# STATUS OF OPEN FABRICS OVER VERBS BASED FABRICS

Sayantan Sur, Intel

Presenting work done by Arun Ilango, Dmitry Gladkov, Dmitry Durnov and Sean Hefty and others in the OFIWG community

6th Annual MVAPICH User Group (MUG) 2018

# Legal Disclaimer & Optimization Notice

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.  For more complete information visit www.intel.com/benchmarks.

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

**Optimization Notice**

# Open Fabric Interfaces

**User-centric interfaces lead to innovation and adoption**

**Open Source**

**Inclusive development effort**

• App and HW developers

**User-Centric**

**Software interfaces aligned with user requirements**

• Careful requirement analysis

*Open Fabric Interfaces*

**Scalable**

**Optimized SW path to HW**

• Minimize cache and memory footprint
• Reduce instruction count
• Minimize memory accesses

**Implementation Agnostic**

**Good impedance match with multiple fabric hardware**

• InfiniBand, iWarp, RoCE, raw Ethernet, UDP offload, Omni-Path, GNI, BGQ, …
• Works on Linux, Windows and MacOS

# OFI – State of the Union

**OFI Insulates applications from wide diversity of fabrics underneath**

| Intel® MPI Library | MPICH* | Charm++* | Open MPI* | GASNet* | Sandia SHMEM* | NetIO* | Intel® MLSL# |
|---|---|---|---|---|---|---|---|

**libfabric Enabled Middleware**

OFI

*Advanced application oriented semantics*

| Tag Matching | Scalable memory registration | Triggered Operations | Remote Completion Semantics | Multi-Receive buffers | Shared Address Vectors | Unexpected Message Buffering |
|---|---|---|---|---|---|---|

| Streaming Endpoints | Reliable Datagram Endpoints |
|---|---|

| Sockets TCP, UDP | Verbs | Cisco usNIC* | Intel® OPA PSM | Cray GNI* | Mellanox* | IBM Blue Gene* | Exciting new providers in development! |
|---|---|---|---|---|---|---|---|

# Exploration

(intel)

# OFI Implementation Update

✓ **OFI Provider Infrastructure**
✓ **OFI API Exploration**
✓ **Companion APIs (Bonus!)**

**1.5 API Updates**
- RxM provider
- SOCK endpoint types
- Memory registration
- API optimizations

**2017** v1.4.0.. ..1.4.2 v1.5.0.. ..1.5.3

**2018** v1.6.0.. v1.6.1 v1.6.2 v1.7.0

**1.6 Provider Enhancements**
- PSM2 – native
- RxM performance
- SHM – shared memory support
- Persistent memory

**1.7 Predictions**
- New providers
  - RxD, multi-rail, new vendors
- SHM – xpmem support
- API enhancements

# Provider Infrastructure Updates

# RXM – Reliable Datagrams over Connections



MPI / SHMEM

OFI

RxM

RDM

OFI

MSG | MSG | MSG | MSG

MSG | MSG | MSG | MSG

RDM | RDM | RDM | RDM

Connection multiplexing

Primary path for HPC apps accessing verbs hardware

Verbs

NetworkDirect

TCP

TCP will replace sockets

Optimizes for hardware features

- **Strong MPI performance**
- **Evaluating tighter provider coupling**

(intel)

# MPI Critical Path Software overhead Analysis

# MPI Critical path software overhead



Critical SW Code Path for Ping Pong Test

Number of Instructions

| | Send | Recv |
|---|---|---|
| MVAPICH (RDMA Fast Path) | 361 | 844 |
| MVAPICH (SRQ) | 364 | 907 |
| MPICH-CH4-OFI (SRQ) | 313 | 625 |

■ Send ■ Recv

MVAPICH 2.3 (default configuration)
GCC 4.8.5, OFI master (@585919d)
-O3, -DNDEBUG

Gains in total code path primarily coming from combination MPICH-CH4 and OFI RXM provider

Instruction counts are an indirect measure help us gauge semantic fit

Ongoing optimization

• Aiming to reduce send path to about 250 instructions, and receive path to 450-480 instructions

