2023 OFA Virtual Workshop
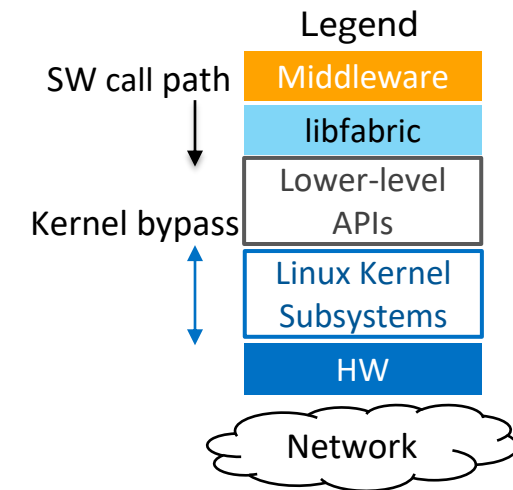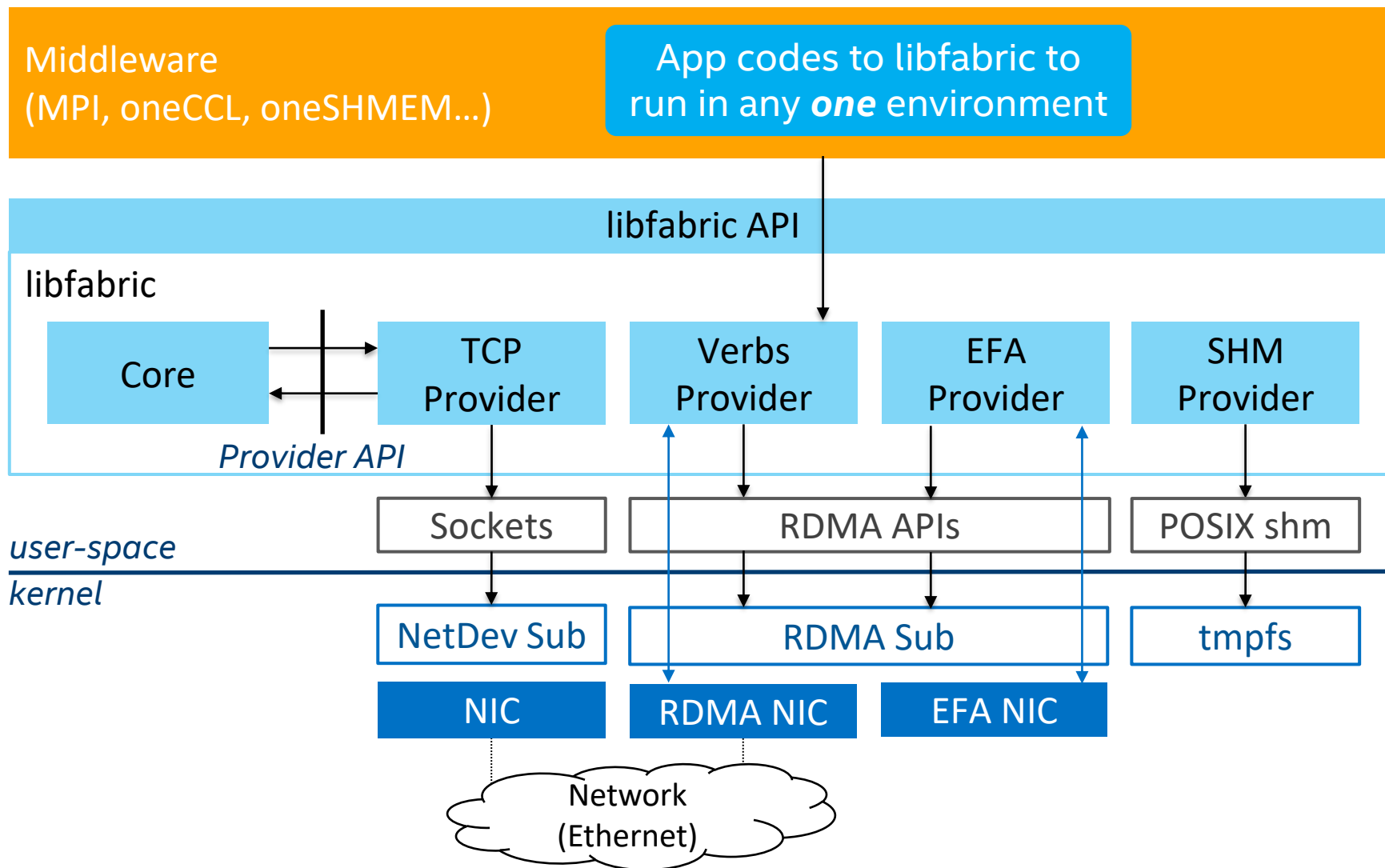
# libfabric Composability:
## Peer Provider Architecture

