



OPENFABRICS
ALLIANCE

12th ANNUAL WORKSHOP 2016

INTRODUCING SANDIA* OPENSHMEM WITH MULTI-FABRIC SUPPORT USING LIBFABRIC:

Kayla Seager, Software Engineer, [Intel](#)

Ryan Grant, Senior Member of Technical Staff, [SNL](#)

Special thanks to Sayantan Sur and Jim Dinan, [Intel](#)

[April 5th, 2016]

OUTLINE

■ **Background**

- Libfabric and SHMEM co-design
- Sandia SHMEM support and software structure

■ **Scalable Design**

- Elements leveraged from Libfabric

■ **Portability**

- SOS support for many fabrics

■ **Results**

- Initial performance results



OPENFABRICS
ALLIANCE

BACKGROUND

OFIWG LIBFABRIC DESIGN

■ OFA's OFIWG Goals for Libfabric:

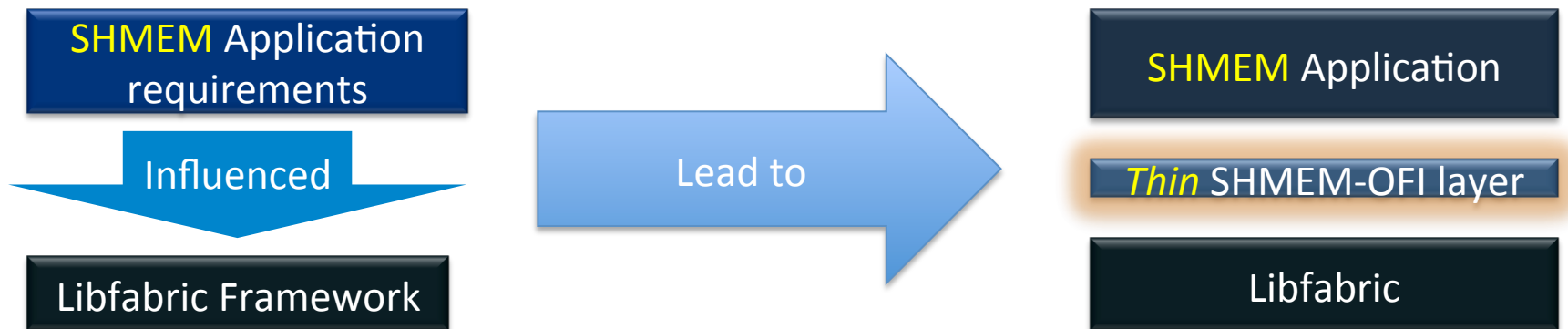
- “tight semantic map between applications and underlying fabric services” – OFIWG <https://ofiwg.github.io/libfabric/>



App Centric Design = Performance, Scalability, Portability

SANDIA-SHMEM-LIBFABRIC METHODOLOGY

Application mapping example embodies Libfabric's goal:



