# Synergistic Co-design:
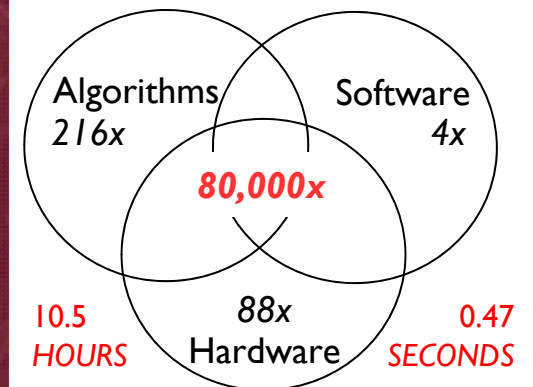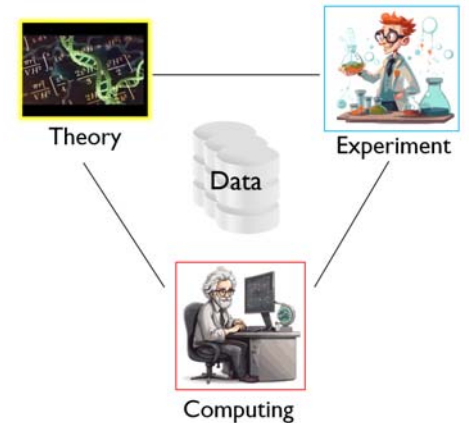# A Multi-Faceted Perspective

Wu Feng
wfeng@vt.edu

QCD at the Femtoscale in the Era of Big Data
Institute of Nuclear Theory



Theory

Data

Experiment

Computing

Algorithms
216x

Software
4x

**80,000x**

10.5
*HOURS*

88x
Hardware

0.47
*SECONDS*

http://www.youtube.com/watch?v=zPBFenYg2Zk

SyNeRG

http://synergy.cs.vt.edu

VIRGINIA TECH.

# A Little Bit About Me …

- Education
  - Ph.D., Computer Science,
    U. Illinois at Urbana-Champaign, 1996

- Professional
  - Current Appointments
    - Professor and Elizabeth & James Turner Fellow; Departments of Computer Science, Electrical & Computer Engg., and Health Sciences; Virginia Tech
    - Director, **SyNeRG** Laboratory (http://synergy.cs.vt.edu/) → SEEC Center
    - Site Director, Center for Space, High-performance, and Resilient Computing (SHREC)
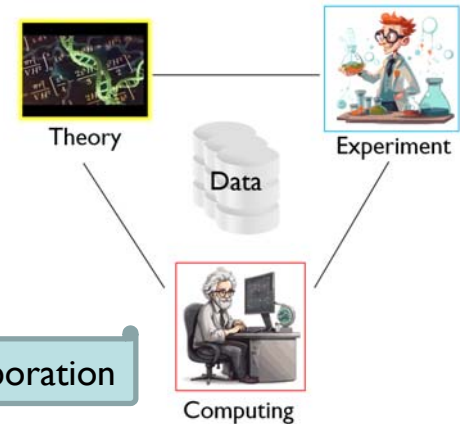  - Previous Appointments & Professional Stints
    - *Academia:* Ohio State U. ('00-'03), Purdue U. ('98-'00), U. of Illinois at Urbana-Champaign ('96-'98)
    - *Government:* Los Alamos Nat'l Lab ('98-'06), NASA Ames Research Ctr ('93)
    - *Industry:* IBM T.J. Watson Rsch ('90), Vosaic ('97), Orion Multisystems ('04-'05), EnergyWare ('08-'10)

# Summary: Re-visiting the Third Pillar of Science

1. The third pillar of science is simply COMPUTING, encompassing simulating physical reality **and** computing on the data. *... for the technical layperson*

2. Synergistic co-design of algorithms, software, and hardware can massively accelerate discovery. *... for scientific collaboration*

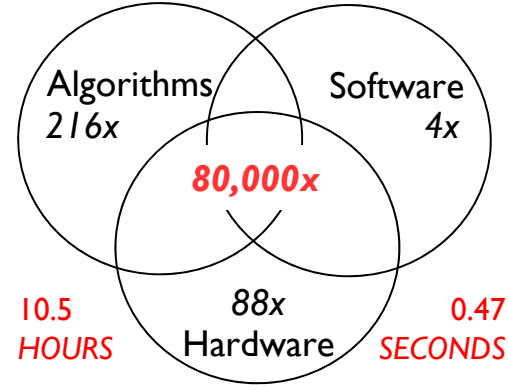3. Don't fool yourself and, in turn, fool the masses. *... for computer scientists*

Theory — Experiment

Data

Computing

## 12 Ways to Fool the Masses
(David H. Bailey, NASA & LBL, 1991)

1. Quote only 32-bit performance results, **not** 64-bit results. → 2x speedup
2. Present performance figures for an inner kernel and then represent these figures as the performance of the entire application.

7. When direct run-time comparisons are required, compare with an old code on an obsolete system.
8. If Mflop/s rates must be quoted, base the operation count on the parallel [version], **not** on the best sequential [version].
9. Quote performance in terms of processor utilization, parallel speedups, or Mflop/s per dollar.
10. Mutilate the algorithm used in the parallel implementation to match the architecture.
11. Measure parallel run-times on a dedicated system but measure conventional run times in a busy environment.
12. If all else fails, show pretty pictures and videos, and don't talk about performance.

### Debunking the 100X GPU vs. CPU Myth:
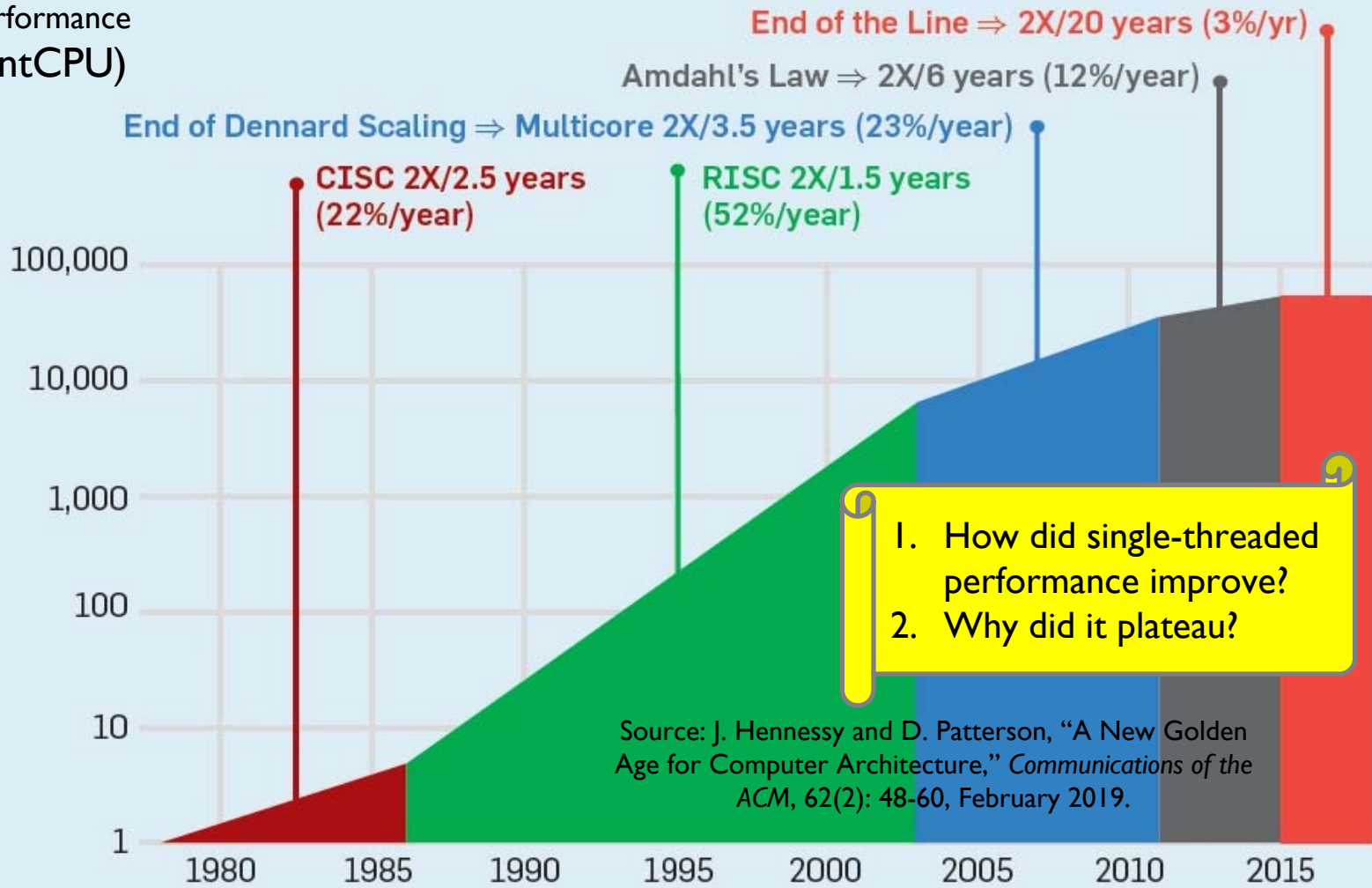An Evaluation of Throughput Computing on CPU and GPU

Victor W Lee†, Changkyu Kim†, Jatin Chhugani†, Michael Deisher†,
Daehyun Kim†, Anthony D. Nguyen†, Nadathur Satish†, Mikhail Smelyanskiy†,
Srinivas Chennupaty⋅, Per Hammarlund⋅, Ronak Singhal⋅ and Pradeep Dubey†

victor.w.lee@intel.com

†Throughput Computing Lab,        ⋅Intel Architecture Group,
Intel Corporation                 Intel Corporation

Algorithms 216x — Software 4x

**80,000x**

10.5 HOURS

88x Hardware

0.47 SECONDS

http://www.youtube.com/watch?v=zPBFenYg2Zk

SyNeRG
synergy.cs.vt.edu

Serial Performance
(SPECintCPU)

