

Genomic Analysis at Scale: Mapping Irregular Computations to Advanced Architectures

Kathy Yelick

Robert S. Pepper Distinguished Professor of EECS

Associate Dean for Research

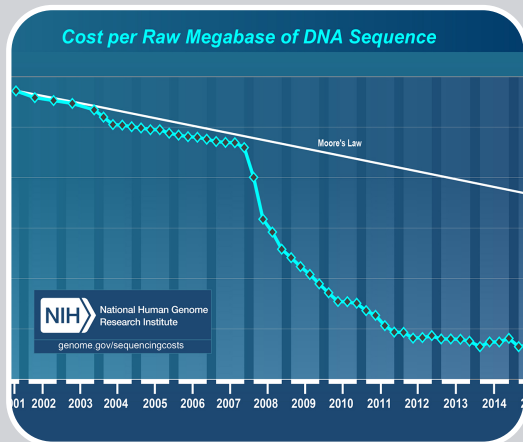
Division of Computing, Data Science, and Society

UC Berkeley

Senior Advisor on Computing

Lawrence Berkeley National Laboratory

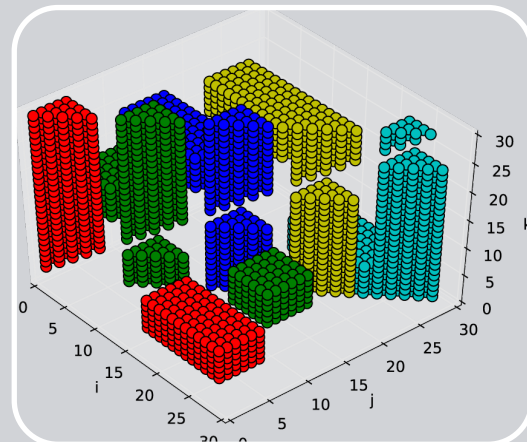
Genomic Analysis at Scale



Big Data

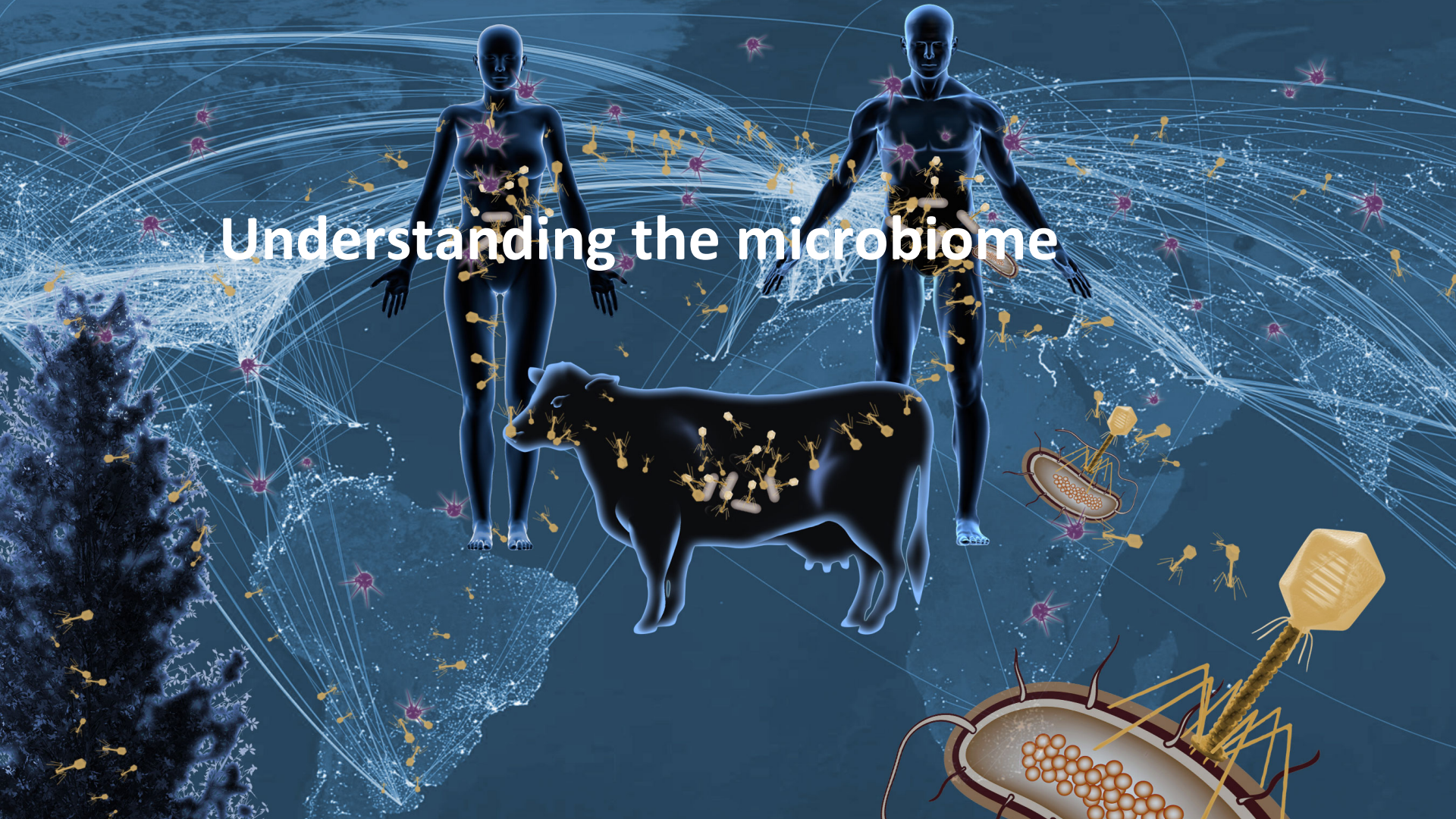


Big
Machines



Scalable
Algorithms

Understanding the microbiome

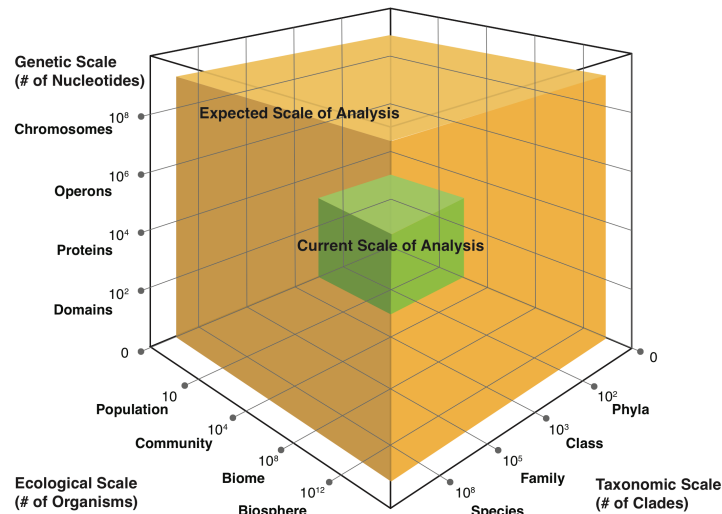
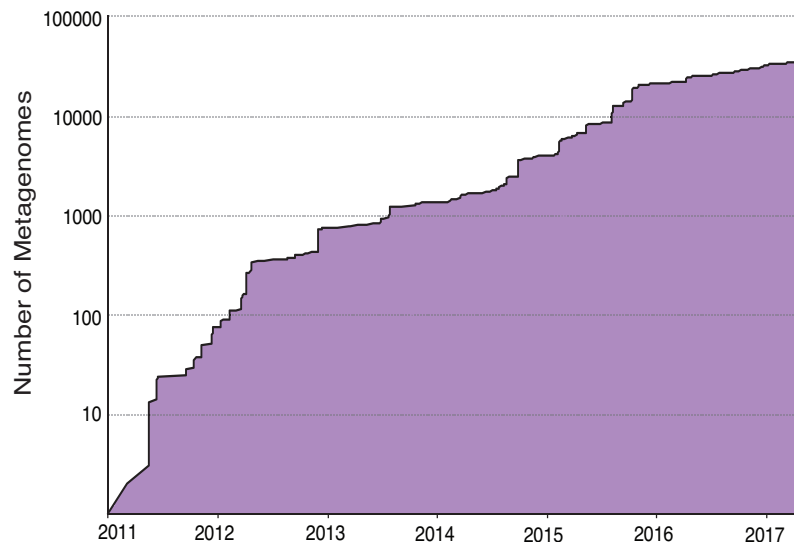




Understanding the microbiome

Who, what, why, how?

New Era of Microbiome Data Science



NMDC launched in 2019 by DOE to collect, analyze, and serve microbial data.

Metagenome Complexity

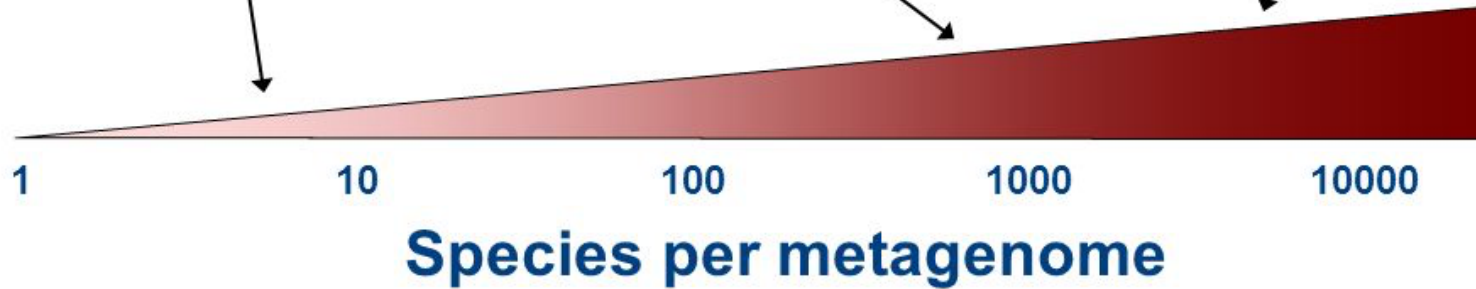
Acid mine



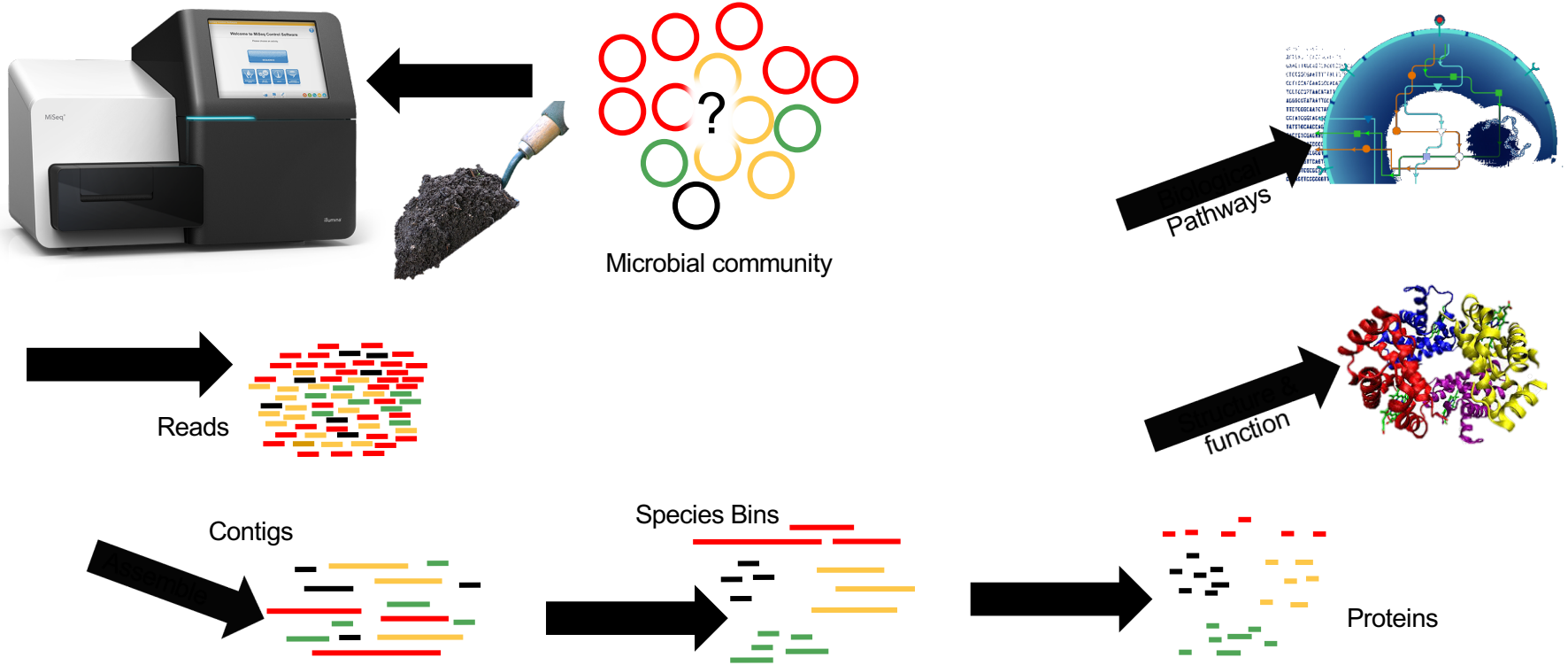
Cow rumen



Soil



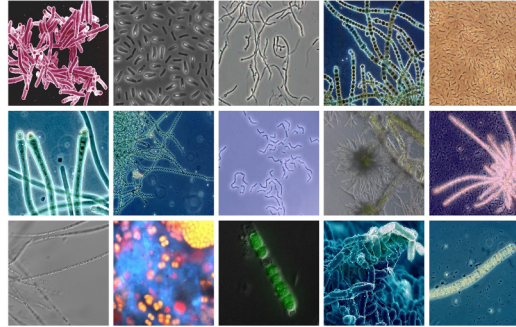
Microbiome analysis: metagenome



Big Science Questions



What happens to microbes after a wildfire? (1.5TB)



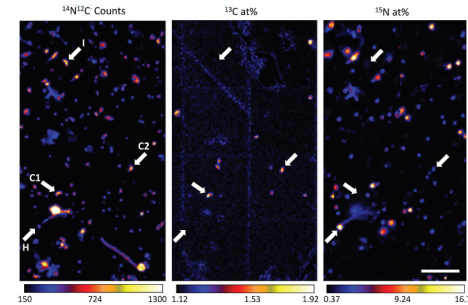
What are the microbial dynamics of soil carbon cycling? (3.3 TB)



How do microbes affect disease and growth of switchgrass for biofuels (4TB)

