



 **Los Alamos**
NATIONAL LABORATORY
EST. 1943

Delivering science and technology
to protect our nation
and promote world stability

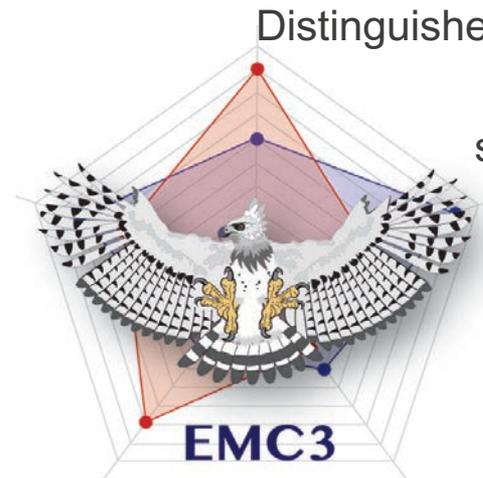
- Responsible for care for the majority of the US Nuclear Stockpile
- A Serious Enduring Mission
- Not a solved problem space or anywhere close

2019-UCF-UCX-F2F

2019-F2F
Austin Tx.

Stephen Poole

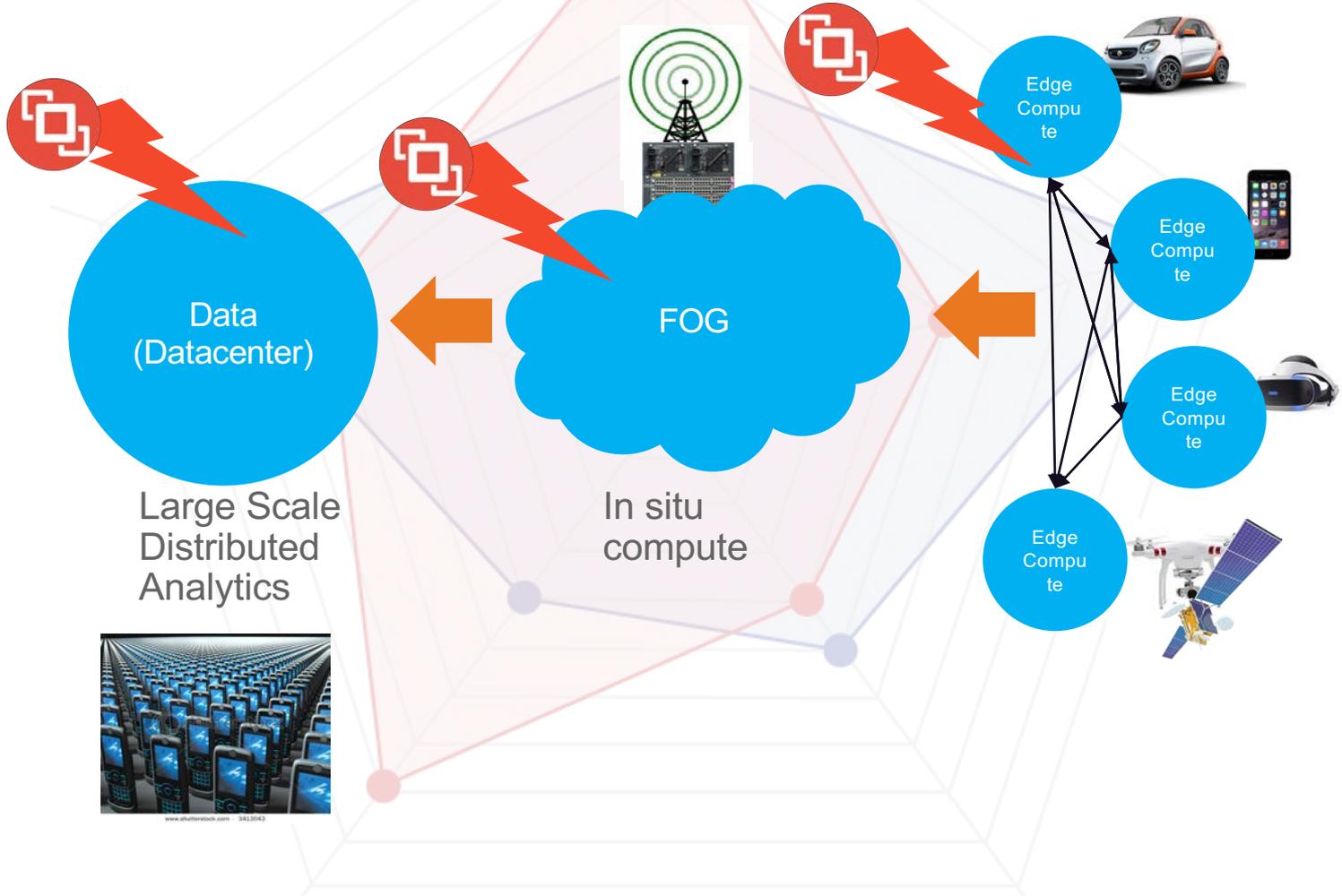
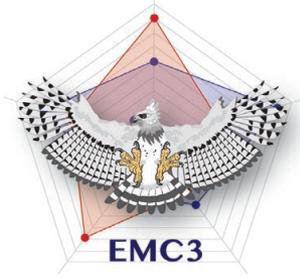
Distinguished Senior Scientist
Chief Architect
swpoole@lanl.gov