**End of the Line ⇒ 2X/20 years (3%/yr)**

**Amdahl's Law ⇒ 2X/6 years (12%/year)**

**End of Dennard Scaling ⇒ Multicore 2X/3.5 years (23%/year)**

**CISC 2X/2.5 years (22%/year)**

**RISC 2X/1.5 years (52%/year)**

Performance vs. VAX11-780

100,000

10,000

1,000

100

10

1

1980   1985   1990   1995   2000   2005   2010   2015
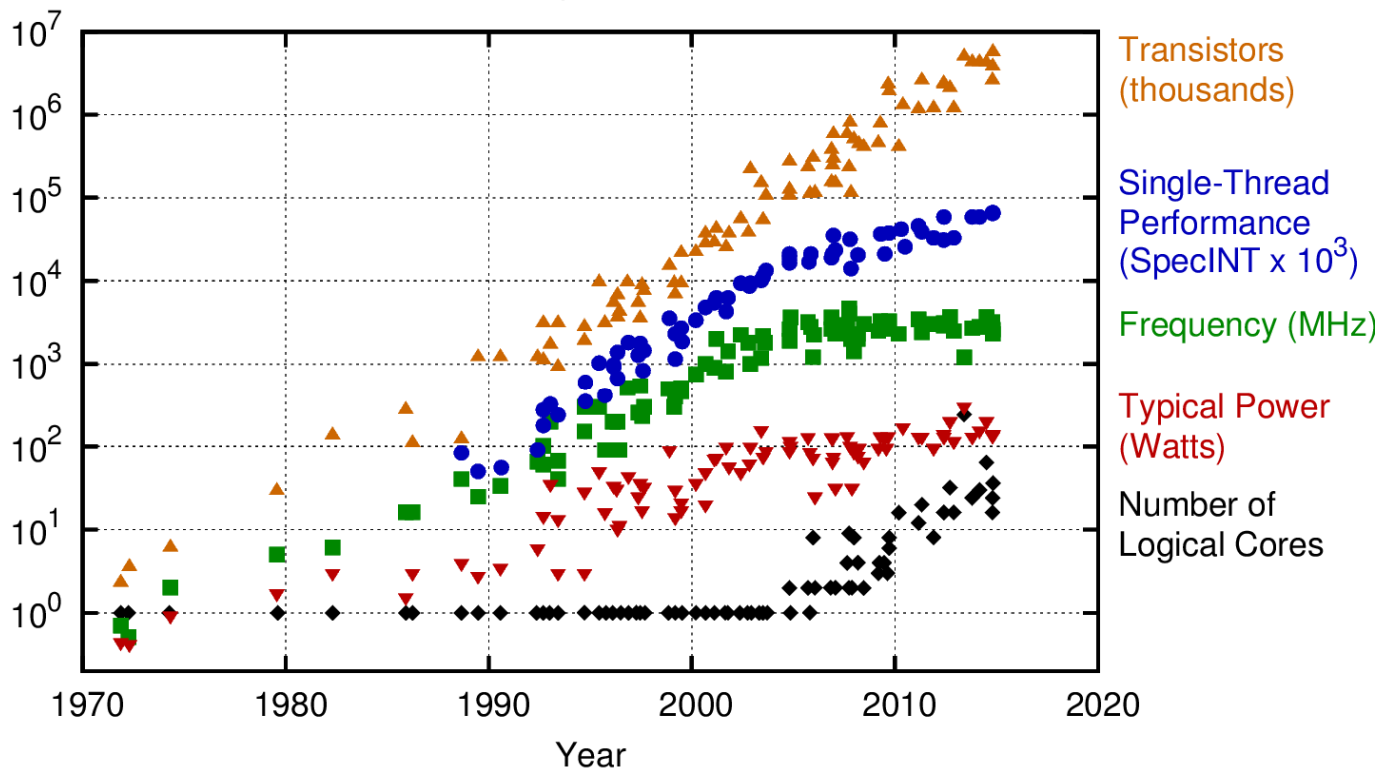
1. How did single-threaded performance improve?
2. Why did it plateau?

Source: J. Hennessy and D. Patterson, "A New Golden Age for Computer Architecture," *Communications of the ACM*, 62(2): 48-60, February 2019.

## 40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)
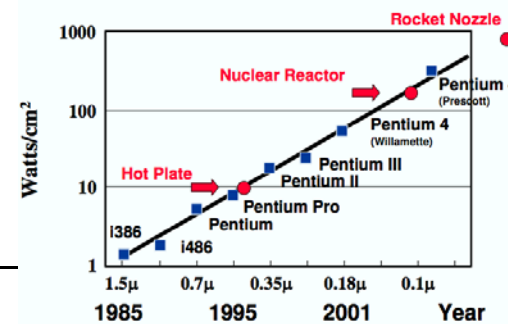
Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

https://www.karlrupp.net/wp-content/uploads/2015/06/40-years-processor-trend.png

**Intel's Paxville: too slow, too hot, too dumb, 2005.**

W. Feng, wfeng@vt.edu, 540.231.1192
QCD at the Femtoscale in the Era of Big Data

VIRGINIA TECH™

Source: Intel

# Growth of DNA Sequencing



Can a similar growth trend be constructed for quantum physics?

Source: Z. Stephens, S. Lee, F. Faghri, R. Campbell, C. Zhai, M. Efron, R. Iyer, M. Schatz, S. Sinha, G. Robinson, "Big Data: Astronomical or Genomical?" *PLoS Biol* 13(7): e1002195, July 2015
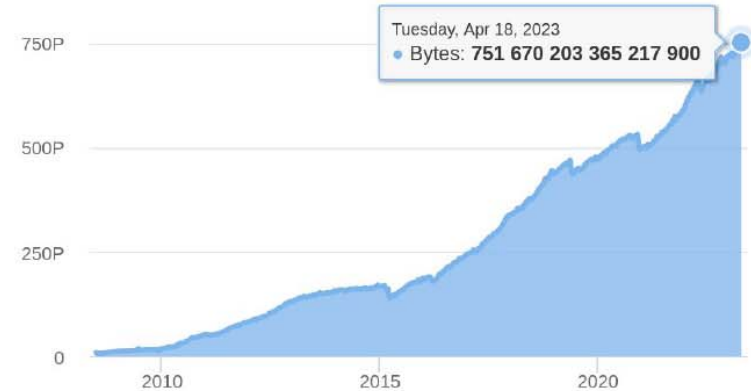
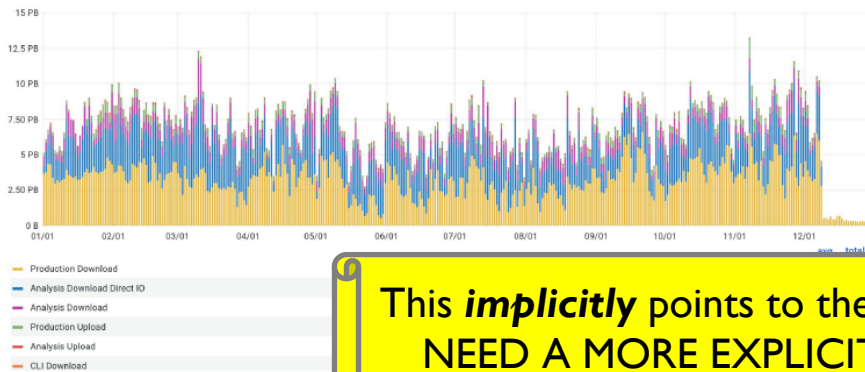# ATLAS Distributed Computing today

- A few numbers showing the scale of ATLAS data
  - 1B+ files, 750+ PB of data, 400+ Hz interaction
  - 120 data centres, 5 HPCs, 3 clouds, 1000+ users
  - 1.2 Exabytes/year transferred
  - 2.7 Exabytes/year uploaded & downloaded

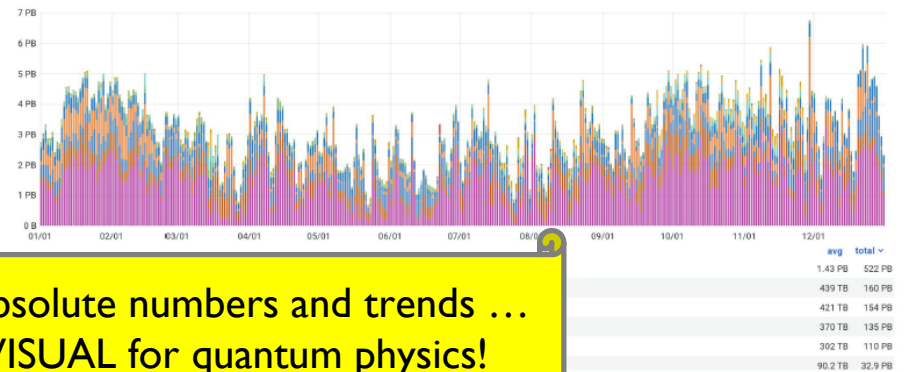- Expect an increase of at least one order of magnitude for the HL-LHC

ATLAS data registered in Rucio

Tuesday, Apr 18, 2023
- Bytes: **751 670 203 365 217 900**

5+ PB/day data access for computation

- Production Download
- Analysis Download Direct IO
- Analysis Download
- Production Upload
- Analysis Upload
- CLI Download

2+ PB/day data transfers between storage

| avg | total |
|---|---|
| 1.43 PB | 522 PB |
| 439 TB | 160 PB |
| 421 TB | 154 PB |
| 370 TB | 135 PB |
| 302 TB | 110 PB |
| 90.2 TB | 32.9 PB |

This *implicitly* points to the absolute numbers and trends ...
NEED A MORE EXPLICIT VISUAL for quantum physics!

D. South. ATLAS Distributed Computing Evolution. CHEP 23, Norfolk, VA, USA. May 2023

4

# Challenge



- The rate of growth in **big data** is *far outstripping* the rate at which computing can (brute-force) **compute** on the data.

# Challenge



- The rate of growth in **big data** is *far outstripping* the rate at which computing can (brute-force) **compute** on the data.
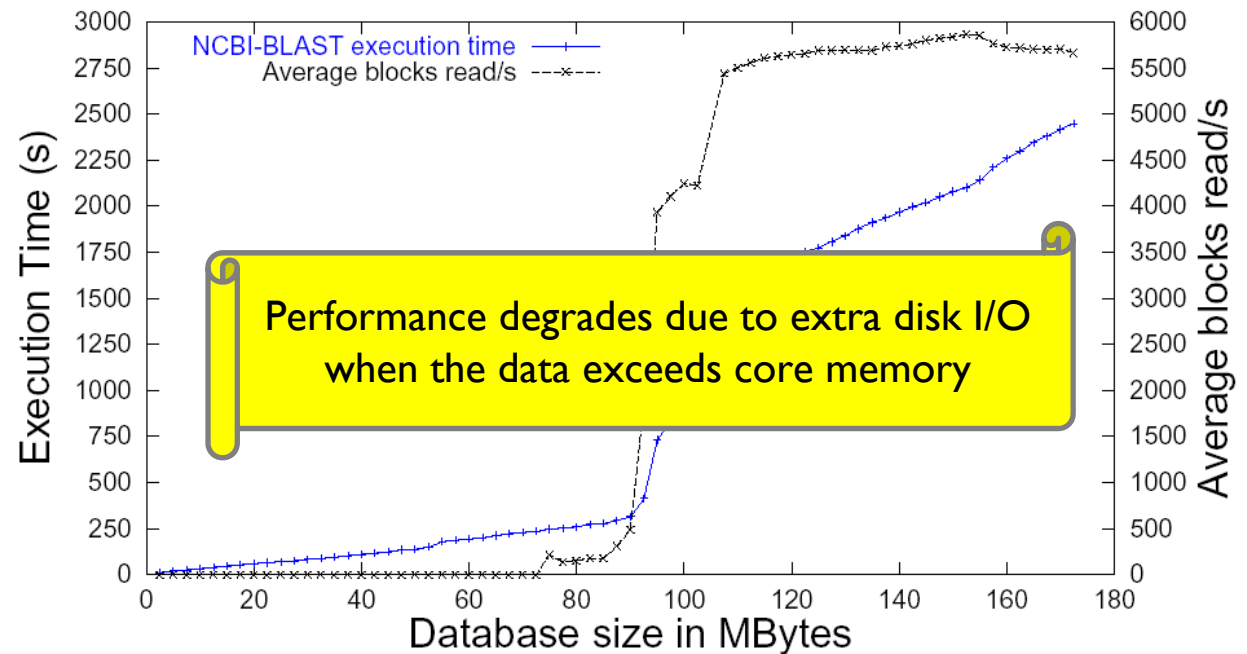
# Approach



- Synergistic co-design of architecture, software, and in particular, algorithms to more *efficiently* and *intelligently* compute on the data.

# Importance of Trend Graphs:  Compute, Data, Compute/Data

- Impacts how programs should be written, e.g., 2004:  BLAST → mpiBLAST

Standard BLAST pairwise
sequence alignment
(128MB RAM)



Performance degrades due to extra disk I/O
when the data exceeds core memory

# Computational Science (→ OpenDwarfs → Berkeley Dwarfs)

**Figure 3. The color of a cell (for 12 computational patterns in seven general application areas and five Par Lab applications) indicates the presence of that computational pattern in that application; red/high; orange/moderate; green/low; blue/rare.**

| | Embed | SPEC | DB | Games | ML | CAD | HPC | Health | Image | Speech | Music | Browser |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. Finite State Mach. | red | red | red | orange | orange | orange | blue | blue | blue | blue | blue | red |
| 2. Circuits | red | blue | green | blue | olive | blue | blue | blue | blue | blue | blue | red |
| 3. Graph Algorithms | red | orange | orange | orange | red | red | blue | red | blue | red | green | green |
| 4. Structured Grid | red | red | blue | orange | blue | blue | red | blue | red | blue | blue | blue |
| 5. Dense Matrix | red | red | orange | red | red | red | blue | red | blue | red | red | blue |
| 6. Sparse Matrix | orange | blue | blue | red | red | red | red | red | blue | blue | red | blue |
| 7. Spectral (FFT) | orange | blue | blue | blue | red | red | blue | blue | green | red | red | red |
| 8. Dynamic Prog | orange | blue | red | blue | red | red | blue | blue | blue | orange | blue | red |
| 9. Particle Methods | blue | orange | blue | red | blue | red | red | blue | blue | blue | blue | blue |
| 10. Backtrack/B&B | blue | blue | orange | red | red | red | blue | blue | blue | blue | orange | blue |
| 11. Graphical Models | blue | blue | orange | red | red | blue | blue | blue | blue | blue | red | blue |
| 12. Unstructured Grid | blue | blue | orange | orange | red | red | red | red | blue | blue | blue | blue |

W. Feng, wfeng@vt.edu, 540.231.1192
QCD at the Femtoscale in the Era of Big Data

VIRGINIA TECH.