Similar optimizations are possible in MVAPICH

(intel)

# MPI Performance Analysis - Latency

## OSU Latency (Relative Performance)



Lower is better

Y-axis: Relative Performance (0 to 1.8)
X-axis: Message Size (Bytes): 0, 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1k, 2k, 4k, 8k, 16k, 32k, 64k, 128k, 256k, 512k, 1M

Legend: MVAPICH2-2.3, IMPI 2019 Beta U1 (OFI Master)

## Platform:

Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz
Infiniband controller: Mellanox Technologies MT27700 Family
[ConnectX-4]
intel_pstate on/turbo on
RHEL 7.4
mlnx1-OFED.4.3.0.2.1.43101.x86_64

## Run details:

$ mpirun -hosts nnlmpibdw01,nnlmpibdw02 -n 2 -ppn 1 numactl
--physcpubind=7 osu_latency -i 40000

$ FI_OFI_RXM_SAR_LIMIT=8192
FI_VERBS_MR_CACHE_ENABLE=1 mpirun -hosts
nnlmpibdw01,nnlmpibdw02 -n 2 -ppn 1 numactl --
physcpubind=7 osu_latency -i 40000

(intel)

# MPI Performance Analysis – Message Rate

## OSU Messaging Rate (Relative Performance)



Higher is better

Message Size (Bytes)

— MVAPICH2-2.3    — IMPI 2019 Beta U1 (OFI Master)

## Platform:

Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz
Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4]
intel_pstate on/turbo on
RHEL 7.4
mlnx1-OFED.4.3.0.2.1.43101.x86_64

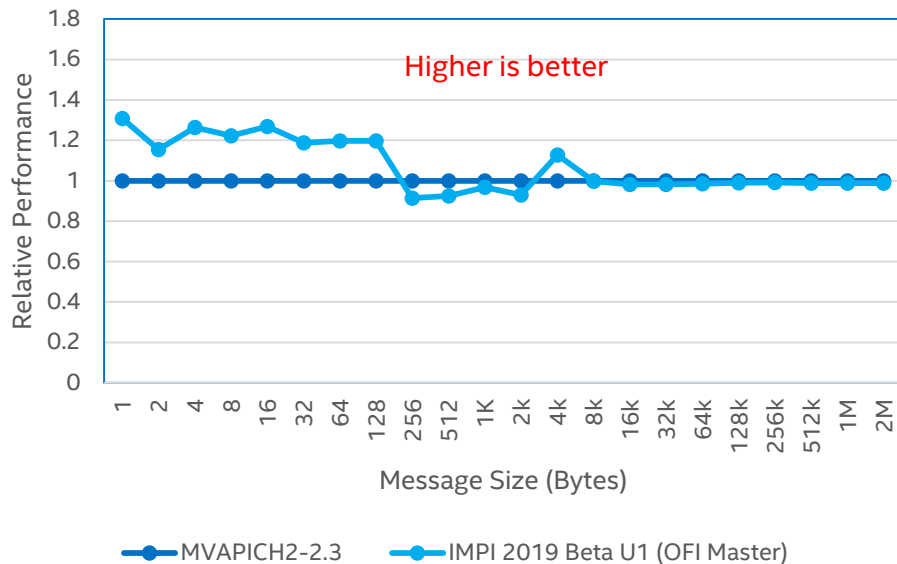## Run details:

$ mpirun -hosts nnlmpibdw01,nnlmpibdw02 -n 2 -ppn 1 numactl --physcpubind=7 osu_mbw_mr

$ FI_OFI_RXM_SAR_LIMIT=8192
FI_VERBS_MR_CACHE_ENABLE=1 mpirun -hosts nnlmpibdw01,nnlmpibdw02 -n 2 -ppn 1 numactl --physcpubind=7 osu_mbw_mr

intel

# RXD – Reliable Datagram over Unreliable Datagram



OFI

MPI / SHMEM

HPC scalability

OFI

RxD    RDM

DGRAM    DGRAM

Offload large transfers

DGRAM
RDM

DGRAM
RDM

DGRAM
RDM

DGRAM
RDM

Reliability, segmentation, and reassembly

Verbs UD

usNIC

UDP

Raw Ethernet

Other..?