LA-UR-19-32281



Early idea for graphic

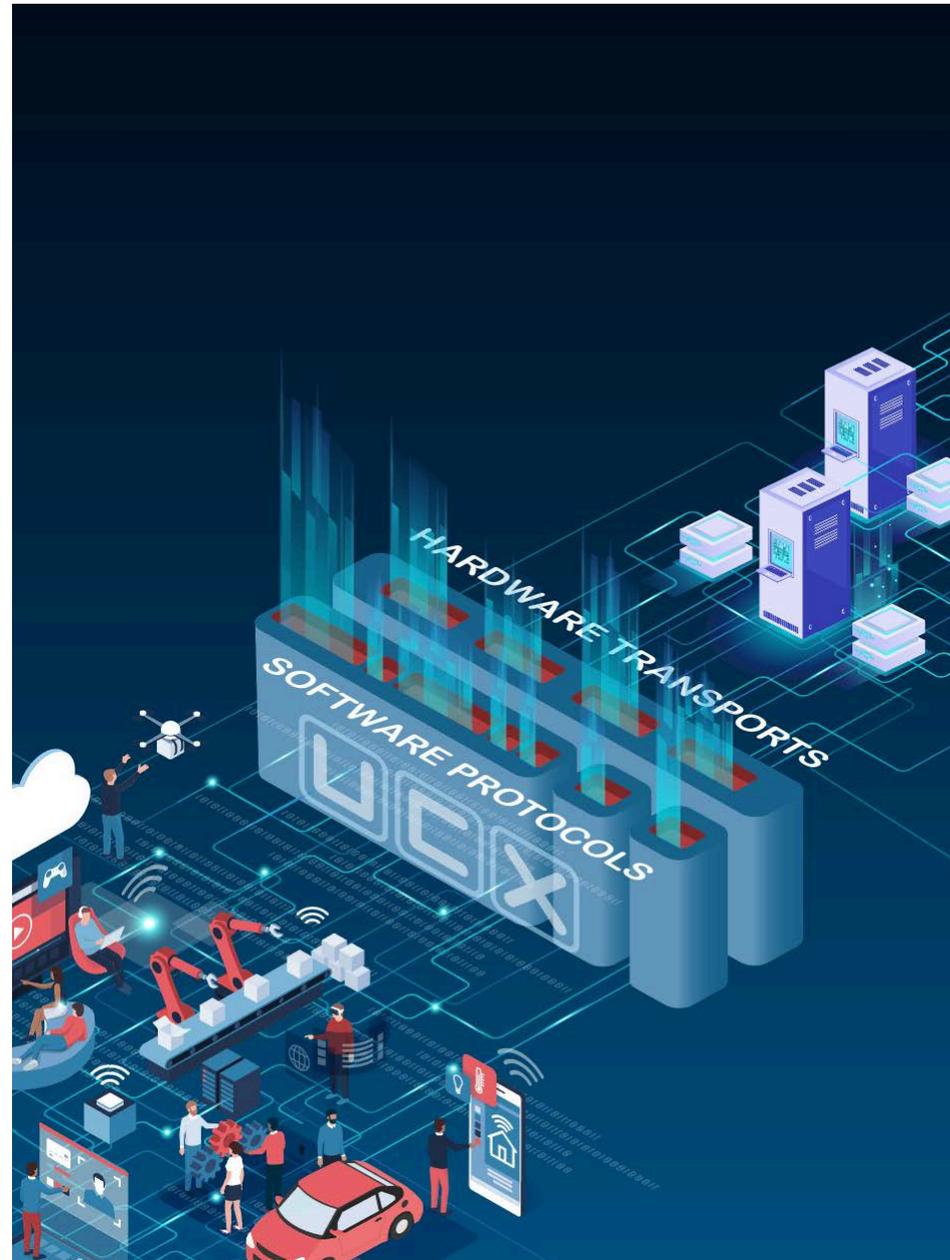




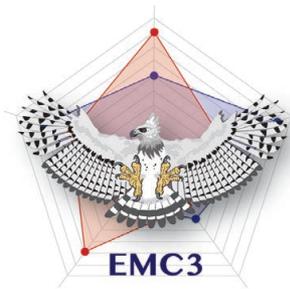
Unified Communications X

An open-source, exascale-ready communications framework

- Solves decades-old problem in high-performance computing (HPC)
- Frees developers from hardware-specific implementations and laborious porting efforts
- Simplifies deployment of advanced research tools, regardless of system complexity
- Advances fields of artificial intelligence, machine learning, deep learning, and internet of things



Advanced Micro Devices,
Argonne National Laboratory,
Arm Ltd., Mellanox Technologies,
NVIDIA, Stony Brook University,
Oak Ridge National Laboratory,
Rice University



Los Alamos has been at the forefront of developing high-speed network interconnects/technologies

- Los Alamos developed one of the first networks, named Hydra, to allow common access to the five CDC machines, 1972
- High-Speed Parallel Interface (HSPI) for inter-computer communication, 50 Mbit/s, 1979-1982
- High-Performance, Parallel Interface (HIPPI), the first gigabit network, 1987
- Gigabyte System Network - GSN 1990's
- Infiniband interconnect came out of ASCI work in the late 1990's
- Optical interconnects started by ASCI ~2000
 - Analysis of optical switches
- QKD – Early 2000's (Spun off)
 - Free space quantum optics
- **In Situ using COTS**
 - Prime candidate for UCX
 - SNIC API (SNAPI)

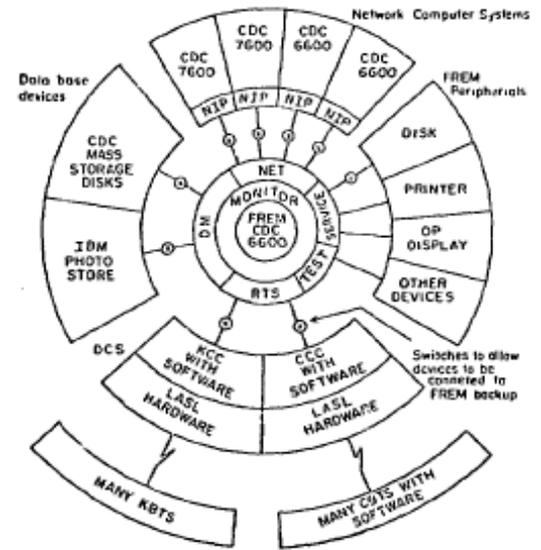
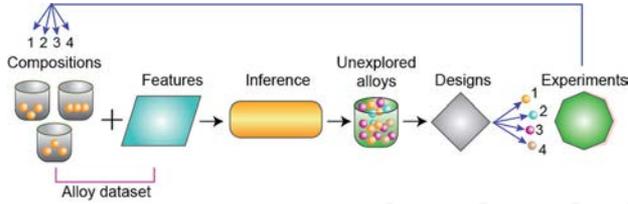
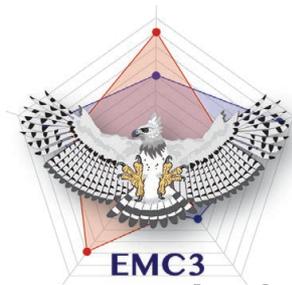


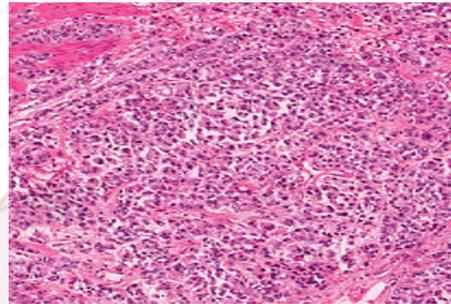
Fig. 1. Hydra system diagram.
Hydra network design, 1972

High-speed interconnects are mostly a standard commercial off-the-shelf technology now

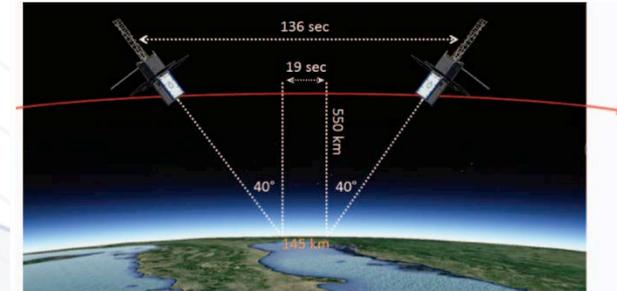
Data analytics is emerging organically and rapidly across many programs



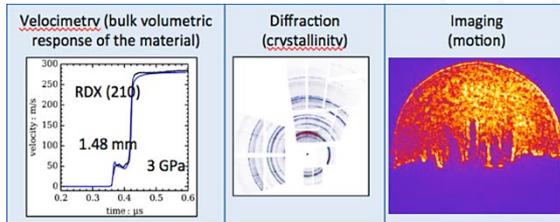
Machine Learning
Accelerating Discovery
of New Materials, old
Materials



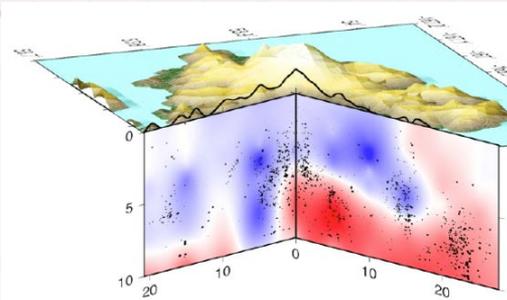
Bioinformatics/
Emergent Diseases



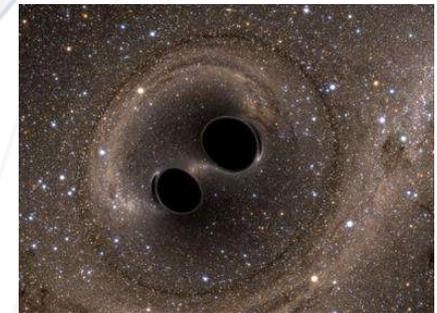
Constellation of CubeSats,
Carrying Ultra-Compact
Spectral Sensors



