# Exploring Mellanox Bluefield SmartNICs as Accelerators for Heterogeneous Architectures

Brody Williams*, Liliana Aguirre-Esparza†, Wendy Poole‡, Stephen Poole‡

*Texas Tech University †New Mexico State University ‡Los Alamos National Laboratory

## Motivation

Limited advances in processor technology in recent years have forced researchers to explore alternative techniques to provide continued system performance improvements and facilitate further scaling. Many resulting approaches have shifted away from the strictly CPU-centric approach used in the past in favor of more heterogeneous architectures. These architectures often employ added computational components to support offloading computation from the CPU. Frequently, these components are also paired with their own local memory in order to minimize performance degradations associated with data movement.

In this work, we propose an extension to the heterogenous architecture paradigm using Mellanox Bluefield SmartNICs. These devices combine state of the art network controllers together with 16 ARM cores into a device that provides unique potential. Herein, we explore the feasibility of utilizing these SmartNICs as accelerators capable of offloading both communication routines, as well as computational kernels, from the CPU.

## Bluefield System on Chip Architecture

16 ARMv8 Cortex-A72 Cores
- Three-level cache hierarchy
- SkyMesh coherent interconnect
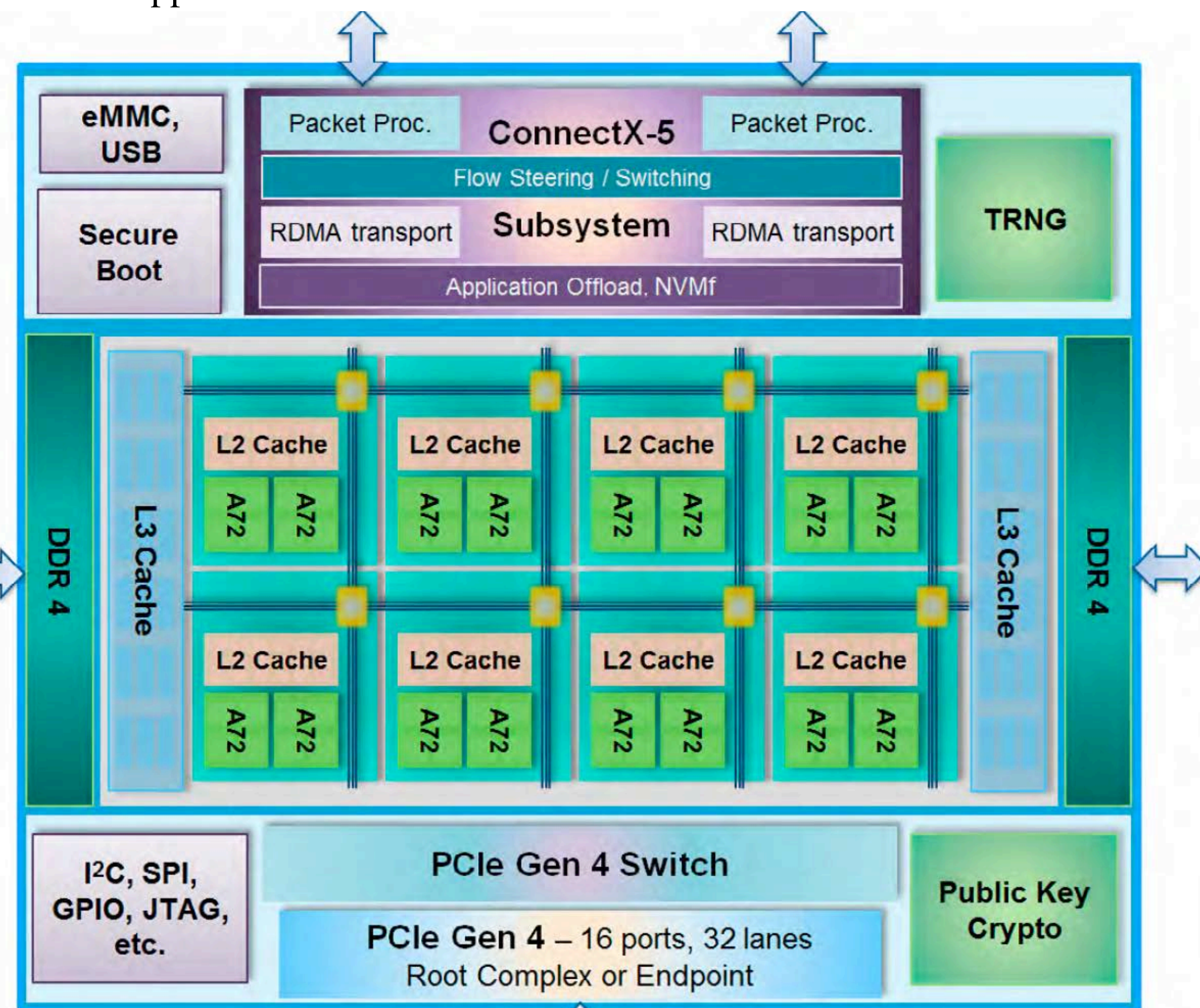- 128b ARM Neon SIMD execution unit per core

