# 64-Bit Floating-Point Accelerators for HPC Applications

**John L. Gustafson, Ph.D.**

**Chief Technology Officer, HPC**

**ClearSpeed Technology, Inc.**

# Outline

- **Acceleration issues in general: host-accelerator bandwidth/latency**

- **Accelerator performance analysis examples**

- **Measured and expected application acceleration:**
  - Molecular Dynamics: AMBER and NAB
  - Quantum Chemistry: GAUSSIAN, Qbox, PARATEC
  - Monte Carlo models for PDEs
  - LS-DYNA and ANSYS
  - PAM-CRASH
  - MATLAB applications
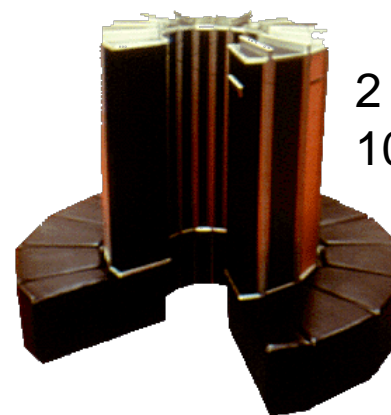
- **Summary**

www.clearspeed.com

# Thesis

- **Performance analysis for accelerator cards is like analysis for message-passing parallelism, but with more levels of memory and communication.**

- **Application porting success depends heavily on attention to memory bandwidths, but (surprisingly) not so much the host-accelerator bandwidth.**

www.clearspeed.com

# The accelerator idea is as old as supercomputing itself



3 MB/s ←→

2 MB
10x speedup

**General-purpose computer
Runs OS, compilers, disk,
printers, user interface**

**Attached vector processor
accelerates certain
applications, but not all**

Even in 1977, HPC users faced issues of when it makes
sense to use floating-point-intensive vector hardware.

"History doesn't repeat itself, but it does rhyme."
—*Mark Twain*

www.clearspeed.com

# IDC survey of planned HPC accelerator use

## Study Results: Accelerators

### Use of Applications Accelerators

Q: Do you have any plans to use applications accelerators?
Multiple responses allowed.

| Accelerator | Number of Mentions | Percentage of Responding Sites Mentioning |
|---|---|---|
| FPGAs | 28 | 90.3% |
| Vector coprocessors | 13 | 41.9% |
| GPUs | 10 | 32.3% |
| Other | 1 | 3.2% |

n = 31

IDC
Analyze the Future
Source: IDC, 2006

©2007 IDC  24

www.clearspeed.com

# Is this trip necessary? Bandwidth issues

| Node memory | | Accelerator memory |
|---|---|---|

Bandwidth = $B$

| Node | Accelerator |
|---|---|

- **Acceleration software tests candidates for work on the board. If too small, it leaves them on the host.**
- **Performance claims *must* assume host-resident data. Beware of benchmarks that leave out the time to move the data to accelerator memory.**
- **Remember Bailey's *12 Ways…***

