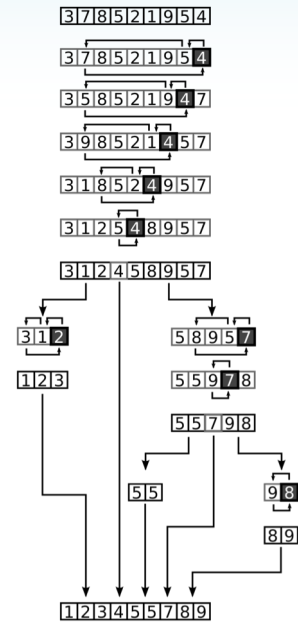


# 6.886: Algorithm Engineering



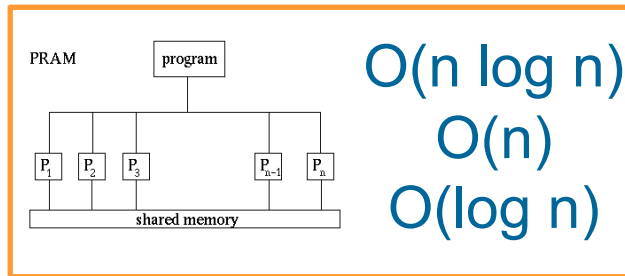
## LECTURE 1 Introduction

Julian Shun  
*February 5, 2019*

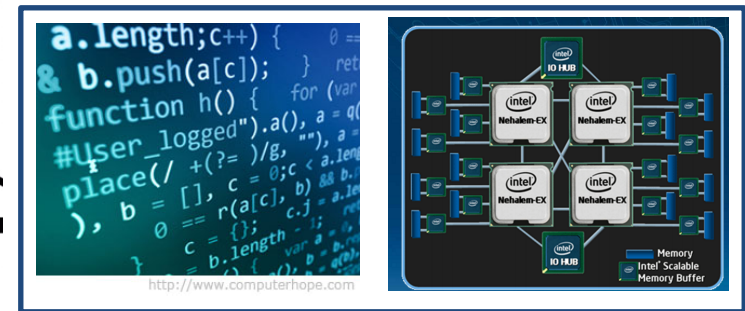
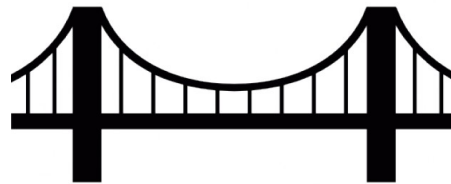


# What is Algorithm Engineering?

- Algorithm design
- Algorithm analysis
- Algorithm implementation
- Optimization
- Profiling
- Experimental evaluation

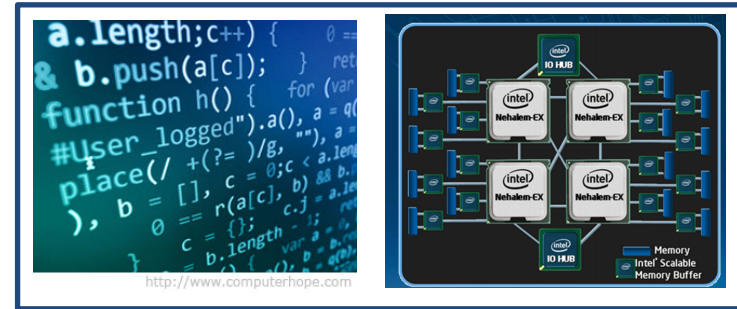
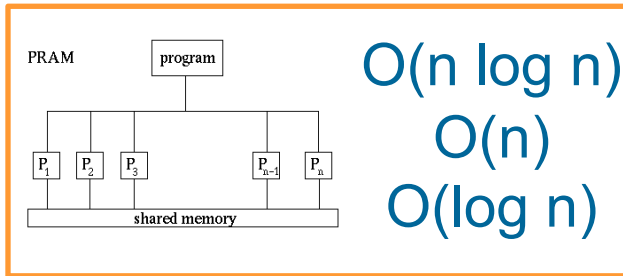


*Theory*



*Practice*

# Bridging Theory and Practice



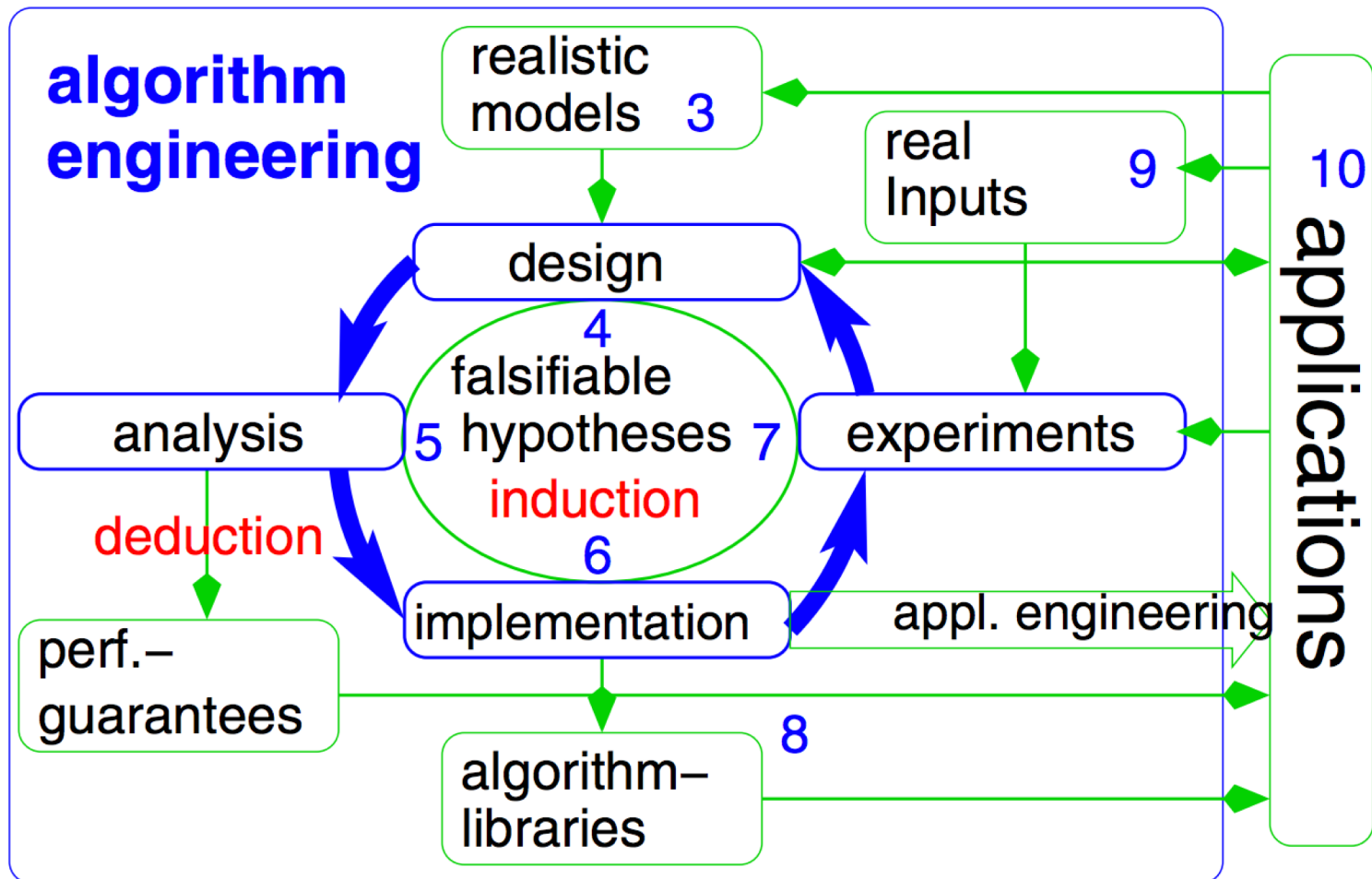
- Good empirical performance
- Confidence that algorithms will perform well in many different settings
- Ability to predict performance (e.g. in real-time applications)
- Important to develop theoretical models to capture properties of technologies

*Use theory to inform practice and practice to inform theory.*

# Brief History

- In early days, implementing algorithms designed was standard practice
- 1970s–1980s: Algorithm theory is a subdiscipline in CS mostly devoted to "paper and pencil" work
- Late 1980s–1990s: Researchers began noticing gaps between theory and practice
- 1997: First Workshop on Algorithm Engineering (WAE) by P. Italiano (now part of ESA)
- 1998: Meeting on Algorithm Engineering & Experiments (ALENEX)
- 2003: annual Workshop on Experimental Algorithms (WEA), now Symposium on Experimental Algorithms (SEA)
- Nowadays many conferences have papers on algorithm engineering

# What is Algorithm Engineering?



Source: "Algorithm Engineering – An Attempt at a Definition", Peter Sanders

# Models of Computation

- Random-Access Machine (RAM)
  - Infinite memory
  - Arithmetic operations, logical operations, and memory accesses take  $O(1)$  time
  - Most sequential algorithms are designed using this model (6.006/6.046)
- Nowadays computers are much more complex
  - Deep cache hierarchies
  - Instruction level parallelism
  - Multiple cores
  - Disk if input doesn't fit in memory

# Algorithm Design & Analysis

	<u>Algorithm 1</u>	<u>Algorithm 2</u>
Complexity	$N \log_2 N$	$1000 N$

- Constant factors matter!
- Avoid unnecessary computations
- Simplicity improves applicability and can lead to better performance
- Think about locality and parallelism
- Think both about worst-case and real-world inputs
- Use theory as a guide to find practical algorithms
- Time vs. space tradeoffs

# Implementation

- Write clean, modular code
  - Easier to experiment with different methods, and can save a lot of development time
- Write correctness checkers
  - Especially important in numerical and geometric applications due to floating–point arithmetic, possibly leading to different results
- Save previous versions of your code!
  - Version control helps with this



# Experimentation

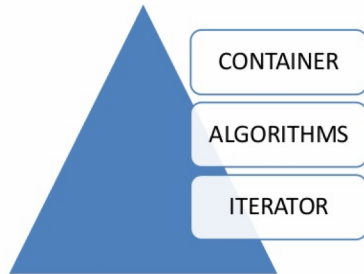
- Instrument code with timers and use performance profilers (e.g., perf, gprof, valgrind)
- Use large variety of inputs (both real-world and synthetic)
  - Use different sizes
  - Use worst-case inputs to identify correctness or performance issues
- **Reproducibility**
  - Document environmental setup
  - Fix random seeds if needed
- Run multiple timings to deal with variance

# Experimentation II

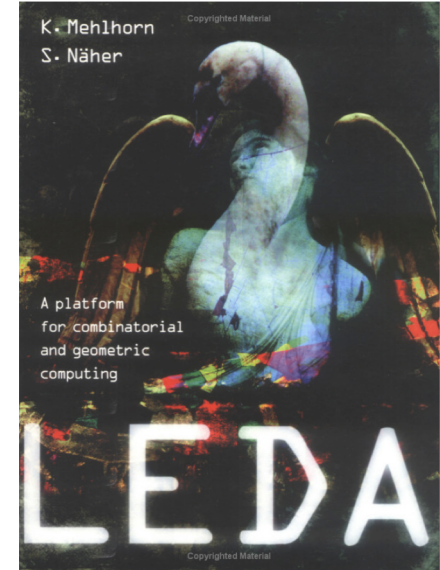
- For parallel code, test on varying number of processors to study scalability
- Compare with best serial code for problem
- For reproducibility, write deterministic code if possible
  - Or make it easy to turn off non-determinism
- Use numactl to control NUMA effects on multi-socket machines
- Useful tools: CilkScale, CilkSan

