# *CAC 2007 Panel Position:*

# *Accelerators In Cluster Communication*

Greg Pfister
Distinguished Engineer
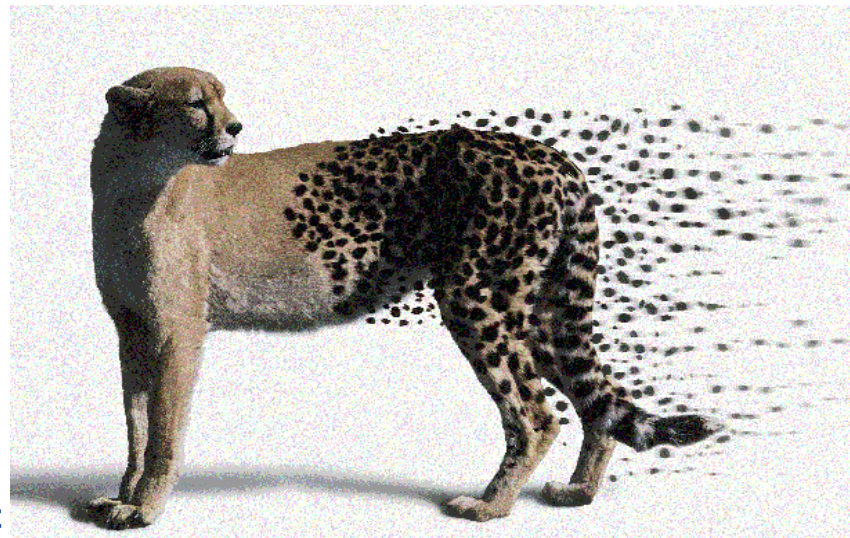IBM Systems and Technology Group
Austin, TX

Any relationship between the opinions expressed here and any official position of the IBM Corporation…

$f$ or of anyone else in the IBM Corporation

...is purely coincidental.

And all copyrighted words are the property of their respective owners.

**Random gratuitous clipart**

# The Questions, and Some of My Answers

| | |
|---|---|
| Will full or partial offload be the norm for Ethernet? | Partial |
| Will any collective functions *typically* be offloaded to the NIC in the near future? | No |
| What about collective support in switches **{*typically* }**? | Heck no |
| Should NICs assist in "read-modify-write" capability to support remote updates? | No, but. |

## *Required*

| | |
|---|---|
| What communication functions should **always** be offloaded to accelerators? | later |
| Which should **never** be offloaded to accelerators? | later |
| Are there cases for which it makes sense to offload to another core rather than to an accelerator? (For SMT cores, I suppose one might also consider offloading to other threads of the same core?) | later, but of course. |

