2022 OFA Virtual Workshop

# IN-NETWORK COLLECTIVE COMMUNICATION ACCELERATIONS
## OFI COLLECTIVES

Sean Hefty

Intel Corporation

April 2022

# HOW COLLECTIVES DIFFER
## MULTICAST VS COLLECTIVE BROADCAST

**Collective broadcast**

Message

data

ack

RDM endpoint

join collective()
broadcast send()

Knows who receivers are

Collective Network

**In-network support**
- *Message* replication
- Reliable delivery – ack coalescing

RDM endpoint

join collective()
broadcast recvfrom()

Recv matched with send

RDM endpoint

join collective()
broadcast recvfrom()

Joins coordinated

RDM endpoint

join collective()
broadcast recvfrom()

Knows who sender is

RDM – Reliable Datagram

# COLLECTIVE OPERATIONS
## CONCEPTUAL: "MULTICAST ATOMICS"

**Definitions by example**

| Peer 0 |
|--------|
| 1 |
| 5 |
| 9 |

← data array

**broadcast**

| Peer 1 | Peer 2 |
|--------|--------|
| 1 | 1 |
| 5 | 5 |
| 9 | 9 |

**barrier** (global sync)

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 1 | 2 | 3 |
| 5 | 6 | 7 |
| 9 | 10 | 11 |

**all to all**

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 1 | 5 | 9 |
| 2 | 6 | 10 |
| 3 | 7 | 11 |

**Same data types and operations as atomic APIs**

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 1 | 1 | 1 |
| 5 | 5 | 5 |
| 9 | 9 | 9 |

**sum**

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 3 | 3 | 3 |
| 15 | 15 | 15 |
| 27 | 27 | 27 |

**all reduce**

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 1 | 5 | 9 |

**all gather**

| Peer 0 | Peer 1 | Peer 2 |
|--------|--------|--------|
| 1 | 1 | 1 |
| 5 | 5 | 5 |
| 9 | 9 | 9 |

**Collectives not appearing on stage**: gather, scatter, reduce, reduce-scatter

**Additional in-network support**
- *Data* replication
- Computation – data format aware
- Data coalescing and distribution

**1. Identify collective membership**

Select participating peers

Local operation – address vector sets

**2. Setup communication groups**

Coordinated join among members

Network operation (maybe) – 2 supported models

**3. Invoke collective**

Collective data transfer operation

# IDENTIFY COLLECTIVE MEMBERSHIP
## ADDRESS VECTOR SETS

**Existing API**

**AV**
Represents the peer universe

Application visible addresses of peer endpoints

Translated addresses

Unicast address used by API

**e.g.: IP : Port**

| e.g.: IP : Port |
|---|
| 10.0.0.1 : 7000 |
| 10.0.0.1 : 7001 |
| 10.0.0.2 : 7000 |
| 10.0.0.3 : 7003 |
| ... |

map

| Address Vector | |
|---|---|
| **fi_addr_t** | **Fabric Address** |
| 0 | 100:3:50 |
| 1 | 100:3:51 |
| 2 | 101:3:83 |
| 3 | 102:3:64 |
| ... | ... |