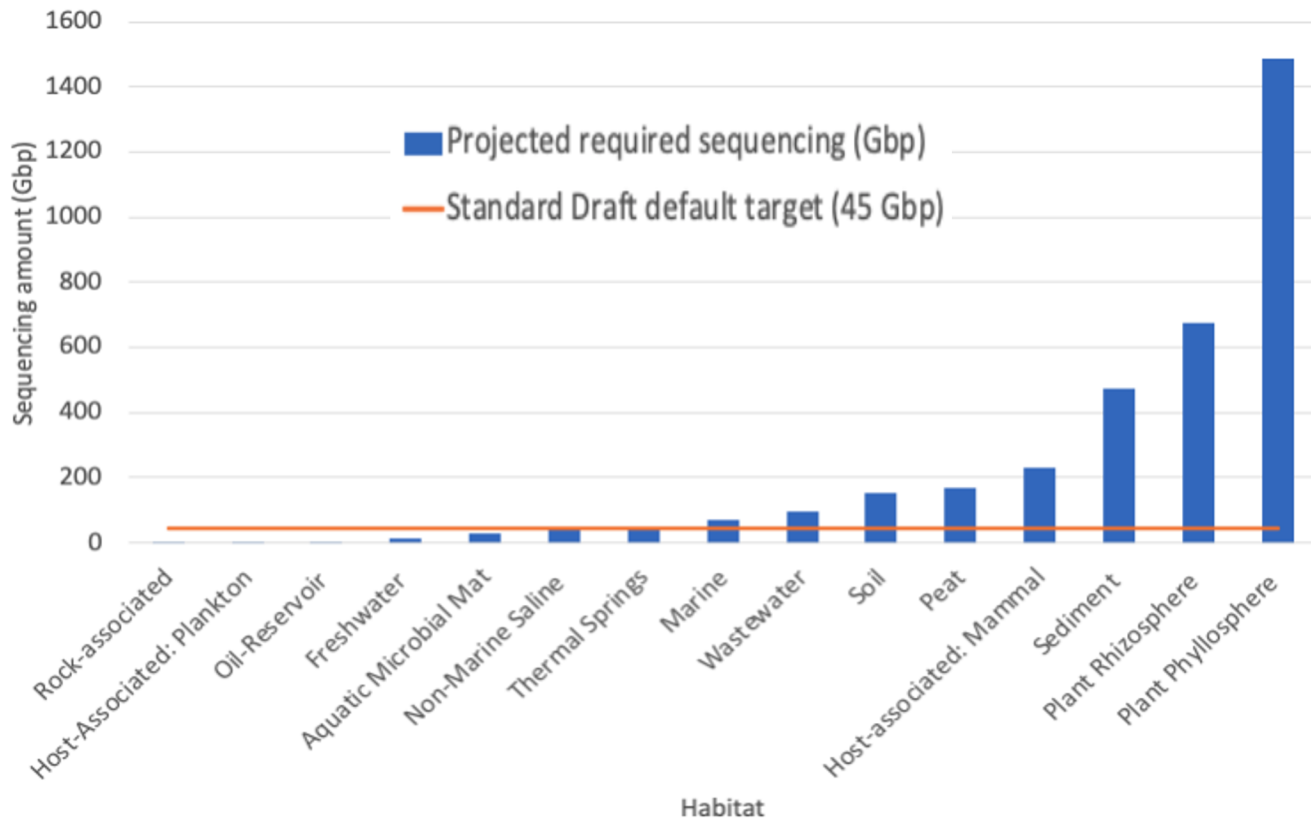


What are the seasonal fluctuations in a wetland mangrove? (1.6 TB)



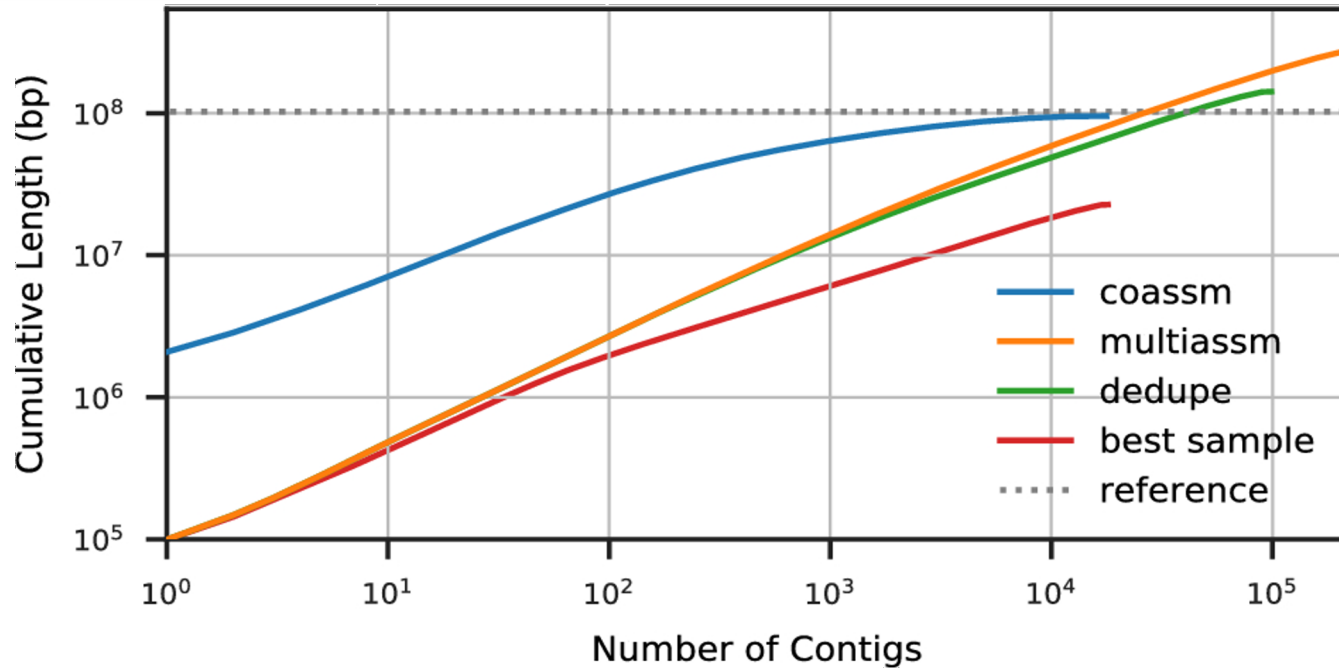
Combine genomics with isotope tracing methods for improved functional understanding (8TB)

Complex metagenomes require terascale data



Coverage analysis done by Kelly Cobaugh using Nonpareil 3 (Rodriguez-R, et al. mSystems 2018)

Big Data, Big Iron → Better Science



Multiassembly: assembling many samples separately

Coassembly: assembling many samples together

Clustering Huge Protein Datasets

- **Protein families via clustering**
- **Functional diversity**
 - Oceans vs human microbiome
- **New genes and proteins, e.g.**
 - Novel CRISPR/Cas genes
 - Gene clusters encoding antibiotics

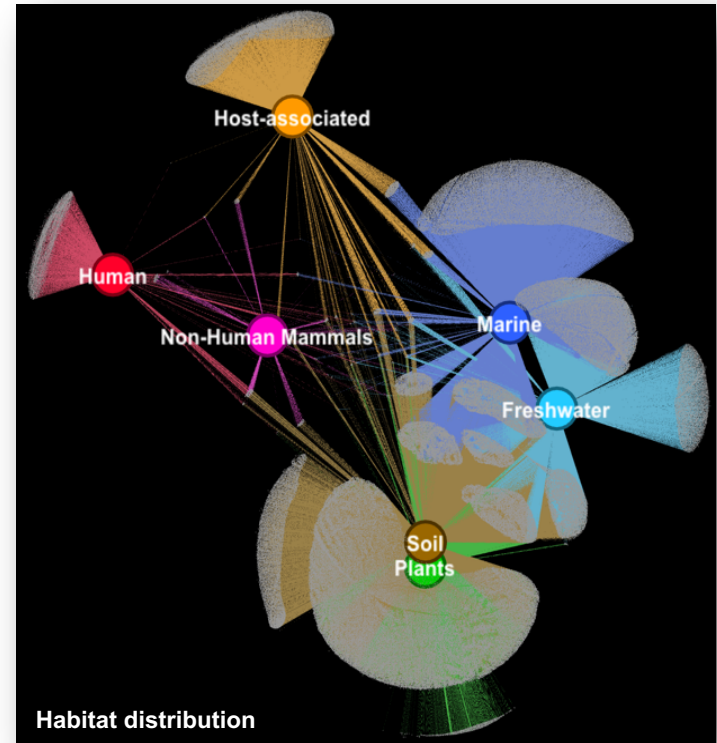


Image: G. Pavlopoulos and N. Kyrpides

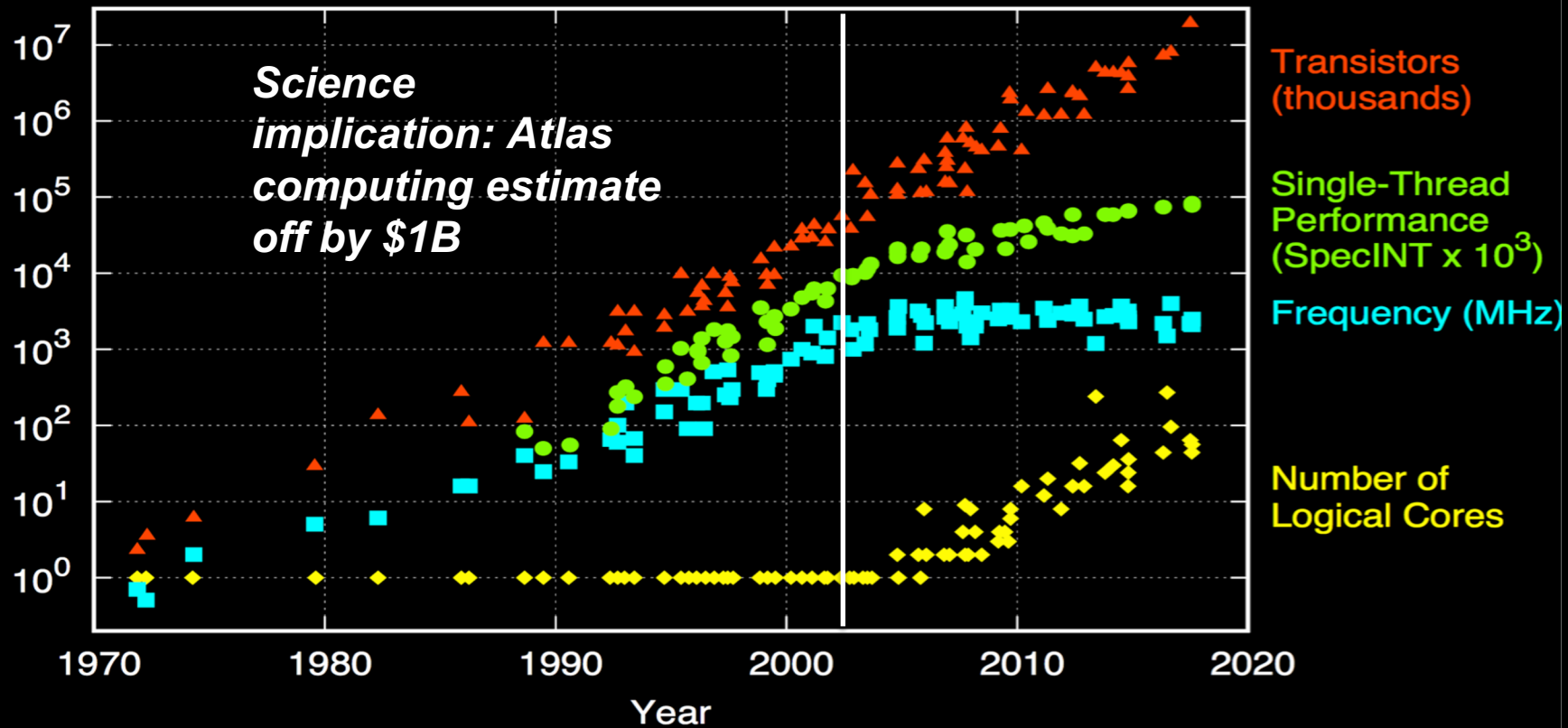


Moore's Law

It's hard to think exponentially

But it's also hard to stop

Dennard Scaling is Dead; Moore's Law Will Follow

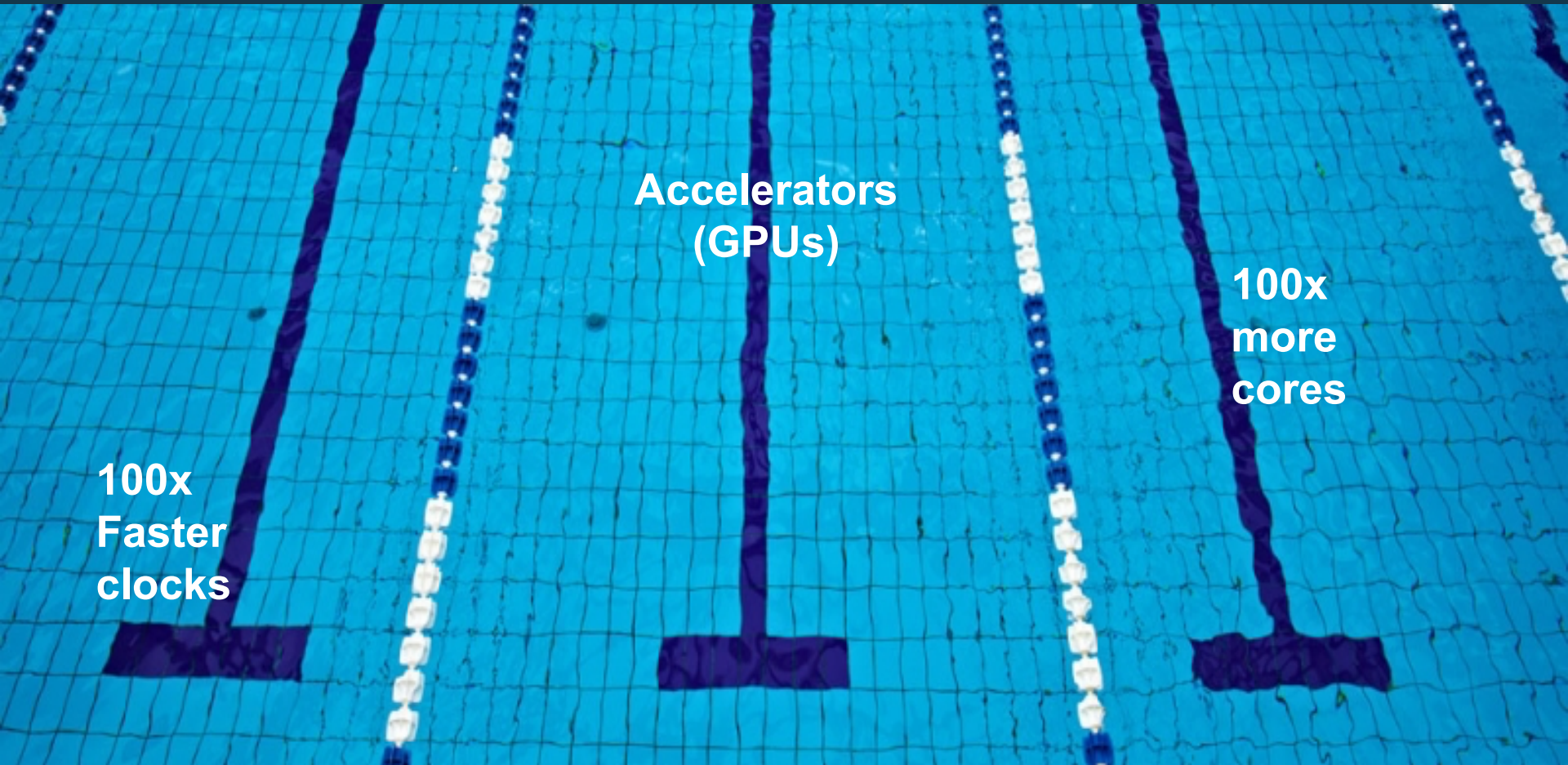


Exascale Architecture Plans (2008)

