



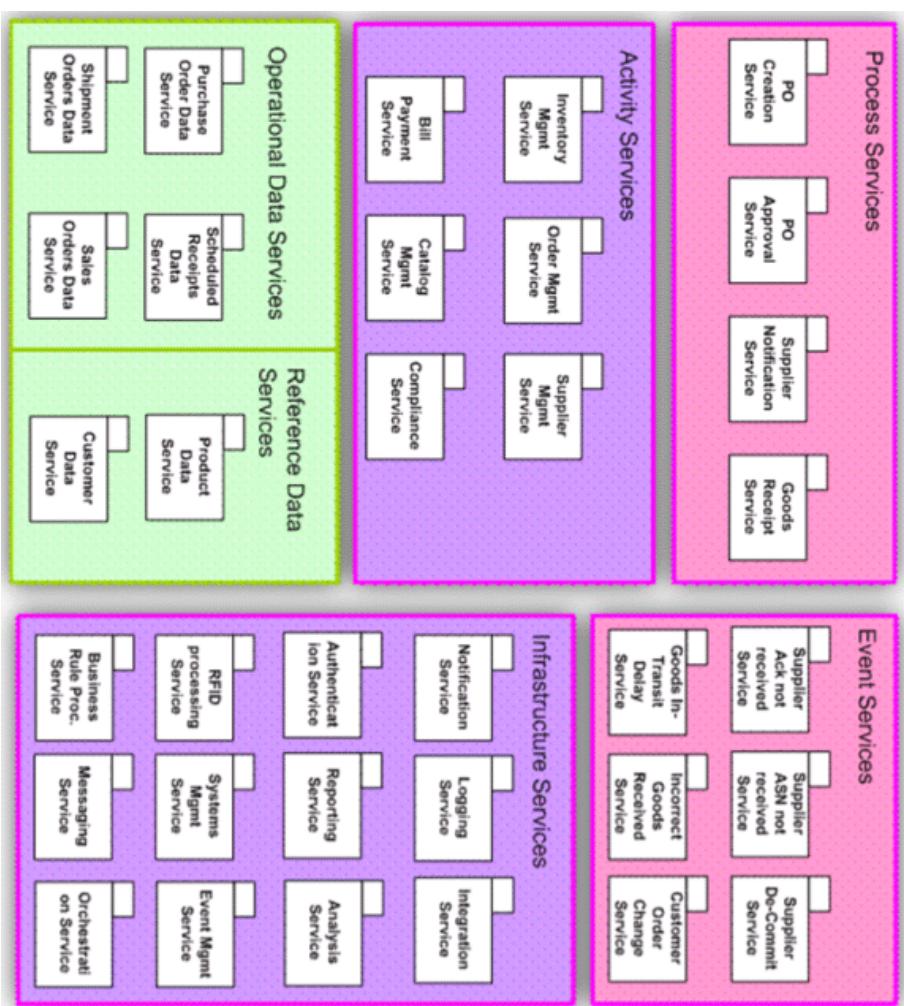
Systems Analysis and Specification

prof. dr. ir. Pieter Simoens



Hello
World!

Real-World software



possibly (likely?) distributed
continuously evolving
developed by a team

... And subject to many challenges!

Latency costs money

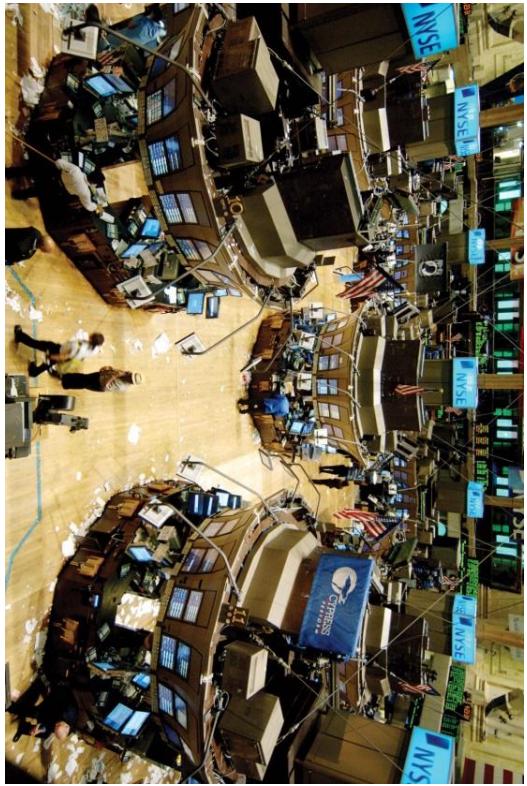


Every 100 ms of latency costs 1 % in sales



400 ms of additional latency to show 30 instead of 10 results decreases searches/user by 0.76 %

A lot of money...



in high-frequency trading, each ms is
worth \$ 100 million per year

a technology firm spent \$ 300 million to realize

a straight line and shave 3 ms of the

communication time with the NASDAQ server

engineers add extra lengths of cables of a few feet

to equalize the runs among all algo-computers
inside the room with the central NYSE server

