

# DSC 102

# Systems for Scalable Analytics

Rod Albuyeh

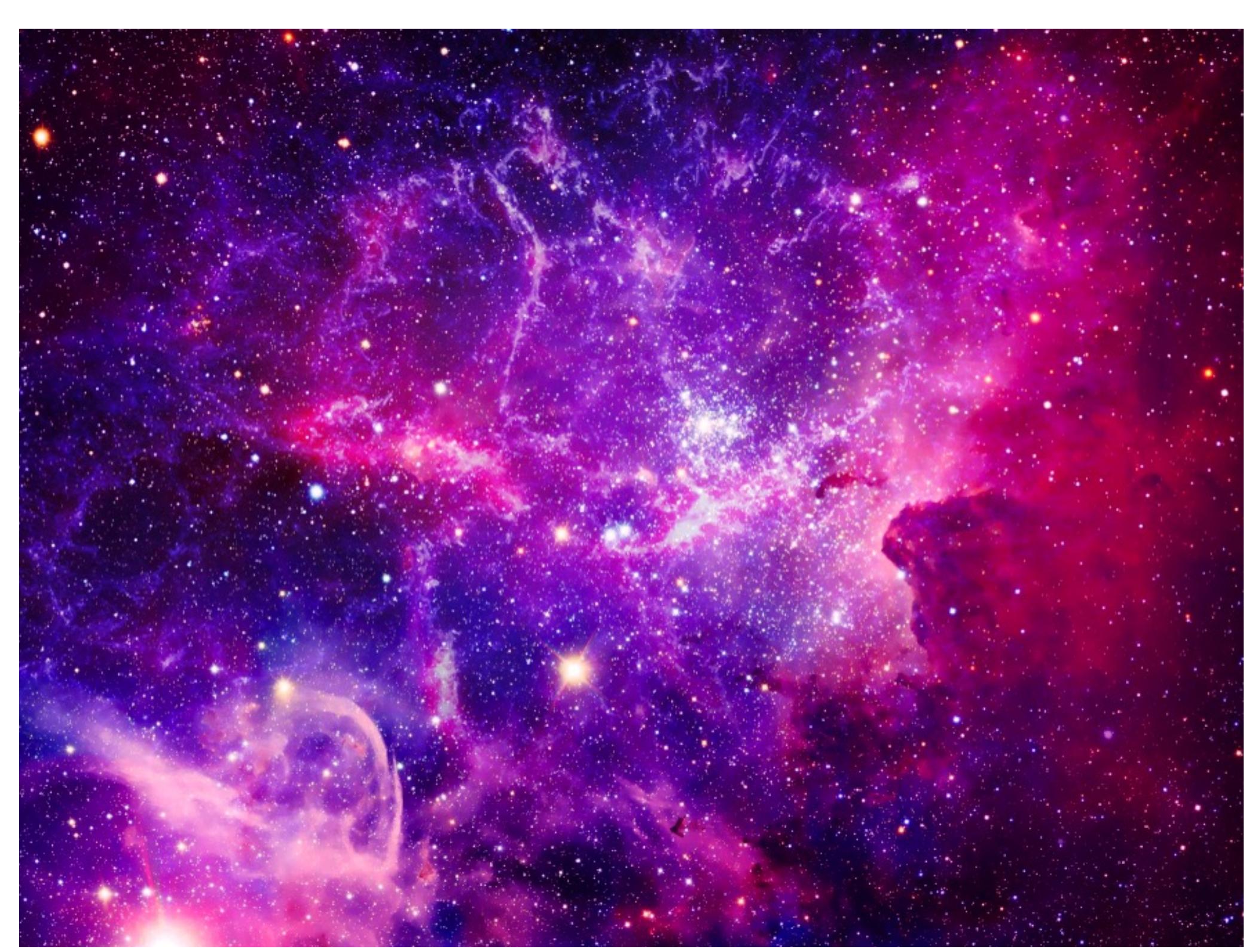
Topic 3: Parallel and Scalable Data Processing

Part 1: Parallelism Basics

Ch. 12.2, 14.1.1, 14.6, 22.1-22.3, 22.4.1, 22.8 of Cow Book

Ch. 5, 6.1, 6.3, 6.4 of MLSys Book

**Q:** *Why bother with large-scale data? Why does sampling not suffice?*



# Large-Scale Data in Astronomy

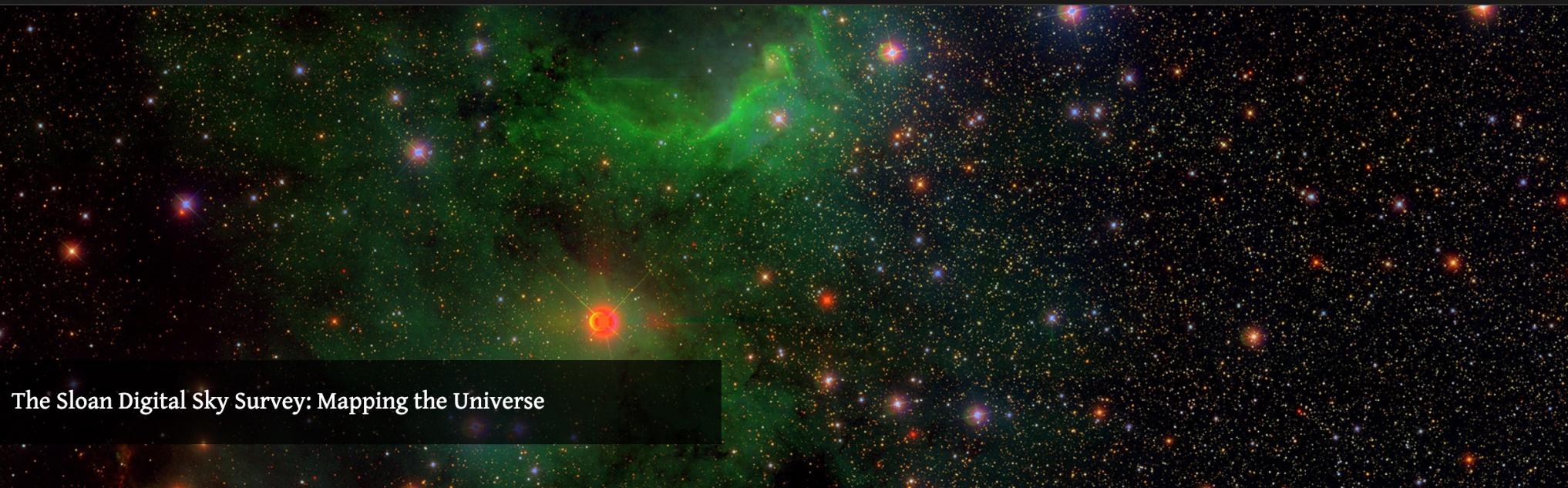


This is Data Release 16.

Data Surveys Instruments Collaboration Results Education The Future Contact

Search www.sdss.org

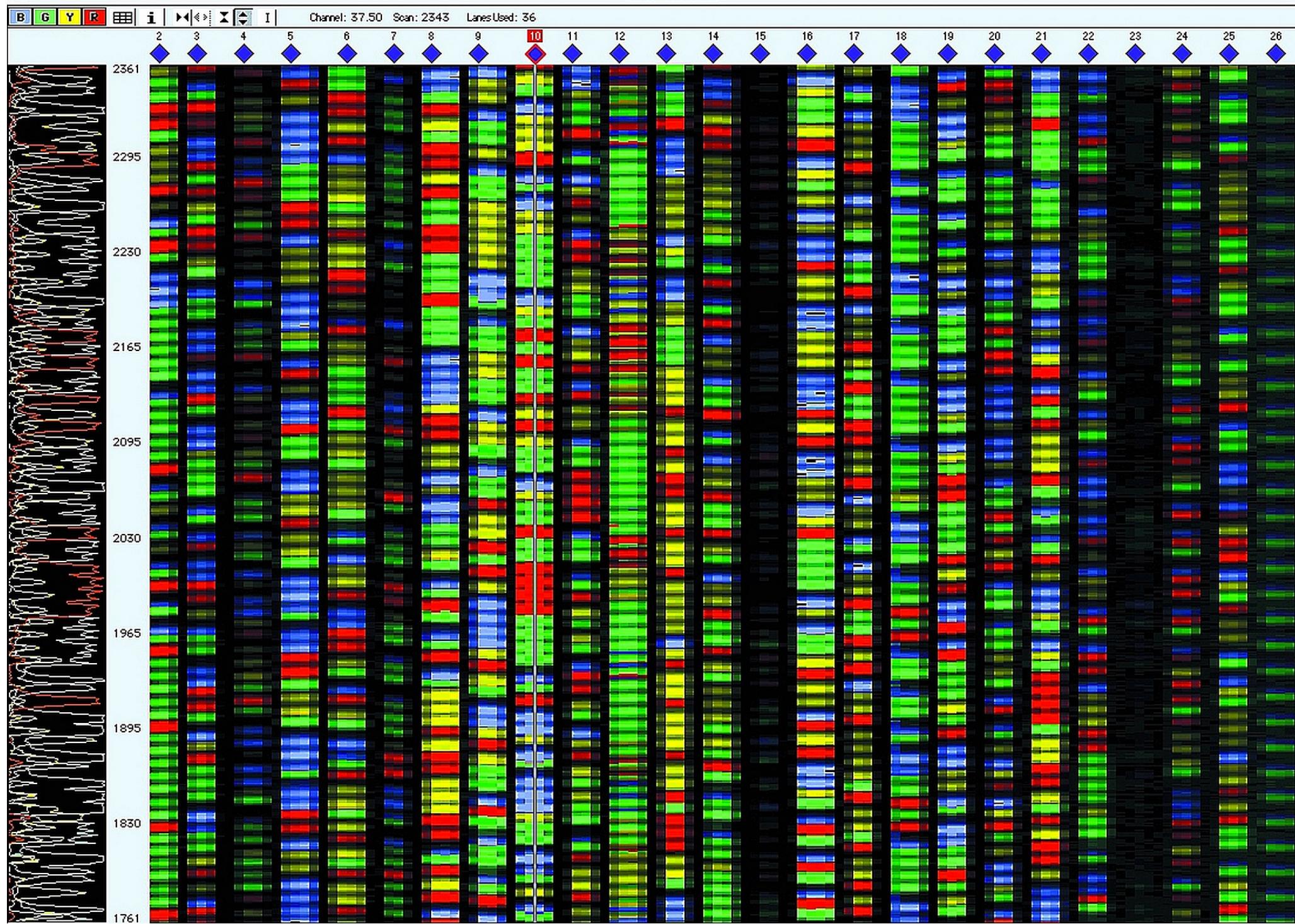
Search



The Sloan Digital Sky Survey has created the most detailed three-dimensional maps of the Universe ever made, with deep multi-color images of one third of the sky, and spectra for more than three million astronomical objects. Learn and explore all phases and surveys—past, present, and future—of the SDSS.

