



Fabric Interfaces Architecture

Sean Hefty - Intel Corporation

Changes

- v2
 - Remove interface object
 - Add open interface as base object
 - Add SRQ object
 - Add EQ group object

Overview

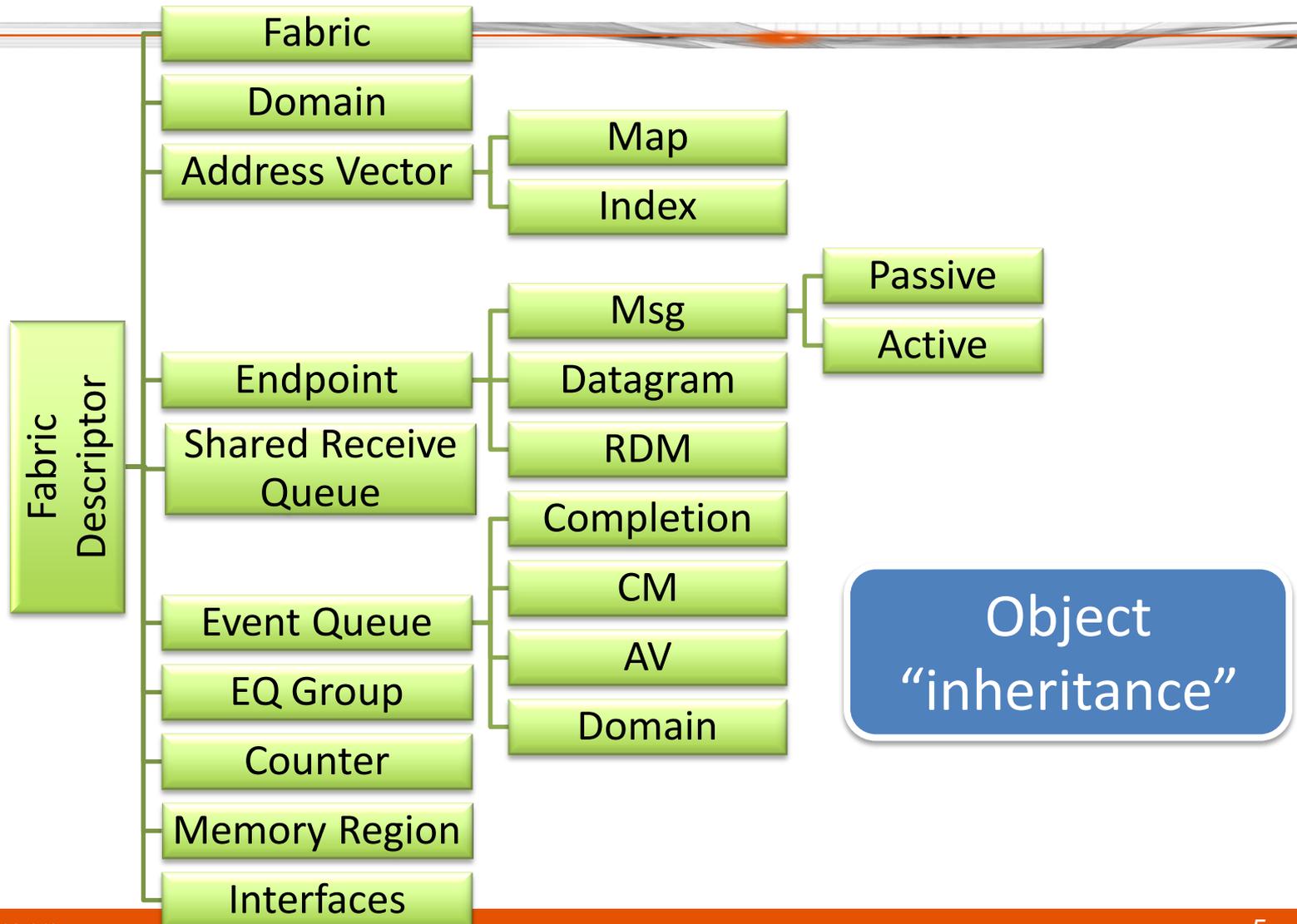
- Object Model
 - Do we have the right type of objects defines?
 - Do we have the correct object relationships?
- Interface Synopsis
 - High-level description of object operations
 - Is functionality missing?
 - Are interfaces associated with the right object?
- Architectural Semantics
 - Do the semantics match well with the apps?
 - What semantics are missing?