Go viral... not down



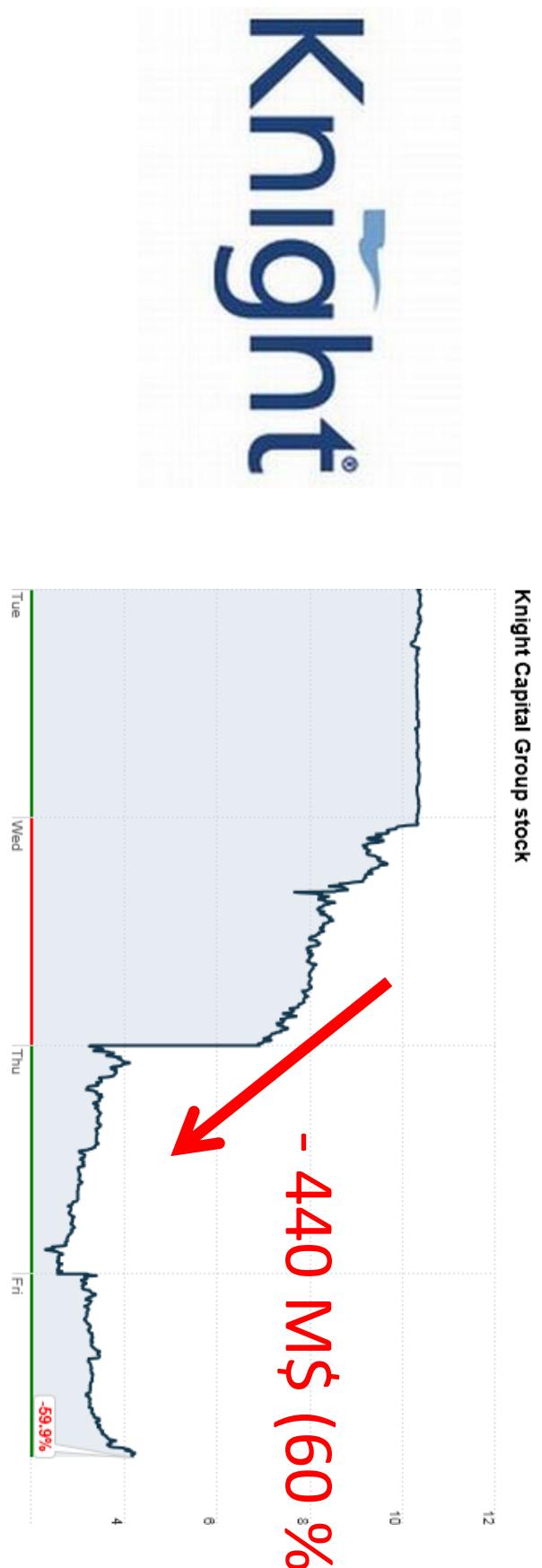
Realo.com

Wie de [site](#) zelf eens wil uitproberen, kan daarbij moeilijkheden ondervinden. Coppens benadrukt dat de site niet offline is gehaald, maar net problemen ondervindt door alle media-aandacht. De site zit namelijk nog in een testfase. "Maar we proberen de problemen zo snel mogelijk op te lossen."

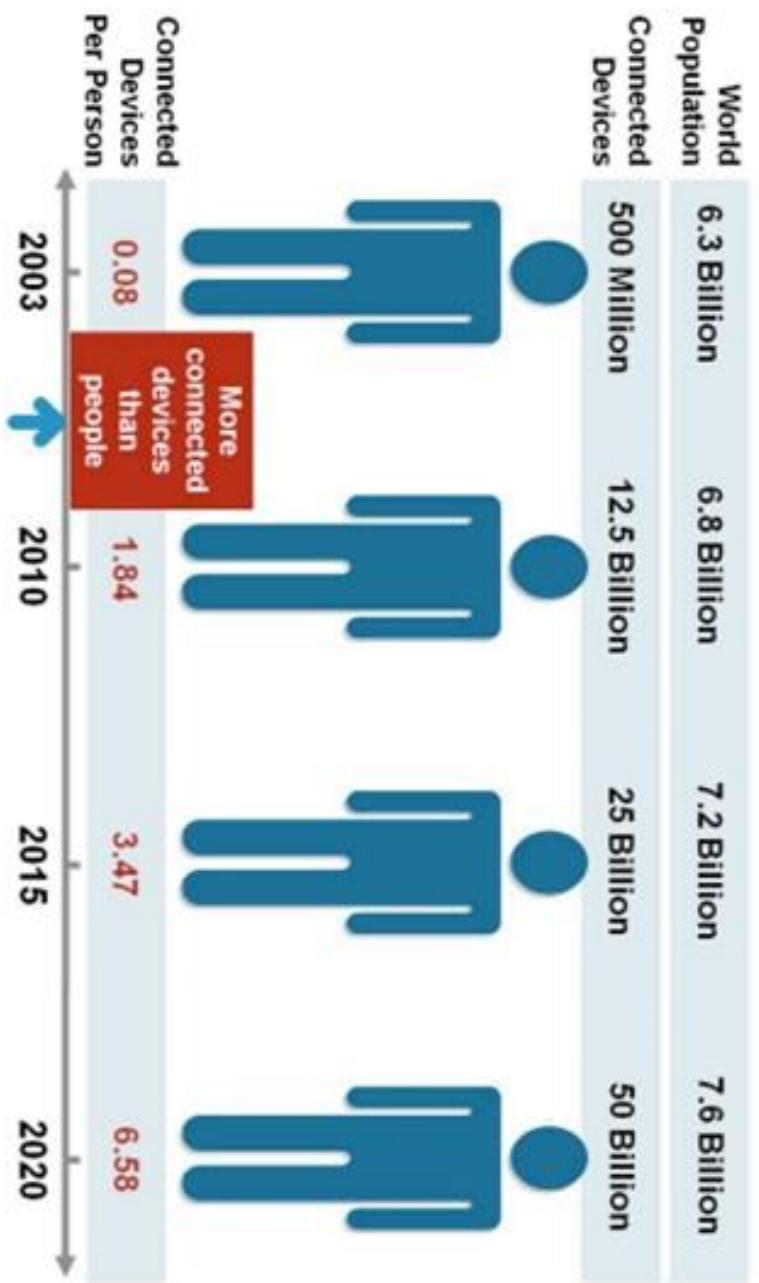


Eiffel Tower's website crashed on March 31 when featured in a Google Doodle

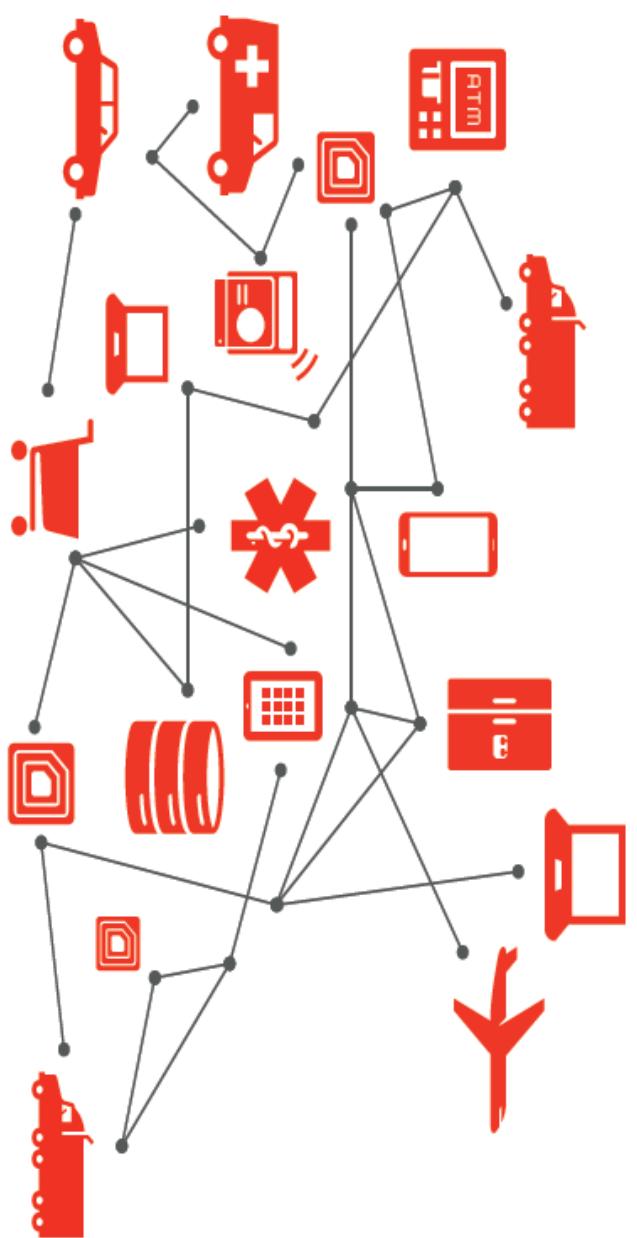
Undisciplined DevOps



- Deploy new code in algorithmic router that sends orders into the stock market
- Update was intended to replace unused code in the order router
- Knight stopped using the old code in 2003
- Repurposed a flag formerly used to active the unused code
- New code was only copied to 7 of the eight servers, so one server executed the old code (flag repurposing!)