# Libraries and Frameworks

## Components of STL



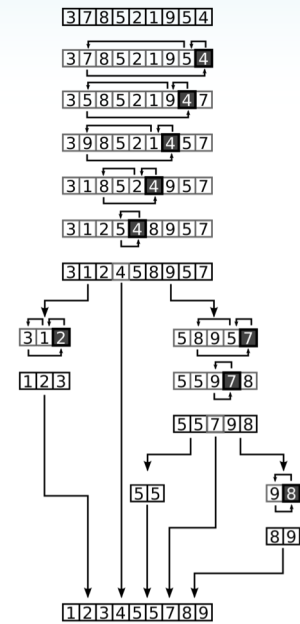
CGAL



## Problem Based Benchmark Suite

[Home](#) [Benchmarks/Code](#) [Inputs](#) [License](#) [People](#) [Publications](#)

- Use efficient building blocks from existing library/frameworks when appropriate
- Develop your own to help others and improve applicability



# COURSE INFORMATION



# Course Information

- Graduate-level class
  - Undergraduates who have taken 6.046 and 6.172 are welcome
- Lectures: Tuesday and Thursday 2:30–4pm
- Instructor: Julian Shun
- Units: 3–0–9
- We will use Piazza for communication
  
- This course will cover various ideas in algorithm engineering, with an emphasis on parallelism and graph problems

# Course Website

<https://people.csail.mit.edu/jshun/6886-s19/>

## Schedule (tentative)

Date	Topic	Speaker	Required Reading	Optional Reading
Tuesday 2/5	Course Introduction	Julian Shun	<a href="#">Algorithm Engineering - An Attempt at a Definition</a> <a href="#">A Theoretician's Guide to the Experimental Analysis of Algorithms</a>	<a href="#">Algorithm Engineering: Bridging the Gap Between Algorithm Theory and Practice</a> <a href="#">A Guide to Experimental Algorithmics</a> <a href="#">Algorithm engineering: an attempt at a definition using sorting as an example</a> <a href="#">Algorithm Engineering for Parallel Computation</a> <a href="#">Distributed Algorithm Engineering</a> <a href="#">Experimental algorithmics</a> <a href="#">Programming Pearls</a> <a href="#">Smoothed analysis of algorithms: Why the simple algorithm usually takes polynomial time</a>
Thursday 2/7	Parallel Algorithms	Julian Shun	<a href="#">Parallel Algorithms</a> <a href="#">Thinking in Parallel: Some Basic Data-Parallel Algorithms and Techniques (Chapters 4-8)</a> CLRS Chapter 27	<a href="#">Prefix Sums and Their Applications</a> <a href="#">Algorithm Design: Parallel and Sequential</a> <a href="#">Introduction to Parallel Algorithms</a>
Tuesday 2/12	Parallel Graph Traversal		<a href="#">Direction-Optimizing Breadth-First Search*</a> <a href="#">A Faster Algorithm for Betweenness Centrality</a> <a href="#">The More the Merrier: Efficient Multi-Source Graph Traversal*</a>	<a href="#">A Work-Efficient Parallel Breadth-First Search Algorithm (or How to Cope with the Nondeterminism of Reducers)</a> <a href="#">Internally Deterministic Parallel Algorithms Can Be Fast</a> <a href="#">SlimSell: A Vectorizable Graph Representation for Breadth-First Search</a> Chapter 3.6 of <a href="#">Networks, Crowds, and Markets</a> (describes Betweenness Centrality with an example) <a href="#">Better Approximation of Betweenness Centrality</a> <a href="#">ABRA: Approximating Betweenness Centrality in Static and Dynamic Graphs with Rademacher Averages</a> <a href="#">KADABRA is an ADaptive Algorithm for</a>

# Grading

Grading Breakdown	
Paper Reviews	15%
Paper Presentations	20%
Research Project	60%
Class Participation	5%

*You must complete all assignments to pass the class.*

# Paper Presentations

- This is a research-oriented course
- Cover content from 2–3 research papers each lecture
- 25–30 minute student presentation per paper
  - Discuss motivation for the problem solved
  - Key technical ideas
  - Theoretical/experimental results
  - Related work
  - Strengths/weaknesses
  - Directions for future work
  - Include several questions for discussion
  - Presentation should cover necessary background to understand paper (you may have to read related papers)
  - Make slides but may use the board for theoretical proofs
- Sign up for presentations this week in Google doc
- Would be helpful to sign up even if listening



# Paper Reviews

- Submit one paper review each week on a paper that will be covered that week
  - Cover motivation, key ideas, results, novelty, strengths/weaknesses, your ideas for improving the techniques or evaluation, any open problems or directions for further work
  - Submit on Learning Modules by Monday 11:59pm each week (before we cover the papers)
  - Reviews will be made viewable to class (anonymously)
  - Read them before the lecture to help prepare for the discussions

# Research Project

- Open-ended research project to be done in groups of 1–3 people
- Some ideas
  - Implementation of non-trivial algorithms
  - Analyzing/optimizing performance of existing algorithms
  - Designing new theoretically and/or practically efficient algorithms
  - Applying algorithms in the context of larger applications
  - Improving or designing new algorithm frameworks or libraries
  - Any topic may involve parallelism, cache-efficiency, I/O-efficiency, and memory-efficiency
- Must contain an implementation component
- Can be related to research that you are doing
- On Tuesday 3/5, you can pitch any project ideas you have and find group members

# Project Timeline

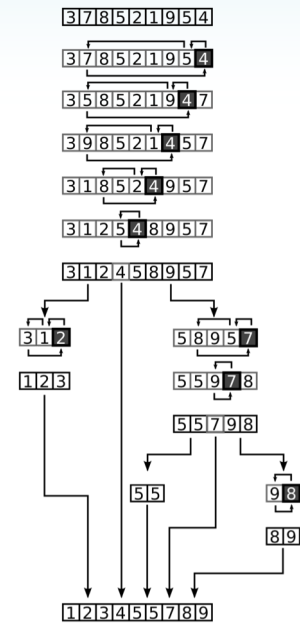
Assignment	Due Date
Project idea pitches	3/5
Pre-proposal meeting	3/12
Proposal*	3/15
Weekly progress reports	3/22, 4/5, 4/12, 4/26, 5/3, 5/10
Mid-term report*	4/19
Project presentations	5/16
Final report	5/16

- Pre-proposal meeting
  - 15-minute meeting to run idea by instructor
- Talk to instructors if you need computing resources for the project
  - We may have some AWS credits

\*You can submit it 3 days later if you go to the Comm Lab at least one day before original deadline.



# PARALLELISM



# Parallelism

*Data is becoming very large!*



41 million vertices  
1.5 billion edges  
(6.3 GB)

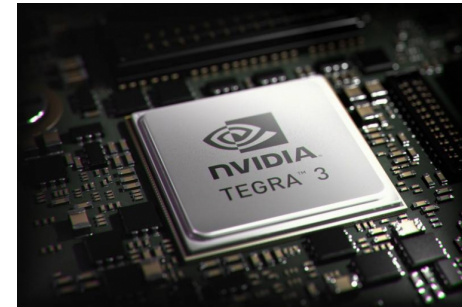


1.4 billion vertices  
6.6 billion edges  
(38 GB)



3.5 billion vertices  
128 billion edges  
(540 GB)

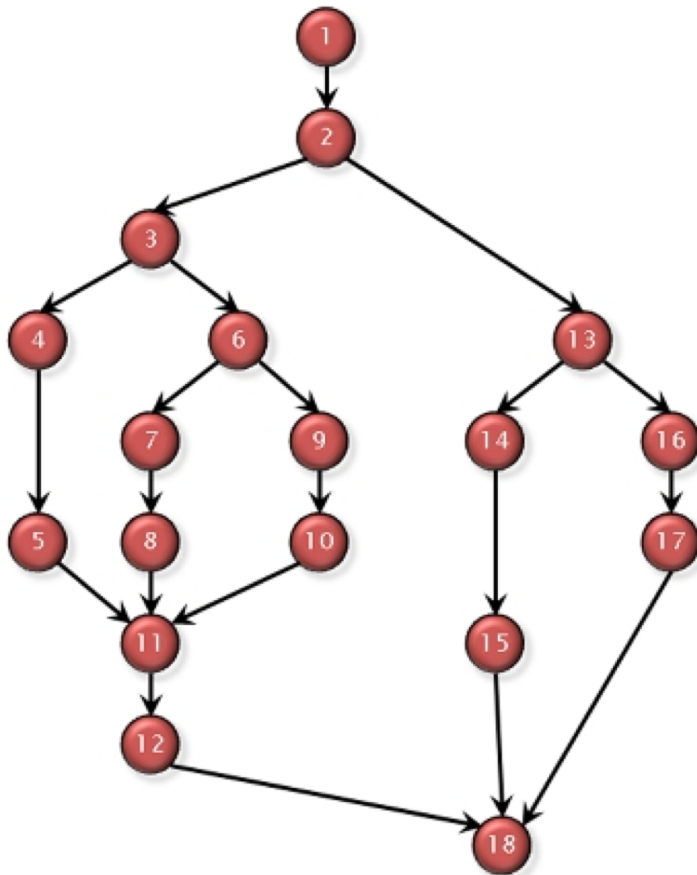
*Parallel machines are everywhere!*



*Can rent machines on AWS with 72 cores  
(144 hyper-threads) and 4TB of RAM*

# Parallelism Models

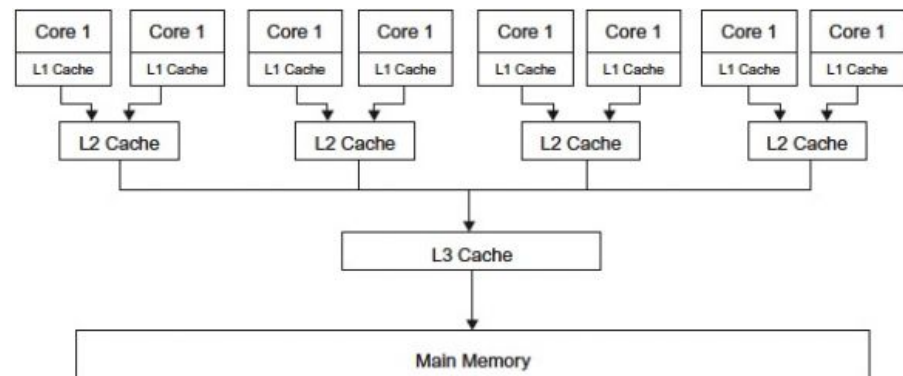
## Computation graph

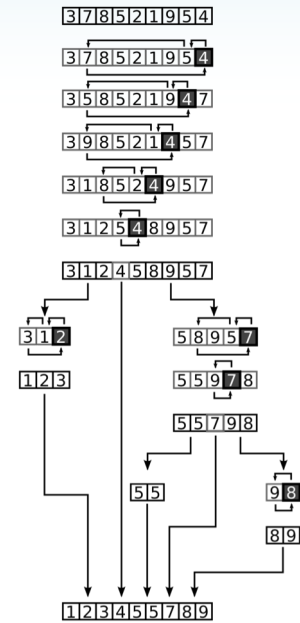


- **Work** = number of vertices in graph (number of operations)
- **Depth (Span)** = longest directed path in graph (dependence length)
- **Parallelism** =  $Work / Depth$