Fast development path for hardware support

Extend features of simple RDM provider

- **Re-designing for performance and scalability**
- **Analyzing provider specific optimizations**

intel

# Shared Memory Provider



SHM Provider

Shared Memory Region

SMR    SMR    SMR

One-sided and two-sided transfers

CMA (cross-memory attach) for large transfers

Version
Flags
PID
Region Size
Lock

Command Queue

Response Queue

Inject Buffers

Peer Address Map

Shared memory primitives

Single command queue

xpmem support under development

# Memory Monitor and Registration Cache



**Provider**

Internal API

Driver notification, hook alloc/free, provider specific

Notification Queue

Memory Monitor Core

Monitor 'Plug-in'

subscribe

events

Get/put MRs

A generic solution is desired here

Callbacks to add/delete MRs

**Registration Cache**
LRU List
Custom Limits
Usage Stats

**MR Map**

MR   MR   MR

Tracks active usage

Merges overlapping regions

(intel)

# Performance Monitoring

Ex: Sample CPU instructions for various code paths

Linux RDPMC

## Performance Data Set

Event Data
Event Data
Event Data
Count
Count
Count
Sum
Sum
Sum

Inline performance tracking

## Performance Management Unit

**CPU** → Cycles Instructions

**Cache** → Hits Misses

**NIC** → ?

Performance 'domains'

# Hooking Provider



User

OFI

OFI Core

Hook

Core/Util Provider

Always available – release and debug builds

Zero-impact unless enabled

Intercept calls to any provider

Debugging, performance analysis, feature enhancements, testing

# Multi-rail provider



User

OFI

mRail

EP    Rail selection 'plug-in'

OFI

EP 1    EP 2

EP 1    EP 2    EP 1    EP 2

RDM    RDM

Application or admin configured

Multiple EPs, ports, NICs, fabrics

One fi_info structure per rail

Isolate rail selection algorithm

Require variable message support

TBD: recovery fallback

Increase bandwidth and message rate

Failover

Active

Standby

(intel)