Connecting people and things has enabled new types of applications



Big Data, IoT, mobile...

Deliver on time, on budget, on value

IT executives identify 4 groups of issues that cause most project failures.

Rough distribution by cause of the 45% of IT projects that experience cost overruns (for those with budgets >\$15 million in 2010 dollars), %

Missing focus

- Unclear objectives
- Lack of business focus

Content issues

- Shifting requirements
- Technical complexity

Skill issues

- Unaligned team
- Lack of skills

Execution issues

- Unrealistic schedule
- Reactive planning

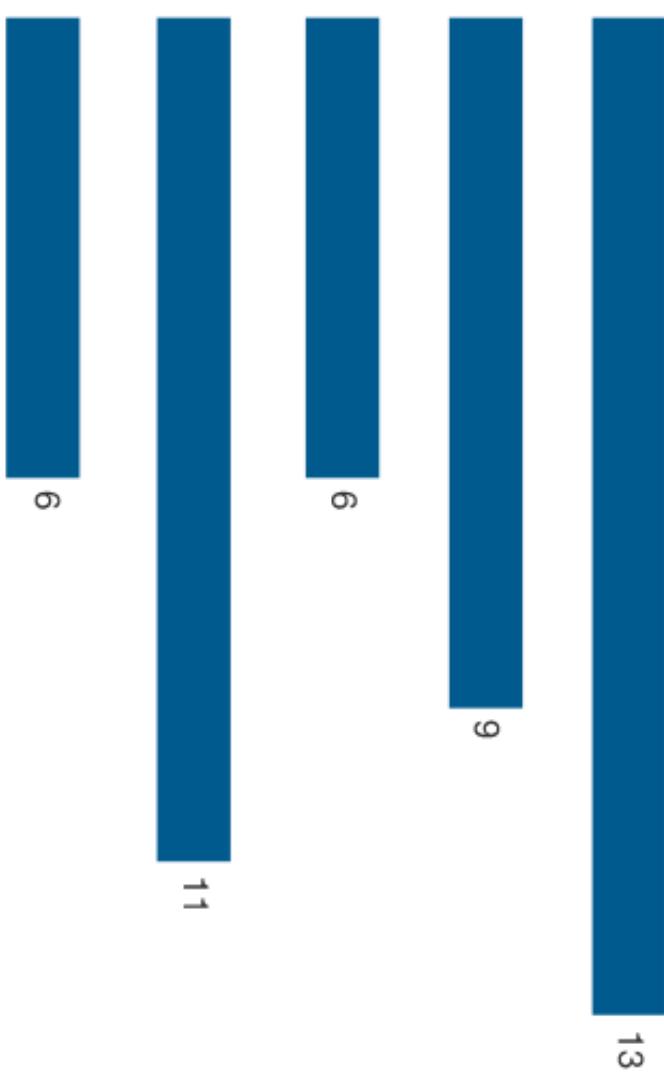
Unexplained causes

IT projects with budgets >\$15 million

Cost overrun, 45%

Schedule overrun, 7%

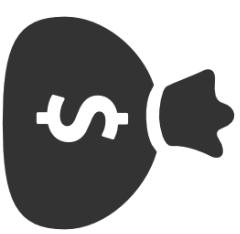
Benefits shortfall, -56%



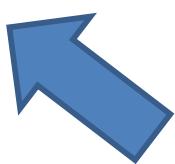
The programmer today...



business requirements
(latency, stability, world-scale problem)



skilled individual



... in an organized team



Pick the appropriate tools:

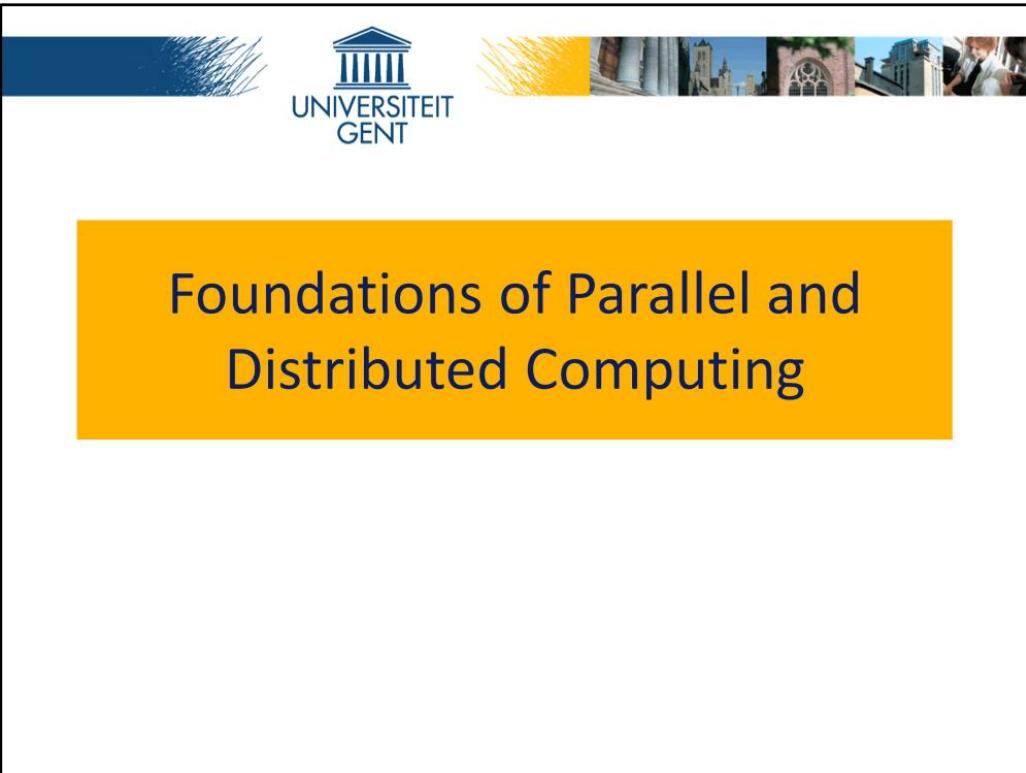
- * cloud computing
- * big data processing and storage
- *
- ...

Organize the work:

- DevOps
- Project Management
- ...

Systeemanalyse en -specificatie

- Course content:
 - Central question: how to develop complex software?
 - Topics:
 - fundamentals of parallel and distributed computing
 - big data processing and storage
 - cloud computing
 - project management and software development
- Weekly, 15:45 – 17:45 @ B2.018
 - mixture of ex cathedra, lab sessions and guest lectures
- Exam
 - written



Parallel vs Distributed

Parallel

focus on **fast** solving compute-intensive
(large) problems
multi-core CPU, supercomputers
tight coordination between processors
shared memory



Distributed

focus on information and resource sharing
P2P, client-server, cloud computing
loose coupling between processes
message passing



The difference between parallel and distributed computing is often vague. Many systems leverage on principles from both paradigms simultaneously. In broad terms, the goal of parallel processing is to employ all processors to perform one large task. In contrast, each processor in a distributed system generally has its own semi-independent agenda, but for various reasons, including sharing of resources, availability, and fault tolerance, processors need to coordinate their actions.