Object “Class” Model

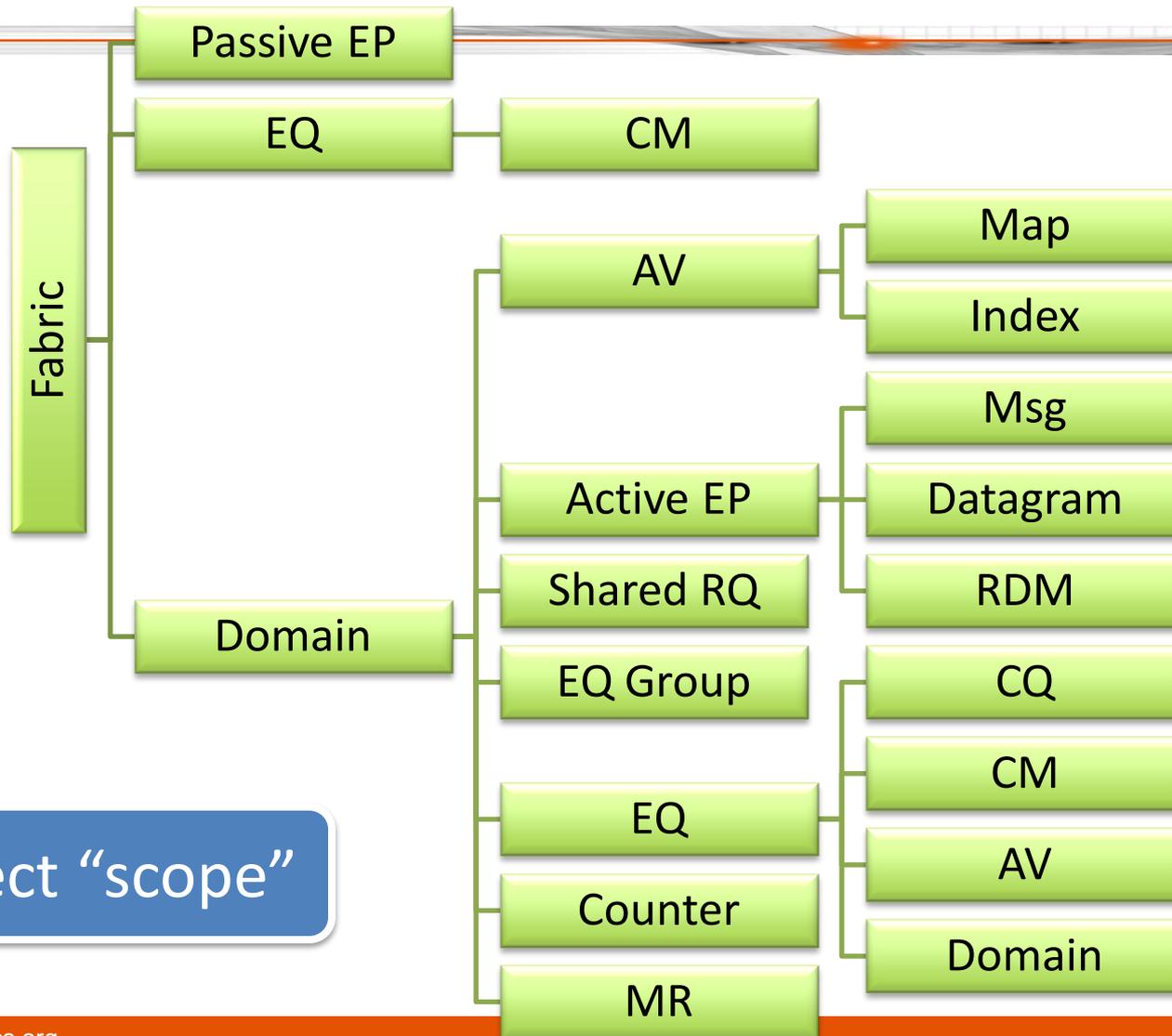
- Objects represent collection of attributes and interfaces
 - I.e. object-oriented programming model
- Consider architectural model only at this point

