



OPENFABRICS  
ALLIANCE

13<sup>th</sup> ANNUAL WORKSHOP 2017

# ADVANCING OPEN FABRICS INTERFACES

Sean Hefty, OFIWG Co-Chair  
Intel Corporation

March, 2017

# MOVING LIBFABRIC FORWARD

## Enhancements

- Deferred and new features

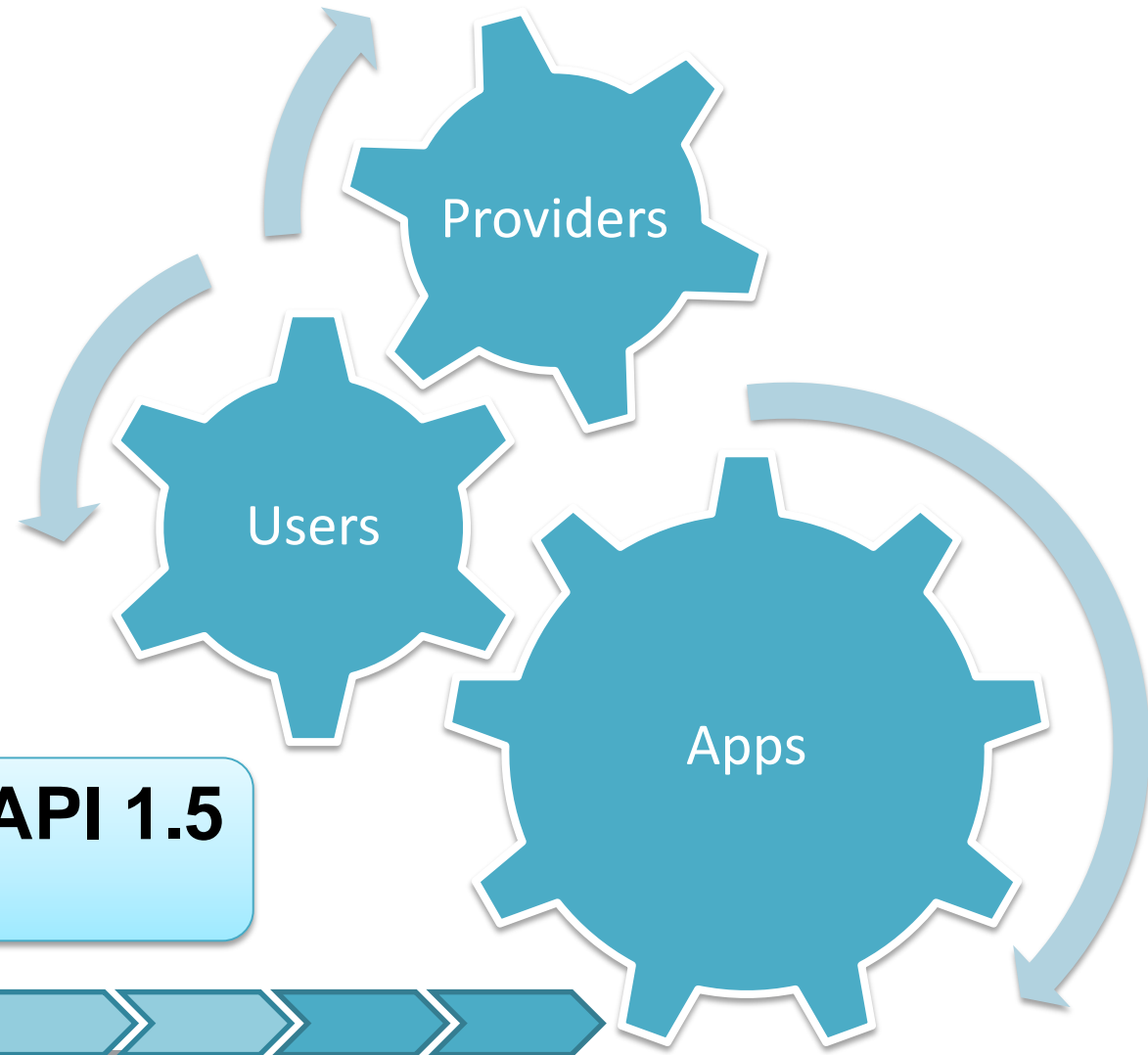
## API adjustments

- Developer feedback

## Performance optimizations

- Provider implementation

**Many features require migration to API 1.5**  
**Existing apps work unmodified**



**New!**

# OFI DEVELOPER GUIDE



*What a novel concept!*

Broad coverage of libfabric

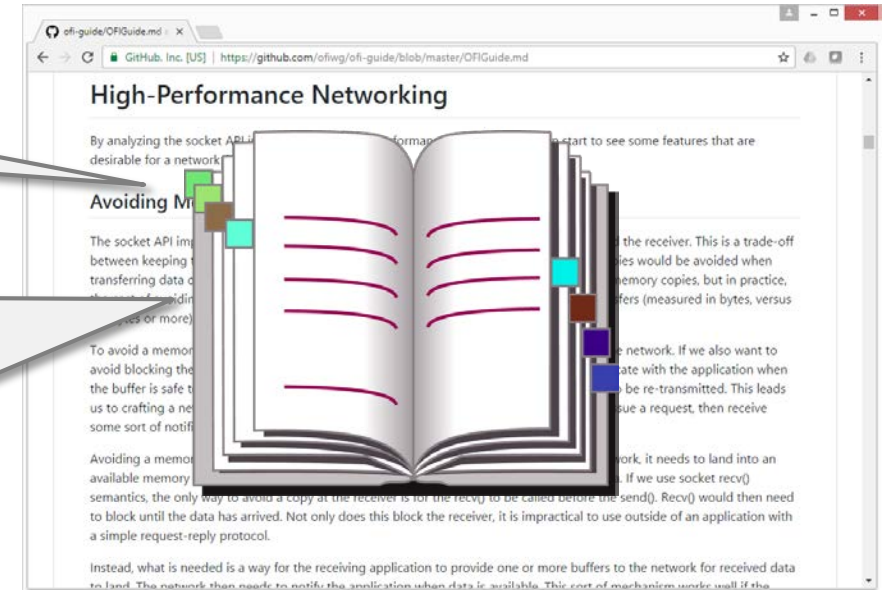
Design motivations, architecture, usage guidance

Assumes minimal background  
Reviews sockets API for context

Analyzes HPC requirements and their impact on an API

Memory copies, network buffering, asynchronous operations, interrupts, direct HW access, ...