**Sean Hefty**

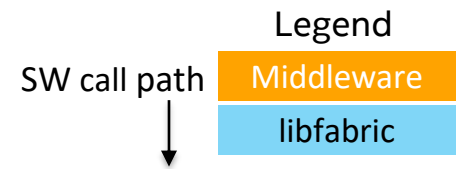Intel Corporation

# libfabric *Classic* Architecture
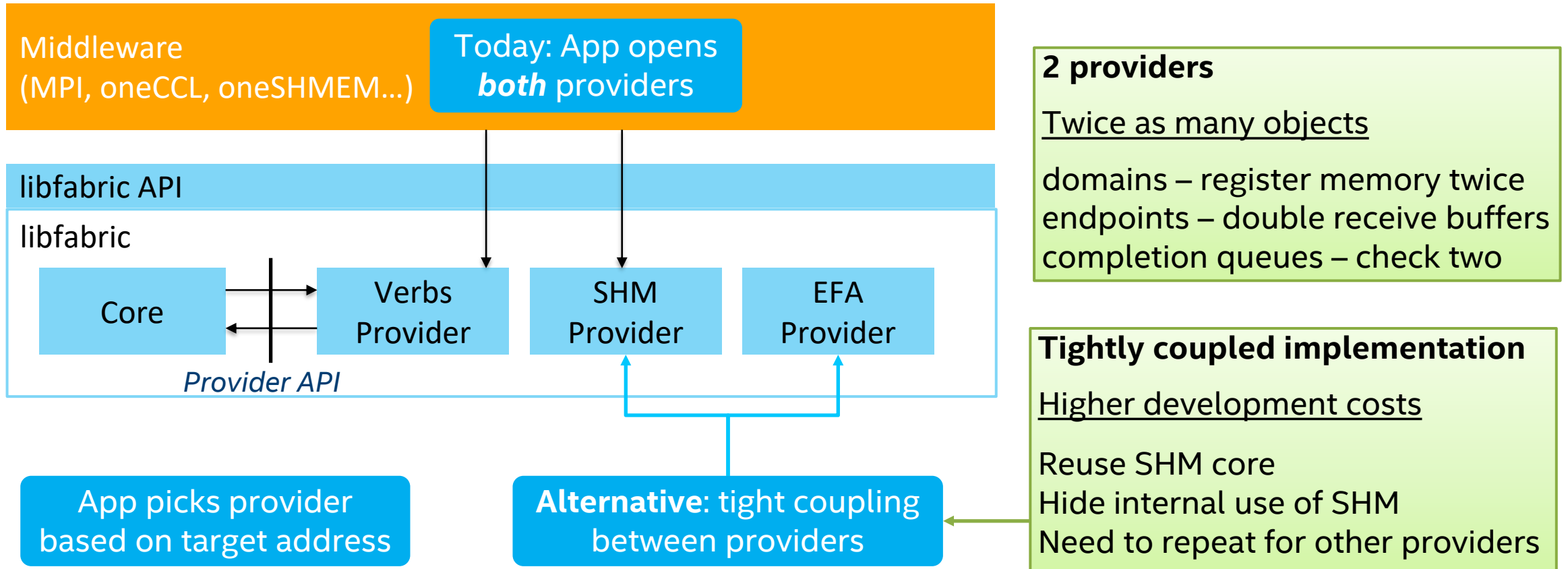## Illustrative Components

**Legend**

SW call path

| | |
|---|---|
| | Middleware |
| | libfabric |
| | Lower-level APIs |

Kernel bypass

| |
|---|
| Linux Kernel Subsystems |
| HW |

Network

---

**Middleware (MPI, oneCCL, oneSHMEM...)**

App codes to libfabric to run in any *one* environment

**libfabric API**

**libfabric**

| Core | | TCP Provider | Verbs Provider | EFA Provider | SHM Provider |
|---|---|---|---|---|---|

*Provider API*

| Sockets | RDMA APIs | POSIX shm |
|---|---|---|

*user-space*

*kernel*

| NetDev Sub | RDMA Sub | tmpfs |
|---|---|---|

| NIC | RDMA NIC | EFA NIC |
|---|---|---|

Network (Ethernet)

# Using Multiple Providers
## Example: Shared Memory + Network Provider
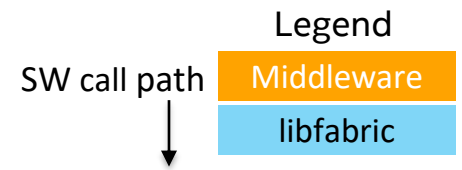
**Legend**
- Middleware
- libfabric

What if the app wants to use 2 providers?

E.g. overall performance will improve if the app can leverage shared memory for intranode communication
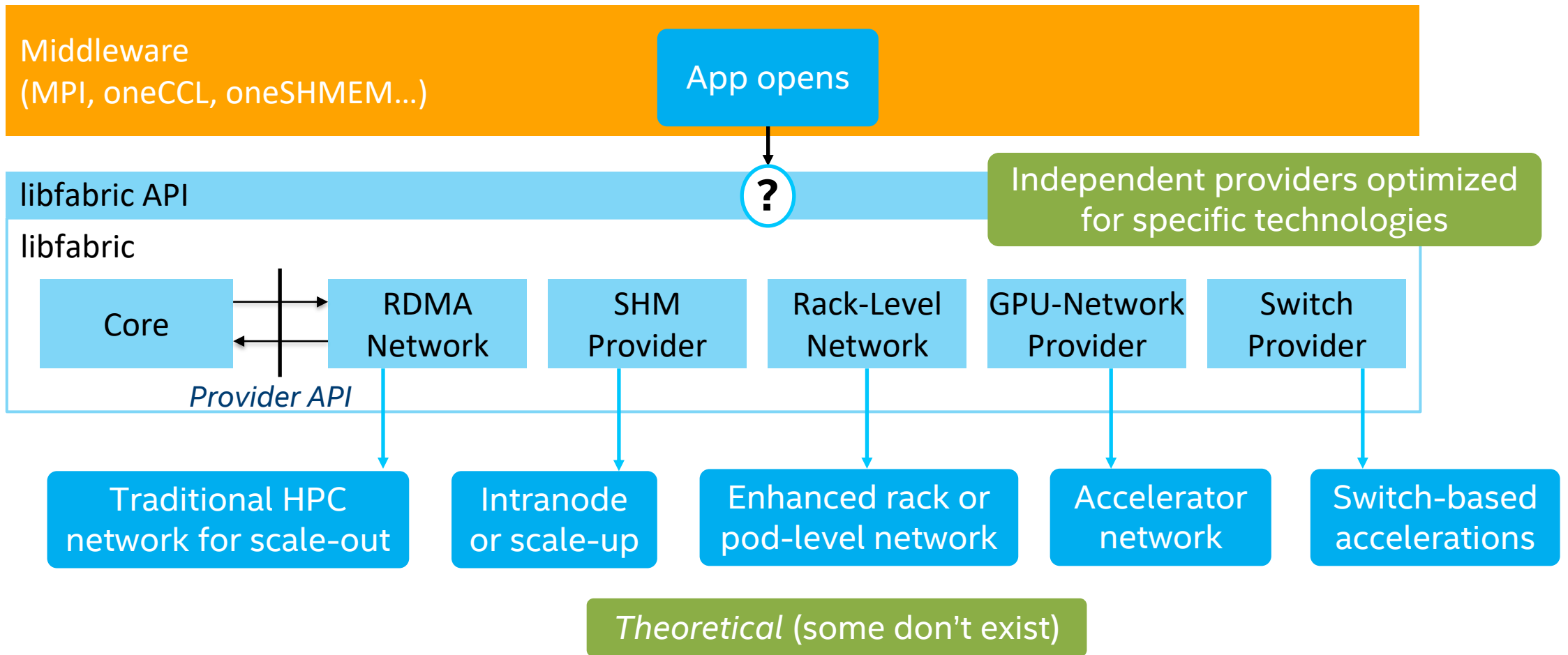
Middleware (MPI, oneCCL, oneSHMEM...)

