

# 12th ANNUAL WORKSHOP 2016 NEW FEATURES AND CAPABILITIES OF PSN2

Ravi Murty Intel

April 7<sup>th</sup>, 2016



## OUTLINE

- What is Performance Scaled Messaging (PSM)?
- New Features in PSM2
- Compatibility with PSM1
- **PSM Architecture**
- PSM and OFI
- PSM Open Source Strategy

# **INTEL® OMNI-PATH ARCHITECTURE OPTIMIZED HOST IMPLEMENTATION**



Host Strategy: Leverage OpenFabrics Alliance\* (OFA) Host Strategy: Leverage OpenFabrics

- Alliance\* (OFA)
- OpenFabrics Alliance compliant: Off-the-shelf application compatibility
- Provides an extensive set of mature upper layer protocols
- Integrates 4<sup>th</sup> generation proven, scalable PSM capability for HPC
- OpenFabrics Interface (OFI) API aligned with application requirements
  - AccessooppenpapurceskeybElements Fabrics Managereand Ce Via OFA
- Channels in Integrate into Linux\*ools, Fabric Manager, and GUI Channels: Integrate into Linux\*
  - **Distributions**
  - Intel® Omni-Path Architecture support • included in standard distributions



# **PSM – A QUICK OVERVIEW**

# WHAT IS PSM?

- A scalable messaging library for large clusters and millions of ranks
  - PSM is carefully designed to match the semantics needed by compute middlewares such as MPI
- End-point communication model
  - Manage flow state for reliability with support for dynamic addition of end-points
- Matched Queue (MQ) component
  - Semantically matched to the needs of MPI using tag matching
  - Provides calls for communication progress guarantees
  - MQ completion semantics (standard vs. synchronized)
- Data transfer strategies optimized for latency and bandwidth
  - Adaptively takes advantage of on-load and off-load mechanisms in HFI, transparent to library user
- Optimized inter and intra node communications encapsulated in the API
- Error handling scheme
  - Per Endpoint and Global error handler and user-defined error handlers.

## **PSM API**

- Global tag matching API with 64-bit tags
- Scale up to 64K processes per job
- MQ APIs provide point-to-point message passing between endpoints
  - e.g. psm\_mq\_send, psm\_mq\_irecv
- No "recvfrom" functionality needed by some applications

# **PSM2 API**

- 4<sup>th</sup> generation mature API with forward compatible extensions to existing PSM API
- Tag matching improvement
  - Increased tag size to 96 bits
  - Fundamentally ((stag ^ rtag) & rtagsel) == 0
  - Supports wildcards such as MPI\_ANY\_SOURCE or MPI\_ANY\_TAG using zero bits in rtagsel
- Allows for practically unlimited scalability
  - Up to 64M processes per job
- Added "recvfrom" API
  - Allows caller to specify the *source* from which to receive message
     PSM: Apprendice API with forward compatible extensions to existing
- Tag matching improvement
  - Increased tag size to 96 bits

# **PSM2 TAG MATCHING**

```
#define PSM_MQ_TAG_ELEMENTS 3
typedef struct psm2_mq_tag {
    union {
        uint32_t tag[PSM_MQ_TAG_ELEMENTS] __attribute__((aligned(16)));
        struct {
            uint32_t tag0;
            uint32_t tag1;
            uint32_t tag2;
        };
    };
} psm2_mq_tag_t;
```

- Application fills 'tag' array or 'tag0/tag1/tag2' and passes to PSM
- Both tag and tag *mask* use the same 96 bit tag type

# **PSM2 SEND SIDE FUNCTIONS**

```
psm_error_t
psm2_mq_isend2(psm2_mq_t mq, psm2_epaddr_t dest,
    uint32_t flags, psm2_mq_tag_t *stag,
    const void *buf, uint32_t len, void *context,
    psm2_mq_req_t *req);
```

```
psm_error_t
psm2_mq_send2(psm2_mq_t mq, psm2_epaddr_t dest,
    uint32_t flags, psm2_mq_tag_t *stag,
    const void *buf, uint32_t len);
```

