# Opportunities and Challenges for Accelerated Network Interfaces in HPC
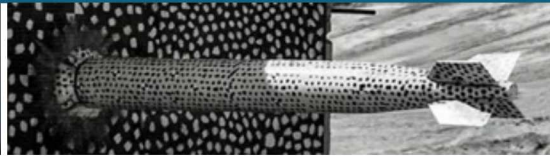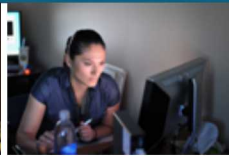
*Department 1423 - Scalable System Software*
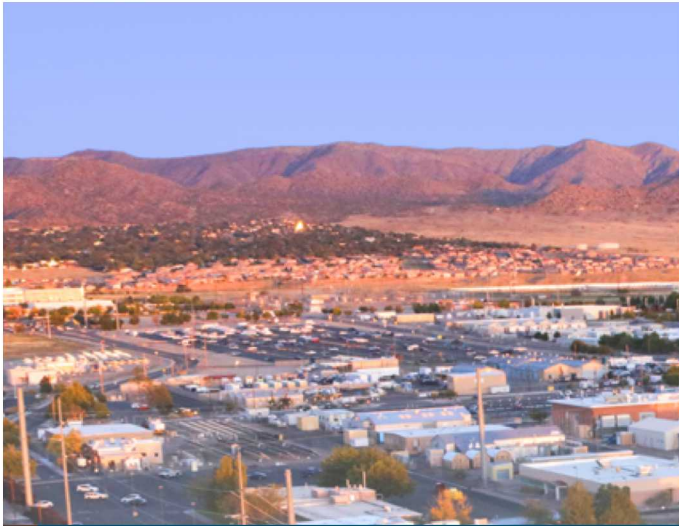
Ron Brightwell, R&D Manager

# Outline

- Past
  - ◦ Brief history of network offload in HPC

- Present
  - ◦ SmartNICs

- Future
  - ◦ Recent work on network offload technologies
  - ◦ Disaggregated system architecture

# A Brief History of
# Network Offload in HPC

# Intel Paragon Node (Circa 1993)



**Message Processor.** When an application decides to send a mesage, the node's i860 XP message processor handles message-protocol processing and frees the application processor to continue with numeric computing. Messaging software is executed from the message processor's internal cache, enabling overlapped communication and application processing to occur without incurring expensive context-switching delays. The message processor is also used to implement efficient global operations such as synchronization, broadcasting and global reduction calculations (e.g., global sum).

**Message Routing.** The actual transmission of messages is carried out by an independent routing system of custom-designed Mesh Router Controllers (MRCs), one for each node, arranged in a two-dimensional mesh. These fixed-function

*General-Purpose Node.* Each GP node dedicates one i860 XP processor to user applications and one to message processing. The GP node's expansion port allows the addition of an I/O or networking interface.

# Intel ASCI Red (TFLOPS) Compute Node



Figure 3: The ASCI Option Red supercomputer Kestrel Board. This board includes two compute nodes chained together through their NIC's. One of the NIC's connects to the MRC on the backplane through the ICF Link.

# Paragon/TFLOPS Network Interface Controller (NIC)

- Attached to the memory bus

- Cache coherent with the processor(s)

- Programmed by the operating system (OS)
  - Device driver was embedded in the OS
  - Driver consisted of programming DMA engines and memory-mapped registers

- Interrupt-driven
  - An interrupt would be generated for:
    - Arrival of an incoming message
    - Completion of an incoming message
    - Completion of an outgoing message

- Messages initiation via system call trap

- Source-routed, circuit-switched, wormhole routed network
  - Message header contained route to destination
  - Message body was one contiguous block

# Basic Assumptions About Networking for Massively Parallel Systems

- A single low-level network API is needed
  - Compute node may not have a TCP/IP stack
  - System is space-shared
    - Compute node application should own all network resources

- Applications will use multiple protocols simultaneously
  - Can't focus on just MPI
  - Runtime system, system call forwarding, I/O protocols too

- Need to support communication between unrelated processes
  - Client/server communication between application processes and system services

- Need to support general-purpose interconnect capabilities
  - Can't assume special collective network hardware