Goal 1: work-efficient and low (polylogarithmic) depth algorithms

Goal 2: simple, practical, and cache-friendly

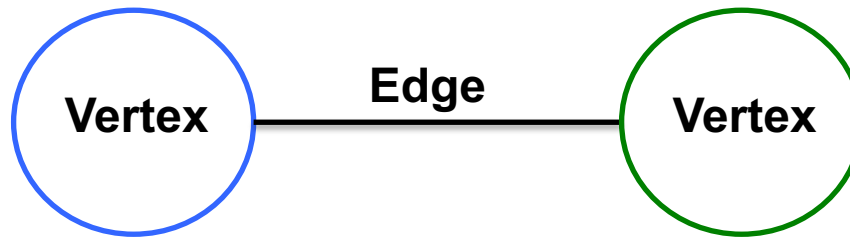




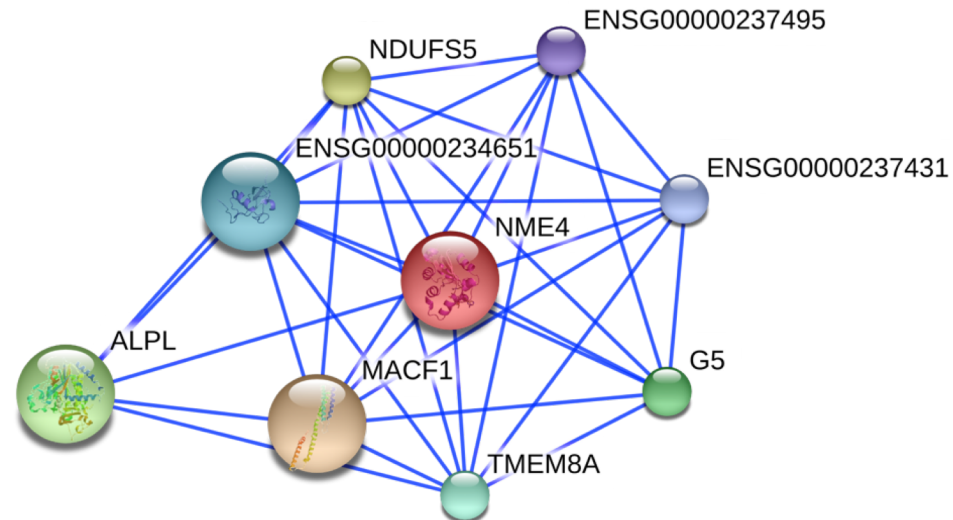
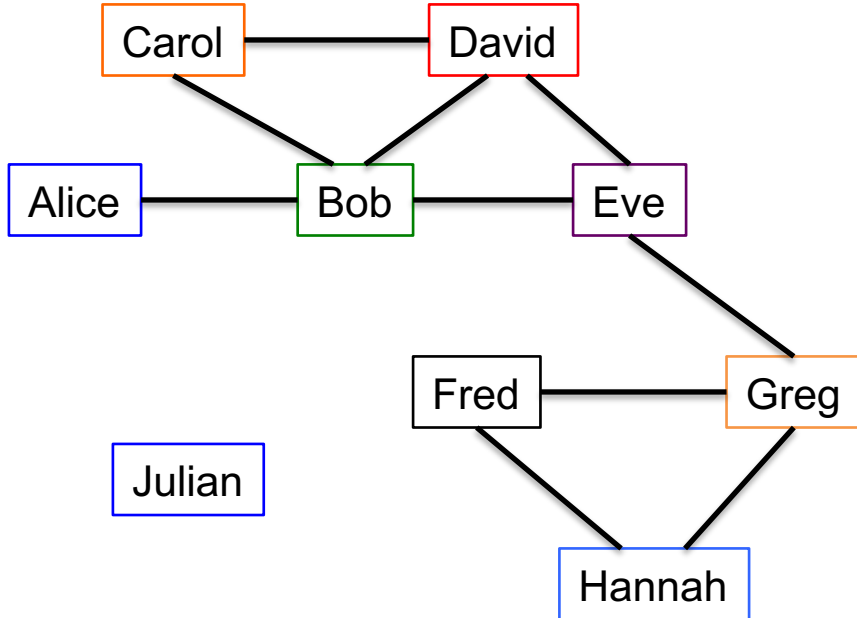
# GRAPHS



# What is a graph?



- Vertices model objects
- Edges model relationships between objects



[https://commons.wikimedia.org/wiki/File:Protein\\_Interaction\\_Network\\_for\\_TMEM8A.png](https://commons.wikimedia.org/wiki/File:Protein_Interaction_Network_for_TMEM8A.png)



# Graph Representations

- Vertices labeled from 0 to  $n-1$

	0	1	2	3	4	
0	0	1	0	0	0	(0,1)
1	1	0	0	1	1	(1,0) (1,3) (1,4)
2	0	0	0	1	0	(2,3)
3	0	1	1	0	0	(3,1)
4	0	1	0	0	0	(3,2) (4,1)

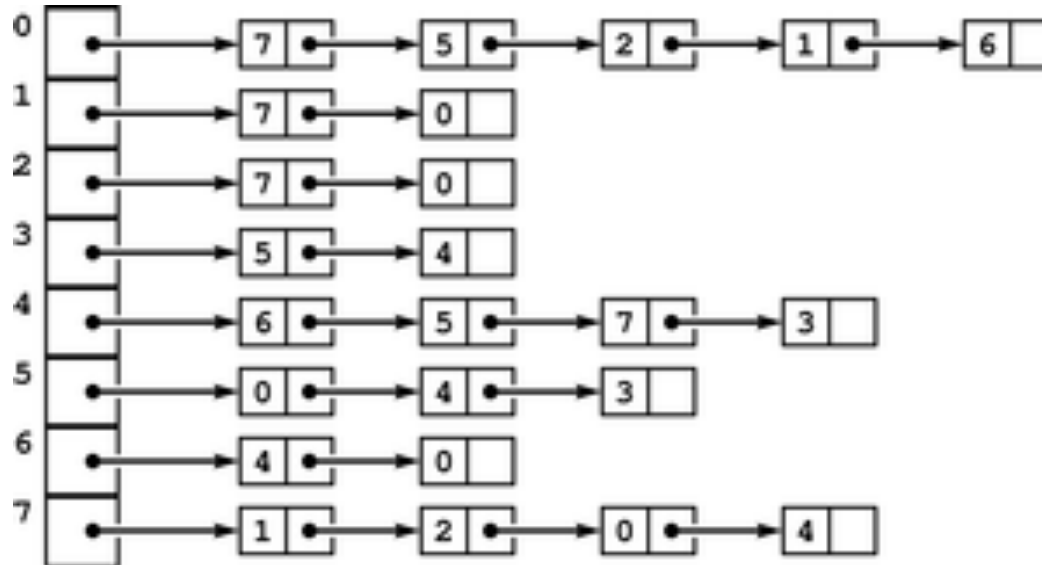
Adjacency matrix  
("1" if edge exists,  
"0" otherwise)

Edge list

- $O(n^2)$  space for adjacency matrix
- $O(m)$  space for edge list

# Graph Representations

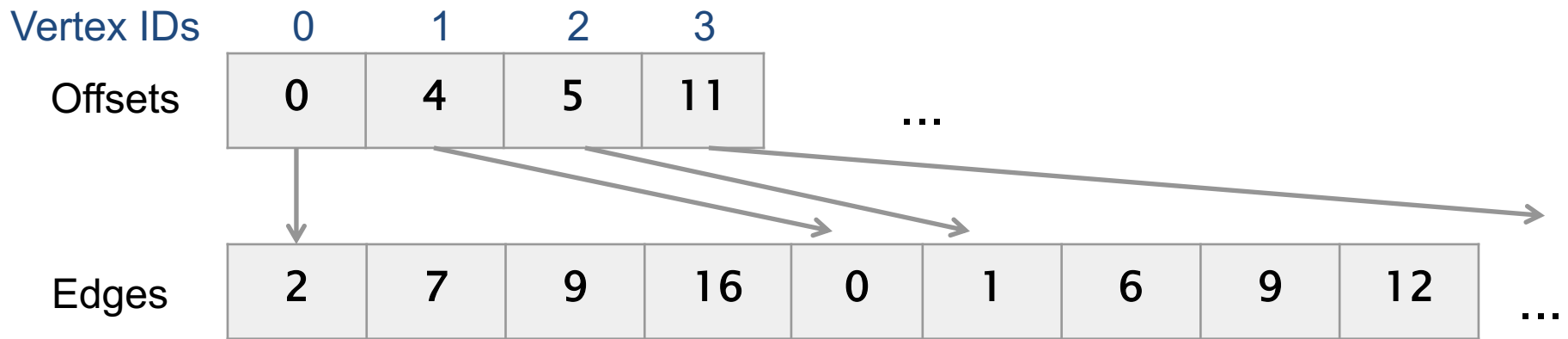
- Adjacency list
  - Array of pointers (one per vertex)
  - Each vertex has an unordered list of its edges



- Space requirement is  $O(n+m)$
- Can substitute linked lists with arrays for better cache performance
  - Tradeoff: more expensive to update graph

# Graph Representations

- Compressed sparse row (CSR)
  - Two arrays: **Offsets** and **Edges**
  - **Offsets**[i] stores the offset of where vertex i's edges start in **Edges**



- How do we know the degree of a vertex?
- Space usage is  $O(n+m)$
- Can also store values on the edges with an additional array or interleaved with **Edges**

# Tradeoffs in Graph Representations

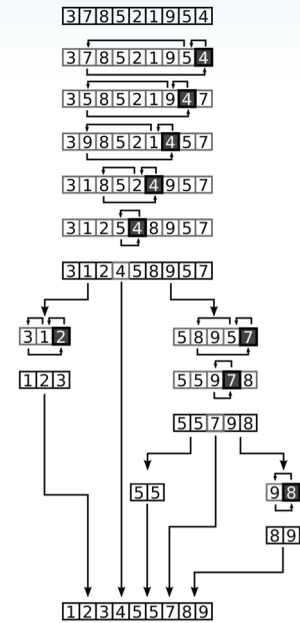
- What is the cost of different operations?

	Adjacency matrix	Edge list	Adjacency list (linked list)	Compressed sparse row
Storage cost / scanning whole graph	$O(n^2)$	$O(m)$	$O(m+n)$	$O(m+n)$
Add edge	$O(1)$	$O(1)$	$O(1)$	$O(m+n)$
Delete edge from vertex $v$	$O(1)$	$O(m)$	$O(\text{deg}(v))$	$O(m+n)$
Finding all neighbors of a vertex $v$	$O(n)$	$O(m)$	$O(\text{deg}(v))$	$O(\text{deg}(v))$
Finding if $w$ is a neighbor of $v$	$O(1)$	$O(m)$	$O(\text{deg}(v))$	$O(\text{deg}(v))$

- There are variants/combinations of these representations



# BREADTH-FIRST SEARCH



# Breadth-First Search (BFS)

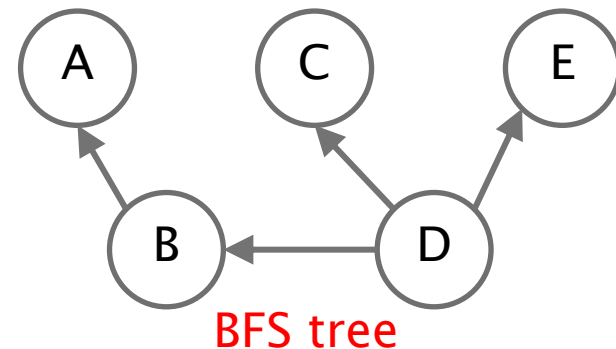
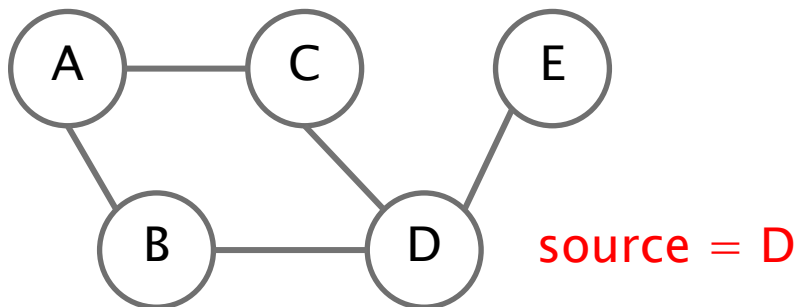
- Given a source vertex  $s$ , visit the vertices in order of distance from  $s$
- Possible outputs:

- Vertices in the order they were visited
  - D, B, C, E, A
- The distance from each vertex to  $s$

A	B	C	D	E
2	1	1	0	1

- A BFS tree, where each vertex has a parent to a neighbor in the previous level

Applications
Betweenness centrality
Eccentricity estimation
Maximum flow
Web crawlers
Network broadcasting
Cycle detection
...