SyNeRG
synergy.cs.vt.edu

# Molecular Dynamics → Cosmology → ???

- Primary computational dwarf? *N*-body method → particle method
- A computational dwarf (or pattern) describes a program's machinery, flow of resources, and outputs.

Algorithms
216x

Software
4x

**80,000x**

10.5
HOURS

88x
Hardware

0.47
SECONDS

Awesome for the domain scientist!
- Can run "what-if" simulations for rational drug design on a GPU server in his office, but

  80,000x relative to *serial*. What if "level playing field?"
  - 80,000 / 216 = 371 (algo. refactor)
  - 371 / 16-core CPU = 23 (1 → 16 cores)
  - 23 / 2 = 11 (DP → SP)
  - 11 / "2" = 5 (calc → lookup)

12 Ways to Fool the Masses
(David H. Bailey, NASA & LBL, 1991)

Debunking the 100X GPU vs. CPU Myth:
An Evaluation of Throughput Computing on CPU and GPU

...nance results, *not*
...up
...es for an inner

7. When direct run-time comparisons are required, compare with an old code on an obsolete system.
8. If Mflop/s rates must be quoted, base the operation count on the parallel [version],

**VIRGINIA TECH**™

# Race to Sequence the Human Genome

- **Theory & Experiment (Collins@NIH)**
  - Goal: Complete in 15 years
    1990 – 2005
  - Cost: $3,000M (1990-2000/2003)

- **Computing (Venter@Celera)**
  - Goal: Complete in 3 years & cheaper
    1998 - 2001
  - Cost: $300M (1998-2000/2003)



June 26, 2000

**mpiBLAST**

"String matching" →
dynamic programming
dwarf

# Pairwise Sequence Alignment (Smith-Waterman Algorithm)

- Performs local sequence alignment by identifying similar regions between two strings of nucleic acid sequences or protein sequences.

## Algorithm [edit]

Let $A = a_1 a_2 \ldots a_n$ and $B = b_1 b_2 \ldots b_m$ be the sequences to be aligned, where $n$ and $m$ are the lengths of $A$ and $B$ respectively.

1. Determine the substitution matrix and the gap penalty scheme.
   - $s(a, b)$ - Similarity score of the elements that constituted the two sequences
   - $W_k$ - The penalty of a gap that has length $k$

2. Construct a scoring matrix $H$ and initialize its first row and first column. The size of the scoring matrix is $(n + 1) * (m + 1)$. The matrix uses 0-based indexing.

$$H_{k0} = H_{0l} = 0 \quad for \quad 0 \le k \le n \quad and \quad 0 \le l \le m$$

3. Fill the scoring matrix using the equation below.

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ \max_{k \ge 1}\{H_{i-k,j} - W_k\}, \\ \max_{l \ge 1}\{H_{i,j-l} - W_l\}, \\ 0 \end{cases} \quad (1 \le i \le n, 1 \le j \le m)$$



Scoring method of the Smith–Waterman algorithm

$$H_{i-1,j-1} + s(a_i, b_j) \qquad \max_{k \ge 1}\{H_{i-k,j} - W_k\}$$
$$\max_{l \ge 1}\{H_{i,j-l} - W_l\} \quad H_{i,j}$$

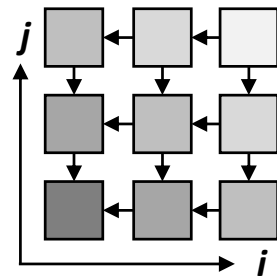**Fill the scoring matrix**

# Wavefront Loops

- Update each entry of a grid based on already-updated values from its neighbors
- Used in many scientific applications, e.g., PDE solver, sequence alignment tools, etc.
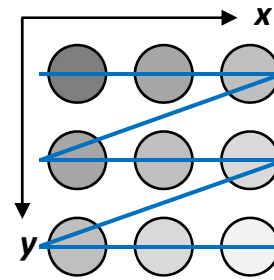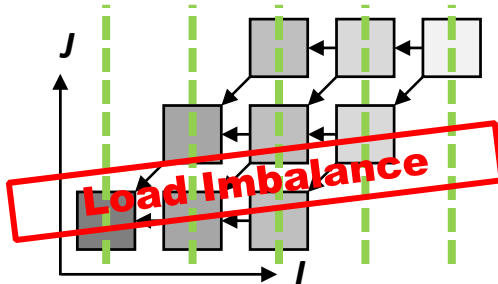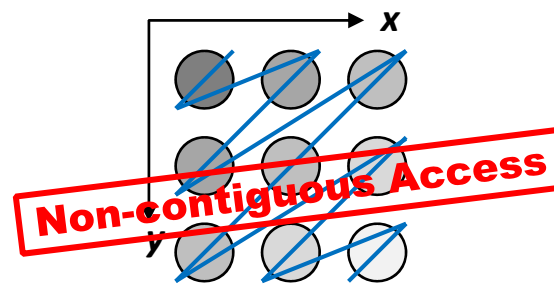
**Example: a wavefront loop (2D matrix)**

**Neither loop can be parallelized.**

```
for(int i = 0; i < m; i++)
  for(int j = 0; j < n; j++)
    A[i][j] = A[i][j-1] * 0.5 + A[i-1][j] * 0.5;
```

**Data Dependence (Iteration Space)**     **Memory Access (Memory Space** A[y][x]**)**

# Wavefront Loops

- Update each entry of a grid based on already-updated values from its neighbors
- Used in many scientific applications, e.g., PDE solver, sequence alignment tools, etc.

**Example: a wavefront loop (2D matrix) -- Tra**

> **J-loop can be parallelized.**

```
for(int I = 0; I < m+n-1; I++)
  for(int J = max(0, I-n+1); J < min(m, I+1); J++)
    A[J][I-J] = A[J][I-J-1] * 0.5 + A[J-1][I-J] * 0.5;
```
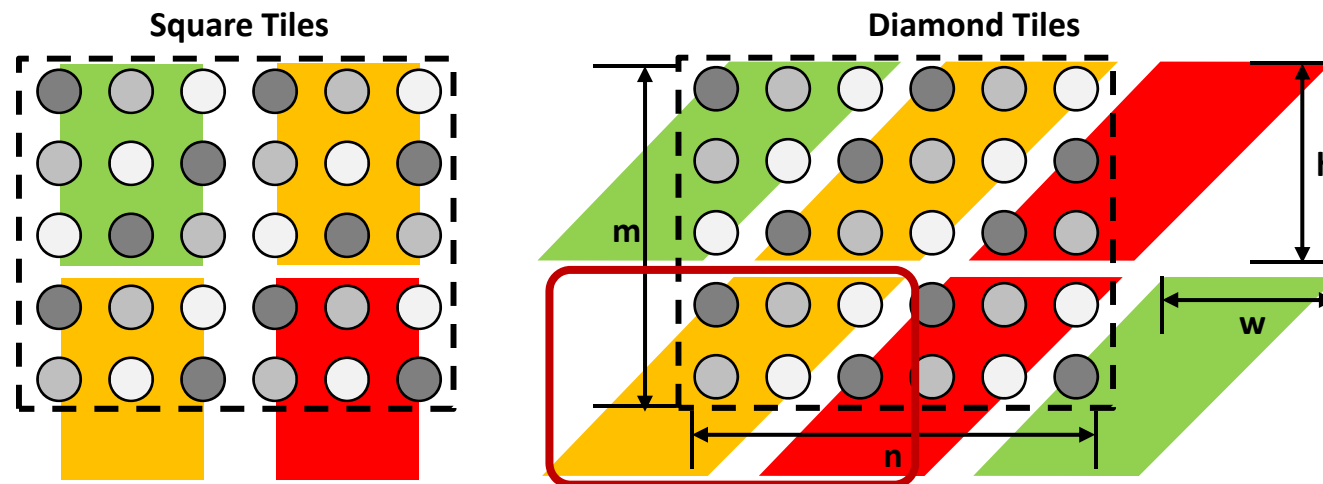
**Data Dependence (Iteration Space)**

**Memory Access (Memory Space** A[y][x]**)**



*Load Imbalance*

*Non-contiguous Access*

# Existing Parallel Solutions

- Tiling-based solutions and their limitations
  - Problem 1: Wasted memory and computing resources
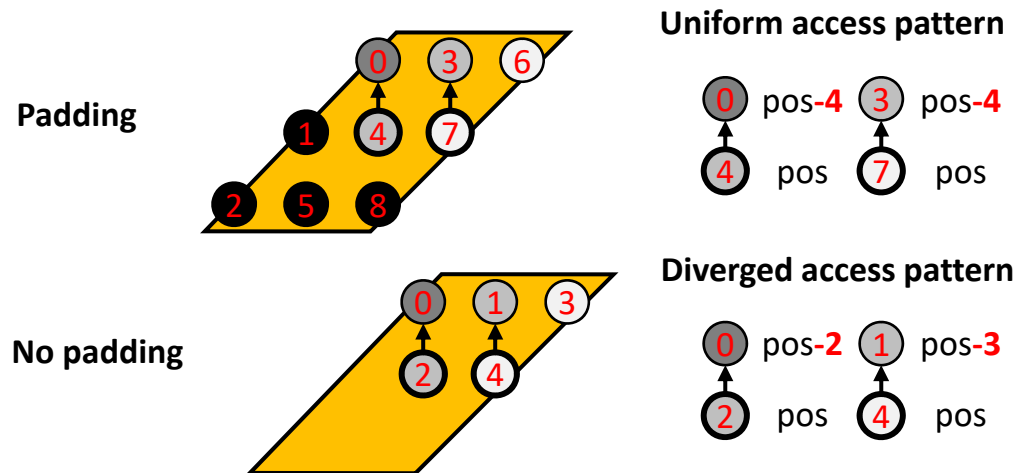
**Square Tiles**

**Diamond Tiles**

Tiles with same color can be executed in parallel

**Non-contiguous memory access still exists**

**Much memory space will be wasted (The rate of effective memory usage $\approx n/(n+h)$)**

# Existing Parallel Solutions

- Tiling-based solutions and their limitations
  - Problem 1: Wasted memory and computing resources



**Padding**

**No padding**

**Uniform access pattern**

**Diverged access pattern**

**Padding-free strategy may greatly increase the complexity of indexing and lead to more branches in GPU kernels**

# Existing Parallel Solutions

- Tiling-based solutions and their limitations
    - Problem 1: Wasted memory and computing resources
    - Problem 2: Layout transformation overhead
    - Problem 3: Task scheduling

**For some workloads, sufficient parallelism can be exposed**

**For other workloads, insufficient parallelism will be met**

**For some workloads, tiling-based solution may lose efficiency because of the small amount of tiles along anti-diagonals**
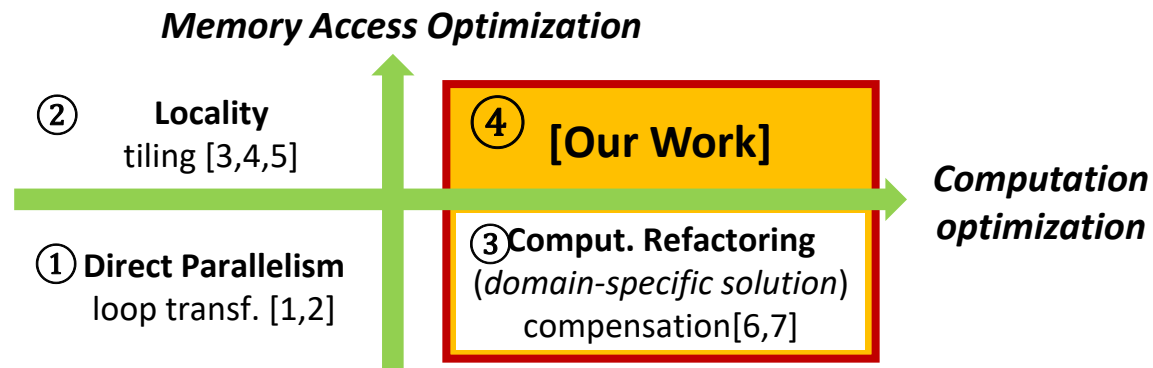
Tiles with same color can be executed in parallel