- Interconnect hardware limitations can't be fixed in software

# Key Network Capabilities

- Independent progress
  - Data should move without requiring polling from user-level library
  - Adhere to the strong progress rule interpretation of MPI

- Overlap
  - Decouple the host processors from the network as much as possible
  - Enable overlap of computation and communication as well as communication and communication

- Scalable use of memory resources
  - Buffer space for MPI unexpected messages
  - Memory use should be independent of the number of peers

- High performance
  - Maximize bandwidth by avoiding memory-to-memory copies
  - Minimize latency by avoiding OS interaction

# Programmable User-Level Networks Enabled API Exploration

- Myrinet (~1994)
  - First commercially available Gb/s standalone network
  - Based on technology developed for Intel MPP networks
  - Initially available for Sun SPARC SBus, later for PCI-based PCs
  - Custom embedded MIPS-based programmable processor (LANai)
  - Myrinet Control Program (MCP) software development environment
  - Destination routed, maximum message size (packets)
  - Numerous APIs and MCPs: AM, FM, GM, PM, MX

- Quadrics QSNet (Elan + Elite) (~2001)
  - Outgrowth of technology developed for Meiko MPP networks
  - Offered several different APIs for user-level networking
  - Provided a development environment for running user-level functions on NIC

Processor    Memory

Memory Bus

Bridge

PCI Bus

NIC

Network

# Portals Interconnect Programming Interface
## A Vehicle for Hardware/Software Co-Design

- Initially developed primarily by Sandia and the University of New Mexico

- Deployed on several production massively parallel processing (MPP) and cluste
  - 1993: 1800-node Intel Paragon (SUNMOS)
  - 1997: 10,000-node Intel ASCI Red (Puma/Cougar)
  - 1999: 1800-node Cplant cluster (Linux)
  - 2005: 10,000-node Cray Sandia Red Storm (Catamount)
  - 2009: 18,688-node Cray XT5 – ORNL Jaguar (Linux)
  - 2017: Bull BXI interconnect

- Focused on providing
  - Lightweight "connectionless" model for massively parallel systems
  - Low latency, high bandwidth
  - Independent progress
  - Overlap of computation and communication
  - Scalable buffering semantics
  - Protocol building blocks to support higher-level application protocols and libraries and system services

http://www.cs.sandia.gov/Portals/

# Fixing Semantic Mismatch Between Layers

**Majority of interconnect software R&D is spent on dealing with the semantic mismatch between what the upper-layer protocols need and what the low-level network software and the underlying hardware provide**

**RDMA (e.g. InfiniBand Verbs)**

One-sided
- Allows process to read/write remote memory implicitly

Zero-copy data transfer
- No need for intermediate buffering in host memory

Low CPU overhead
- Decouples host processor from network

Fixed memory resources
- No unexpected Messages

Supports unstructured, non-blocking data transfer
- Completion is a local event

**MPI Point-to-Point**

Two-sided
- Short messages are copied
- Long messages need rendezvous

CPU involved in every message
- Message matching

Unexpected messages
- Need flow control

Completion may be non-local
- Need control messages

# Motivation for Low-Level Transport APIs

- Targeting a single programming model target is too limiting (top down approach)
  - MPI - MPICH, OpenMPI
  - PGAS - GasNet, OpenSHMEM
  - I/O – LNet, Mercury, Nessie

- Desire to reduce development costs
  - Provide one network abstraction for all ULPs
  - Large porting effort is a strong indication of the semantic mismatch

- "Vendor differentiation
  - Which really defeats the portability goal

- Vehicle for hardware/software co-design R&D

# Red Storm – Prototype for Cray XT Series

- Architected by Sandia, engineered jointly with Cray in 2003
  - Sandia contributed to the design of the SeaStar network interface and router

Sandia also developed
  - Lightweight kernel compute node OS
  - Scalable parallel job launching system
  - Portals high-performance interconnect programming interface
  - SeaStar firmware

- 140+ systems to 80 different customers worldwide
  - Including ORNL, NERSC, and LANL

- Following Red Storm, Cray's market share rose from 6% in 2002 to 21% in 2007*

- Revenue of $1B +

- Basis of Cray's business today

- "Virtually everything we do at Cray - each of our three business units - comes from Red Storm. It spawned a company around it, a historic company struggling as to where we would go next. Literally, this program saved Cray." – Pete Ungaro, CEO, Cray Inc.