# Sample of Industry Acceleration Activity

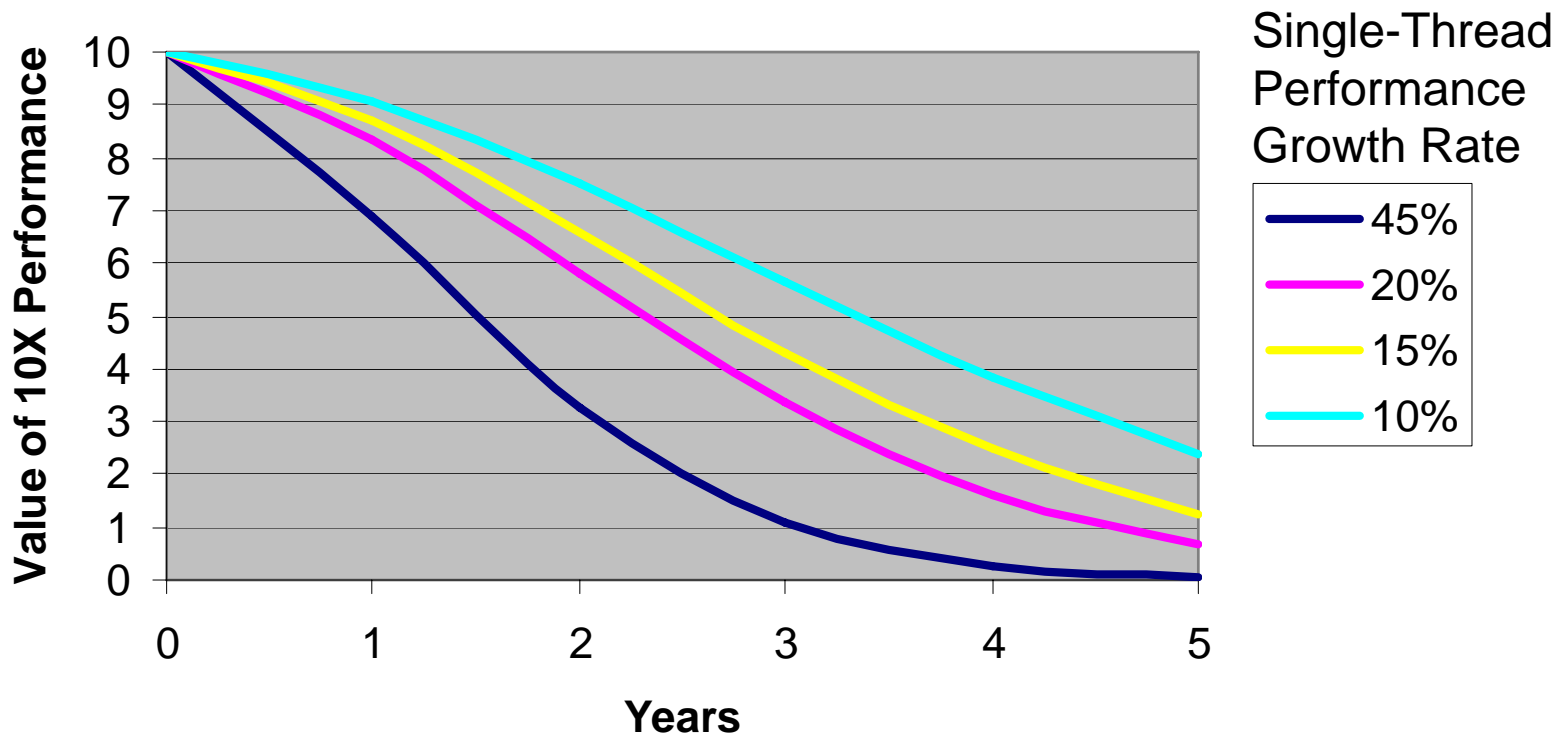| Category | Subcategory | Companies |
|---|---|---|
| XML, Java | XML Processing | Actional, Cisco, Conformative, DataPower, Reactivity, Sarvega, Tarari |
| | Java Virtual Machine | Azul |
| Networks: Storage, Cluster, LAN | TCP/IP + Ethernet | Adaptec, Allied Telesyn, Amasso, Brocade, Chelsio, Cisco(Topspin), Citrix, Crescendo Networks, Enigma Semiconductor, Infrant, NetEffect, NextIO, Nortel, Precision I/O, Silverback, Sensory Networks, Tehuti, Toplayer , Voltaire |
| | InfiniBand | Cisco (Topspin), Mellanox, SilverStorm, Voltaire |
| Application Acceleration | | Cisco(AON), Citrix, F5 Networks, Juniper Networks(Redline), Radware |
| Security, Privacy, Rights Management | Trusted Environment | STMicroelectronics, Intel, Infineon |
| | Cryptographic Functions | Actional, Allied Telesyn, Broadcom, Check Point, 3Com, Cipheroptics, Cisco, DataPower, Enterasys, Forum Systems, HP, Intel, Kasten Chase, Lucent, nCipher, Nokia, Nortel, F5 Networks, Radware, Reactivity, Sonicwall, Sun, Tarai, Vormetric |
| Real-time Analytics | Data Warehouse Based | Cogent Systems, DATAllegro, Netezza |
| | Data Stream Based | Tarari |
| Collaboration & Information Mgmt | Audio, video, data, IM Fusion | ClearSpeed, Tarari |
| | Media Distribution | Tarari, SeaChange |
| | Enterprise Search | Netapp, Google, Search Cacher, Thunderstone |
| High Performance Computing | Floating Point & Integer Comp | ClearSpeed, Cray, Mitrionics, SRC, Tarari, TimeLogic |
| | Messaging | Cisco (Topspin), Mellanox, Myricom, Quadrics, QLogic, Voltaire, … |
| | Reconfigurable Systems & SW | Cray, SRC, Nallatech, Tarari, Cadence, Mitrion, Koan, Synopsys, Celoxica, Impulse, Starbridge, SGI |
| Intelligent Storage Network | Storage Virtualization | Acopia, Brocade, EMC, HP, Hitachi, Index Engines, IBM, NeoPath, Sun, Troika, |
| | Storage Services | |
| Accelerator Technology | System-on-Chip | Arteris S.A., Bay Microsystems, Broadcom, Cavium, Freescale Semiconductor, Infineon, LSI Logic, Rapport (Kilocore), Raza Microelectronics, STMicroeletronics, Teja, Tensilica |
| | ASIC | Advanced Architectures, Britestream, Cavium, Critical Blue, Elixent, Forte, Freescale Semiconductor (Seaway Networks), IP Fabrics, LSI Logic, Mellanox, nCipher, Propulsion Networks, STMicroeletronics, Xelerated |
| | FPGAs | Xilinx, Altera, Lattice, Acte |