100x
Faster
clocks

Accelerators
(GPUs)

100x
more
cores



Exascale Architecture Plans (2021)

US DOE Office of Science Systems



Exascale
HPE AMD+AMD

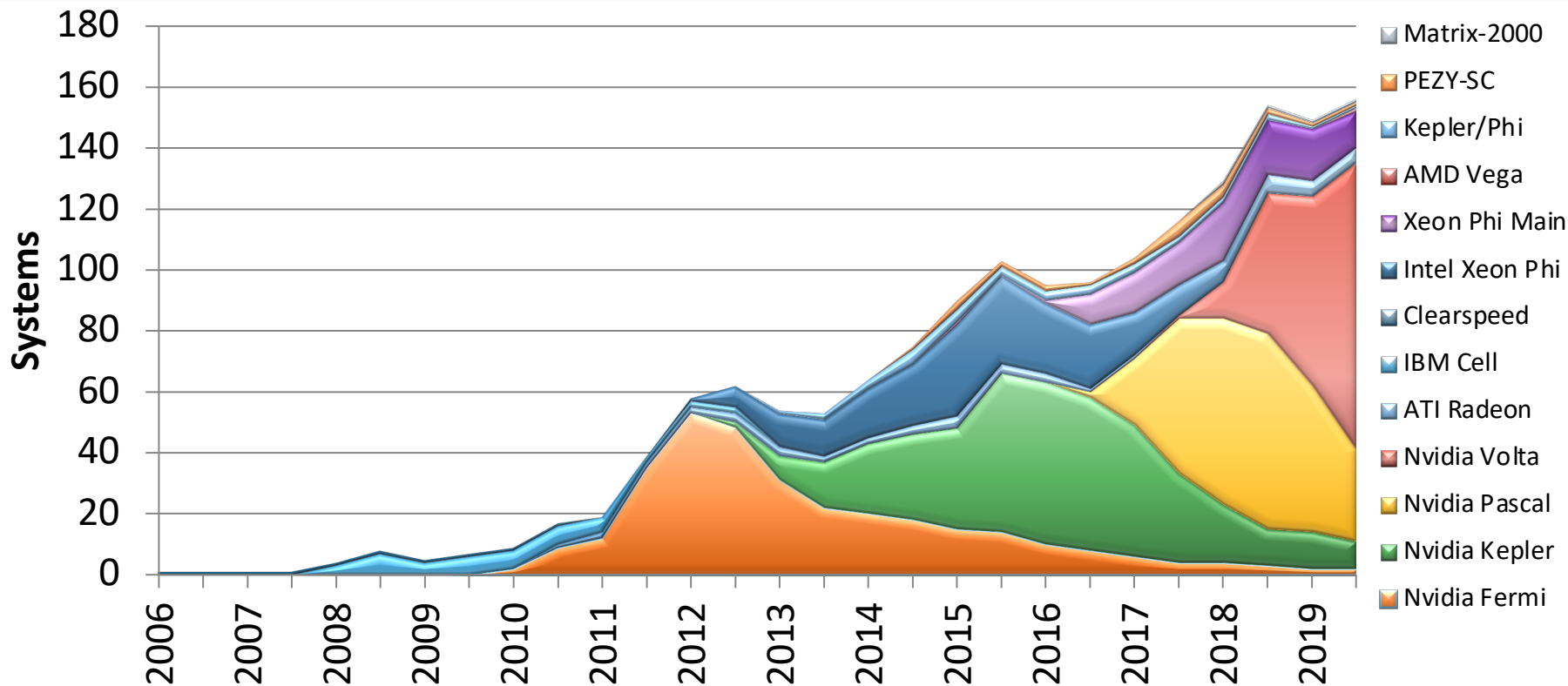


Exascale
HPE Intel+Intel

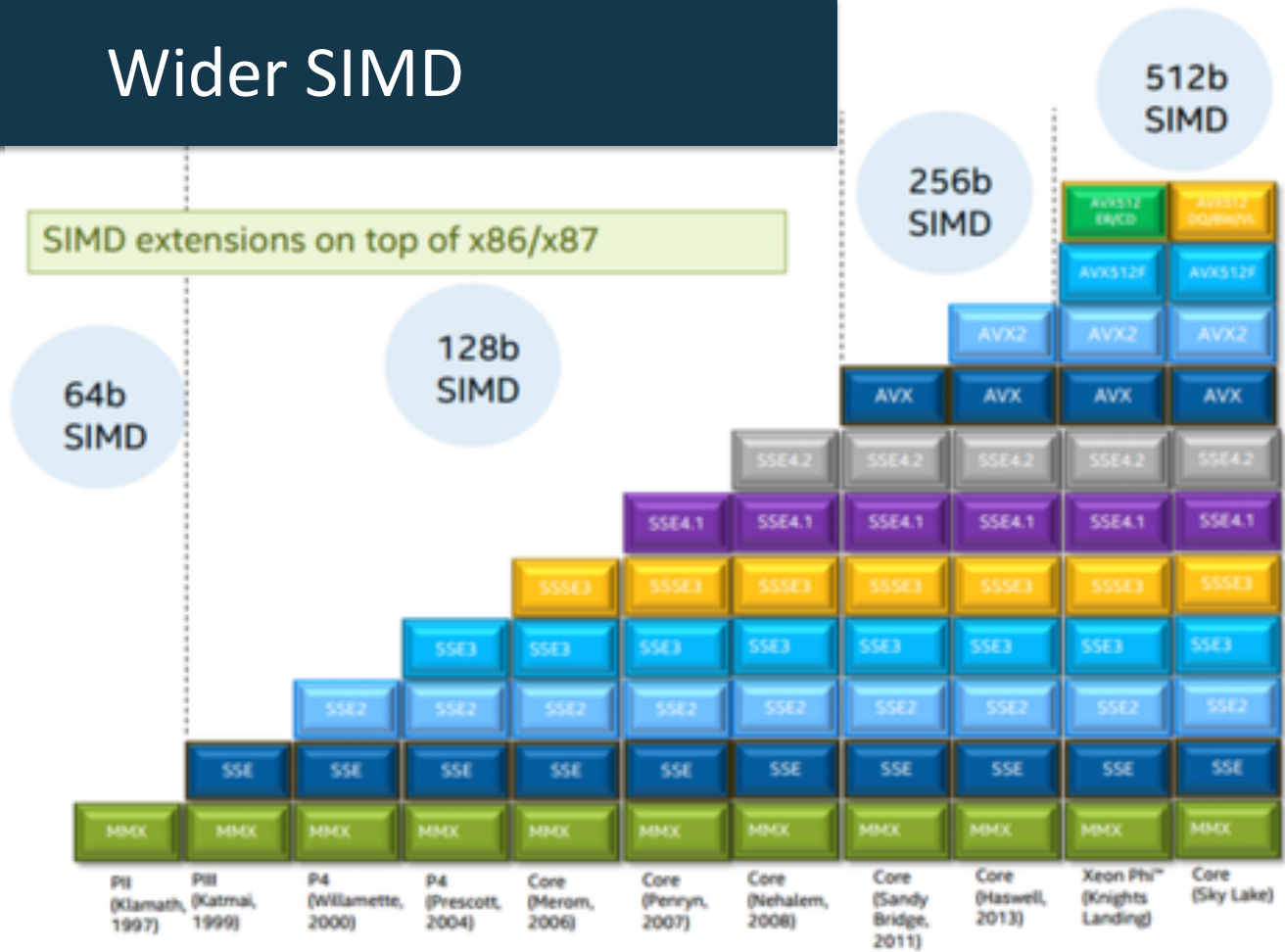


Pre-exascale
HPE AMD+NVIDIA

Accelerators

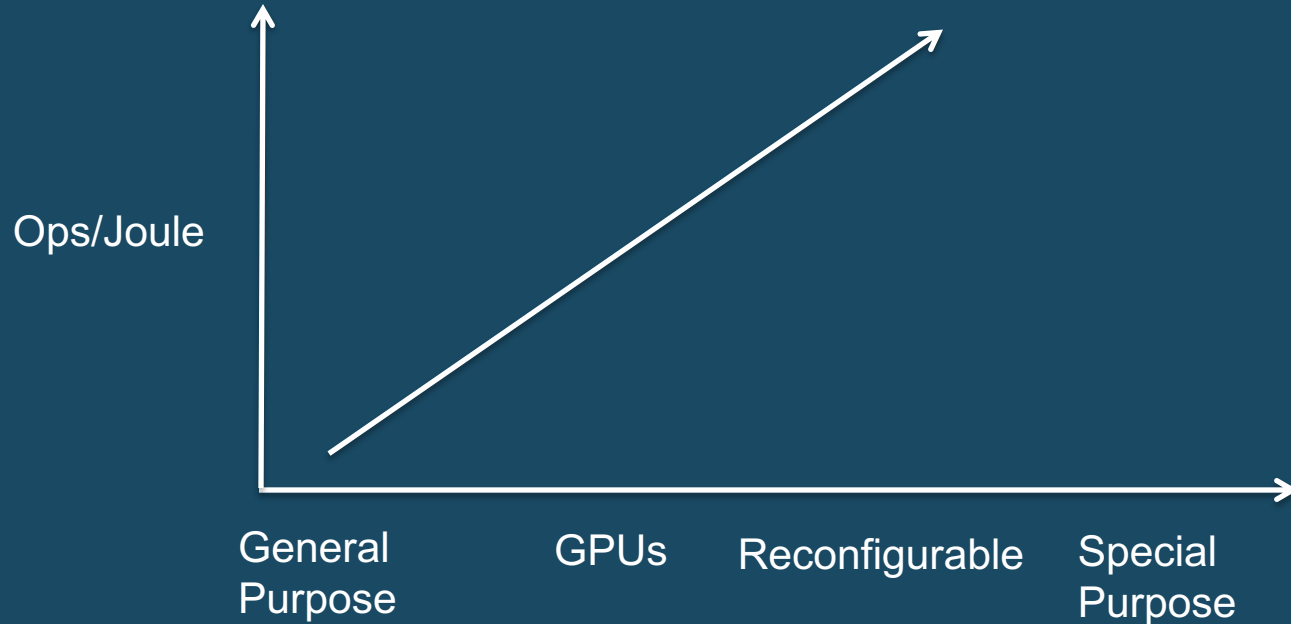


Wider SIMD

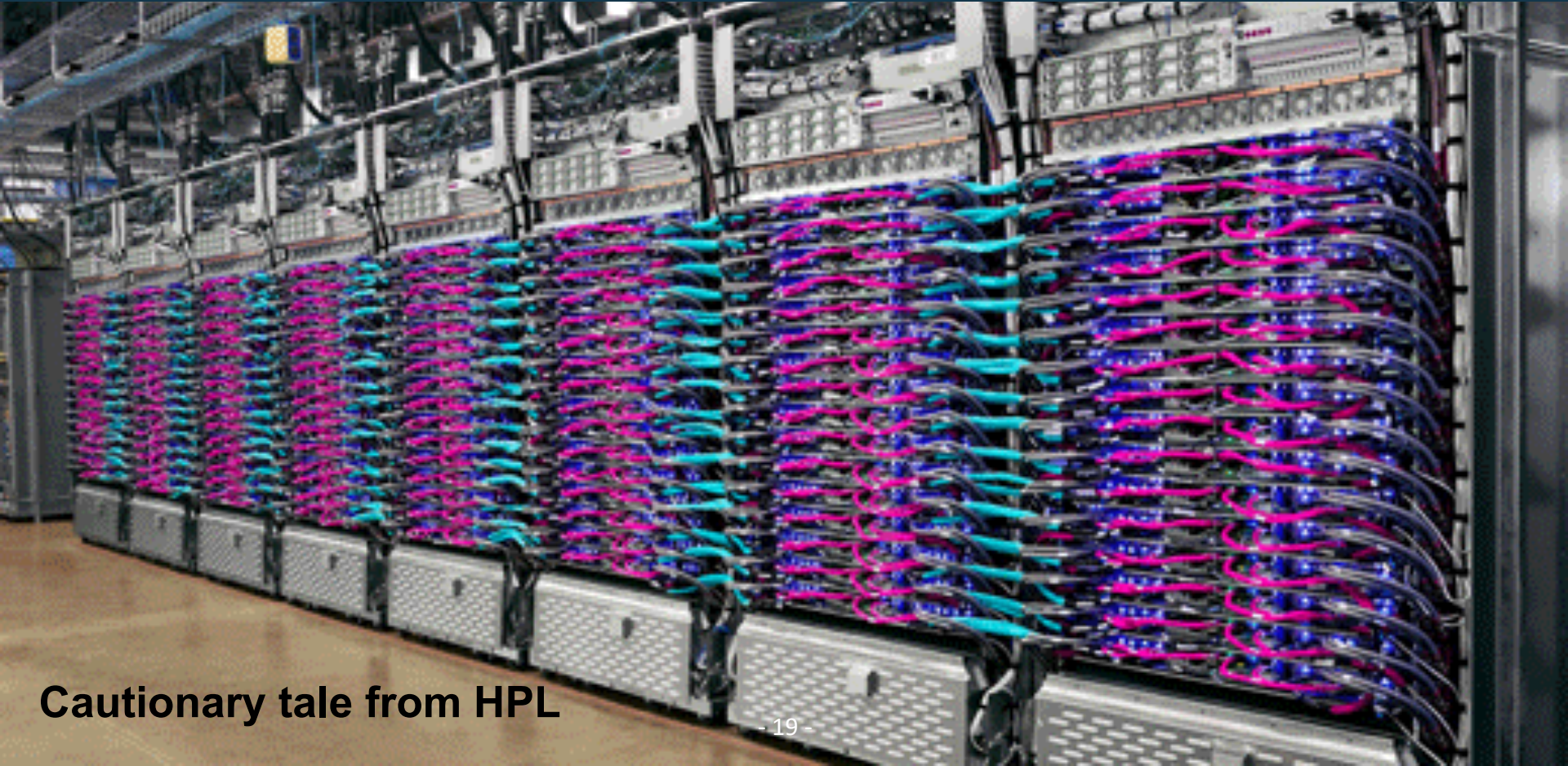


Top 500
#1 Fugaku
ARM w/ 512b

Specialization: End Game for Moore's Law



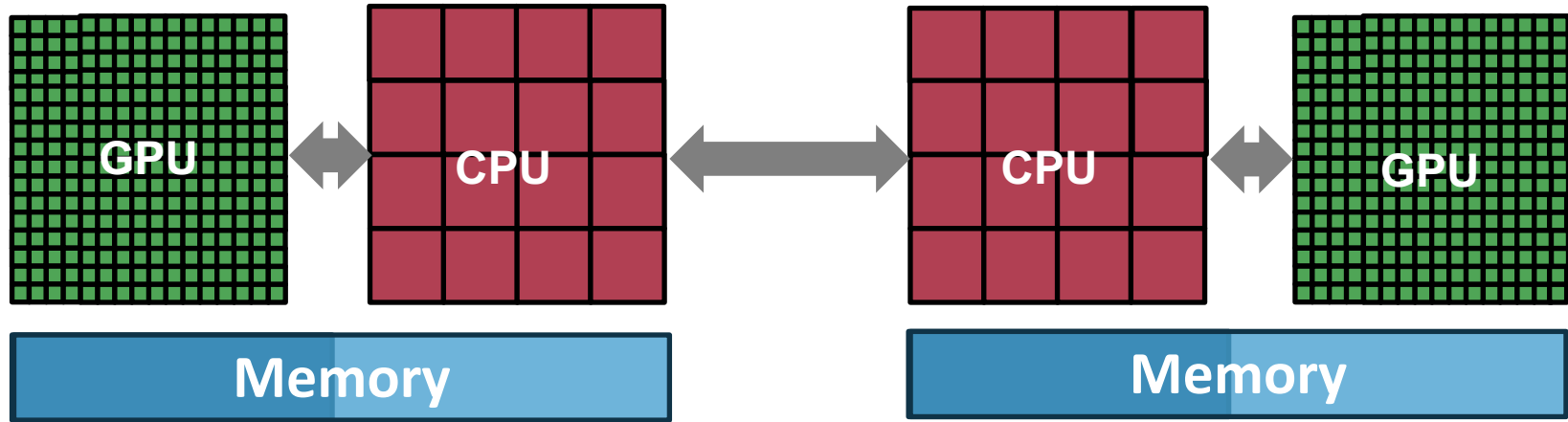
Is deep learning the only application?