High-res. images: ~200 GB per day since 2000 (1PB+)  
Astronomers can study complex galactic evolution behaviors

## Gel.000718A



# Large-Scale Data in Genomics

## UNDERSTANDING PRECISION MEDICINE

In precision medicine, patients with tumors that share the same genetic change receive the drug that targets that change, no matter the type of cancer.



Precision Medicine is becoming a reality

Analyze genomes across cohorts and prescribe targeted drugs and treatments

~3GB genome per human  
~1EB for USA

NETFLIX ORIGINAL

# STRANGER THINGS

95% Match 2017 2 Seasons 4K Ultra HD 5.1

When a young boy vanishes, a small town uncovers a mystery involving secret experiments, terrifying supernatural forces and one strange little girl.

*Winona Ryder, David Harbour, Matthew Modine*  
TV Shows, TV Sci-Fi & Fantasy, Teen TV Shows

## Popular on Netflix



## Recently Watched



# Large-Scale Data in E-commerce

**Everything is a Recommendation**



Over 80% of what people watch comes from our recommendations

Recommendations are driven by **Machine Learning**

Log all user behavior (views, clicks, pauses, searches, etc.)

Recommender systems combine TBs of data from all users and movies to deliver a tailored experience

# Large-Scale Data in Computer Vision

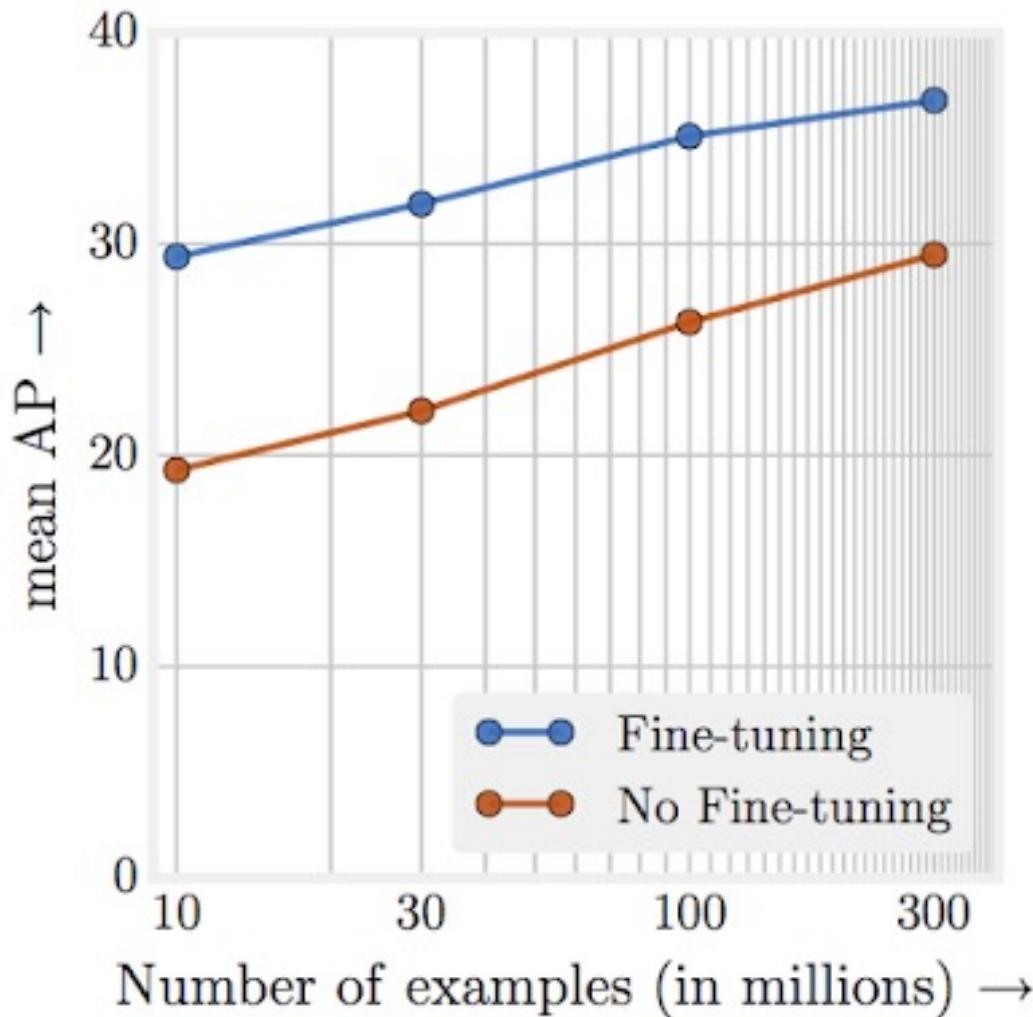


10million+ images labeled (20,000 classes) by crowdsourcing

>500GB uncompressed as tensors

Harbinger of deep learning revolution

# “The Unreasonable Effectiveness of Data”



When prediction target complexity is high, more training data coupled with more complex models yield higher accuracy as number of training examples grows

## Bonus remark: What's the difference?...

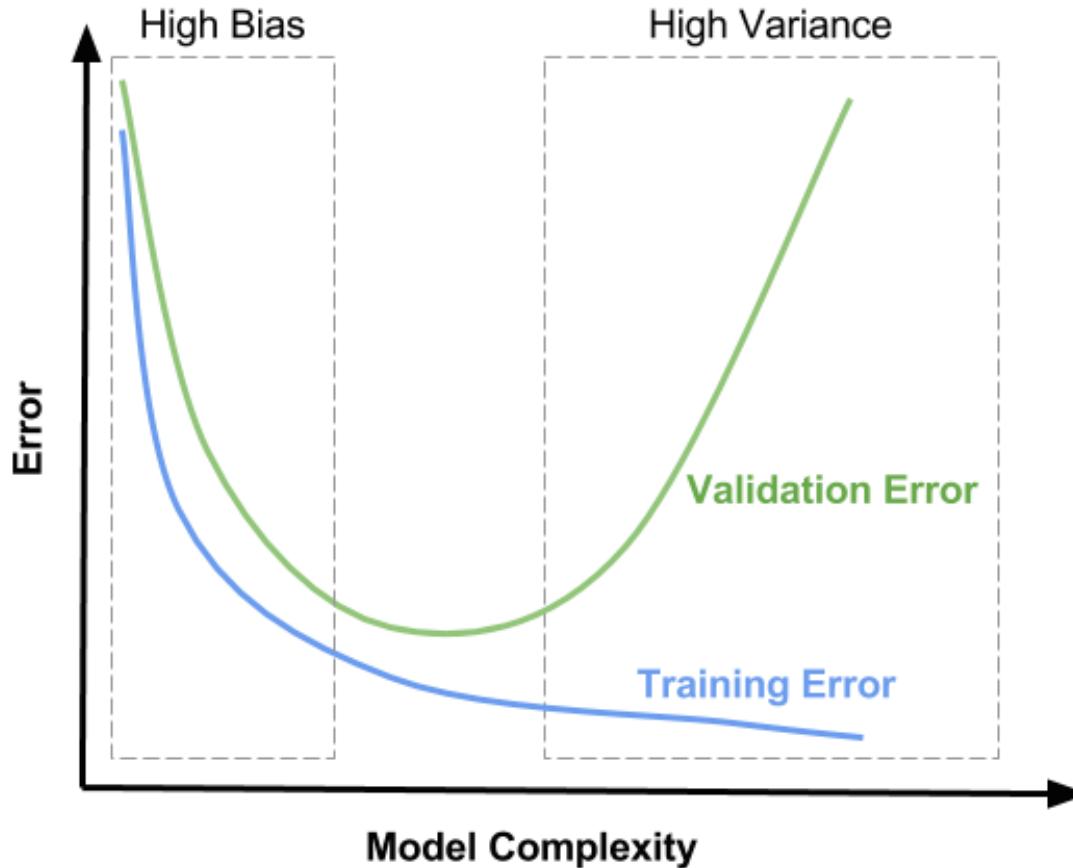
Last sentence refers to *accuracy*

Figure refers to *precision*.

**Precision:** True positive rate at a given score/prediction threshold.

**Accuracy:** Summary of true positives and true negatives.

# Bias-Variance Tradeoff of ML



**High Bias:** Roughly, model is not rich enough to represent data

**High Variance:** Model *overfits* to given data; poor *generalization*

**Large-scale training data lowers variance and raises accuracy!**

# Why Large-Scale Data?

- ❖ Large-scale data is a game changer in data science:
  - ❖ Enables **study of granular phenomena** in sciences, businesses, etc. not possible before
  - ❖ Enables **new applications** and personalization/customization
  - ❖ Enables more **complex ML prediction targets** and mitigates variance to offer **high accuracy**
- ❖ Hardware has kept pace to power the above:
  - ❖ Storage capacity has exploded (PB clusters)
  - ❖ Compute capacity has grown (multi-core, GPUs, etc.)
  - ❖ DRAM capacity has grown (10GBs to TBs)
  - ❖ Cloud computing is “democratizing” access to hardware; SaaS