Connect X-5 Subsystem
- Dual Virtual Protocol Interconnect (VPI) ports
- Ethernet/Infiniband at 100Gbps per port
- RDMA & NVMe-oF support

Integrated PCIe Switch
- 32 bifurcated PCI 4.0 lanes
  - Configurable as 2x16/4x8/8x4/16x2
  - Speeds up to 200Gbps

Memory Controllers
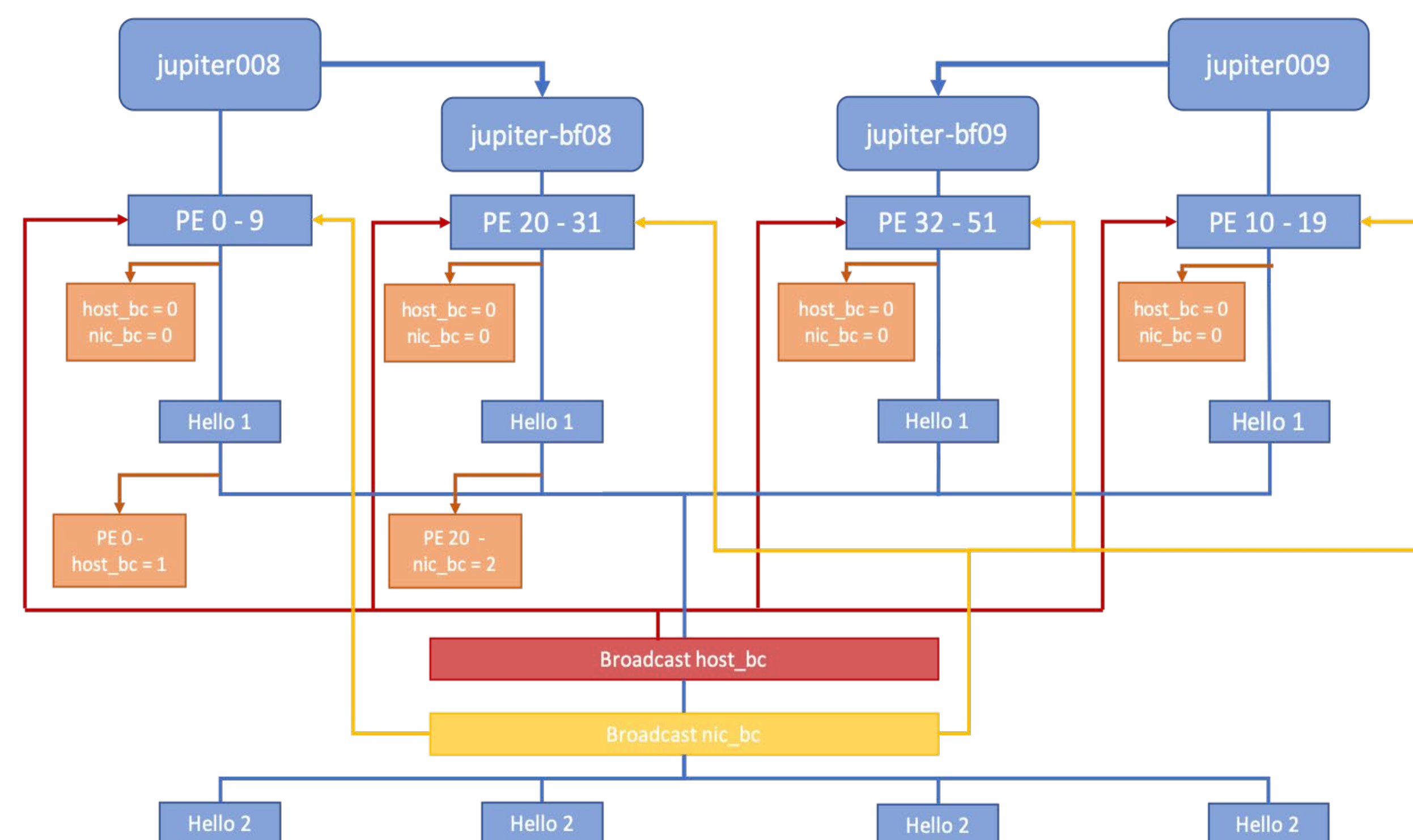- Supports two channels of 256 GB DDR4 DRAM at 1333MHz