- https://www.datacenterdynamics.com/analysis/after-the-storm-the-supercomputer-that-saved-cray/



SeaStar was a PowerPC-Based System-on-a-Chip (SOC)



Cumulative Cabinet Deliveries

*Source: IDC #209251 *Technical Computing Systems: Competitive Analysis*, November 2007

# Significant Vendor Impact of Sandia's Portals Networking Technology

All of these production vendor-supported systems used Portals as the network hardware programming interface.
Portals enabled the first TeraFLOPS platform (ASCI Red) and the first non-accelerated PetaFLOPS platform (Jaguar).


Intel Paragon
**Portals 0**


Intel ASCI Red
**Portals 2**


Cray Red Storm
**Portals 3**


Cray XT3, XT4, XT5
**Portals 3**


Atos Tera1000
**Portals 4**

Unlike other low-level network programming interfaces, Portals is intended to enable co-design rather than serve as a portability layer.
The influence and impact of Portals can be seen in vendor co-design activities, other low-level network programming interfaces, and emerging network hardware.

**AMD FastForward Project based on Portals 4**



**OFI Libfrabric API based on Portals 4**



**Atos Bull eXascale Interconnect (BXI) based on Portals 4**



**Lustre File System network based on Portals 4**



**Mellanox ConnectX-5 MPI tag matching in hardware**



**Cray Slingshot Supports Portals 4 header**

# SmartNICs

# Onload Versus Offload Argument (~2005)

- Why design a custom NIC for offload?

- It is too expensive
  - ◦ HPC market is too small to justify a custom ASIC for an HPC NIC

- Just dedicate a core
  - ◦ A 3 GHz Xeon will outperform a 500 MHz embedded processor on network protocol processing
  - ◦ Overhead cost will go down as core count increases
  - ◦ Cores won't be getting slower, right?

# Cray Core Specialization

- Cray's chose a complete re-design for Gemini NIC
  - Gemini was delivered six years after the first SeaStar chip

- Dedicate "OS" cores to handle MPI progress
  - MPI progress threads run on a dedicated set of cores



**S3D Time Step Summary**

| # Application Threads | Progression disabled | Progression enabled |
|---|---|---|
| 14 | 4.77 | 3.93 |
| 15 | 4.68 | 4.05 |
| 16 | 4.59 | 4.06 |

**MILC Run Time Summary(secs)**

| # Run Type | 4096 ranks | 8192 ranks |
|---|---|---|
| No progression | 2165 | 1168 |
| Progression (phase 1) | 2121 | 1072 |
| Progression (phase 2) | 3782 | 2138 |
| Progression (phase 1) no reserved cores | 3560 | 2210 |
| Progression (phase 1) reserve core but no corespec | 2930 | 2070 |

# HPC NIC Timeline

**2001**
- Mellanox SDR

**2003**
- Quadrics Elan3

**2005**
- Cray SeaStar
- Quadrics Elan4

**2005**
- Mellanox DDR

**2008**
- Mellanox QDR

**2009**
- Quadrics goes out of business

**2010**
- Cray Gemini

**2011**
- Mellanox FDR

**2012**
- Cray Aries
- Intel acquires Cray's network assets

**2013**
- Myrinet goes out of business

# 2014 - 2017

- NDAs prohibit me from saying much about this time period

- Two offload capable HPC NIC in design/development
  - Bull BXI

- Offload-capable and programmable NICs aren't fashionable

# 2019 - SmartNICs are Everywhere

# What Happened?

- Core frequency stalled and signaling rate kept going
  - A single core doing network protocol processing can't keep up

- Cloud computing with virtual machines
  - SDN, Open vSwitch, DPDK

- Main barrier to developing offload capable NICs was not technical
  - Economics drive innovation