*speed*

break-even

**accelerator**

**node**

*time*
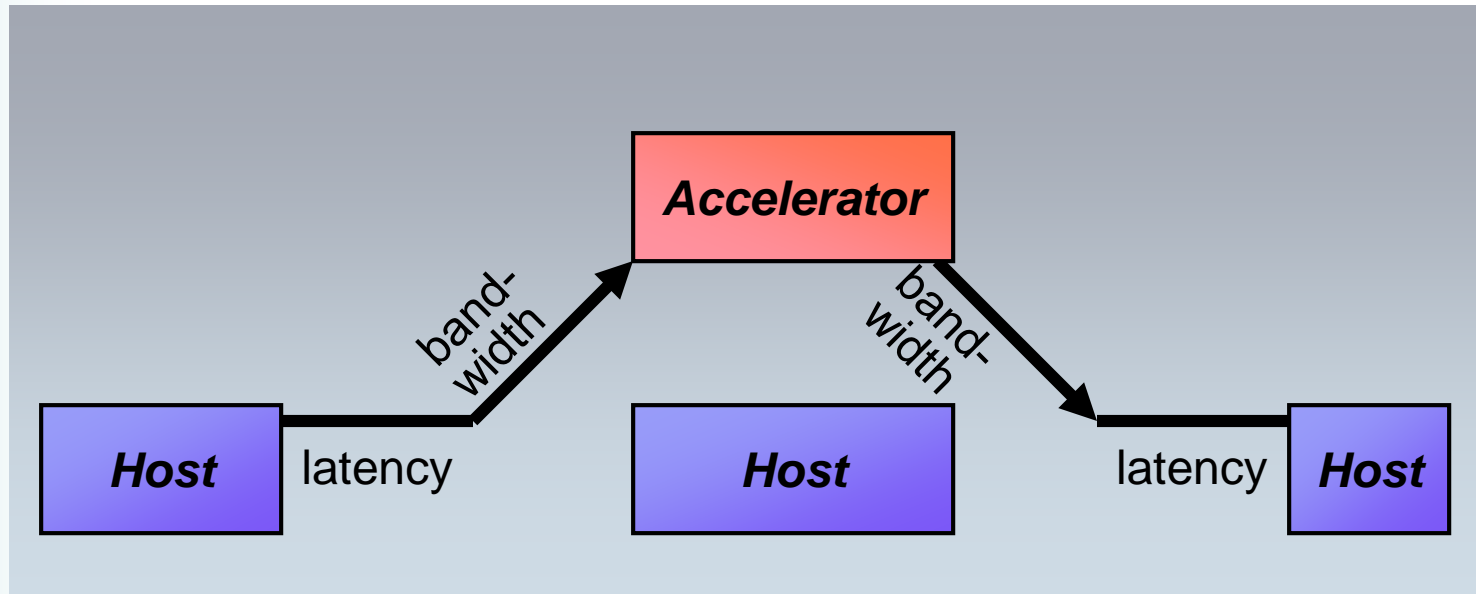*(larger problem size)*

**www.clearspeed.com**
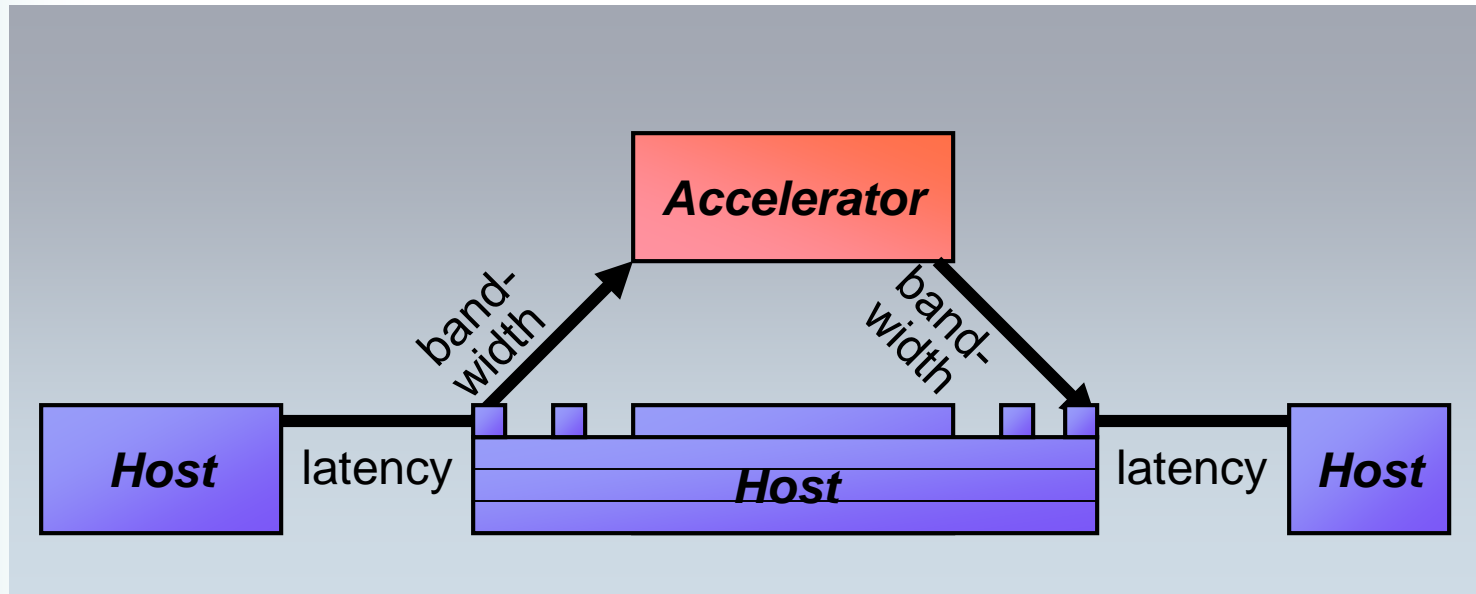
# Simple offload model is out of date



- **Accelerator must be quite fast for this approach to have benefit**
- **This "mental picture" may stem from early days of Intel 80x87, Motorola 6888x math coprocessors**

www.clearspeed.com

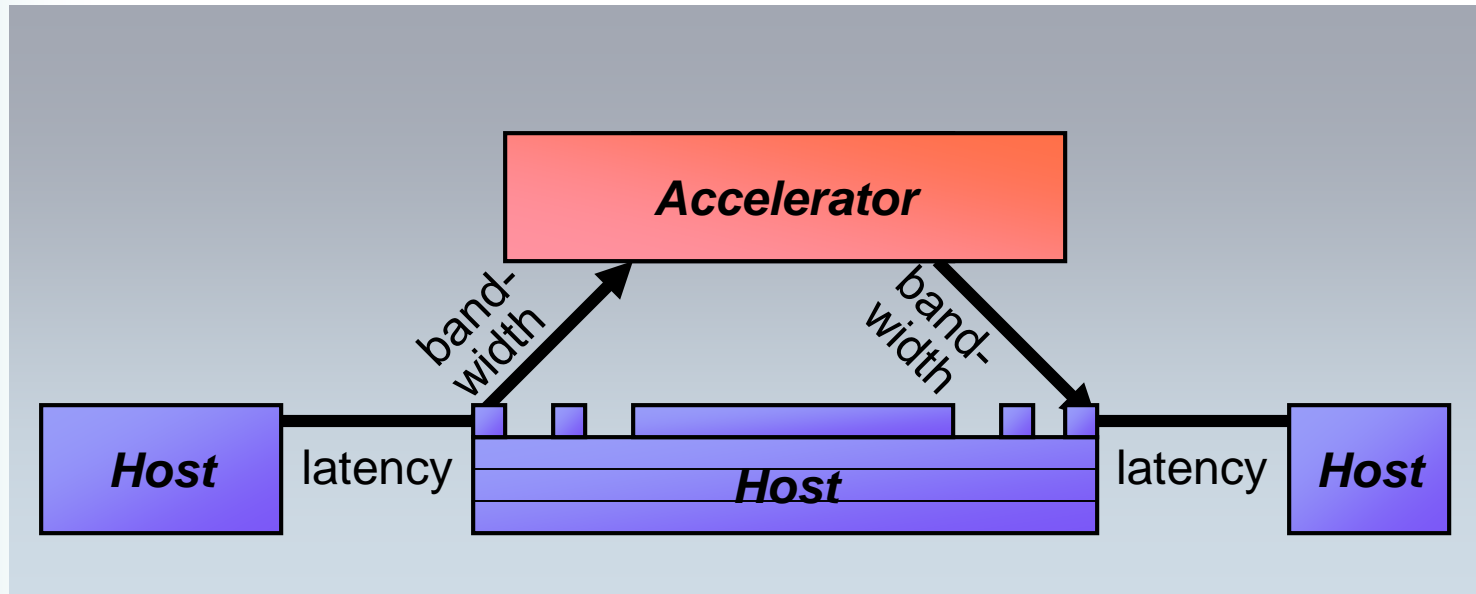# Acceleration model: Host continues working