Real-time Adaptive
Acceleration of Dynamic
Experimental Science

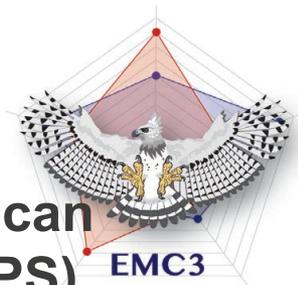


Critical Stress
in Subsurface
Energy Dynamics

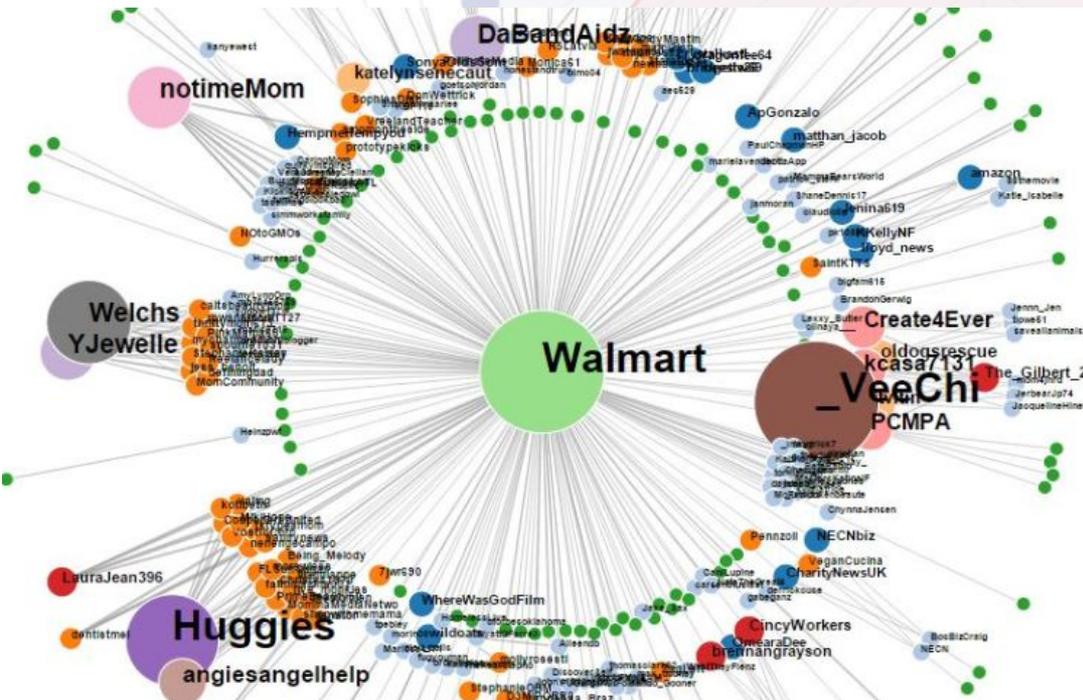


Gravitational Wave
Emissions from
Colliding Black Holes
(LIGO)

Graph Analytics and Event Simulation

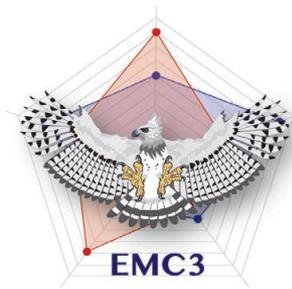


- Communications intensive, event/data driven, working set can exceed main and often have short words, mostly read. (GUPS)
 - Worst case is every lookup is not cached, in some random memory location somewhere in the cluster/system
 - Non-cached latency is memory latency ~40-100 cycles
 - Worse is across a low latency network ~1000 cycles (non-deterministic)



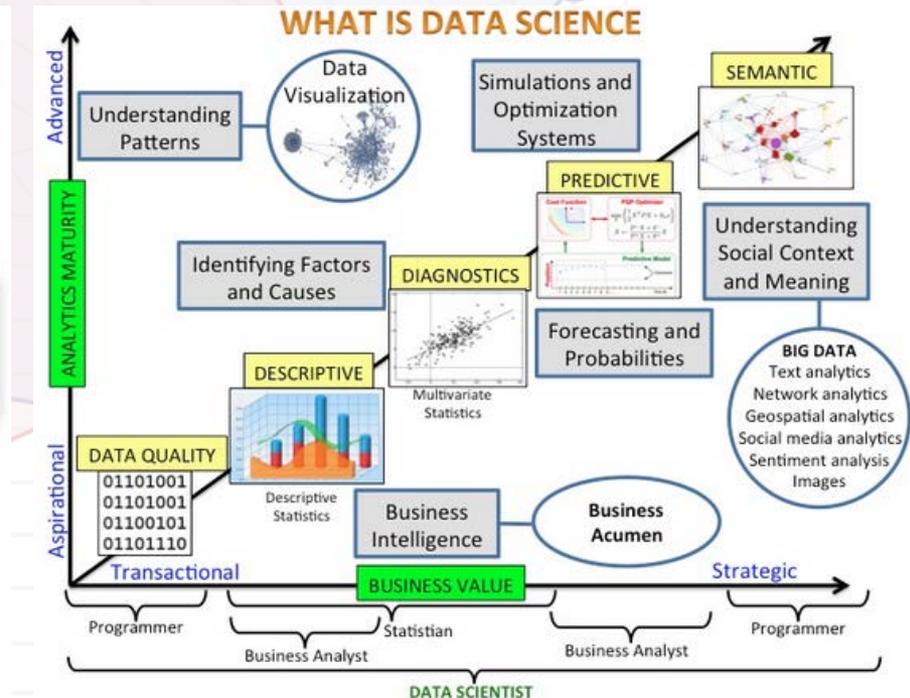
Follow relationships,
needles in haystacks, etc.

Many forms of Data Science: Analytics, Streaming Analytics, and some forms of Machine Learning/Deep Learning/AI (potential in situ)



- Often IO Intensive, Data Parallel works very well, often with small word sizes, mostly read dominant
 - Usually parallelizes well
 - IO and Network Bandwidth are sometimes limiters
 - We are working in these areas, so HW has to address issues.

Machine Learning basics



Dated results from Grand Unified File Index has improved by about 5X since this was done

Trinity



Query

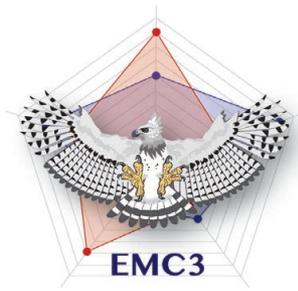


	Find all files in scratch1 as uid 67890		Find all files in lustre scratch1		Find all files in scratch1 and NFS home as uid 67890	
	POSIX	GUFI	POSIX	GUFI	POSIX	GUFI
Files	22,771,329	22,509,652	119,296,067	118,509,899	-	22,522,140
Dirs	240,736	237,759	5,541,230	5,523,153	-	239,603
Time (s)	531.6	14.5	11,309	134.2	-	14.9
Files/s	42,835	1,553,956	10,548	883,413	-	1,511,553
	36x		84x			

Index load

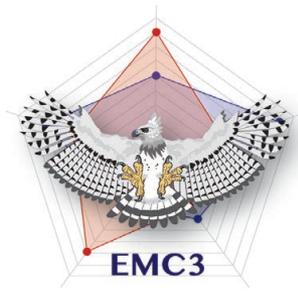
- 118,509,899 files from lustre filesystem into GUFI tree: 148.9s
 - 795,902 files/s
- 13,229,405 files from NFS home filesystem into GUFI tree: 38.4s
 - 344,515 files/s

Why does LANL care about UCX?