**Lower the barrier to adoption**  
**Many concepts are generic**

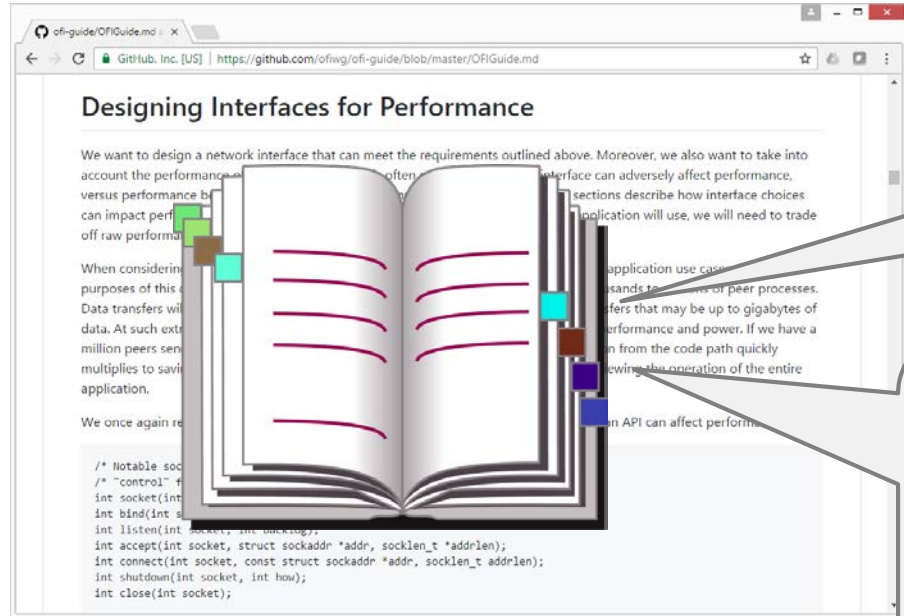


**New!**

# OFI DEVELOPER GUIDE



*Owl read  
this later*



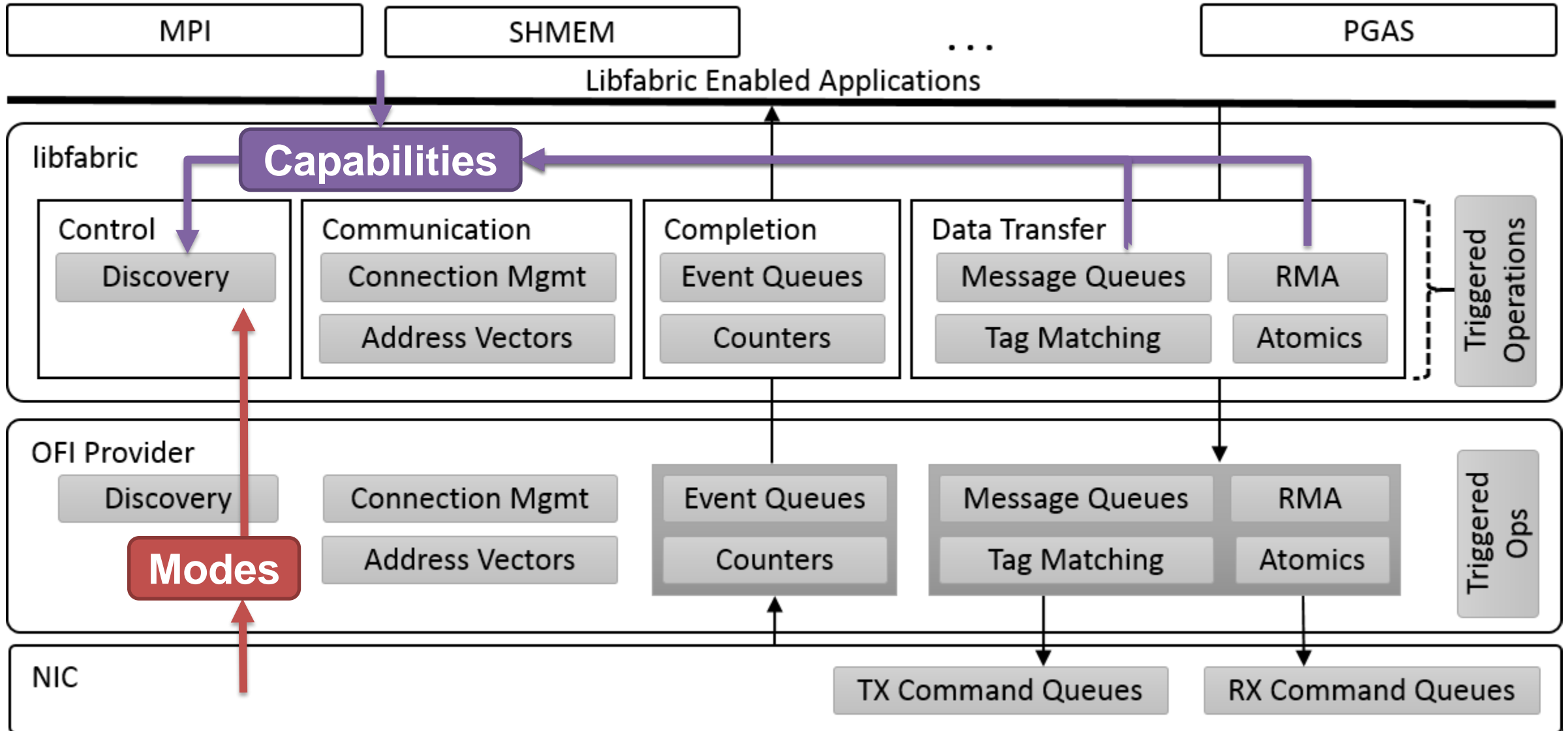
Examines API trade-offs  
Call setup, branches, memory footprint, ...

Defines communication primitives  
Command queues, memory buffers, completion notification, progress models, data and message ordering, ...