Today: App opens *both* providers

libfabric API

libfabric

Core

Verbs Provider

SHM Provider

EFA Provider

*Provider API*

App picks provider based on target address

**Alternative**: tight coupling between providers

**2 providers**

<u>Twice as many objects</u>

domains – register memory twice
endpoints – double receive buffers
completion queues – check two

**Tightly coupled implementation**

<u>Higher development costs</u>

Reuse SHM core
Hide internal use of SHM
Need to repeat for other providers
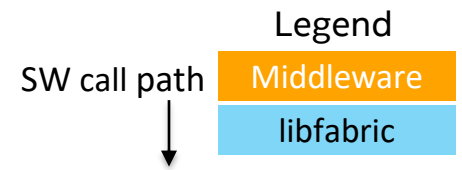
# Using Multiple Providers
## Complex, Theoretical Scenarios

What if multiple providers are needed for optimal performance?

Tight coupling becomes impractical

Middleware
(MPI, oneCCL, oneSHMEM...)

App opens

libfabric API

libfabric

Independent providers optimized for specific technologies

Core

RDMA Network

SHM Provider

Rack-Level Network

GPU-Network Provider

Switch Provider

*Provider API*

Traditional HPC network for scale-out

Intranode or scale-up

Enhanced rack or pod-level network

Accelerator network

Switch-based accelerations

*Theoretical* (some don't exist)

# Using Multiple Providers
## Example: Shared Memory + Network Provider

Needed: *efficient* cooperation of *independent* providers

➡ Peer APIs

Middleware
(MPI, oneCCL, oneSHMEM...)

App opens **one** provider, but has access to both

libfabric API

libfabric

Link Provider (Optional)

Core

Verbs Provider

SHM Provider

*Provider API*

*Peer API*

Orchestration handled by core provider or (future) link provider

Keep providers highly focused

Allow independent development

Easy for providers to adopt

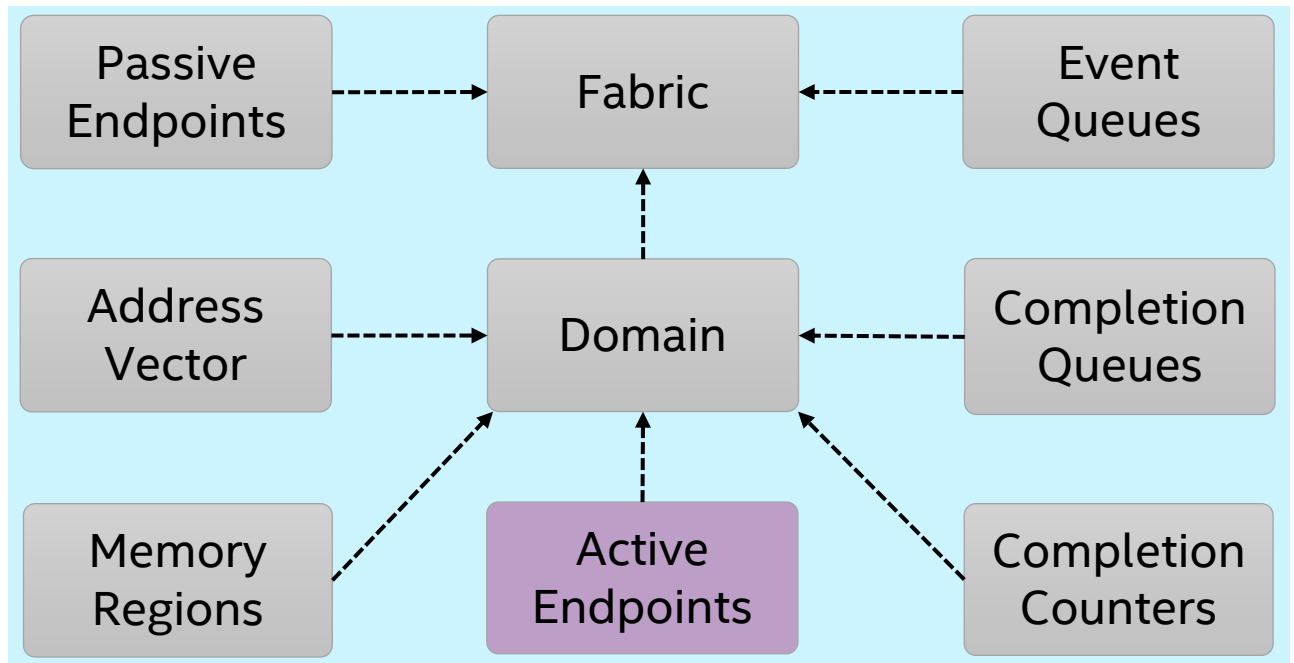Once a provider is enabled for peer API support, can swap in another to its right or left

# Review: libfabric API

API defines *user* interface to objects

| Passive Endpoints | Fabric | Event Queues |
|---|---|---|
| Address Vector | Domain | Completion Queues |
| Memory Regions | Active Endpoints | Completion Counters |

Example: active endpont

```
struct fid_ep {
    …
    struct fi_ops_msg *msg;
    struct fi_ops_rma *rma;
    …
};

static inline
fi_send(ep, buf, len, …)
{
    return ep->msg->send(ep, …)
}
```
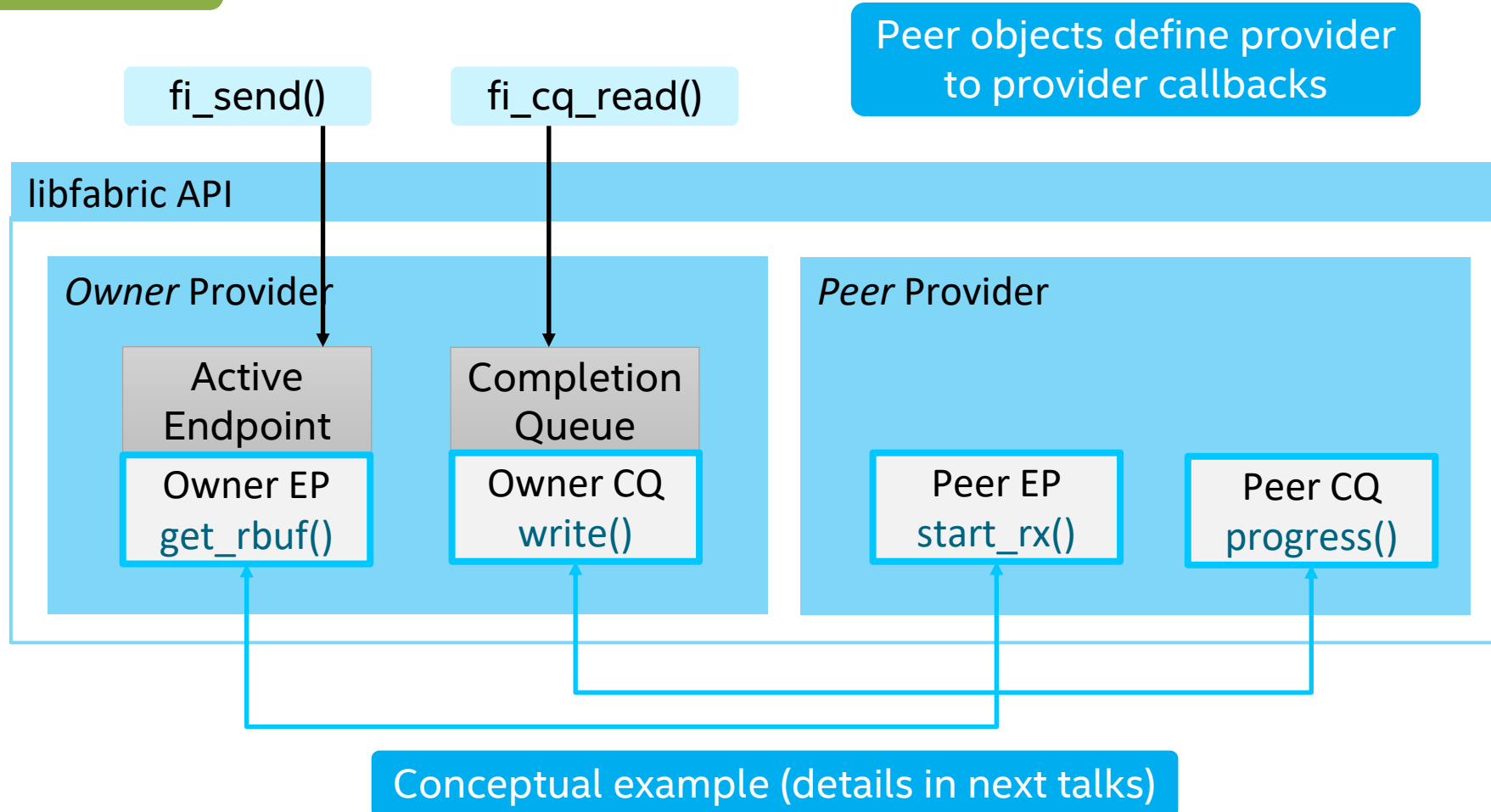
User invokes direct call on object