Parallelism is generally concerned with accomplishing a particular computation as fast as possible, exploiting multiple processors. The scale of the processors may range from multiple arithmetical units inside a single processor, to multiple processors sharing memory, to distributing the computation on many computers. On the side of models of computation, parallelism is generally about using multiple simultaneous threads of computation internally, in order to compute a final result.

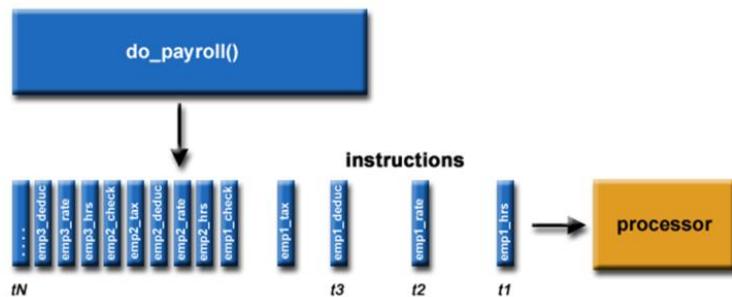
Distributed computing studies separate processors connected by communication links. Whereas parallel processing models often (but not always) assume shared memory, distributed systems rely fundamentally on message passing. Individual entities are often loosely coupled. The model is automatically used when independent entities cooperate (e.g. a web shop provider contacts a payment provider when a customer is finalizing his order). It is also used in cloud computing, where different application components reside in separate virtual environments. Failure (of processor nodes or communication links) is a normal situation.

Outline

- Parallel computing
 - serial vs parallel
 - the need for speed: why parallelism?
 - theoretical boundaries to maximum speed-up
- Distributed computing
 - what is distributed computing
 - CAP theorem and its extensions
 - RAFT consensus algorithm

Serial computing

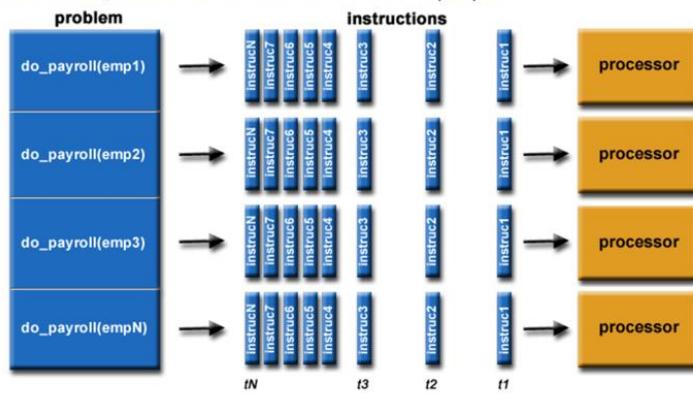
- Discrete series of instructions executed sequentially
- Executed on a single processor
- Only one instruction executed at any time



Parallel computing

Simultaneous use of multiple compute resources
to solve a computational problem

- Problem decomposed in parts that can be solved concurrently
- Instructions from each part execute simultaneously and are combined afterwards
 - Note: concurrent computing: IPC *during* task execution
- An overall control/coordination mechanism is employed



Why parallelism?

- ... it incurs additional troubles:
 - need to rewrite programs
 - synchronization and control
- Three walls to serial performance (=single CPU speed)
 - Memory wall
 - ILP wall
 - Power wall
- Not walls, but increasingly steep hills to climb

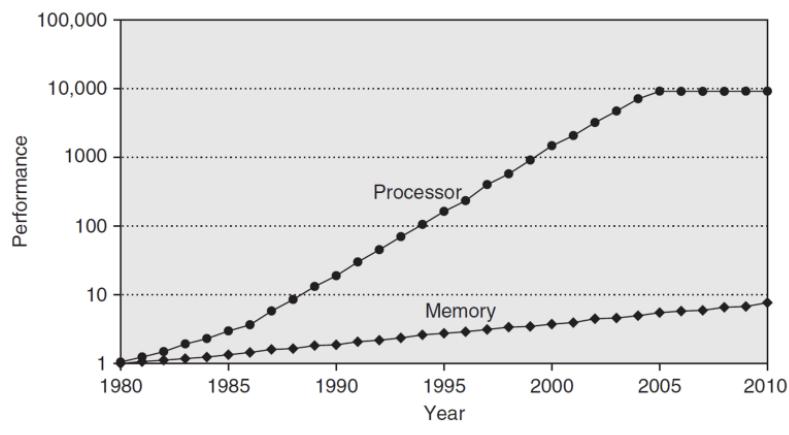


For general-purpose processors, much of the motivation for multi-core processors comes from greatly diminished gains in processor performance from increasing the operating frequency. This is due to three primary factors:

- The *memory wall*; the increasing gap between processor and memory speeds. This, in effect, pushes for cache sizes to be larger in order to mask the latency of memory. This helps only to the extent that memory bandwidth is not the bottleneck in performance.
- The *ILP wall*; the increasing difficulty of finding enough parallelism in a single instruction stream to keep a high-performance single-core processor busy.
- The *power wall*; the trend of consuming exponentially increasing power with each factorial increase of operating frequency. This increase can be mitigated by “shrinking” the processor by using smaller traces for the same logic. The *power wall* poses manufacturing, system design and deployment problems that have not been justified in the face of the diminished gains in performance due to the *memory wall* and *ILP wall*.

Memory speed

Growing disparity between CPU and memory outside the CPU chip



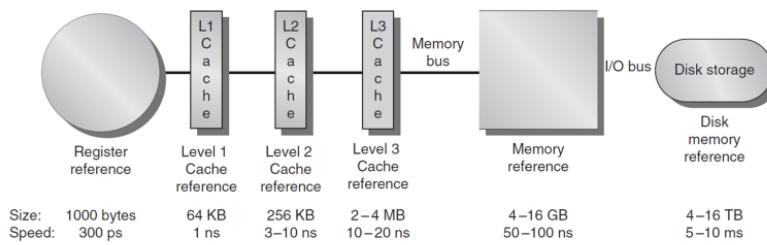
© Computer Architecture (5th edition), Hennessy

The advance in performance of processors has increased the importance of the memory wall. The graph plots single processor performance projections against the historical performance improvement in time to access main memory. The processor line shows the increase in memory requests per second on average (i.e., the inverse of the latency between memory references), while the memory line shows the increase in DRAM accesses per second (i.e., the inverse of the DRAM access latency).

Starting with 1980 performance as baseline, the gap in performance as the difference in time between processor memory requests (for a single processor or core) and the latency of a DRAM access is plotted over time. Note that the vertical axis is on a logarithmic scale.