***Then it goes into detail about OFI***  
**<https://github.com/ofiwg/ofi-guide>**

**Old!**

# ARCHITECTURE



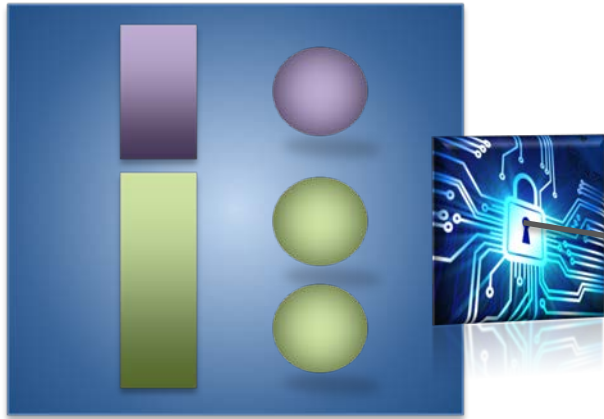


*Job Keys*

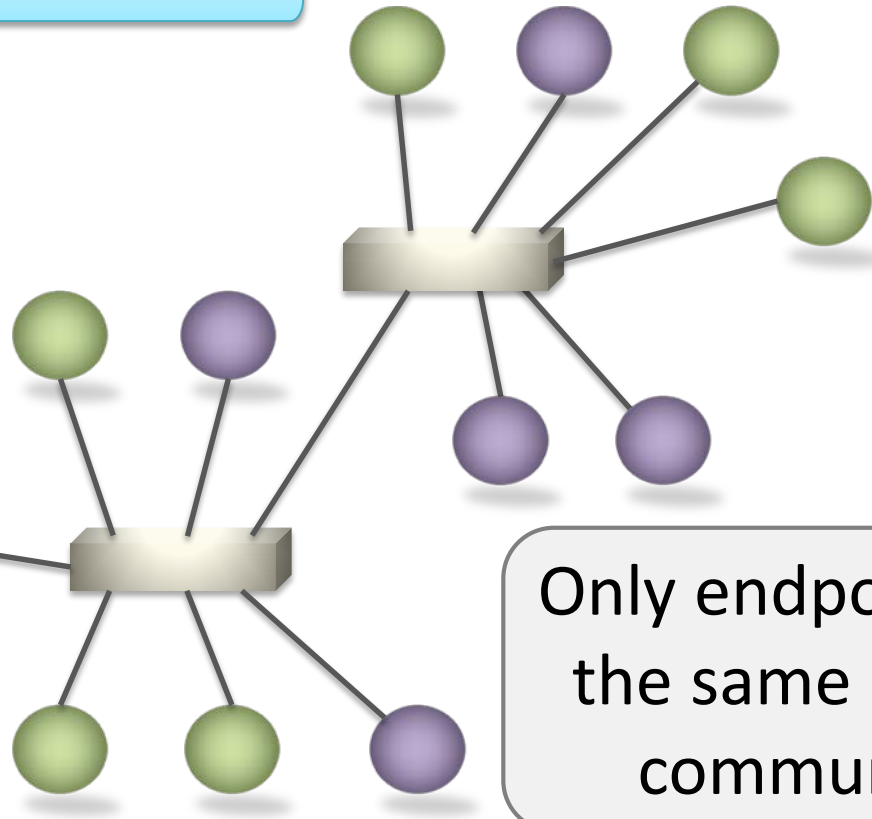
# AUTHORIZATION KEYS

**Isolate traffic between virtual sets of peers**

Associated with endpoints and memory regions



Memory regions are only accessible through endpoints with matching keys



Only endpoints with the same key may communicate

***At Long Last***

# MULTICAST

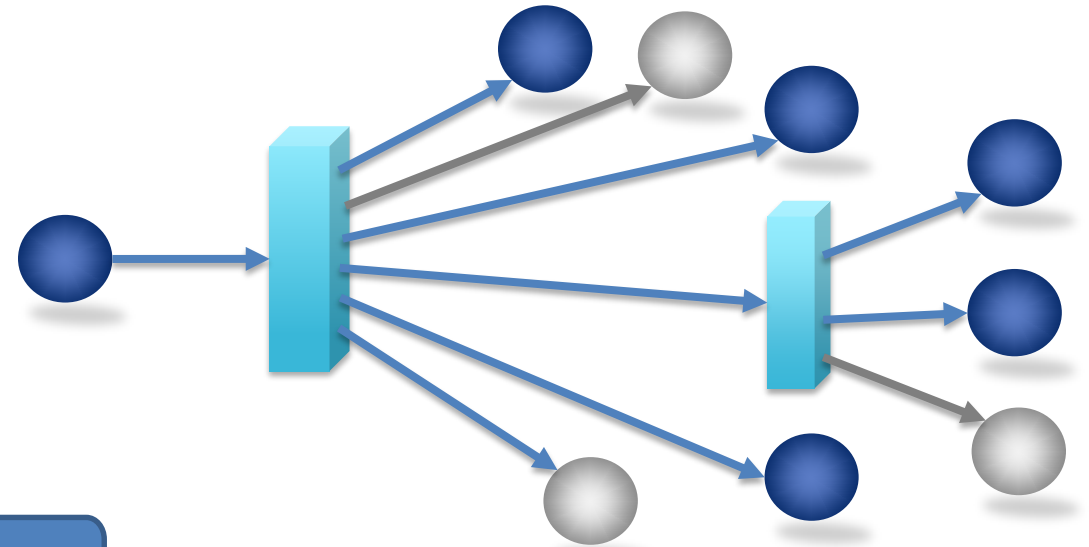
**Datagram endpoint capability**

**New fi\_join call to connect to group**

- Supports send-only joins
- Returns fabric address (fi\_addr\_t) for group

**Uses existing data transfer APIs**

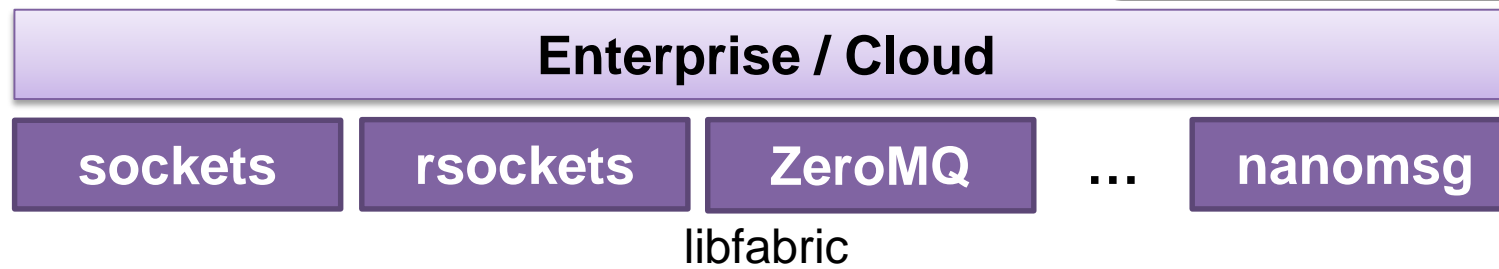
- Only message queue interface currently defined
- Multicast sends require use of FI\_MULTICAST flag
  - Distinguish between unicast and multicast address
- Received messages also marked with FI\_MULTICAST