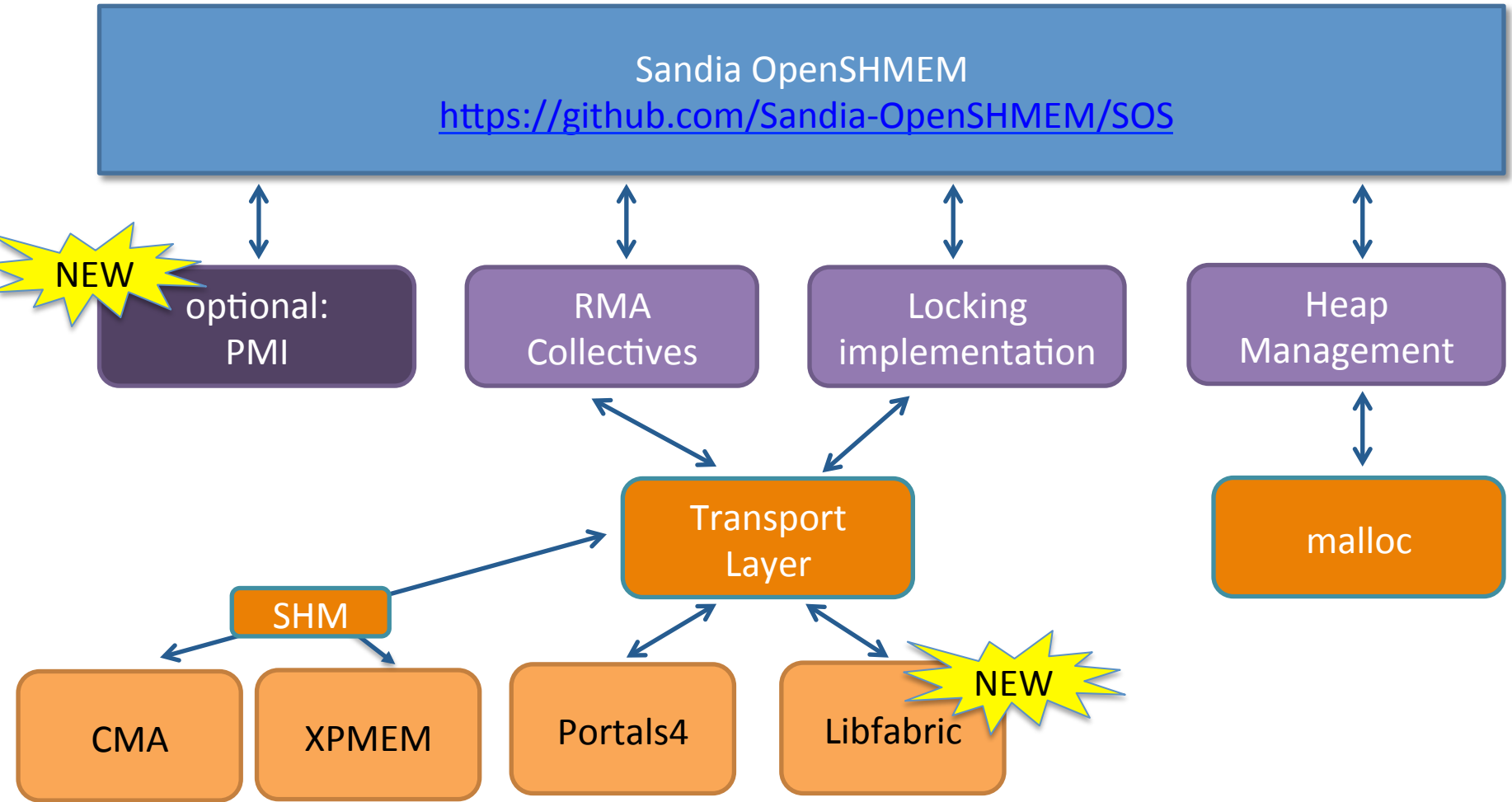
SHMEM-Libfabric = Performance, Scalability, & Portability

WHAT IS OPENSHMEM?

- **HPC Communication Programming Model API**
 - RMA & Atomic Pt-Pt
 - Distributed shared memory model (symmetric addressing)
 - Collectives
 - barrier, broadcast, reduce, all-to-all, strided all-to-all

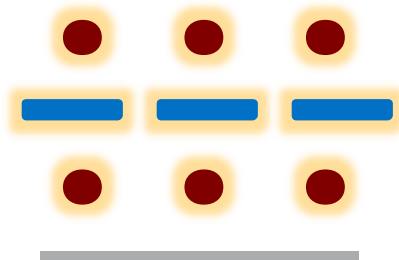


SANDIA OPENSHEMEM



SANDIA-SHMEM CODEBASE PART 2

- **OpenSHMEM 1.2.0 compliant release**
- **OpenSHMEM 1.3.0 compliant release (1st)**
 - includes true non-blocking API and all-to-all collectives
- **Open Source BSD license**
- **Vehicle to...**
 - establish interface between Libfabric and OpenSHMEM
 - Drive both specifications



PR TESTING FRAMEWORK: TRAVIS CI



■ Software Quality

- Continuous integration testing
- Rich test suite x Multiple interfaces
- Supports independent run environments
- Easy issue tracking -> bug resolution

PR TESTING FRAMEWORK: TRAVIS CI



Sandia-OpenSHMEM / SOS build passing

Current Branches Build History Pull Requests More options

✓ Pull Request #142 Merge v1.3.0 development branch #219 passed

Merge pull request #136 from jdinan/topic/finalize-errcheck
Fix handling of erroneous finalize calls