# Peer Object Model
## Sharable Fabric Identifiers (FIDs)

Define objects to share between **providers**

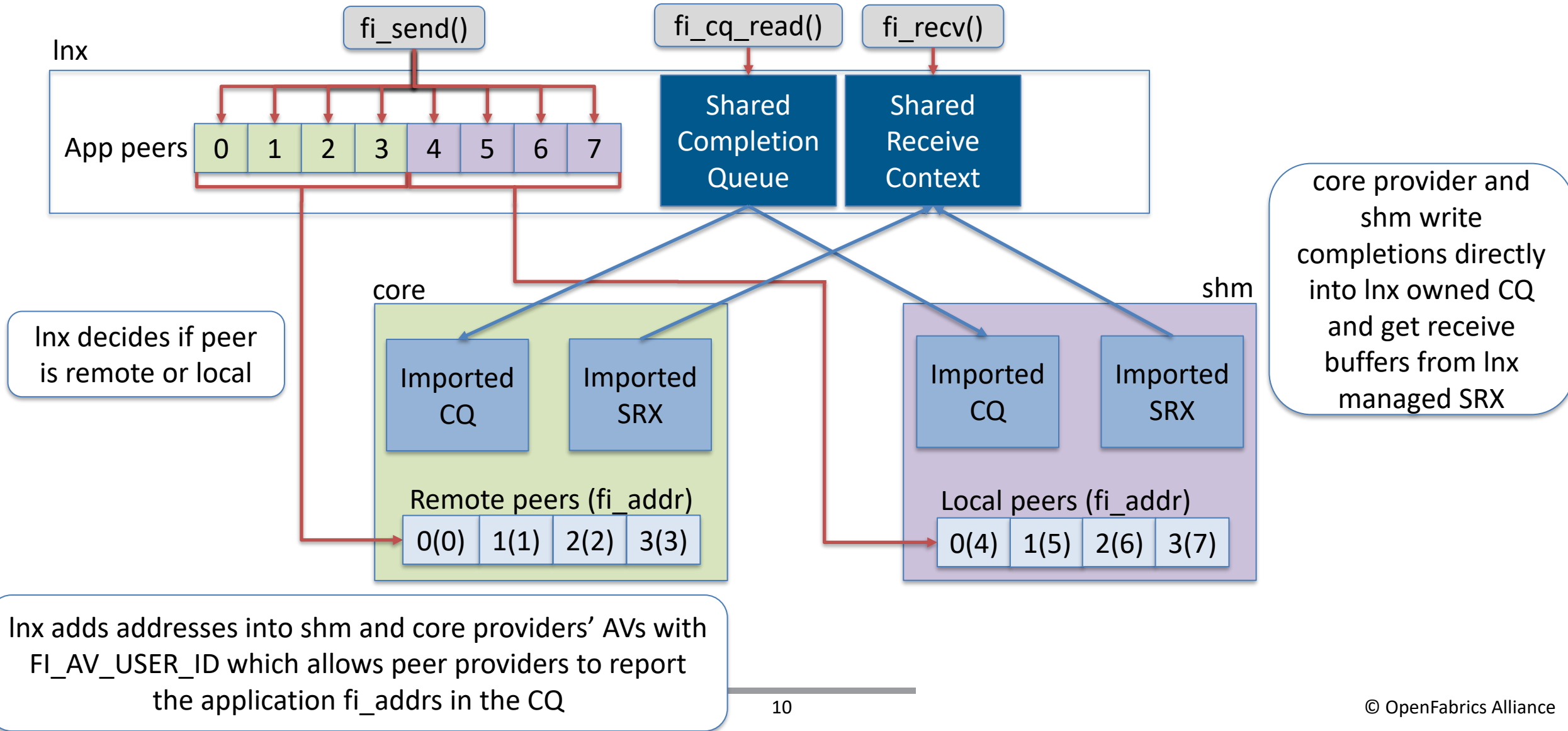Peer objects define provider to provider callbacks

fi_send()

fi_cq_read()

**libfabric API**

*Owner* Provider

| Active Endpoint |
| --- |
| Owner EP get_rbuf() |

| Completion Queue |
| --- |
| Owner CQ write() |

*Peer* Provider

| Peer EP start_rx() |
| --- |

| Peer CQ progress() |
| --- |

Conceptual example (details in next talks)