# Sequential BFS Algorithm

```
Breadth-First-Search(Graph, root):
```

```
  for each node n in Graph:  
    n.distance = INFINITY  
    n.parent = NIL
```

Source: [https://en.wikipedia.org/wiki/Breadth-first\\_search](https://en.wikipedia.org/wiki/Breadth-first_search)

- BFS requires  $O(n+m)$  work on  $n$  vertices and  $m$  edges

# Sequential BFS Algorithm

- Assume graph is given in compressed sparse row format
  - Two arrays: **Offsets** and **Edges**
  - $n$  vertices and  $m$  edges (assume  $\text{Offsets}[n] = m$ )

```
int* parent =
  (int*) malloc(sizeof(int)*n);
int* queue =
  (int*) malloc(sizeof(int)*n);

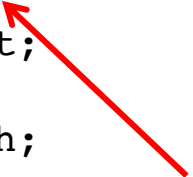
for(int i=0; i<n; i++) {
  parent[i] = -1;
}

queue[0] = source;
parent[source] = source;

int q_front = 0, q_back = 1;

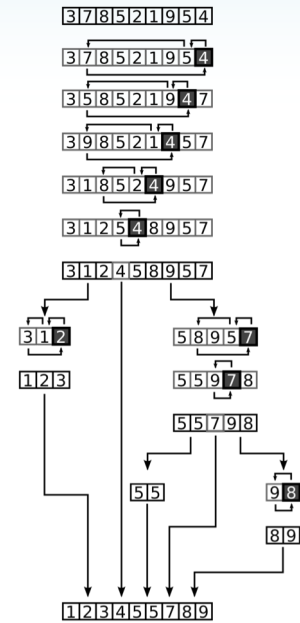
//while queue not empty
while(q_front != q_back) {
  int current = queue[q_front++]; //dequeue
  int degree =
    Offsets[current+1]-Offsets[current];
  for(int i=0; i<degree; i++) {
    int ngh = Edges[Offsets[current]+i];
    //check if neighbor has been visited
    if(parent[ngh] == -1) {
      parent[ngh] = current;
      //enqueue neighbor
      queue[q_back++] = ngh;
    }
  }
}
```

Total of  $m$   
random accesses



- What is the most expensive part of the code?
  - Random accesses cost more than sequential accesses





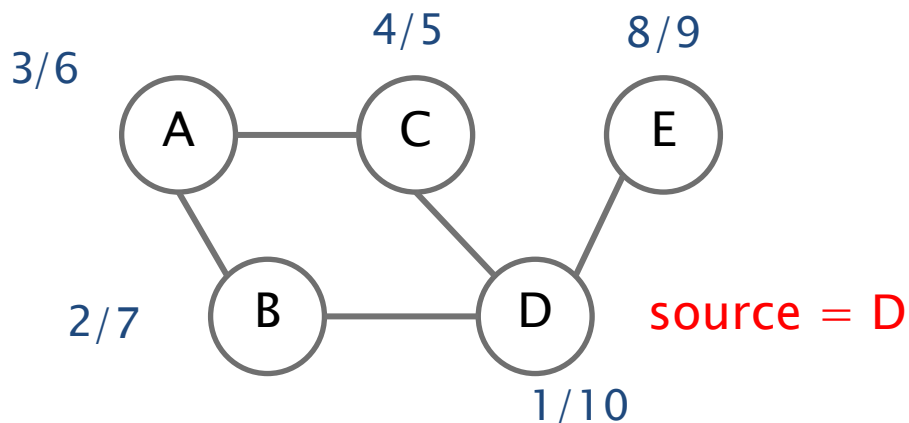
# DEPTH-FIRST SEARCH



# Depth-First Search (DFS)

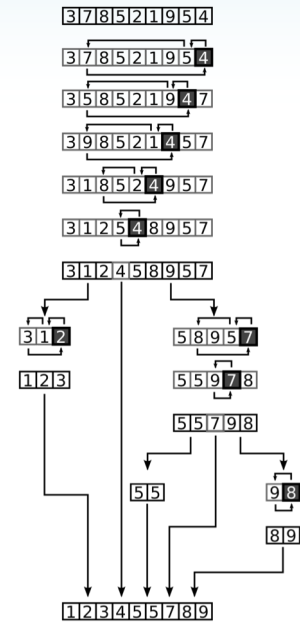
- Explores edges out of the most recently discovered vertex
- Possible outputs:
  - Depth-first forest
  - Vertices in the order they were first visited (preordering)
  - Vertices in the order they were last visited (postordering)
  - Reverse postordering

Applications
Topological sort
Solving mazes
Biconnected components
Strongly connected components
Cycle detection
...



Preorder: D, B, A, C, E  
Postorder: C, A, B, E, D  
Reverse postorder: D, E, B, A, C

*DFS requires  $O(n+m)$  work on  $n$  vertices and  $m$  edges*

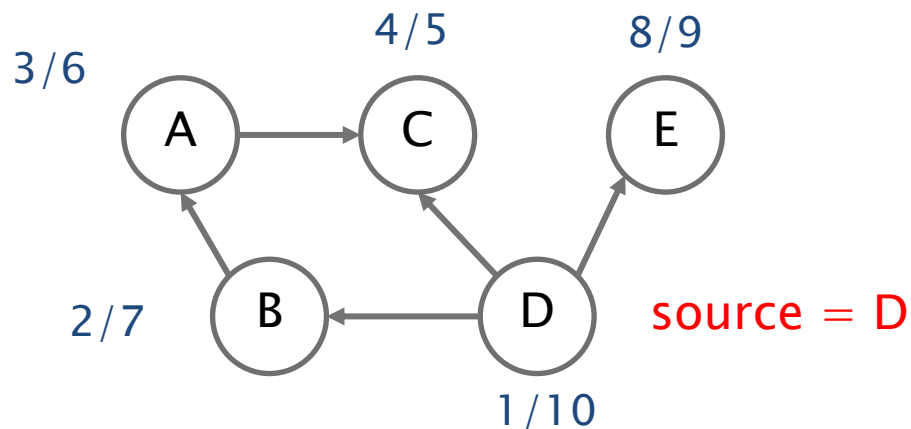


# TOPOLOGICAL SORT



# Topological Sort

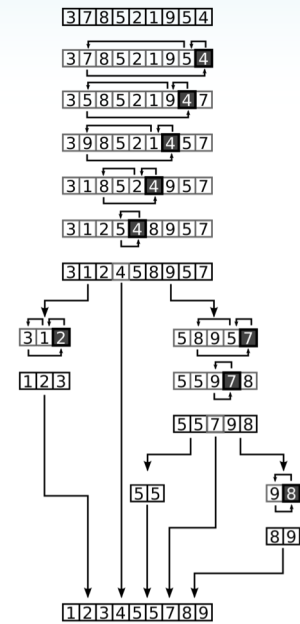
- Given a directed acyclic graph, output the vertices in an order such that all predecessors of a vertex appear before it
  - Application: scheduling tasks with dependencies (e.g. parallel computing, Makefile)
- Solution: output vertices in reverse postorder in DFS



Reverse postorder: D, E, B, A, C



# SHORTEST PATHS



# Single-Source Shortest Paths

- Given a weighted graph and a source vertex, output the distance from the source vertex to every vertex
- Non-negative weights
  - Dijkstra's algorithm
  - $O(m + n \log n)$  work using Fibonacci heap
- General weights
  - Bellman-Ford algorithm
  - $O(mn)$  work

# Dijkstra's Algorithm

```
1 function Dijkstra(Graph, source):  
2     dist[source] ← 0                                // Initialization  
3  
4     create vertex set Q  
5
```

- $O((m+n)\log n)$  work using normal heap
- $O(m + n\log n)$  work using Fibonacci heap
  - Extract-min takes  $O(\log n)$  work but decreasing priority only takes  $O(1)$  work (amortized)

# Bellman-Ford Algorithm

Bellman-Ford( $G$ , source):

ShortestPaths =  $\{\infty, \infty, \dots, \infty\}$  //size  $n$ ; stores shortest path distances

ShortestPaths[source] = 0

for  $i=1$  to  $n-1$ :

    for each vertex  $v$  in  $G$ :

        for each  $w$  in neighbors( $v$ ):

            if(ShortestPaths[ $v$ ] + weight( $v,w$ ) < ShortestPaths[ $w$ ]):

                ShortestPaths[ $w$ ] = ShortestPaths[ $v$ ] + weight( $v,w$ )

    if no shortest paths changed:

        return ShortestPaths

report “negative cycle”

- At most  $n$  rounds, each doing  $O(n+m)$  work
- Total work =  $O(mn)$



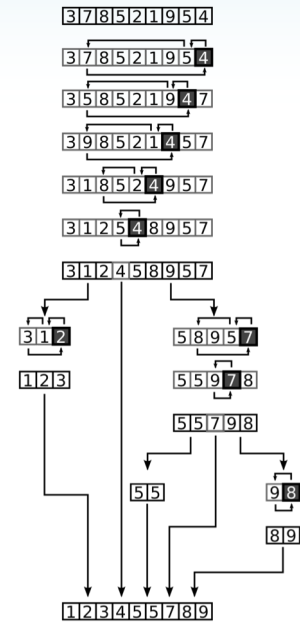
# More Graph Algorithms

- We will study algorithms for particular problems
  - Parallelism, cache-efficiency, I/O-efficiency, dynamic updates

Breadth-first search	Betweenness centrality
PageRank	Triangle Computations
Low-diameter decomposition	SSSP
Connected components	Maximal independent set
K-core decomposition	Multi-BFS
Minimum spanning forest	Spanning forest
Maximal matching	Set cover
Eccentricity estimation	Subgraph matching



# GRAPH PROCESSING FRAMEWORKS



# Graph Processing Frameworks

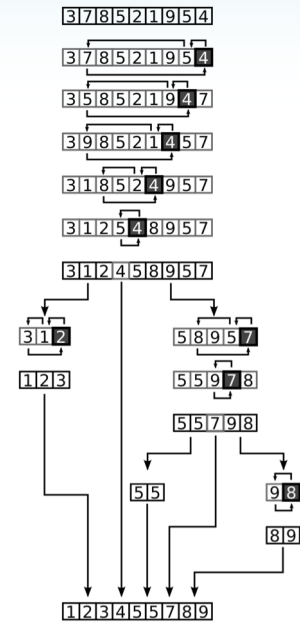
- Provides high level primitives for graph algorithms
- Reduce programming effort of writing efficient parallel graph programs