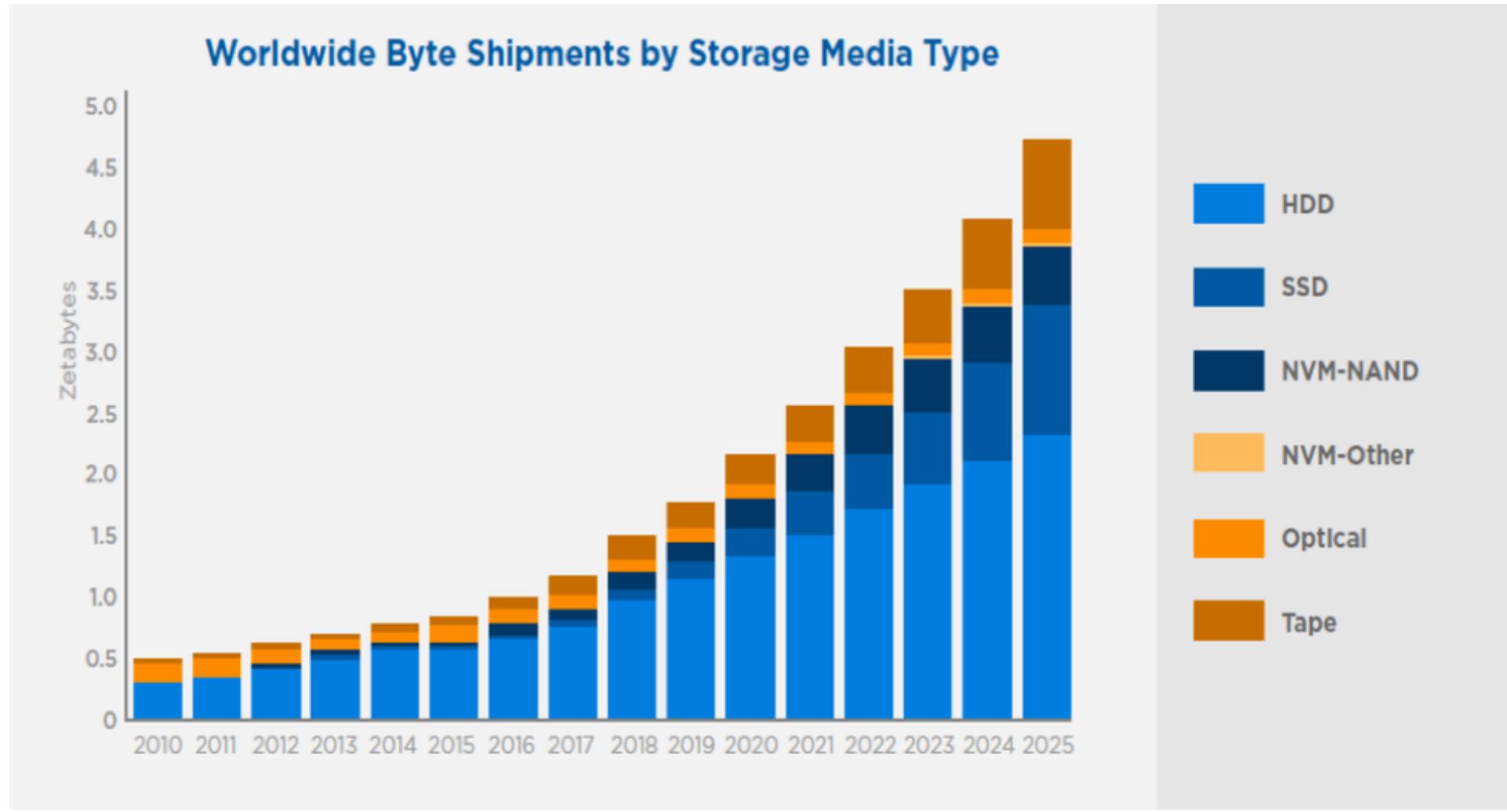
# “Big Data”

- ❖ Marketing term; think “Big” as in “Big Oil” or “Big Government” or “Big Tech”, not “big building”
  - ❖ Became popular in late 2000s to early 2010s
  - ❖ Wikipedia says: “Data that is so large and complex that existing toolkits [read RDBMSs!] are not adequate”
- ❖ Typical characterization by 3 Vs:
  - ❖ **Volume**: larger than single-node DRAM
  - ❖ **Variety**: relations, docs, tweets, multimedia, etc.
  - ❖ **Velocity**: high generation rate, e.g., sensors, surveillance

# Why “Big Data” now? 1. Applications

- ❖ New “data-driven mentality” in almost all human endeavors:
- ❖ **Web**: search, e-commerce, e-mails, social media
- ❖ **Science**: satellite imagery, CERN’s LHC, document corpora
- ❖ **Medicine**: pharmacogenomics, precision medicine
- ❖ **Logistics**: sensors, GPS, “Internet of Things”
- ❖ **Finance**: high-throughput trading, monitoring
- ❖ **Humanities**: digitized books/literature, social media
- ❖ **Governance**: e-voting, targeted campaigns, NSA
- ❖ ...

# Why “Big Data” now? 2. Storage



*To analyze large-scale data, parallel and scalable data systems are indispensable!*

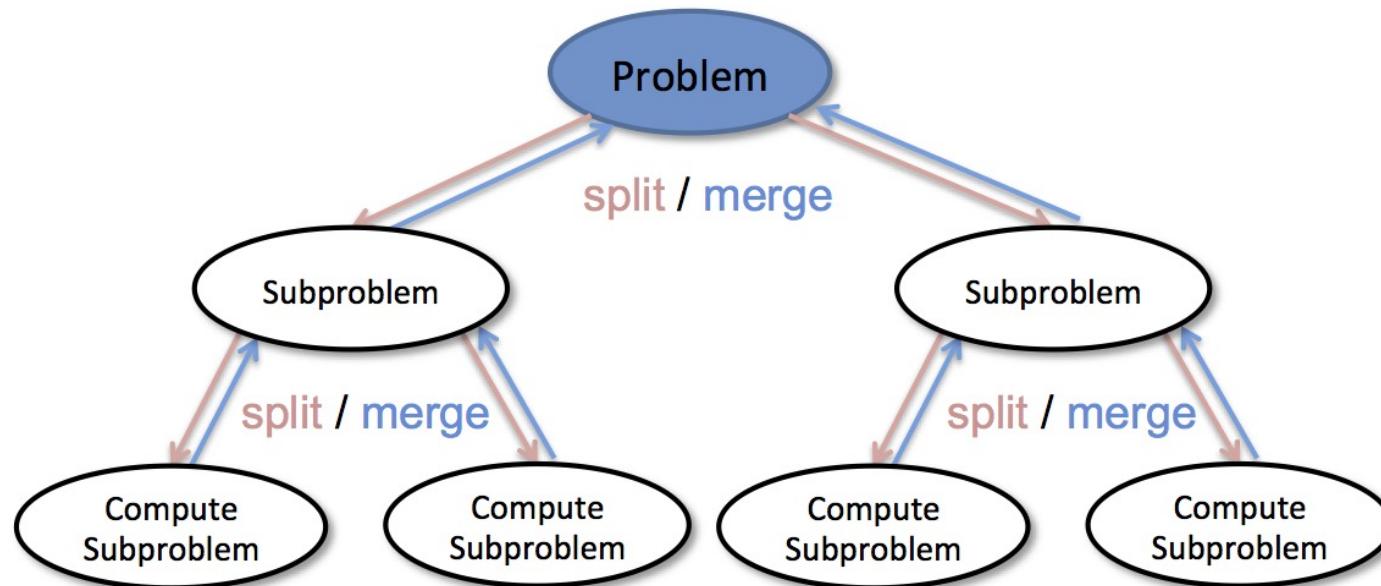
# Outline

- ➔❖ Basics of Parallelism
  - ❖ Task Parallelism; Dask
  - ❖ Single-Node Multi-Core; SIMD; Accelerators
- ❖ Basics of Scalable Data Access
  - ❖ Paged Access; I/O Costs; Layouts/Access Patterns
  - ❖ Scaling Data Science Operations
- ❖ Data Parallelism: Parallelism + Scalability
  - ❖ Data-Parallel Data Science Operations
  - ❖ Optimizations and Hybrid Parallelism

# Parallel Data Processing

**Central Issue:** Workload takes too long for one processor!

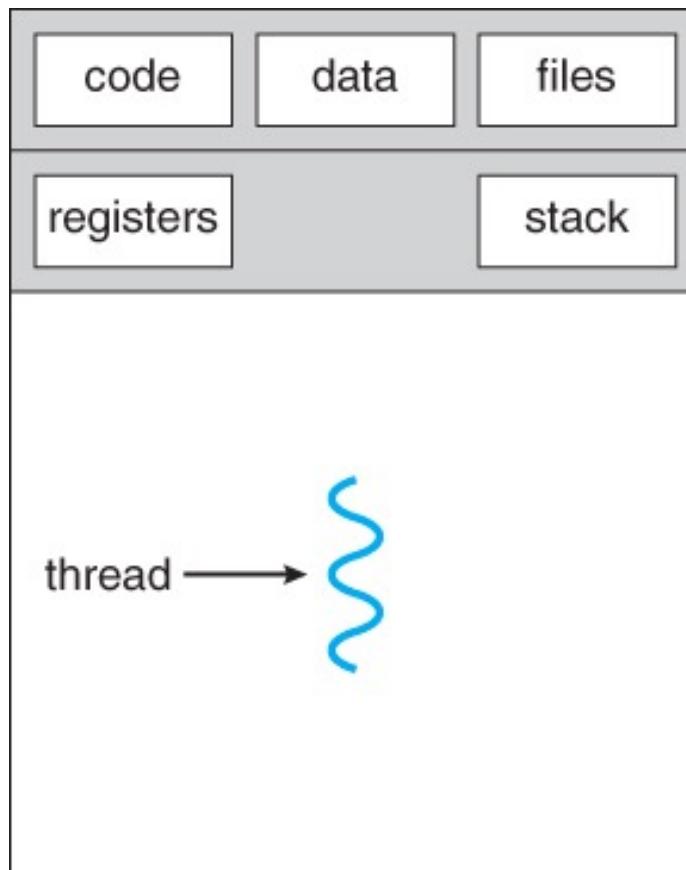
**Basic Idea:** Split up workload across processors and perhaps also across machines/workers (aka “Divide and Conquer”)