Bluefield System on Chip Architecture.
http://www.mellanox.com/related-docs/npu-multicore-processors/PB_Bluefield_SoC.pdf

## Communication Experiment

In order for the SmartNICs to perform effectively as accelerators at any useful scale, they need to be able to communicate with other devices. Therefore, as a necessary prerequisite to any application performance optimization attempt, we first conducted an experiment to determine compatibility between the ARM cores onboard each SmartNIC and prominent distributed-address space programming paradigms.

### Mellanox Testbed – jupiter007 - jupiter010
- Intel Xeon E5-2680 v2 10-core processors
- 64 GB of memory
- Paired Bluefield SmartNIC with 16GB onboard memory
- CentOS 7
- OpenMPI 4.0.1 with Unified Communication X (UCX) 1.6



Experiment Program Flow

Output from our experiment shows that processor cores are properly utilized across multiple nodes and SmartNICs. Further, correct broadcast variable values demonstrate that proper inter-device communication and synchronization takes place despite buffered print statements. (Note that the output has been simplified for presentation purposes.)
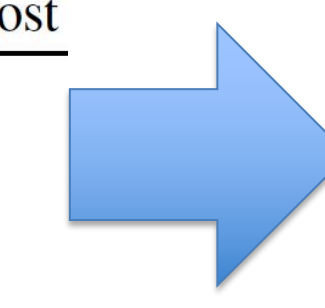


Output Screenshot

## Acceleration Opportunities

Asynchronous Execution of Task-Parallel Code Segments
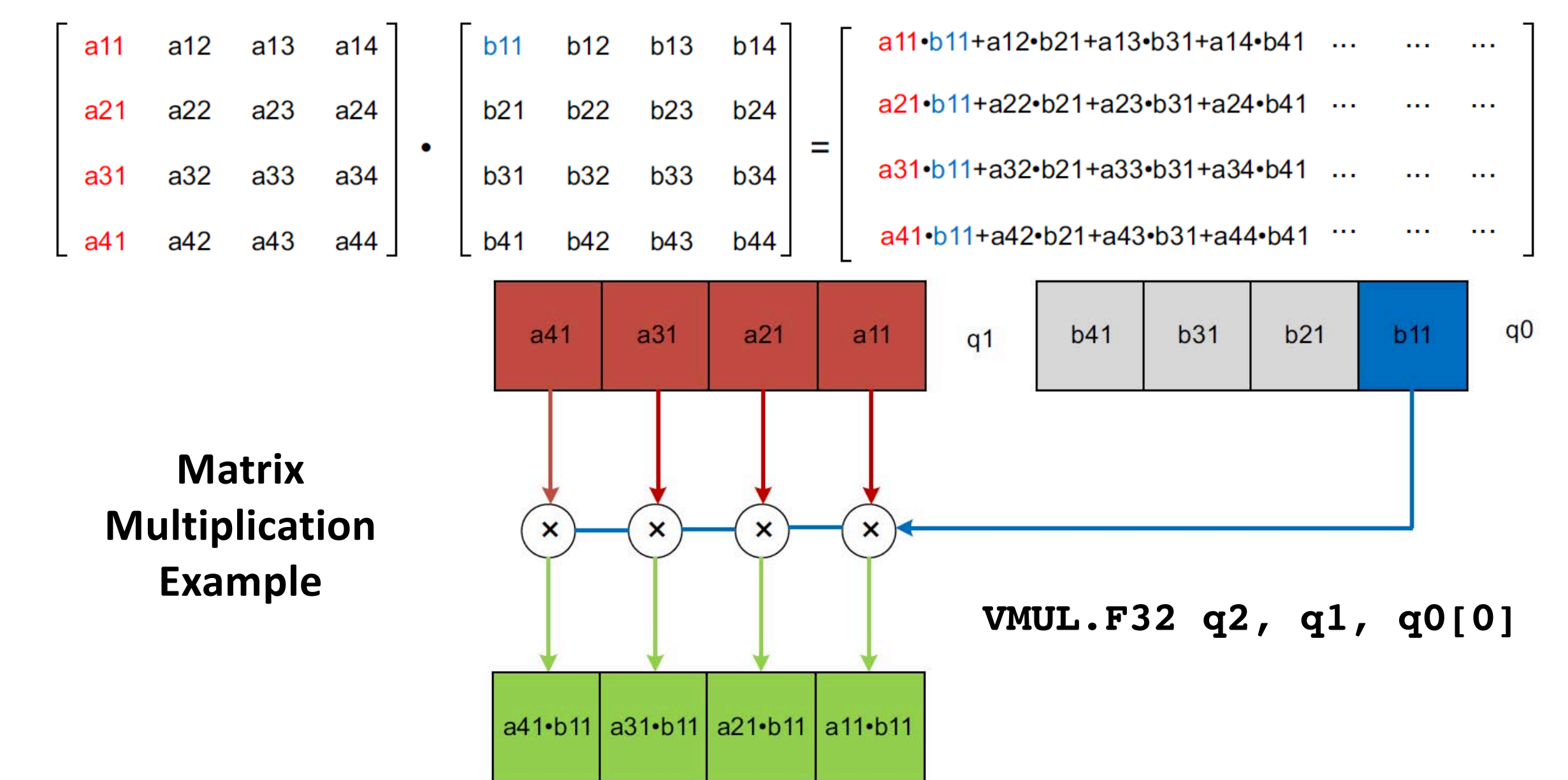- Map host and SmartNIC cores to orthogonal tasks