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

**Azure Accelerated Networking: SmartNICs in the Public Cloud**

Daniel Firestone    Andrew Putnam    Sambhrama Mun
Mike Andrewartha    Hari Angepat    Vivek Bhanu
Harish Kumar Chandrappa    Somesh Chaturmohta    Matt
Fengfen Liu    Kalin Ovtcharov    Jitu Padhye    Gauthan
Mark Shaw    Gabriel Silva    Madhan Sivakumar    Nisheeth S
Deepak Bansal    Doug Burger    Kushagra Vaid
Microsoft

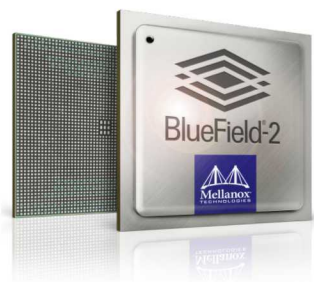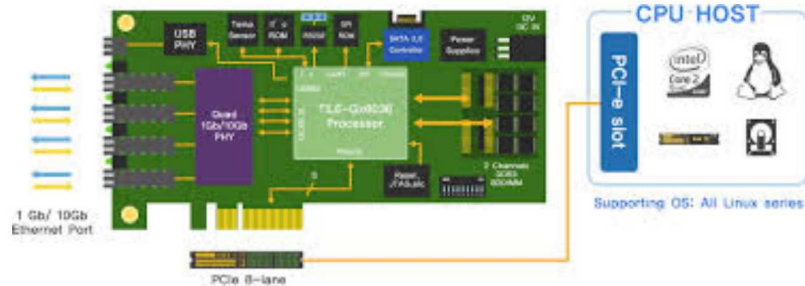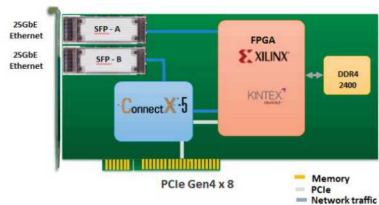Even considering that some fraction of cores are unsold at any time and that clouds typically offer customers a discount for committed capacity purchases, using even one physical core for host networking is quite expensive compared to dedicated hardware. Our business fundamentally relies on selling as many cores per host as possible to customer VMs, and so we will go to great lengths to minimize host overheads. Thus, running a high-speed SDN datapath using host CPU cores should be avoided.

- **Better perf means better reliability**. One of the biggest benefits of AccelNet for VMs is that the network datapath no longer shares cores or resources with the rest of the host, and is not subject to transient issues — we've seen much more reliable performance and lower variance as a result.

# What is a SmartNIC?

- Designed for virtualization environments and SDN

- Need significant more processing for network protocol offload

- Programmable from the network protocol perspective
  ◦ Not really intended to offload upper-level protocol functionality

- Move parallelization into the NIC because they can't get it in the host
  ◦ Flows can be parallelized more easily below the VM

- Tradeoffs between ASIC, FPGA, SOC
  - FPGAs provide a more flexible platform for rapidly evolving TLAs
  - FPGAs provide more reliable performance and are not subject to variance



NIC Implementation Comparison

Implementation does **Not** make a NIC a _SmartNIC_

**ASIC Based**
- Excellent price-performance
- Vendor development cost high
- Programmable and extensible
  - Easy to program but flexibility is limited to pre-defined capabilities

**FPGA Based**
- Good performance but expensive
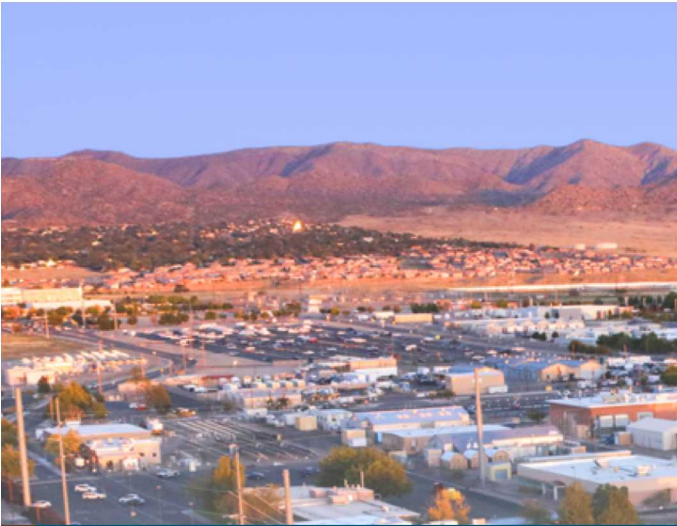- Very difficult to program
- Workload specific optimization