# EXAMPLE - OWNER: LNX



fi_send()

fi_cq_read()

fi_recv()

lnx

App peers: 0 1 2 3 4 5 6 7

Shared Completion Queue

Shared Receive Context

core provider and shm write completions directly into lnx owned CQ and get receive buffers from lnx managed SRX

lnx decides if peer is remote or local

core

Imported CQ

Imported SRX

Remote peers (fi_addr): 0(0) 1(1) 2(2) 3(3)

shm

Imported CQ

Imported SRX

Local peers (fi_addr): 0(4) 1(5) 2(6) 3(7)

lnx adds addresses into shm and core providers' AVs with FI_AV_USER_ID which allows peer providers to report the application fi_addrs in the CQ

10

# SHARED COMPLETION QUEUE API

1. Owner allocates a peer cq and defines peer CQ write ops

```
struct fid_peer_cq {
    struct fid fid;
    struct fi_ops_cq_owner *owner_ops;
};
        struct fi_ops_cq_owner {
            ssize_t (*write)();
            ssize_t (*writeerr)();
        };
```

3. Peer calls imported peer_cq->owner_ops in order to write an entry to the shared CQ

2. Owner calls fi_cq_open, passing in the peer_cq via context indicating a peer with attr->flags | FI_PEER

```
fi_cq_open(peer_domain, &attr, &peer_cq, peer_context);
                    struct fi_peer_cq_context {
                        struct fid_peer_cq *cq;
                    };
```

```
struct fi_peer_srx_context {
    struct fid_peer_srx *srx;
};
```

> 1. Owner creates peer_srx_context and sets owner ops

```
struct fid_peer_srx {
    struct fid_ep ep_fid;
    struct fi_ops_srx_owner *owner_ops;
    struct fi_ops_srx_peer *peer_ops;
};
```

> 2. Owner imports SRX into peer by calling fi_srx_context passing in the peer_cq via context indicating a peer with attr->flags | FI_PEER. Peer sets peer_ops

```
fi_srx_context(peer_domain, &attr, &srx_fid, peer_srx_context);
```

```
struct fi_ops_srx_owner {
    int (*get_msg)();
    int (*get_tag)();
    int (*queue_msg)();
    int (*queue_tag)();
    void (*free_entry)();
};
```

> Peer calls owner ops to get, queue, and free messages

```
struct fi_ops_srx_peer {
    int (*start_msg)();
    int (*start_tag)();
    int (*discard_msg)();
    int (*discard_tag)();
};
```

> Owner calls peer ops to start and discard unexpected messages

# EXAMPLE SRX FLOW

**OWNER**

**PEER**

SRX

`fi_recv()`

`fi_srx_context()`

`start_msg()`

`fi_recv()`

msg

`get_msg()`
`free_entry()`

FI_ENOENT    msg

`get_msg()`
`queue_msg()`

`free_entry()`

- **API summary**
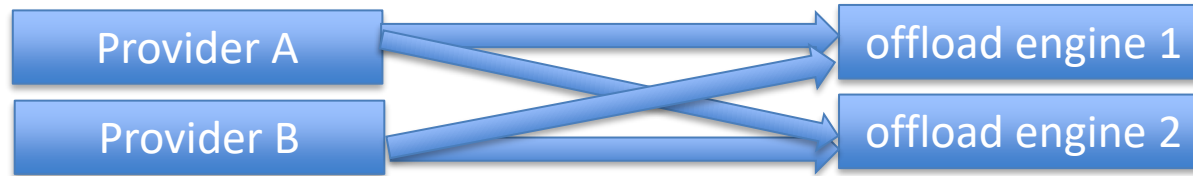  - Asynchronous
  - Defined in `<rdma/fi_collective.h>`
  - Supported ops: barrier, broadcast, alltoall, allreduce, allgather, reduce, reduce_scatter, scatter, gather
  - Wrapper functions: `fi_barrier(), fi_broadcast(), ……`
  - Collective groups: av_set
    - A set of addresses (`fi_addr_t`) representing group members
    - Can perform set operations: insert, remove, intersect, union, diff
    - Similar to multicast group, join via the same `fi_join()` call, but with `FI_COLLECTIVE` flag

- **Collective ops can be defined for each endpoint**

```
struct fid_ep {
    ……
    struct fi_ops_collective collective;
}
```

# IMPLEMENTATION CONSIDERATIONS

- **Goal: Efficiently enable multiple providers over multiple collective offload engines**



| Provider A | → offload engine 1 |
| Provider B | → offload engine 2 |

  - An example of offload engines is switch with collective support
- **Option 1 -- fully independent implementations**
  - Each provider implements collective ops for each offload engine
  - Pros: good separation between providers and between offload engines
  - Cons: a lot of duplicated efforts
- **Option 2 -- collective functions as utility code**
  - Pros: reduce code duplication
  - Cons: utility code enforce common basic data structures (domain, ep, cq, etc) to be used by providers
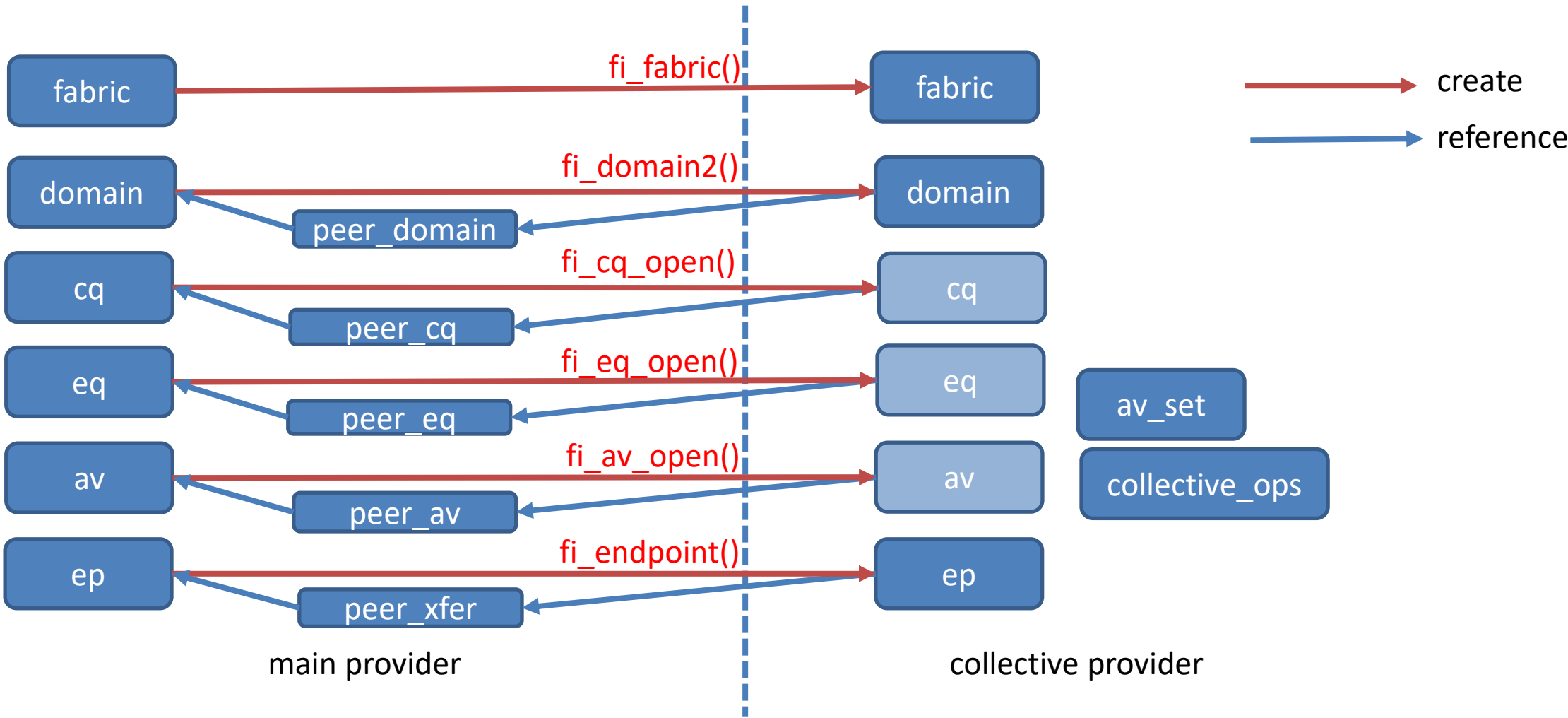- **Peer-provider provides a better option**