# NEW ENDPOINT TYPES

Apps/middleware designed around socket semantics

Synchronous completions  
No completion queue  
Application owns buffers



Enable provider specific optimizations  
RMA, circular receive buffers  
Minimize data copies, use offload

Associated with file descriptor  
Support for select/poll





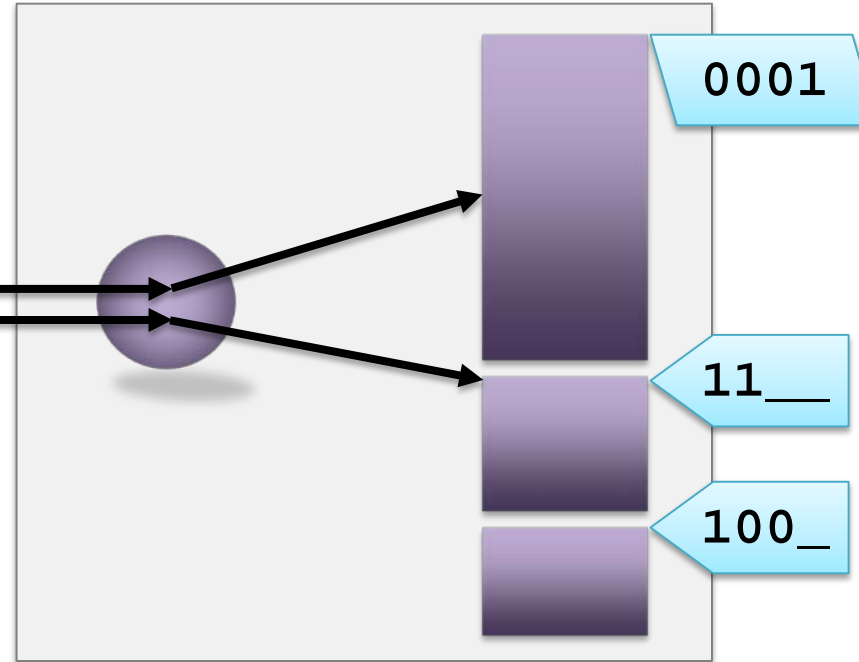
**Free:  
100% More**

# EXPANDING ATOMICS

*RMA key (not address)  
identifies the buffer*

**(RMA) Atomic**

**Tagged Atomic**



**FI\_TAGGED op flag**  
- Replace key with tag

**fi\_atomic\_query()**  
- Domain level call  
- Also reports datatype size

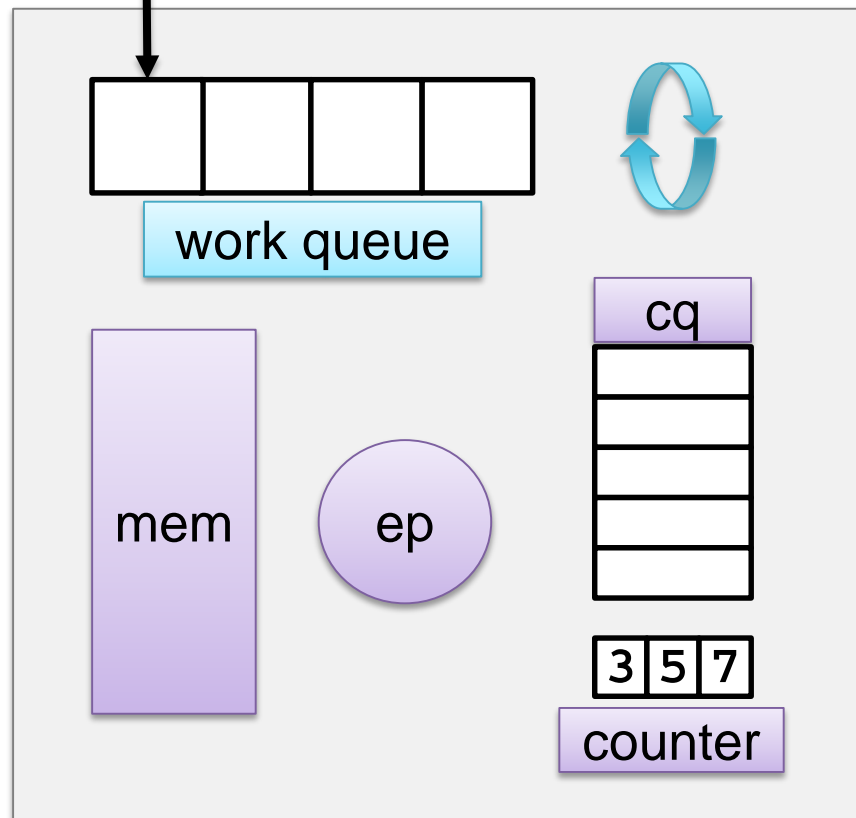
**Experimental**

# DEFERRED WORK QUEUES

Generalizes  
triggered operations

work request

**Conceptual  
model only**



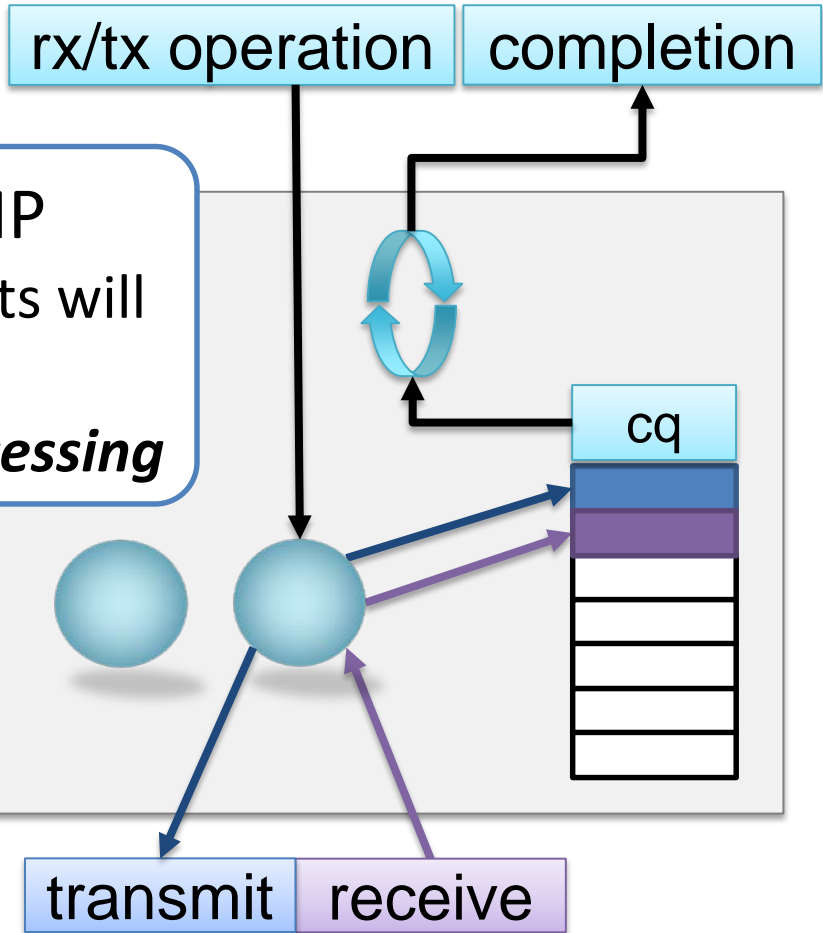
Conceptual domain  
level work queue

Allows (theoretical)  
associations between  
any domain level object  
Definition is limited

Primitives to develop and  
evaluate collective APIs  
e.g. directed acyclic graphs

**Reducing  
Overhead**

# OPTIMIZING COMPLETIONS



**FI\_RESTRICTED\_COMP**  
Only similar endpoints will share CQs / counters  