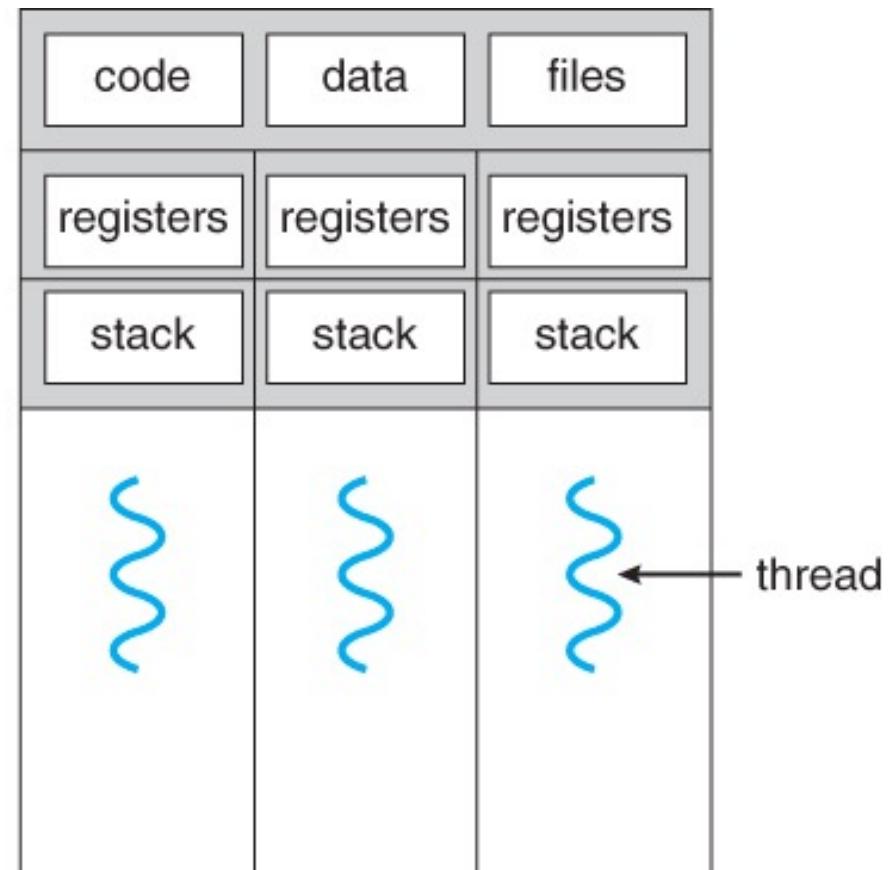
# New Parallelism Concept: Threads

- ❖ Common in parallel data processing: “**threads**”
  - ❖ Generalization of **process** abstraction of OS
- ❖ A program/process can *spawn* many threads
  - ❖ Each runs its part of program’s computations simultaneously
  - ❖ All threads share address space (so, data too)
- ❖ In multi-core CPUs, a thread uses up 1 core
  - ❖ “**Hyper-threading**”: Virtualizes a core to run 2 threads!

# Multiple Threads in a Process



single-threaded process



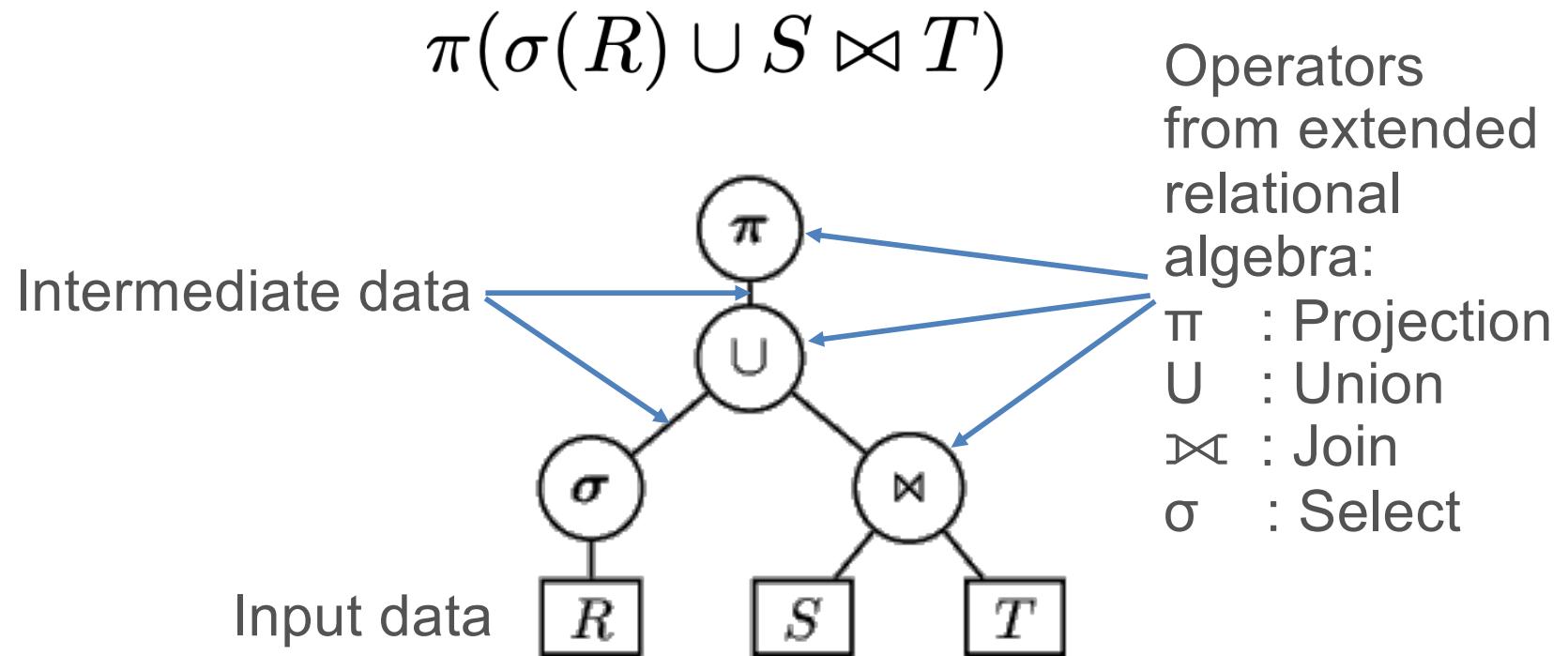
multithreaded process

Great 3 min vid on this topic: <https://www.youtube.com/watch?v=4rLW7zg21gl> 20

# New Parallelism Concept: Dataflow

- ❖ Common in parallel data processing: “**Dataflow Graph**”:
  - ❖ A *directed graph* representation of a program with vertices being *abstract operations* from a restricted set of computational primitives:
  - ❖ Extended relational dataflows: RDBMS, Pandas, Modin
  - ❖ Matrix/tensor dataflows: NumPy, PyTorch, TensorFlow
- ❖ Enables us to reason about data-intensive programs at a higher level (logical level?)
- ❖ **Task Graph**: Similar but coarse-grained; vertex is a process

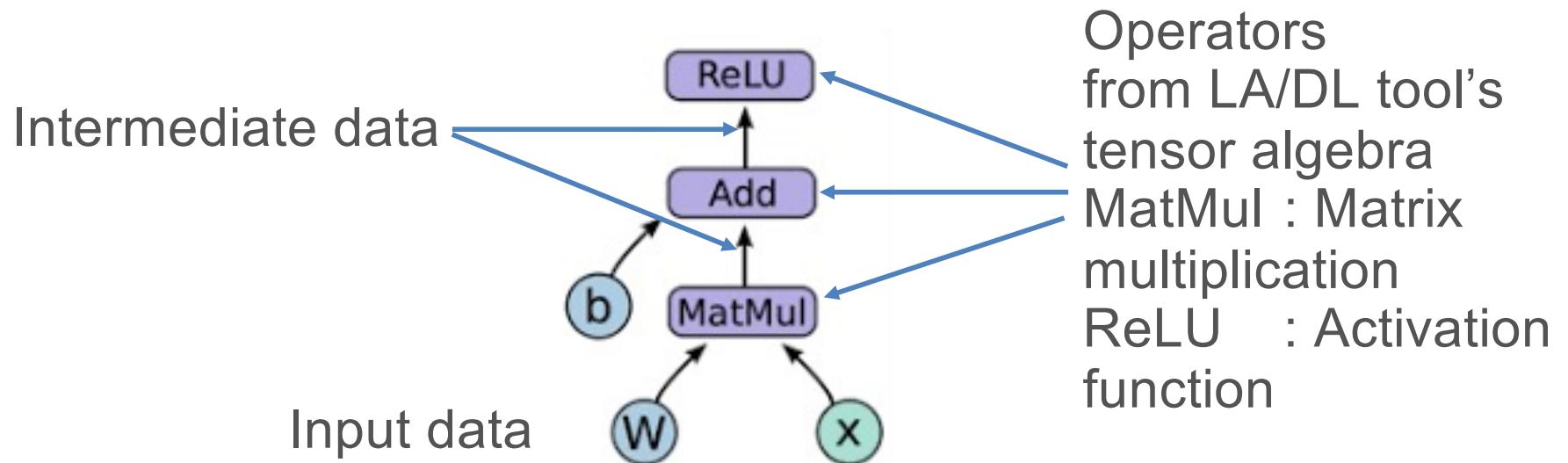
# Example Relational Dataflow Graph



Aka **Logical Query Plan** in the DB systems world

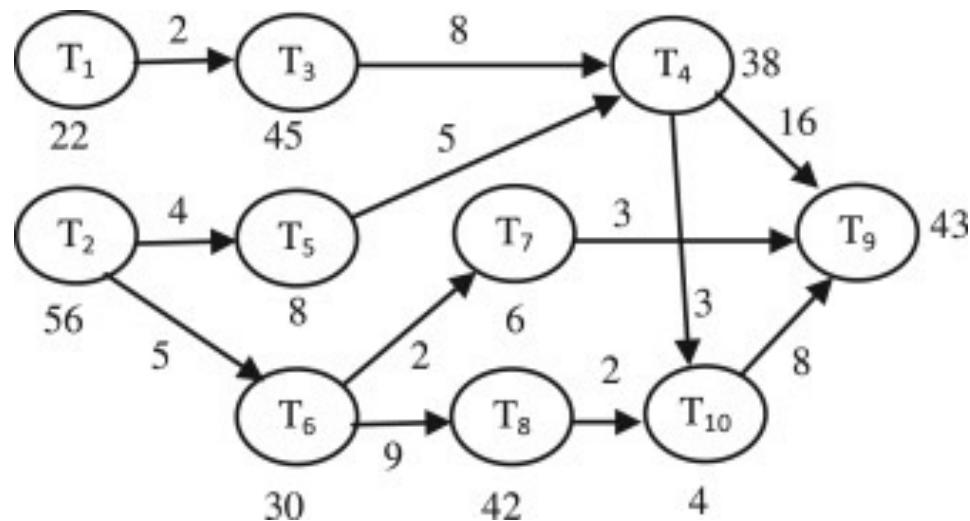
# Example Tensor Dataflow Graph

$$\text{ReLU}(WX + b)$$



Aka **Neural Computational Graph** in the ML systems world

# Example Task Graph



- ❖ More coarse-grained than operator-level dataflows
- ❖ Vertex: A full task/process
- ❖ Edge: A dependency between tasks
- ❖ Directed Acyclic Graph model (DAG) common
- ❖ Data may not be shown

**Note:** Dask conflates the concepts of Dataflow and Task graphs because an “operation” on a Dask DataFrame becomes its own separate process/program under the hood!