# Existing Parallel Solutions

- Compensation-based solutions and their limitations
  - Problem 1: Global synchronizations
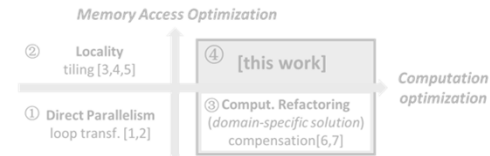  - Problem 2: Limited usage in sequence alignment algorithms

**❶** Compute partial results by ignoring horizontal dependency

**❷** Compensate the partial results

**❸** Combine the results from **❶** and **❷**



**Multiple expensive global synchronizations are required for processing each row; the compensation-based solution works well for string-matching operations**

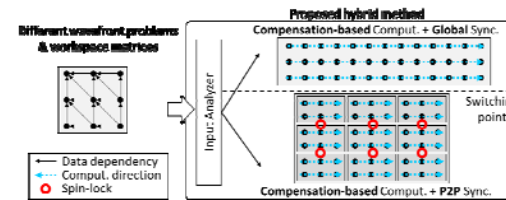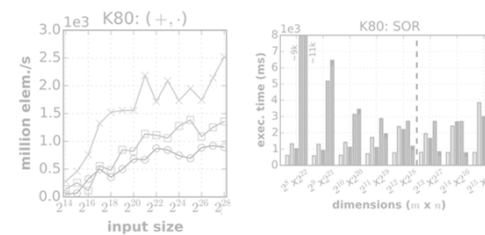# Our Highly *Efficient* Wavefront Parallelism

Diamond Tiles

**Wasted Memory**

**Sequential Scheduling**

**Sync. Overhead**

*Memory Access Optimization*

② **Locality**
tiling [3,4,5]

④ **[Our Work]**

*Computation optimization*

① **Direct Parallelism**
loop transf. [1,2]

③ **Comput. Refactoring**
(*domain-specific solution*)
compensation[6,7]

# Outline

- Introduction

- Motivation



- Our Method
  - Compensation-based Method
  - GPU Implementation
  - Hybrid Parallel Strategy



- Evaluation
  - Weighted-scan Kernel Performance
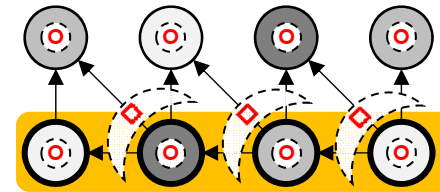  - Wavefront Kernel Performance

# Compensation-based Method

- ## Wavefront Pattern

$$A_{i,j} = (A_{i,j-1} \circ b_0) \diamond (A_{i-1,j} \circ b_1) \diamond (A_{i-1,j-1} \circ b_2)$$

- ○ *generic distribution operator (for adding weights)*
- ◇ *generic accumulation operator (for adding neighbors)*

- ## Compensation-based Method

**Step 1:** $\tilde{A}_{i,j} = (A_{i-1,j} \circ b_1) \diamond (A_{i-1,j-1} \circ b_2)$

**Step 2:** $B_{i,j} = \begin{cases} \sum_{u=0}^{j-1}(\tilde{A}_{i,u} \circ \prod_{v=u}^{j-1} b_0) & \text{when } \circ \neq \diamond \\ \sum_{u=0}^{j-1}(\tilde{A}_{i,u} \diamond b_0) & \text{when } \circ = \diamond \end{cases}$

**Step 3:** $A_{i,j} = \tilde{A}_{i,j} \diamond B_{i,j}$

This is valid when (1) ○ has the distributive property over ◇ ; (2) ○ is same with ◇. *
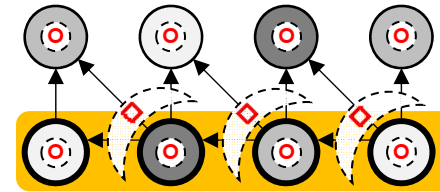
* The mathematical proof is included in our paper.
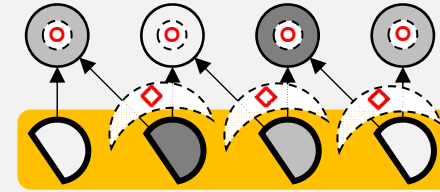
# Compensation-based Method

- ## Wavefront Pattern

$$A_{i,j} = (A_{i,j-1} \circ b_0) \diamond (A_{i-1,j} \circ b_1) \diamond (A_{i-1,j-1} \circ b_2)$$

- ○ *generic distribution operator (for adding weights)*
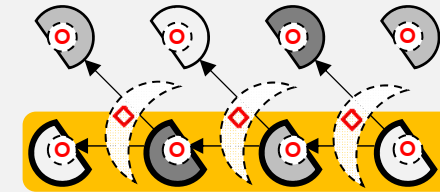- ◇ *generic accumulation operator (for adding neighbors)*

- 

**Step 1:** $\tilde{A}_{i,j} = (A_{i-1,j} \circ b_1) \diamond (A_{i-1,j-1} \circ b_2)$

**Step 2:** $B_{i,j} = \begin{cases} \sum_{u=0}^{j-1}(\tilde{A}_{i,u} \circ \prod_{v=u}^{j-1} b_0) & \text{when } \circ \neq \diamond \\ \sum_{u=0}^{j-1}(\tilde{A}_{i,u} \diamond b_0) & \text{when } \circ = \diamond \end{cases}$

**Step 3:** $A_{i,j} = \tilde{A}_{i,j} \diamond B_{i,j}$

# Compensation-based Method

- Wavefront loops can be expressed as compensation-based parallelism patterns
- SOR (*Successive Over-Relaxation*) Solver:

  ( ◇ , ○ ) maps to (+, ·)

  ```
  A[i][j] = (A[i][j] + A[i][j-1] + A[i-1][j] +
             A[i+1][j] + A[i][j+1]) / 5;
  ```

- SW (Smith-Waterman):

  ( ◇ , ○ ) maps to (max, +)

  ```
  A[i][j] = max(A[i][j-1] - 2, A[i-1][j] - 2,
               A[i-1][j-1] + s(i,j), 0);
  ```
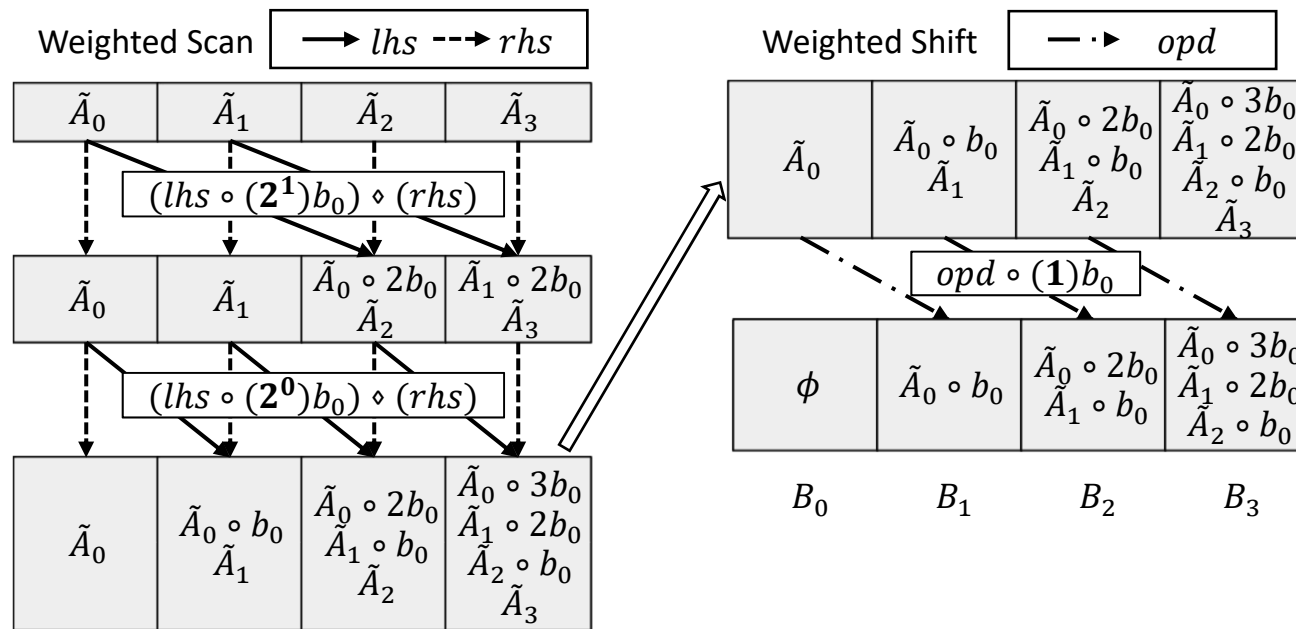
- SAT (Summed-Area Table):

  ( ◇ , ○ ) maps to (+, +)

  ```
  A[i][j] = p[i][j] + A[i][j-1] + A[i-1][j] - A[i-1][j-1];
  ```

# GPU Implementation

- Step 2 of the compensation-based method is the critical part: "***Weighted Scan***"*


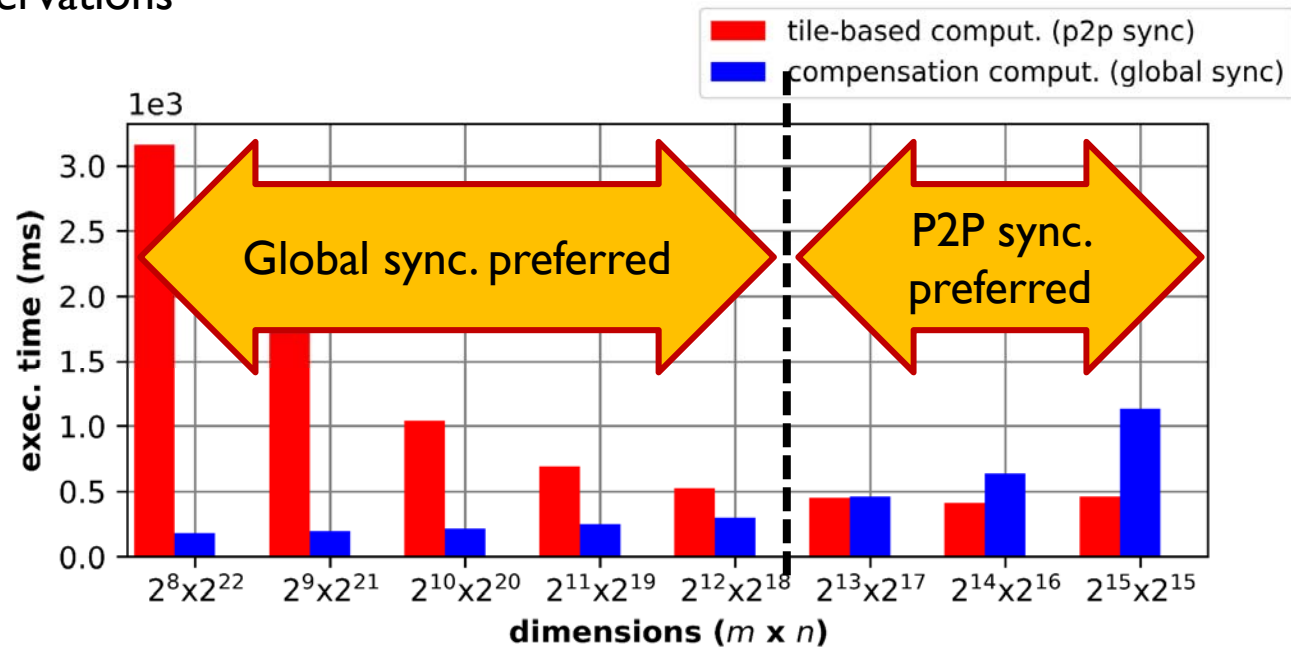
* which also includes a weighted shift operation