***Eliminates post-processing***

**New mode bits**

**FI\_NOTIFY\_FLAGS\_ONLY**

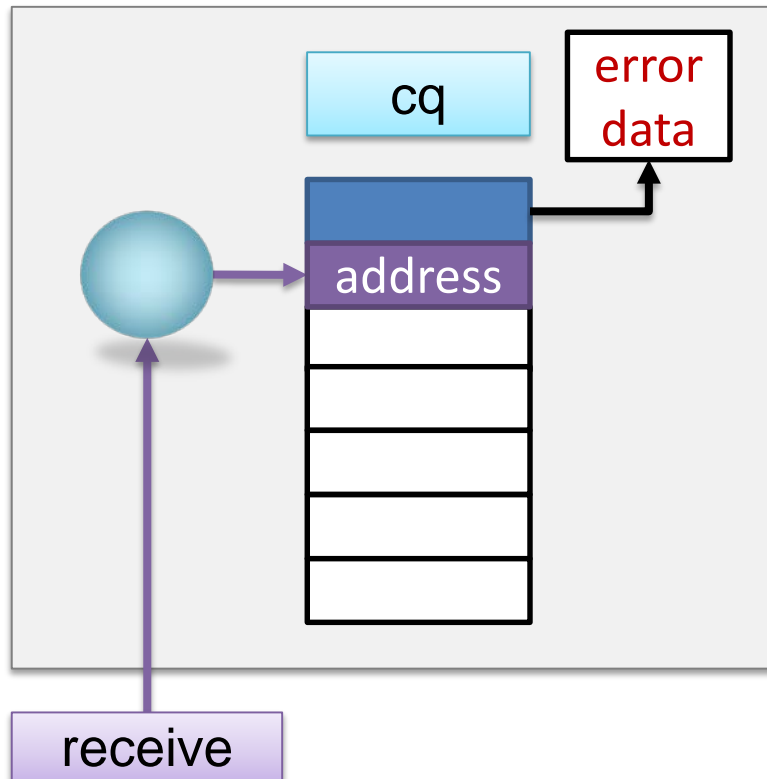
- Drops operation flag
- Retains notification flags:
  - Remote CQ data
  - Multi-recv buffer freed

***Avoids op save/restore***

**App  
Request**

# COMPLETION DATA

**Multi-  
Threading**



## Reporting errors

- Provider specific data reported for debugging
  - Data is opaque to application
  - Storage managed by provider
- Report maximum size of error data
- Application provides buffer

***Improves multi-threaded completion processing***

## FI\_SOURCE\_ERR capability

- Validate source address against local AV
- Report raw fabric address if not found
  - EADDRNOTAVAIL completion error

***Pushes 'connect' address exchange to provider***



## New attributes

- Fabric: api version
- Domain: counter limits, memory region iov limits

## String based addressing

- FI\_ADDR\_STR address format
- Generic string format for any address
  - Application does not *need* to interpret
  - Well-known format for sockaddr
- Pass directly to `fi_getinfo()` or AV insertion