Background table (partially visible):

| XML, Java | XML Processing | Actional, Cisco, Conformative, DataPower, Reactivity, Sarvega, Tarari |
| | Java Virtual Machine | Azul |
| Networks: Storage, Cluster, LAN | TCP/IP + Ethernet | ...asso, Brocade, Chelsio, Cisco(Topspin), Citrix, Crescendo Networks, Enigma Semiconductor, Infrant, NetEffect, NextIO, Nortel, Precision IO, Silverpack, Sensory Networks, Tehuti, Toplayer , Voltaire |
| | InfiniBand | ...SilverStorm, Voltaire |
| Application Acceleration | | ...rix, F5 Networks, Juniper Networks(Redline), Radware |
| | Trusted Environment | STMicroelectronics, ... |
| Security, Rights Management | Cryptographic Functions, etc. | Actional, AlliedTelesyn, Broadcom, CheckPoint, Ciperoptics, Cisco, DataPower, Enterasys, Forum Systems, HP, Intel, Kasten Chase, Lucent, nCipher, Nokia, Nortel, F5 Networks, Radware, Reactivity, Sonicwall, Sun, ... |
| Real-time Analytics | Data Warehouse Based | Google(SYSense, DATAllegro) Netezza ... |
| Collaboration & Information Mgmt | Audio, Video, data; IM, etc. | ...ClearSpeed, Tarari |
| | Enterprise Search | Netapp, Google, Search Cacher, Thunderstone |
| High Performance Computing | Messaging | Cisco (Topspin), Mellanox, Myricom, Quadrics, QLogic, Voltaire, ... |
| | Reconfigurable Systems | ...ence, Mitrion, Koan, Synopsys, Celoxica, Impulse, Starbridge, SGI |
| Intelligent Network | Storage Services | ...h, Sun, Troika, |
| | System-on-Chip | ...Broadcom, Cavium, Freescale Semiconductor, Infineon, LSI Logic, Rapport (Kilocore), Raza Microelectronics, STMicroelectronics, ... |
| Accelerator Technology | ASIC | ...ctures, Britestream, Cavium, Critical Blue, Elixent, Forte, Freescale Semiconductor (Seaway Networks), IP Fabrics, LSI Logic, Mellanox, nCipher, Propulsion Networks, STMicroelectronics, Xelerated |
| | FPGAs | Xilinx, Altera, Lattice, Acte |