# GPU Implementation

- Step 2 of the compensation-based method is the critical part: "**Weighted Scan**"
- Our algorithm handles the changing weights during each stages of the operations
- A hierarchical design is used for GPUs
    - *Register level*: compute how the preceding neighbor affects the current one via <u>data shuffle instructions</u>
    - *Shared memory level*: compute how the preceding **"warp"**\* of neighbors affect the current one via <u>shared memory access</u>
    - *Global memory level*: compute how the preceding **"block"**\* of neighbors affect the current one via <u>global memory access</u>

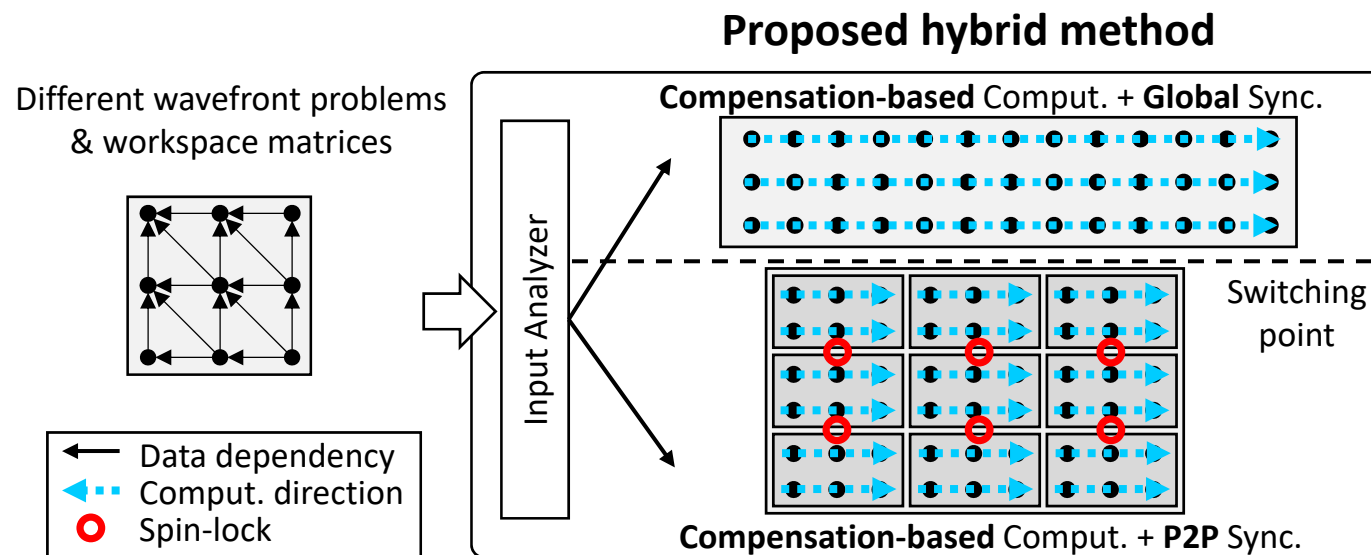\* which are thread organization units in NVIDIA GPU terminology

# Hybrid Parallel Strategy

- *Is the compensation-based method sufficient for any types of workloads?*
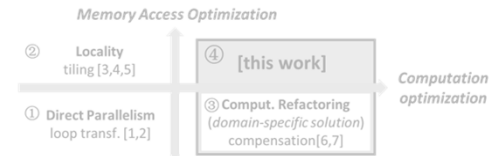- Observations

# Hybrid Parallel Strategy

- Our hybrid design switches to the appropriate parallel method, based on the input workload
- All the computation follows the compensation-based parallelism pattern

**Proposed hybrid method**



Different wavefront problems & workspace matrices

**Compensation-based** Comput. + **Global** Sync.

Switching point

**Compensation-based** Comput. + **P2P** Sync.

Input Analyzer

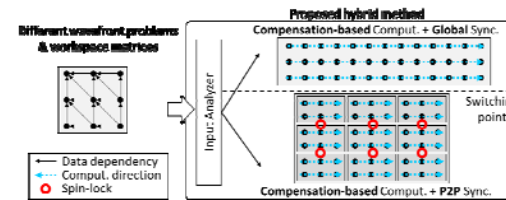Data dependency
Comput. direction
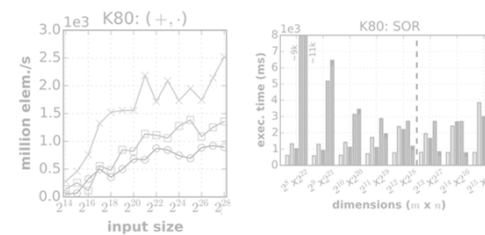Spin-lock

# Outline

- Introduction

- Motivation

- Our Method
  - Compensation-based Method
  - GPU Implementation
  - Hybrid Parallel Strategy

- Evaluation
  - Weighted-scan Kernel Performance
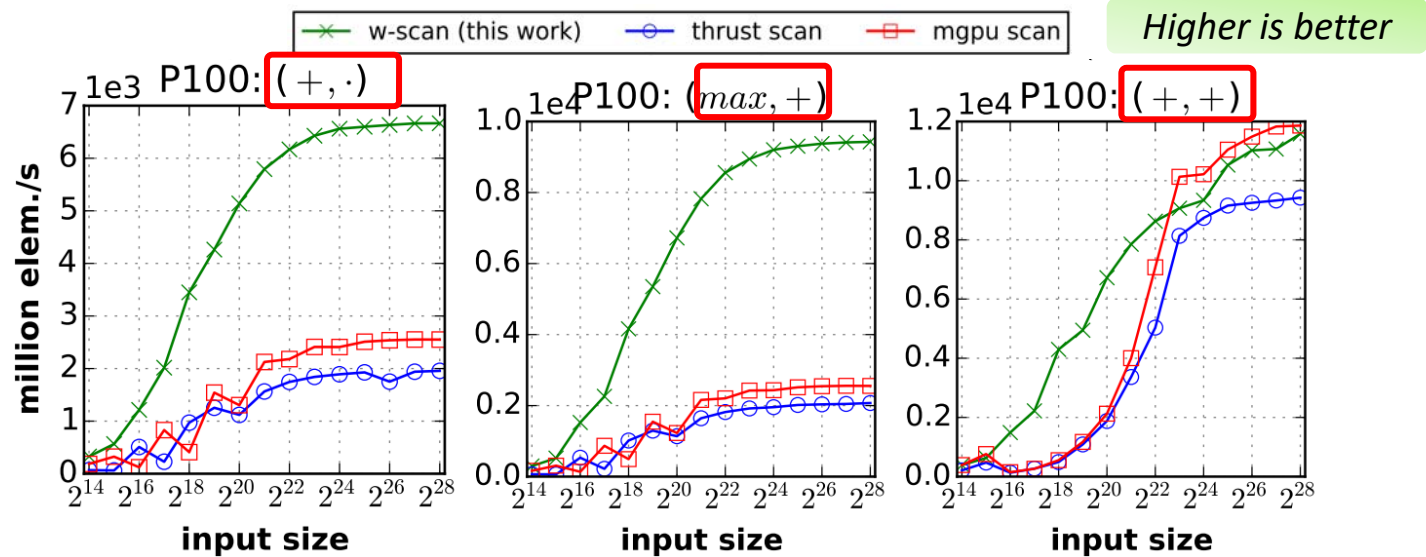  - Wavefront Kernel Performance

# Experimental Platforms

- **<u>nVidia Tesla K80 (Kepler-K80)</u>**, 2496 CUDA cores @ 824 MHz, 240 GB/s bandwidth

- **<u>nVidia Pascal P100 (Pascal-P100)</u>**, 3584 CUDA cores @ 405 MHz, 720 GB/s bandwidth *

- Our **<u>Weighted Scan</u>** vs. other tools
    a. Thrust v.1.8.1 (thrust::exclusive_scan w/ custom comparator)
    b. ModernGPU v.2.0 (mgpu::scan w/ custom comparator)
        - Using 1D array of data to mimic different rows

- Our **<u>Hybrid Wavefront</u>** kernel vs.
    a. Tile-based methods [`15] (incl. square & diamond tiles)
    b. Compensation-based methods [`12, `16, this work]
        - Using 2D array of data to mimic different workloads

\* We only show the performance results of P100 GPU here.
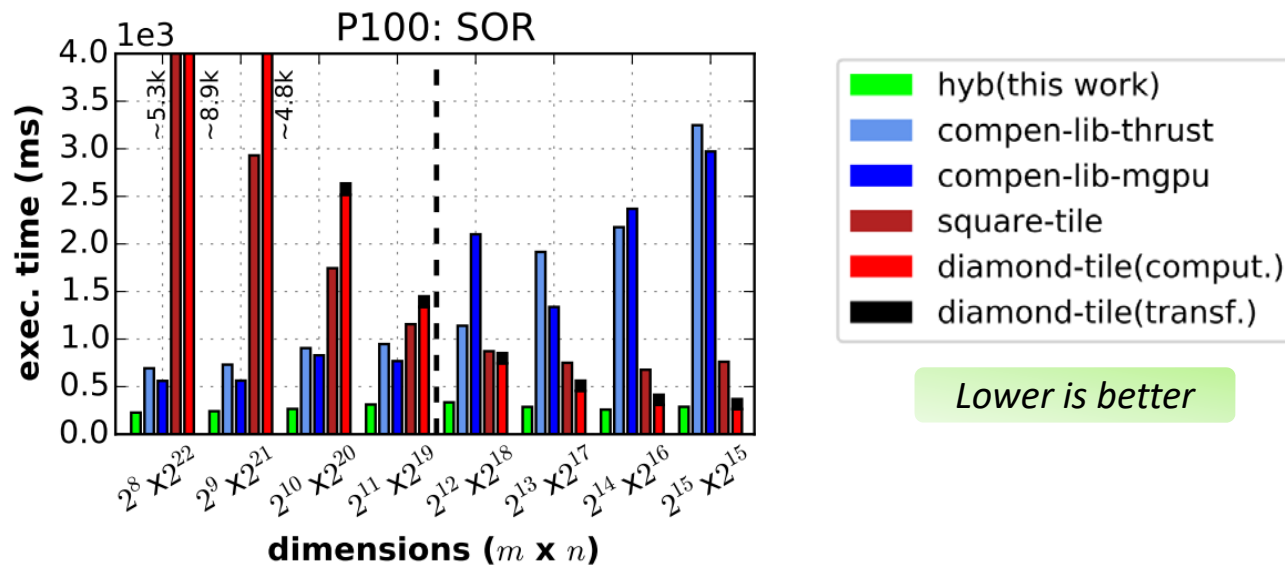
# Weighted-scan Kernel Performance