- **LANL is actively engaged in developing “in-network computing”.**
- **All of the previously mentioned apps can take advantage of UCX**
- **UCX offers a continuum of deployment spaces.**
 - NIC/HCA
 - Phones/IoT
 - Standard applications
 - OpenSHMEM, MPI, PGAS
 - KVS
 - I/O applications
- **UCX offers portability across a very diverse portfolio of systems**
 - We have lots
- **UCX is easily adapted to new and development HW**
 - We are developing new systems and applications/analytics
- **We are currently porting some of the LANL mini-apps to UCX**
- **LANL and key USG partners are supporting UCX for external developers**
- **LANL is working with O&G customers wrt UCX (See EMC³)**
- **LANL has on-going development projects with Mellanox**
- **LANL is fully supportive of an Open SmartNIC API (OpenSNAPI)**
 - Arm, Mellanox, Broadcom, others?

EMC³ Areas of Interest (where are we doing these things?)



• **Networking**

- Next Gen Network Requirements/Design
- Apps/computing/programmability in the network

• **I/O & Storage futures**

- Data protection at extremes
- IO Forwarding
- Data motion innovation

• **Resilient Computing**

- Characterization/prediction

• **High Performance Data Analytics Systems**

- True performance benchmarks/measurement of HPDA systems

• **Inexact computing**

- Characterization for exploration

• **System SW**

- Next Gen Systems mgmt. (boot/launch/manage/etc.)
- Leveraging container tech

• **HPC Environments**

- Launch, Run Time, Monitoring, Tools innovation

• **Application of ML/DL/AI to HPC**

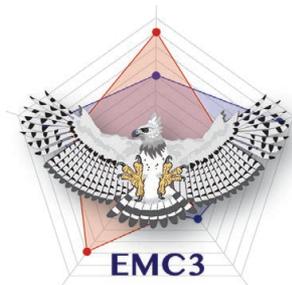
• **Processor/Memory Complex**

- Balanced Application focused
- Power requirements
- Balanced Performance
- Scaling

• **Benchmarking and Simulation of Systems**

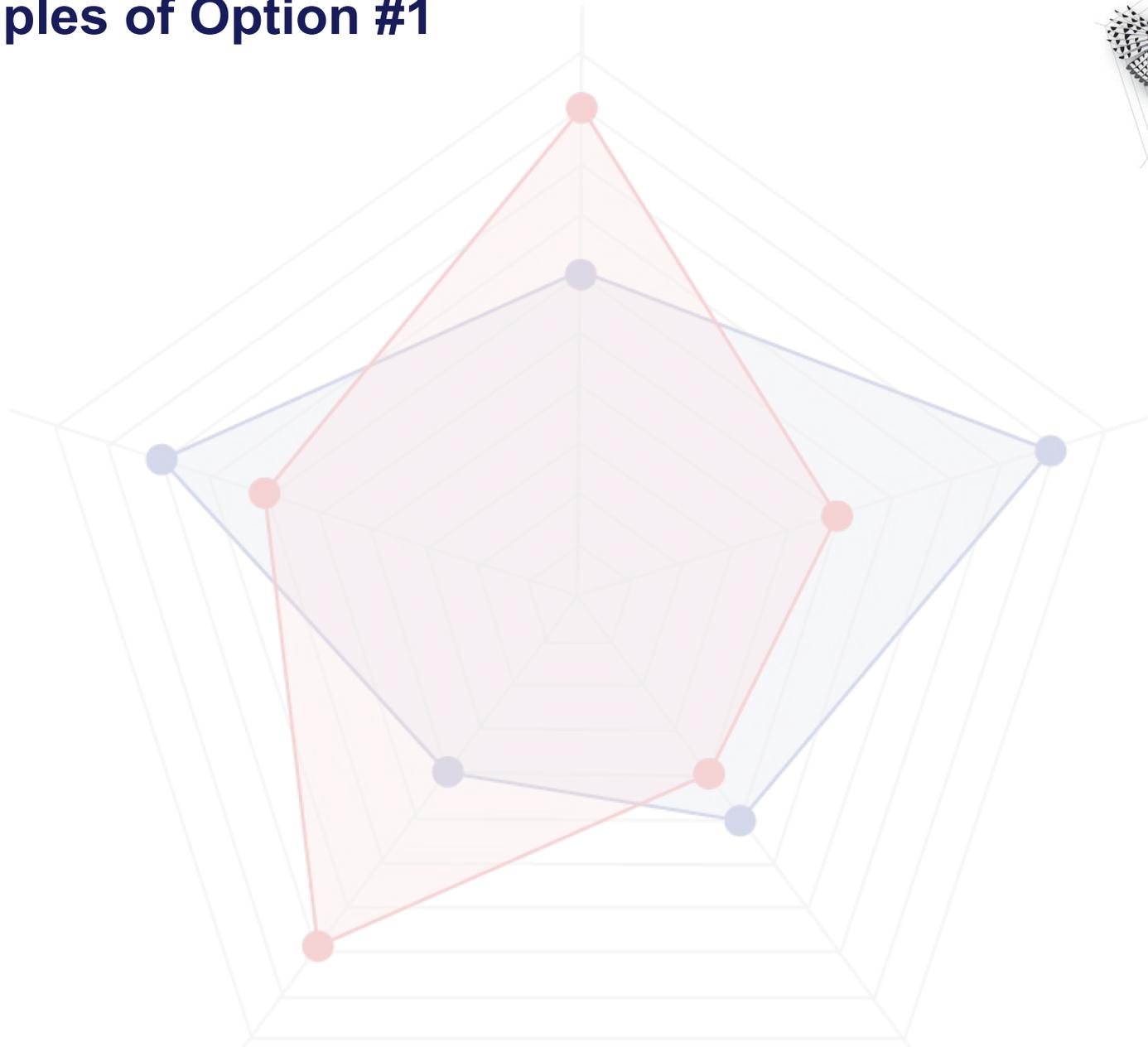
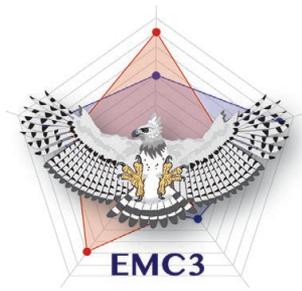
- Use of Benchmarks/Simulation against codes to achieve balanced forward progress

OpenSNAPI (potential options/ideas)

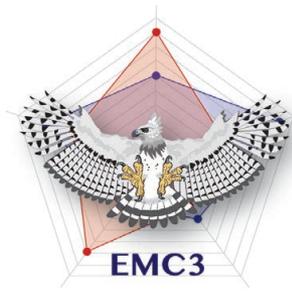


- **Option 1 – Compile and deploy application/parts on SmartNIC**
 - Already doable. (Summer students)
 - Assumes the NIC has an OS, memory, processor. (Just another node)
- **Option 2 – Direct calls into “resident” library**
 - Communication via PCIe/Network interface
 - UCX under OpenSHMEM, OpenFAM
 - CUDA like?
- **Option 3 – Direct “put” of function calls into NIC Space**
 - FPGA like.
 - #PRAGMA and function calls (Xilinx)
 - Assumes driver can set up drop/return spaces.
 - IXP2800 like
 - PCIe or Network? Both?
- **Option 4 – New language?**
 - P4
 - <http://lccn.cs.technion.ac.il/wp-content/uploads/2019/02/HeavyHitter-Detection-with-P4-ver-1.2.pdf>
 - Micro-C
 - Define one of our own? (LONG PROCESS)
- **Which one do we start with first?**
- **Do we need to assume the NIC is “host” resident? (Sever or serverless)**
 - Yes
 - PCIe or Network comms
 - No
 - Network comms
- **What about programs/functions to the Switch**