Cautionary tale from HPL

Specialization, Yes

Accelerators, No!



More
cores

More data
parallelism

Narrow
data types

More
memory
spaces

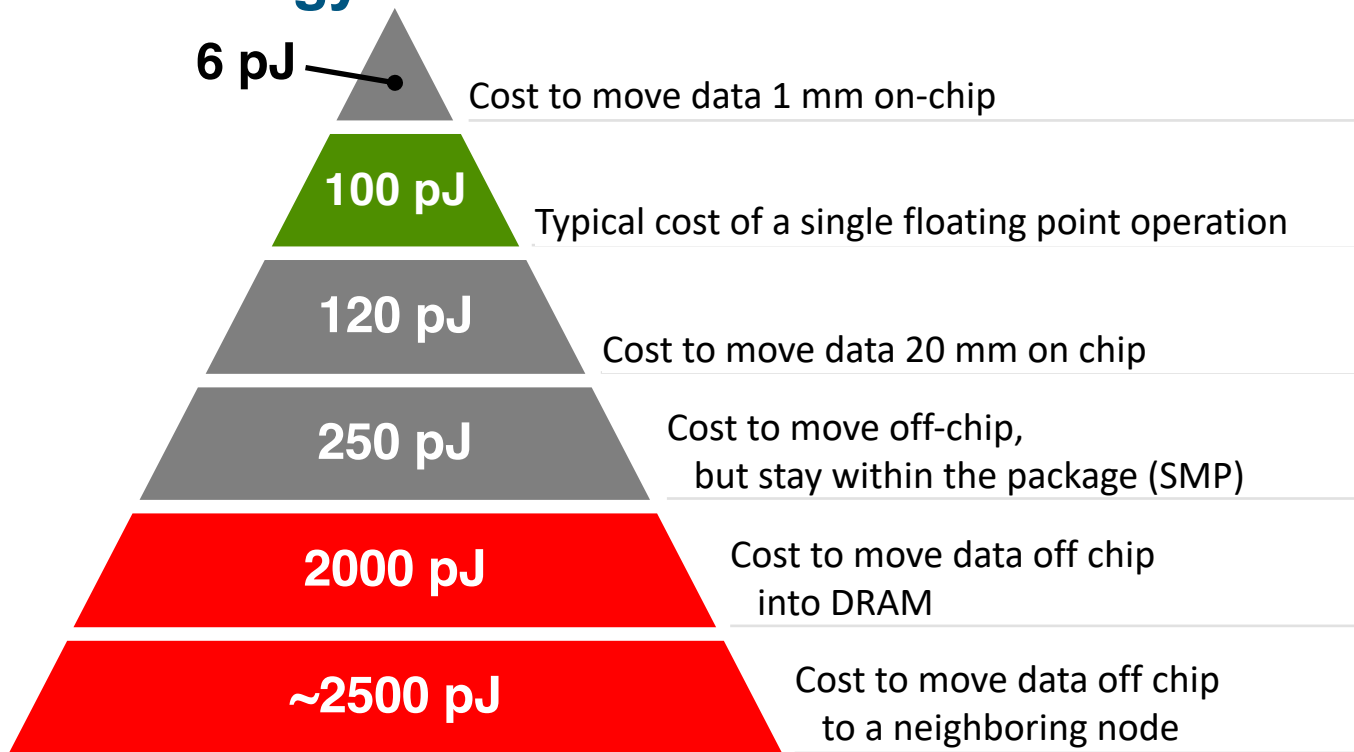
CPUs in
control

CPUs
communicate

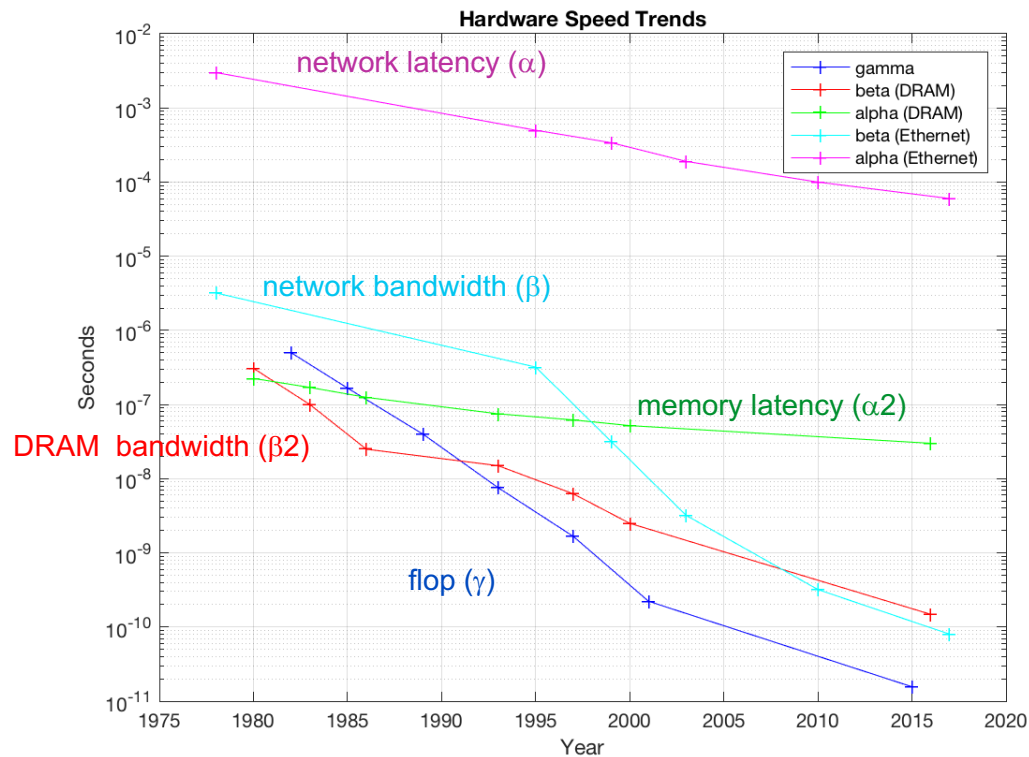


Data Movement is Expensive

Hierarchical energy costs.



Communication Dominates: Dennard was too good



Time =

$$\# \text{ flops} * \gamma +$$

$$\# \text{ message} * \alpha +$$

$$\# \text{ bytes comm} * \beta +$$

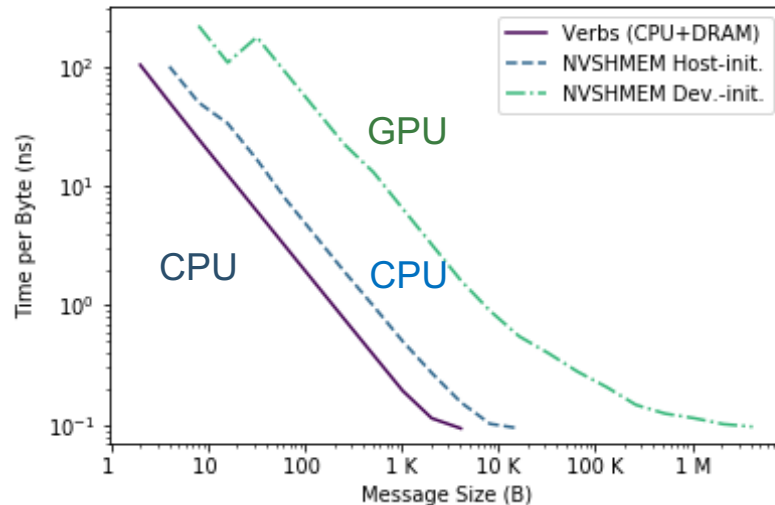
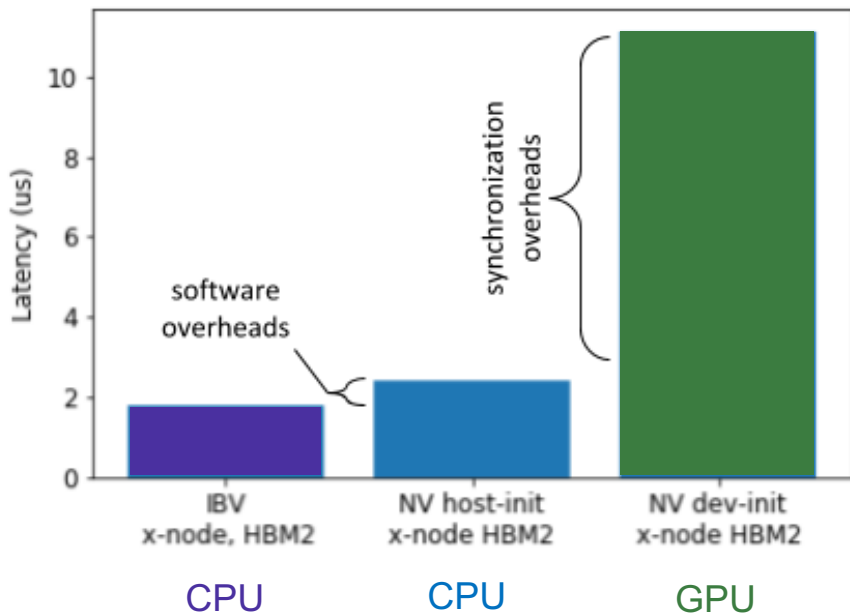
$$\# \text{ diff memory locs} * \alpha^2 +$$

$$\# \text{ memory words} * \beta^2$$

Data from Hennessy / Patterson, Graph from Demmel

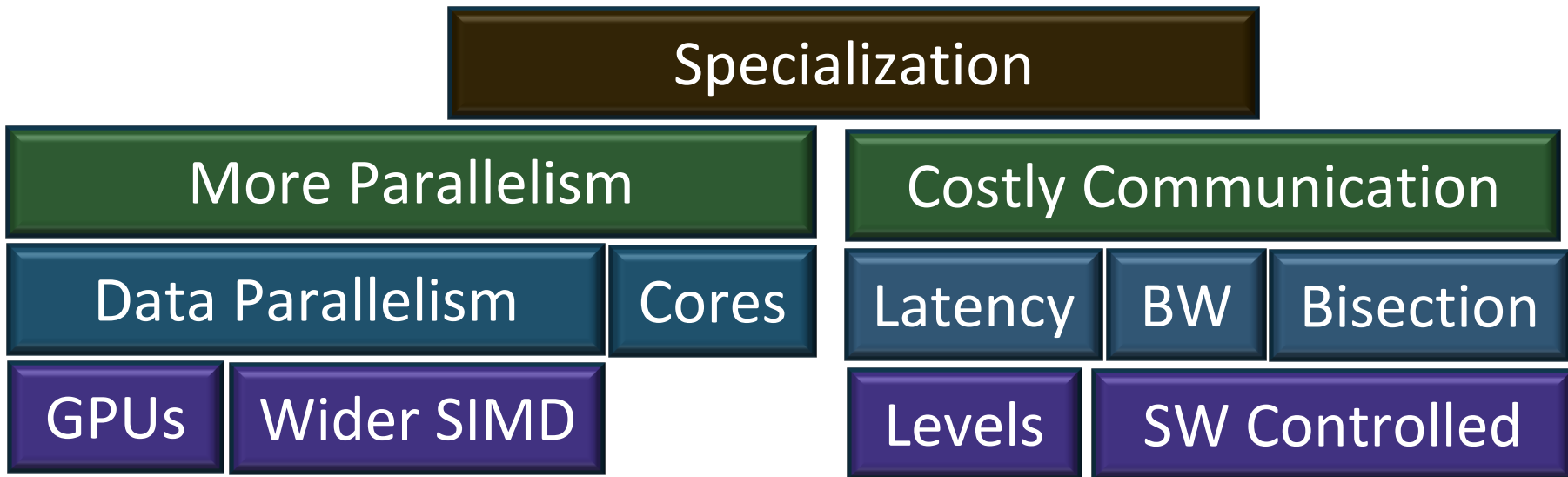
Put Accelerators in Charge of Communication

Architecture and software are not yet structured for accelerated-initiated communication (Summit with NVLink between Power9 CPUs and NVIDIA GPUs)



Taylor Groves et al

Hardware Trends



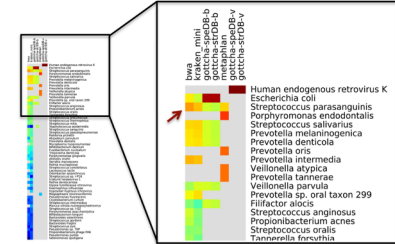
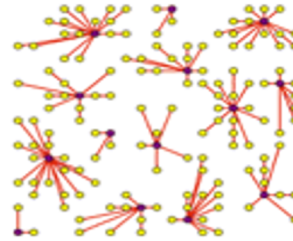
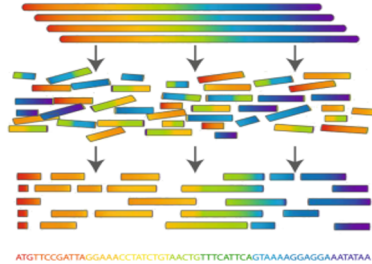
Tradeoffs in integration (faster communication) vs scale (amount of fast memory) and flexibility