- Processing a row of data with variable sizes



- For $\circ \neq \bullet$, our method delivers significant performance benefit (*mainly because we can calculate the distance-related weights more **efficiently** in the kernel*)
- For $\circ = \diamond$, our method reduces to an ordinary scan kernel
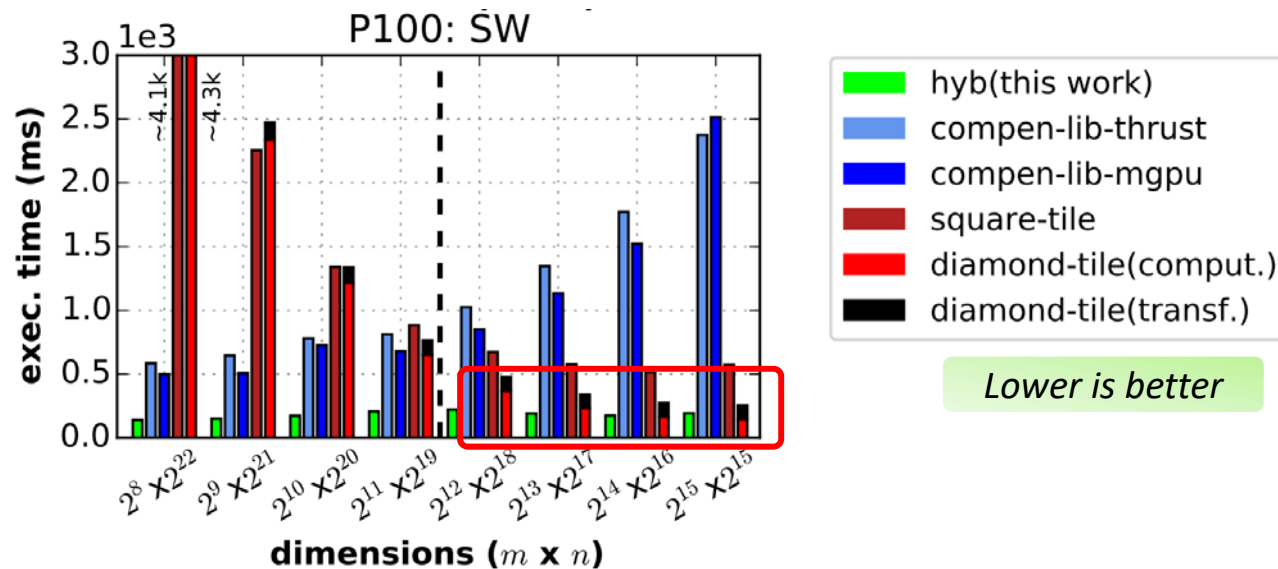
# Wavefront Kernel Performance

- Using SOR, SW, and SAT as representative wavefront kernels
- Processing 2D matrices of data with variable dimensions



*Lower is better*

- Our method always delivers better performance than previous solutions
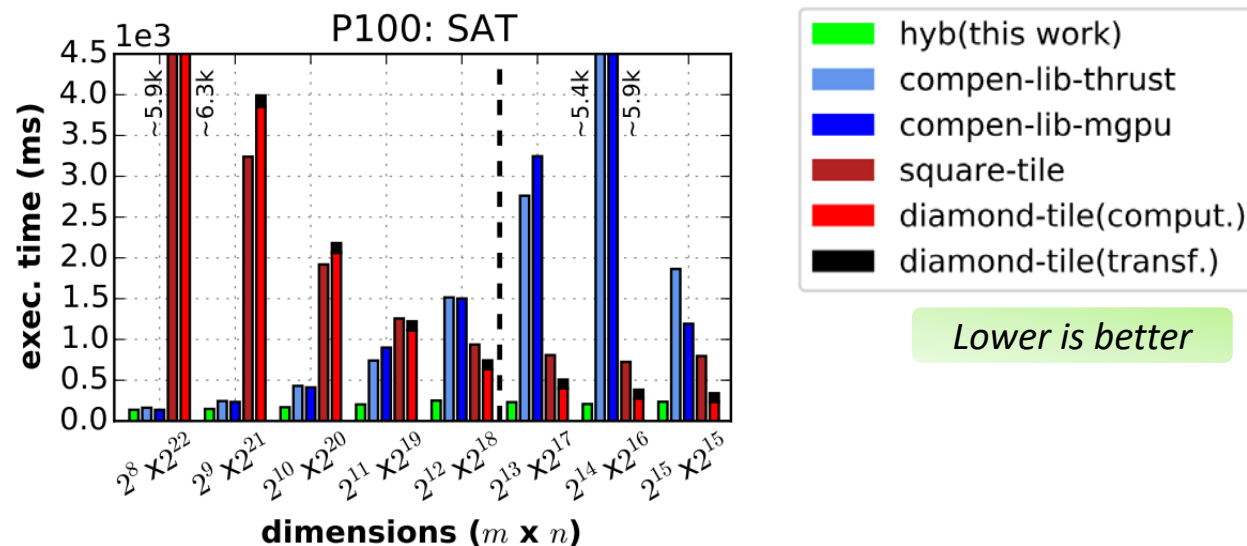
# Wavefront Kernel Performance

- Using SOR, SW, and SAT as representative wavefront kernels
- Processing 2D matrices of data with variable dimensions



- The transformation overhead becomes non-negligible for the diamond-tile method
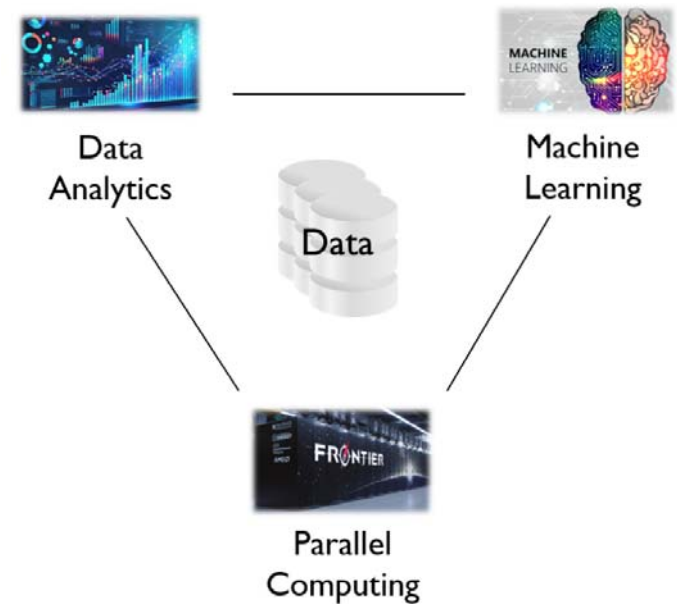
# Wavefront Kernel Performance

- Using SOR, SW, and SAT as representative wavefront kernels
- Processing 2D matrices of data with variable dimensions



- – Our hybrid method exhibits superior performance *regardless* of the workloads and wavefront types

# At the Synergistic Intersection of Parallel Computing, Data Analytics, and Machine Learning

**Wu Feng**

wfeng@vt.edu



Data Analytics

Machine Learning

Data

Parallel Computing

W. Feng, wfeng@vt.edu, 540.231.1192
QCD at the Femtoscale in the Era of Big Data

# Challenge

- The rate of growth in **big data** is *far outstripping* the rate at which computing can (brute-force) **compute** on the data.

# Approach

- Synergistic co-design of architecture, software, and in particular, algorithms to more *efficiently* and *intelligently* compute on the data.

# At the Synergistic Intersection of Parallel Computing, Data Analytics, and Machine Learning

- Systems-Oriented
  - Automated GPU Blocksize Tuning via Iterative Machine Learning (Cui)
  - Scalable I/O for Deep Learning (Pumma)

- Applications-Oriented
  - SparkLeBLAST: High-Productivity DNA Sequence Search (Youssef)
  - Visual Data Analytics (Dash, in collaboration with C. North CS@VT)
  - Data-Oriented Computational Fluid Dynamics (Cui)
  - Understanding Carcinogenesis (Dash, in collaboration with VCOM)
  - Graph Analytics (Wanye, in collaboration with MIT LL)
  - Biomedical Imaging (Goel et al., in collaboration with BEAM@VT)

| Sequence Alignment | Molecular Dynamics |
|---|---|
| Earthquake Modeling | Neuro-informatics |
| CFD for Mini-Drones | Cyber-Security |

# Automated GPU Blocksize Tuning via Iterative ML

- **Problem**

  *Many parameters to tune to achieve best performance*
  - ✓ Thread block size
  - ✓ # streams
  - ✓ Register usage
  - ✓ Compiler optimization flags
  - ✓ *… and so on*

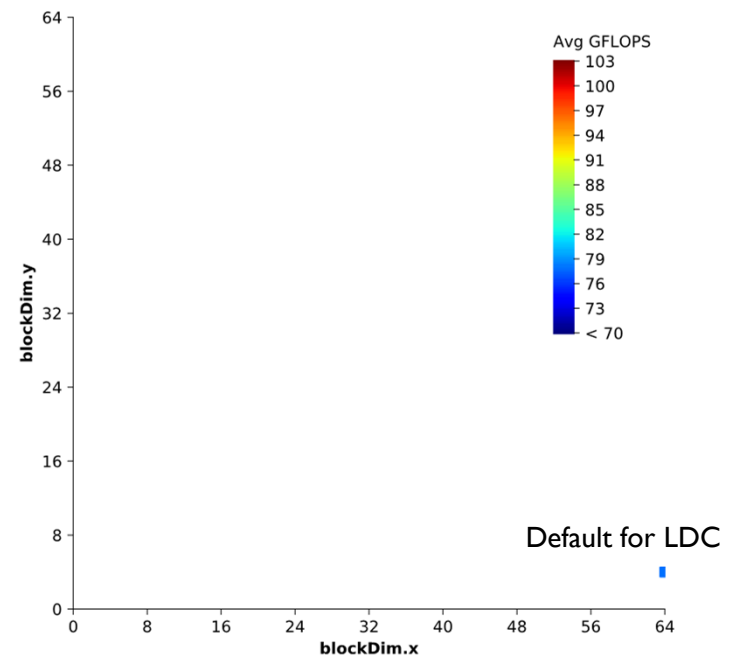  **O(millions) potential software configurations**
  **for the *same code***

- Our Focus
  - ✓ ***Thread block size***

- Example
  - ✓ Lid-driven cavity (LDC) code with
    varying GPU thread block size (NVIDIA K20m GPU)
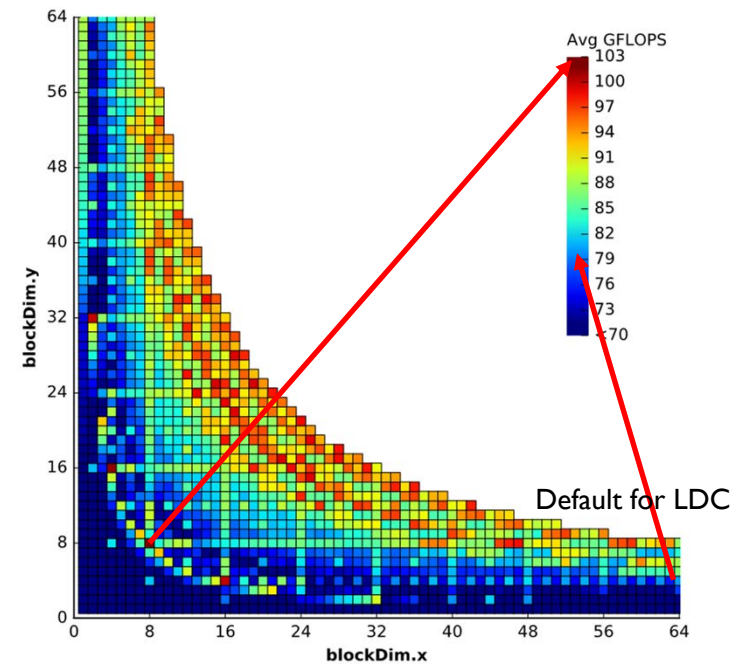
$$1 <= \{blockDim.x, \ blockDim.y\} <= 1024$$
$$1 <= blockDim.x * blockDim.y <= 1024$$

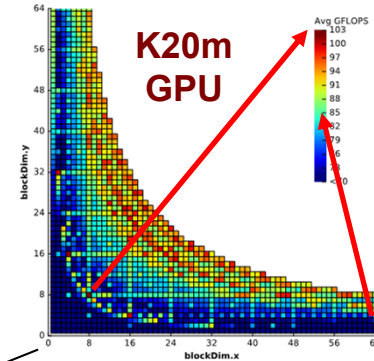# Automated GPU Blocksize Tuning via Iterative ML

- ***Challenge: Huge search space***
  (even when considering only **one** parameter,
  i.e., thread block size)

- ***What should I set my thread block size to?***
  - Brute-force search (7262 runs)
    - Takes more than a day to search
  - Reliance on developer experience?
    - Recommended block size
      - 8x8, 8x16, 16x8, 16x16
      (see next slide)

$$1 <= \{blockDim.x, blockDim.y\} <= 1024$$
$$1 <= blockDim.x * blockDim.y <= 1024$$