Examples of Option #1



TEST PROGRAM #1



```
#include <stdlib.h>
#include "getcpu.h"

#define DEBUG 1

#define HOST_PROCS 20
#define NIC_PROCS 32
#define HOST_START 0
#define HOST_END (HOST_PROCS - 1)
#define NIC_START (HOST_END + 1)
#define NIC_END (NIC_START + NIC_PROCS)

long psync(SHMEM_REDUCE_SYNC_SIZE);

int main(int argc, char** argv)
{
    shmem_init();
    int np = shmem_n_pes();
    int id = shmem_my_pe();
    int *host_bc = (int*) shmem_malloc(sizeof(int));
    int *nic_bc = (int*) shmem_malloc(sizeof(int));
    char host_name[100];
    gethostname(host_name, 100);
    int cpu_id = get_cpu_id();

    #ifdef DEBUG
    int i, j;
    if(host_bc == NULL || nic_bc == NULL)
    {
        printf("Symmetric pointer == NULL\n\n");
        exit(1);
    }
    for(i = 0; i < np; i++)
    {
        if(shmem_pe_accessible(i) != 1)
            printf("PE %d is unable to access PE %d\n", id, i);
        if(shmem_addr_accessible(nic_bc, i) != 1)
            printf("PE %d is unable to access nic_complete PE %d\n", id, i);
        if(shmem_addr_accessible(host_bc, i) != 1)
            printf("PE %d is unable to access host_complete PE %d\n", id, i);
    }
    shmem_barrier_all();
    #endif

    printf("Hello #1 from process %03d out of %03d, hostname %s, cpu_id %d, host_bc %d, nic_bc %d\n",
        id, np, host_name, cpu_id, *host_bc, *nic_bc);

    shmem_barrier_all();

    if(id == HOST_START)
    {
        *host_bc = 1;
    }
    else if(id == NIC_START)
    {
        *nic_bc = 2;
    }

    shmem_barrier_all();

    shmem_broadcast32(host_bc, host_bc, 1, HOST_START, HOST_START, 0, np, psync);
    shmem_broadcast32(nic_bc, nic_bc, 1, NIC_START, NIC_START, 0, np, psync);

    printf("PE %d completed broadcasts\n", id);

    shmem_barrier_all();

    printf("Hello #2 from process %03d out of %03d, hostname %s, cpu_id %d, host_bc %d, nic_bc %d\n",
        id, np, host_name, cpu_id, *host_bc, *nic_bc);

    shmem_free(host_bc);
    shmem_free(nic_bc);
    shmem_finalize();

    return 0;
}
```

The Jupiter nodes have 10 cores each. Each core has 2 threads, resulting in a total of 20 Pes per node. Each Jupiter node also has a Bluefield card, which has 16 cores. For this test program we used 10 cores on Jupiter nodes 8 and 9, and all 16 in each of the Bluefield cards, for a total of 52 PEs.

The goal of the program was to explore communication amongst all PEs, understand the relationship between the host nodes and Bluefield cards, and figure out the mapping of PE ranks to physical resources.

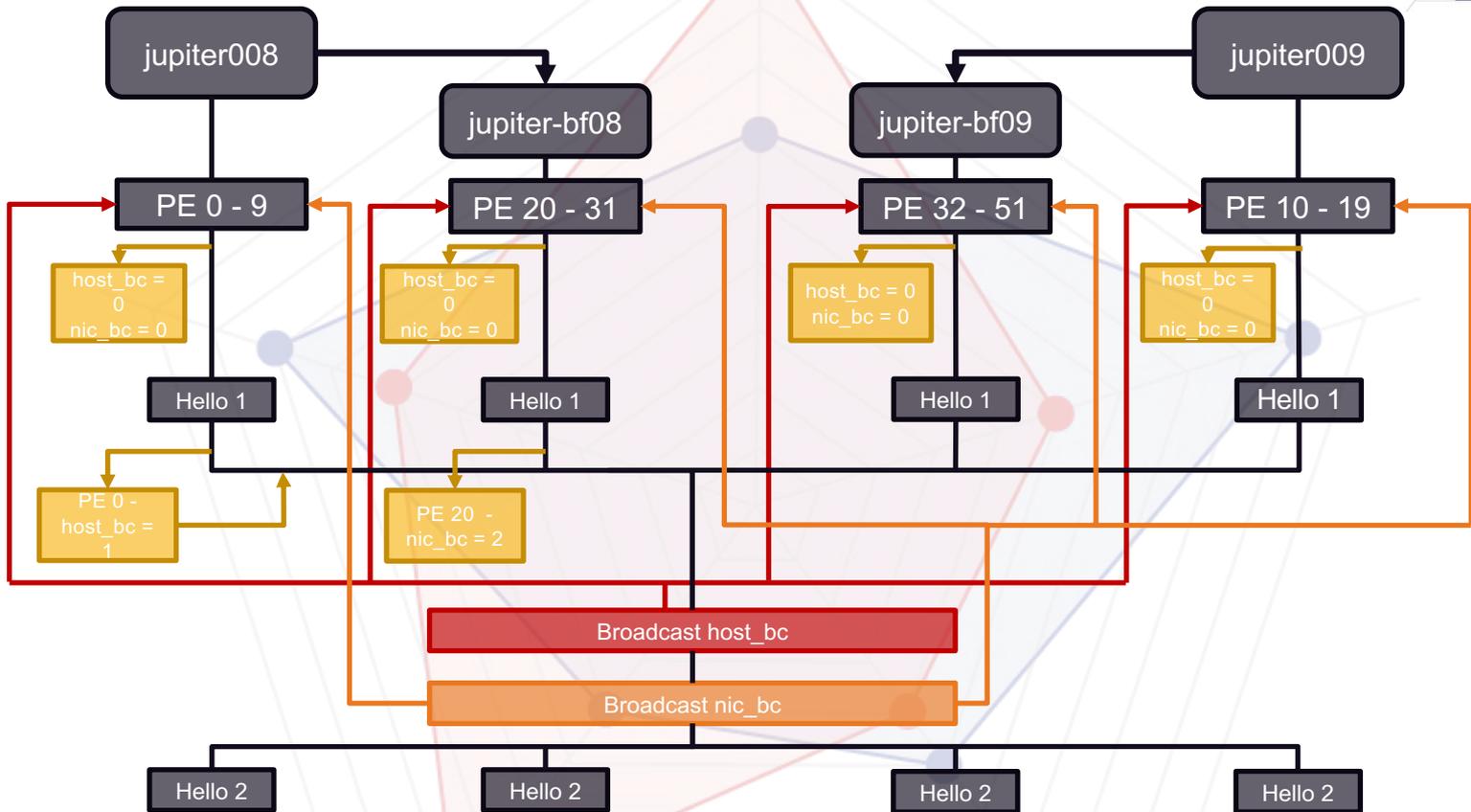
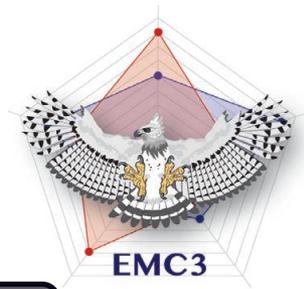
Note that in order to run the program in parallel, the code must be on both the host and BF the cards. The program was run using an appfile, which specified how many cores each node and BF card should use. It was noted that the order of the entries in the appfile, determines the PE ranks.

In the program itself, two symmetric objects were declared: host_bc, and nic_bc, both initialized to 0. These variables were utilized to check for proper synchronization and communication between PE's. A message was printed from each PE showing the initial variable values. The variable host_bc was set to 1 by PE 0 (the first host PE), and the variable nic_bc was set to 2 by PE 20 (the first BF PE).

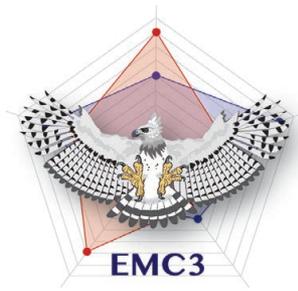
A broadcast was then sent to update the values of the two variables on all the PEs. Finally, once the broadcasts were complete another message was printed to verify that communication was successful.