Overlay box:

## This is a zoo

- There is no single accelerator architecture.
  - ƒ Some are totally unique hardware.
  - ƒ Many are general-purpose systems with a special algorithm, security chip, NIC, etc.
- There is no single accelerator programming model.
  - ƒ Subroutine, coroutine, SOA, SQL, passthrough, shared memory, message-based, system access, user-mode access, …
- There is no single accelerator attachment technique.
  - ƒ Ethernet, IO bus adapter, part of chipset, etc.

**The strategy of acceleration is:** *specialize to the application.*
  - ƒ To optimize performance, power, area, cost
- The tactics – designs – are completely variable.

# Accelerator Longevity (Controversial!)



**Single-Thread Performance Growth Rate**

- —— 45%
- —— 20%
- —— 15%
- —— 10%

(Chart: Value of 10X Performance vs. Years, 0 to 5)

- **Frequency growth slowdown ⇨ enhanced business case:**
  - ƒ Past (45%): After 3 years – useless.
  - ƒ Now (<20%): as much 5 years useful lifetime

- **Controversy?** *Aggregate* chip performance still at 45%
  - ƒ There will be lots of cores and threads.
  - ƒ Must use them all (nontrivial), but accelerators often parallel too.

# What is an Accelerator? (vs. an Appliance)

- **NOT a general-purpose system.**

  - *ƒ* The strategy is specialization
  - *ƒ* "GP accelerator" is an oxymoron

- **Accelerator:** a device optimized to enhance the performance or function of a computing system.

  - *ƒ* *Does not function on its own.*
    - n Requires invocation from host programs.
    - n Intention & design optimization, not physics
    - n May contain GP system parts (like CPU) and be substantially software or firmware.
  - *ƒ* May contain other accelerators
    - n E.g., crypto in protocol offload in XML processing.

- **Appliance**: a device that performs a complete customer-visible function with substantially reduced management & programming requirements.
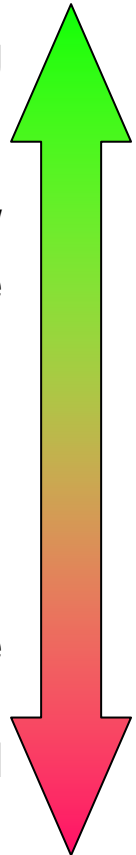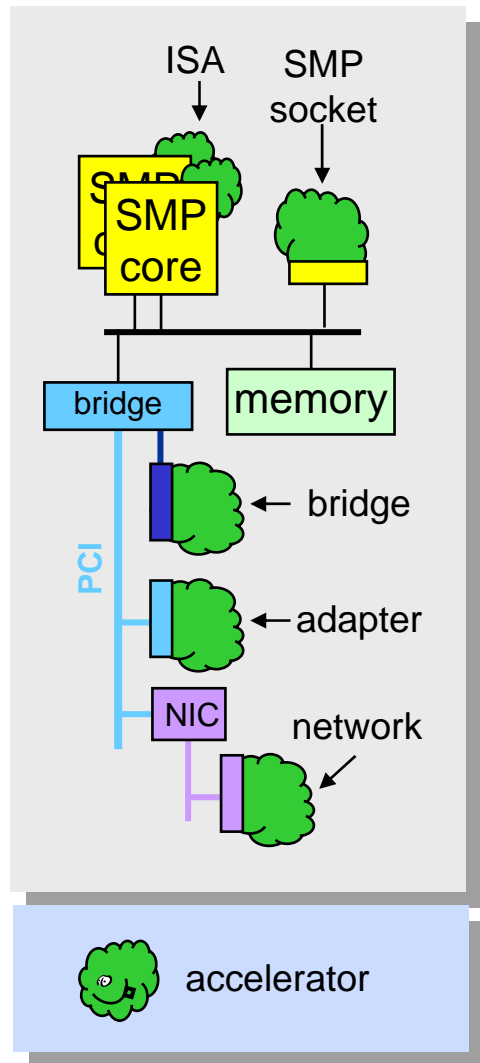
  - *ƒ* Functions on its own.

