

“Physics will converge our programming models” (nVIDIA promise/threat 😊)

- **Clock rates cease to increase while arithmetic capacity continues to increase dramatically w/concurrency, consistent with Moore’s Law**
 - **Storage capacity diverges exponentially below arithmetic capacity**
 - **Transmission capacity diverges exponentially below arithmetic capacity**
 - **Mean time between hardware interrupts shortens**
-

Apps programming model convergence

- **Billions of dollars of scientific software hang in the balance while better formulations and algorithms evolve to span the architectural gap**
 - **Historically, in the 40 apps surveyed by Beckman for the San Francisco IESP meeting, explicit message passing is effectively hidden from most developers by well developed libraries**
 - **Similarly, it will eventually be easier to program hybrid, heterogeneous architectures than to do user-managed data placement**
 - **In the meantime...**
-

Applications Agenda

- **New hardware-tolerant formulations with**
 - ◆ **greater arithmetic intensity (flops per bytes moved into and out of registers and upper cache)**
 - ◆ **reduced communication**
 - ◆ **reduced synchronization**
 - ◆ **assured accuracy with (adaptively) less floating-point precision**
 - **Quantification of trades between limiting resources**
 - ***Plus* enhancing applications to obtain all of the exciting payoffs that exascale is meant to exploit**
-

Payoffs

- **Resolve extreme scales**
 - **Accommodate full physics in full dimensions**
 - **Combine multiple complex models**
 - **Solve an inverse problem, or perform data assimilation**
 - **Perform optimization or control**
 - **Quantify uncertainty**
 - **Validation and verification**
 - **Solve stochastic problems**
 - **Data analytics (including high dimensional tensors)**
-

Why exa- is different...

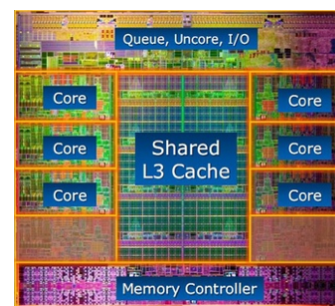
Which steps of FMADD take more energy?

64-bit floating-point fused multiply add

or

moving four 64-bit operands 20 mm across the die

$$\begin{array}{r} 934,569.299814557 \quad \text{input} \\ \times \quad 52.827419489135904 \quad \text{input} \\ \hline = 49,370,884.442971624253823 \\ + \quad 4.20349729193958 \quad \text{input} \\ \hline = 49,370,888.64646892 \quad \text{output} \end{array}$$



20 mm

(Intel Sandy Bridge, 2.27B transistors)

Going across the die requires up to an order of magnitude more !

DARPA study predicts that by 2019:

- ◆ Double precision FMADD flop: 11pJ
- ◆ cross-die per word access (1.2pJ/mm): 24pJ (= 96pJ overall)

A community working model

- Mathematical and modeling ideas should be welcomed from many sources, including nontraditional
 - It is important to show pay-offs on skeleton applications before these ideas *can, should, or will* be adopted
 - Applications communities have detailed roadmaps between now and 2020 apart from adapting to emerging architectures
 - Expertise from a small number of “enlightened” computational mathematicians and computer scientists will have to be packaged for wider adoption, since not every applications community can find or afford all the computational experts it otherwise needs
-

Algorithmic Path Forward

- **Requirements (?) of algorithm and software developers**
 - New depth of knowledge about what the hardware is doing
 - New exercise of control over what the hardware is doing
 - Performance models that are “good enough” to inform decisions
 - More levels of abstraction in representing ideas and coding for multiple architectures
 - Better understanding of numerical error propagation to exploit reduced precision
- **Modes of adaptation to emerging architectures (incl. some opposing tendencies)**
 - Concentrate locality for efficient access vs. relax locality for load-balancing flexibility
 - Aggregate to reduce overhead vs. disaggregate to hide latency
 - redundant or extra work to avoid communication or synchronization
 - Catch and conditionally tolerate errors in user space
 - Checkpoint in user space
 - Apply machine learning to optimize ordering and layout

Algorithmic Path Forward

- **Other ideas still to be redacted from our discussions... stay tuned tomorrow**