# New Accelerations for Parallel Programming

**HPC Advisory Council – Spain 2012**
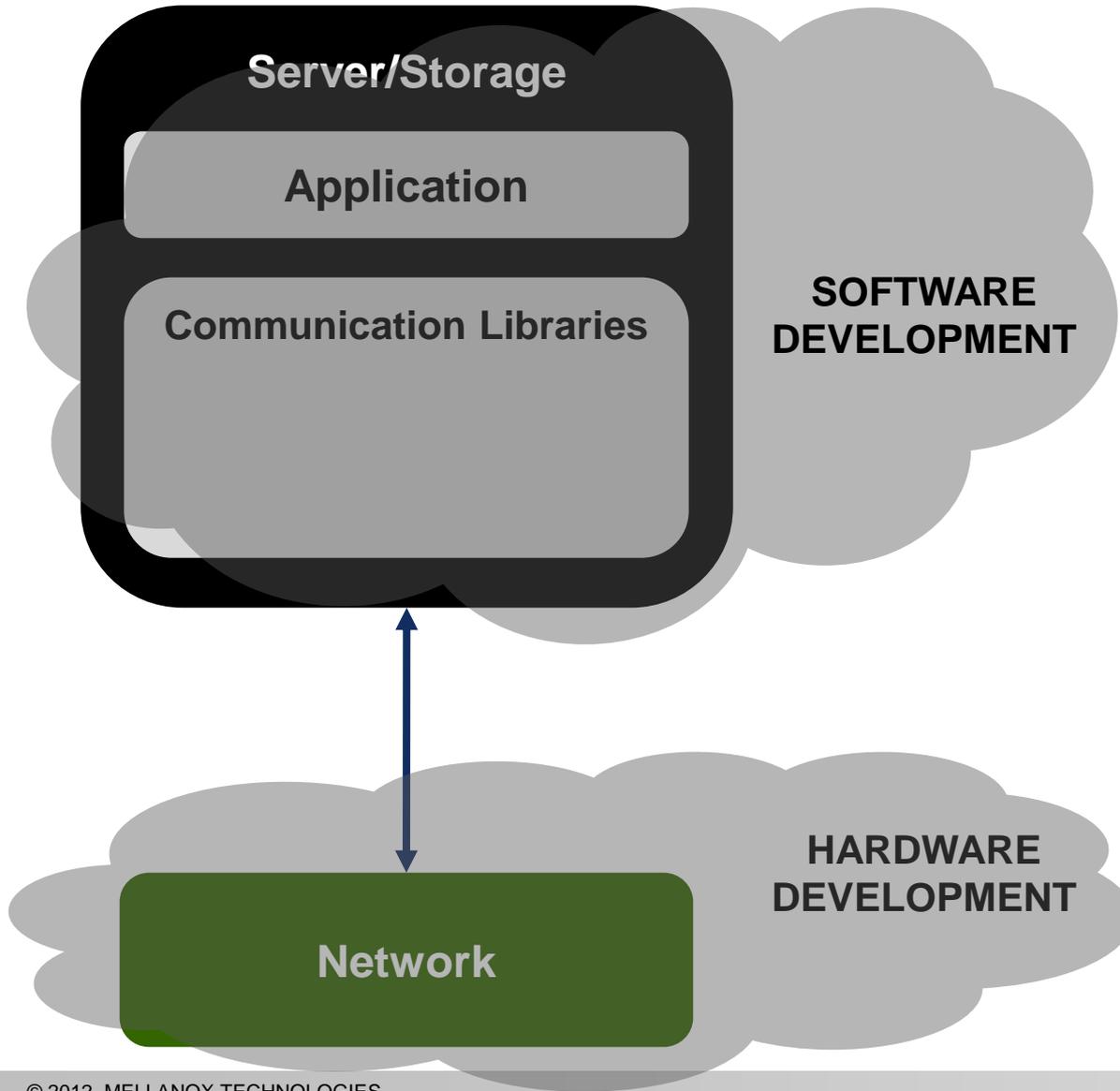