Memory wall

- A system (CPU-memory duo) can't move at CPU top speed
- Improvements?
 - memory bandwidth: more pins (32 → 64 bit →...)
 - memory latency: even more challenging!
 - caches are fast, but limited in size



Even when the clock frequency of CPU cores keeps improving, we cannot fully utilize this because of the memory wall. The memory wall is caused by two factors:

- **Memory bandwidth:** how much data can be transferred per second between CPU and memory. Improving the bandwidth requires to add more “pipes”, e.g. more pins that come out of the chip for the DRAM for example. This is challenging in terms of technology since the real estate on a chip is limited.
- **Memory latency:** the amount of time it takes for an operation to complete. For example: the time needed for a CPU to complete a 32 bit write/read operation (to off-chip memory). This is even harder to improve than bandwidth. While improving bandwidth can (in principle) be realized by doing more of the same; improving latency requires breakthroughs in material sciences (e.g. optical communication).

Instruction Level Parallelism wall

- Basic idea
 - Overlap execution of independent instructions
 - Work on many instructions during the same clock cycle
- How?
 - Instruction pipelining, out-of-order execution, speculative execution, superscalar execution
- But... acceleration using ILP is plateauing
 - You need large blocks of instructions that can run in parallel
 - “speculation” success difficult to predict
 - super-linear increase in complexity without linear speedup

```
for(int i=0; i < 1000; i++)  
    x[i] = x[i] + y[i];
```

```
e = a + b;  
f = c + d;  
g = e*f;
```

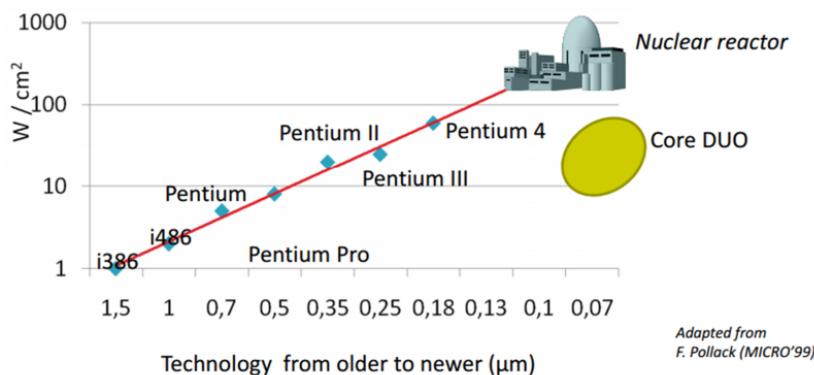
Instruction Level Parallelism aims to increase the throughput of the number of instructions executed per unit of time (clock cycle) by a single core. Various techniques are used:

- Instruction pipelining: divide each instruction into series of sub-steps (micro-operations) so that the execution of multiple instructions can be partially overlapped
- Out-of-order execution: instructions execute in any order but without violating data dependencies
- Speculative execution: allows the execution of complete instructions or parts of instructions before being sure whether this execution is required
- Superscalar execution: multiple execution units are used to execute multiple instructions in parallel

However, the upper limits of the acceleration through ILP exploitation are almost reached. First, for ILP to be efficient, you need large blocks of instructions that can be [attempted to be] run in parallel. This is not always the case and you cannot go beyond your critical path (see later). Second, ILP is often based on speculation: if you guessed wrong, you throw away that part of your result. Third, data dependencies may prevent successive instructions from executing in parallel, even if there are no branches.

Power wall

- Static power consumption
 - leakage per transistor
 - worse as transistor gates get smaller
- Dynamic power consumption
 - proportional to clock frequency

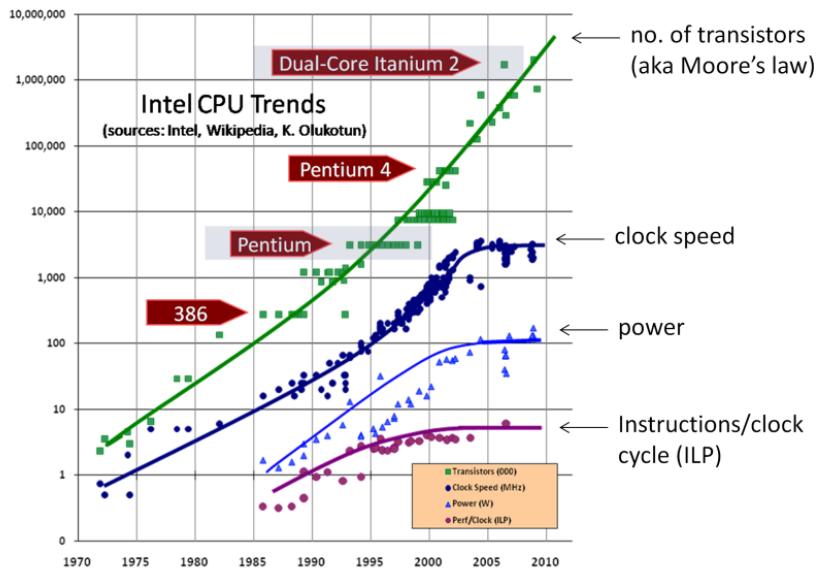


Power, and not manufacturing, limits improvements in the speed of individual cores. The power consumption of a core is the sum of two contributions:

- **Static power consumption:** this is mostly leakage, it is the power dissipated by a transistor whose gate is intended to be off. Leakage power dissipation gets worse as transistor gates get smaller, because gate dielectric thickness must proportionally decrease.
- **Dynamic power consumption:** this is the part that is related to the amount of work executed by a single core. The dynamic power consumption is proportional to (a.o.) the clock frequency.

For this reason, the increase in INTEL CPU clock speed was already stopped in 2007. From then on, they switched to multiple cores.

Sequential execution has lost steam



The three walls discussed in the previous slides indicate that the model of sequential execution has lost momentum. As chip geometries shrink and clock frequencies rise, the transistor leakage current and the dynamic power increases, leading to excess power consumption and heat. Also, the advantages of higher clock speeds are in part negated by memory latency, since memory access times have not been able to keep pace with increasing clock frequencies.

On the above graph, there is however one bright spot: the number of transistors per unit area is still increasing. This is the well-known Moore's law.



If one ox could not do the job they did not try to grow a bigger ox, but used two oxen.

When we need greater computer power, the answer is not to get a bigger computer, but...to build systems of computers and operate them in parallel.

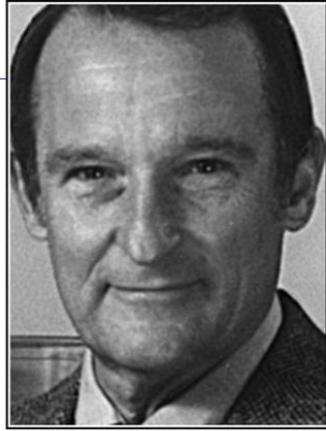
(Grace Hopper)

izquotes.com

Sacrificing uniprocessor performance for power savings can save you a lot

Example:

- Scenario One: one-core processor with power budget W
 - Increase frequency/ILP by 20%
 - Substantially increases power, by more than 50%
 - But, only increase performance by 13%
- Scenario Two: Decrease frequency by 20% with a simpler core
 - Decreases power by 50%
 - Can now add another core (one more ox!)