Note that the print statements are buffered by the filesystems. However, the broadcast values in the output prove that in reality the four nodes are working in parallel. Therefore, it was concluded that all PEs remained synchronized throughout the process.

The following graphic demonstrates the program's flow with regard to individual processes.



Test Program #2



```
#include <stdio.h>
#include <unistd.h>
#include <shmem.h>

int main(void)
{
    int me, npes;
    shmem_init();
    int *sum = (int*) shmem_malloc(sizeof(int));
    me = shmem_my_pe();
    npes = shmem_n_pes();
    char hostname[100];
    gethostname(hostname, 100);

    if (me > 0) {
        shmem_int_atomic_add(sum, me, 0);
    }

    shmem_barrier_all();

    if (me == 0) {
        printf("From hostname %s: Sum from 1 to %d = %d\n", hostname, npes - 1, *sum);
    }

    shmem_finalize();

    return 0;
}

"sum2n.c" 82L, 2655C
```

For the second test program we used 20 cores on both Jupiter nodes 8 and 9, and all 16 in each of the Bluefield cards, for a total of 72 PEs.

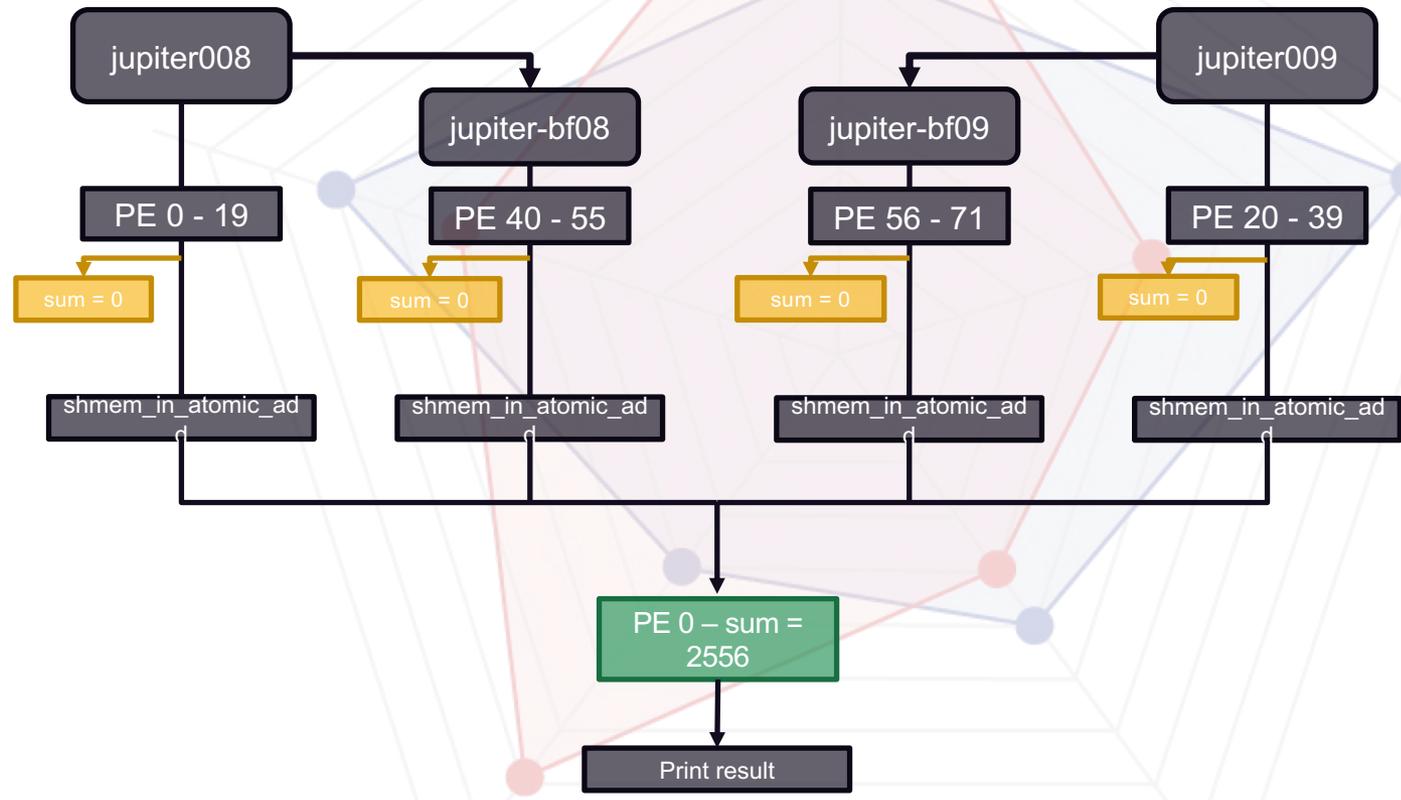
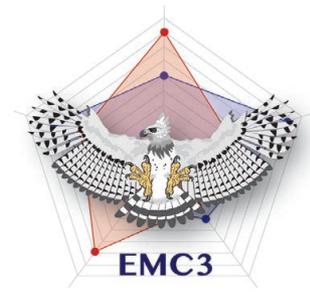
The goal of this simple program was to test the use of atomic operations across PEs on the Jupiter hosts and the Bluefield cards.

Just as the previous test, this program was run using an appfile.

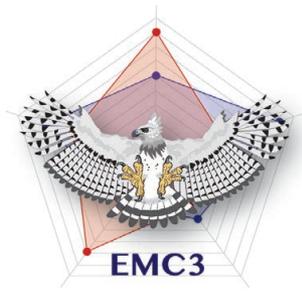
This program atomically adds all PE rank values and stores the result into the symmetric object “sum” on PE 0. Once all rank values have been added, the result is printed.

Based on the output, it was concluded that atomic operations are correctly executed across all PEs.

The following graphic demonstrates the program’s flow with regard to individual processes.