Objects do not necessarily map directly to hardware or software objects

Conceptual Object Hierarchy

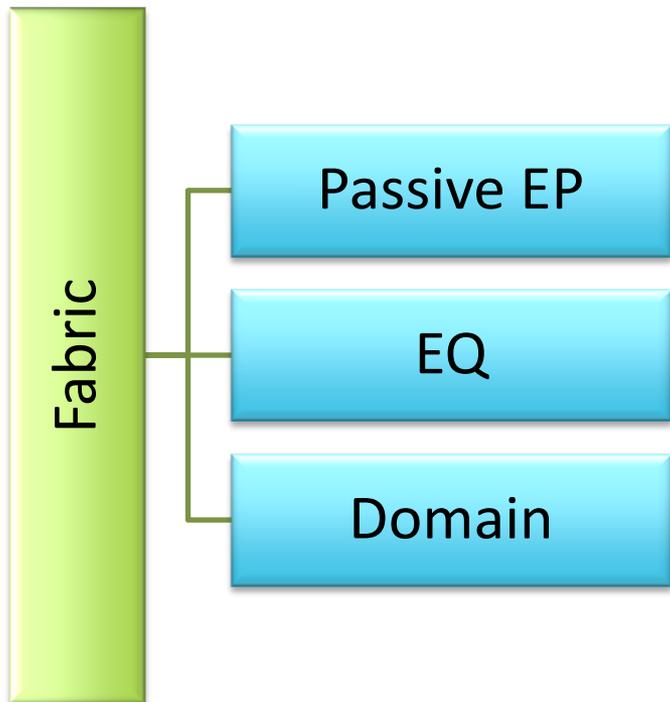


Object Relationships



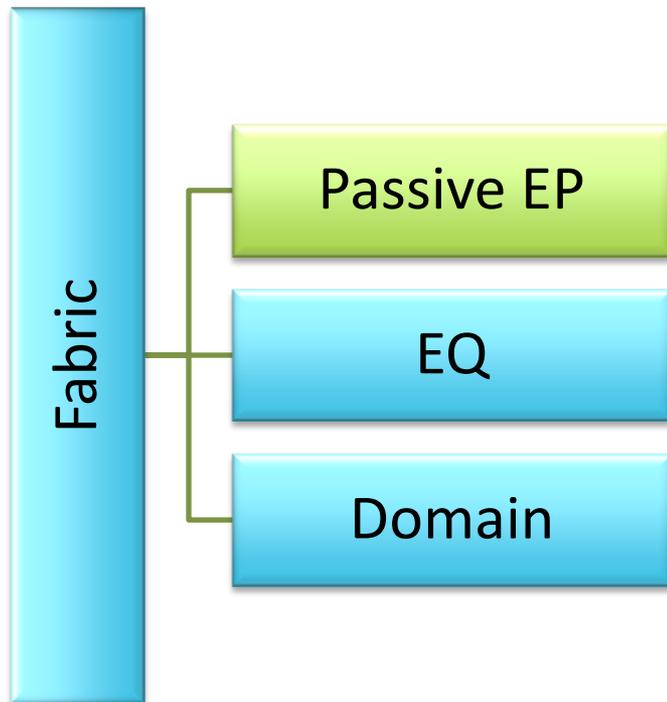
Object "scope"