Default for LDC
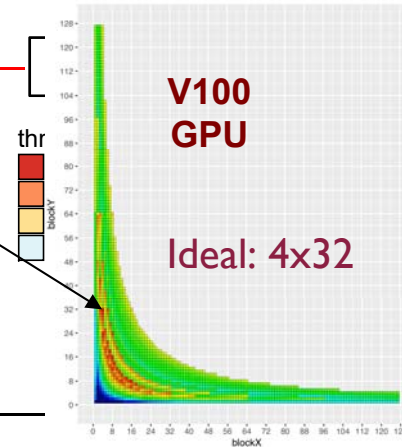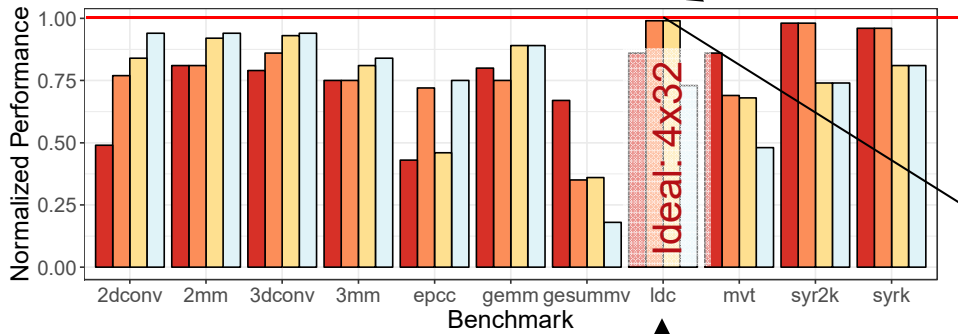
# Automated GPU Blocksize Tuning via Iterative ML



- Parameter tuning the GPU
  - Reliance on end-user experience (or intuition)
    - Typically chosen block sizes?  64, 128, and 256
  - Best parameter setup varies …
    - Between applications
    - Between different devices

**K20m GPU**

Ideal: 8x8

**V100 GPU**

Ideal: 4x32

Ideal: 4x32

LDC code from upper right
and from previous slide

# Automated GPU Blocksize Tuning via Iterative ML

- **Parameter tuning the GPU**
    - Reliance on end-user experience (or intuition)
    - **Statistical methods**
        - Build a model *a priori* based on a (required) **large training set**
        - Predict the best parameter(s) based on real-time profiling data and model
        - May perform poorly for **new** algorithms on **new** devices or systems
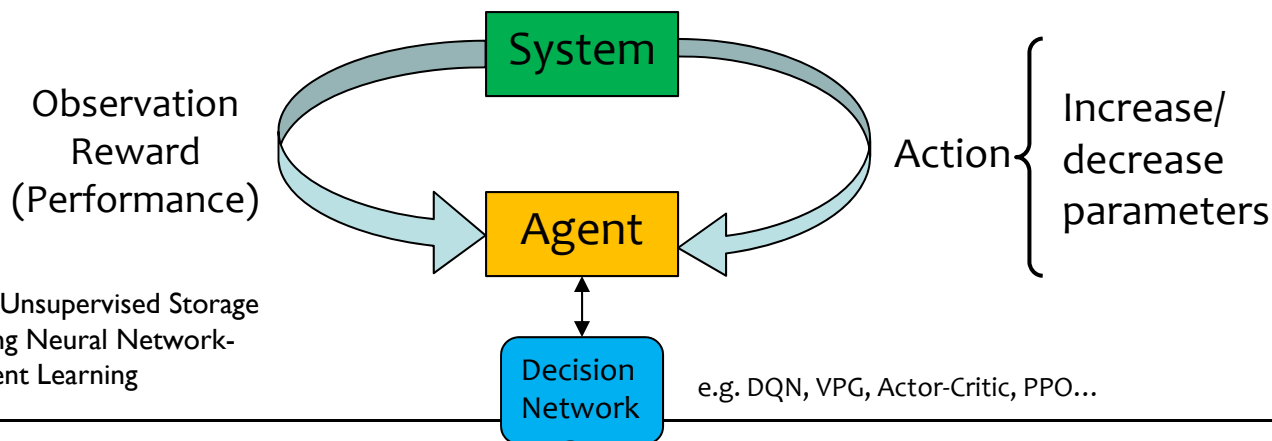
# Automated GPU Blocksize Tuning via Iterative ML

- Parameter tuning the GPU
  - Reliance on end-user experience (or intuition)
  - Statistical methods
  - Reinforcement learning methods
    - Requires **no prior knowledge** of the target system
    - Run continuously to adapt and dynamically update parameters
    - May take the decision system significant time to converge



Observation
Reward
(Performance)

System

Action {
Increase/
decrease
parameters

Agent

Example: SC17 : CAPES: Unsupervised Storage
Performance Tuning Using Neural Network-
Based Deep Reinforcement Learning

Decision
Network

e.g. DQN, VPG, Actor-Critic, PPO…

# Automated GPU Blocksize Tuning via Iterative ML

- **Our Challenge:  How to deploy**

  … new applications

  … new algorithms

  … new systems

  … new hardware accelerators (e.g., GPUs, FPGAs, etc.)

  ***and*** tune the multi-dimensional search space of hardware/software/algorithmic parameters to optimize applications

X. Cui and W. Feng, "Iterative Machine Learning (IterML) for Effective Parameter Pruning and Tuning in Accelerators," *16th ACM International Conference on Computing Frontiers*, April-May 2019.

- **Our Goal**

  - Intelligently tune parameters with *no prior knowledge* and *no pre-trained model*
    - Deliver near-optimal performance with the least amount of effort and domain knowledge.

# Automated GPU Blocksize Tuning via Iterative ML

- Iterative Machine Learning (IterML)
  - Uses samples from one iteration to then look for potentially *better samples* in subsequent iterations.

Pick ratio:
- sample ratio in each iteration

Cut ratio:
- ratio of the space pruned each iteration



X. Cui and W. Feng, "Iterative Machine Learning (IterML) for Effective Parameter Pruning and Tuning in Accelerators," *16th ACM International Conference on Computing Frontiers*, April-May 2019.
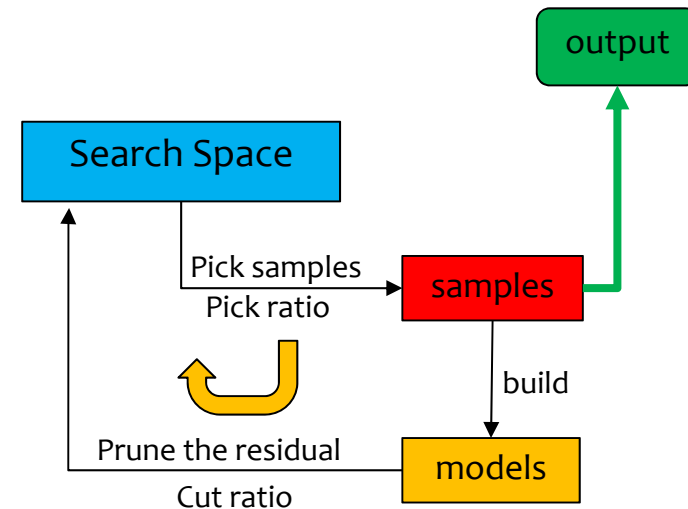
# Automated GPU Blocksize Tuning via Iterative ML

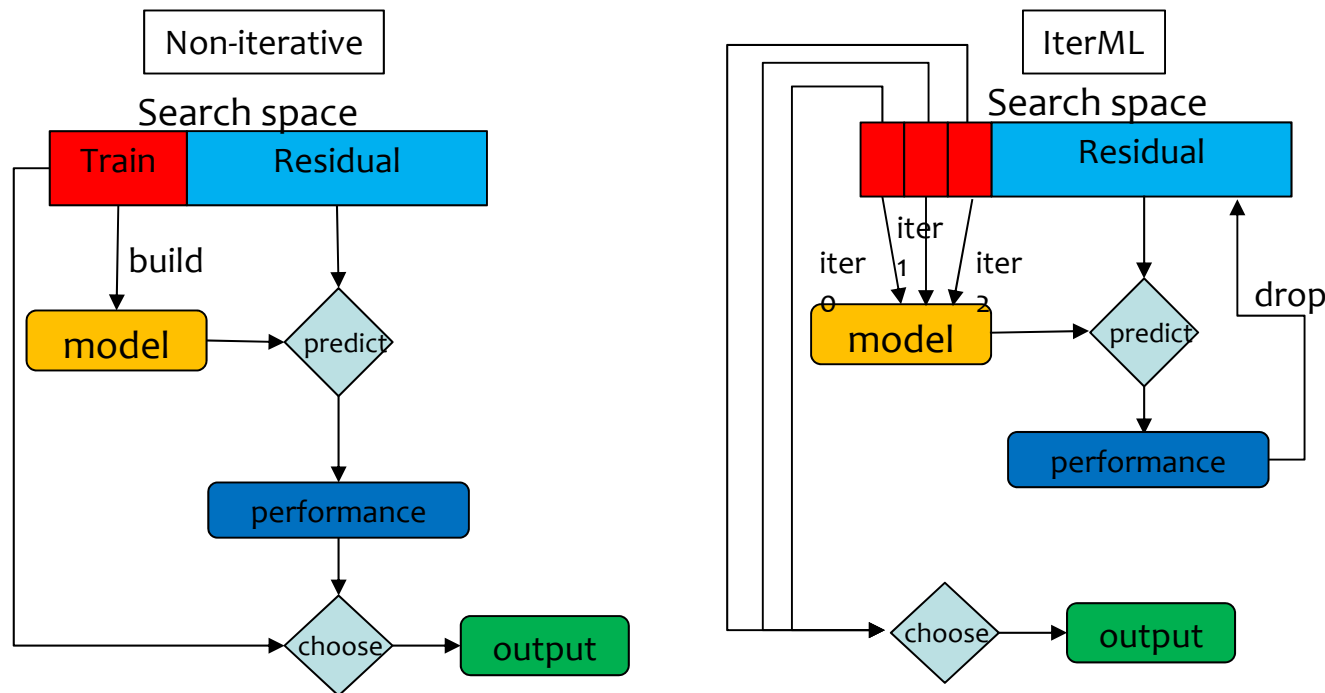- Non-Iterative vs. Iterative Machine Learning (IterML)
  - With no prior knowledge or pretrained model



\* X. Cui and W. Feng, "Iterative Machine Learning (IterML) for Effective Parameter Pruning and Tuning in Accelerators," *16th ACM Int'l Conf. on Computing Frontiers*, April-May 2019.
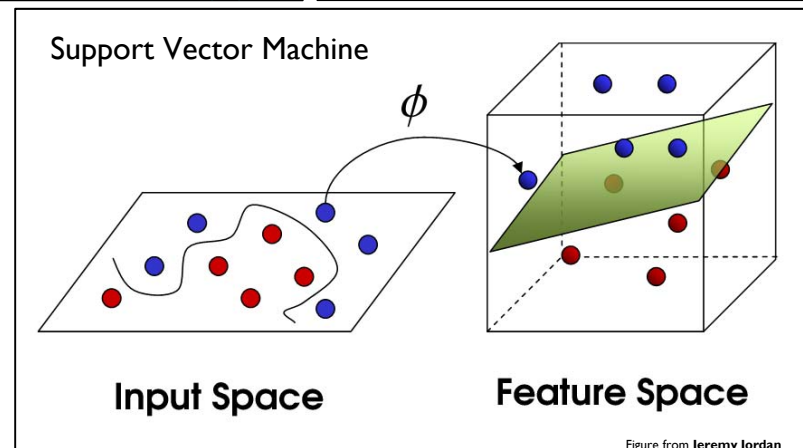
Machine Learning Models

Classification & Regression Trees

yes    sex male?    no

Age>9.5    survived

died    sibsp>2.5

died    survived

Random Forest

Figure from machinelearning-blog.com

K-Nearest Neighbors

Figure from MathWorks

Input Layer    Hidden Layer    Output Layer

Input 1
Input 2
Input 3
Input 4

Output

Multilayer Perceptron

Support Vector Machine

$\phi$

Input Space    Feature Space

Figure from Jeremy Jordan

# Takeaway:
# Automated GPU Blocksize Tuning via Iterative ML

- Iterative machine-learning (IterML) approach
  … to prune the massive parameter search space