# API Exploration

# Persistent Memory



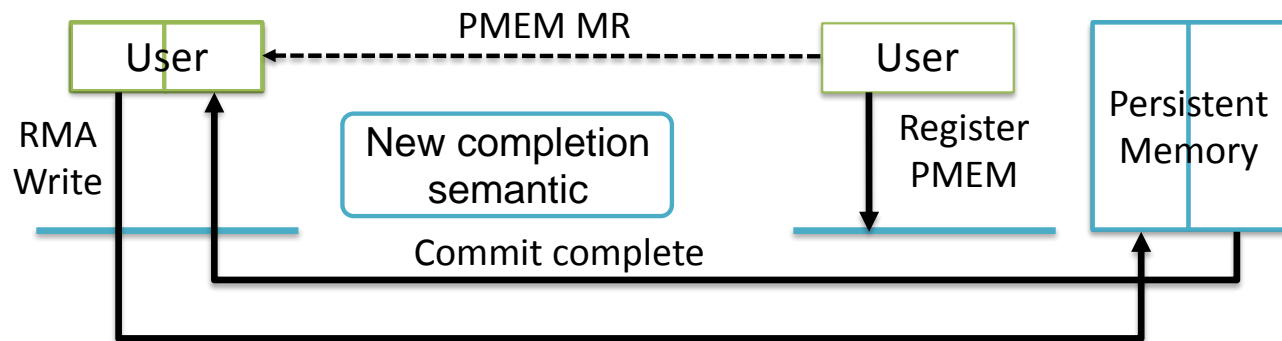Work with SNIA (Storage Networking Industry Association)
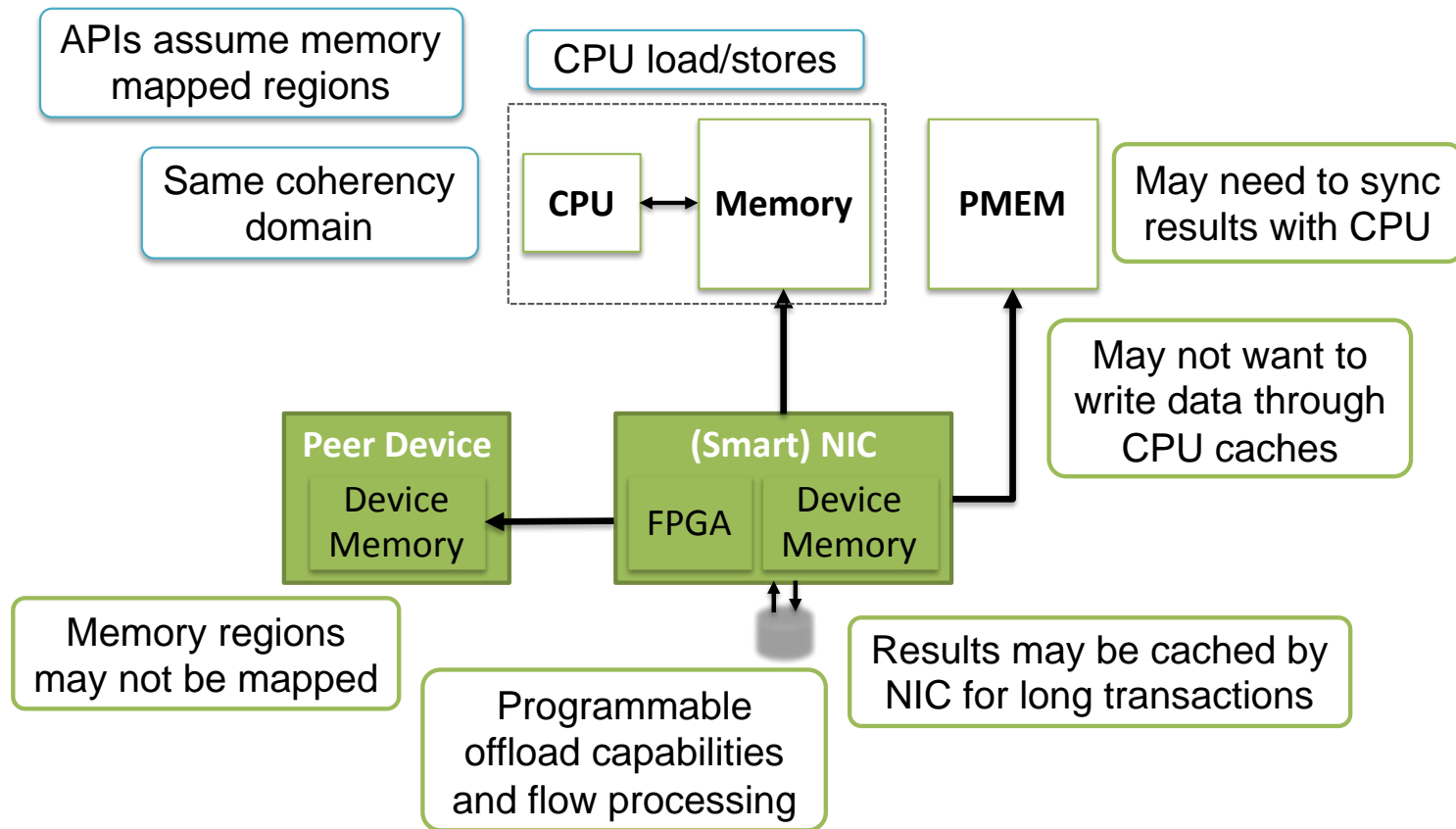
Evolve APIs to support other usage models

- **Exploration**
  - Byte addressable or object aware
  - Single or multi-transfer commit
  - Advanced operations (e.g. atomics)

- **Keep implementation agnostic**
  - Handle offload and on-load models
  - Support multi-rail
  - Minimize state footprint

# Data Domains

APIs assume memory mapped regions

Same coherency domain

CPU load/stores

CPU ◄──► **Memory**

**PMEM**

May need to sync results with CPU

May not want to write data through CPU caches

**Peer Device**
Device Memory

**(Smart) NIC**
FPGA | Device Memory

Memory regions may not be mapped

Programmable offload capabilities and flow processing

Results may be cached by NIC for long transactions

(intel)