**SOC Based**
System on Chip - NIC + CPU
- Good price-performance
- C Programmable Processors
- Highest Flexibility
- Easiest programmability

| Capability | Workloads Accelerated | Basic NIC | Cloud NIC | SmartNIC |
|---|---|---|---|---|
| **Basic Virtualization and Data Movement** | | | | |
| TCP/IP Acceleration | Enterprise workloads | ✔ | ✔ | ✔ |
| SRIOV (NIC-level virtualization) | Enterprise workloads | ✔ | ✔ | ✔ |
| VXLAN, NVGRE (Network-level virtualization) | Multi-Tenant workloads | ✔ | ✔ | ✔ |
| **Data Transport Acceleration** | | | | |
| RDMA and RoCE | Big Data, Storage, AI/ML, Virtual Machines | | ✔ | ✔ |
| ACLs and QoS support | Web serving, Content Distribution (CDN) | | ✔ | ✔ |
| OVS hardware acceleration | Efficient, Scalable Virtualized Apps | | ✔ | ✔ |
| Storage functions acceleration (NVMe-oF, RAID, T10) | High-performance, low-latency Storage | | ✔ | ✔ |
| Flow match/action engine | Software Defined Networking (SDN) | | ✔ | ✔ |
| Flow monitoring/reporting | Visibility, Network Packet Broker, IBN | | ✔ | ✔ |
| **Smart Networking, Security & Virtualization** | | | | |
| Stateful n-tuple / ACL filtering | Load Balancing, IPS/IDS, Unified Threat Management | | | ✔ |
| Isolated OVS control plane, Fault Isolation & Response | Bare Metal Cloud, High-Availability Datacenter | | | ✔ |
| Analytics Engine | DPI, Network Monitoring and Diagnostics | | | ✔ |
| On-board Security VNFs | Firewall, Anti-DDoS, Anti-Malware, Host Introspection | | | ✔ |
| Datapath encryption/decryption | Secure data-at-rest or data-in-flight | | | ✔ |
| Public Key crypto, hardware RNG | Secure Authentication & Key Exchange | | | ✔ |

# Recent Work on Network Offload Technologies

CCR
**Center for Computing Research**

# Active Messages (AM)

- T. von Eicken, et al.: "Active Messages: A Mechanism for Integrated Communication and Computation" (1992)

- Lots of different flavors
  - Pure active messages
    - Origin sends message to target containing code and data
    - Target invokes code on that data
  - Generalized active messages
    - Origin sends message to target containing function id and data
    - Function id maps to existing code in target's address space