Commit 422b060
#142: Merge v1.3.0 development branch
James Dinan authored and committed

Testing for every incoming PR

Build Jobs

✓ # 219.1	</> Compiler: gcc	no environment variables set
✓ # 219.2	</> Compiler: gcc	SOS_BUILD_OPTS="--disable-mr-scalable" 8 min 45 sec
✓ # 219.3	</> Compiler: gcc	SOS_BUILD_OPTS="--enable-cma --enable-er 8 min 37 sec
✓ # 219.4	</> Compiler: gcc	SMA_BOUNCE_SIZE=0 SOS_BUILD_OPTS="--c 9 min 45 sec
✓ # 219.5	</> Compiler: gcc	SMA_BARRIER_ALGORITHM=auto SMA_BCAST 9 min 13 sec
✓ # 219.6	</> Compiler: gcc	SMA_BARRIER_ALGORITHM=linear SMA_BCAST 8 min 29 sec
✓ # 219.7	</> Compiler: gcc	SMA_BARRIER_ALGORITHM=tree SMA_BCAST 7 min 55 sec

Tests multiple environment combinations

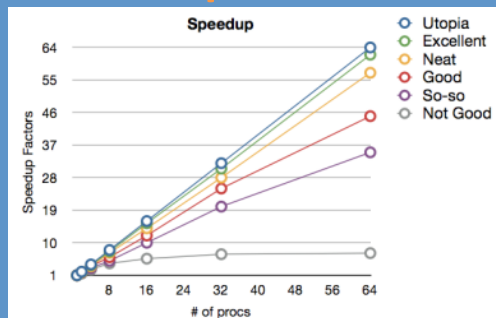


OPENFABRICS
ALLIANCE

SCALABILITY AND PERFORMANCE

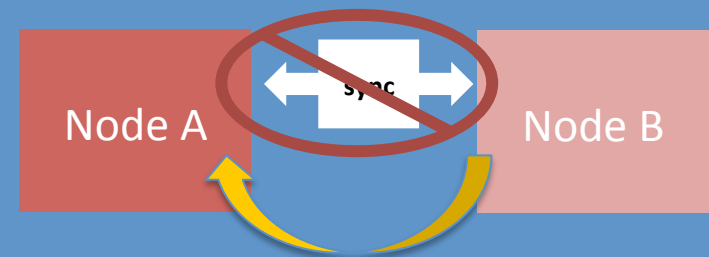
OFI WG'S SHMEM REQUIREMENTS BAKED INTO OFI USED BY SHMEM-OFI

Scalable Endpoint Enumeration



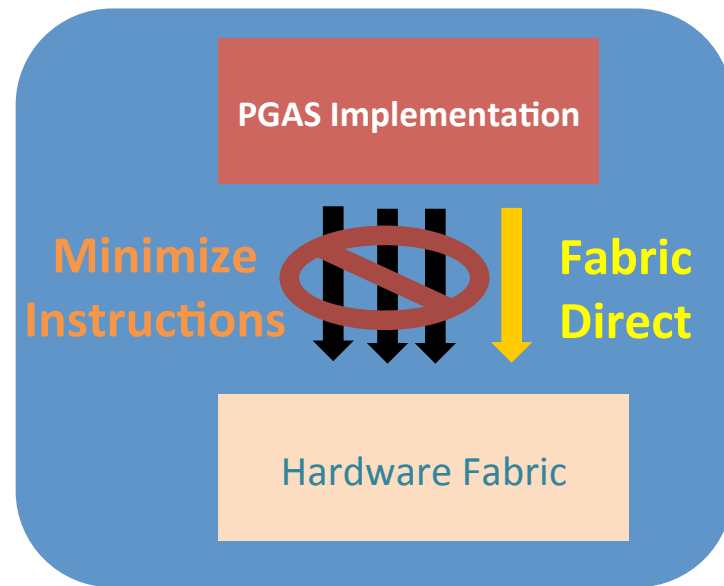
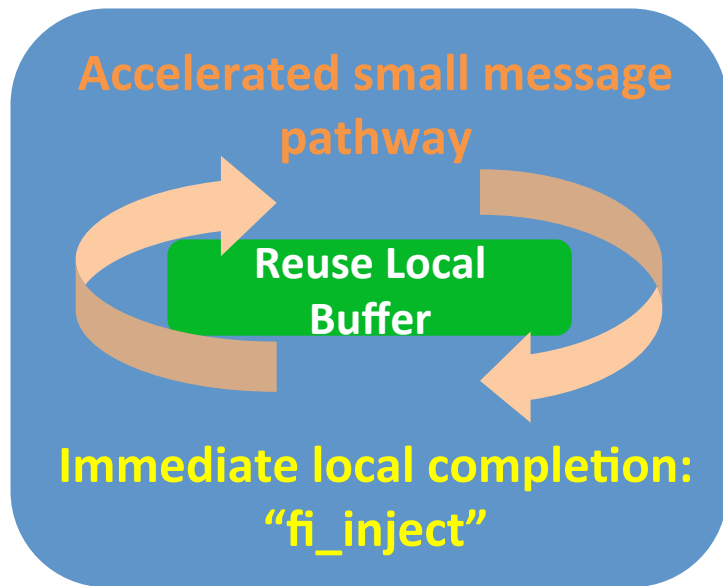
Reliable connectionless endpoints with logical addressing