- **Accelerator needs only be fast enough to make up for time lost to bandwidth + latency**
- **Easiest use model: host and accelerator share the same task, load balanced to complete at same time**
- **More flexible: Host, accelerator *specialize* what they do**

**www.clearspeed.com**
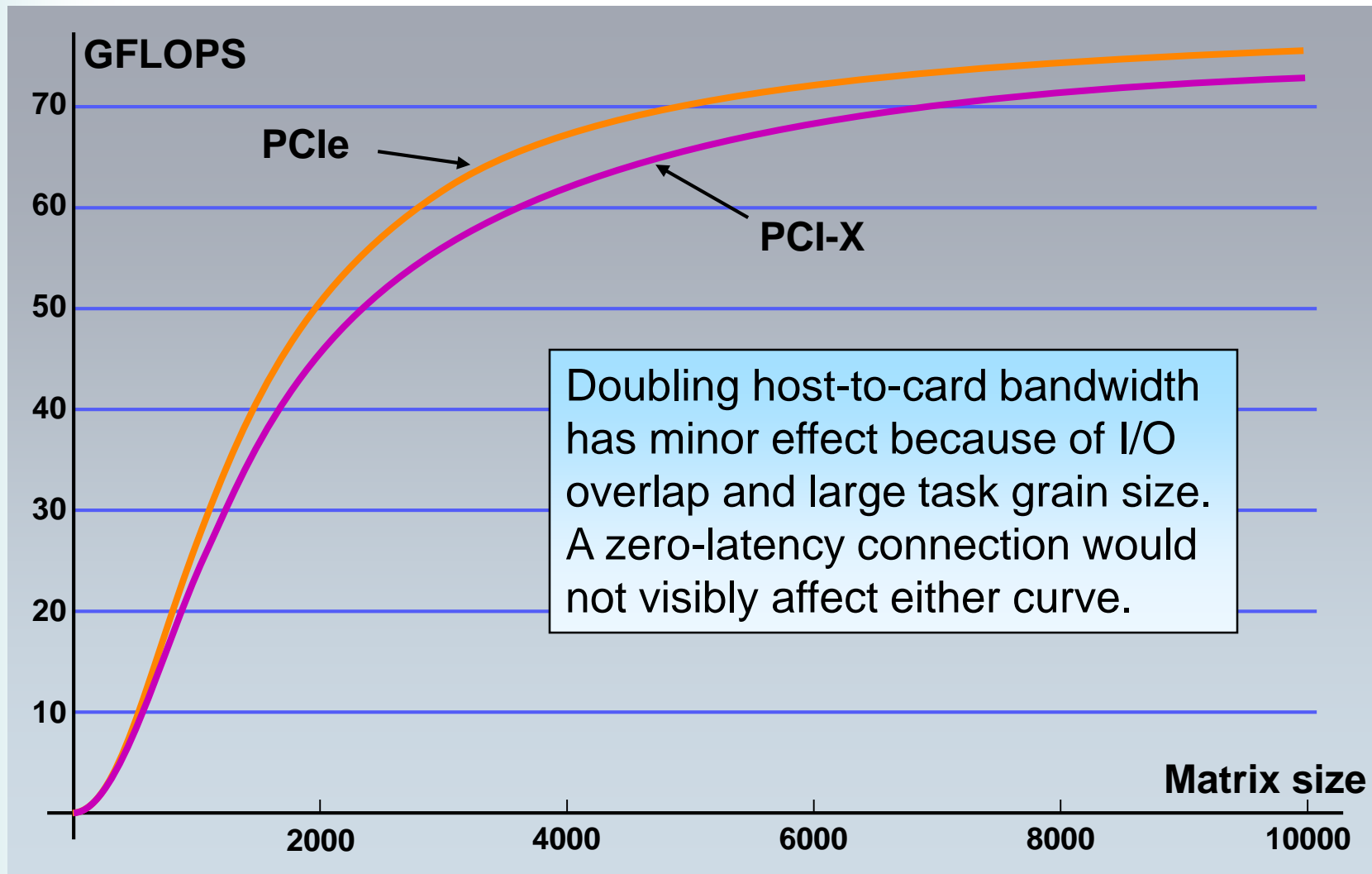
# Host can work while data is moved



- **PCI transfers might burden a *single x86 core* by 60%**
- **Other cores on host continue productive work at full speed**

# Card need not wait for *all* data before starting

**Accelerator**

band-width

band-width

**Host**        latency        **Host**        latency        **Host**

- **In practice, latency is *microseconds*; the accelerator task takes *seconds*. Latency gaps above would be microscopic if drawn to scale.**

