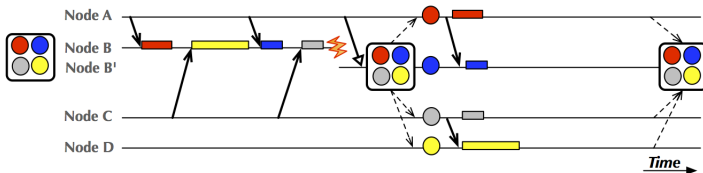
CeNAT

# Parallel Recovery

Migrate objects to speed up recovery

Resilient Parallel
Objects

Esteban Meneses,
Laxmikant V. Kalé

Parallel Objects

Prevention

Recovery

Detection

Containment
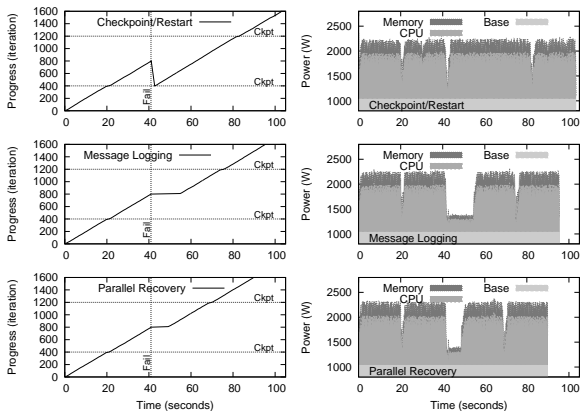
Sayantan Chakravorty and Laxmikant Kalé. **A Fault Tolerance Protocol with Fast Fault Recovery**. IEEE International Parallel and Distributed Processing Symposium (IPDPS), 2007.

11

CeNAT

# Advantages of Message Logging

Faster and greener

Resilient Parallel Objects

Esteban Meneses, Laxmikant V. Kalé

Parallel Objects

Prevention

**Recovery**

Detection

Containment

Esteban Meneses, Osman Sarood, and Laxmikant Kalé. **Energy Profile of Rollback-Recovery Strategies in High Performance Computing**. Parallel Computing (ParCo), 2014.

CeNAT

# ACR: Automatic Checkpoint/Restart

Soft and hard error protection

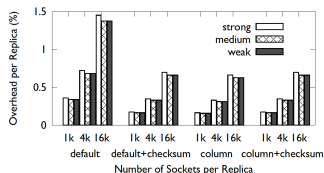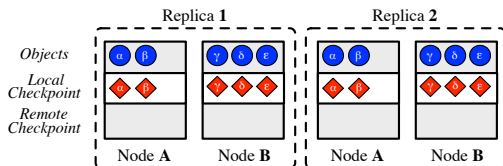Resilient Parallel Objects

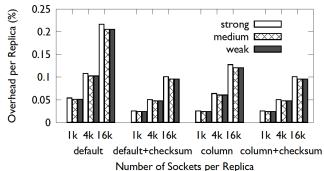Esteban Meneses, Laxmikant V. Kalé

Parallel Objects

Prevention

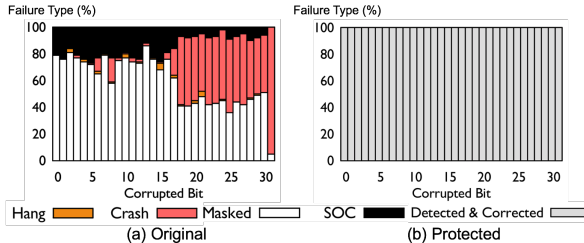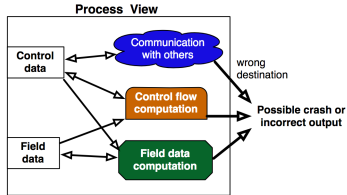Recovery

Detection

Containment

(a) Jacobi3D Charm++

(b) LeanMD

Xiang Ni, Esteban Meneses, Nikhil Jain, and Laxmikant Kalé. **ACR: Automatic Checkpoint/Restart for Soft and Hard Error Protection**. International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2013.

CeNAT

# FlipBack

Automatic targeted protection against silent data corruption

Resilient Parallel Objects

Esteban Meneses, Laxmikant V. Kalé

Parallel Objects

Prevention

Recovery

Detection

Containment

Xiang Ni and Laxmikant Kalé. **FlipBack: Automatic Targeted Protection Against Silent Data Corruption**. International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2016.

14

# Acknowledgments

# Concluding Remarks

- Parallel objects provide a fertile ground to **enhance resilience techniques**
- **Adaptivity and introspection** at the core of novel strategies
- Objects are natural **failure containment units**

**Thank You!**
**Q&A**