If you were plowing a field, which would you rather use? Two strong oxen or 1024 chickens?

— Seymour Cray —

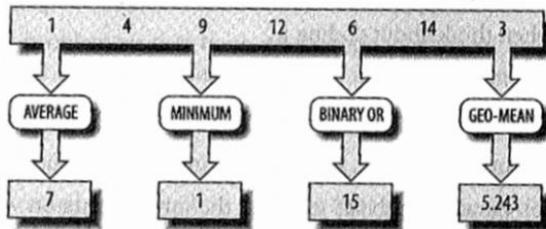
AZ QUOTES

For certain classes of applications (not including field plowing...) you can run many cores at lower frequency and come ahead (big time) at the speed game.

Types of parallelism

- Task parallelism
 - entirely different calculations on either the same or different sets of data
 - allocate subtasks to a processor
- Data parallelism
 - same calculation is performed on the same or different sets of data
 - allocate subset of data to process
- Pipelining: hybrid data/task parallelism
 - a parallel pipeline of tasks, each of which might be data parallel
 - e.g. image processing

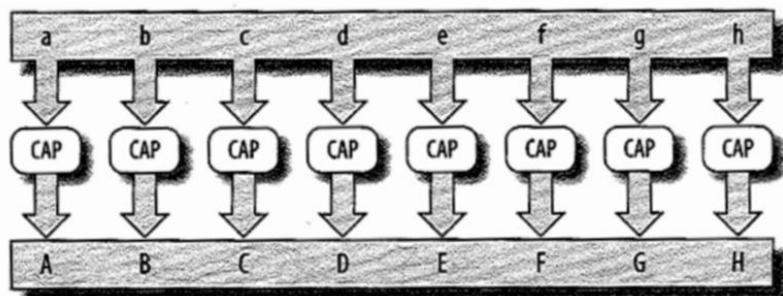
Task parallelism



- Several functions on the same data
- No dependencies between tasks, all can run in parallel

An example of task parallelism is shown here. Several functions are calculated on the same data: average, minimum, binary or geometric mean.

Data parallelism



- Can divide parts of the data between different tasks and perform the tasks in parallel
- Key: no dependencies between the tasks that cause their results to be ordered

This is an example of data parallelism: all characters in a text file must be converted to upper case.

Pipeline parallelism

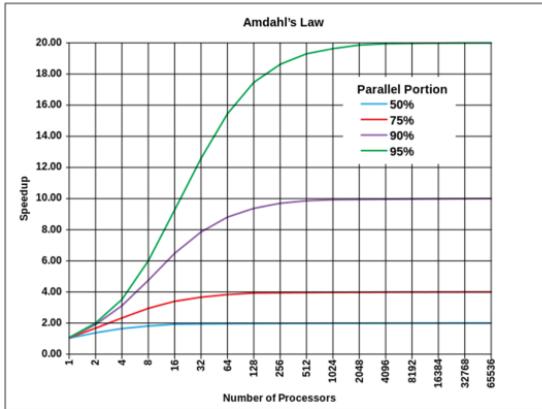


- Output of one task is the input to the next one
- Each task can run in parallel
- Throughput impacted by the longest-latency element in the pipeline

Amdahl's law

Determines the maximum speed-up S by dividing an amount of work with parallel fraction P over N units

$$S \leq \frac{1}{(1-P) + P/N}$$

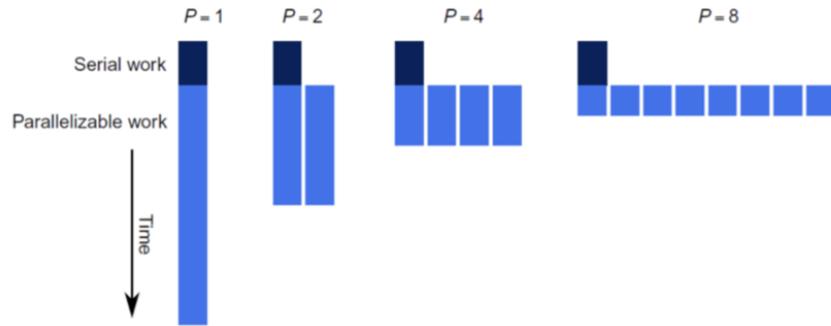


Speedup is limited by the time needed for the sequential fraction of the program

Amdahl's law prescribes the upper limit on the speedup that can be achieved by dividing a program (work load) over N processors.

Each program can be divided into a fraction P that is parallelizable and a portion $(1-P)$ that is intrinsically serial. The time needed to execute the serial portions is unaffected by the number of parallel processing units. The law gives an upper bound that can only be realized by ignoring the overhead due to message passing, gathering of results, etc...

Philosophical question

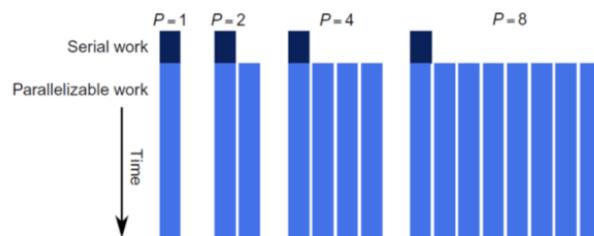


If the maximum speed-up is bounded by Amdahl's law, then why are we building supercomputers with 1000s of nodes?

Gustafson's law

- As processor power increases, programmer's tend to increase the size of the problem
 - set the problems to be solved within a practical fixed time
 - e.g. more pixels, larger scale, smaller time steps...
- Key assumption: total amount of work varies linearly with the number of processors

$$T_{seq} = a + b \cdot N$$



Gustafson's law

$$S = \frac{T_{seq}}{T_{par}} = \frac{a+b \cdot N}{a+b}$$

$$S = N - \alpha \cdot (N - 1) \quad \alpha = a/(a + b)$$

- α = sequential fraction of the total execution time on a parallel session
- if α is small, then the speed-up is almost linear with the amount of processors

Gustafson's law indicates that when the available computational power increases with a factor N , you can handle workloads that are similarly scaled while keeping the total time constant.

The difference between Amdahl's law and Gustafson's law lies in whether you want to make a program run faster with the same workload or run in the same time with a larger workload. History clearly favors programs attacking and solving larger, more complex problems, so Gustafson's observations fit the historical trend. Nevertheless, Amdahl's Law still haunts you when you need to make an application run faster on the same workload to meet some latency target.