- **The accelerator can be *slower* than the host, and still add performance!**

**www.clearspeed.com**

# Square DGEMM speeds as of December 2006

Note: curve only samples integer multiples of vector size

www.clearspeed.com

# Accelerator memory hierarchy
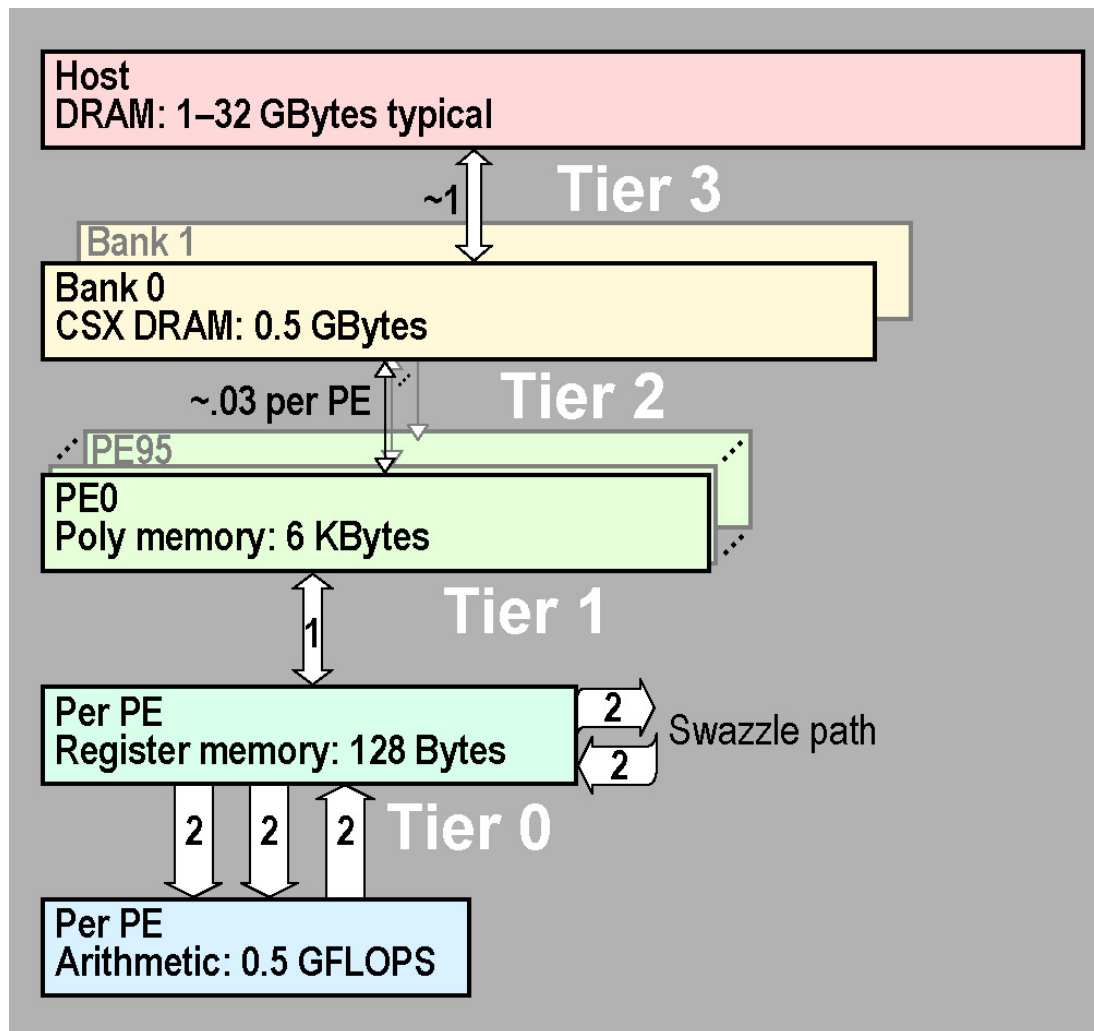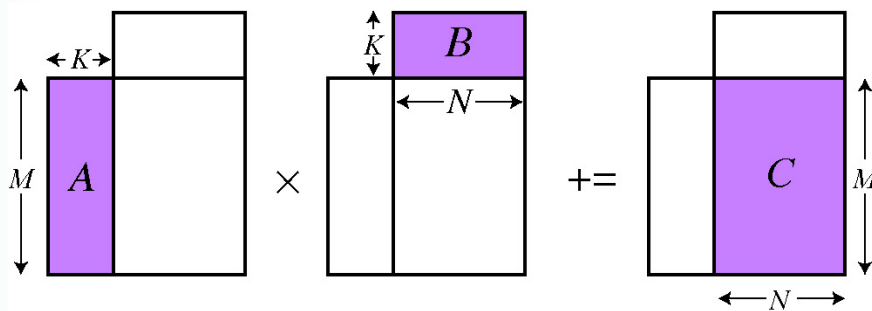
Total: 1.0 GB

Total: 6.4 GB/s

Total: 1.1 MB

Total: 192 GB/s

Total: 24 KB

Total: 2 TB/s

Total: 96 GFLOPS
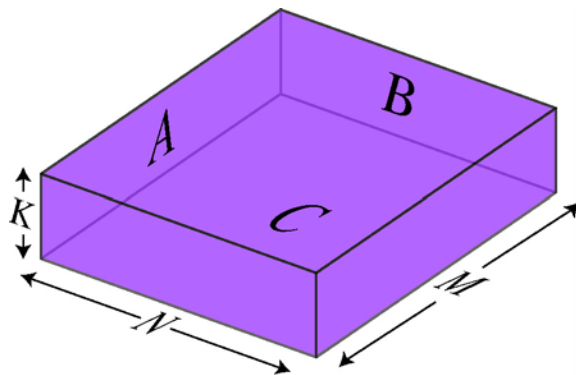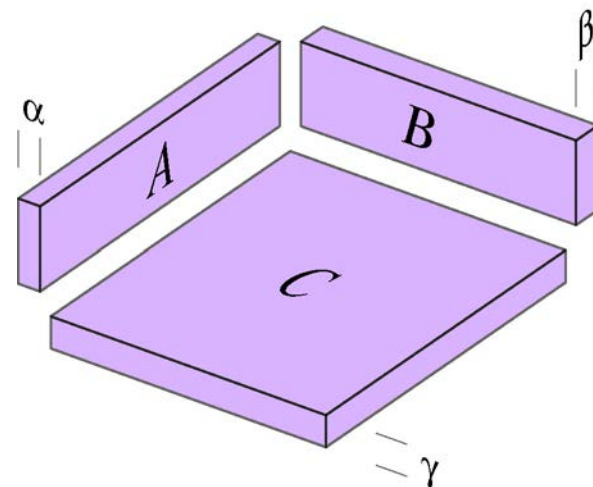(but only 25 watts)

**Host**
**DRAM: 1–32 GBytes typical**

~1    **Tier 3**

Bank 1

**Bank 0**
**CSX DRAM: 0.5 GBytes**

~.03 per PE    **Tier 2**

PE95

**PE0**
**Poly memory: 6 KBytes**

1    **Tier 1**

**Per PE**
**Register memory: 128 Bytes**    2 / 2    Swazzle path

2  2  2    **Tier 0**

**Per PE**
**Arithmetic: 0.5 GFLOPS**

www.clearspeed.com

Matrix multiply (DGEMM) is a perfect analog to a folded box.