Algorithms and Software

ExaBiome project overview

Exascale algorithms & systems for previously intractable problems



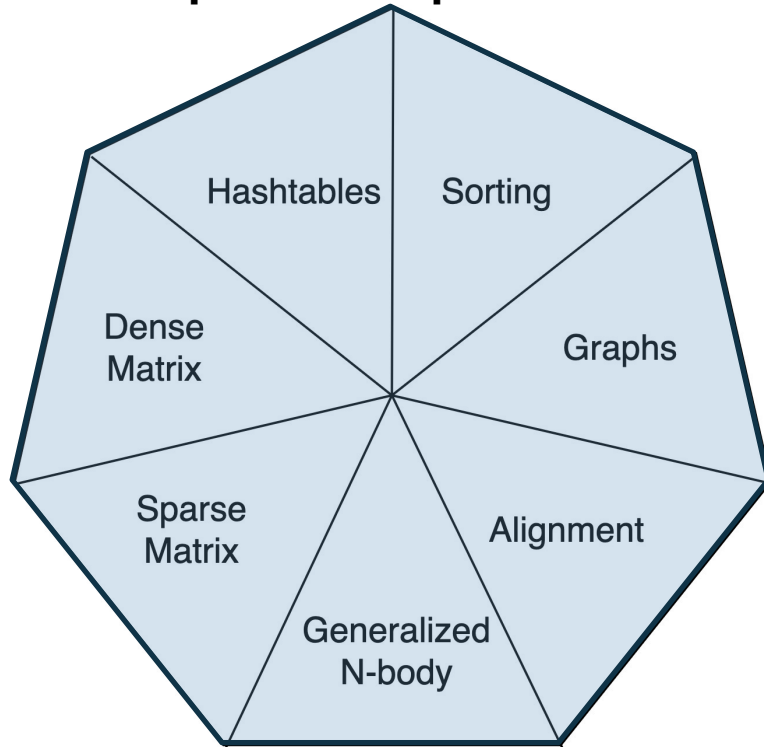
	Metagenome Assembly	Protein Clustering and Annotation	Comparative Metagenome Analysis
Science Need	Find species, genes and relative abundance in microbial communities	Improve understanding of tree of life for microbes; aid in identifying gene function	Track microbiome over time or space, changes in environment, climate, etc.
Computing Technique	hash tables, alignment, k-mer counts, graph walks	direct tables, alignment, k-mer counts, sparse matrices, ML (clustering, GNNs)	hash tables, alignment, k-mer counts, ML (dimensionality reduction)



<http://exabiome.org>

Motifs of Genomic Data Analysis

These computational patterns dominate ExaBiome Project experience



Application problems

- Assemble genomes
- Compute distances
- Cluster (contigs, proteins,...)
- Annotate

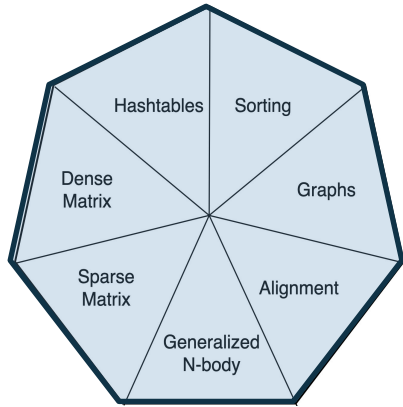
Distributed memory platforms open up new approaches and science questions

Analytics vs. Simulation Kernels:

7 Giants of Data	7 Dwarfs of Simulation
Basic statistics	Monte Carlo methods
Generalized N-Body	Particle methods
Graph-theory	Unstructured meshes
Linear algebra	Dense Linear Algebra
Hashing	Sparse Linear Algebra
Sorting	Spectral methods
Alignment	Structured Meshes

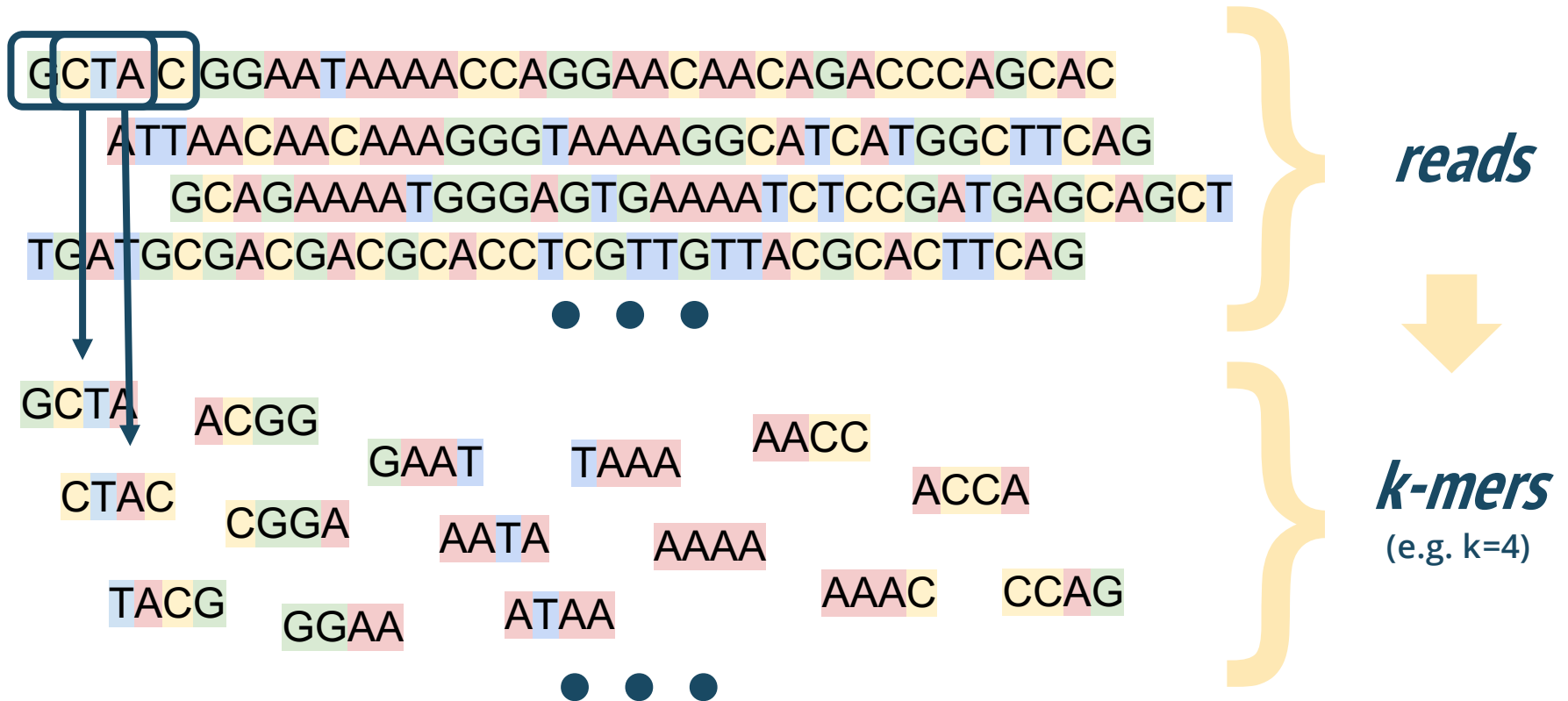
NRC Report + our paper

Phil Colella



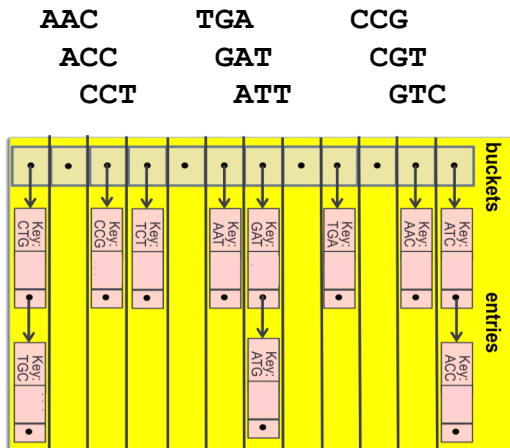
Hashing

Common Technique: Analyze K-mers



Distributed Hash Tables of K-Mers

Make hash table of k-mers



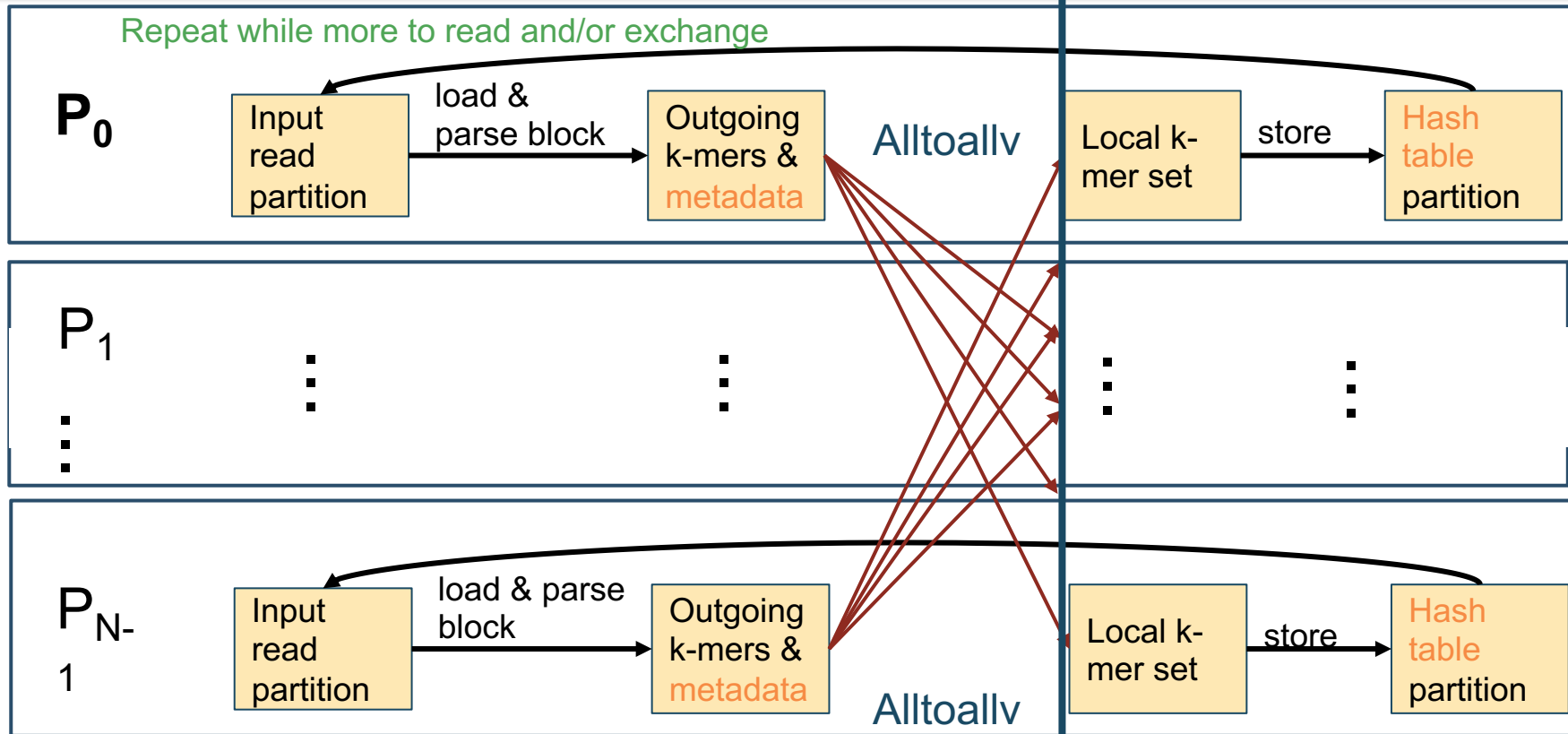
1-sided comm or irregular all-to-all + memory

Keys are fixed-length strings:

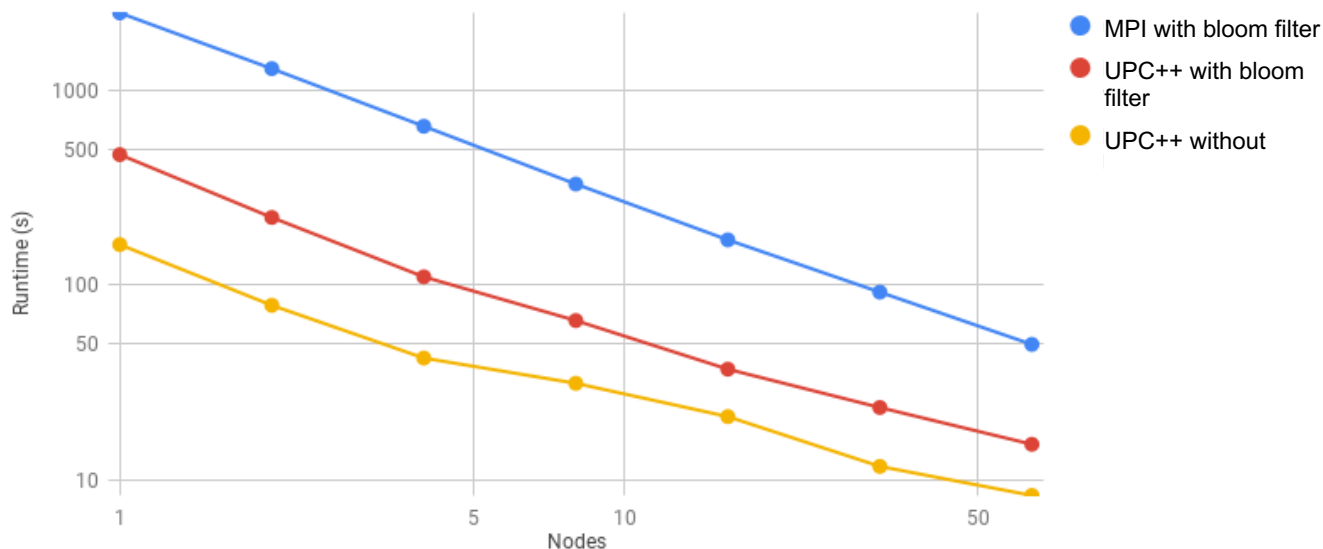
Values depend on application:

- Count and remove errors (singletons)
- Find strings with matching k-mer
- Connected components
- Use histogram as approximation

Distributed Hashing / Histogramming



K-mer counting: All the Wires All the Time

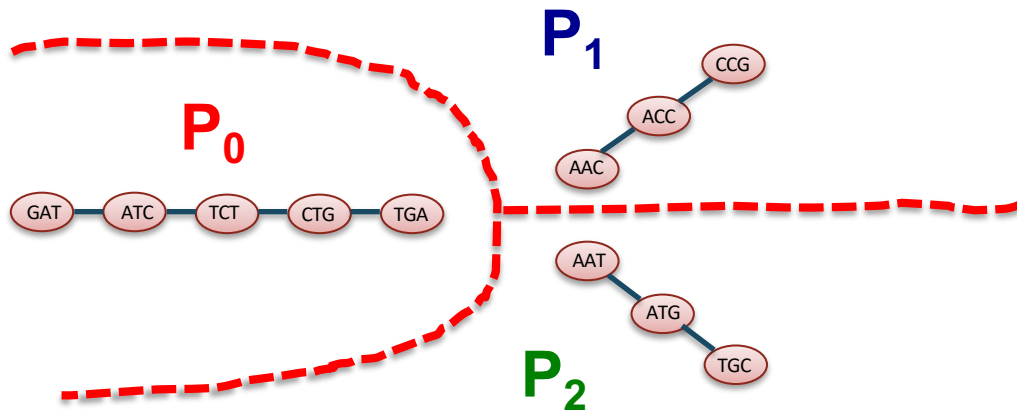


- Used to be bulk-synchronous MPI
- UPC++ communication is asynchronous and 1-sided
- UPC++ version is faster, avoids barriers, saves memory (one runtime)
- And it's simpler!