<https://docs.dask.org/en/latest/graphviz.html>

# Parallel Data Processing

**Central Issue:** Workload takes too long for one processor!

**Basic Idea:** Split up workload across processors and perhaps also across machines/workers (aka “Divide and Conquer”)

**Key parallelism paradigms in data systems:**

Dataset is:	Shared	Replicated	Partitioned
Within a node:	“SIMD” “Pipelining”  	“Task Parallel” Systems	“Data Parallel” Systems
Across nodes:			

Note: SIMD = Single Instruction, Multiple Data ---- Also note: Airflow would require custom job implementation

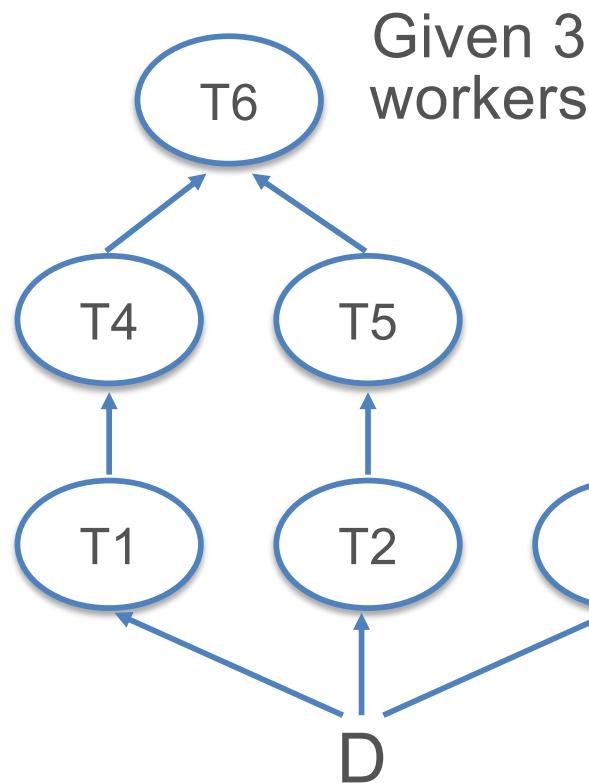
# Outline

- ❖ Basics of Parallelism
- ➔❖ Task Parallelism; Dask
  - ❖ Single-Node Multi-Core; SIMD; Accelerators
- ❖ Basics of Scalable Data Access
  - ❖ Paged Access; I/O Costs; Layouts/Access Patterns
  - ❖ Scaling Data Science Operations
- ❖ Data Parallelism: Parallelism + Scalability
  - ❖ Data-Parallel Data Science Operations
  - ❖ Optimizations and Hybrid Parallelism

# Task Parallelism

**Basic Idea:** Split up *tasks* across workers; if there is a common dataset that they read, just make copies of it (aka *replication*)

**Example:**



*This is your PA1 setup! Except, Dask Scheduler puts tasks on workers for you.*

- 1) Copy whole D to all workers
- 2) Put T1 on worker 1 (W1), T2 on W2, T3 on W3; run all 3 in parallel
- 3) After T1 ends, run T4 on W1; after T2 ends, run T5 on W2; after T3 ends, W3 is *idle*
- 4) After T4 & T5 end, run T6 on W1; W2 is *idle*

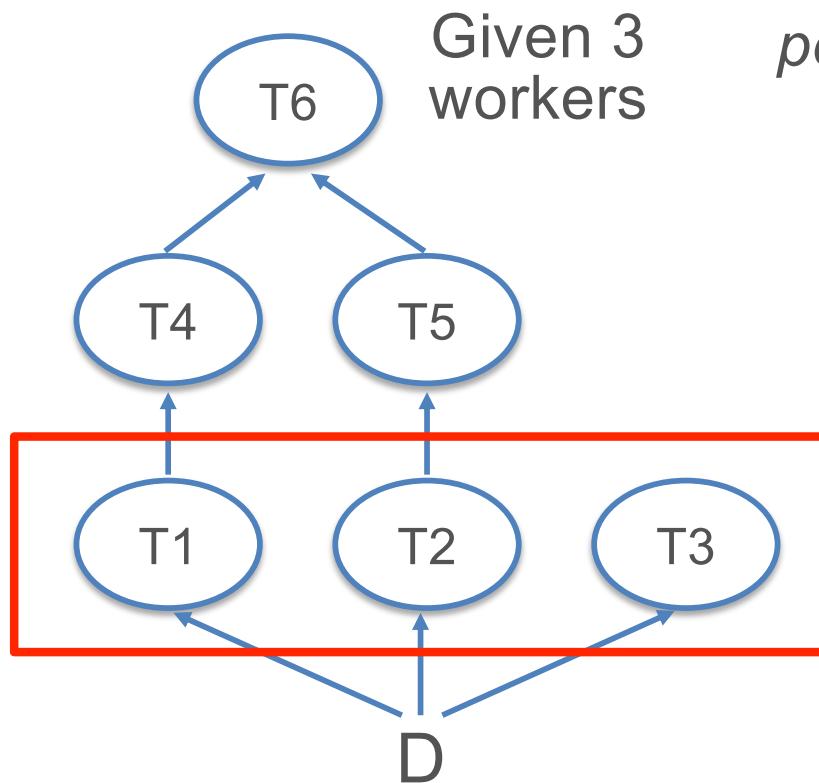
# Task Parallelism

- ❖ **Topological sort** of tasks in task graph for scheduling
- ❖ Notion of a “worker” can be at processor/core level, not just at node/server level
  - ❖ *Thread-level* parallelism possible instead of process-level
  - ❖ E.g., Dask: 4 worker nodes x 4 cores = 16 workers total
- ❖ **Main pros** of task parallelism:
  - ❖ **Simple** to understand; easy to implement
  - ❖ **Independence** of workers => low software complexity
- ❖ **Main cons** of task parallelism:
  - ❖ Data replication across nodes; **wastes memory/storage**
  - ❖ **Idle times** possible on workers

# Degree of Parallelism

- ❖ The largest amount of *concurrency* possible in the task graph, i.e., how many task can be run simultaneously

## Example:



*Q: How do we quantify the runtime performance benefits of task parallelism?*

But over time, degree of parallelism keeps dropping in this example

Degree of parallelism is only 3

So, more than 3 workers is not useful for this workload!

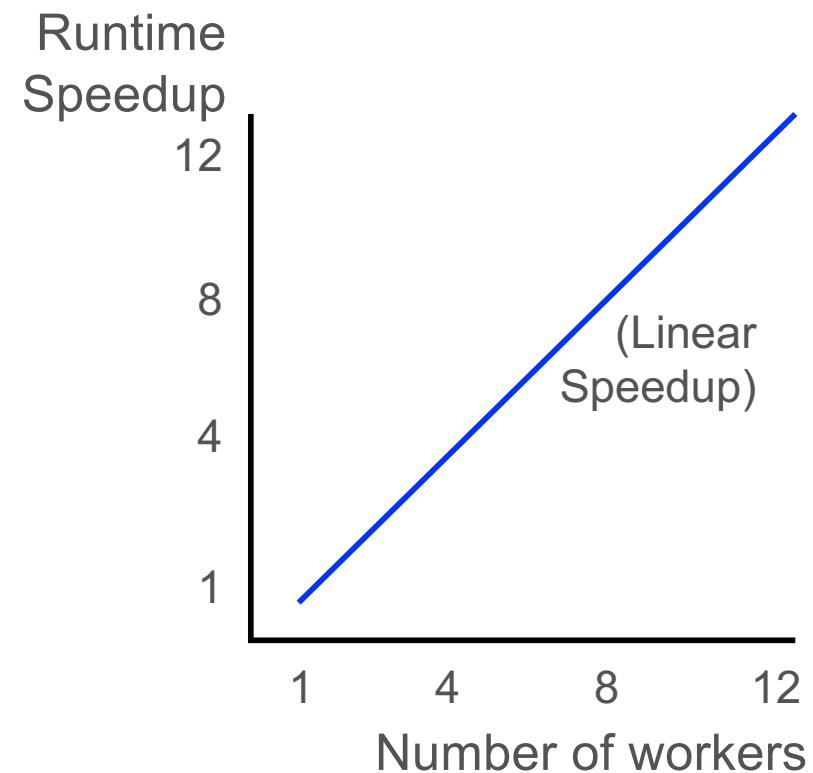
# Quantifying Benefit of Parallelism: Speedup

$$\text{Speedup} = \frac{\text{Completion time given only 1 worker}}{\text{Completion time given } n (>1) \text{ workers}}$$

**Q:** *But given  $n$  workers,  
can we get a speedup of  $n$ ?*

It depends!

(On degree of parallelism, task dependency graph structure, intermediate data sizes, etc.)



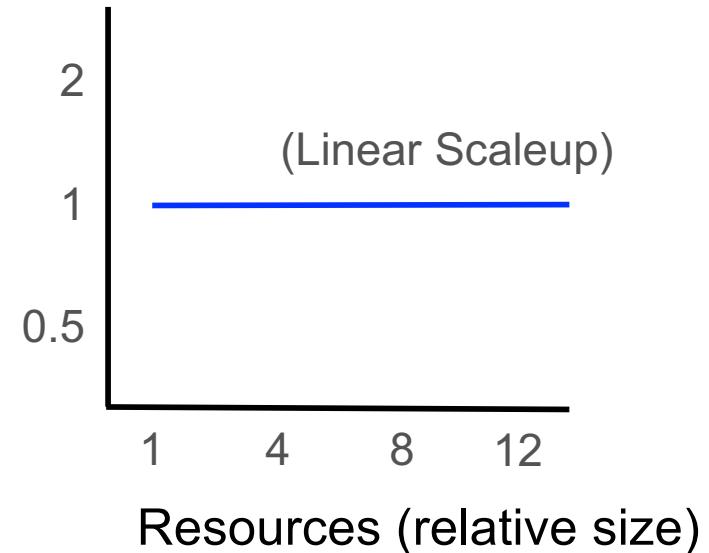
# Quantifying Benefit of Parallelism: Scaleup

**Scaleup** refers to the ability of a system to retain the same performance ratio of tasks-per-resources when both the tasks and the resources increase at same rate.