Fabric



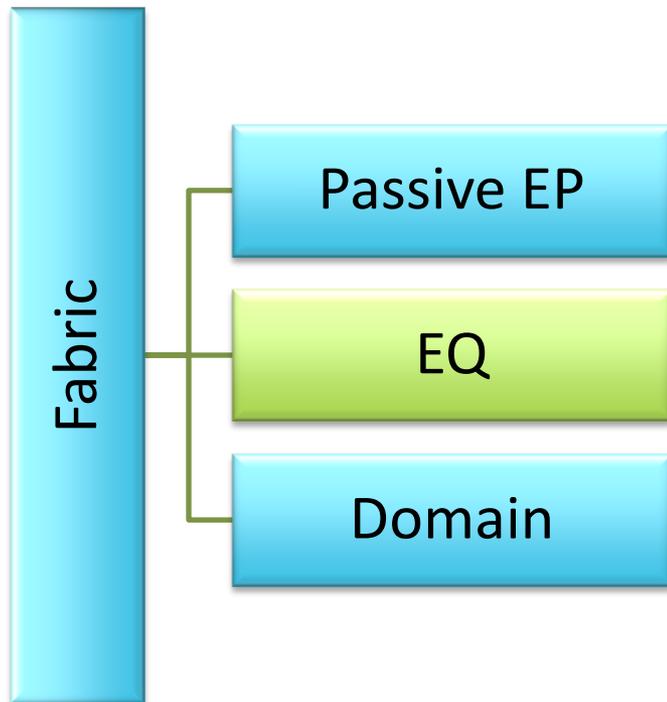
- Represents a communication domain or boundary
 - Single IB or RoCE subnet, IP (iWarp) network, Ethernet subnet
- Multiple local NICs / ports
- Topology data, network time stamps
- Determines native addressing
 - Mapped addressing possible
 - GID/LID versus IP

Passive (Fabric) EP



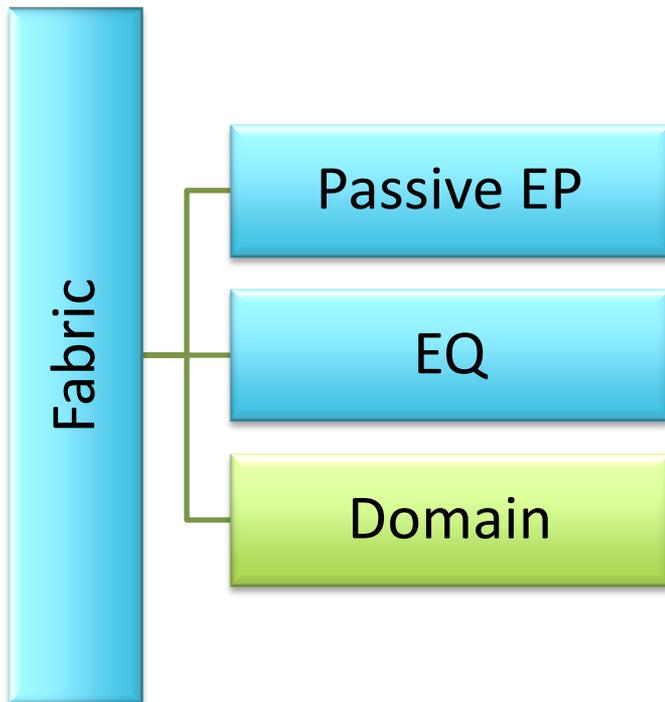
- Listening endpoint
 - Connection-oriented protocols
- Wildcard listen across multiple NICs / ports
- Bind to address to restrict listen
 - Listen may migrate with address

Fabric EQ



- Associated with passive endpoint(s)
- Reports connection requests
- Could be used to report fabric events