Communication-Avoiding hash table lookup

Caching for temporal locality (reuse): if few large items, so lookups will repeat

Layout for spatial locality: if we have an “oracle” that approximate final genome

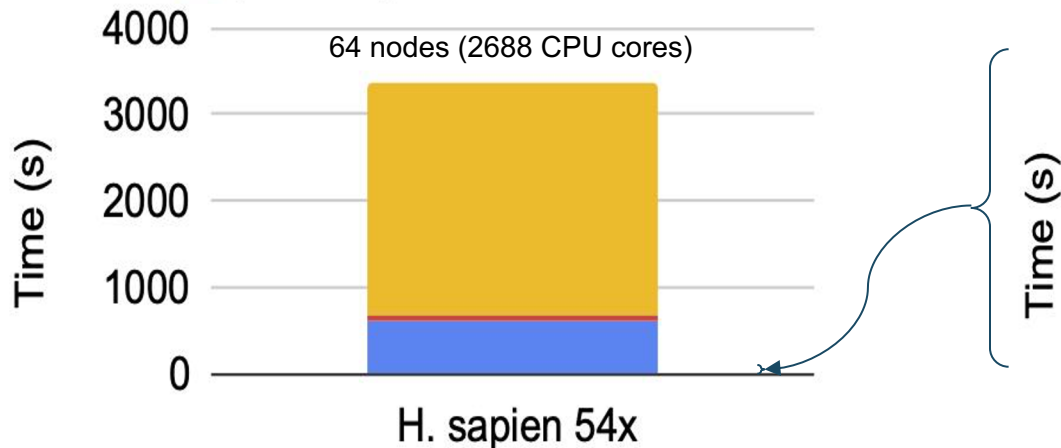


Traversal is up to 2.8x faster!
Up to 76% reduction of off-node communication !

K-mer Counting: Finding Data Parallelism

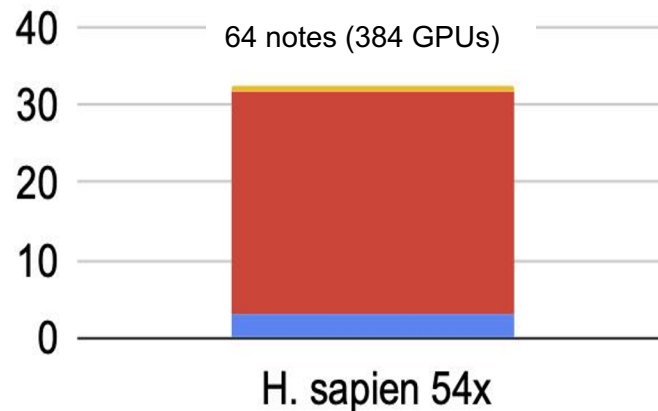
■ kmer counter ■ exchange (incl. MPI call)

■ parse & process kmers



■ kmer counter ■ exchange (incl. MPI call)

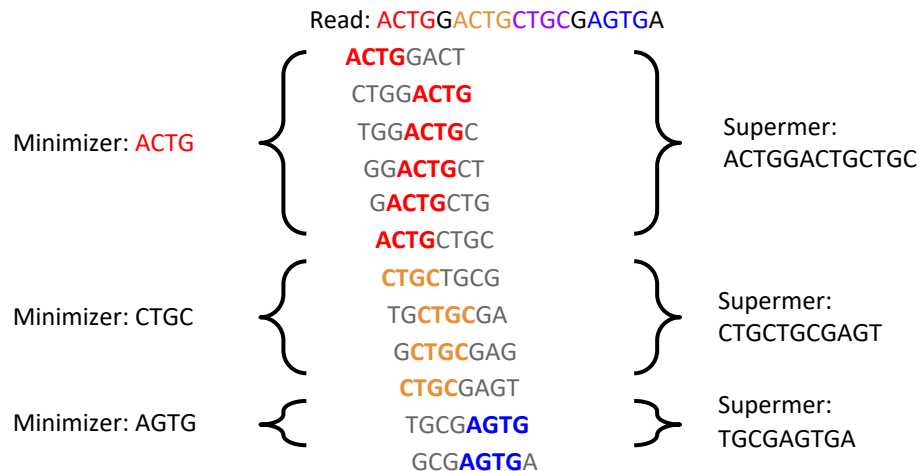
■ parse & process kmers



- K-mer counter on Summit. (Note scales -- red k-mer exchange time is roughly equal.)
- Reduce CPU/GPU communication by parsing as well as processing on GPU

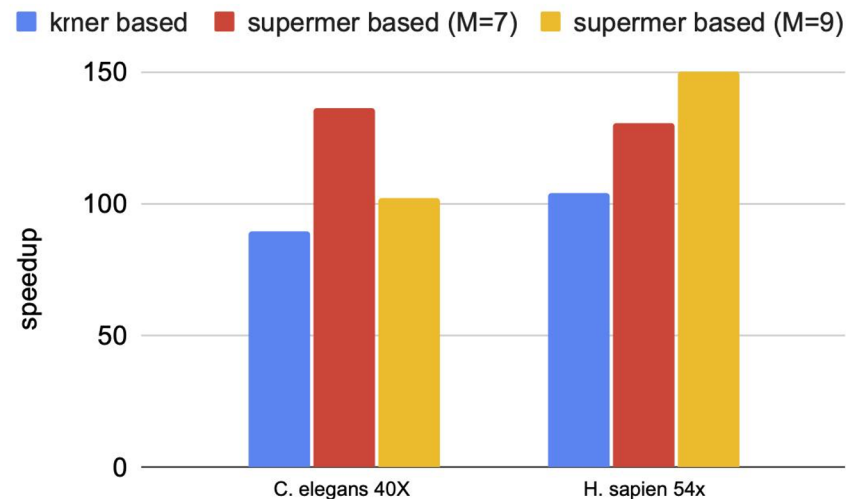
Over 100x speedup!!

K-mer Counting: Reducing Communication



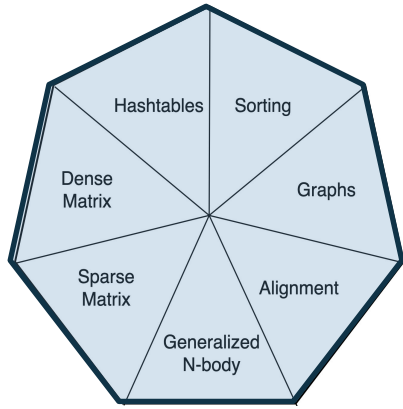
Reduce communication with “Supermers”

- Multiple contiguous k-mer
- map to the same process ID with minimizer-based hashing
- Saves volume (bandwidth) and number of messages (latency)



Speedup on 64 Summit nodes

- 6 GPUs / node
- baseline: 42 cores / node



Alignment

Smith-Waterman: Dynamic Programming

	_	G	A	T	C	A	G	C	T
_	0	0	0	0	0	0	0	0	0
G	0	1	0	0	0	0	0	0	0
A	0	0	2	0	0	1	0	0	0
T	0	0	0	3	1	0	0	0	1
A	0	0	1	1	2	2	0	0	0
G	0	0	0	0	2	1	3	1	0
C	0	0	0	0	1	1	2	4	2
C	0	0	0	0	1	0	2	3	3

GATCACCT
GAT_ACCC

Scoring
insert/delete = -2
match = 1
mismatch = -1.

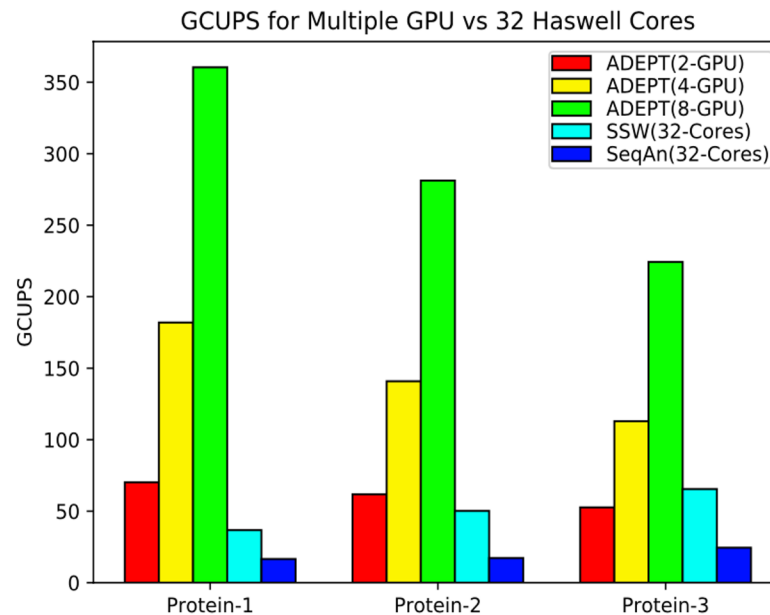
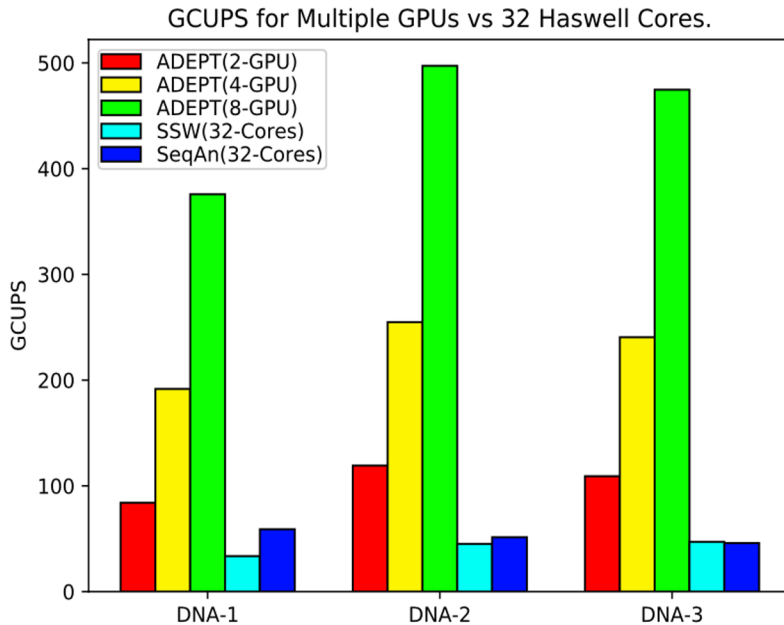
Options to search matrix

- Full search (Smith-Waterman)
- Banded (only search near diagonal)
- X-Drop stop search when the score drops by more than X

Variations for local / global alignment, per character penalties, seeding, etc.

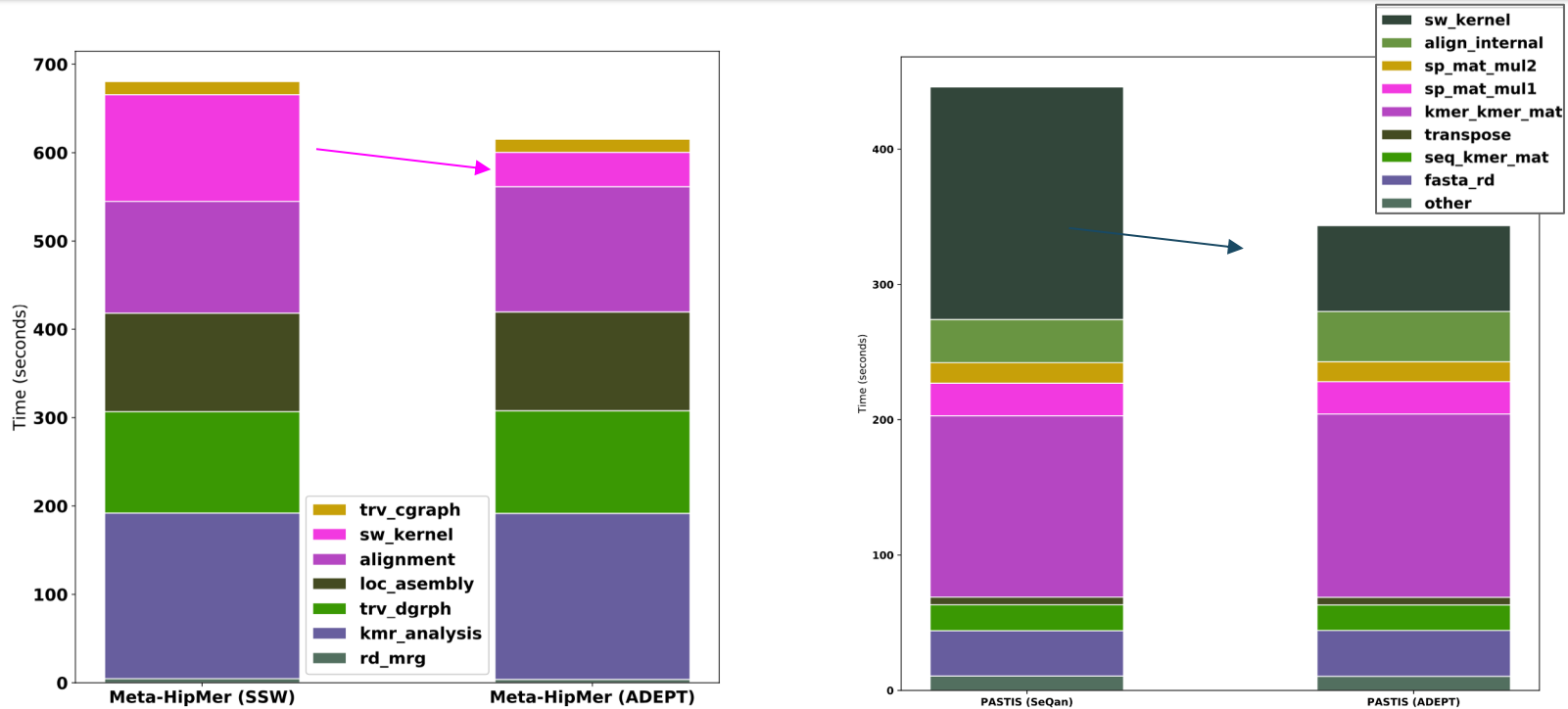
ADEPT: Batch Alignment on GPUs

GCUPS = Giga-Cell Updates per second

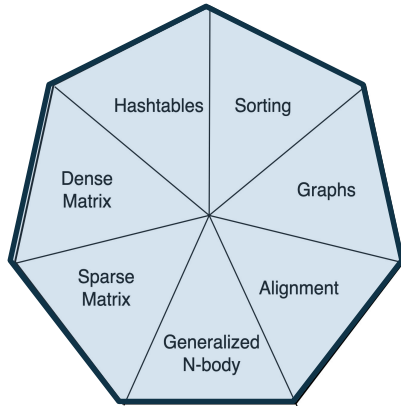


Adept is designed for relatively short, low-error sequences, both DNA (left) and proteins (right)
SSW and SeqAn are vectorized implementations of Smith-Waterman Algorithm on CPU.