- Similar to remote procedure invocation without returning a result to origin

- Semantically equivalent to blocking a thread on an incoming message and invoking a handler when the message arrives
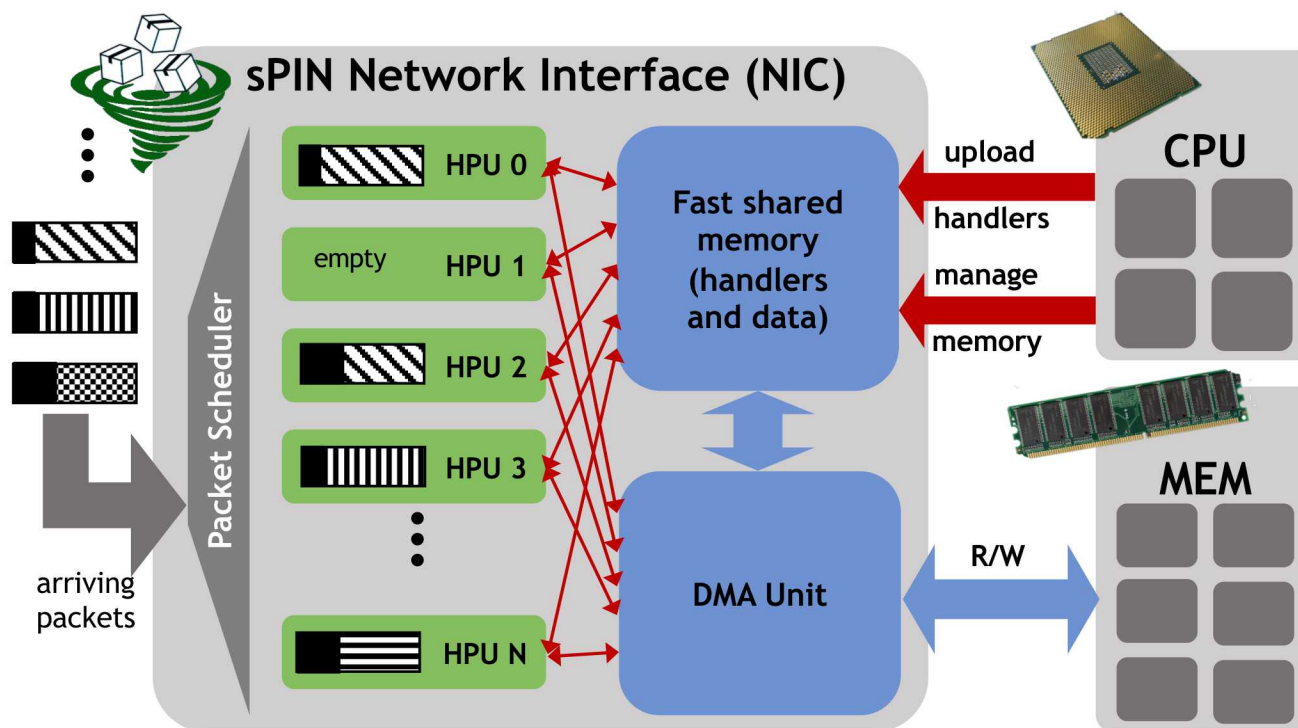
# Issues with Active Messages

- Data delivery
  - Who determines where the data goes – origin or target?
  - How much data can be delivered?

- Handlers
  - When are resources allocated – a priori or on arrival?
  - What can be called?
  - Where do they run (context)?
  - When do they run relative to message delivery?
  - How long do they run?
  - Why?
    - One-sided messages decouple processor from network
    - Active messages tightly couple processor and network
      - Active messages aren't one-sided
      - Memory is the endpoint, not the cores
    - Lightweight mechanism for sleeping/waking thread on memory update
      - Why go through the network API for this?

- Scheduling lots of unexpected thread invocations leads to flow control issues

# Is There a Better Way to Get AM Semantics?

- Cores are slower, more energy-efficient
  - Modern cores require 15-20 ns to access L3 cache
    - Haswell – 34 cycles
    - Skylake - 44 cycles

- Terabit per second networks are coming
  - 400 Gib/s can deliver a 64-byte message every 1.2 ns

- Need to remove processor from network processing path (offload)

- RDMA only supports data transfer between virtual memory spaces
  - Data is placed blindly into memory
  - Need varying levels of steering the data at the target

streaming Processing In the Network (sPIN)



sPIN Network Interface (NIC)

HPU 0
empty HPU 1
HPU 2
HPU 3
HPU N

Packet Scheduler

arriving packets

Fast shared memory (handlers and data)

DMA Unit

upload
handlers
manage
memory

CPU

R/W

MEM

Hoefler, Di Girolamo, Taranov, Grant, Brightwell. "sPIN: High-Performance Streaming Processing in the Network," in Proceedings of SC"17, November 2017.

# sPIN is not Active Messages

- Tightly integrated NIC packet processing

- AMs are invoked on full messages
  - sPIN works on packets
  - Allows for pipelining packet processing

- AM uses host memory for buffering messages
  - sPIN stores packets in fast buffer memory on the NIC
  - Accesses to host memory are allowed but should be minimized

- AM messages are atomic
  - sPIN packets can be processed atomically

# sPIN Approach

- Handlers are executed on NIC Handler Processing Units (HPUs)

- Simple runtime manages HPUs

- Each handler owns shared memory that is persistent for the lifetime of a message
  - Handlers can use this memory to keep state and communicate

- NIC identifies all packets belonging to the same message

- Three handler types
  - Header handler – first packet in a message
  - Payload handler – all subsequent packets
  - Completion handler – after all payload handlers complete