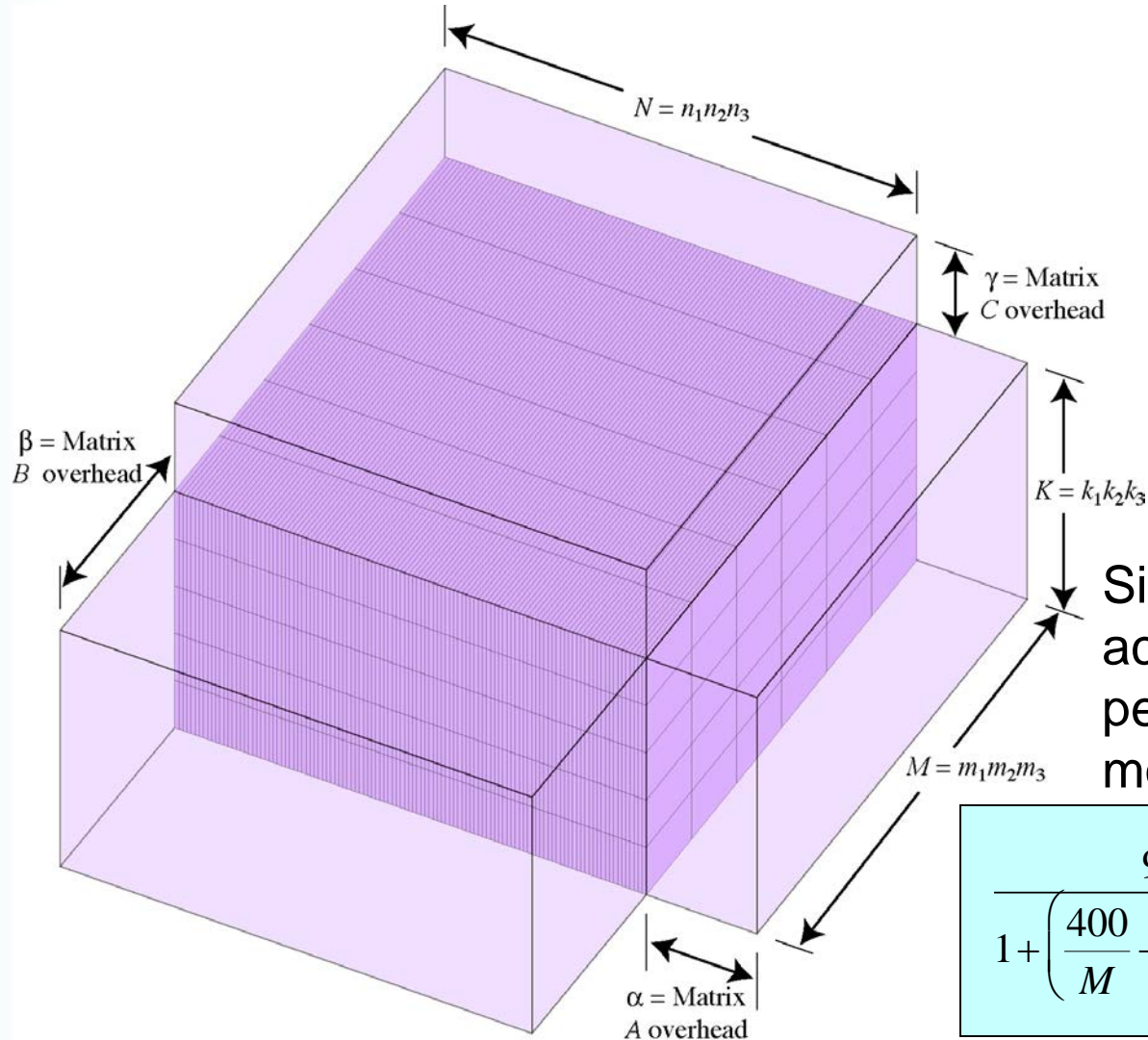Volume is the number of multiply-adds.