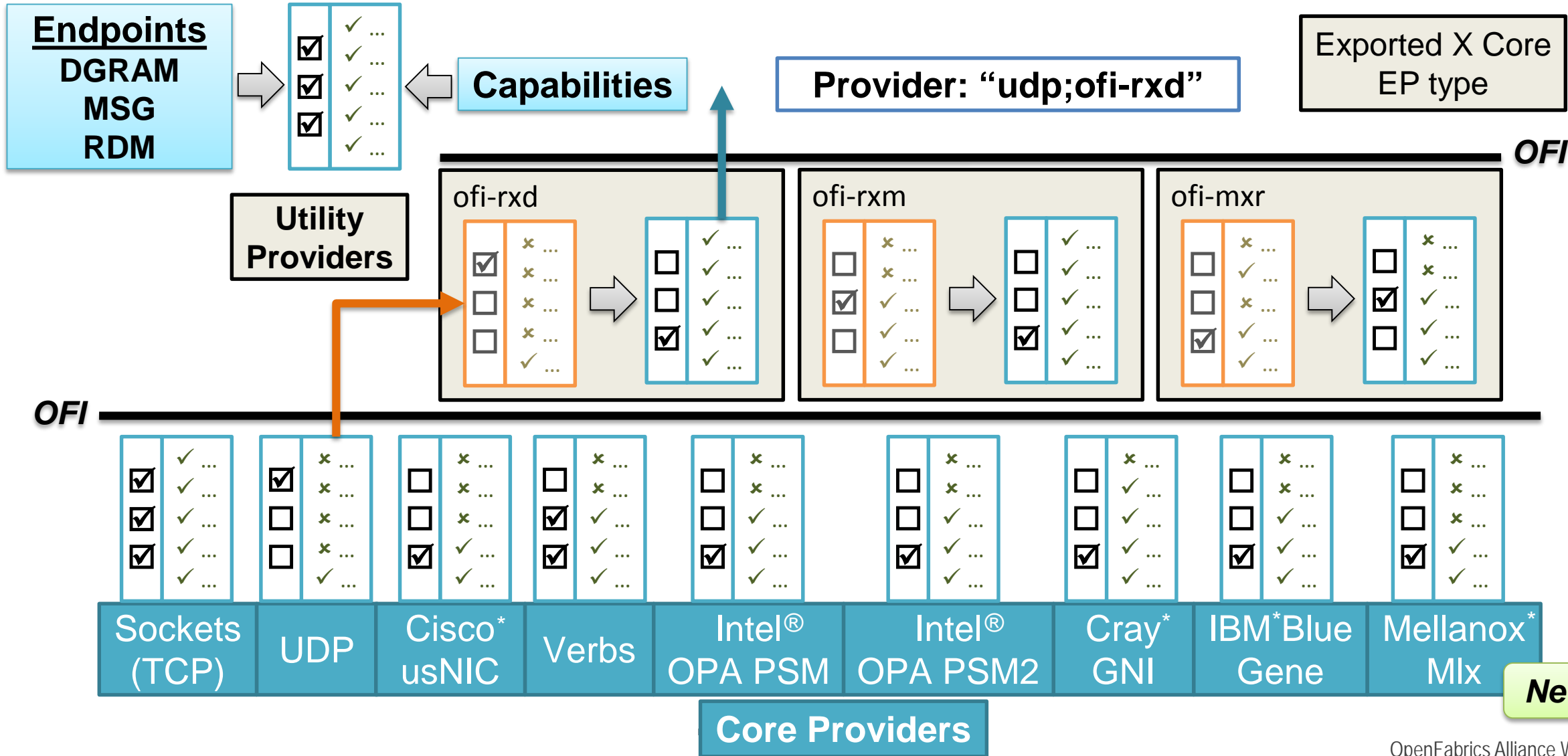
## Communication scope

- Local: within the local node
  - Shared memory, most NICs
- Remote: with external nodes
  - Any NIC
- Evading more refined definitions



**Simplify  
Development**

# PROVIDER UPDATES



**Exposing  
the Pain**

# MEMORY REGISTRATION

**FI\_MR\_BASIC**  
Traditional RDMA fabrics



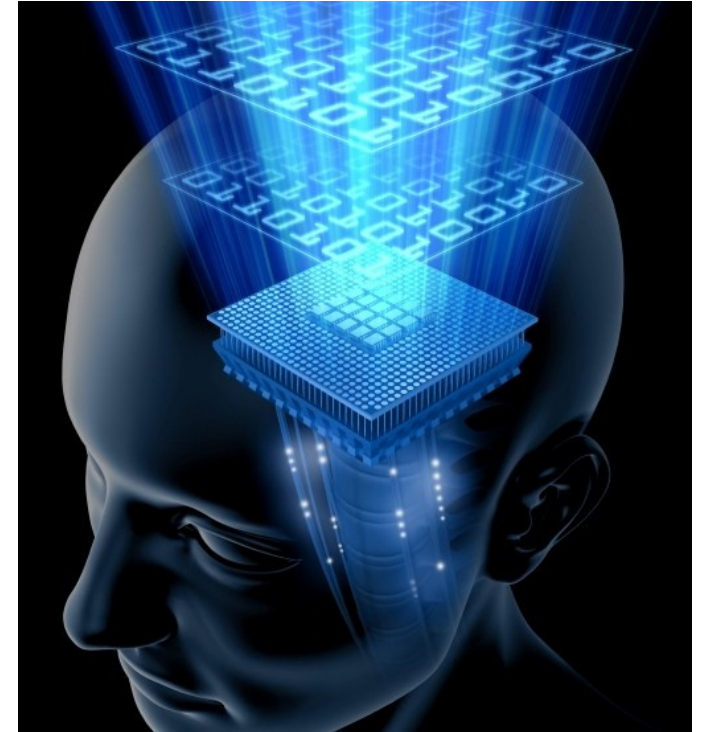
**FI\_MR\_SCALABLE**  
Application driven capability