  - Performance evaluation of traditional non-iterative ML vs. our IterML

  - Empirical demonstration that IterML with the random forest (RF) model reduces search effort by 40%~80%

  - Random forest (RF) produces better and more stable results than other popular ML models

  X. Cui and W. Feng, "Iterative Machine Learning (IterML) for Effective Parameter Pruning and Tuning in Accelerators," *16th ACM International Conference on Computing Frontiers*, April-May 2019.
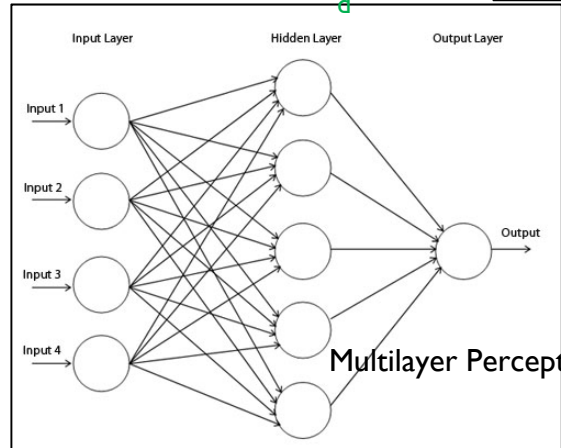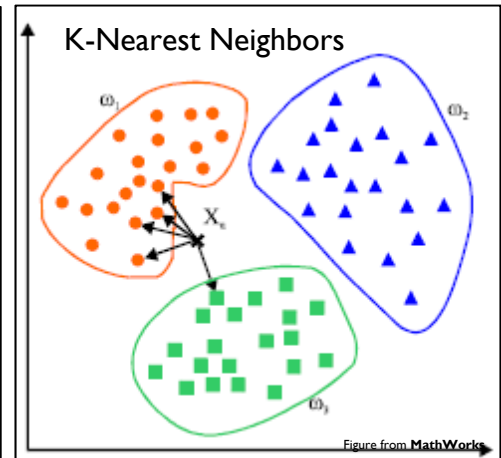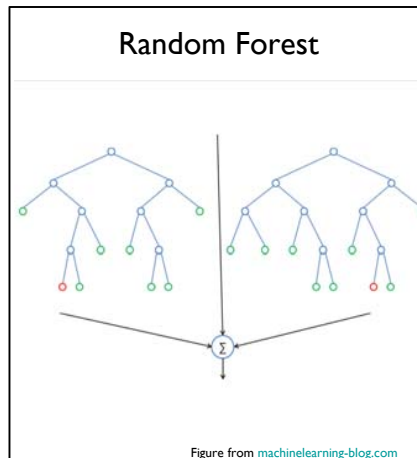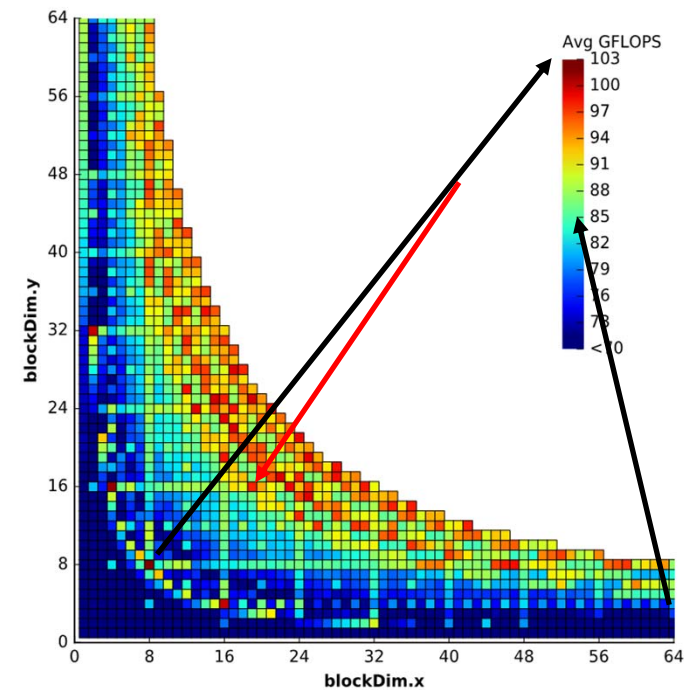
# Challenges: Scalable I/O in Large-Scale Deep Learning

***Not all large-scale deep learning is compute-bound…***

File I/O-bound

Large per batch data size & DNN with non-dense weight layers

## Semi-supervised Bounding Box Prediction[1]
Extreme weather detection
Input size: 768 x 768 x 16
By NERSC, LBNL, DOE, Stanford, Intel

[1] Kurth et al. Deep learning at 15pf: Supervised and semi-supervised classification for scientific data. SC'17.

## Image Classification
Lung cancer detection
Input size: 253x235x240

From Kaggle Data Science Bowl 2017

Compute-bound

High-resolution input data

Network I/O-bound

Large number of trainable parameters

## Unsupervised Image Feature Extraction[2]
LLNL's DNN with **15 billion** parameters

[2] Ni et al. Large-scale deep learning on the YFCC100M dataset. arXiv preprint arXiv:1502.03409, 2015.

# I/O Scaling of Deep Learning

Overall Training Time

660x worse than ideal

— Caffe/LMDB    ···· Ideal



Training Time Breakdown

- Read time
- Transform time
- Total forward time
- Total backward time
- I/O skew time
- Param sync time
- Param calculation time
- Param update time

# Scalable I/O in Deep Learning: Parallel Data Reading

- Lightning Memory-Mapped Database (LMDB)
  - Widely used in deep-learning frameworks, e.g., Caffe (default), Caffe2, TensorFlow, Keras-TensorFlow
  - Uses mmap internally (memory-mapped file I/O)
  - Database layout: B+ tree
- No collaboration between readers
  - Each reader opens the LMDB database in its virtual memory space
  - In each iteration, each reader reads B/NP samples of data via LMDB's API in a strided manner (B = batch size, NP = number of processes)



*Example of 3 processes performing strided data access*

# Scalable I/O in Deep Learning

*Our Solution: LMDB-IO* (Lightning Memory-Mapped Database – I/O)

An optimized I/O subsystem for large-scale deep learning
LMDB-IO optimizations are divided into <span style="color:red">three classes</span>

## Intra-node I/O

**Problem**

High inter-process contention in multi-reader environment due to **mmap**

**Solution**

- Localizing mmap (using one reader per node)
- Using MPI-3 shared buffer to share data between processes

## Speculative Distributed I/O (Inter-node)

**Problem**

High I/O skew due to indeterministic DB layout
(direct DB access is not allowed)

**Solution**

- Collaborating between reader processes to reduce I/O skew
- Speculatively reading reading data in parallel

## Direct I/O

**Problem**

**mmap** is highly inefficient.
The user has no control over I/O operations

**Solution**

- Replacing mmap with Posix I/O (direct I/O)
- Using several techniques to improve direct I/O performance

CIFAR10-Large & AlexNet Scaling

Total Execution Time
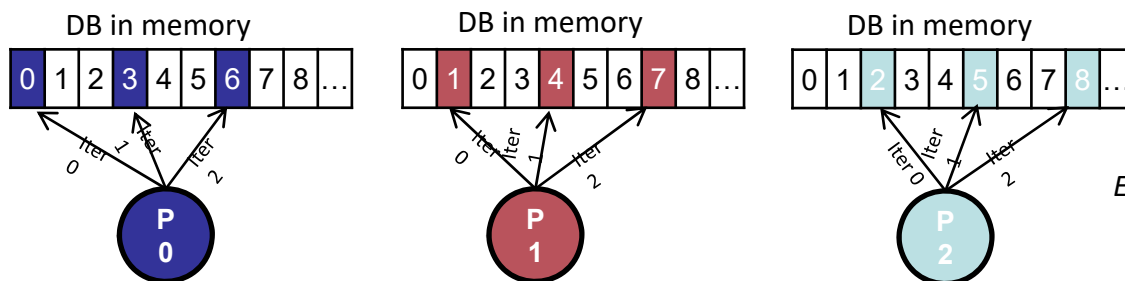
**Dataset:** CIFAR10-Large (50M 3KB images, 10 classes, 190 GB)
**DNN:** AlexNet (13 layers, 89K parameters)
**Batch size:** 18,432 **Training iterations:** 512
**Framework:** Caffe **Testbed/Storage:** LCRC Bebop/GPFS
(Each node: 36 cores Intel Broadwell, 128 GB memory)

— LMDB
— LMDBIO-LMM
— LMDBIO-LMM-DIO

Factor of Improvement over LMDB

LMDBIO-LMM
LMDBIO-LMM-DIO
LMDBIO-LMM-DM
LMDBIO-LMM-DIO-PROV
LMDBIO-LMM-DIO-PROV-COAL
LMDBIO-LMM-DIO-PROV-COAL-STAG

# ImageNet-Large & CaffeNet Scaling

**Total Execution Time**

| | | | | | |
|---|---|---|---|---|---|
| Time (s): 10000, 1000, 100 | | | | | |

**Dataset:** ImageNet-Large (6M 192KB images, 1000 classes, 1.1 TB)
**DNN:** CaffeNet (22 layers, 60M parameters)
**Batch size:** 18,432   **Training iterations:** 32
**Framework:** Caffe   **Testbed/Storage:** LCRC Bebop/GPFS
(Each node: 36 cores Intel Broadwell, 128 GB memory)

LMDB

LMDBIO-LMM

X-axis: 576   1152   2304   4608   9216

**Factor of Improvement over LMDB**

- LMDBIO-LMM
- LMDBIO-LMM-DIO

Y-axis: Factor of Improvement over LMDB — 5.0, 4.0, 3.0, 2.0, 1.0, -

576: 3.7
1152: 2.5
2304: 2.2
4608: 1.5
9216: 1.6

Number of Cores

W. Feng, wfeng@vt.edu, 540.231.1192
QCD at the Femtoscale in the Era of Big Data

VIRGINIA TECH

SyNeRG
synergy.cs.vt.edu

# Conclusion

- Synergistic co-design of algorithms, software, and hardware can massively accelerate discovery.

# What's Next?

- Case Studies on Synergistic Co-Design of Algorithms, Software, and Hardware
    - Brain Tomography on GPU.  Carcinogenesis: Weighted Set Cover vs. Graph Cluster. […]
- HPC Systems
    - IterML: Iterative Machine Learning (AFOSR & DOD)
        - Context:  Computational fluid dynamics (CFD) →
            OpenDwarfs, i.e., fundamental "DNA" building blocks for scientific computing
    - CoreTSAR: <u>Core</u> <u>T</u>ask-<u>S</u>ize <u>A</u>dapting <u>R</u>untime <u>S</u>ystem (DOE & NSF)
        - Context:  Initially, discrete CPU+GPU systems w/ discrete memory
            Now, also "fused" co-located CPU+GPU systems w/ shared memory
            See Aurora @ ANL with PVC & El Capitan @ LLNL with MI-300a
    - Scalable Deep Learning (with ANL → Meta & Llama-3)
        - Context:  Caffe, Caffe2, and Tensorflow
        - Takeaway:  Large-scale multi-node DL does *NOT* scale.

# An (Intra-Node) Ecosystem for Heterogeneous Parallel Computing



Sequence Alignment | Molecular Dynamics | Earthquake Modeling | Neuro-informatics | CFD for Mini-Drones | Biomedical Imaging | Cybersecurity

Software Ecosystem

**Design & Compile Time**

Computational & Communication Patterns: 13 Dwarfs

Source-to-Source Translation & Optimization Framework

Architecture-Aware Optimizations

MACHINE LEARNING

**Run Time**

Affinity Cost Models

Task Scheduling System

Performance & Power Models

… inter-node with applications to BIG DATA (StreamMR)

Computing Platforms

VIRGINIA TECH™

SyNeRG
synergy.cs.vt.edu