## Graph processing frameworks/libraries

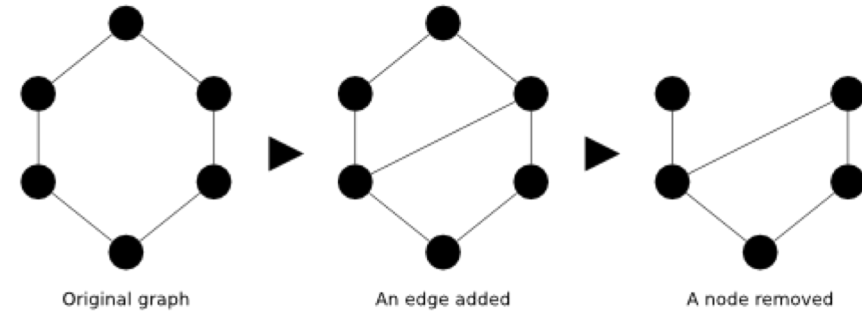
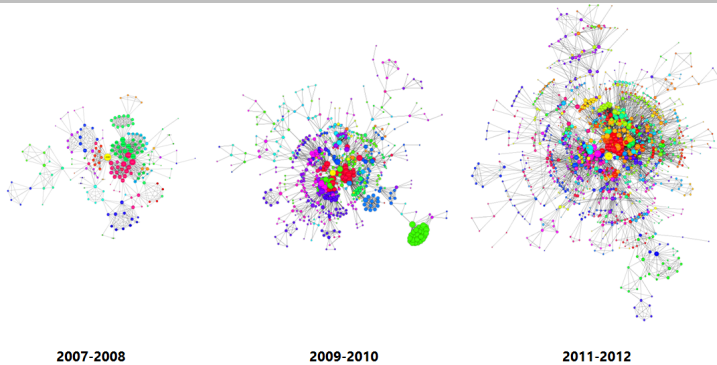
Pregel, Giraph, GPS, GraphLab, PowerGraph, PRISM, Pegasus, Knowledge Discovery Toolbox, CombBLAS, GraphChi, GraphX, Galois, X-Stream, Gunrock, GraphMat, Ringo, TurboGraph, TurboGraph++, FlashGraph, Grace, PathGraph, Polymer, GPSA, GoFFish, Blogel, LightGraph, MapGraph, PowerLyra, PowerSwitch, Imitator, XDGP, Signal/Collect, PrefEdge, EmptyHeaded, Gemini, Wukong, Parallel BGL, KLA, Grappa, Chronos, Green-Marl, GraphHP, P++, LLAMA, Venus, Cyclops, Medusa, NScale, Neo4J, Trinity, GBase, HyperGraphDB, Horton, GSPARQL, Titan, ZipG, Cagra, Milk, Ligra, Ligra+, Julienne, GraphPad, Mosaic, BigSparse, Graphene, Mizan, Green-Marl, PGX, PGX.D, Wukong+S, Stinger, cuStinger, Distinguer, Hornet, GraphIn, Tornado, Bagel, KickStarter, Naiad, Kineograph, GraphMap, Presto, Cube, Giraph++, Photon, TuX2, GRAPE, GraM, Congra, MTGL, GridGraph, NXgraph, Chaos, Mmap, Clip, Floe, GraphGrind, DualSim, ScaleMine, Arabesque, GraMi, SAHAD, Facebook TAO, Weaver, G-SQL, G-SPARQL, gStore, Horton+, S2RDF, Quegel, EAGRE, Shape, RDF-3X, CuSha, Garaph, Totem, GTS, Frog, GBTL-CUDA, Graphulo, Zorro, Coral, GraphTau, Wonderland, GraphP, GraphIt, GraPu, GraphJet, ImmortalGraph, LA3, CellIQ, AsyncStripe, Cgraph, GraphD, GraphH, ASAP, RStream, and many others...



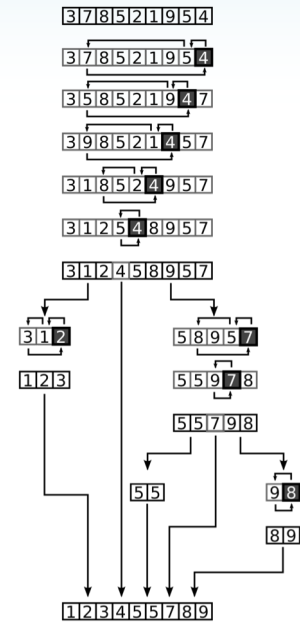
# DYNAMIC GRAPHS



# Dynamic Graphs



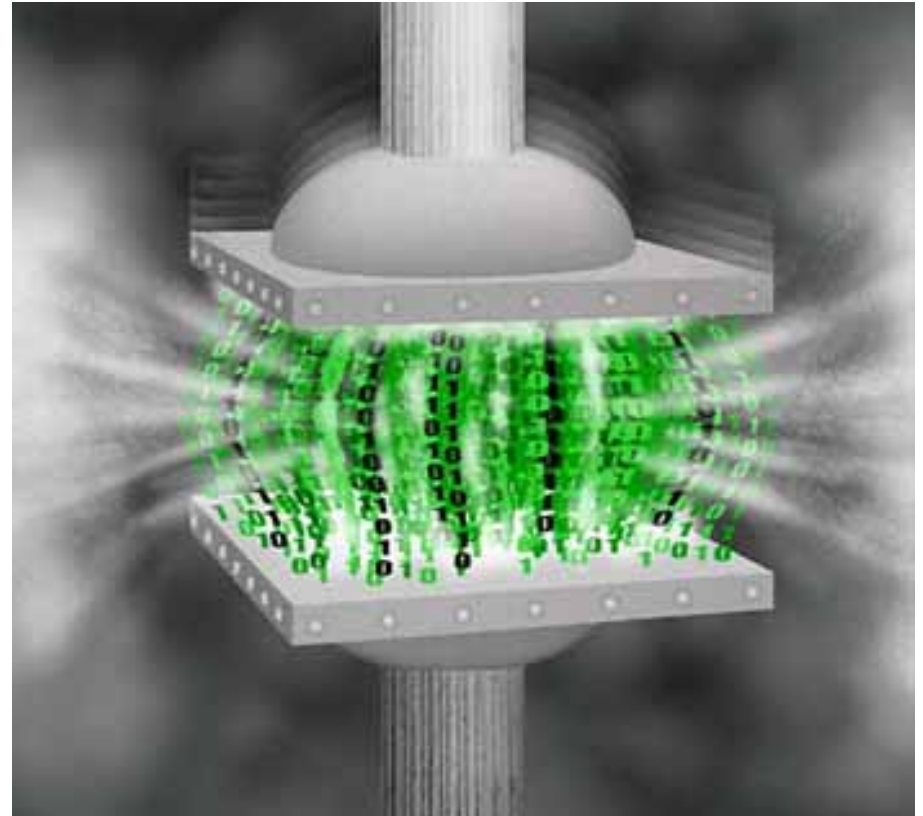
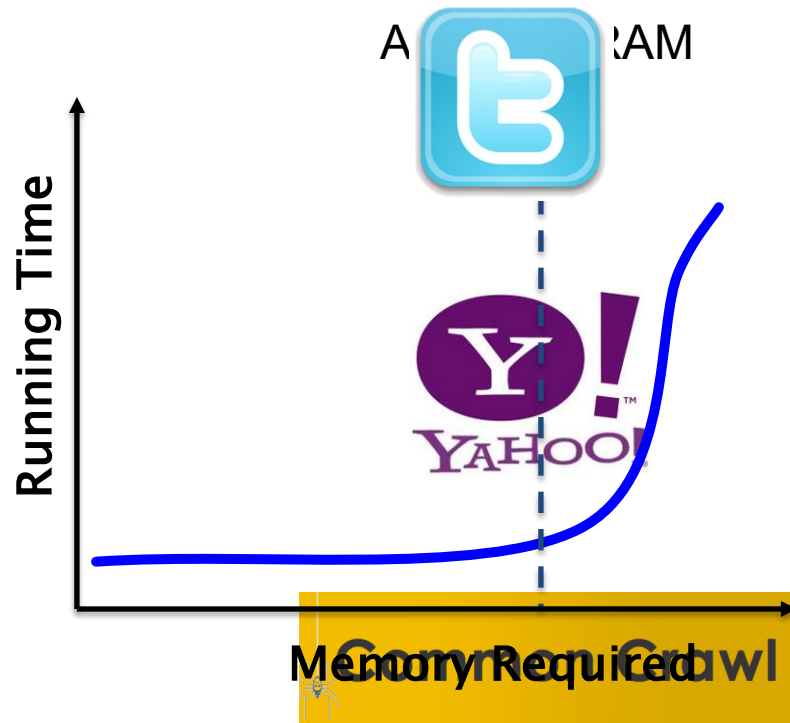
- Many graphs are changing over time
  - Adding/deleting connections on social networks
  - Traffic conditions changing
  - Communication networks (email, IMs)
  - World Wide Web
  - Content sharing (Youtube, Flickr, Pinterest)
- Need graph data structures that allow for efficient updates (in parallel)
- Need (parallel) algorithms that respond to changes without re-computing from scratch



# COMPRESSION



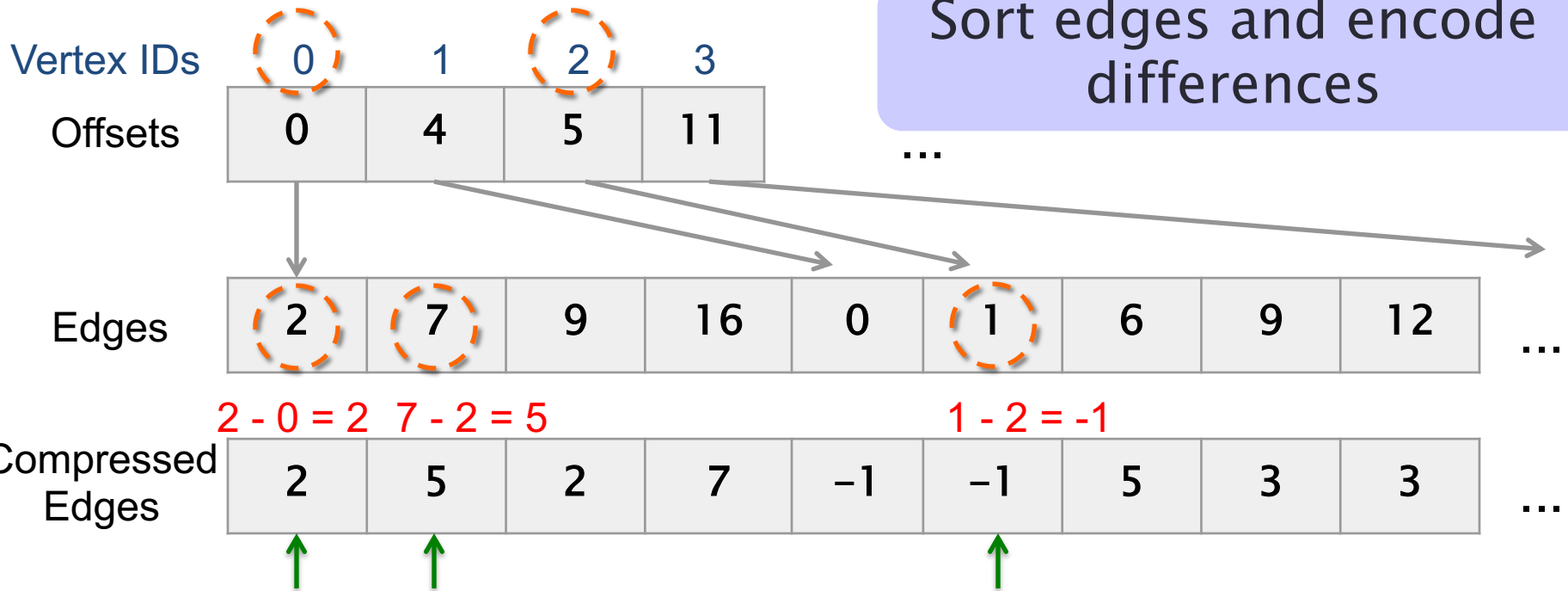
# Large Graphs



- What if you cannot fit a graph on your machine?
- Cost of machines increases with memory size