Bibliography: parallel computing

- D. Ernst, Introduction to Accelerators and GPGPU
- D. Negrut, High Performance Computing for Engineering Applications,
<http://sbel.wisc.edu/Courses/ME964/2013/Lectures/lecture0918.pdf>
- B. Barney, Introduction to Parallel Computing
https://computing.llnl.gov/tutorials/parallel_comp
- <http://www.cs.umd.edu/class/fall2013/cmsc433/lectures/concurrency-basics.pdf>
- M. Gillespie, Amdahl's Law, Gustafson's Trend, and the Performance Limits of Parallel Applications
- <http://www.drdobbs.com/parallel/amdahls-law-vs-gustafson-barsis-law/240162980?pgno=2>

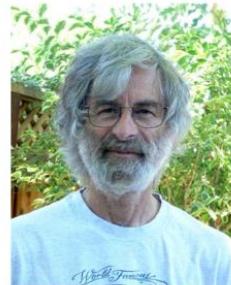
Outline

- Parallel computing
- Distributed computing
 - what is distributed computing
 - CAP theorem and its extensions
 - RAFT consensus algorithm

Distributed system

“A distributed system is one where I can’t get the job done
because a computer I never heard of drashed”

Leslie Lamport



Distributed system

Distributed system

System where hardware and software components are located at networked computers. Components communicate and coordinate their actions only by passing messages.

“The network is the computer”

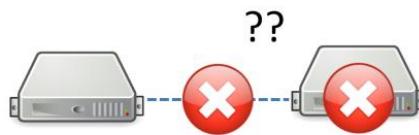
WWW, grids,
P2P, cloud computing



Consequences of distributed systems



no obvious spatial limit to the extent of the system



failure of remote node or of interconnecting network?



no global time notion



inconsistency due to partial failures and/or concurrent execution

There are a number of natural consequences of distributed systems.

- 1) There is no spatial limit on the number of nodes in the distributed systems. Nodes of a single distributed system can be located on the same server, on different servers in the same datacenters, or on servers in different locations.
- 2) When communication with a remote node fails, it is hard for the requesting node to decide if the failure is the consequence of a network disruption or a crash of the node itself.
- 3) Each node of the distributed system has its own clock. This means you can not use locally generated timestamps for ordering distributed events.
- 4) Components are likely to execute concurrently. Partial failures are likely to happen, while the others continue to work. This results in inconsistent views and inconsistencies between the nodes.

Fallacies of distributed computing

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology does not change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous



There is no such thing
as a free lunch.

The fallacies of distributed computing are a set of assumptions that programmers new to distributed applications invariably make. These assumptions ultimately prove false, resulting either in the failure of the system, a substantial reduction in system scope, or in large, unplanned expenses required to redesign the system to meet its original goals.

1. The network is reliable – Hardware failures can never be excluded: power failures, someone tripping over a network cord, etc... On the infrastructure side, this means you need to think about hardware redundancy and weigh the risks of failure vs. the required investment (e.g. buying a spare switch). On the software side, you need to think about messages/calls getting lost. For one, you can use full reliable messaging.
2. Latency is zero – Latency can be relatively good on a LAN but deteriorates quickly when you move to WAN scenarios or internet scenarios. Even when you work on a LAN with Gigabit Ethernet you should bear in mind that the latency is much bigger than accessing local memory. Assuming the latency is zero you can be easily tempted to assume making a call over the wire is almost like making a local call. This means you should strive to make as few as possible calls and assume you have enough bandwidth to move data in/out in these calls.
3. Bandwidth is infinite – Although network bandwidth is continuously growing, so does the amount of information we try to squeeze through it. Also, in the WAN,

packet loss and latency will limit your throughput (e.g. over TCP). While the previous fallacy stimulates to send fewer, but larger messages, this fallacy states that you should strive to limit the size of individual messages as well.

4. Network is secure – Build security in your software from Day 1. Security is usually a multi-layered solution that is handled on the network, infrastructure and application levels.
5. Topology does not change – Usually, the network topology is out of your control. The cloud infrastructure provider may add and remove servers or make other changes to the network. Try not to depend on specific endpoints or routes. Provide location transparency (e.g. multicast), discovery services and abstract the physical structure of the network (e.g. using DSN names instead of IP addresses)
6. There is one administrator – Your company might collaborate with external entities (e.g. hosting provider), your application might consume external services, etc... The DevOps strategy can help, because system administrators then become part of the development team. But with external parties, this is not always possible.
7. Transport cost is zero – Marshaling (serializing) information across the different network stack layers takes computer resources and adds to the latency. The second way to interpret this statement is that the monetary cost for running a network is not free.
8. The network is homogeneous – Most networks are not homogeneous. Even in a home network, you might encounter Linux devices, Windows PCs, NAS servers and mobile devices. At the network layer, this fallacy does not cause too much trouble since every device will talk IP. At the application level, you have to assume interoperability will be needed sooner or later. Do not rely on proprietary protocols and use standard technologies.

Why distributed systems?

- Many challenges
 - difficult to manage
 - no global time notion
 - no global state
 - almost always concurrent execution
 - (partial) failures likely to will to happen
- Many advantages
 - resource sharing
 - information (e.g. music files, group documents)
 - hardware (e.g. collective storage system, cloud infrastructure)
 - scalability: “easy” to adapt to larger user base
 - fault tolerance: “easy” to cope with failures

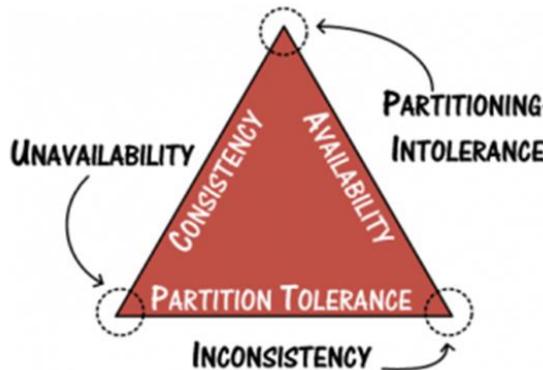
Despite the intrinsic difficulties you will have to cope with in a distributed system, there are many very important advantages.

- It allows for resource sharing: you can access information and/or services provided by 3rd parties
- It allows for scalability: if user demand increases, you can add more servers to distribute computation and data
- It allows for redundancy: replicating data or computation on multiple nodes allows that the overall system stays alive, even during failures.

Outline

- Parallel computing
- Distributed computing
 - what is distributed computing
 - CAP theorem and its extensions
 - Definition
 - Coping with CAP in real-world applications
 - RAFT consensus algorithm

CAP Theorem



© StackExchange

- describes trade-offs involved in distributed system
- impossible to provide the following *three guarantees at the same time*
 - Consistency: all nodes see the same data at the same time
 - Availability: all non-failing nodes are available for queries
 - Partition-tolerance: underlying system can be split in non-communicating groups

The CAP theorem was conjectured by Prof. Eric Brewer during a keynote talk at the PODC (Principle of Distributed Computing) conference in 2000. In 2002, other authors published a formal proof of Brewer's conjecture, rendering it into a theorem.