ADEPT: Impact on ExaBiome Applications



Soil assembly SW time: 2.8 node hours on Cori, 0.1 node hours on Summit (hidden behind CPU work)

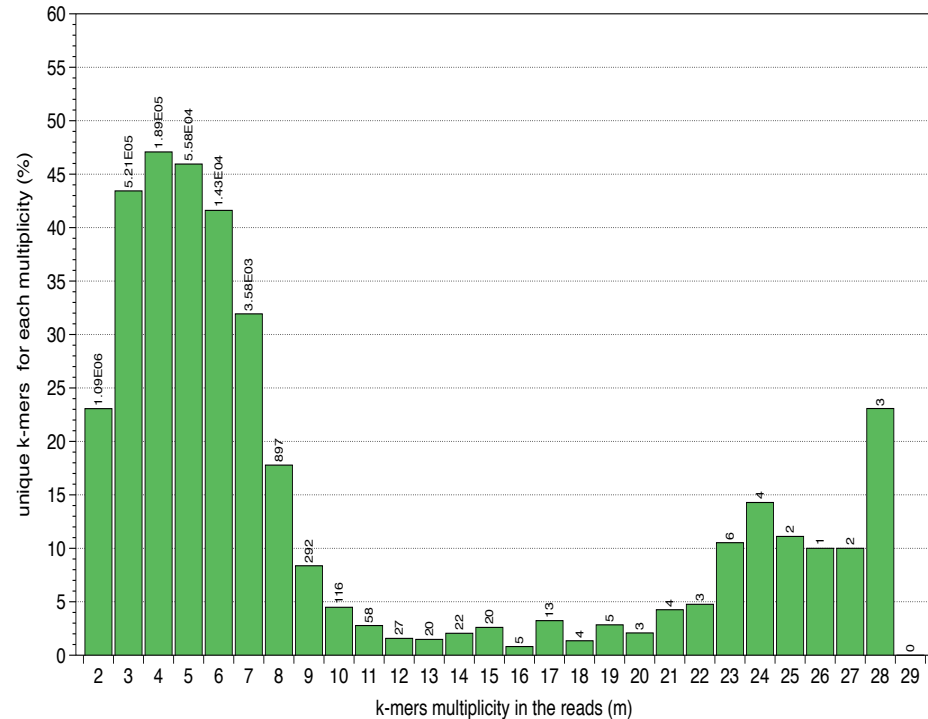


Generalized N-Body

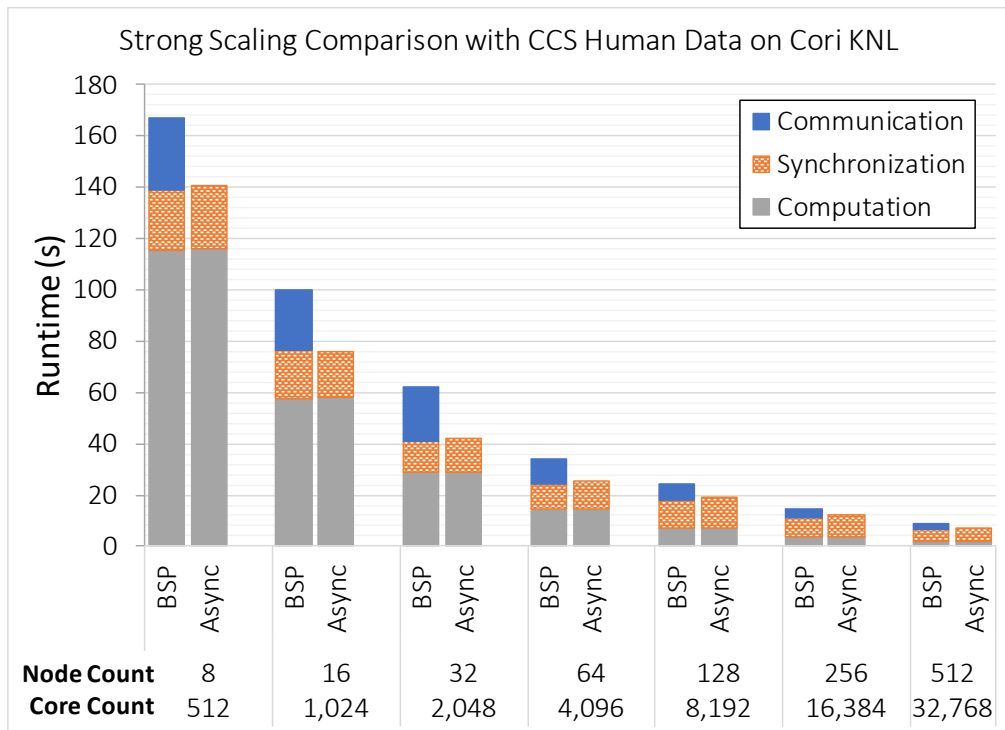
diBELLA: Towards a Long Read Assembler

Long reads (PacBio, etc.)

- Longer alignments
 - More compute-intensive
 - More GPU friendly
- ## No need for De Bruijn graph
- Pairwise alignments
 - Filtered k-mers



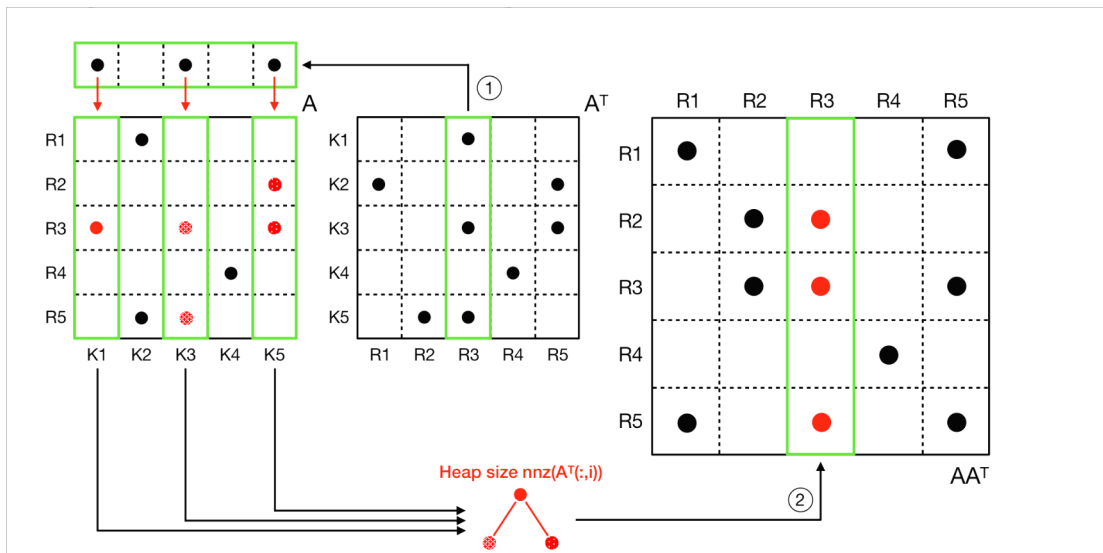
Bulk-Synchronous vs 1-sided Asynchronous



Asynchronous communication hides latency and uses less memory in general

Set Alignment is a Sparse All-to-All

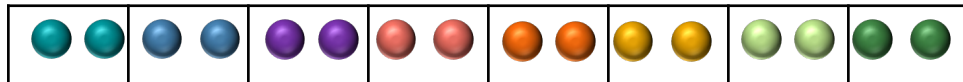
Run expensive alignment on all pairs with a common k-mer



Avoid Communication, Maximize Parallelism

Compute on all pairs of particles or strings, or...

Obvious solution

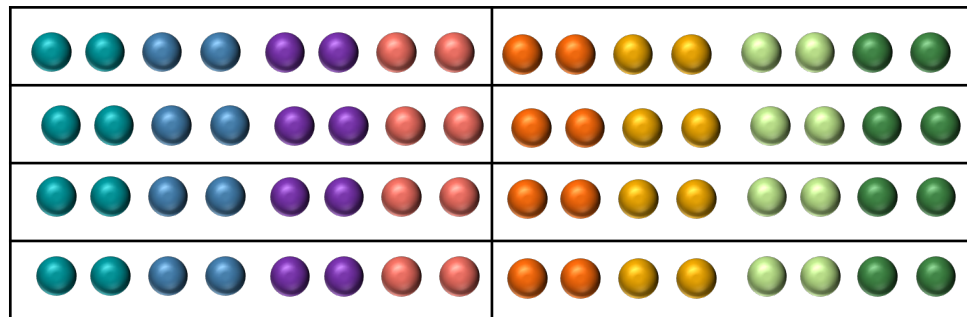


16 particles on 8 processors
Pass all particles around (p steps)

Decreases

- #messages by factor c^2
- #volume sent by factor c

Better solution

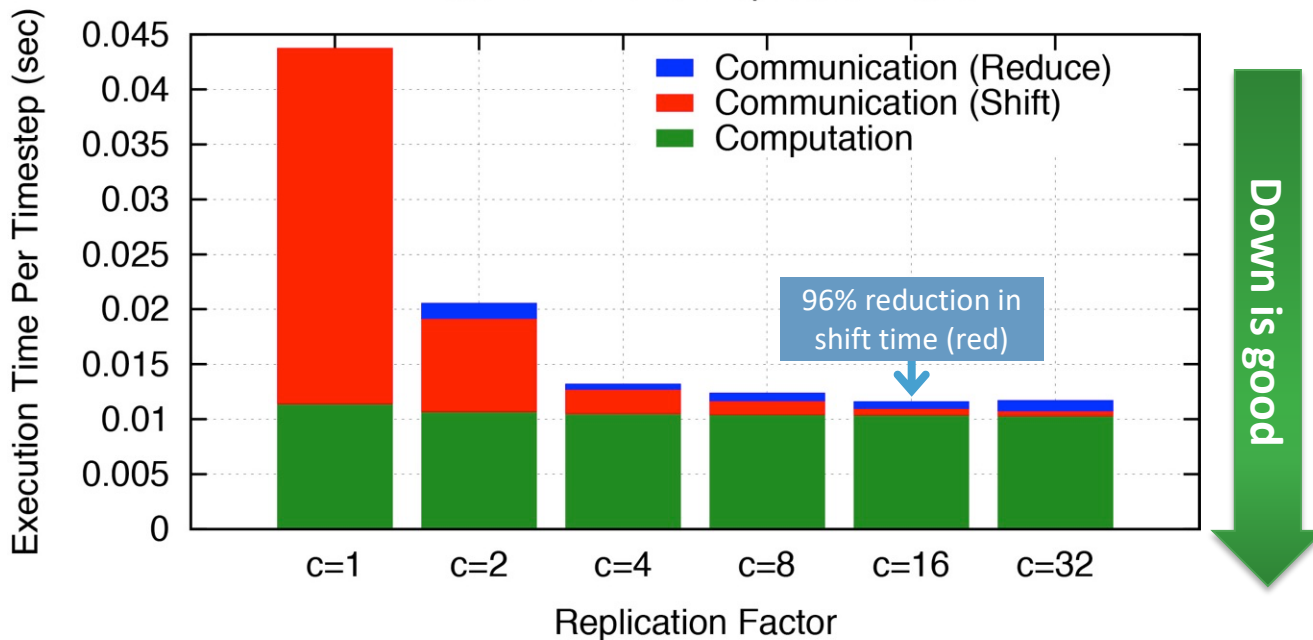


$c = 4$ copies of particles
8 particles on each

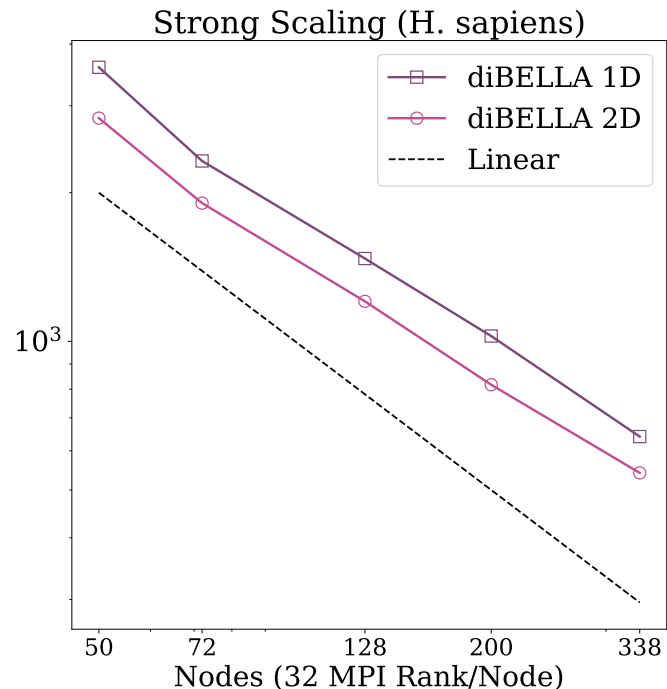
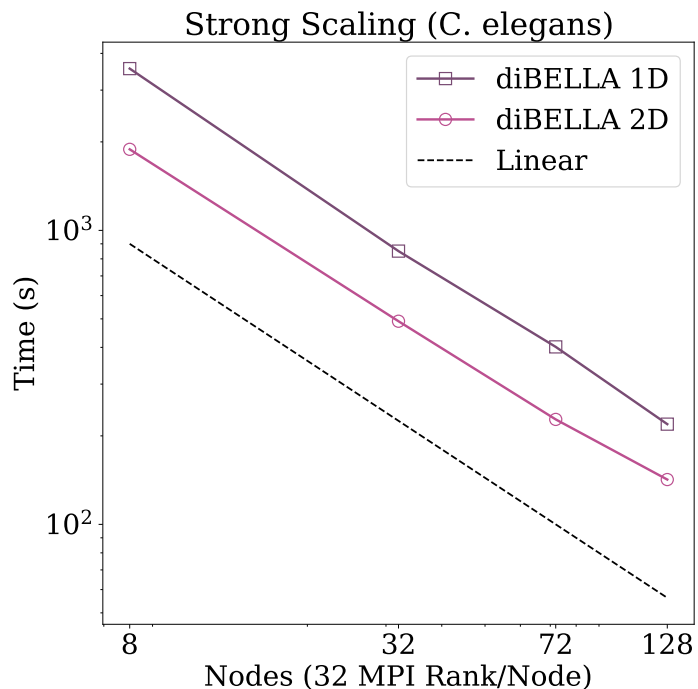
Less Communication..