**Todd Wilde – Director of Technical Computing and HPC**
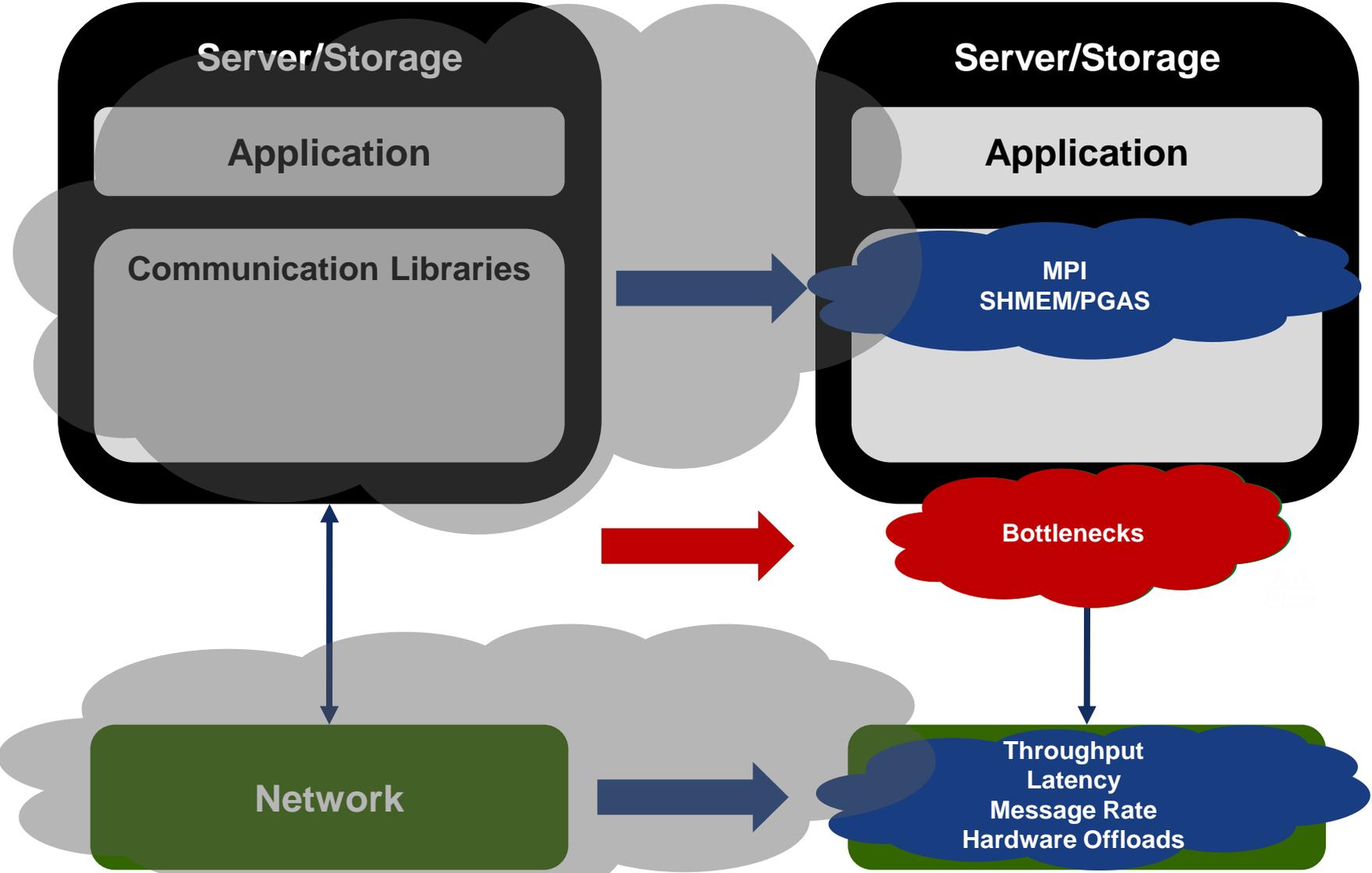
**todd@mellanox.com**