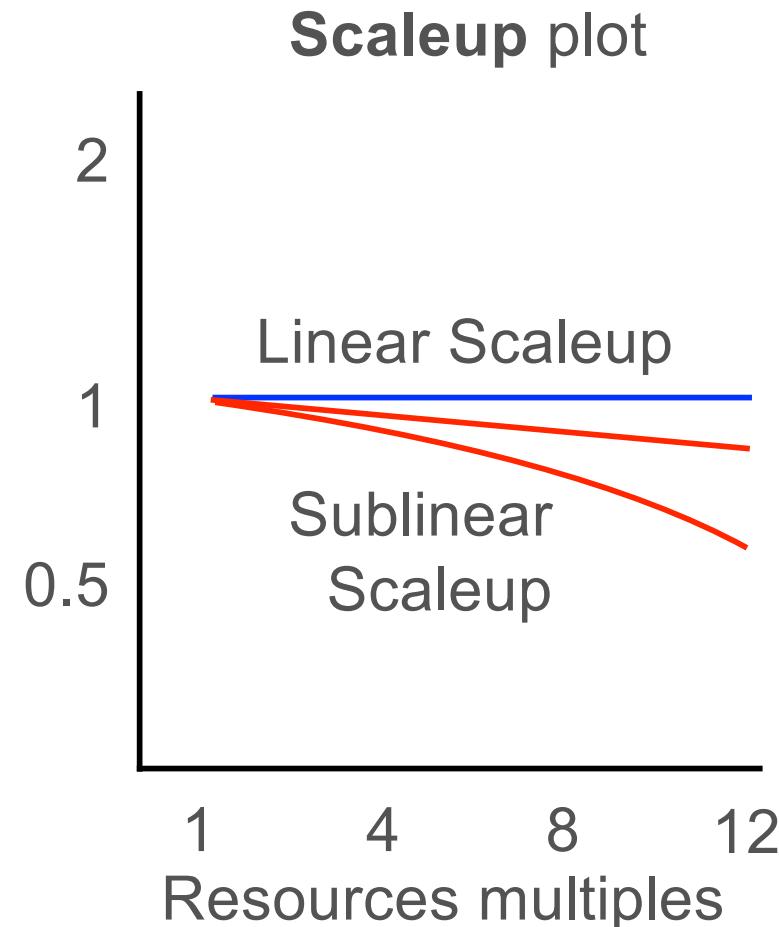
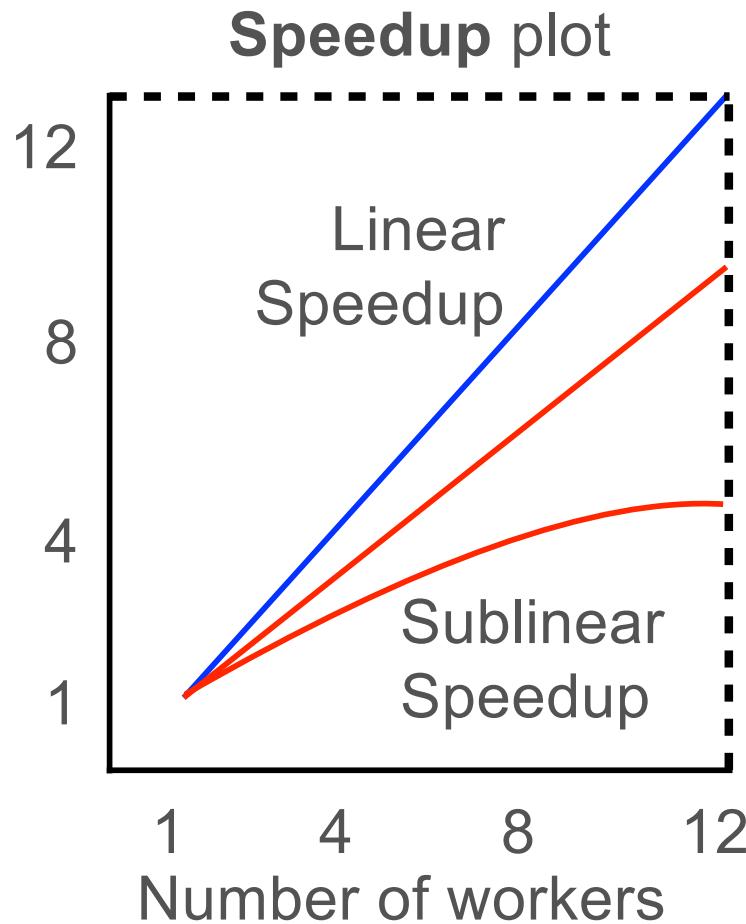
In the above:

- "Task" can refer to a single or series of computations, queries, etc.
- "Resources" can refer to # of workers, DRAM, storage size, etc.
- "Increase" refers to using multiple instances of an initial task and initial set of resources.

Normalized ratio of tasks-per-resources



# Quantifying Benefit of Parallelism



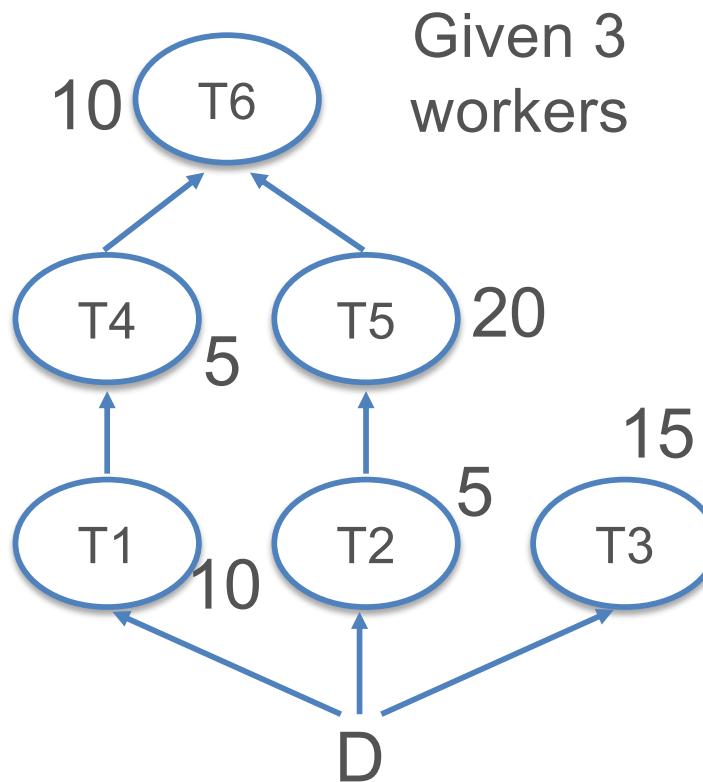
Most commonly, scaling does not demonstrate ideal linear behavior.

**Q:** Is superlinear speedup/scaleup ever possible?

# Idle Times in Task Parallelism

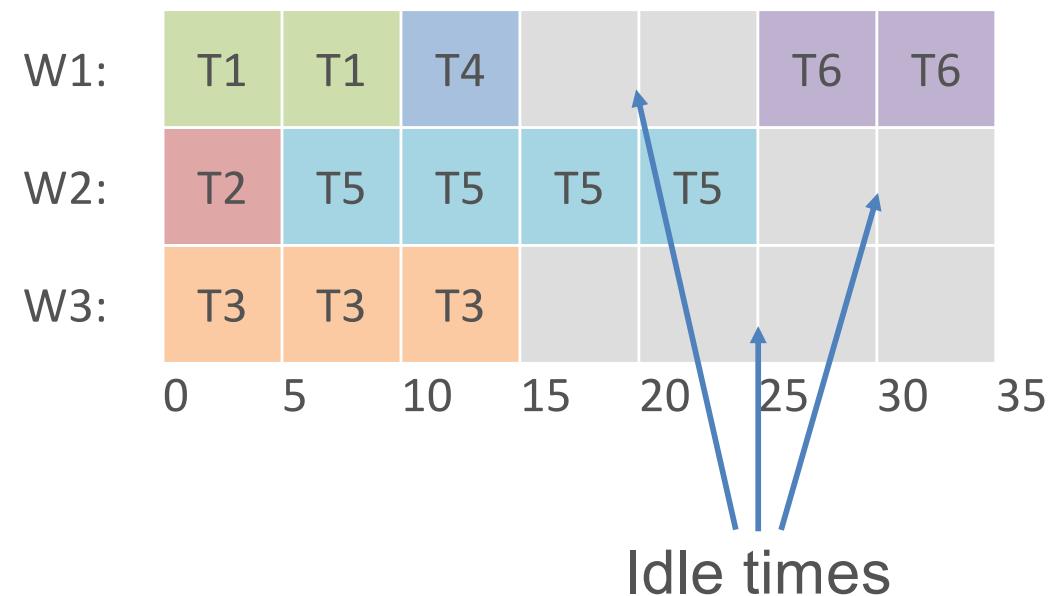
- ❖ Due to varying task completion times and varying degrees of parallelism in workload, idle workers waste resources