# Variable Length Messages



- **Eager message→    ←rendezvous→**
  - RMA read or tagged message
- **MTU→    ←ack    remaining transfer→**
  - RMA write, tagged send, send
- **RTS→    ←CLS    transfer→**

# Variable Length Messages (continued)

No change at sender… *maybe*

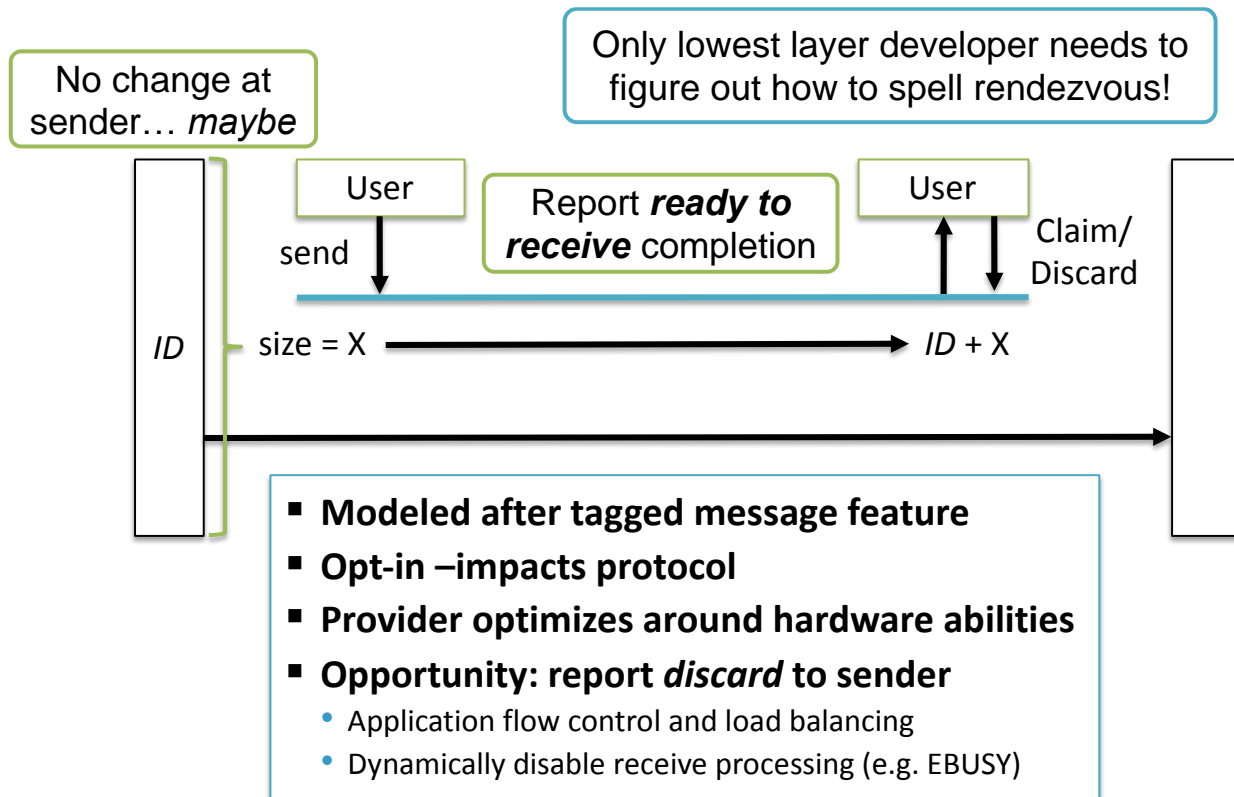Only lowest layer developer needs to figure out how to spell rendezvous!