# Agenda

- The Co-Design Architecture for Parallel Programming Languages

- An Introduction to PGAS Languages

- FCA – Fabric Collective Accelerations

- Mellanox/HP Collaboration On InfiniBand Scalability for One-Sided Communication
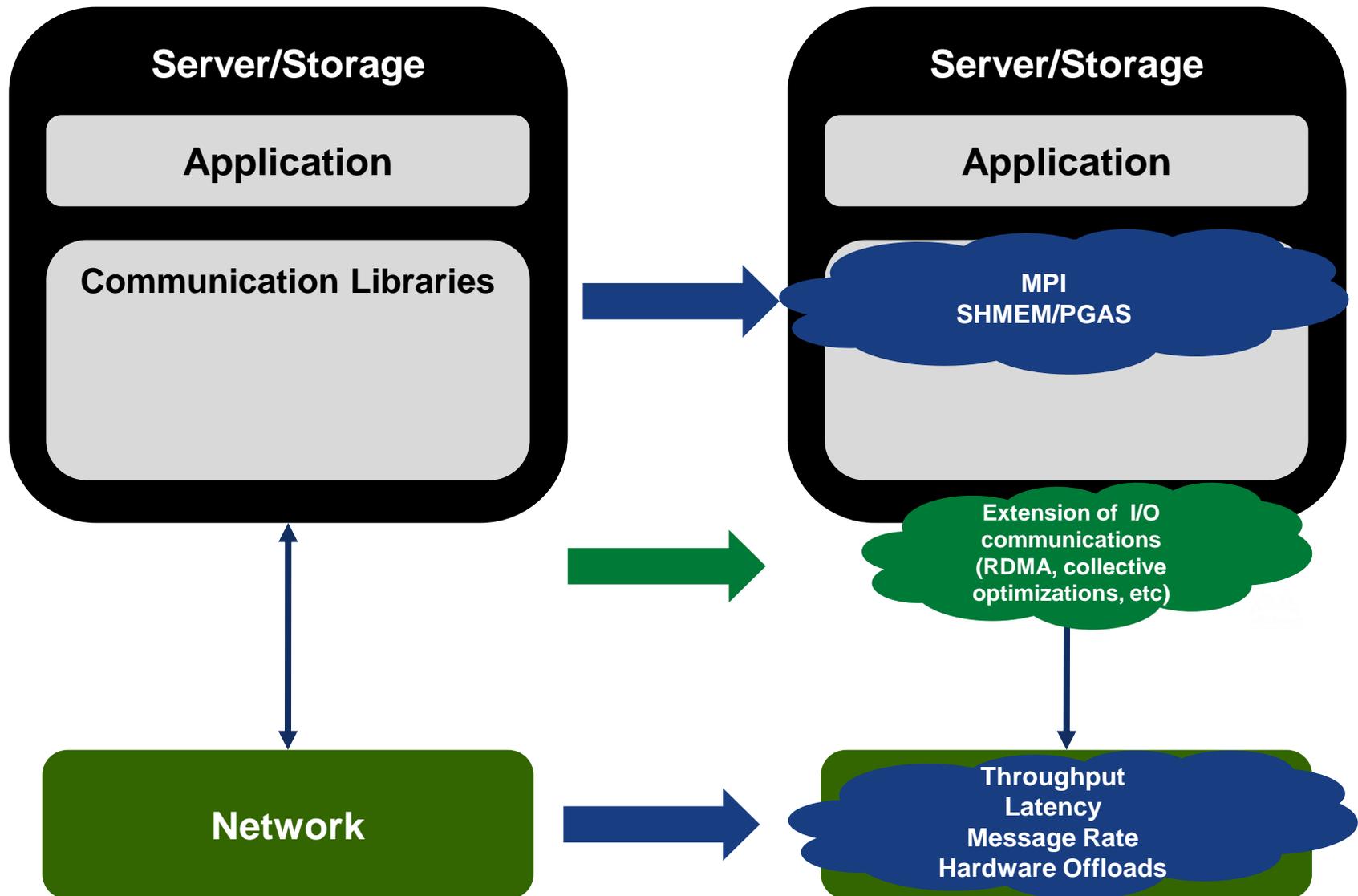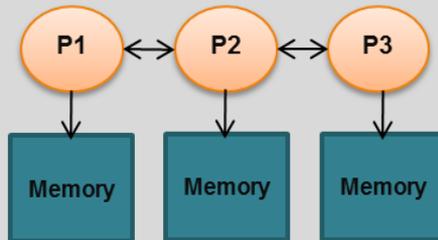
# The Co-Design Architecture
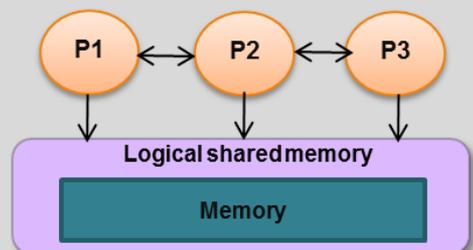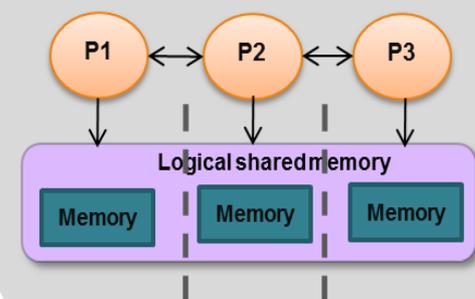
# The Co-Design Architecture