Efficient Remote Completion for Put/Get

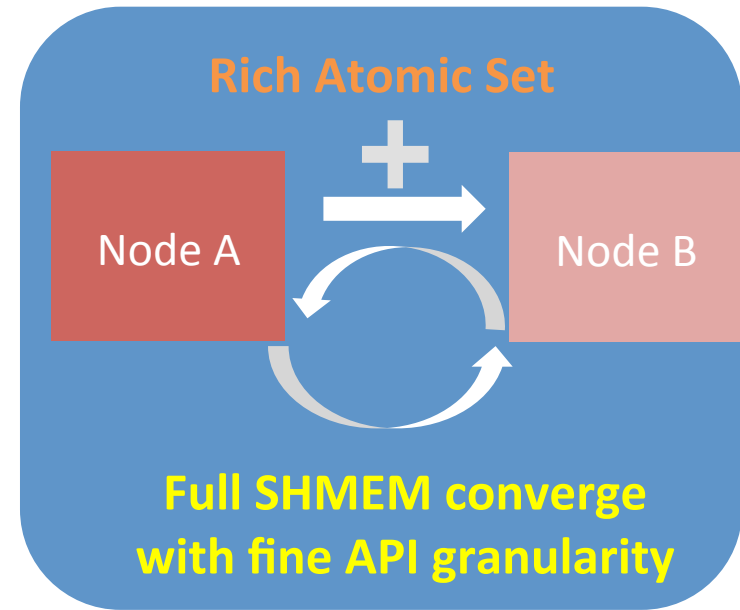
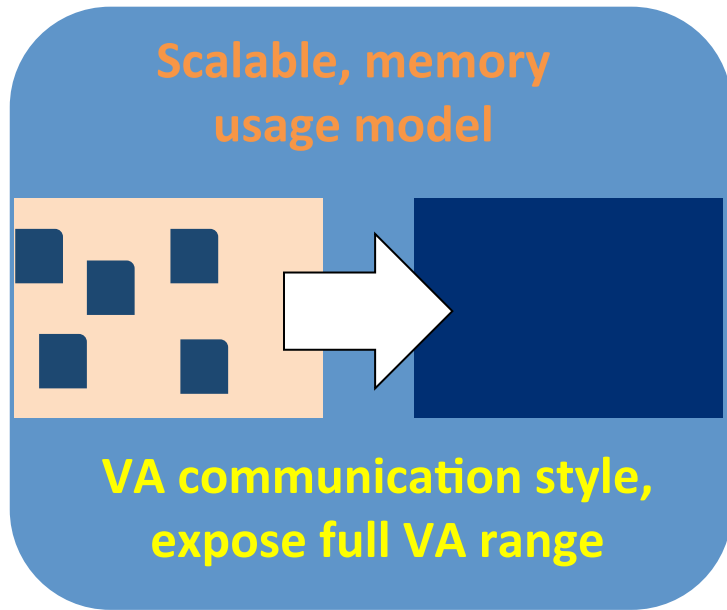