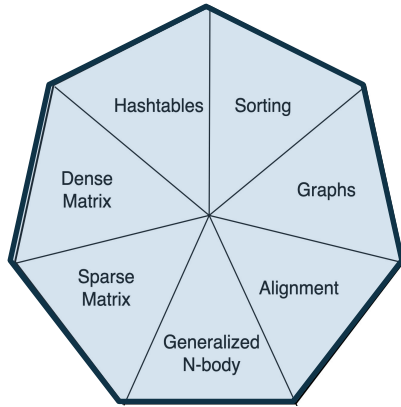
Cray XE6; n=24K particles, p=6K cores

Execution Time vs. Replication Factor



1D vs 2D Algorithm on DNA “overlap”



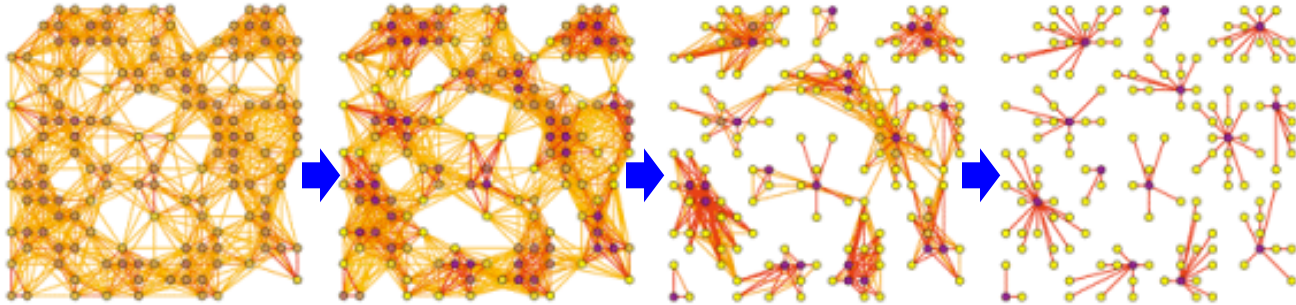


Graphs and Sparse Matrices (unsupervised learning)

Protein Clustering with Sparse Matrices

Input: Adjacency matrix A (sparse)

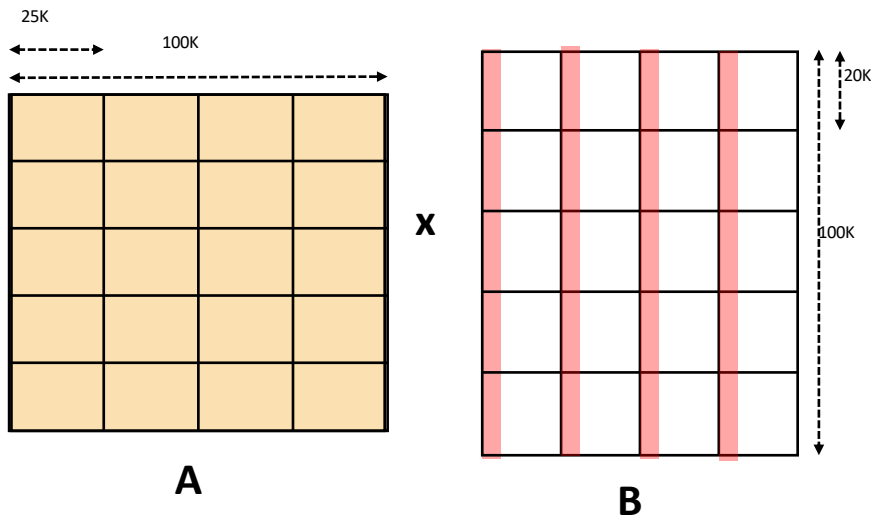
Image source: <http://micans.org/mcl/>



- **Similarity Matrix:** “Many-to-many” protein alignment
- **Expansion:** Square matrix, pruning small entries, dense columns
- **Inflation:** element-wise powers

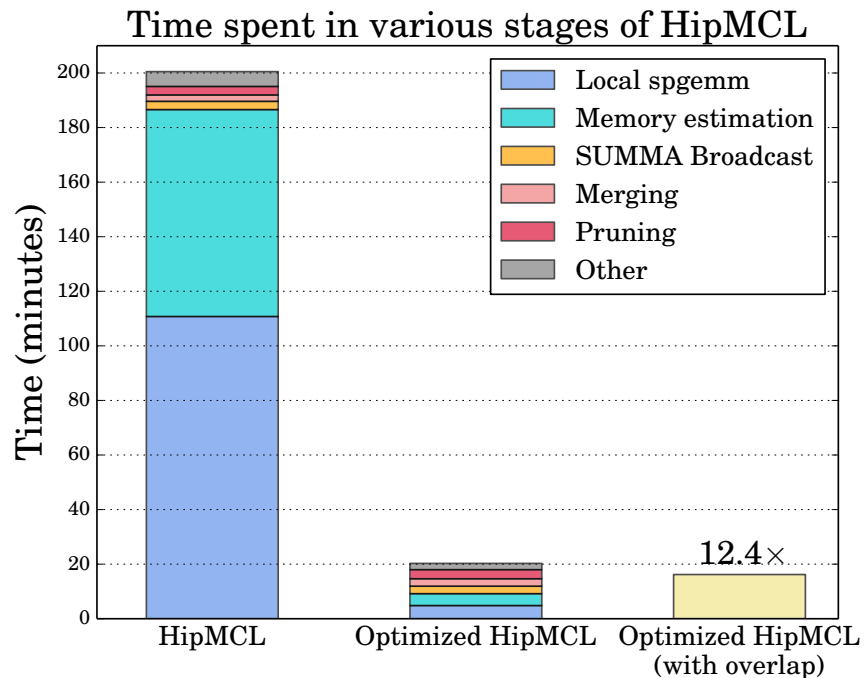
PASTIS + HipMCL

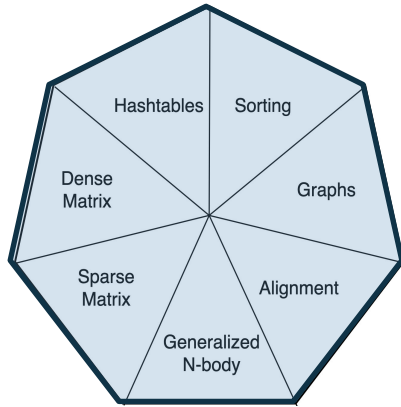
Sparse Matrix Algorithms



Distributed memory enabled new science

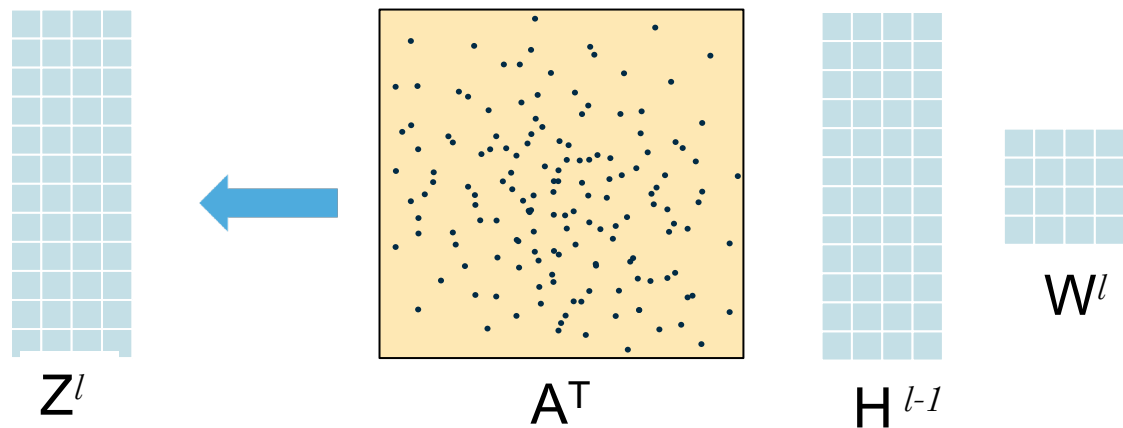
12.4× faster with GPUs!





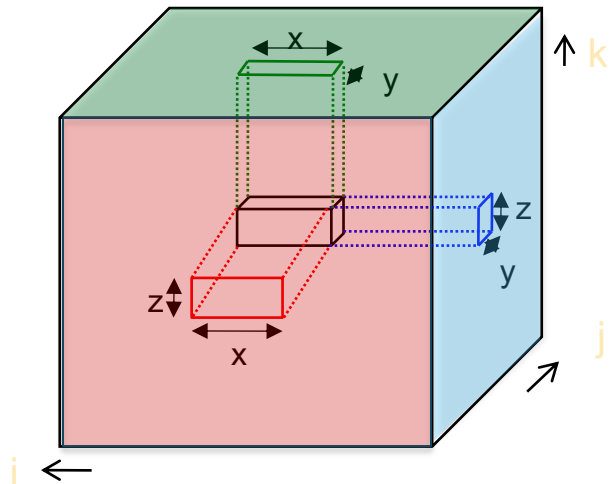
Graphs and Sparse and Dense Matrices (supervised learning)

Bottleneck in GNN Training



- $A^T H^{l-1}$ sparse-dense matmul (SpMM)
- $(A^T H^{l-1}) W^l$ dense-dense matmul (DGEMM)
- **SpMM is the bottleneck, not DGEMM!**

Communication-Avoiding Matrix Multiply

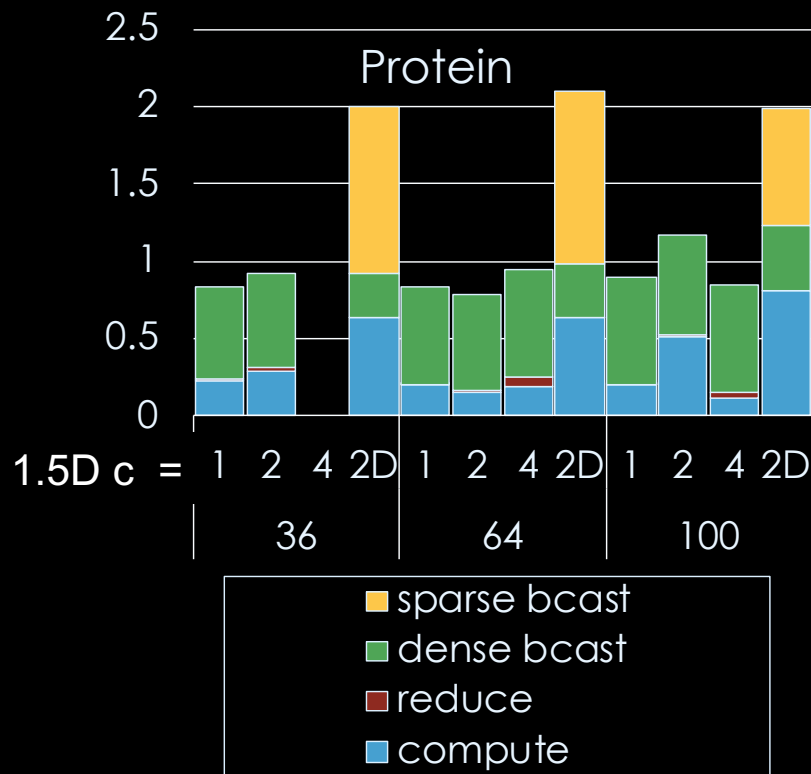
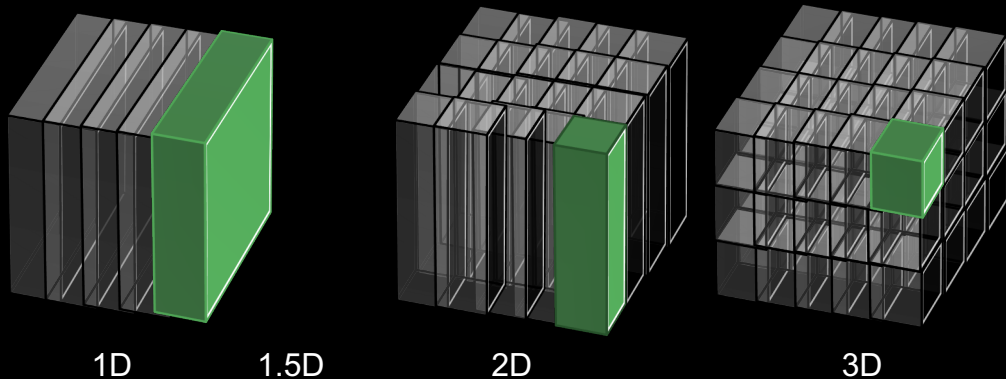


- 2D algorithm: never chop k dim
- 3D: Assume + is associative; chop k, which is \rightarrow replication of C matrix

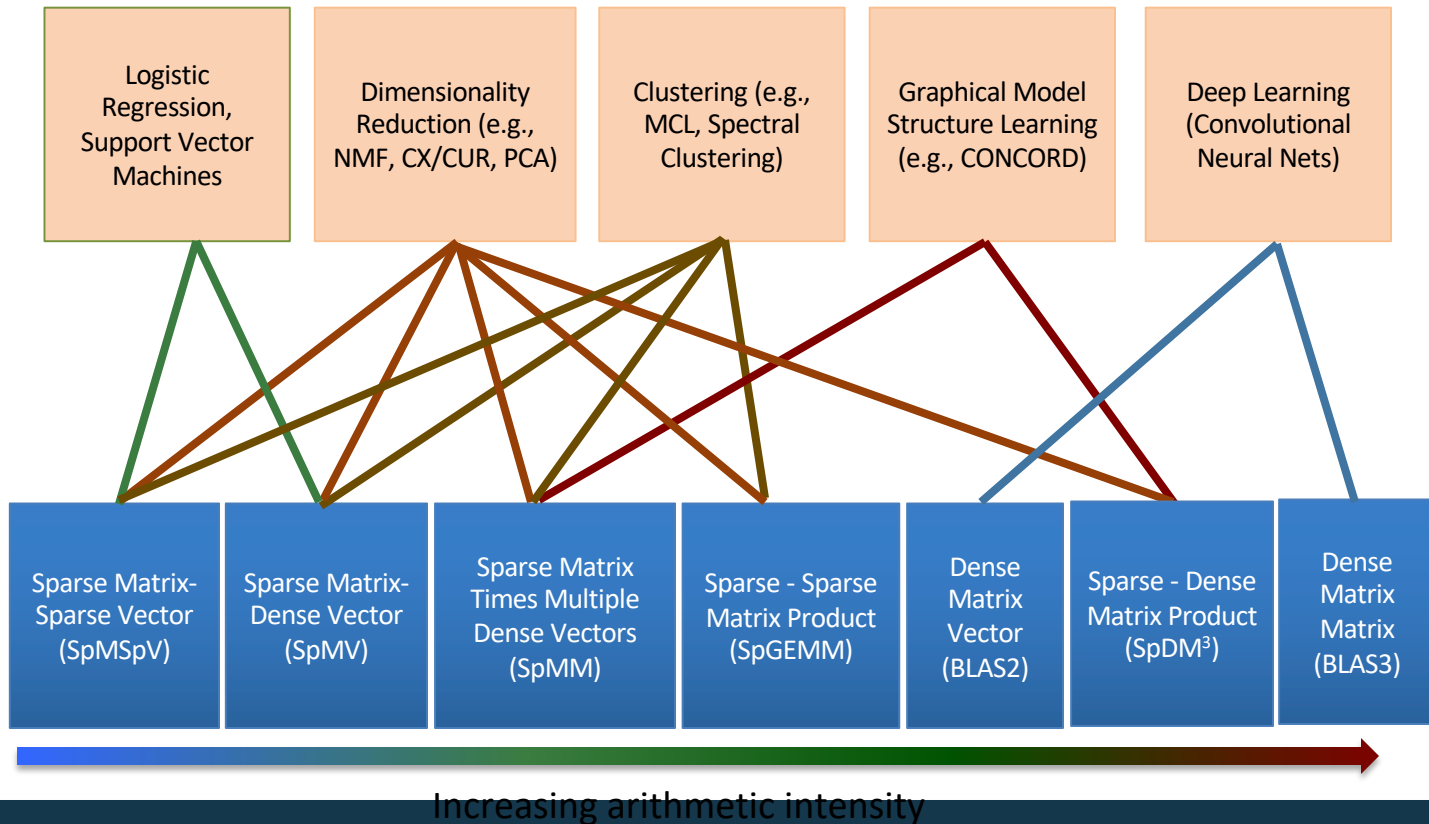
Matrix Multiplication code has a 3D iteration space
Each point in the space is a constant computation (*/+)

```
for i
  for j
    for k
      C[i,j] ... A[i,k] ... B[k,j] ...
```

Avoiding Communication in GNNs



Machine Learning Mapping to Linear Algebra



Take-Aways

- **Applications**
 - More data, more computing can reveal new insights
 - Genomics problems dominated by ~7 motifs
- **Architectures**
 - Specialization and data parallelism will be increasingly important
 - Communication will (still) dominate
 - Need better integration, lower overheads mechanism
- **Algorithms**
 - Irregular, fine-grained problems
 - Can map to distributed memory and data parallelism
 - Avoid communication:
 - Hide latency or aggregating messages (can trade off)
 - Reduce bandwidth (volume), and
 - Use all the wires all the time