**accelerator**

**accelerator**

GP CPU?

**accelerator**

**System**

# Must **Attach** to a Computing System

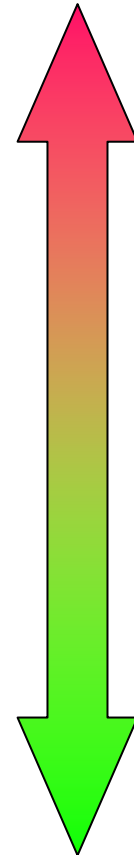- Accelerators have been attached to every conceivable orifice of a computer

Tighter coupling
Lower overhead
Faster synchro.
Integrated SW
Broader use

Narrower use
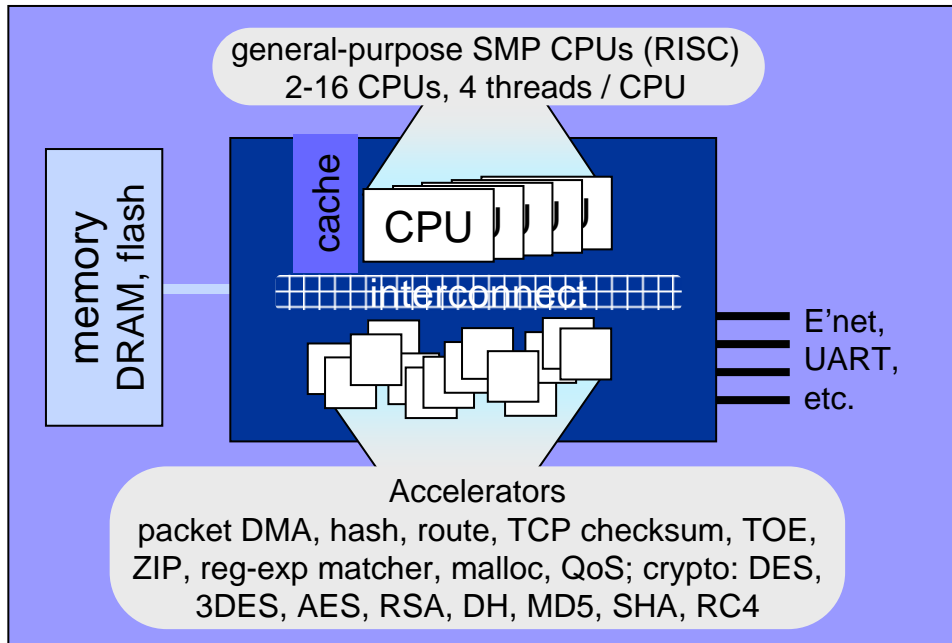Slower synchro.
Higher overhead
Loose coupling