enabled counters for lightweight remote completion

OFI WG'S SHMEM REQUIREMENTS BAKED INTO OFI USED BY SHMEM-OFI



OFI WG'S SHMEM REQUIREMENTS BAKED INTO OFI USED BY SHMEM-OFI

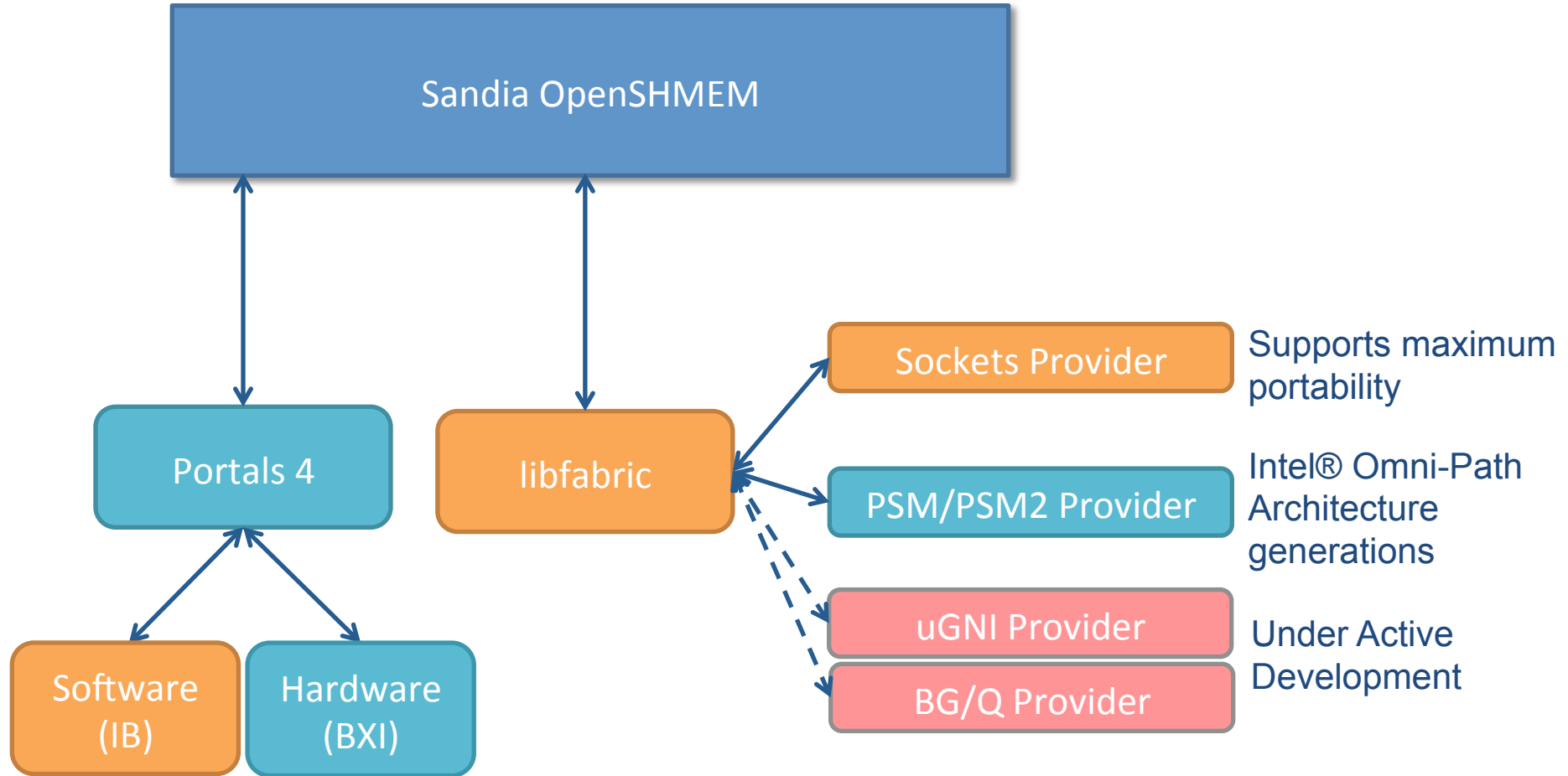




OPENFABRICS
ALLIANCE

PORTABILITY

PORTABILITY SCOPE



MEMORY MODEL BEFORE

Libfabric Semantic	Remote Virtual Addr. (OS support)
<p>MR_SCALABLE</p> <ul style="list-style-type: none">*user defined keys*Full VA Space <p>+no key exchange +single registration for full VA space +minimal footprint (don't track keys)</p>	<p>Symmetric Addressing</p> <p>(Address on Node A = Address on Node B)</p> <p>+full VA space addressing +minimal footprint (don't track addresses)</p>