Surface "padding" shows overheads

www.clearspeed.com

Tier 1

Tier 2

www.clearspeed.com

$N = n_1 n_2 n_3$

$\gamma = $ Matrix $C$ overhead

$\beta = $ Matrix $B$ overhead

$K = k_1 k_2 k_3$

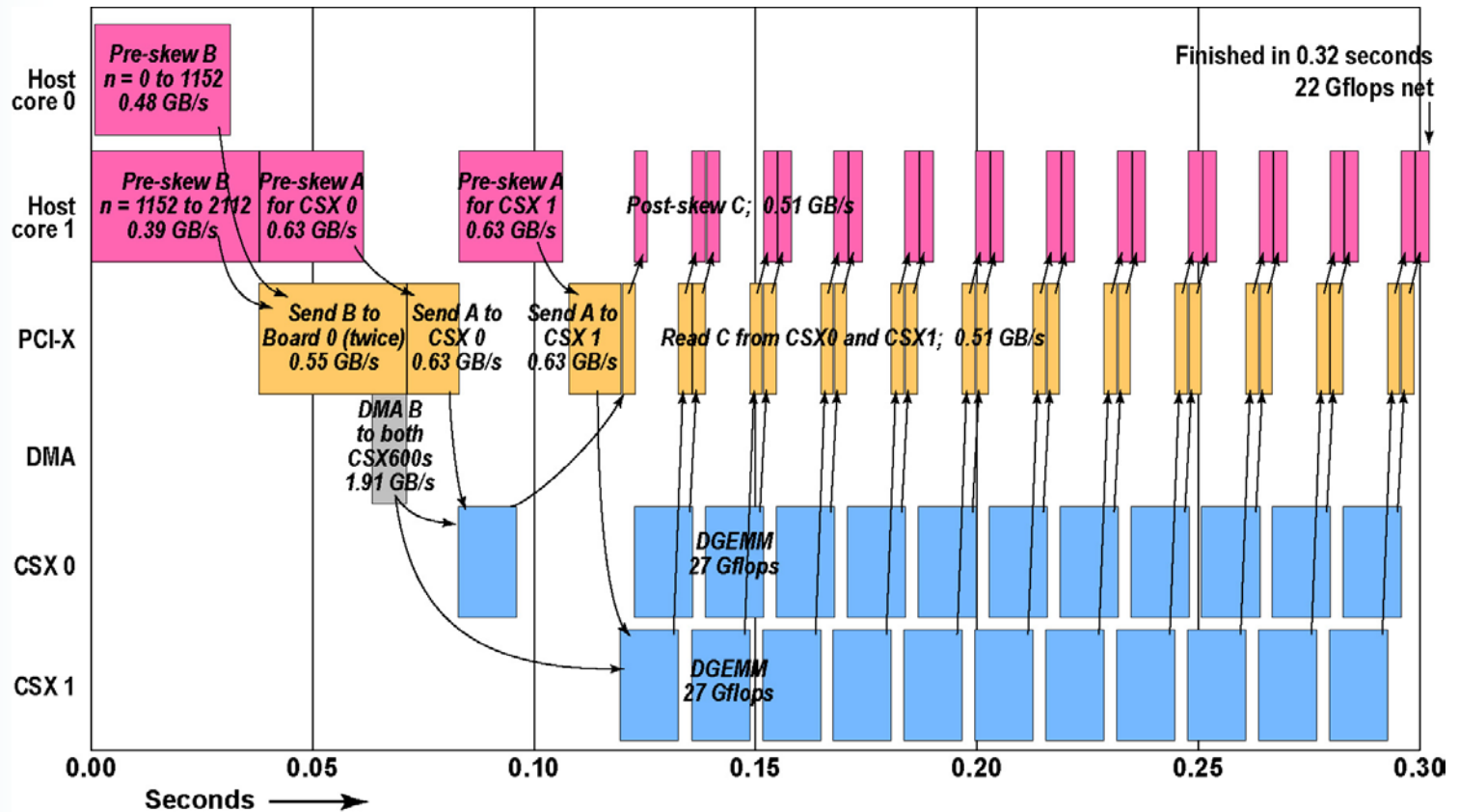$M = m_1 m_2 m_3$

$\alpha = $ Matrix $A$ overhead

Simple but decent accelerator performance model:

$$\frac{90.2}{1 + \left( \dfrac{400}{M} + \dfrac{200}{K} + \dfrac{200}{N} \right)} \text{ GFLOPS.}$$

www.clearspeed.com

# Unoptimized Tier 3 timing

## $K=960$, $M=1920$, $N=1920$

www.clearspeed.com

# Optimized Tier 3 timing



Can now accelerate matrices as small as N=576

www.clearspeed.com

# Almost doubled sustained speed



**GFLOPS**

**Revised DGEMM asymptotic to ~90 GFLOPS**

Does not include host contribution, which adds 5 to 60 GFLOPS depending on host

**Previous DGEMM asymptotic to ~56 GFLOPS**

**Matrix size**
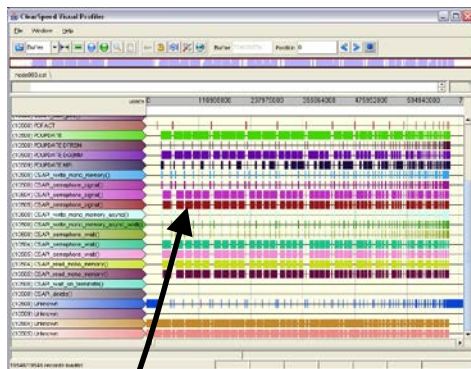
Note: curve only samples integer multiples of vector size

# ClearSpeed™

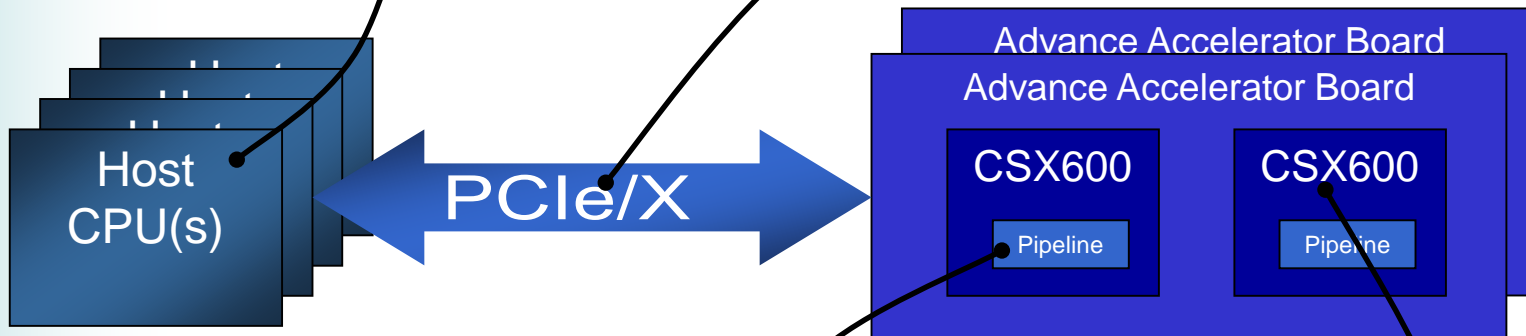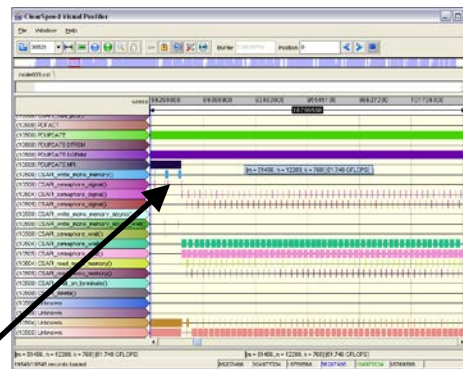# Detailed profiling is essential for accelerator tuning

**HOST CODE PROFILING**

Visually inspect multiple host threads.
Time specific code sections.
Check overlap of host threads

**HOST/BOARD INTERACTION**

View host/board interactions
Measure transfer bandwidth.
Check overlap of host and board compute

Advance Accelerator Board
Advance Accelerator Board

Host
CPU(s)

PCIe/X

CSX600
Pipeline

CSX600
Pipeline

**ACCELERATOR PIPE**

View instruction issue.
Visualize overlap of executing instructions.
Get cycle-accurate timing.
Remove instruction-level performance bottlenecks.