# The Co-Design Architecture

# Mellanox ScalableHPC Accelerate Parallel Applications

# An Introduction to PGAS Languages

# Introduction to PGAS Languages

- PGAS – Partitioned Global Address Space – Best of both worlds
  - Message passing and shared memory methods

- Explicitly-parallel programming model with SPMD parallelism like MPI
  - Fixed at program start-up, typically 1 thread per processor

- Global address space model of memory
  - Allows programmer to directly represent distributed data structures

- Address space is logically partitioned
  - Local vs. remote memory (two-level hierarchy)

- SHMEM is being used/proposed as a lower level interface for PGAS implementations.

- Multiple PGAS languages: UPC (C), CAF (Fortran), Titanium (Java)

# PGAS Language Example - UPC

- UPC – Unified Parallel C
  - Open source compiler from LBNL/UCB
  - Currently operates over Infiniband Verbs RC connections via GASnet interface

- Utilizes a distributed shared memory programming model
  - Similar to traditional shared memory model, but allows for data locality
  - Distributed shared memory is divided into partitions where each $M_i$ is associated with thread $TH_i$.

- Features include:
  - Simple statements for remote memory access
  - Minimization of thread communication overhead by exploiting data locality

**UPC**

$TH_0$   $TH_1$   $TH_{n-1}$

$M_0$   $M_1$   $M_{n-1}$

**Shared Address Space**

# UPC Memory Model

- UPC memory is divided into private and shared space

- Each thread has its own private space in addition to a portion of the shared space

- A UPC *shared* pointer can access any locations in the shared space.   A *private* pointer may reference only addresses in it's private space or local portion of shared space.
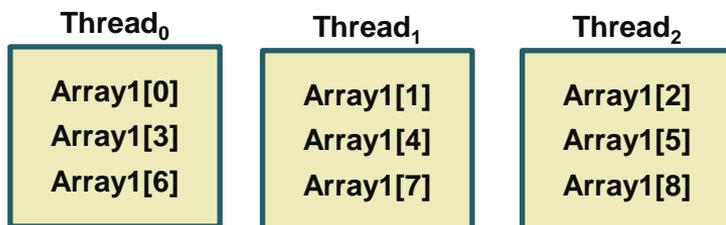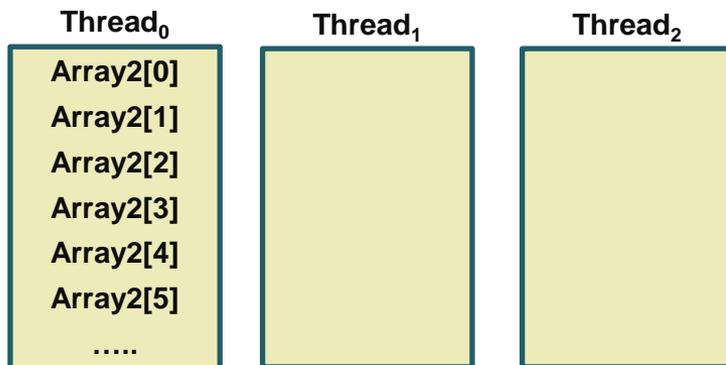
**Affinity to thread 0**

| Thread$_0$ | Thread$_1$ | | Thread$_{N-1}$ |
|---|---|---|---|
| | | Shared Address Space | |
| Private$_0$ | Private$_1$ | …….. | Private$_{N-1}$ |

# UPC Memory Model – Shared Memory Usage

- Must use shared qualifier in variable declaration
  - *shared [block_size] type variable_name* : means that variable_name is distibuted across memory space in the span of block_size per thread.