MR_SCALABLE and symmetric VA = good SHMEM semantic match

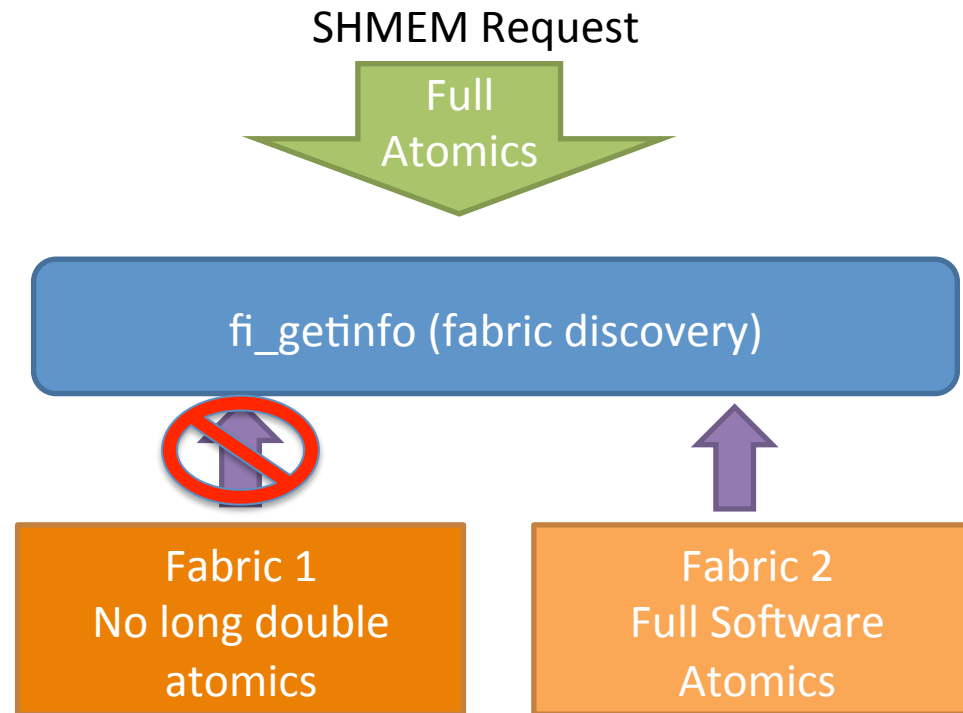
MEMORY MODEL AFTER

OFI Semantic	Remote Offset	Remote Virtual Addr. (OS support)
MR_SCALABLE *user defined keys	Offset addressing, -track and exchange offset	Symmetric addressing
MR_BASIC *fabric defined keys	Offset Addressing -track and exchange offset -track and exchange keys	Symmetric Addressing -track and exchange keys