ISA
SMP socket

SMP core

bridge
memory

PCI

← bridge

← adapter

NIC
network

accelerator

High dev expense
Arch. dependent
Longer lead times
High SW expense
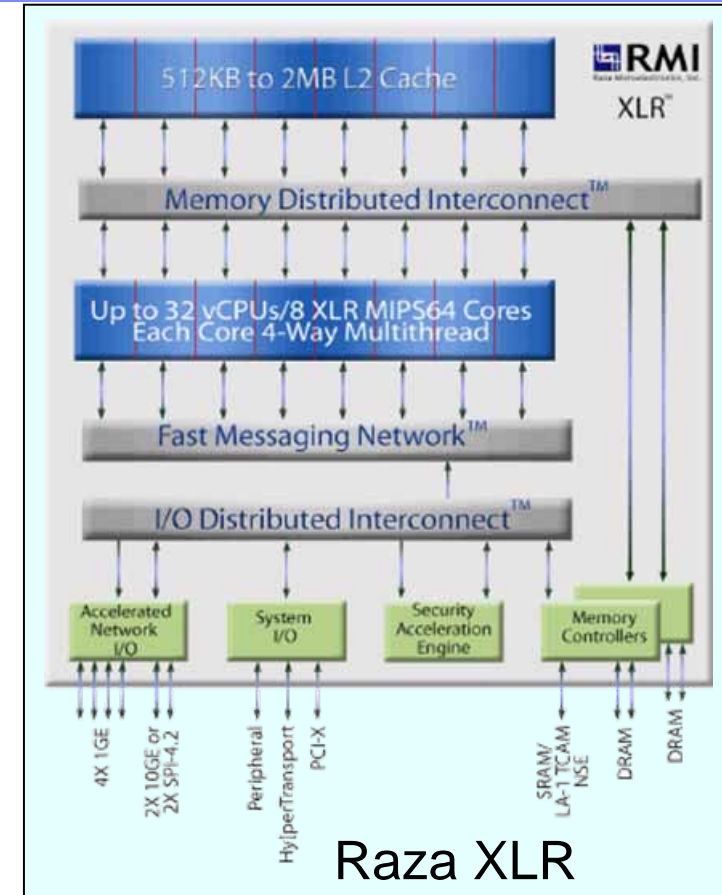3rd party⇒lower RAS
Few entrants
(AMD: 4; Intel: 2)

Many entrants
Lower SW exp
Fast time-to-mkt
Arch. independent
Low dev expense

(memory attach, too – proposals)

# Network Sys/chip: Arms Dealers to Appliances/Accelerators

general-purpose SMP CPUs (RISC)
2-16 CPUs, 4 threads / CPU

memory
DRAM, flash

cache

CPU

interconnect

Accelerators

E'net,
UART,
etc.

Accelerators
packet DMA, hash, route, TCP checksum, TOE,
ZIP, reg-exp matcher, malloc, QoS; crypto: DES,
3DES, AES, RSA, DH, MD5, SHA, RC4



Raza XLR

- Values: low power, size, cost, parts count
  - $f$   Allow appliance/accelerator focus on new value add, not standard stuff
  - $f$   Uses: SSL, VPN, VoIP, virus scanning, etc.
- **Key IP: packet scanning, classification, & ultra-fast dispatch to threads.**
  - $f$   **RTOS, privileged mode code.**
- Why? One reason:
  *40% of TCP packets are 40-byte TCP ACKs.*
  - $f$   At 10Gb/s, one every 32 ns.

**Vendors:** Cavium, P.A. Semi, Raza, Broadcom (SyByte), PMC Sierra, Freescale

# Industry Initiatives for Accelerator Attachment

- **Geneseo**: Extending PCI Express® (ann. 9/26/06)
  - *f* Proposed by Intel and IBM
    - n 30+ companies support at announcement
  - *f* Goal: new open standard for attaching accelerators and co-processors to server platforms
    - n power management, transaction ordering (on PCIe), atomic ops, software overhead, faster data rate
  - *f* Developing through and with the PCI SIG
    - n May ultimately be called "PCI 3.0" or "2.1" or the like.