Resource Domain



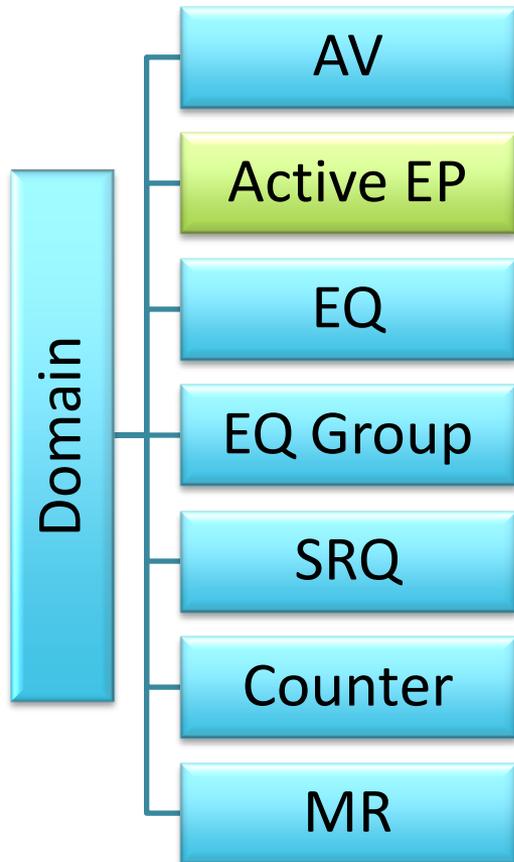
- Boundary for resource sharing
 - Physical or logical NIC
 - Command queue
- Container for data transfer resources
- A provider may define multiple domains for a single NIC
 - Dependent on resource sharing

Domain Address Vectors



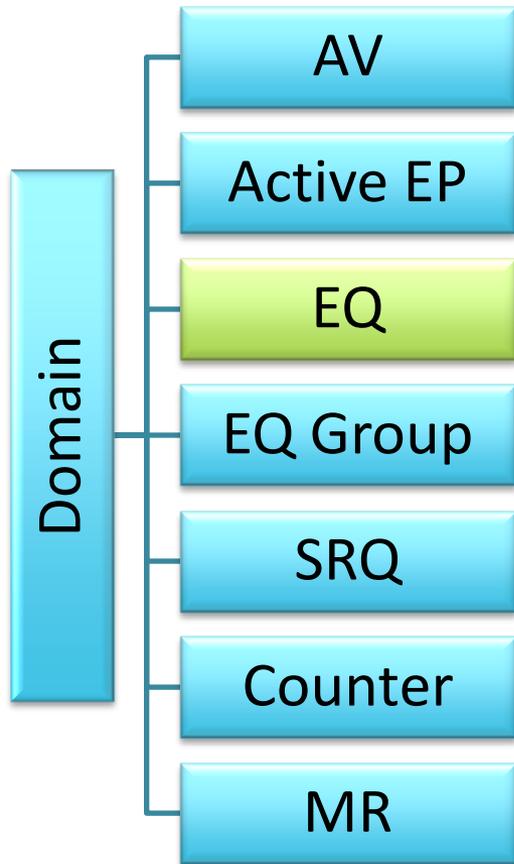
- Maintains list of remote endpoint addresses
 - Map – native addressing
 - Index – ‘rank’-based addressing
- Resolves higher-level addresses into fabric addresses
 - Native addressing abstracted from user
- Handles address and route changes