**Collective Extension**

**Local setup to identify members for new group**

**AV Set**
Set of peers in a collective group

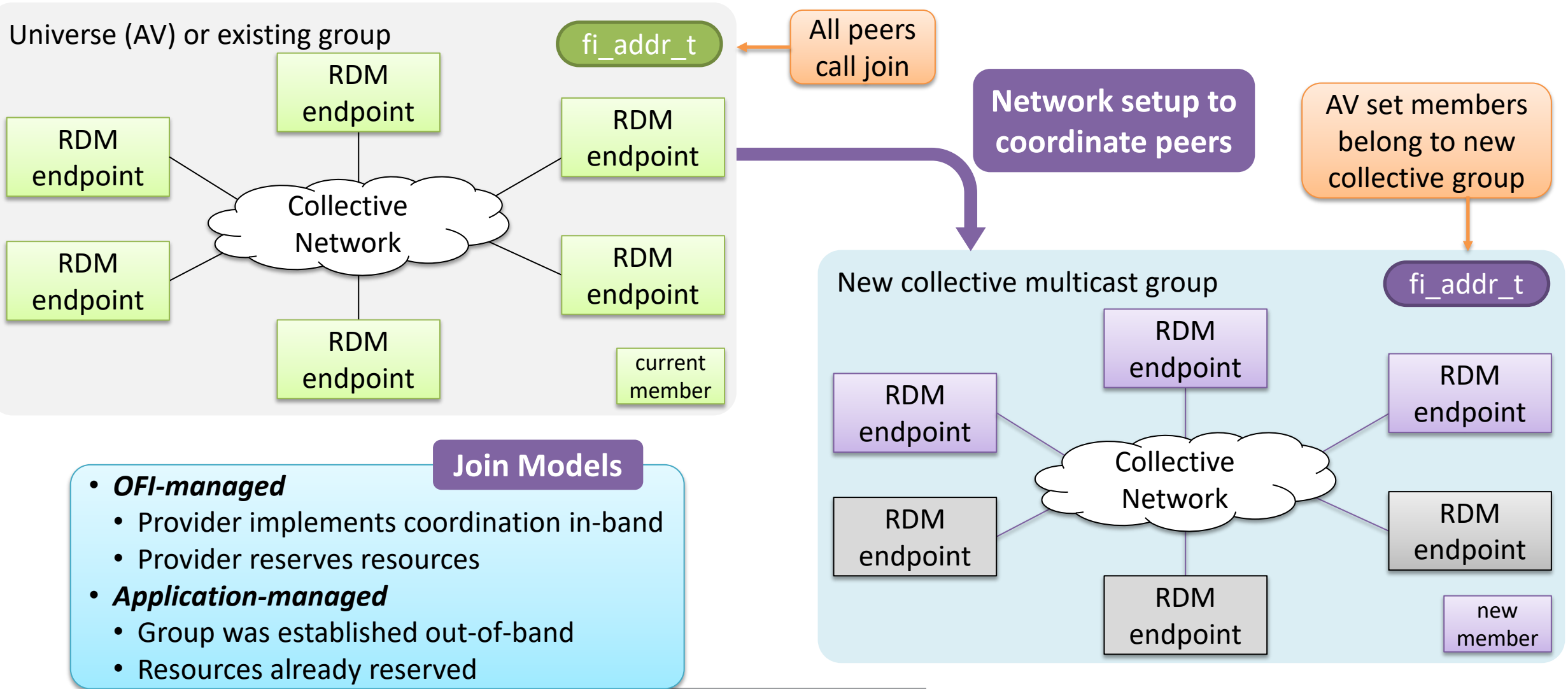Associated collective address

select

| Address Vector Set |
|---|
| **fi_addr_t** |
| 0 |
| 2 |
| 4 |
| 6 |
| ... |

fi_addr_t

**AV Set operations**
- Insert
- Remove
- Union
- Intersect
- Difference

AV – Address Vector

Universe (AV) or existing group

fi_addr_t

All peers call join

RDM endpoint

RDM endpoint

RDM endpoint

Collective Network

RDM endpoint

RDM endpoint

RDM endpoint

current member

**Network setup to coordinate peers**

AV set members belong to new collective group

New collective multicast group

fi_addr_t

RDM endpoint

RDM endpoint

RDM endpoint

Collective Network

RDM endpoint

RDM endpoint

RDM endpoint

new member

**Join Models**

- *OFI-managed*
  - Provider implements coordination in-band
  - Provider reserves resources
- *Application-managed*
  - Group was established out-of-band
  - Resources already reserved

```
struct fi_info *hints, *info;

hints fi_allocinfo();
<format hints>
hints->caps |= FI_COLLECTIVE;

fi_getinfo(FI_VERSION(1,14), hostname, NULL, FI_SOURCE, hints, &info);

<allocate fabric resources>

struct fi_collective_attr attr = {0};
attr.op = FI_SUM;
attr.datatype = FI_FLOAT;

fi_query_collective(domain, FI_ALLREDUCE, &attr, 0);
assert(attr.datatype_attr.count >= 100 && attr.max_members >= 50)
```

Request support for in-network collectives

Verify support for collective that we need

```
struct fid_av_set *av_set;
struct fi_av_set_attr attr = {0};
attr.start_addr = 2;
attr.end_addr = 100;
attr.stride = 2;
fi_av_set(av, &attr, &av_set, NULL);


struct fid_mc *mc;
fi_join_collective(ep, FI_ADDR_UNAVAIL, av_set, 0, &mc, NULL);


struct fi_eq_entry entry;
uint32_t event;
fi_eq_sread(eq, &event, &entry, sizeof(entry), -1, 0);
assert(event == FI_JOIN_COMPLETE);


fi_allreduce(ep, input_array, 100, NULL, result_array, NULL,
        fi_mc_addr(mc), FI_FLOAT, FI_SUM, 0, my_context);
```

Create AV set and identify group members

First join involves all peers

Creates collective multicast group

Join completes asynchronously

Asynchronous all-reduce operation

**Managing in-network resources**

Guarantee resources are available

App may want to prioritize which collectives to accelerate

**API object: collective resource tokens?**

**Priority**

Define impact on active collectives

Preempt possible?  Pause-resume or abort/cancel?

libfabric defines priority at the endpoint level

**Do resource tokens act as a proxy?**

## Reproducibility of results

Order that data is fed into operations can produce different results

Relaxed reproducibility can reduce in-network memory

**Setting: per-operation, group (AV set), resource token?**

## Sparse data

Avoid sending / storing null data

**Define a compact, data aware SGL?**

## Network topology

Query collective support - local vs global?

Peer endpoints relative to switches and accelerators

**Scope of the job or resource manager?**

## Programmable in-network accelerations

Non-collective operations

How does app specify operation and parameters?

Entity responsible for programming switch/FPGA?

THANK YOU