- Examples:
  - shared [1] int array1[N]

| Thread$_0$ | Thread$_1$ | Thread$_2$ |
|---|---|---|
| Array1[0] | Array1[1] | Array1[2] |
| Array1[3] | Array1[4] | Array1[5] |
| Array1[6] | Array1[7] | Array1[8] |

  - shared [N] int array2[N]

| Thread$_0$ | Thread$_1$ | Thread$_2$ |
|---|---|---|
| Array2[0] | | |
| Array2[1] | | |
| Array2[2] | | |
| Array2[3] | | |
| Array2[4] | | |
| Array2[5] | | |
| ….. | | |

# UPC Execution Model

- ## A number of threads working independently in a SPMD fashion

  - Number of threads specified at compile-time or run-time; available as program variable **THREADS**

  - **MYTHREAD** specifies thread index ($0$..$THREADS-1$)

  - **upc_barrier** is a global synchronization: all wait

  - **upc_forall** is similar to for-loop but also indicates which thread will run the loop iteration

- ## There are two compilation modes

  - Static Threads mode:
    - THREADS is specified at compile time by the user
    - The program may use THREADS as a compile-time constant

  - Dynamic threads mode:
    - Compiled code may be run with varying numbers of threads

# My first UPC program

```c
#include <upc_relaxed.h>
#include <stdio.h>

void main()
{
  if (MYTHREAD==0){
    printf("Rcv'd: 'Starting Execution' from THREAD %d\n",MYTHREAD );
  }

  printf("Hello World from THREAD %d (of %d THREADS)\n", MYTHREAD, THREADS);
}
```
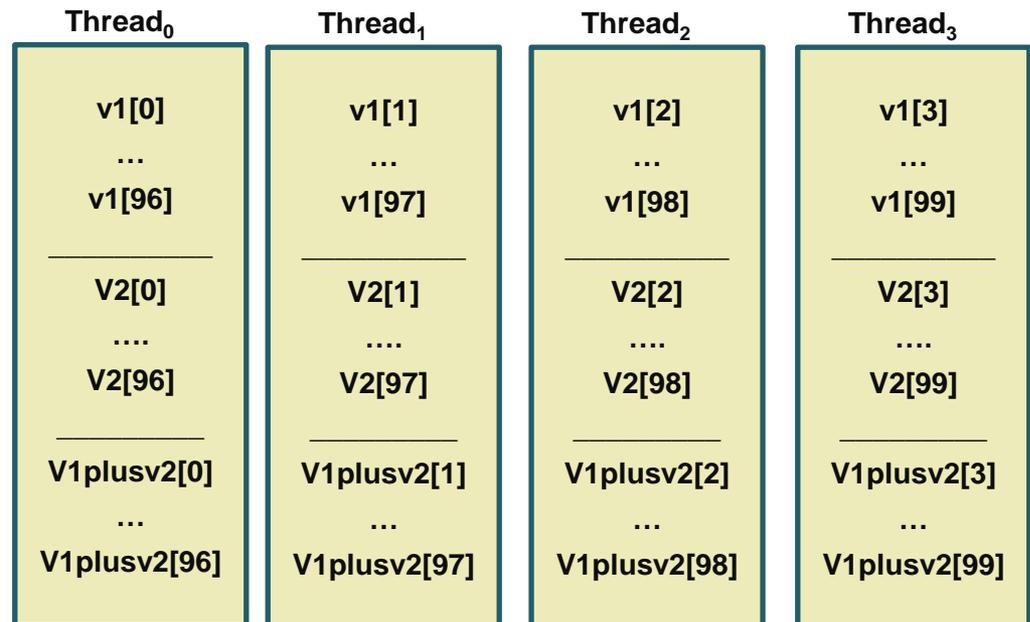
# UPC Programming – Array Copy Example

```
#include<upc_relaxed.h>
#define N 100
shared int v1[N], v2[N], v1plusv2[N];

void main()
{
    int i;
    for(i=0;i<N;i++)
        if(MYTHREAD==i%THREADS)
            v1plusv2[i]=v1[i]+v2[i] ;
}
```