**FI\_MR\_NEITHER**  
What other fabrics support

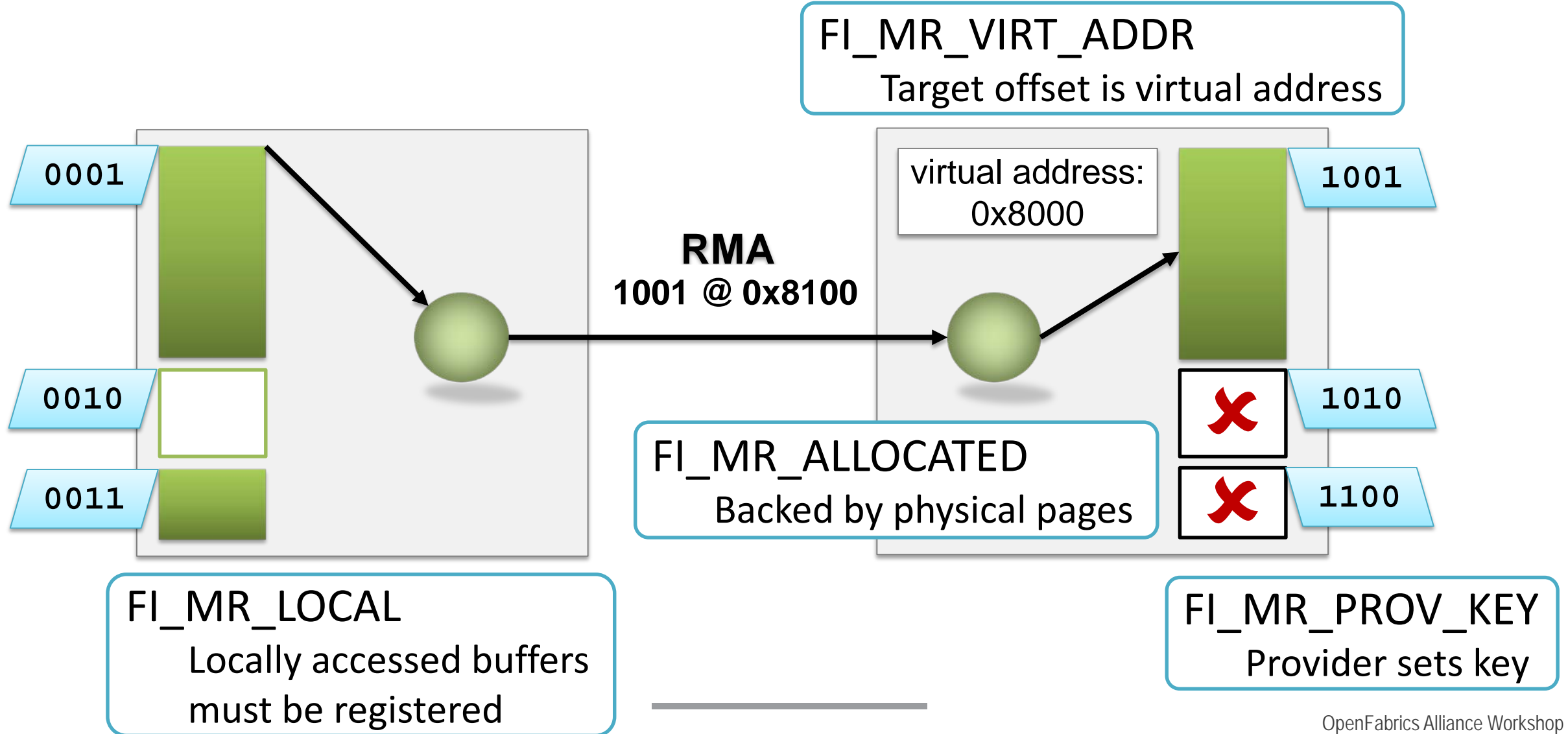
# MR MODE BITS

- **Redefine domain attribute `mr_mode`**
  - enum → integer flags
- **Divide `FI_MR_BASIC` into specific restrictions**
  - `FI_MR_SCALABLE` implied if `mr_mode = 0`
- **`FI_MR_BASIC`, `FI_MR_SCALABLE` kept for compatibility**
  - Behavior dependent on selected API version