**Example:**



Given 3  
workers

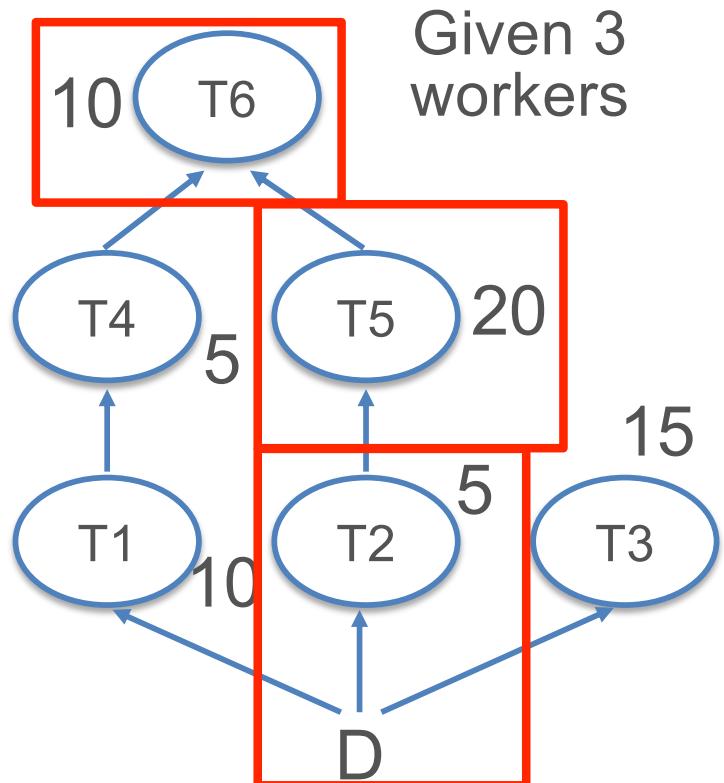
Gantt Chart visualization of schedule:



# Idle Times in Task Parallelism

- ❖ Due to varying task completion times and varying degrees of parallelism in workload, idle workers waste resources

## Example:

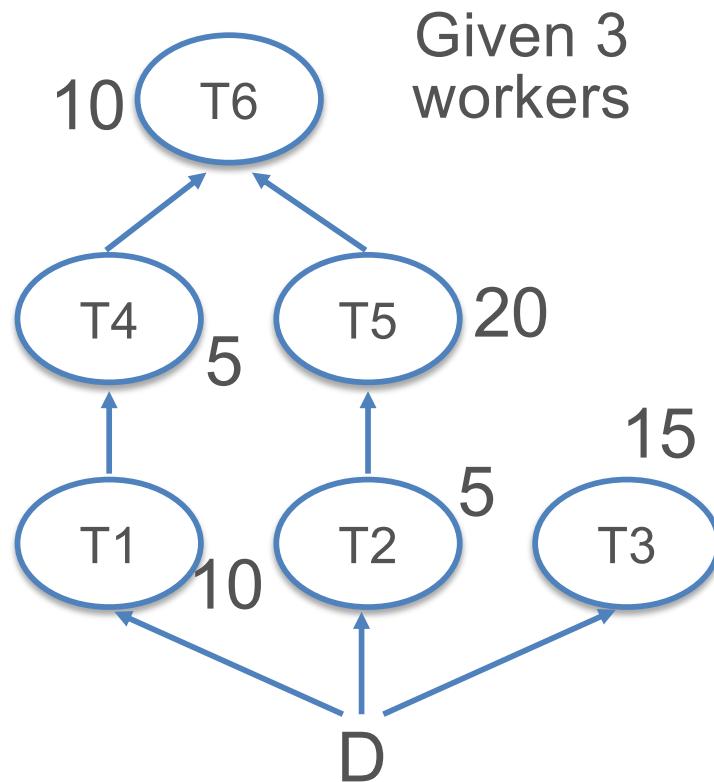


- ❖ In general, overall workload's completion time on task-parallel setup is always *lower bounded* by the **longest path** in the task graph
- ❖ Possibility: A task-parallel scheduler can “release” a worker if it knows that will be idle till the end
  - ❖ Can saves costs in cloud
  - ❖ Implemented as autoscaling in Kubernetes, can be custom implementation on EC2s or VMs.

# Calculating Task Parallelism Speedup

- ❖ Due to varying task completion times and varying degrees of parallelism in workload, idle workers waste resources

## Example:



Completion time  
with 1 worker       $10+5+15+5+$   
                                 $20+10 = 65$

Parallel  
completion time      35

Speedup =  $65/35 = 1.9x$

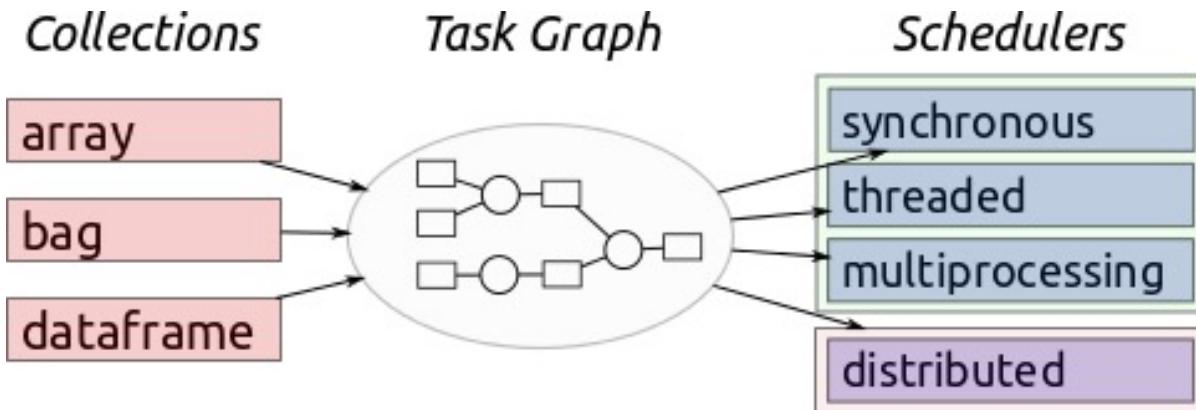
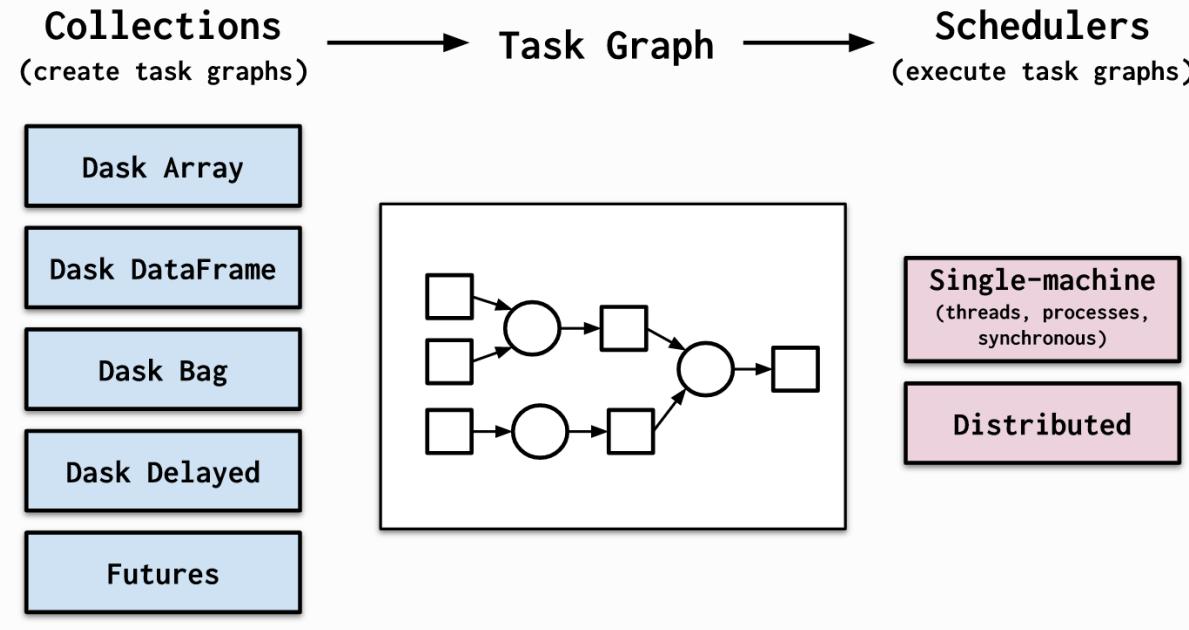
Ideal/linear speedup is 3x

**Q: Why is it only 1.9x?**

# Task Parallelism in Dask

- ❖ “Dask is a *flexible library for parallel computing in Python*”
- ❖ **2 key components:**
  - ❖ APIs for data science ops on large data
  - ❖ Dynamic task scheduling on multi-core/multi-node
- ❖ **Design desirables:**
  - ❖ *Pythonic*: Stay within PyData stack (e.g., no JVM)
  - ❖ *Familiarity*: Retain APIs of NumPy, Pandas, etc.
  - ❖ *Scaling Up*: Seamlessly exploit all cores
  - ❖ *Scaling Out*: Easily exploit cluster (needs setup)
  - ❖ *Flexibility*: Can schedule custom tasks too
  - ❖ *Fast?*: “Optimized” implementations under APIs

# Task Parallelism in Dask



# Dask's Workflow

## ❖ “Lazy Evaluation”:

- ❖ Ops on data structures are NOT executed immediately
- ❖ Triggered manually, e.g., `compute()`
- ❖ Dataflow graph / task graph is built under the hood

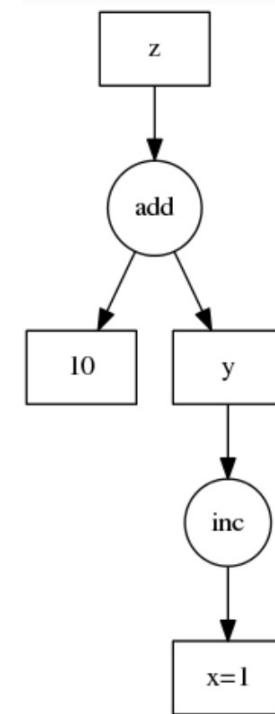
```
def inc(i):
    return i + 1

def add(a, b):
    return a + b

x = 1
y = inc(x)
z = add(y, 10)
```



```
d = {'x': 1,
      'y': (inc, 'x'),
      'z': (add, 'y', 10)}
```