Domain Endpoints



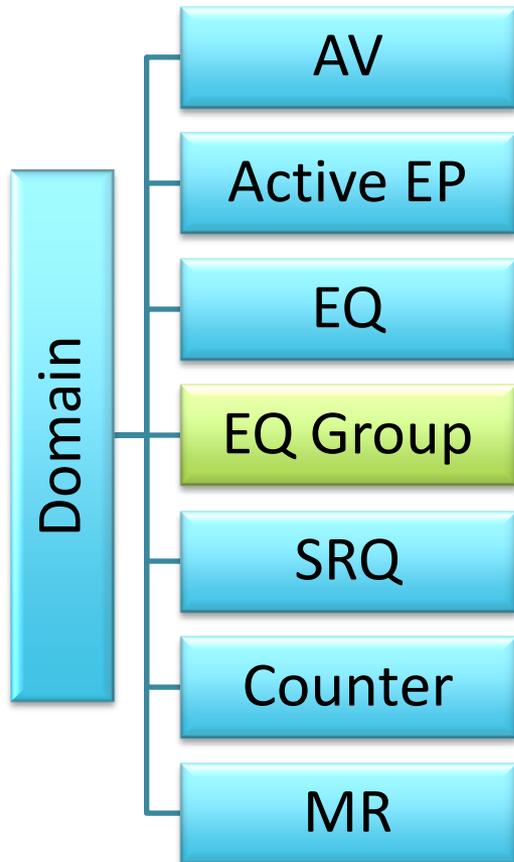
- Data transfer portal
 - Send / receive queues
 - Command queues
 - Ring buffers
- Multiple types defined
 - Connection-oriented / connectionless
 - Reliable / unreliable
 - Message / stream

Domain Event Queues



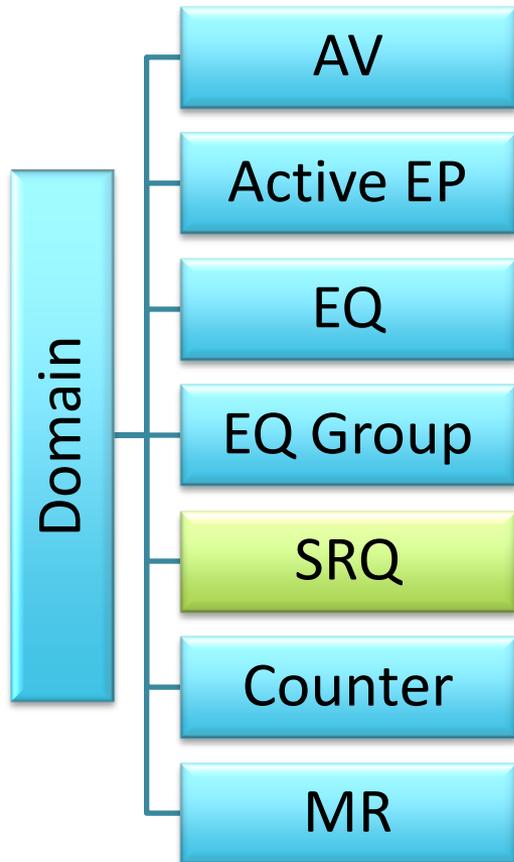
- Reports asynchronous events
- Unexpected errors reported 'out of band'
- Events separated into 'EQ domains'
 - CM, AV, completions
 - 1 EQ domain per EQ
 - Future support for merged EQ domains

EQ Groups



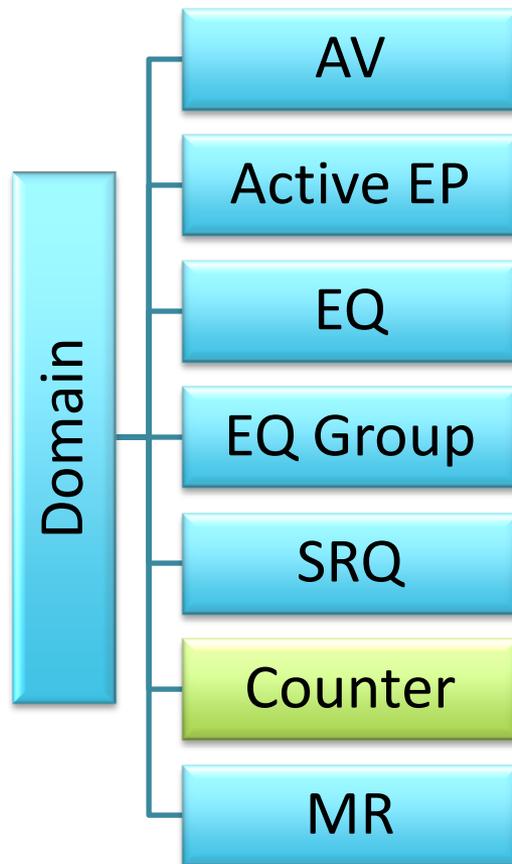
- Collection of EQs
- Conceptually shares same wait object
- Grouping for progress and wait operations