Non-blocking and blocking send PSM API extended to support 96-bit tag

# **PSM2 RECEIVE SIDE STATUS**

```
typedef struct psm2_mq_status2 {
    psm2_epaddr_t msg_peer;
    psm2_mq_tag_t msg_tag;
    uint32_t msg_length;
    uint32_t nbytes;
    psm_error_t error_code;
    void *context;
} psm2 mg_status2_t;
```

- psm2\_mq\_status2\_t is derived from psm2\_mq\_status\_t
  - msg\_tag field changed from 64-bits to the new 96-bit tag
- Additionally, a new field msg\_peer added
  - Provide the source of the message
  - Eliminates application burden to maintain relationship between message and source of message

# **PSM2 RECEIVE SIDE FUNCTIONS**

psm\_error\_t

psm2\_error\_t

psm\_mq\_improbe2(psm2\_mq\_t mq, psm2\_epaddr\_t src, psm2\_mq\_tag\_t \*rtag, psm2\_mq\_tag\_t \*rtagsel, psm\_mq\_req\_t \*req, psm2\_mq\_status2\_t \*status);

- Functions psm2\_mq\_irecv2 and psm2\_mq\_iprobe2
  - Tag and tag mask are 96-bit wide. Also src argument allows receiver to specify the source from which to receive the message
- psm2\_mq\_improbe2 is a new function for MPI3 support

# **PSM2 PROGRESSION API**

psm\_error\_t psm2\_mq\_ipeek2(psm2\_mq\_t mq, psm\_mq\_req\_t \*req, psm2\_mq\_status2\_t \*status);

psm\_error\_t psm2\_mq\_wait2(psm2\_mq\_req\_t \*request, psm2\_mq\_status2\_t
\*status);

psm\_error\_t psm2\_mq\_test2(psm2\_mq\_req\_t \*request, psm2\_mq\_status2\_t
\*status);

• Change in status argument from previous PSM API



# PSM1/PSM2 FORWARD COMPATIBILITY AND CO-EXISTENCE

# **PSM1/PSM2 FORWARD COMPATIBILITY**

- Both PSM and PSM2 are supported on Intel® OPA
- On Intel® OPA PSM1 64-bit tag is zeroextended to 96-bits
- A PSM function only returns the first 64 bits of tag
- A PSM2 function returns all 96 bits of tag
- For a PSM2 receive function, if 'src' argument is NULL, PSM uses only the 96b *tag* to match a message from any source
- status2 will always return the message source epaddr 'msg\_src', independent of the 'src' argument in receiving function



No software application changes

# **PSM1/PSM2 CO-EXISTENCE**

- Co-existence between Intel® True Scale and Intel® Omni-Path Architecture
  - PSM2 renamed the symbol names but ...
- Forward compatibility is provided by psm2-compat RPM:
  - Provides binary compatibility for PSM1 based
  - Set LD\_LIBRARY\_PATH=/usr/lib64/psm2-compat/ prior to application launch.
  - /usr/lib64/psm2-compat/libpsm\_infinipath.so.1, is a very thin wrapper exporting PSM APIs, but uses PSM2 library. Minimal performance impact from native PSM2



# PSM2 OPTIMIZED ON INTEL® OMNI-PATH ARCHITECTURE

# **INTEL® OPA: SEND OPTIONS**

#### 2 Modes of Sending data Independent of Receive mode



#### Send DMA (SDMA)

- Optimizes Bandwidth for Large messages
- 16 SDMA Engines for CPU Offload

#### Programmed I/O (PIO)

 Optimizes Latency and Message Rate for small messages

# **INTEL® OPA: RECEIVE OPTIONS**

#### Eager-Receive

- Received data buffers copied to Application Buffer
- No handshake needed

#### Direct Data placement in Application Buffer

 Data placed directly into Application Memory

#### 2 Modes of Receiving data Independent of Send mode



# INTEL® OPA PROVIDES EFFICIENT DATA TRANSFER MECHANISMS

Message Size	Send Side	Receive Side
Up to 8KB	PIO Send	Eager Receive
> 8KB and < 64KB	SDMA	Eager Receive
64KB or more	SDMA	Expected Receive