# COLLECTIVE OFFLOAD WITH PEER PROVIDER

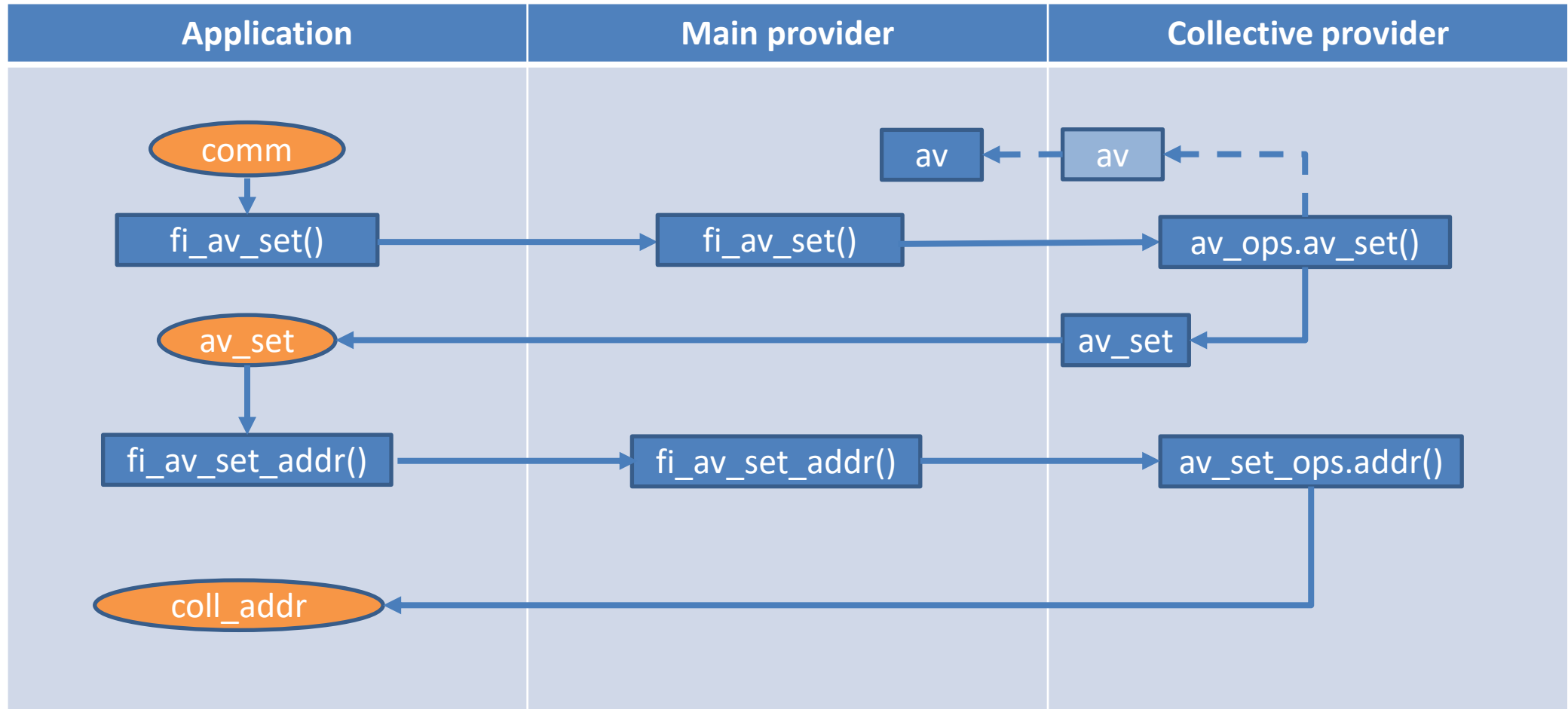- **Implement a collective-only provider for each offload engine**

| | | |
|---|---|---|
| provider A | offload provider 1 | offload engine 1 |
| provider B | offload provider 2 | offload engine 2 |

