# Moving the Abyss:
## Database Management on Future 1000-core Processors

**NFS XPS Workshop**
**June 1-2, 2015**



**Andy Pavlo**
CMU

**Srini Devadas**
MIT

# Future Database Architectures

Develop new DBMS components for future many-core CPU architectures.

- Concurrency Control
- Storage Methods
- Logging / Recovery
- Indexing

# Non-Volatile Memory

# Many-Core Processors

Let's Talk About Storage & Recovery Methods for Non-Volatile Memory Database Systems

Staring into the Abyss: An Evaluation of Concurrency Control with One Thousand Cores

- Evaluate concurrency control schemes for transaction processing on 1000 cores.

- Custom test environment:
  - *DBx1000*
  - *MIT Graphite Simulator*

- No scheme scales due to lock thrashing, memory copying, and timestamp allocation bottlenecks.

# Many-Core Processors

Staring into the Abyss: An Evaluation of Concurrency Control with One Thousand Cores, *Proceedings of VLDB, vol. 8, iss. 3, pp. 209-220, November 2014*

# Non-Volatile Memory



Let's Talk About Storage & Recovery Methods for Non-Volatile Memory Database Systems, *Proceedings of SIGMOD, pp. 707-702, June 2015*

- Evaluate multiple methods with NVM-only storage hierarchy using a single test-bed platform.

- Three different architectures:
  - *In-place, Copy-on-Write, Log-structured Updates.*

- NVM-optimized components that use persistent pointers to reduce write-amplification.

# Next Steps

Software/Hardware Co-Designs

Push key DBMS components into hardware extensions.

# http://cmudb.io/1000cores