- **Torrenza**: Extending HyperTransport™ (ann, 6/1/06)
  - *f* Proposed by AMD
    - n 6+ companies support at announce, plus consortia
  - *f* Goal: improve support for the integration of specialized coprocessors in systems based on AMD Opteron microprocessors
    - n both I/O (HT) and SMP socket (cHT).
  - *f* Allied with HyperTransport Consortium and OpenFPGA consortium

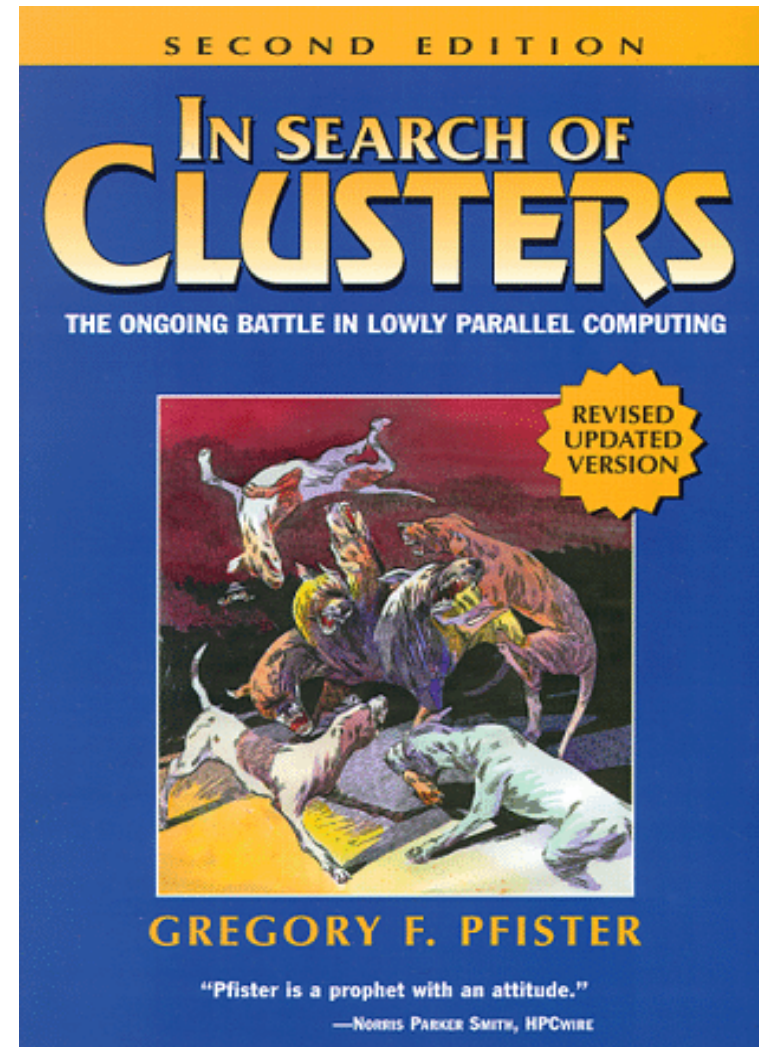| | |
|---|---|
| Will full or partial offload be the norm for Ethernet? | Partial |
| Will any collective functions *typically* be offloaded to the NIC in the near future? | No |
| What about collective support in switches {*typically* }? | Heck no |
| Should NICs assist in "read-modify-write" capability to support remote updates? | No, but. |

## *Required*

| | |
|---|---|
| What communication functions should **always** be offloaded to accelerators? | ACK $\Rightarrow$ RDMA, packet classification, work "dispatch" |
| Which should **never** be offloaded to accelerators? | Everything else. |
| Are there cases for which it makes sense to offload to another core rather than to an accelerator? (For SMT cores, I suppose one might also consider offloading to other threads of the same core?) | That's where the "everything else" goes. |

- # Thank you for listening.

- # Any (more) questions?

**Just in case any of you were wondering...**

**(No, I can't give a presentation without plugging my book.)**

**SECOND EDITION**

**IN SEARCH OF CLUSTERS**

**THE ONGOING BATTLE IN LOWLY PARALLEL COMPUTING**

REVISED UPDATED VERSION

**GREGORY F. PFISTER**

"Pfister is a prophet with an attitude."
—Norris Parker Smith, HPCwire

**Extremely nonrandom clipart**