*Graph Compression*

# Graph Compression on CSR



- For each vertex  $v$ :
  - First edge: difference is  $\text{Edges}[\text{Offsets}[v]] - v$
  - $i$ 'th edge ( $i > 1$ ): difference is  $\text{Edges}[\text{Offsets}[v] + i] - \text{Edges}[\text{Offsets}[v] + i - 1]$
- Want to use fewer than 32 or 64 bits per value
- Compression can improve running time

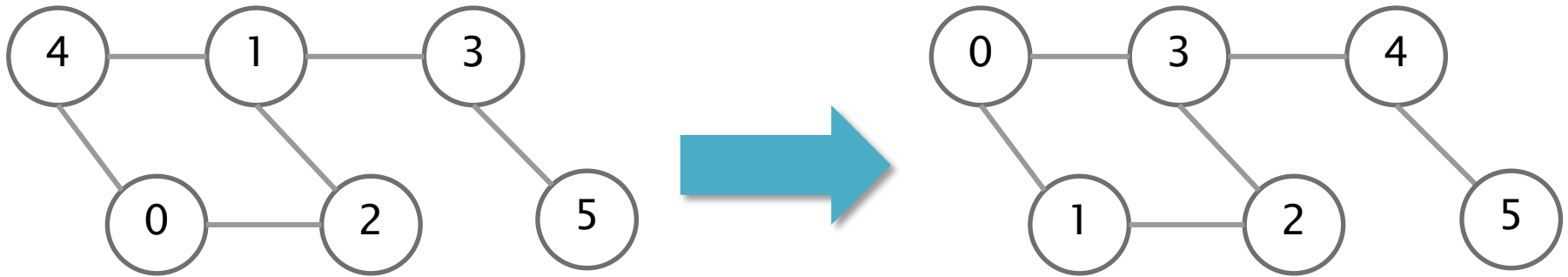


# Fast Compression Schemes

- Study speed and space tradeoffs in compression schemes for integer sequences
- Study how compression has been used to speed up sparse matrix–vector multiplication and graph processing
- Also useful in storing inverted lists for information retrieval

# Graph Reordering

- Reassign IDs to vertices to improve locality
  - Goal: Make vertex IDs close to their neighbors' IDs and neighbors' IDs close to each other



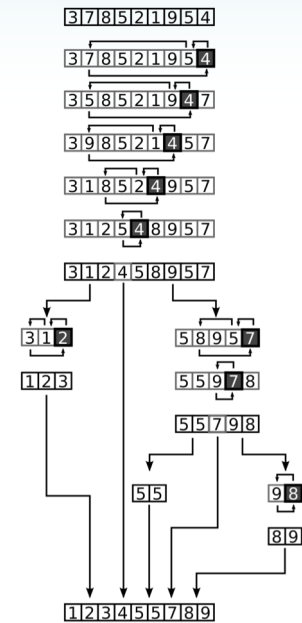
Sum of differences = 21

Sum of differences = 19

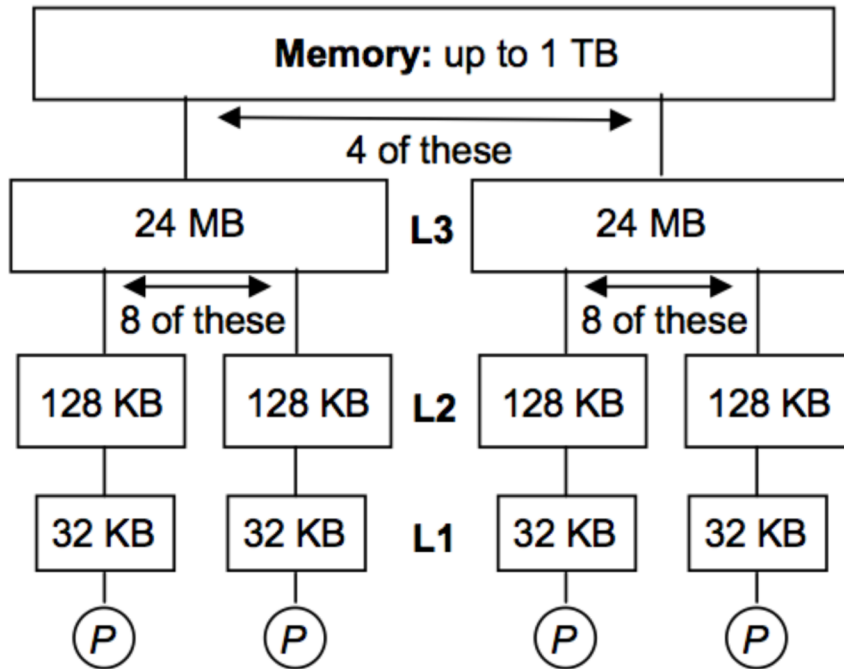
- Can improve compression rate due to smaller “differences”
- Can improve performance due to higher cache hit rate
- Various methods: BFS, DFS, METIS, degree, etc.



# CACHING AND NON-UNIFORM MEMORY ACCESS



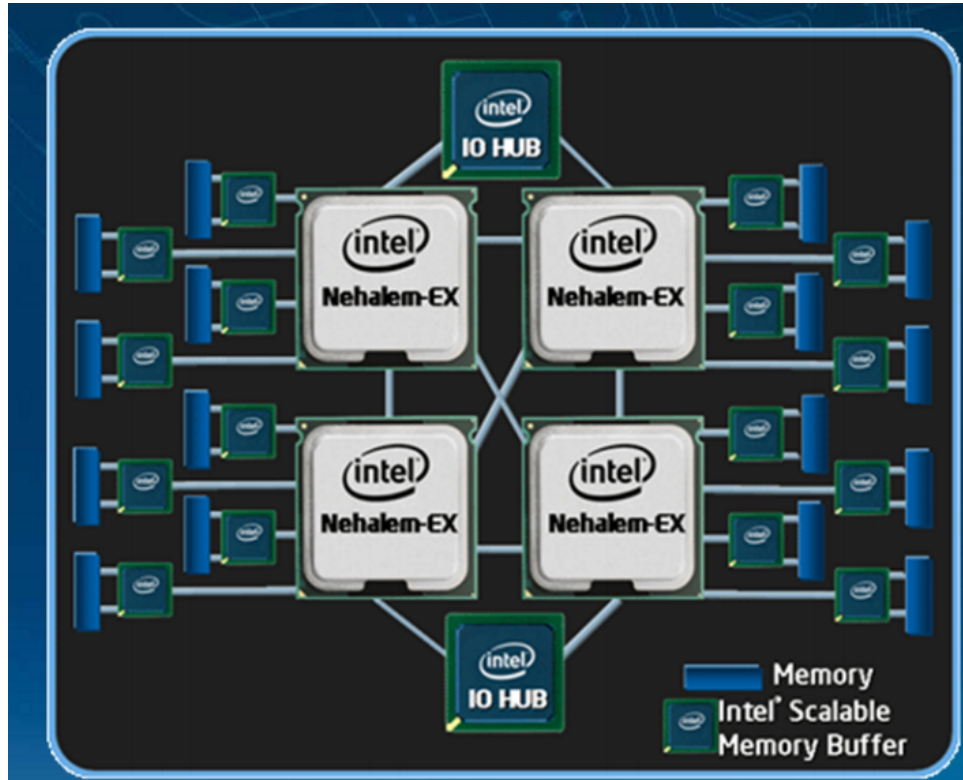
# Cache Hierarchies



Design cache-efficient and cache-oblivious algorithms to improve locality

Memory level	Approx latency
L1 Cache	1-2ns
L2 Cache	3-5ns
L3 cache	12-40ns
DRAM	60-100ns

# Non-uniform Memory Access (NUMA)

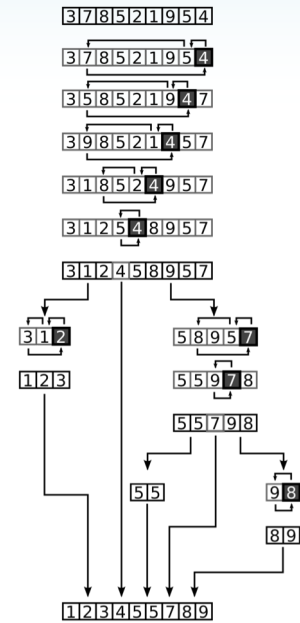


Design NUMA-aware algorithms to improve locality

- Accessing remote memory is more expensive than accessing local memory of a socket
  - Latency depends on the number of hops



# I/O EFFICIENCY



# I/O Efficiency



- Need to read input from disk at least once
- Need to read many more times if input doesn't fit in memory

Memory	Latency	Throughput
DRAM	60–100 ns	Tens of GB/s
SSD	Tens of $\mu$ s	500 MB–2 GB/s (seq), 50–200 MB/s (rand)
HDD	Tens of ms	200 MB/s (seq), 1 MB/s (rand)

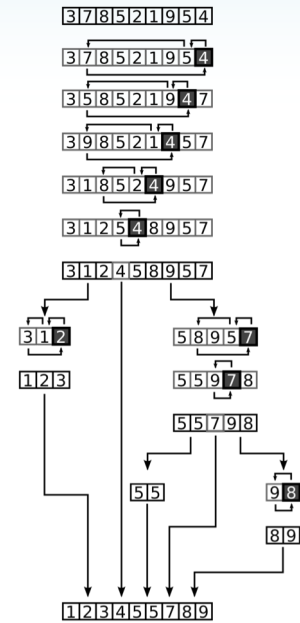
# I/O Efficiency

- For graphs larger than main memory, disk-based computing can be competitive with distributed clusters
- GraphChi: Large-Scale Graph Computation on Just a PC (OSDI 2012)

Application & Graph	Iter.	Comparative result	GraphChi (Mac Mini)	Ref
Pagerank & domain	3	GraphLab[30] on AMD server (8 CPUs) <b>87 s</b>	<b>132 s</b>	-
Pagerank & twitter-2010	5	Spark [45] with 50 nodes (100 CPUs): <b>486.6 s</b>	<b>790 s</b>	[38]
Pagerank & V=105M, E=3.7B	100	Stanford GPS, 30 EC2 nodes (60 virt. cores), <b>144 min</b>	approx. <b>581 min</b>	[37]
Pagerank & V=1.0B, E=18.5B	1	Piccolo, 100 EC2 instances (200 cores) <b>70 s</b>	approx. <b>26 min</b>	[36]
Webgraph-BP & yahoo-web	1	Pegasus (Hadoop) on 100 machines: <b>22 min</b>	<b>27 min</b>	[22]
ALS & netflix-mm, D=20	10	GraphLab on AMD server: <b>4.7 min</b>	<b>9.8 min</b> (in-mem) <b>40 min</b> (edge-repl.)	[30]
Triangle-count & twitter-2010	-	Hadoop, 1636 nodes: <b>423 min</b>	<b>60 min</b>	[39]
Pagerank & twitter-2010	1	PowerGraph, 64 x 8 cores: <b>3.6 s</b>	<b>158 s</b>	[20]
Triange-count & twitter- 2010	-	PowerGraph, 64 x 8 cores: <b>1.5 min</b>	<b>60 min</b>	[20]

- Lots of follow-up work on disk-based computing that we will study
- External-memory algorithms to minimize I/O's