**Algorithm 1: Sample Execution on Host**
$$\text{for } i \leftarrow 0 \text{ to } N \text{ do}$$
$$\quad \text{Host\_Proc(Compute}(i))$$
$$\quad \text{Host\_Proc(Verify}(i))$$
$$\text{end}$$

**Algorithm 2: Sample Execution on Host + SmartNIC**
$$\text{for } i \leftarrow 0 \text{ to } N + 1 \text{ do}$$
$$\quad \text{if } Host\_Proc \text{ then}$$
$$\quad\quad \text{if } i \neq (N + 1) \text{ then}$$
$$\quad\quad\quad \text{Host\_Proc(Compute}(i))$$
$$\quad\quad \text{else if } SmartNIC\_Proc \text{ then}$$
$$\quad\quad\quad \text{if } i \neq 0 \text{ then}$$
$$\quad\quad\quad\quad \text{SmartNIC\_Proc(Verify}(i-1))$$
$$\text{end}$$

Vectorization of SIMD Operations using Neon Units
- Individual Neon unit per ARM core prevents resource contention
- Similar to acceleration using GPUs

**Matrix Multiplication Example**

```
VMUL.F32 q2, q1, q0[0]
```

http://infocenter.arm.com/help/topic/com.arm.doc.dui0489f/DUI0489F_arm_assembler_reference.pdf

Inter-Process Communication Calls
- Offload routines to SmartNIC cores
- Prevent blocking of host CPU cores
- Perform buffering and collective computation locally on SmartNICs

## Future Work

Our preliminary work for this project indicates that acceleration using Mellanox Bluefield SmartNICs is feasible. Further, our communication experiment demonstrates that the SmartNICs are fully compatible with the already widely adopted OpenSHMEM and MPI standards. Our future work will focus on determining whether or not these SmartNICs perform effectively as accelerators, and, if so, how best to optimize and deploy code for SmartNIC acceleration.

In particular, we plan to first optimize several provided Department of Defense benchmarks using the methods discussed above. We are also interested in investigating the use of active messages to pass fully dependence-free code segments to the SmartNICs for execution. Finally, when more familiar with writing code for SmartNIC acceleration, we plan to explore developing a library that abstracts SmarNIC acceleration away from the software developer while providing optimal performance.

## Acknowledgements