- Act as a peer provider to the "main" provider
- The main provider shares necessary data structure (domain, cq, eq, av, etc) via the peer-provider API
  - Eliminate the needs of creating duplicated queues / tables
  - The collective provider reports completions / events directly to the main provider
- Pros:
  - Reduce code duplication
  - Separation between the main provider and the offload provider – interface via peer-provider API only
- Cons:
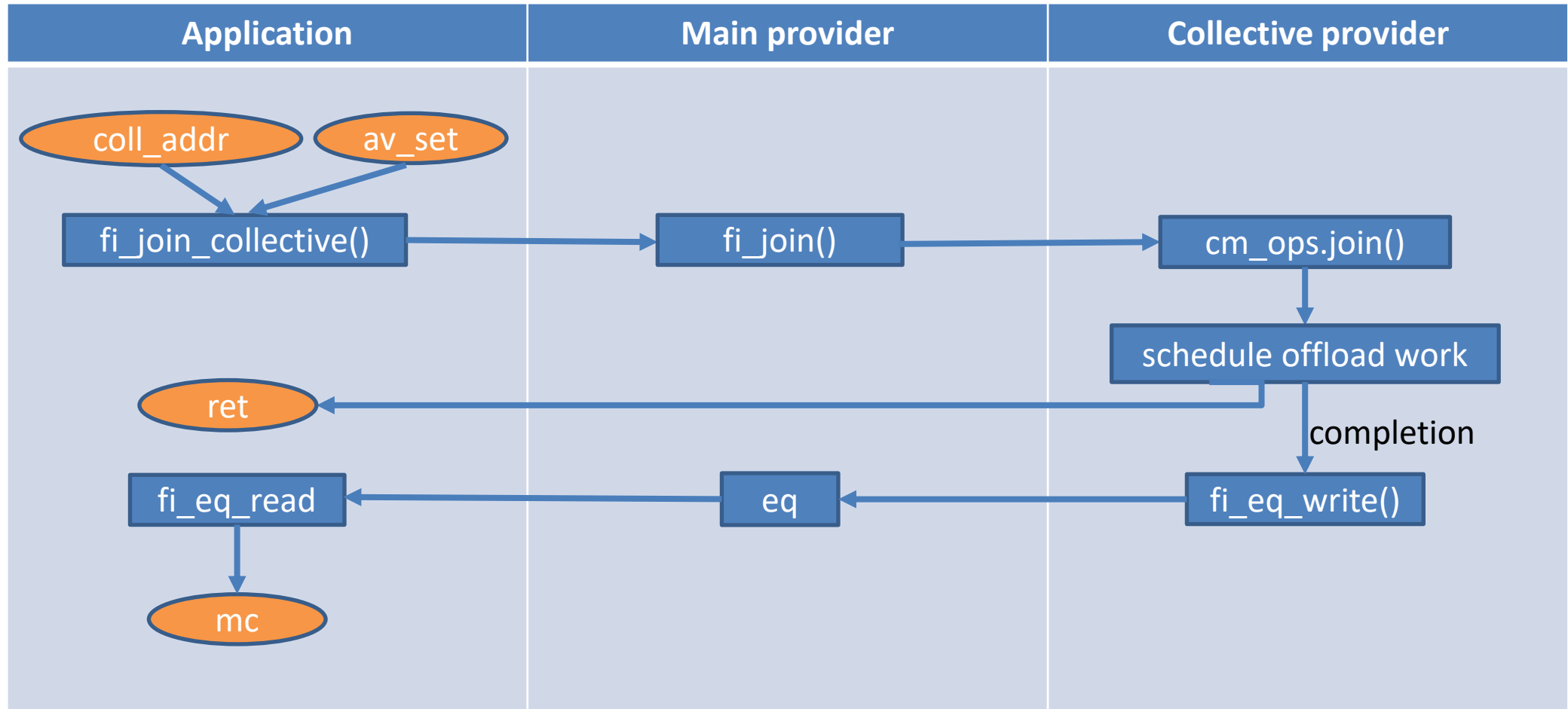  - The provider-to-provider workflow must be coordinated and well-defined
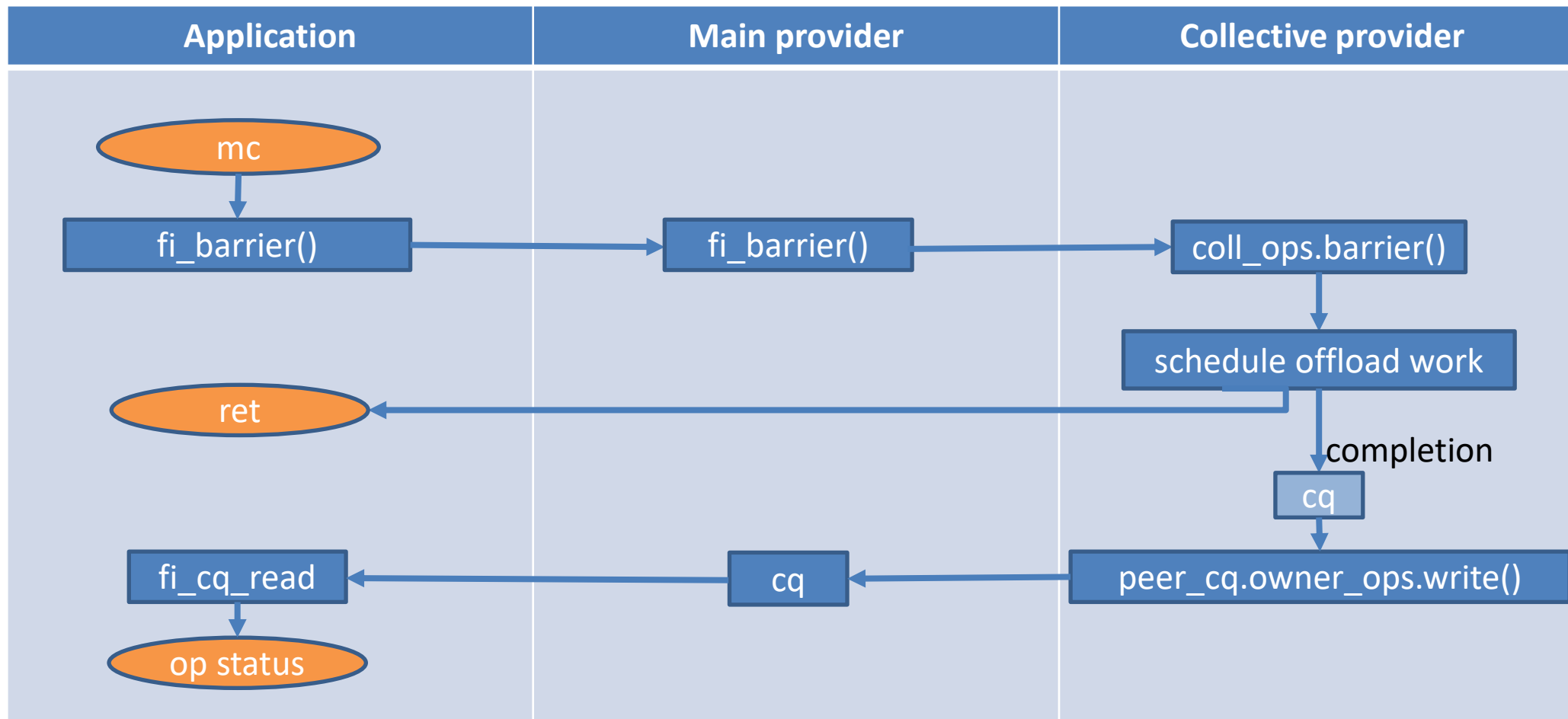
# DESIGN OVERVIEW