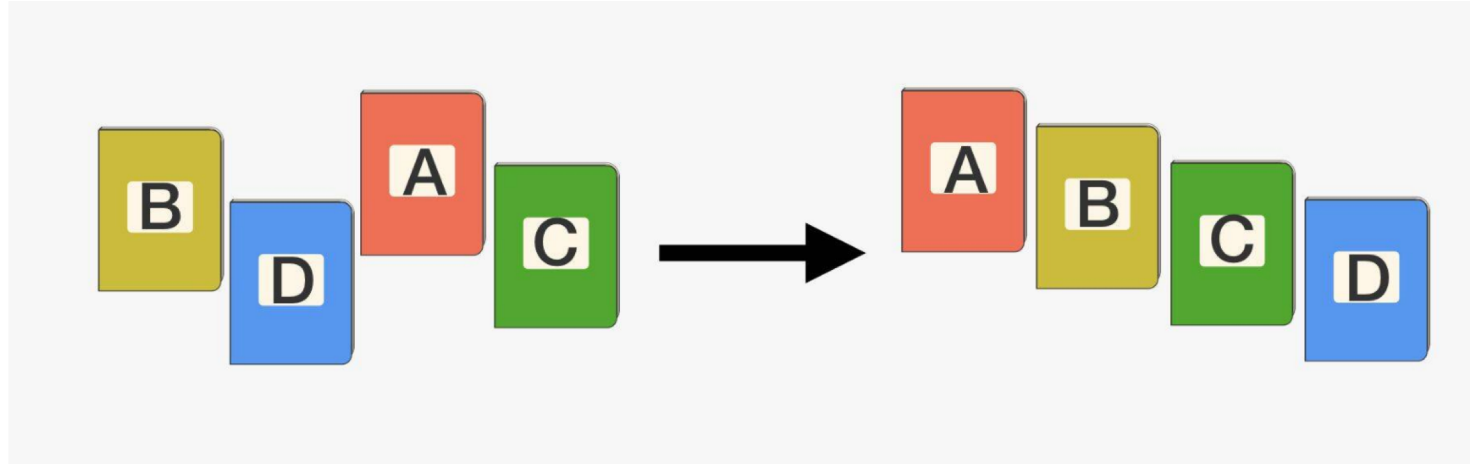




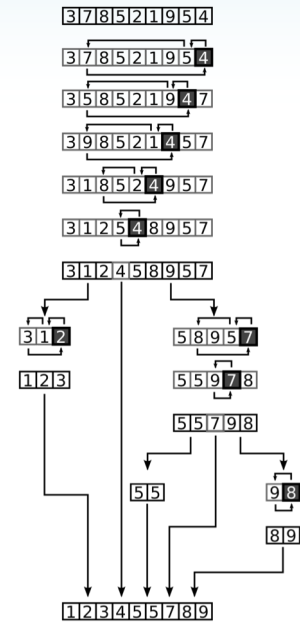
# SORTING ALGORITHMS



# Sorting



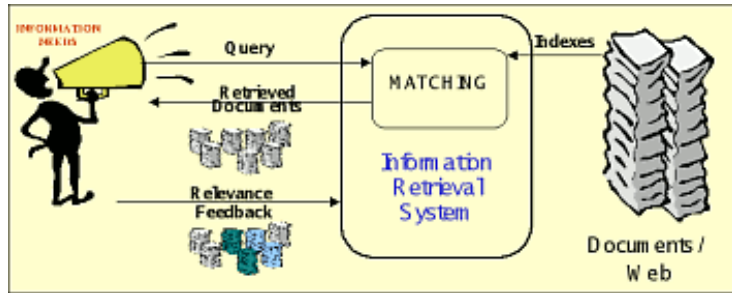
- Lots of research on engineering sorting algorithms
- Will study parallel comparison sorting and radix sorting algorithms
- <http://sortbenchmark.org/>



# STRING ALGORITHMS

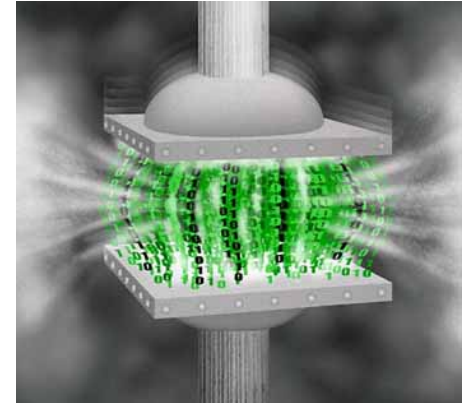


# String Algorithms

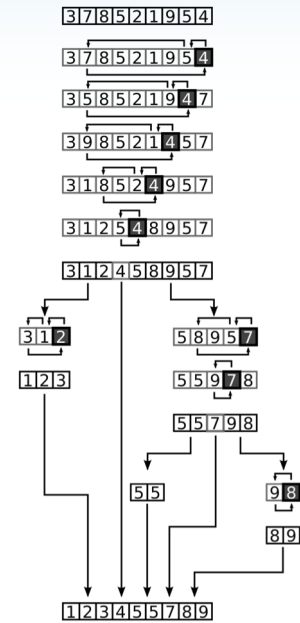


A	B	C	D	E	F	G	H	I	J	K	L	M
ROT13	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
N	O	P	Q	R	S	T	U	V	W	X	Y	Z

H	E	L	L	O
ROT13	↓	↓	↓	↓
U	R	Y	Y	B



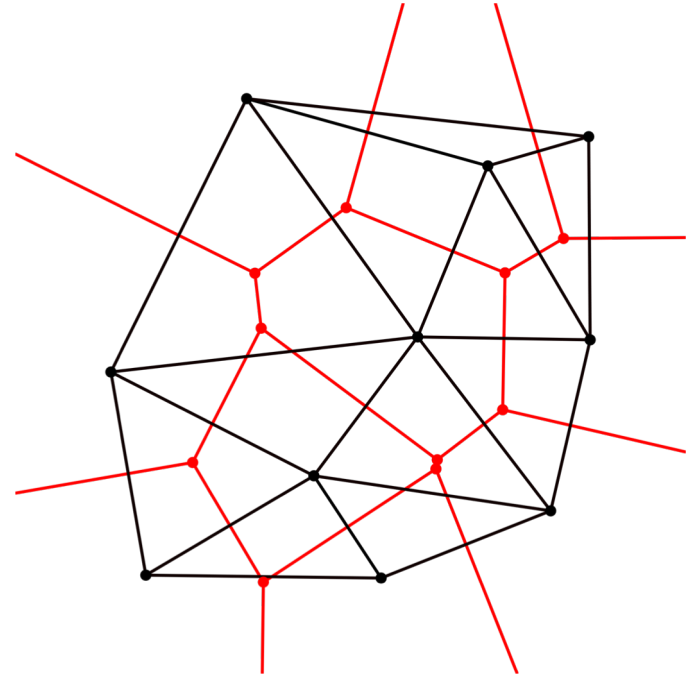
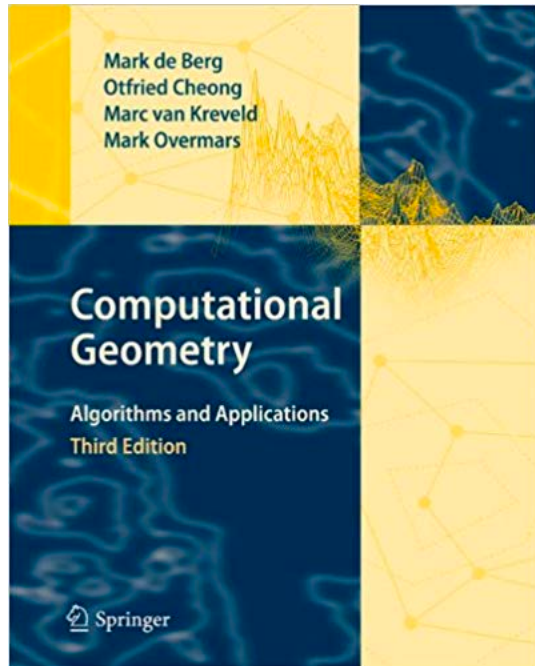
- We will study algorithms for efficiently constructing suffix arrays and suffix trees
- Many other interesting problems (edit distance, Lempel–Ziv compression, approximate string matching, alignment, etc.)



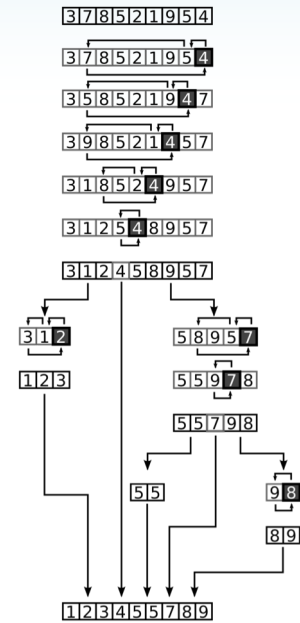
# GEOMETRY ALGORITHMS



# Computational Geometry



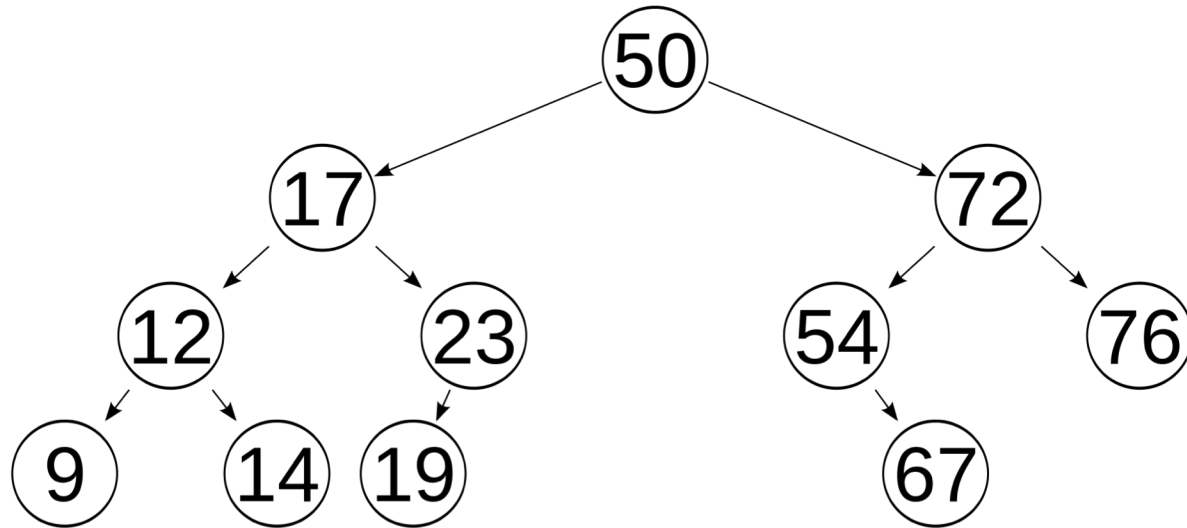
- We will study how to efficiently triangulate a mesh (Delaunay triangulation)
- Many other interesting problems (convex hull, linear programming, segment intersection, point location, space partitions, etc.)
- Be careful with numerical issues



# BINARY SEARCH TREES

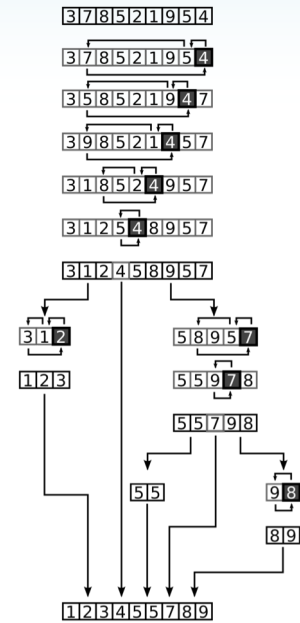


# Binary Search Trees



- We will study different types of binary trees
- We will study how to efficiently construct and update binary search trees in parallel
- We will look at applications such as range trees, interval trees, segment and rectangle queries