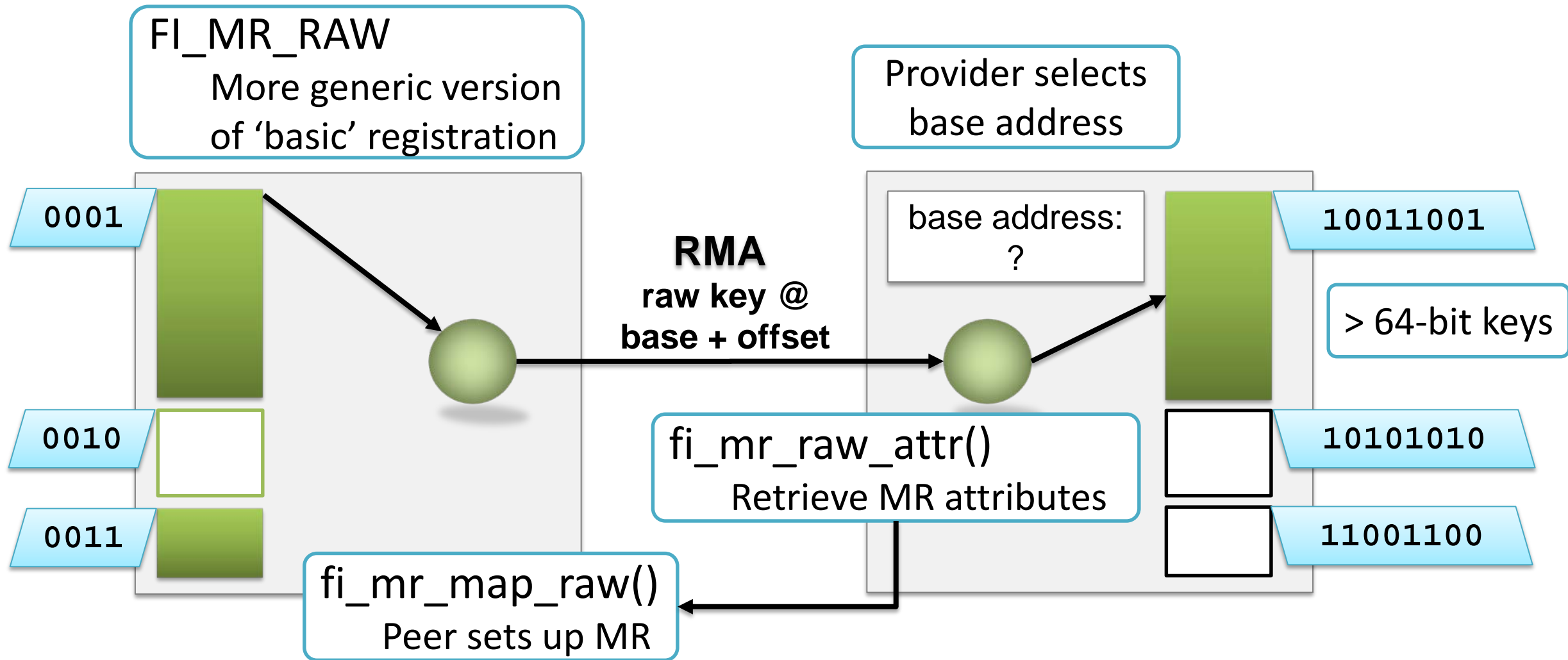


# 'BASIC' MR MODE BITS

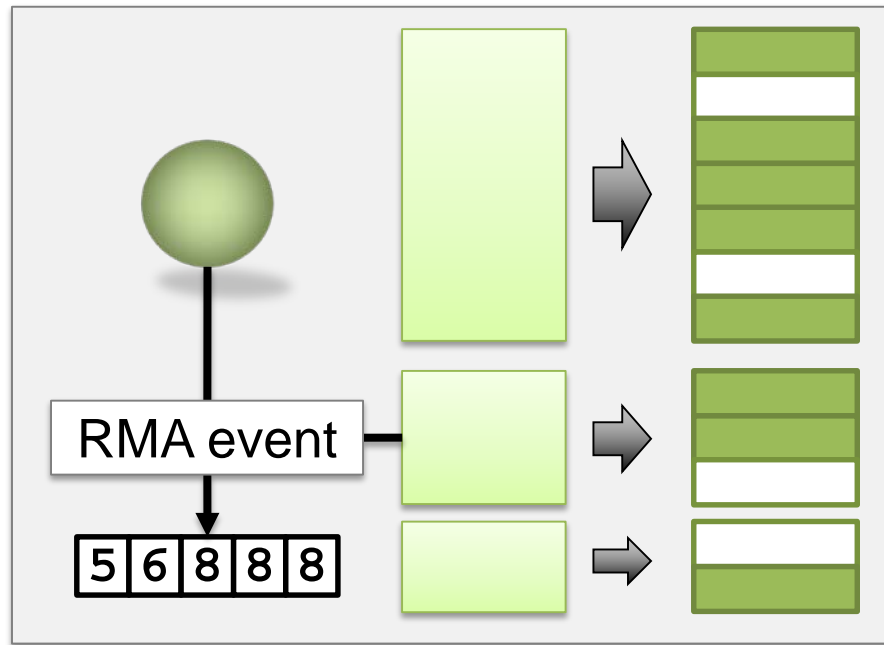


*It's only a  
flesh wound*

# RAW MEMORY REGIONS



# NON-OMNIPOTENT PROVIDERS



## FI\_MR\_MMU\_NOTIFY

- App notifies provider when physical pages backing MR change
- I.e. almost scalable, but NIC not linked into MMU
- Only those addresses accessed need to be backed

## FI\_MR\_RMA\_EVENT

- App indicates if MR will be used with RMA event counting
- MRs must be enabled after resource binding



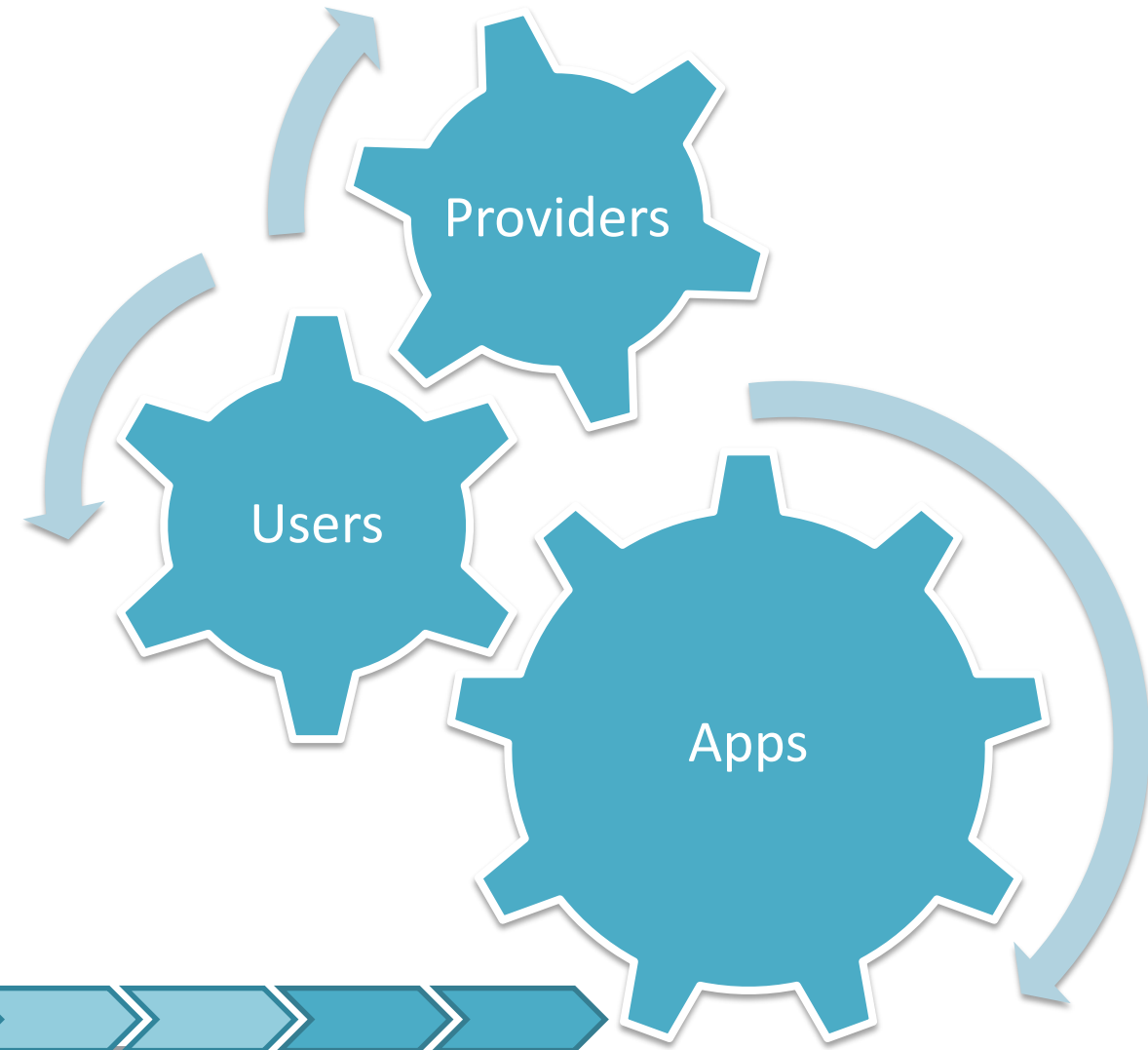
# LOOKING BACK

Balance application needs, developer requests, and provider capabilities

New features for expanding use cases

Developer assistance

Sensible performance optimizations





13<sup>th</sup> ANNUAL WORKSHOP 2017

**THANK YOU**

Sean Hefty, OFIWG Co-Chair

Intel