User

send

Report **ready to receive** completion

User

Claim/ Discard

*ID*

size = X

*ID* + X

- **Modeled after tagged message feature**
- **Opt-in –impacts protocol**
- **Provider optimizes around hardware abilities**
- **Opportunity: report *discard* to sender**
  - Application flow control and load balancing
  - Dynamically disable receive processing (e.g. EBUSY)

# Companion APIs

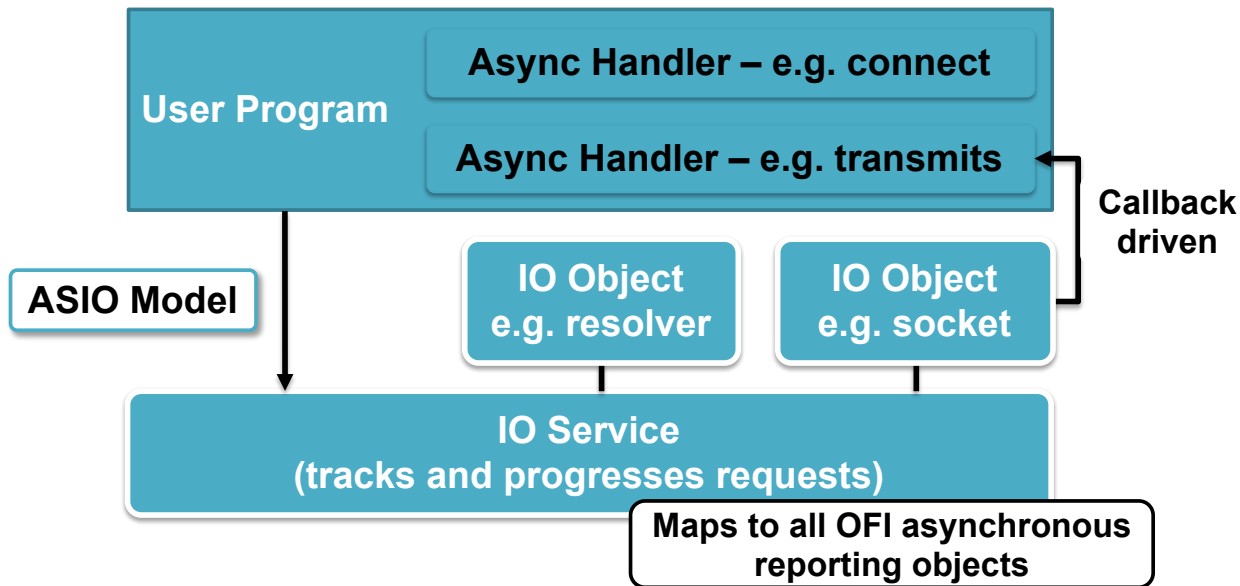(intel)

# C++ Standardization

Add support for fabrics directly to the C++ language

**Feedback from C++ community**
- Implement proposal
- Detail alternatives
- Justify extensions

**Proposal**
- Extend ASIO
- Implement over libfabric



**User Program**

**Async Handler – e.g. connect**

**Async Handler – e.g. transmits**

**Callback driven**

**ASIO Model**

**IO Object e.g. resolver**

**IO Object e.g. socket**

**IO Service (tracks and progresses requests)**

**Maps to all OFI asynchronous reporting objects**

(intel)

# Rsockets

rsockets
(librdmacm)

RDMA CM

Verbs
RC QP
UD QP

Significantly boosts performance versus sockets with HW acceleration

Increase OS & fabric portability

Pursuing OpenJDK integration

rsockets
(librsockets)

OFI

| Verbs | Omni Path | TCP | UDP | Network Direct |
|---|---|---|---|---|
| SOCK DGRAM EP | SOCK DGRAM EP | | SOCK DGRAM EP | |
| SOCK STREAM EP | SOCK STREAM EP | SOCK STREAM EP | | SOCK STREAM EP |

Always available

(intel)

# Summary

Significant software work ongoing to implement full set of OFI features on Fabric providers that lack native support

Components developed are generic and re-usable across Fabrics

Fabric vendors can implement subset of features and get access to wide OFI software ecosystem by leveraging utility components

As newer features are added to OFI, provide a pathway to quickly enable those features in older providers – applications can track latest OFI APIs

Participation in OFIWG is free, simple, no associations or boards to join

http://libfabric.org

intel