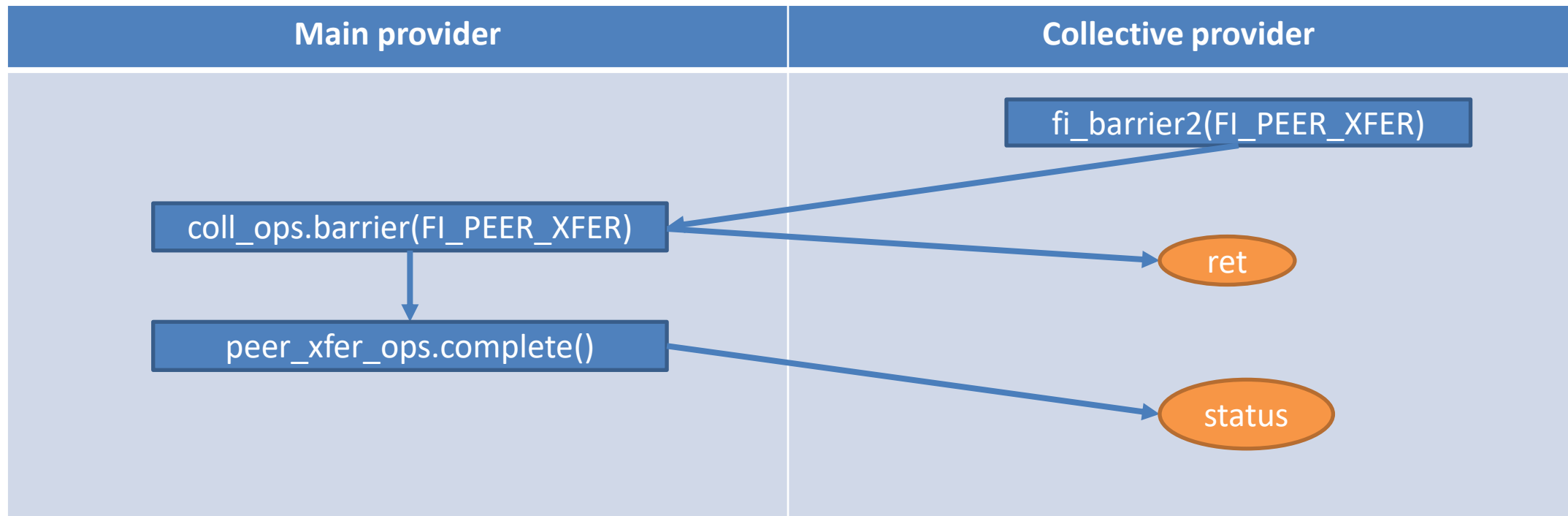# COLLECTIVE GROUP CREATION

# JOIN COLLECTIVE GROUP
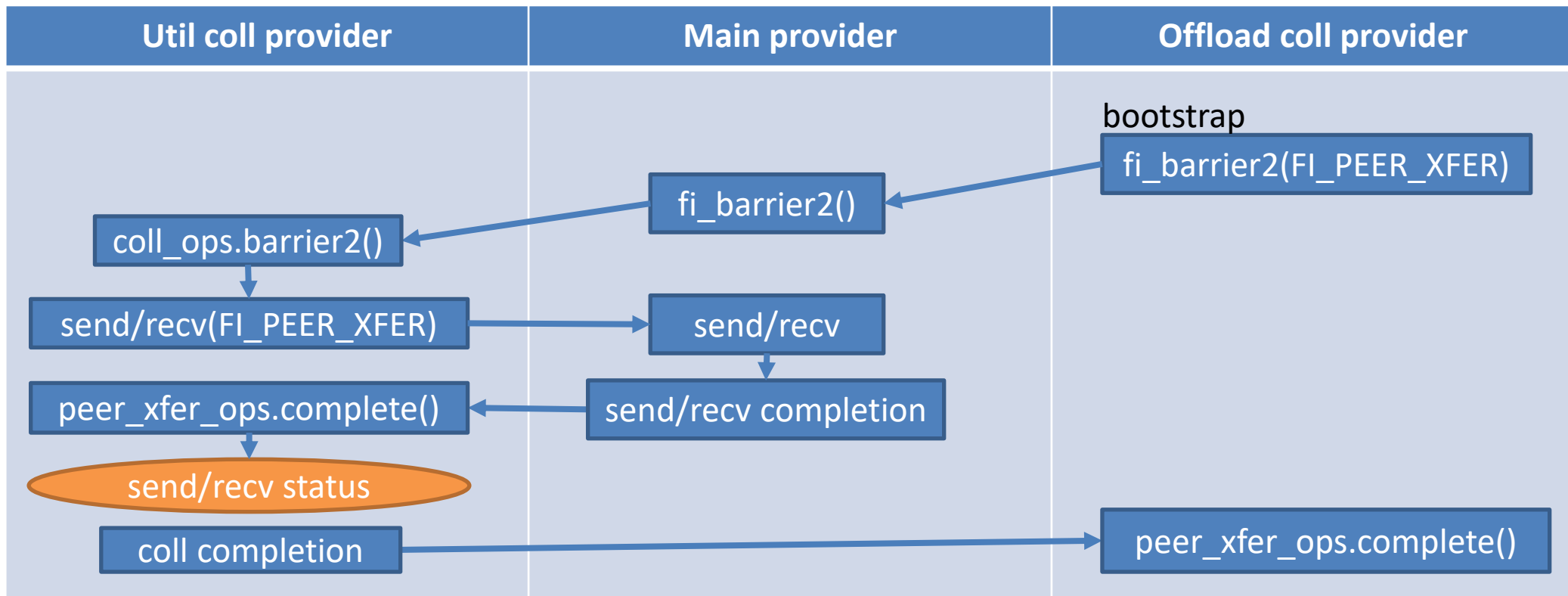
# COLLECTIVE OPS

# BOOTSTRAP COLLECTIVE

- **Offload collective engine may require a small set of out of band collectives for bootstrapping. Can be implemented in the main provider using pt2pt communication**

# UTILITY COLLECTIVE PROVIDER

- **Pt2pt based collectives can be moved to its own provider**

# CONCLUSION AND FUTURE WORK

- **Peer provider provides a mechanism for implementing "functional" providers w/o duplicating important data structures. Collective offload is one such function that suits this model well**

- **As a proof-of-concept, a utility collective provider has been implemented to provide software-based collective functionality.**
  - The rxm provider now uses this utility collective provider for default collective support instead of the old "shared utility code" based implementation.
  - Enables other providers to leverage the pt2pt based collective implementation more easily

- **Future work will have offload collective provider(s) implemented for popular collective offload engine(s). That's when upper layer middleware can start taking advantage of OFI collectives.**