Owner computes

| Thread$_0$ | Thread$_1$ | Thread$_2$ | Thread$_3$ |
|---|---|---|---|
| v1[0] … v1[96] | v1[1] … v1[97] | v1[2] … v1[98] | v1[3] … v1[99] |
| V2[0] …. V2[96] | V2[1] …. V2[97] | V2[2] …. V2[98] | V2[3] …. V2[99] |
| V1plusv2[0] … V1plusv2[96] | V1plusv2[1] … V1plusv2[97] | V1plusv2[2] … V1plusv2[98] | V1plusv2[3] … V1plusv2[99] |

```
#include<upc_relaxed.h>
#define N 100
shared int v1[N], v2[N], v1plusv2[N];

void main()
{
  int i;
  upc_forall(i=0;i<N;i++,i)
    if(MYTHREAD==i%THREADS)
      v1plusv2[i]=v1[i]+v2[i] ;
}
```
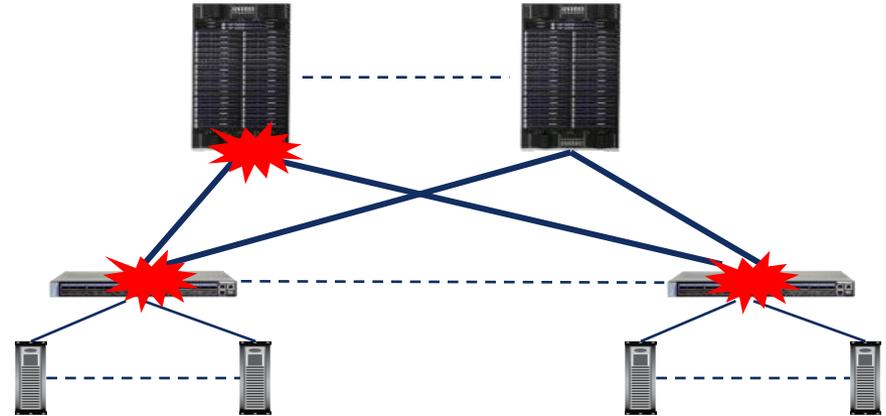
Owner computes

# Fabric Collective Accelerations
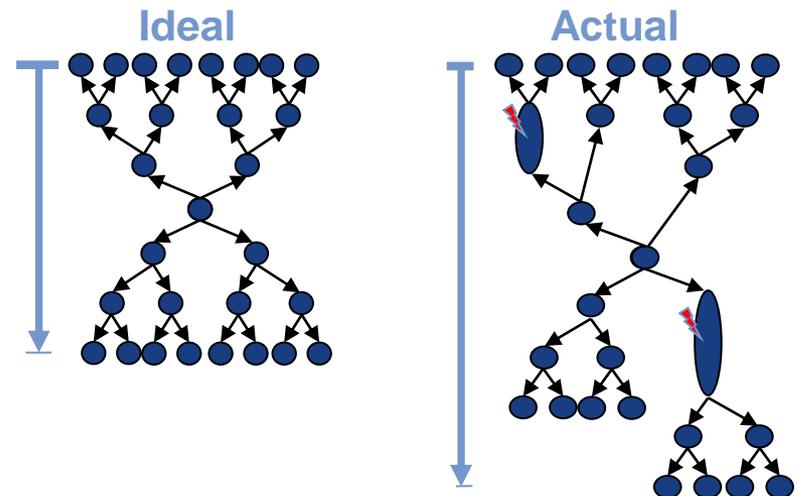
# What are Collective Operations?

- Collective Operations are Group Communications involving all processes in job

- Synchronous operations
  - By nature consume many 'Wait' cycles on large clusters

- Popular examples
  - Barrier
  - Reduce
  - Allreduce
  - Gather
  - Allgather
  - Bcast