Shared Receive Queue



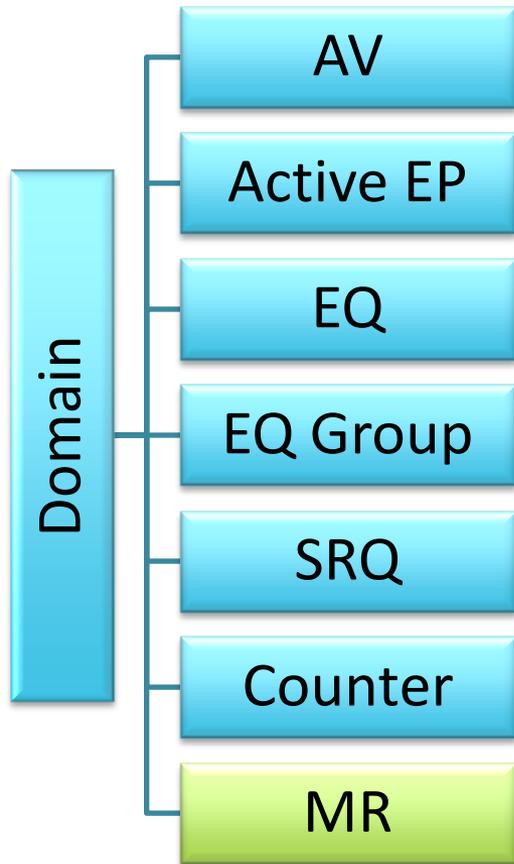
- Shares buffers among multiple endpoints
- Not addressable
 - Addressable SRQs are abstracted within the endpoint

Domain Counters



- Provides a count of successful completions of asynchronous operations
 - Conceptual HW counter
- Count is independent from an actual event reported to the user through an EQ

Domain Memory Regions



- Memory ranges accessible by fabric resources
 - Local and/or remote access
- Defines permissions for remote access

Interface Synopsis

- Operations associated with identified ‘classes’
- General functionality, versus detailed methods
 - The full set of methods are not defined here
 - Detailed behavior (e.g. blocking) is not defined
- Identify missing and unneeded functionality
 - Mapping of functionality to objects

Use timeboxing to limit scope of interfaces to refine by a target date

Base Class

Close	Destroy / free object
Bind	Create an association between two object instances
Sync	Fencing operation that completes only after previously issued asynchronous operations have completed
Control	(~fcntl) set/get low-level object behavior
I/F Open	Open provider extended interfaces

Fabric

Domain	Open a resource domain
Endpoint	Create a listening EP for connection-oriented protocols
EQ Open	Open an event queue for listening EP or reporting fabric events

Resource Domain

Query	Obtain domain specific attributes
Open AV, EQ, EP, SRQ, EQ Group	Create an address vector, event or completion counter, event queue, endpoint, shared receive queue, or EQ group
MR Ops	Register data buffers for access by fabric resources

Address Vector

Insert	Insert one or more addresses into the vector
Remove	Remove one or more addresses from the vector
Lookup	Return a stored address
Straddr	Convert an address into a printable string