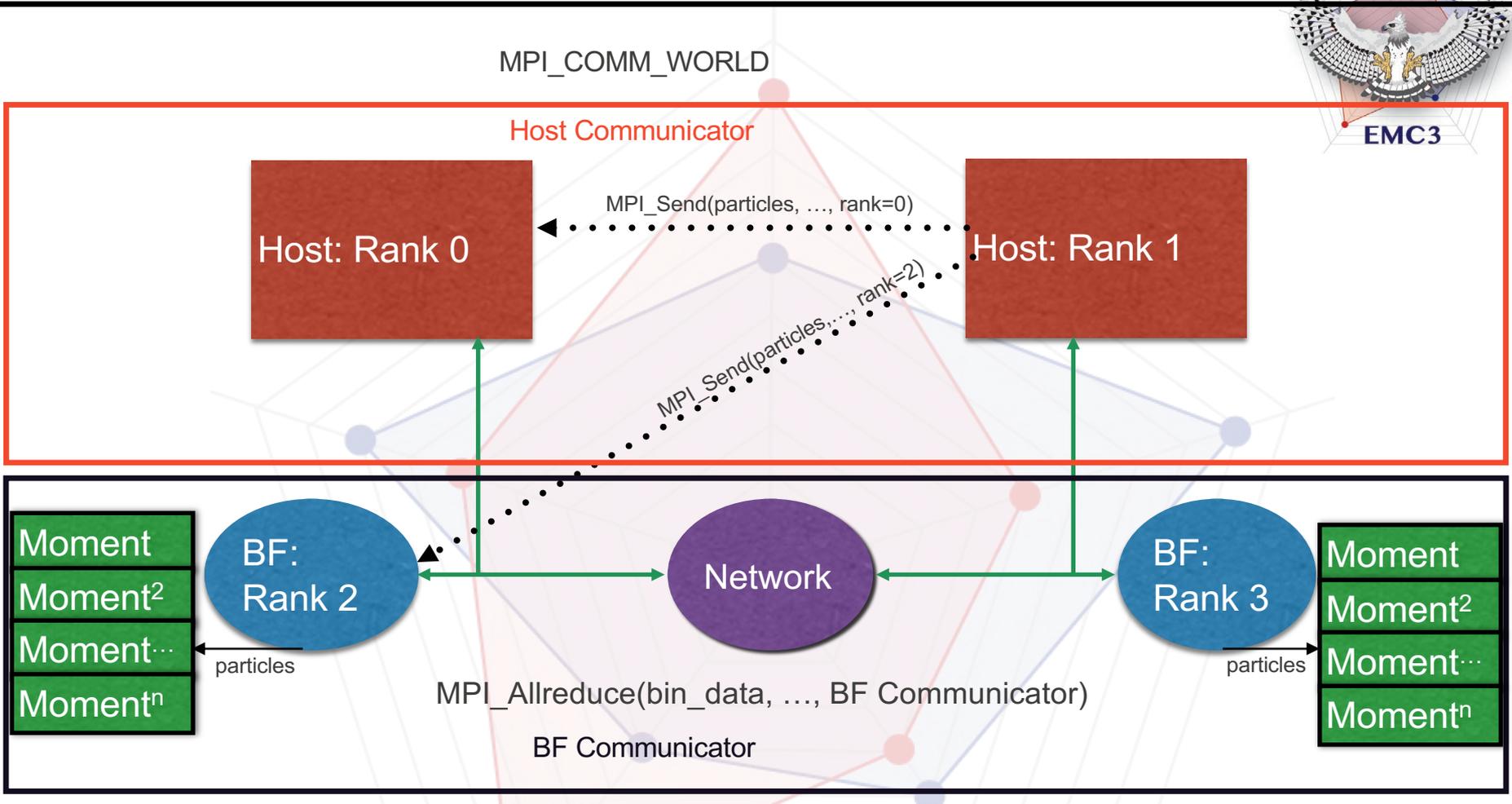
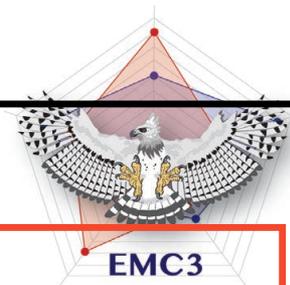


VPIC Problem (Courtesy Brad Settlemyer)



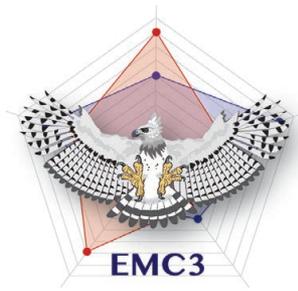
- Problem Statement

- Consider the problem of tracking a spatially correlated energy peak through space, i.e. tracking an energy wave. In a particle-in-cell code, such as LANL's VPIC, the energy values are calculated using trillions of small particles (approximately 32 Bytes each) that move between processes as the simulation progresses. In order to identify all energy peaks at a timestep it is necessary to understand the energy distribution of the particles. The approach most commonly used in practice is to post-process the output data sets and construct energy histograms. The primary challenge in constructing an energy histogram is that the histogram bin widths are typically impossible to estimate in advance. To determine appropriate bin widths during post-processing requires using computational resources similar to the original simulation in order to efficiently post-process all of the simulation output. Another alternative is to simply sample some small percentage of the particles and estimating an empirical distribution based on the sampled particles. This technique is problematic because simulations are typically constrained by memory use, and both the all-to-all communications and buffering required to sample the energy space are expensive to achieve in practice.



- BF's are assigned an MPI Rank like any other host
- Particles are exchanged with MPI when they cross boundaries
- BF is also sent particles, which are then binned
- BF calculates Moments for each particle
- BF reconstructs quantiles of particle distribution from Moments

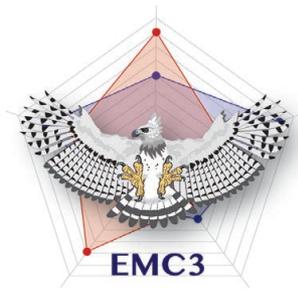
Option #2 (Some potential instructions)



```
SNICresult status = snicInit(0);
status = snicDeviceGet(&dev, 0);
status = snicDeviceComputeCapability(&major,&minor,dev);
snicGetDeviceCount(&device_count);
snicSetDevice( device );
snicGetDevice(&device);
snicGetDeviceProperties(&deviceProp,device);
snicDeviceReset();
snicDriverGetVersion ( int* driverVersion )
snicDeviceGetAttribute ( int* pi, snicDEVICE attrib, snicDEVICE_dev )
snicDeviceGetName ( char* name, int len, snicDEVICE dev )
snicDeviceTotalMem ( size_t* bytes, snicDEVICE dev )
```

NVIDIA has MANY great potential examples for this mode.
I copied these from NVIDIA docs.

Option #3 (Some potential examples)



- <http://xillybus.com/tutorials/vivado-hls-c-fpga-howto-2>

- Can use #pragma

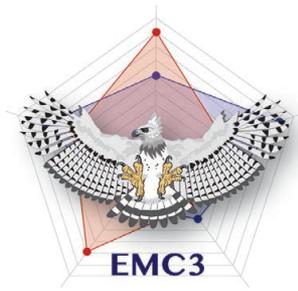
```
void snicBUS_wrapper(int *in, int *out)
{
  #pragma AP interface ap_fifo port=in
  #pragma AP interface ap_fifo port=out
  #pragma AP interface ap_ctrl_none port=return
  snic_puts("Hello, world\n"); // Handle input data
}
```

- Can write directly to device

- `fdr = open("/dev/snicBUS_read_32", O_RDONLY);`
- `fdw = open("/dev/snicBUS_write_32", O_WRONLY);`
- `write(fdw, (void *) &tologic, sizeof(tologic));`
- `read(fdr, (void *) &fromlogic, sizeof(fromlogic));`