**CSX600 SYSTEM**

Trace at system level.
Inspect overlap of compute and I/O.
View cache utilization.
Graph performance.

# Pipeline view of DGEMM inner loop



Profile the code running at the instruction level

See the pipeline performance for each instruction

Tune the instruction scheduling for the application code

www.clearspeed.com

# System level: multiple DGEMM calls



View the DGEMM calls on host and accelerator

Each call syncs up with the host view of accelerator

Much higher level of detail available from the profiler

Dongarra's new LAPACK
will make QR, Cholesky,
LU factoring, etc. much
easier to accelerate.

Volume = ⅓ $N^3$
multiply-adds

DGEMM

DGEMM

DGEMM

DGEMM

*Ax = b*

*N* iterations

*N* equations

*N* unknowns

Excellent test of hardware correctness

www.clearspeed.com

# Accelerated cluster model accurate to about ±5%

$$R_{est} = \cfrac{1}{\cfrac{1}{PQ\gamma} + \cfrac{3\alpha\left[(N_B + 1)\lg P + P\right]}{2N^2 N_B} + \cfrac{3\beta(3P + Q)}{4NPQ}}$$

$N_B$ = Block size, the width of the "panels" used to update the linear system with DGEMM

$N$ = Dimension of the linear system (number of equations to solve)

$P, Q$ = Dimensions of the two-dimensional mapping of computational nodes

$\alpha$ = Effective point-to-point latency of MPI broadcast, in seconds

$\beta$ = Effective point-to-point reciprocal bandwidth of message broadcast, in seconds per datum

$\gamma$ = Effective floating-point operations per second of a node independent of MPI operations

# Memory bandwidth dominates performance model

| Node memory |
|---|

**17 GB/s**    ⟷ (PCIx or PCIe 1 to 2 GB/s) ⟷    | Accelerator DRAM |

| Multicore x86 |
|---|

**6.4 GB/s**

| Accelerator |
|---|

**192 GB/s**

| Accelerator Local RAM |
|---|

- **Apps that can stage into local RAM (Tier 1) can go 10x faster than current high-end Intel, AMD hosts**
- **Apps that must reside in DRAM (Tier 2) will actually run *slower* by about 3x (for fully optimized host code)**
- **Fast Fourier Transforms can go either way!**

# Math functions reside at Tier 1, hence fast

**64-bit Function Operations per Second (Billions)**

Legend:
- 2.6 GHz dual-core Opteron
- 3 GHz dual-core Woodcrest
- ClearSpeed Advance card

Y-axis: 0.0, 0.5, 1.0, 1.5, 2.0, 2.5

X-axis (Function name): Sqrt, InvSqrt, Exp, Ln, Cos, Sin, SinCos, Inv

**Typical speedup of ~8X over the fastest x86 processors, because math functions stay in the local memory on the card**

# Monte Carlo PDE methods exploit Tier 1 bandwidth

**Real apps do work resembling "EP" of NAS Parallel Benchmarks. "Quants" solve PDEs this way for options pricing, Black-Scholes model (a form of the Heat Equation)**

- **No acceleration:** **200M samples, 79 seconds**
- **1 accelerator:** **200M samples, 3.6 seconds**
- **5 accelerators:** **200M samples, 0.7 seconds**

![ClearSpeed]

# Why do EP-type Monte Carlo apps need 64-bit?

- **Accuracy increases as the square root of the number of trials, so five-decimal accuracy takes 10 billion trials.**
- **But, when you *sum* many similar values, you start to scrape off all the significant digits.**
- **64-bit summation needed to get a single-precision result!**

Single precision:
$1.0000 \times 10^8 + 1$
$= 1.0000 \times 10^8$

Double precision:
$1.0000 \times 10^8 + 1$
$= 1.00000001 \times 10^8$

www.clearspeed.com

# We may need to rethink 64-bit flops…

- **Every operation has an optimum number of bits of accuracy**
  - Using too few gives unacceptable errors
  - Using too many wastes memory, bandwidth, joules, dollars.
- **It is unlikely that a code uses *just the right amount* of precision needed.**



Optimum precision

IEEE 754 double precision

All floating-point operations in application →

www.clearspeed.com

# How do HPC programmers pick FP precision?

- **Assume 64-bit is plenty, and use it everywhere.**

- **Use what is imposed by hardware (word size).**

- **Try two precisions; if answers agree, use the less precise one, otherwise use the more precise one.**

- **Compare computed answer for special cases where an analytic answer is known.**

- **Compare computed answer with physical experiment (rare).**

- **Perform careful analysis (very rare).**

**www.clearspeed.com**

## Can a tool estimate joules, W, $, min. precision?

64-bit op
42 bits needed
10 pJ
$1x10^{-15}$

64-bit op
39 bits needed
12 pJ
$1x10^{-15}$

```
a(i,j) -= a(i,k) * a(k,j)
```

Tier 0 read/write
28 pJ
$2x10^{-15}$

Tier 1 read
50 pJ
$8x10^{-14}$

Tier 2 read
1900 pJ
$5x10^{-13}$

Cost and electrical power and precision are almost as important as timing… why not develop analysis tools for them? You can only optimize what you can measure.