Base EP

Enable	Enables an active EP for data transfers
Cancel	Cancel a pending asynchronous operation
Getopt	(~getsockopt) get protocol specific EP options
Setopt	(~setsockopt) set protocol specific EP options

Passive EP

Getname	(~getsockname) return EP address
Listen	Start listening for connection requests
Reject	Reject a connection request

Active EP

CM	Connection establishment ops, usable by connection-oriented and connectionless endpoints
MSG	2-sided message queue ops, to send and receive messages
RMA	1-sided RDMA read and write ops
Tagged	2-sided matched message ops, to send and receive messages (conceptual merge of messages and RMA writes)
Atomic	1-sided atomic ops
Triggered	Deferred operations initiated on a condition being met

Shared Receive Queue

Receive Post buffer to receive data



Event Queue

Read	Retrieve a completion event, and optional source endpoint address data for received data transfers
Read Err	Retrieve event data about an operation that completed with an unexpected error
Write	Insert an event into the queue
Reset	Directs the EQ to signal its wait object when a specified condition is met
Strerror	Converts error data associated with a completion into a printable string

EQ Group

Poll	Check EQs for events
Wait	Wait for an event on the EQ group

Completion Counter

Read	Retrieve a counter's value
Add	Increment a counter
Set	Set / clear a counter's value
Wait	Wait until a counter reaches a desired threshold

Memory Region

Desc	(~lkey) Optional local memory descriptor associated with a data buffer
Key	(~rkey) Protection key against access from remote data transfers

Architectural Semantics

Need refining

- Progress
- Ordering - completions and data delivery
- Multi-threading and locking model
- Buffering
- Function signatures and semantics

Once defined, object and interface semantics cannot change – semantic changes require new objects and interfaces

Progress

- Ability of the underlying implementation to complete processing of an asynchronous request
- Need to consider **ALL** asynchronous requests
 - Connections, address resolution, data transfers, event processing, completions, etc.
- HW/SW mix

All(?) current solutions require significant software components

Progress - Proposal

- Support two progress models
 - Automatic and implicit (name?)
- Separate operations as belonging to one of two progress domains
 - Data or control
 - Report progress model for each domain

Progress - Proposal

- Implicit progress
 - Occurs when reading or waiting *on EQ(s)*
 - Application can use separate EQs for control and data
 - Progress limited to objects associated with selected EQ(s)
 - App can request automatic progress
 - E.g. app wants to wait on native wait object
 - Implies provider allocated threading

Ordering - Completions

- Outbound
 - Is any ordering guarantee needed?
 - ‘Sync’ call completion guarantees all selected, previous operations issued on an endpoint have completed
- Inbound
 - Ordering only guaranteed for message queue posted receives

Ordering – Data Delivery

- Interfaces may imply specific ordering rules
 - E.g. Tagged – “If a sender sends two messages in succession to the same destination, and both match the same receive, then this operation cannot receive the second message if the first one is still pending. If a receiver posts two receives in succession, and both match the same message, then the second receive operation cannot be satisfied by this message, if the first one is still pending.”

Ordering – Data Delivery

- Required ordering specified by application
 - [read | write | send] after [read | write | send]
 - RAR, RAW, WAR, WAW, SAW
- Ordering may differ based on message size
 - E.g. size \geq MTU

Needs more analysis with a
formal proposal

Multi-threading and Locking

- Support both thread safe and lockless models
- Lockless – based on MPI model
 - Single – single-threaded app
 - Funneled – only 1 thread calls into interfaces
 - Serialized – only 1 thread at a time calls into interfaces
 - Are all models needed?
- Thread safe
 - Multiple – multi-threaded app, with no restrictions

Buffering

- Support both application and network buffering
 - Zero-copy for high-performance
 - Network buffering for ease of use
 - Buffering in local memory or NIC
 - In some case, buffered transfers may be higher-performing (e.g. “inline”)
- Registration option for local NIC access
 - Migration to fabric managed registration
- Required registration for remote access
 - Specify permissions