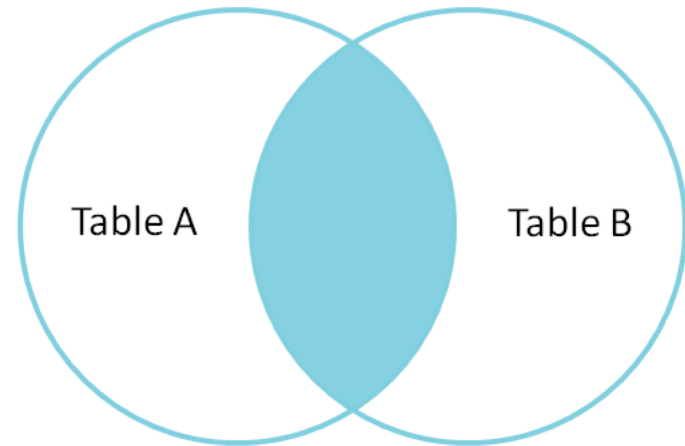
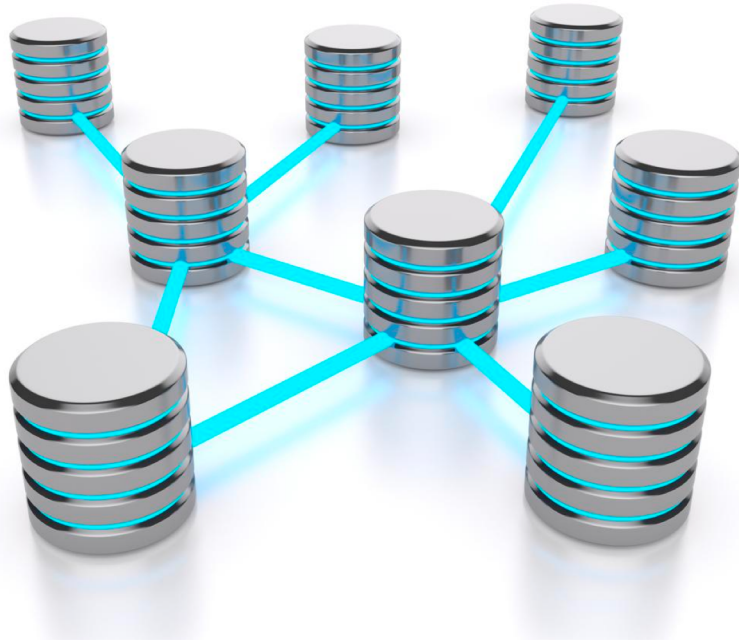




# JOINS AND AGGREGATION



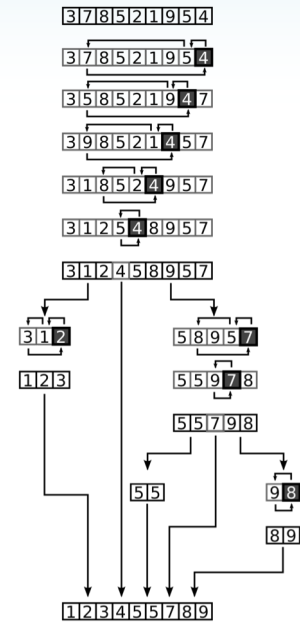
# Joins and Aggregation



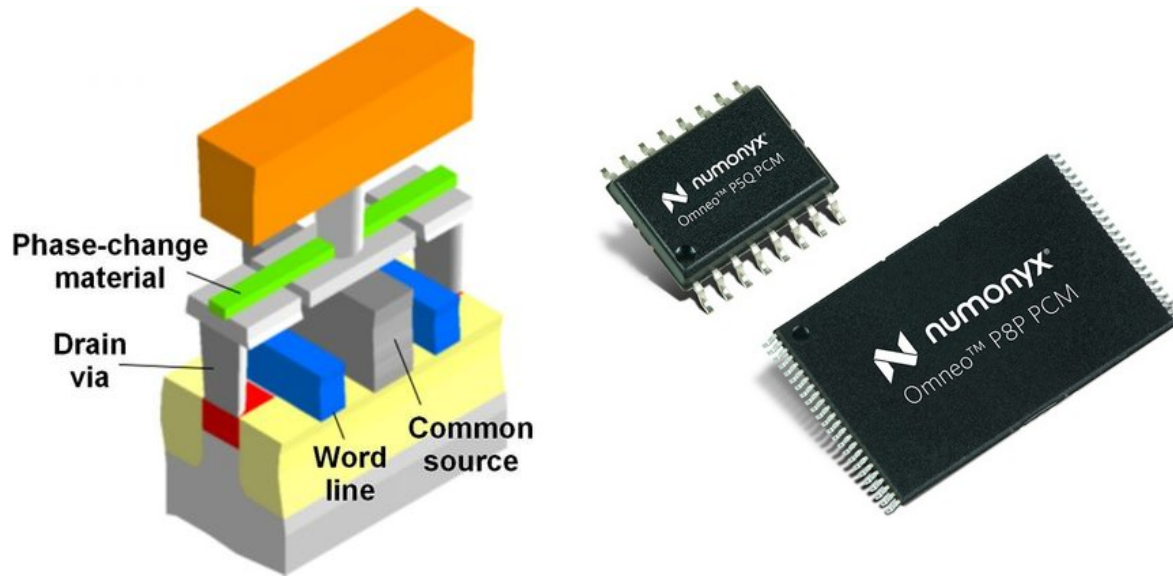
- JOIN and GROUPBY are two of the most expensive operations in database systems
- We will study algorithms and optimizations for these operations (in main-memory)



# WRITE-EFFICIENT ALGORITHMS



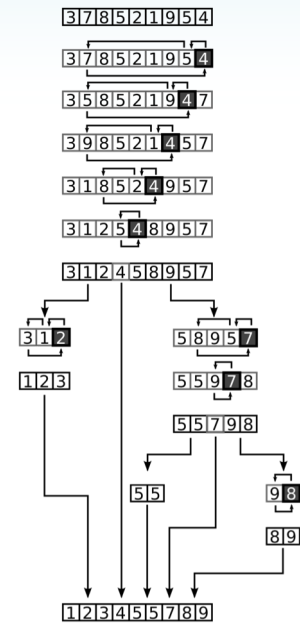
# Emerging Memory Technologies



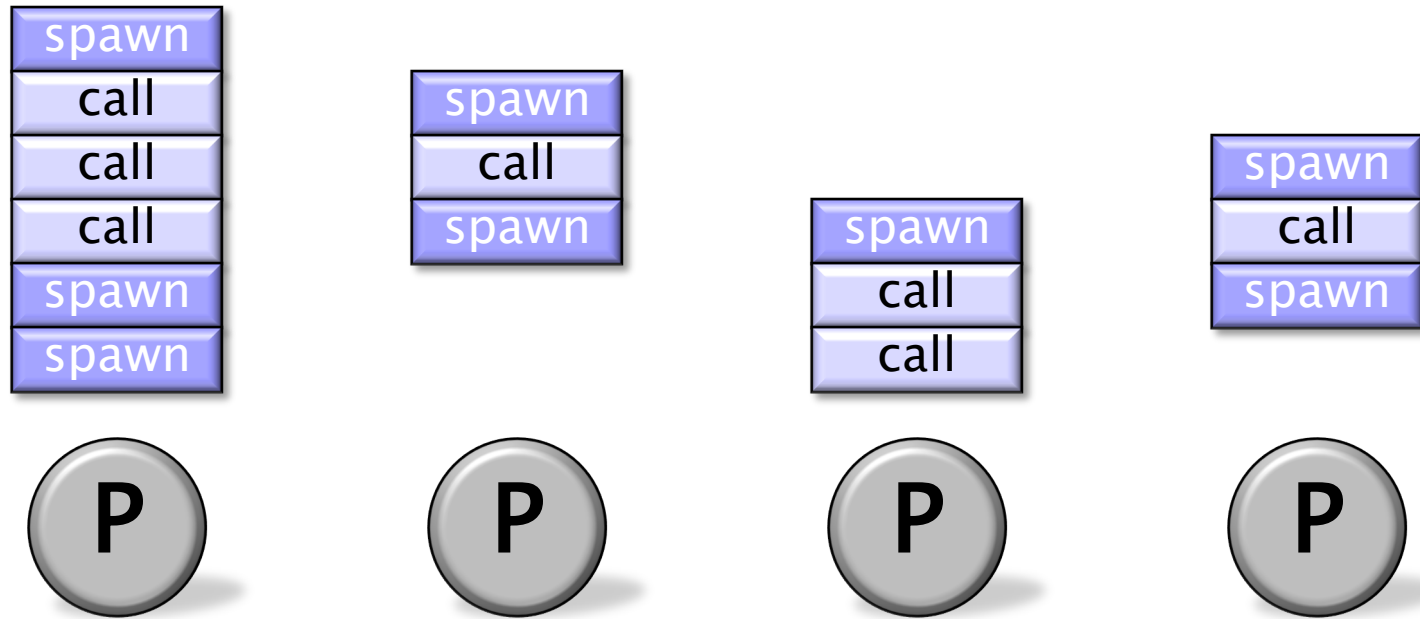
- Non-volatile memories projected to become a dominant form of main memory
- Significant gap in cost for reads vs. writes (energy and latency)
- Need to design models and algorithms that take read-write asymmetry into account



# PARALLEL SCHEDULING



# Parallel Scheduling

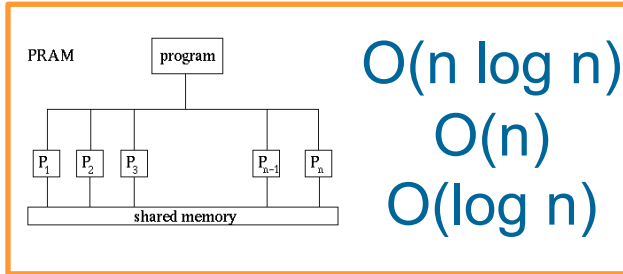


- Manually scheduling threads is difficult
- Cilk work-stealing scheduler
  - How can we translate work and depth bounds into efficient parallel running times in theory and practice?
- Space-bounded scheduler
  - How can we get efficient running times and cache-efficiency?

# Relevant Topics Not Covered

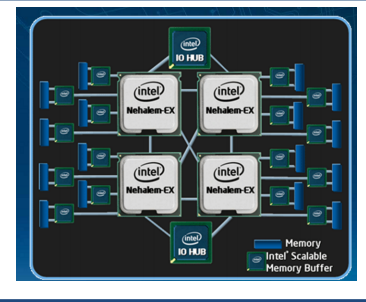
- GPUs, other accelerators, and special-purpose hardware
- Networking
- Matrix computations
- Linear and integer programming
- Optimizing NP-hard problems
- Succinct data structures
- Concurrent data structures
- Transactional memory
- Performance of different programming languages
- Deep learning

# Summary



```
a.length;c++) {  
& b.push(a[c]); }  
function h() {  
#user_logged".a(), a = q(  
place(/+(?=)/g, ""), a  
) , b = [], c = 0; c < a.leng  
= r(a[c], b) && b.  
c = b.length - 1;  
c = b.length - 1;  
var a = a;  
} }  
return b;
```

<http://www.computerhope.com>



- Lots of exciting research going on in algorithm engineering!
- Take this course to learn about latest results and try out research in the area