# Model/measure power use, not just performance

**Standard Cluster Power Consumption / Node**



Base System: HP DL380 G5, Intel Xeon 5160 x 2 @ 3GHz, 14GB

www.clearspeed.com

# Accelerated Cluster Node

**Accelerated Cluster Power Consumption / Node**



**Average Power Consumption: 437.5 watts**

(y-axis: watts, values 250, 300, 350, 400, 437.5, 450, 500, 550)

**20 Minutes**

Base System: HP DL380 G5, Intel Xeon 5160 x 2 @ 3GHz, 14GB
Accelerated System:  As above + 1 ClearSpeed Advance X620

**www.clearspeed.com**

# ClearSpeed™

## *Energy* Usage: Standard vs. Accelerated

### Comparative Cluster Power Consumption / Node

— Standard Cluster Node    — CS Accelerated Node

**53.6% Less Energy Used**

watts: 550, 500, 450, 400, 350, 300, 250

**42 Minutes**

Base System: HP DL380 G5, Intel Xeon 5160 x 2 @ 3GHz, 14GB
Accelerated System:  As above + 1 ClearSpeed Advance X620

**www.clearspeed.com**

# AMBER 9 acceleration

- **Model: memory motion is $k_1N$, operations are $k_2N^2$. Overheads easily overcome for typical $N$.**

- **~4x speedup for pharmaceutical company production runs achieved recently.**

- **225 hour (Opteron) run reduced to 58 hours**

- **Didn't exploit symmetry of $i\,j$ forces, which will give another ~1.5–1.9x speedup (31–39 hours)**

- **Solvation model GB1; 500000 time steps**

- **Correctness checked incrementally throughout conversion process.**

**www.clearspeed.com**
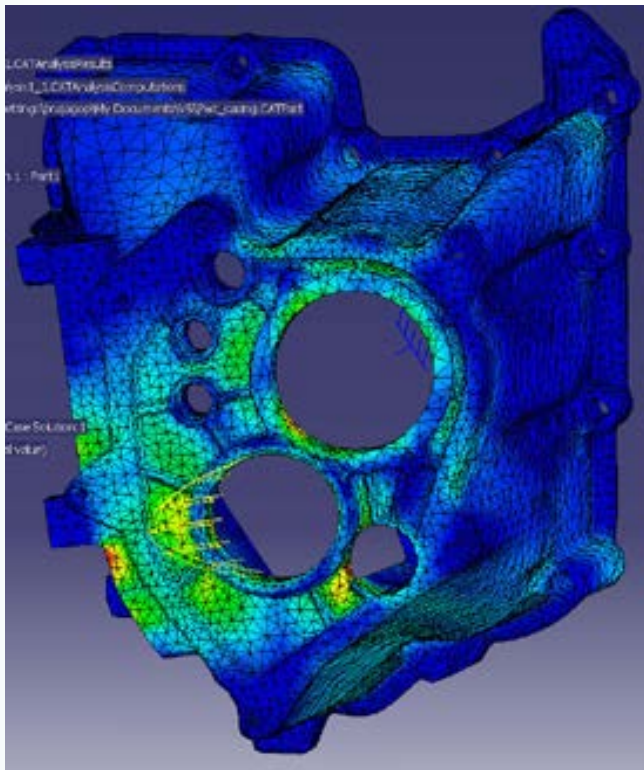
# NAB and AMBER 10 acceleration

- **Newton-Raphson refinement now possible; analytically-computed second derivatives**

- **2.6x speedup obtained for this operation in three hours of effort (no source code changes)**

- **Enables accurate computation of entropy and Gibbs free energy for first time.**

- **Available now in NAB (Nucleic Acid Builder) code. Slated for addition to AMBER 10.**

www.clearspeed.com

# Plug-and-play quantum chemistry acceleration?

- **DGEMM content is 18% to 65% in GAUSSIAN test suite, but typical sizes only ~10 to 100.**

- **No changes to license or to any of the source code. Just invoke dynamic linking option in makefile.**

- **Sample GAUSSIAN tests to date are *too small* to accelerate; below *N* = 576 threshold.**
  - Need larger problem sizes
  - Realistic to be that large? (Lots of occupied orbitals)

- **PARATEC, Qbox much better candidates. Plane wave models are over half DGEMM, huge dimensions**

# The economics of CAE acceleration

**Structural Analysis
ANSYS, LS-DYNA implicit**



- **Each host costs $3,000.**
- **Software license costs ~$30,000 *per core*, which discourages use of multiple cores.**
- **MCAE engineer costs over $200,000/year.**
- **In California, anyway.**
- **Accelerator card would be cost-effective even with a 7% performance boost. Actual performance boost should be more like 260% for large problems.**
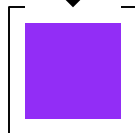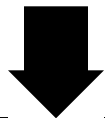
**www.clearspeed.com**

# Accelerating ANSYS, LS-DYNA with Lucas' solver

10 million degrees of freedom (sparse)

becomes…

50,000 *dense* equivalent

Accelerator can solve at over 50 GFLOPS

**Non-solver**
**Solver setup**

**DGEMM**
**on x86**
**host**

**Est. 260% net application acceleration**

**10x**

**Non-solver**
**Solver setup**

**DGEMM**
**with ClearSpeed**

- **Potentially pure plug-and-play**
- **No added license fee**
- **Needs ClearSpeed's 64-bit precision and speed**
- **Enabled by recent DGEMM improvements; still needs symmetric $A^TA$ variant**
- **Could enable some CFD acceleration (for codes based on finite elements, low Reynolds numbers`)**

www.clearspeed.com

## PAM-CRASH and MATLAB acceleration

- **Work done in China shows 1.40x PAM-CRASH speedup using one ClearSpeed accelerator**

- **We await details; this is preliminary**
  - Problem size?
  - Explicit or implicit?
  - What was done to the code?
  - Compared to what host?

- **Japanese industrial researcher got 5x acceleration of MATLAB waveguide model; won't allow publication of results (?)**

- ***sigh* So we continue citing TOP500 results…**

**www.clearspeed.com**

# Summary

- **Accelerator tuning demands attention to memory bandwidth at all levels (bandwidth to host, less so)**

- **Now seeing value for real 64-bit applications in chemistry, electromagnetics, financial modelling, crash codes, etc.**

- **Fit-for-purpose analysis starts with analytical model based on memory tiers, but is verified using detailed performance tools.**

**www.clearspeed.com**