- HPU memory is managed by the host OS

- Host compiles and offloads handler code to the HPU

- Handler code is only a few hundred instructions

# sPIN Approach (cont'd)

- Handlers are written in standard C/C++ code

- No system calls or complex libraries

- Handlers are compiled to the specific Network ISA

- Handler resources are accounted for on a per-application basis
  - Handlers that run too long may stall NIC or drop packets

- Programmers need to ensure handlers run at line rate

- Handlers can start executing within a cycle after packet arrival
  - Assuming an HPU is available

- Handlers execute in a sandbox relative to host memory
  - They can only access application's virtual address space
  - Access to host memory is via DMA

# INCA: In-Network Compute Assistance

- Traditional offloaded network processing:
  - Often task-specific, even when general-purpose hardware is used (CPUs, FPGAs)
  - May sometimes allow users to offload general-purpose kernels (e.g., Myrinet, SPiN)
  - Nonetheless deadline-based:
    - Compute resources must be available for incoming data
      - Limits on number of instructions that can be executed per message or packet (see next slide)
      - Kernels written for one system may not honor deadlines imposed by future systems
      - Desideratum of harvesting compute resources when the network is idle not addressed

- INCA: Enable deadline-free, general-purpose compute offloading

Schonbein, Grant, Dosnjh, Arnold. "INCA: In-Network Compute Assistance," in Proceedings of SC'19, November 2019.

# INCA: In-Network Compute Assistance

- Instruction limits
  - For instance: for a data rate of 200Gb/s, 64B packets, and 32 2.5GHz cores with 1 IPC, a stream-based approach is limited to packet-processing kernels of less than 500 instructions.
  - This limit can exclude:
    - Dot product, matrix transpositions, and Hadamard products involving 1024 elements or more.
    - Significant matrix multiplications (exceeding 16 64bit elements).
    - Simple 3x3 convolutions
    - Linear interpolations exceeding 16 64bit points.

# INCA: In-Network Compute Assistance

# INCA: In-Network Compute Assistance

1. Triggered operation generates message containing 1st argument.

2. Unique matching element specifies buffer containing 2nd argument and atomic.

3. Atomic unit performs specified operation and stores result.

4. Counter incremented

```
→•→ ┌─────────────┐    ┌──────────┐    ┌─────────────┐    ⎛  ⎞
     │  Triggered  │ →  │ Matching │ →  │   Atomic    │ →  │++│
     │  Operation  │    │          │    │  Operation  │    ⎝  ⎠
     └─────────────┘    └──────────┘    └─────────────┘
```

# INCA: In-Network Compute Assistance



Turing Complete*

Triggered Operation Machine

| Triggered Operation | | Matching | | Atomic Operation | | ++ |

*assuming some slight modifications to existing hardware…

# INCA: In-Network Compute Assistance

- Kernel execution is deadline-free:
  - Execution can be suspended (e.g., by incoming network traffic) between triggered messages being generated and matching.

- Kernels are forward compatible:
  - No limits on number of instructions, so faster networks do not imply fewer instructions can be executed.
  - … and faster networks mean INCA kernels will execute faster when they are scheduled.

- Kernel execution can continue when the network is otherwise idle.

# INCA: In-Network Compute Assistance

• INCA: abstract triggered operation machine programs into something more recognizable, e.g.:

---
**Algorithm 2** INCA Dot Product
---
1: PUTL i, 0, i16
2: PUTL $r_0$, i, i16
3: LT $r_0$, $r_0$, b, i16
4: BLEZ $r_0$, 10
5: PUTL $r_1$, A[i], f
6: MUL $r_1$, $r_1$, B[i], f
7: ADD c, c, $r_1$, f
8: ADD i, i, 1, i16
9: JMP 2
10: END
---

# Disaggregated Systems Architecture



*Stolen from John Shalf

# HPE's OpenFAM – Disaggregated Persistent Memory

# Summary

- Innovation and technology are driven more by economics than science

- SmartNICs are fundamental enablers for Cloud Computing
  - Not clear that they will enable HPC-style offload
  - Disaggregated system architecture may eliminate need for virtualization

Questions?