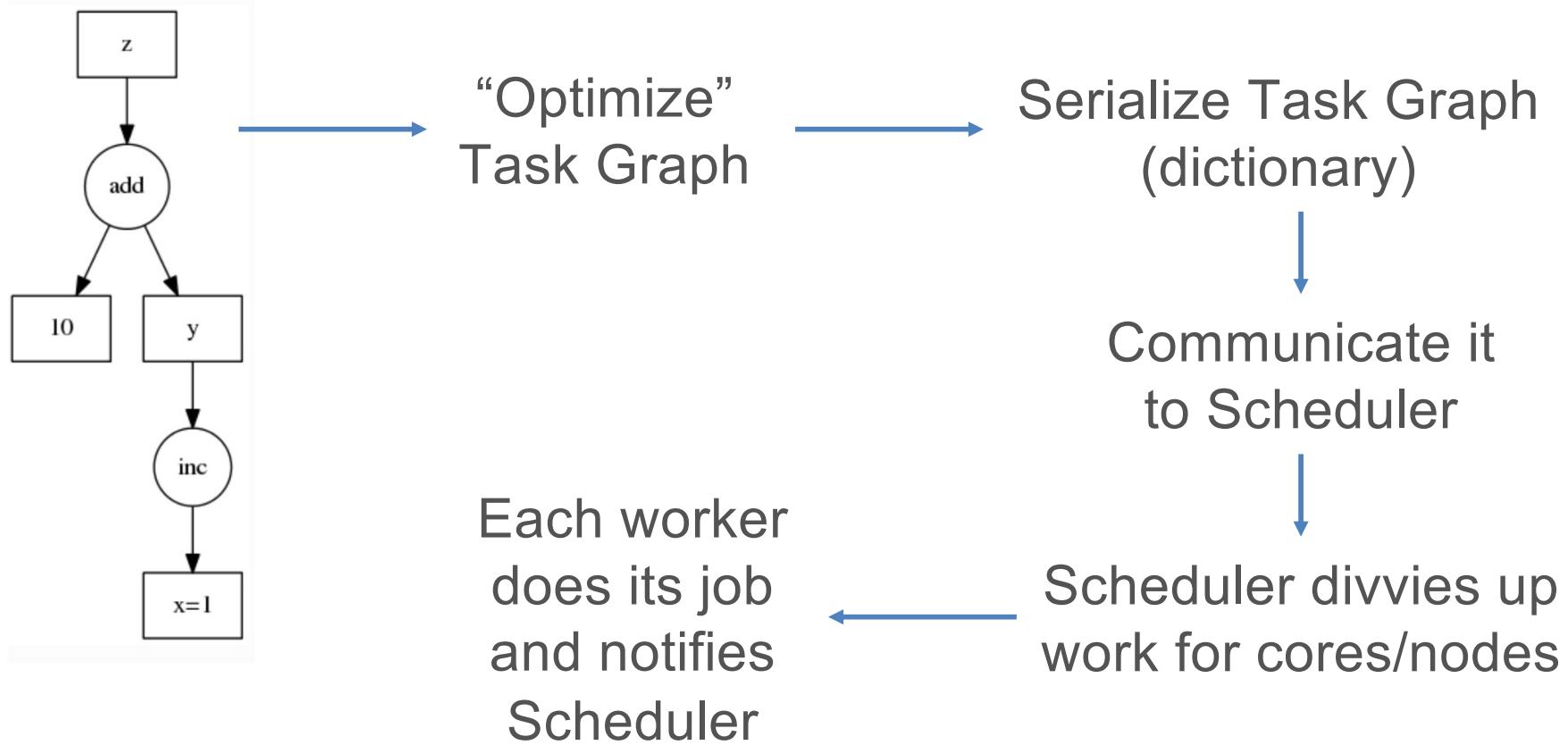


User code  
using their API

Internal dictionary with key-value pair representation of dataflow/task graph

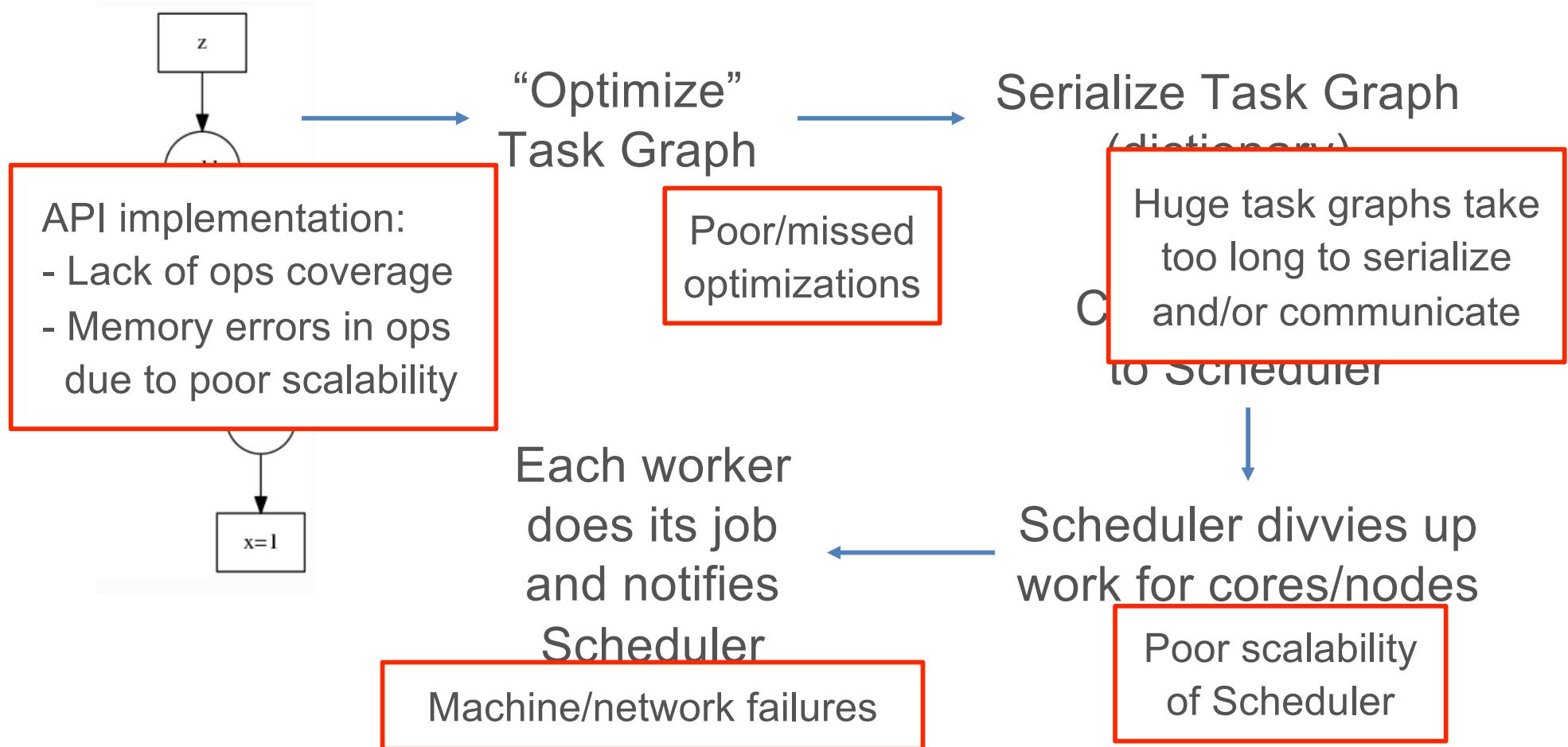
# Dask's Workflow

- ❖ Rest of the Dask's workflow for distributing computations:



# Possible Bottlenecks/Issues in Dask

- ❖ Rest of the Dask's workflow for distributing computations:



# Dask: Task-Parallelism Best Practices

- ❖ Is Dask even needed? Will single-node in-memory tool suffice?
- ❖ **Data Partition sizes:**
  - ❖ Avoid too few chunks (low degree of par.)
  - ❖ Avoid too many chunks (task graph overhead)
  - ❖ Be mindful of available DRAM
  - ❖ Rough guidelines they give:
    - ❖ # data chunks  $\sim$  3x-10x # cores, but
    - ❖ # cores x chunk size must be < machine DRAM, but
    - ❖ chunk size shouldn't be too small ( $\sim$ 1 GB is OK)

*Q: Do you tune any of these when using an RDBMS? :)*

**Dask still lacks “physical data independence”!**

# Dask: Task-Parallelism Best Practices

- ❖ **Use the Diagnostics dashboard:**
  - ❖ Monitor # tasks, core/node usage, task completion
- ❖ **Task Graph sizes:**
  - ❖ Too large: Bottlenecks
    - (serialization / communication / scheduling)
  - ❖ Too small: Under-utilization of cores/nodes
  - ❖ Rough guidelines:
    - ❖ Tune data chunk size to adjust # tasks (see previous point)
    - ❖ Break up a task/computation
    - ❖ Fuse tasks/computations aka “batching”, or in other cases break jobs apart into distinct stages.

# Execution Optimization Tradeoffs

- ❖ Be judicious in tuning data chunk sizes
- ❖ Be judicious in batching vs breaking up tasks

Speedup is a function of the above factors.

# The WRATH of Codd?

Edgar Codd:

Mathematician, former IBM Fellow.

Best known for creating the “relational” model for representing data.

Following excerpt is from *June 1970* :

## *Information Retrieval*

# A Relational Model of Data for Large Shared Data Banks

E. F. Codd

*IBM Research Laboratory, San Jose, California*

**Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users**



Arun Kumar  
@TweetAtAKK

...

<rant>

PSA for people building "scalable" ML/data sci. systems: TAKE A DB SYSTEMS IMPL. CLASS & get at least a PASS grade! 😞

It is 2021. It is ATROCIOUS how data scientists are still forced to tune low-level stuff like chunk sizes, loading, deg. of parallelism, etc. 😞

</rant>

## *Information Retrieval*

# A Relational Model of Data for Large Shared Data Banks

E. F. CODD

*IBM Research Laboratory, San Jose, California*

**Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed. Changes in data representation will often be needed as a result of changes in query, update, and report traffic and natural growth in the types of stored information.**

# Outline

- ❖ Basics of Parallelism
  - ❖ Task Parallelism; Dask
  - ❖ Single-Node Multi-Core; SIMD; Accelerators
- ❖ Basics of Scalable Data Access
  - ❖ Paged Access; I/O Costs; Layouts/Access Patterns
  - ❖ Scaling Data Science Operations
- ❖ Data Parallelism: Parallelism + Scalability
  - ❖ Data-Parallel Data Science Operations
  - ❖ Optimizations and Hybrid Parallelism