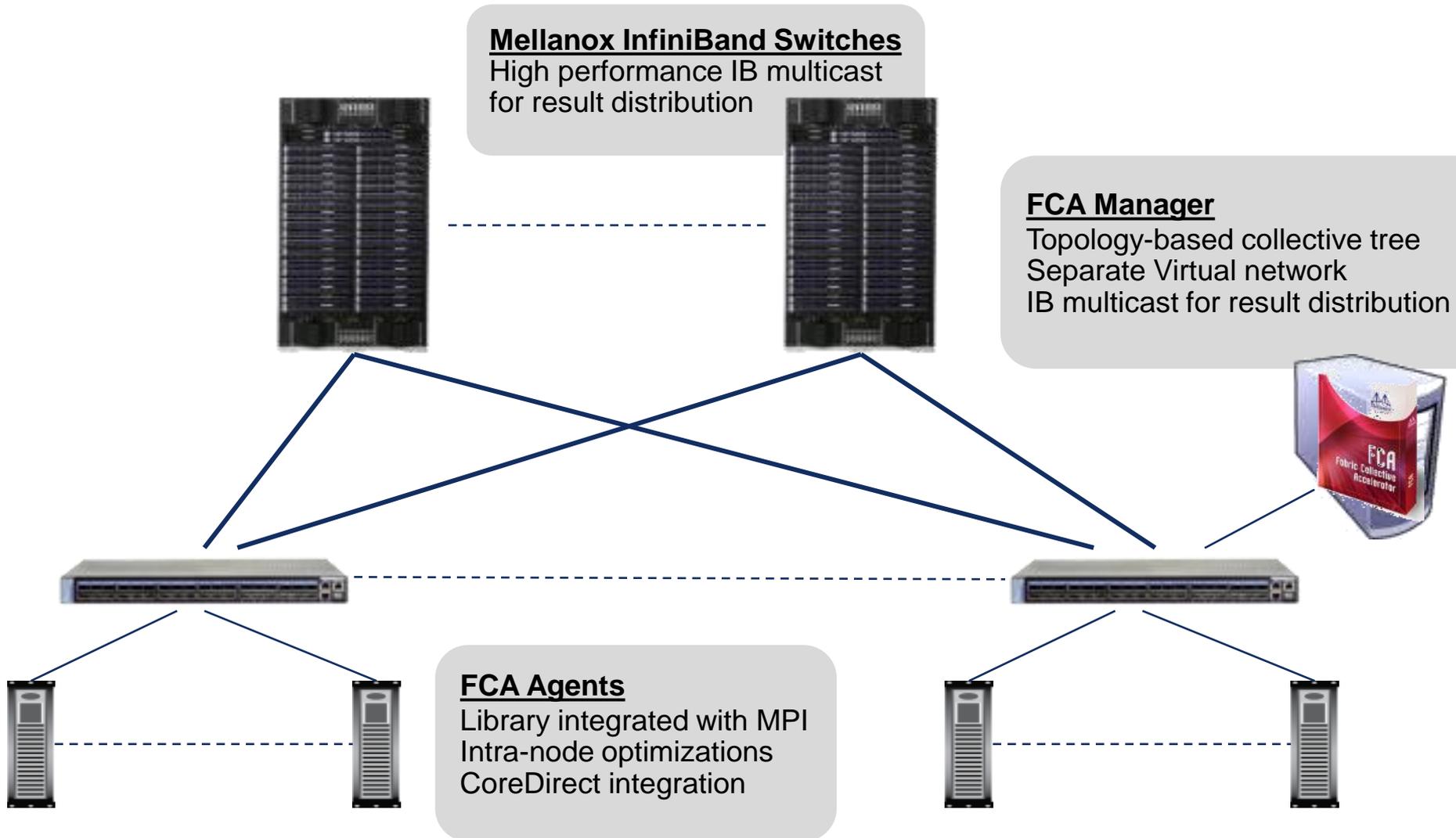- Collective algorithms are not topology aware and can be inefficient



- Congestion due to many-to-many communications

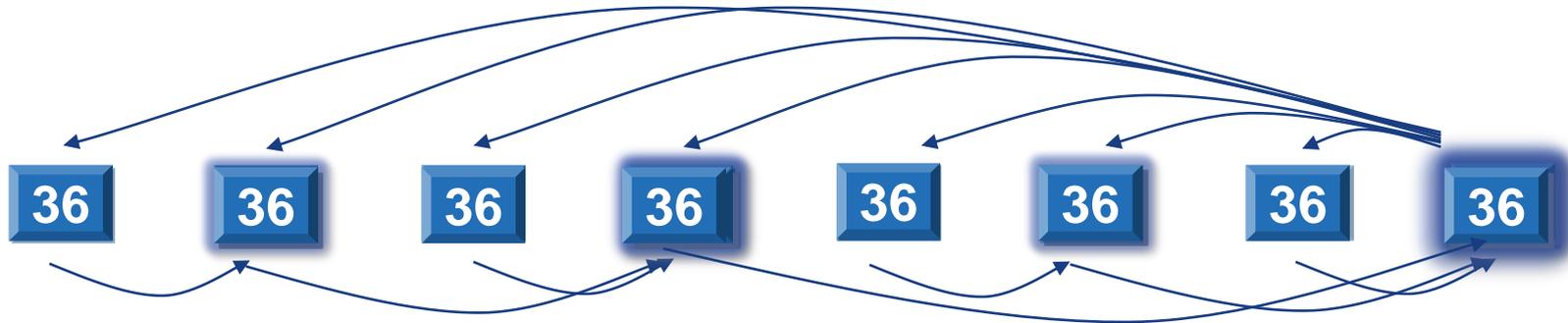- Slow nodes and OS jitter affect scalability and increase variability