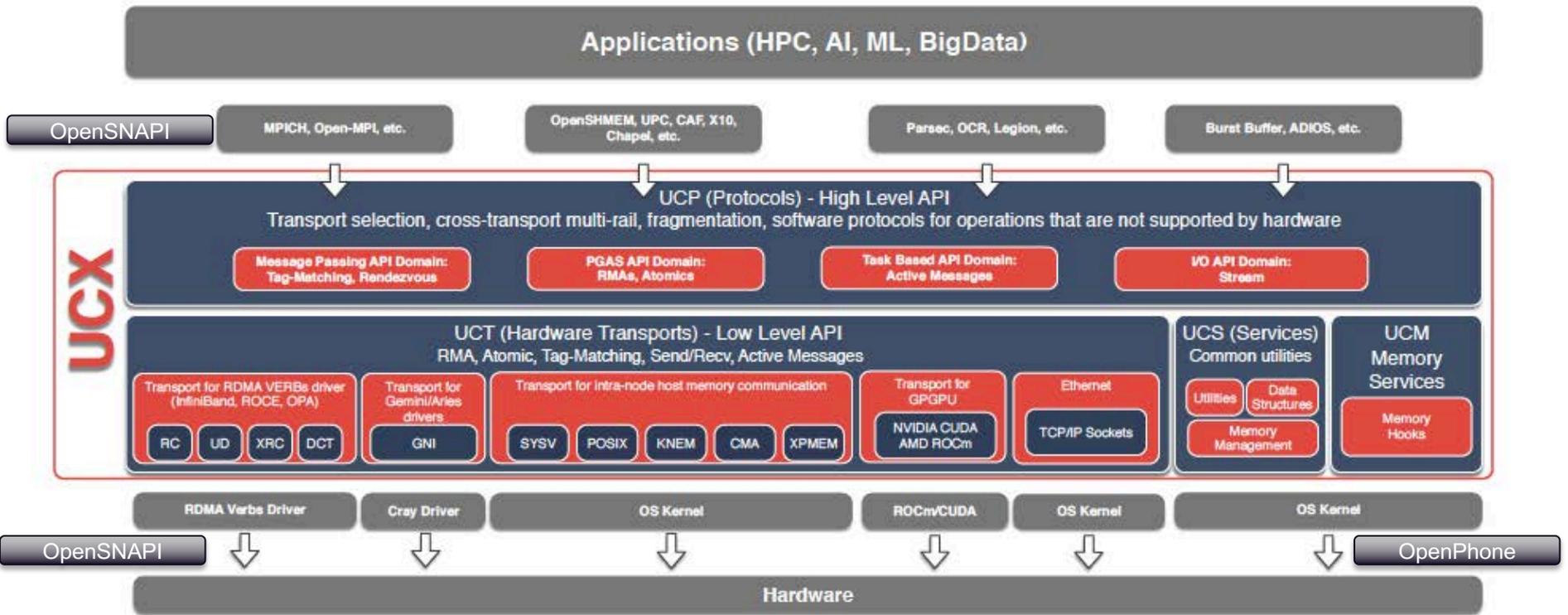
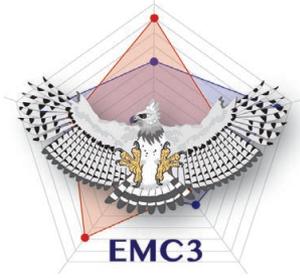
- Mellanox and others have many examples on how to use at this level.

Option #4 (Some potential examples)



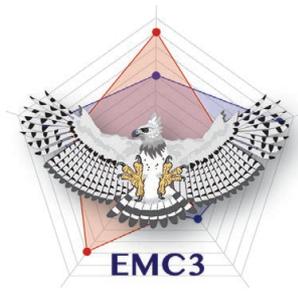
- Do we define our own? (Language)
- Do we adopt something like P4 or Micro-C
- UNO
 - SDN-controlled NF offload architecture (wisc)
 - <https://wiscr.cs.wisc.edu/papers/p506-le.pdf>
- Floem
 - a language, compiler, and runtime — for programming NIC-accelerated applications.
 - <https://www.usenix.org/system/files/osdi18-phothilimthana.pdf>
- What about Python?
- What can we learn from SDN?
 - Can we adopt/adapt a larger communities work?
- How do we make it fit best with UCX?
- <https://blog.mellanox.com/2018/09/why-you-need-smart-nic-use-cases/>
- Fund a summer intern(s) to develop summary paper(s) on:
 - Vendors
 - Languages
 - Functionalities
 - Classifications

Potential future UCX Diagram



Where would OpenSNAPI fit in the UCX “eco-system”?

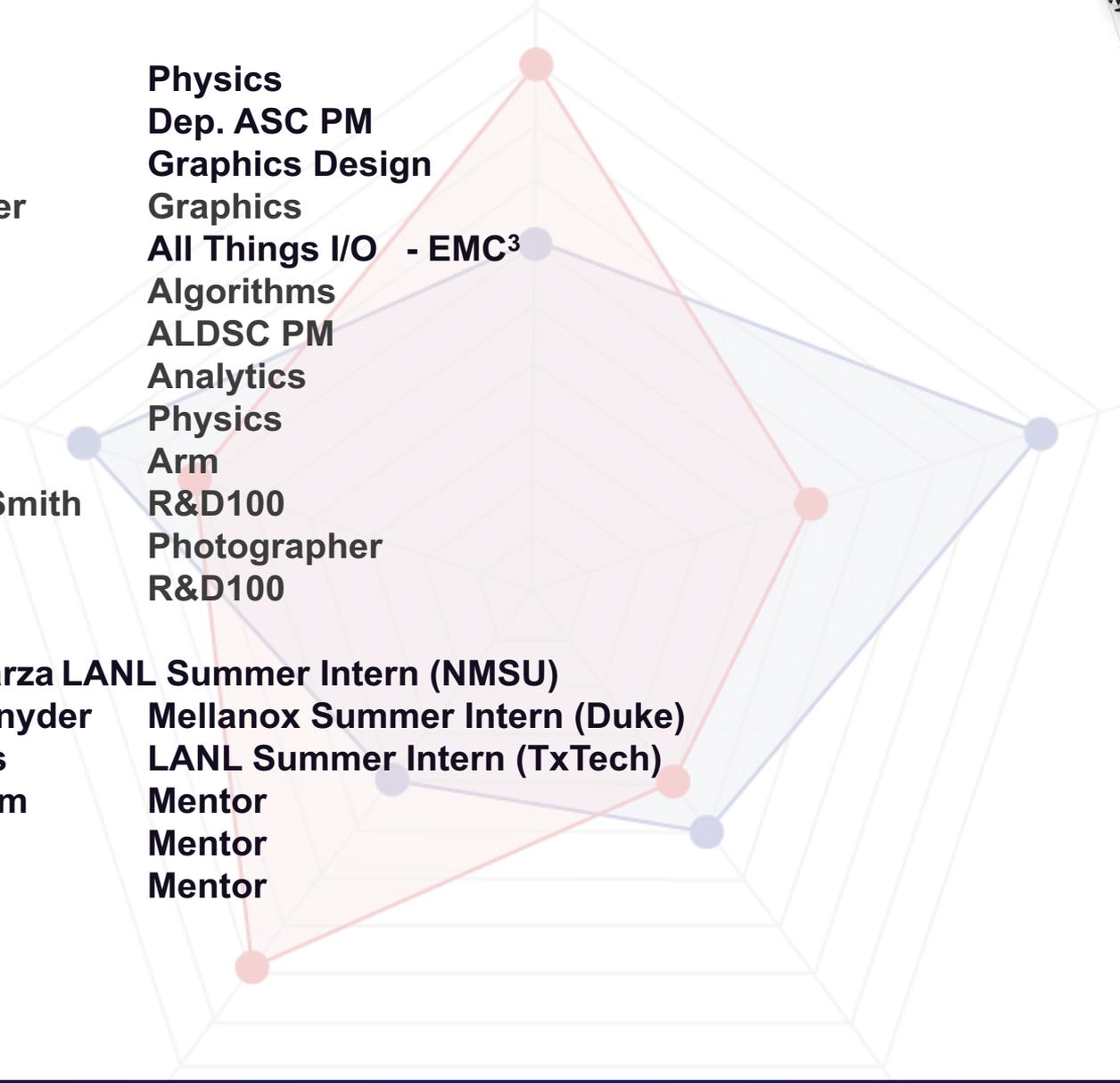
Contributors



- **Bill Archer**
- **Jerry Brock**
- **Jim Cruz**
- **AnnMarie Cutler**
- **Gary Grider**
- **Paul Henning**
- **Beth Kaspar**
- **Wendy Poole**
- **John Sarrao**
- **Pavel Shamis**
- **Janet Mercer-Smith**
- **Kevin Sutton**
- **Justin Warner**

Physics
Dep. ASC PM
Graphics Design
Graphics
All Things I/O - EMC³
Algorithms
ALDSC PM
Analytics
Physics
Arm
R&D100
Photographer
R&D100

- **Liliana A. Esparza LANL Summer Intern (NMSU)**
- **John "Jack" Snyder Mellanox Summer Intern (Duke)**
- **Brody Williams LANL Summer Intern (TxTech)**
- **Richard Graham Mentor**
- **Steve Poole Mentor**
- **Wendy Poole Mentor**



Questions?