 Most efficient data movement method automatically chosen based on message size

# **MPI PERFORMANCE ON INTEL® OPA**



Tests performed on Intel® Xeon® Processor E5-2697v3 dual-socket servers with 2133 MHz DDR4 memory. Turbo mode enabled and hyper-threading disabled. Ohio State Micro Benchmarks v. 4.4.1. Intel OPA: Open MPI 1.10.0 with PSM2. Intel Corporation Device 24f0 – Series 100 HFI ASIC. OPA Switch: Series 100 Edge Switch – 48 port. IOU Non-posted Prefetch disabled in BIOS. Intel® True Scale: Open MPI. QLG-QLE-7342(A), 288 port True Scale switch. 1. osu\_latency 8 B message. 2. osu\_bw 1 MB message. 3. osu\_mbw\_mr, 8 B message (uni-directional), 28 MPI rank pairs





# **PSM PROVIDERS FOR OFI**

- In addition to a standalone library, OFI supports providers for both versions of PSM
- OFI provider for PSM and PSM2: <a href="https://github.com/ofiwg/libfabric">https://github.com/ofiwg/libfabric</a>
  - The psm provider runs over PSM1.x interface supported by Intel® True Scale
  - The psm2 provider runs over PSM2.x interface supported by Intel® OPA
  - The psm2 provider supports a richer set of capabilities due to the additional functionality provided by the PSM2 library described earlier.
- The PSM2 provider allows creation of multiple OFI endpoints without limitations
  - PSM provider has limitations on endpoint capability settings when multiple endpoints are opened.
- The fi\_psm and fi\_psm2 man pages provide additional information including capabilities supported by these providers
  - <u>https://github.com/ofiwg/libfabric/wiki/Provider-Feature-Matrix-v1.3.0</u>
- Status: Both providers are in stable state. On-going work will be mainly bug fixes/code refactoring

# **PSM IS OPEN SOURCE**

- PSM1 has been part of OFED for a very long time.
  - PSM2 is still in the works
- PSM2 sources: <u>https://github.com/01org/opa-psm2</u>
- Distros: PSM2 is accepted to go into RHEL 7.2 and SLES SP2
- Documentation: Additional details on PSM2 and its API are available: <u>http://www.intel.com/content/dam/support/us/en/documents/network/omni-adptr/sb/Intel\_PSM2\_PG\_H76473\_v1\_0.pdf</u>

# CONCLUSIONS

- PSM2 is an optimized user-space library designed to meet the requirements of MPI
- MQ component of PSM2 supports tag matching with 96-bit tags and "recvfrom" functionality, among other improvements
- Existing applications written to the PSM1 API should run on Intel® OPA via the compat library
- PSM support two send mechanism and two receive mechanisms that allow for very efficient messaging from a latency and BW perspective.
- OFI has providers for PSM1 and PSM2
- Finally, PSM2 is open source

## **LEGAL DISCLAIMERS**

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest Intel product specifications and roadmaps.

Any forecasts of goods and services needed for Intel's operations are provided for discussion purposes only. Intel will have no liability to make any purchase in connection with forecasts published in this document.

Cost reduction scenarios described are intended as examples of how a given Intel- based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling, and provided to you for informational purposes. Any differences in your system hardware, software or configuration may affect your actual performance.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at intel.com.

Intel processor numbers are not a measure of performance. Processor numbers differentiate features within each processor family, not across different processor families: Go to: <u>Learn About Intel® Processor Numbers</u>

The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Copies of documents which have an order number and are referenced in this document, or other Intel literature, may be obtained by calling 1-800-548-4725, or go to: <u>http://www.intel.com/design/literature.htm</u>

Intel, Intel Xeon, Intel Xeon Phi<sup>™</sup> are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries.

Copyright © 2016, Intel Corporation

# **OPTIMIZATION NOTICE**

#### **Optimization Notice**

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel.

Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804



12<sup>th</sup> ANNUAL WORKSHOP 2016

**THANK YOU** 

Ravi Murty

Intel