**Ideal**          **Actual**

# Mellanox Fabric Collectives Accelerations (FCA)

**Mellanox InfiniBand Switches**
High performance IB multicast
for result distribution

**FCA Manager**
Topology-based collective tree
Separate Virtual network
IB multicast for result distribution

**FCA Agents**
Library integrated with MPI
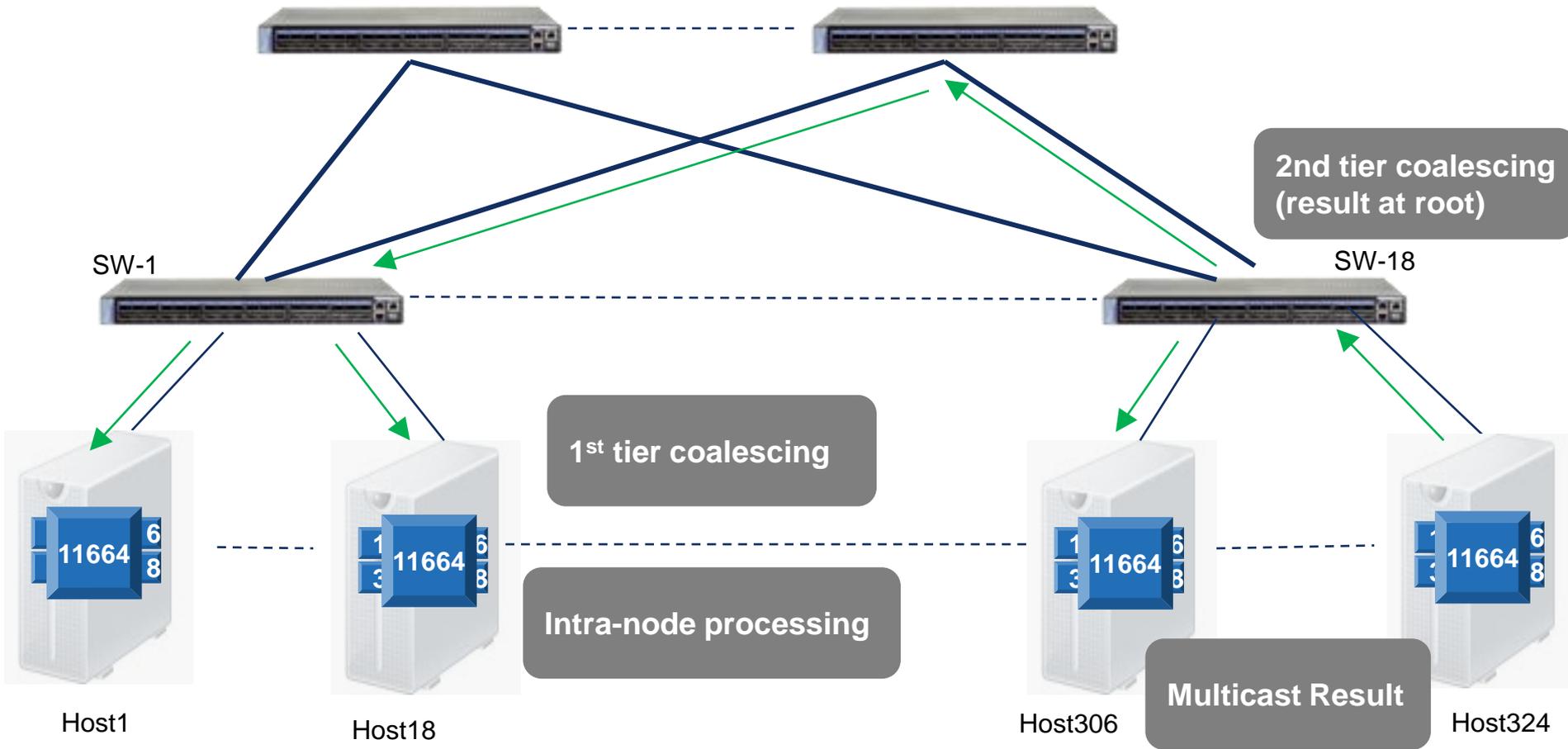Intra-node optimizations
CoreDirect integration

# Collective Example – Allreduce using Recursive Doubling

- Collective Operations are Group Communications involving all processes in job



- A 4000 process Allreduce using recursive doubling is 12 stages

# Scalable Collectives with FCA



SW-1

SW-18

**2nd tier coalescing (result at root)**

**1st tier coalescing**

**Intra-node processing**

**Multicast Result**

Host1

Host18

Host306

Host324

# Thank You
HPC@mellanox.com