Adding MR Basic & offset enables broader provider support

ATOMICS MODEL BEFORE

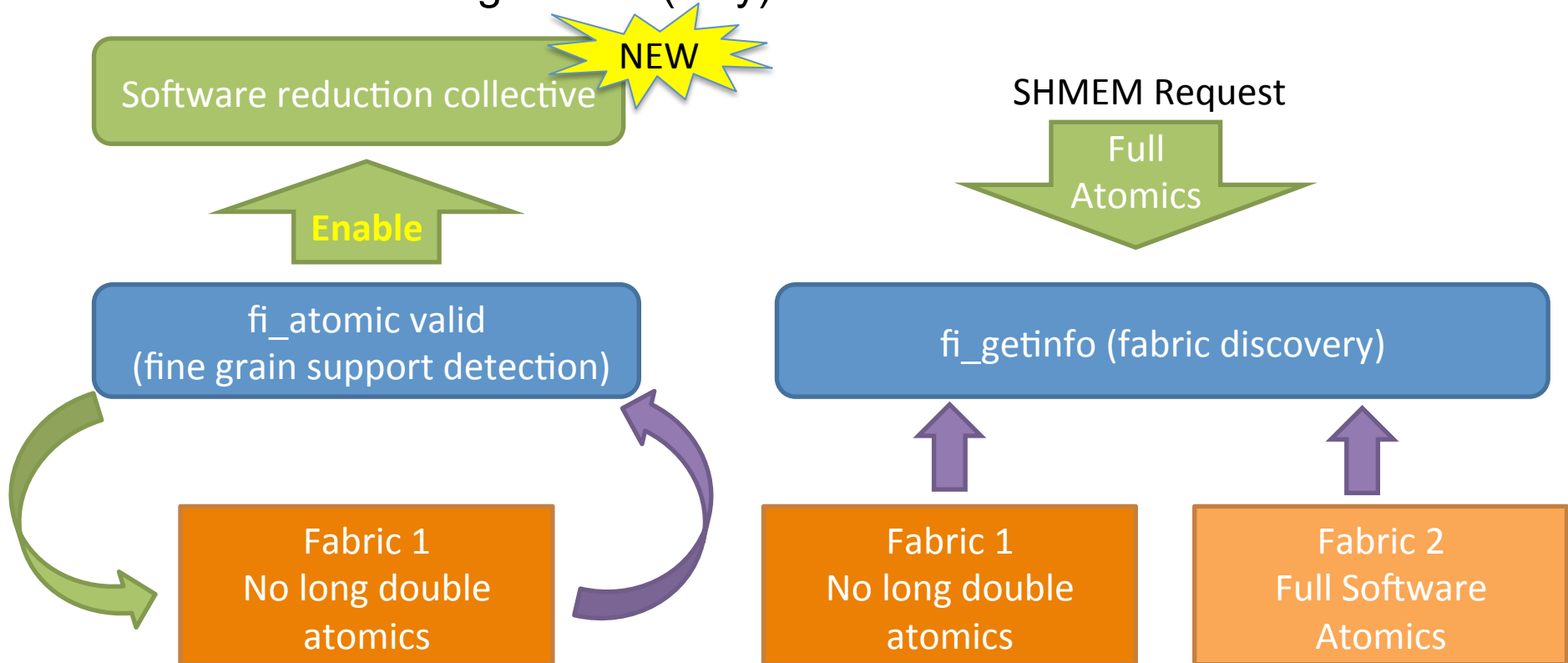
- **Require full atomic fabric support**
 - No current hardware supports long double atomics



ATOMICS MODEL AFTER

▪ Enabled “hybrid” atomic support for portability

- detect and allow long double (only) limitation



EARLY RESULTS

■ Scalability:

- Sandia-SHMEM-OFI-PSM2: scaling to **2,048 PEs** (16ppn)
 - Using ISx (integer sort) with 134,217,728 keys per PE
- Sandia-SHMEM-OFI-Sockets: scaling **512 PEs** (14ppn)
 - Full OSU* SHMEM test-suite

■ Performance:

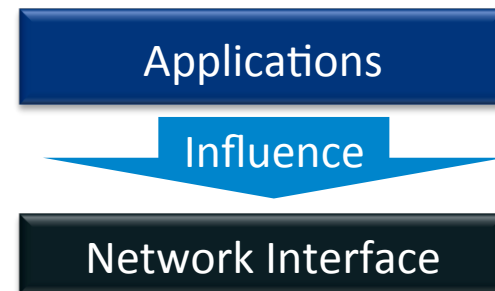
Thin SHMEM-OFI layer

- Shmem_int_p -> fi_inject_write: **16 instructions**
- Shmem_fence -> fi_cntr_wait: **11 instructions**

SANDIA-SHMEM-LIBFABRIC SUMMARY

■ Performance and Scalability

- Enabled by SHMEM requirements baked into libfabric framework
- shmem_int_p: 16 instructions
- scaled to 2,048 processes



■ Portability

- 4 different fabrics/providers
- SHMEM-libfabric changes for portability
 - Flexible memory model (compile time choice)
 - Software long double atomic support



OPENFABRICS
ALLIANCE

12th ANNUAL WORKSHOP 2016

THANK YOU

Kayla Seager, Software Engineer, **Intel**
Ryan Grant, Senior Member of Technical Staff, **SNL**
Special thanks to Sayantan Sur and Jim Dinan



Part of this work was funded through the Computational Systems and Software Environment sub-program of the Advanced Simulation and Computing Program funded by the National Nuclear Security Administration