The theorem states that it is impossible for a distributed computer system to provide all three of the following guarantees:

- Consistency: all nodes see the same data at the same time. In other words, each server gives the correct response to each request
- Availability : a guarantee that every request receives a response about whether it succeeded or failed. Alternative formulations: every request received by a non-failing node in the system must result in a response; all (non-failing) nodes are available for queries)
- Partition tolerance: this refers to the underlying system and not to the service. Servers partitioned to groups that are not able to communicate. The system continues to operate despite arbitrary partitioning due to network failures, at least one partition should remain functioning and accessible to the clients.

CAP Theorem

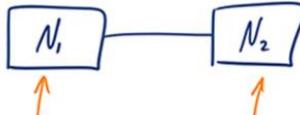
Distributed system = a collection of interconnected nodes that share data

Consistency



All nodes see the same data at the same time

Availability



Every request gets a response on success/failure, regardless of the state of any individual node

Partition Tolerance

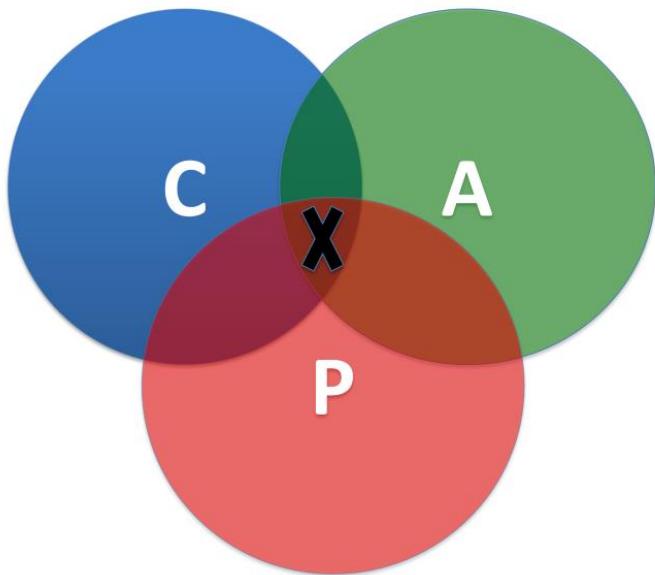


The system continues to function despite message loss between nodes

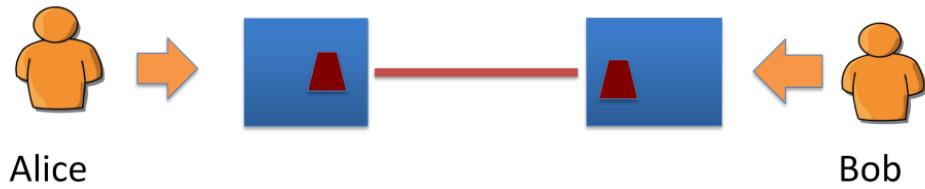
This is a more visual view on the three elements in the CAP theorem.

A distributed system can satisfy any two of these guarantees at the same time **but not all three.**

CAP theorem



Example: hotel booking



How does a system look like
that is AP, CP or CA?

We will use a hotel booking system as a running example to illustrate the CAP theorem. The hotel booking system comprises two servers at physically different locations. Alice and Bob are simultaneously trying to book a room in a hotel. Of course, we want to avoid that a room is double-booked.

(note: the “last” hotel room is simply a metaphor for a specific data item that is read/writable by multiple users)

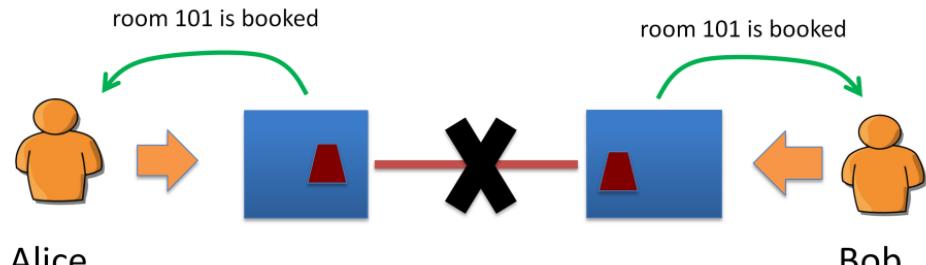
C+A



In a distributed system, you cannot **not** choose partition tolerance

CA is only possible as a monolithic, single server database. As soon as you have two or more servers, you introduce a network link and you basically *require* “partition tolerance”. In a distributed system, you cannot **not** choose partition tolerance.

P+A

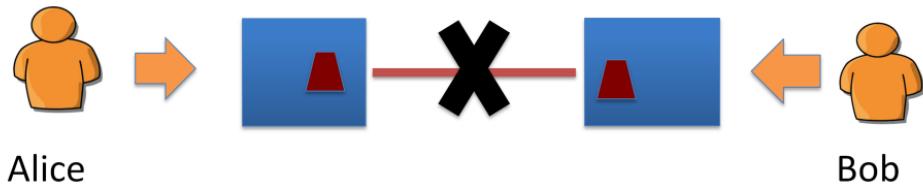


All rooms can be booked and servers work independently

Not consistent: view of both servers on room availability can be out-of-sync!

As soon as we have more than 1 server, we tolerate partitions. According to the CAP theorem, this means we must choose between A and C. In a P+A model, the system keeps functioning even during a network partition, or when a node fails. Clearly, consistency is sacrificed since a user can book a room on one server without the other server knowing this.

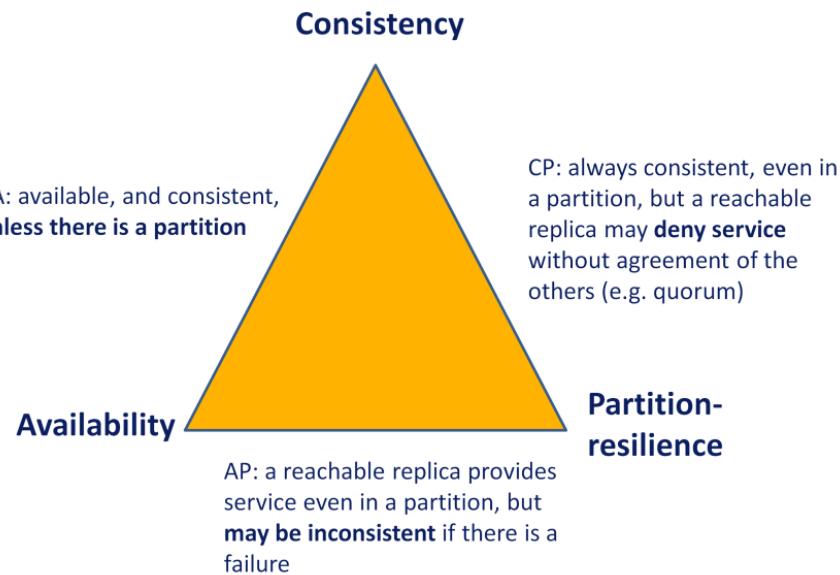
P+C



If we require consistency, then there is no other option than take at least one of the two servers offline during a partition, although all servers are working properly

If we don't allow the inconsistency of the previous slide, then we must take at least one of the two servers offline during network partitions.

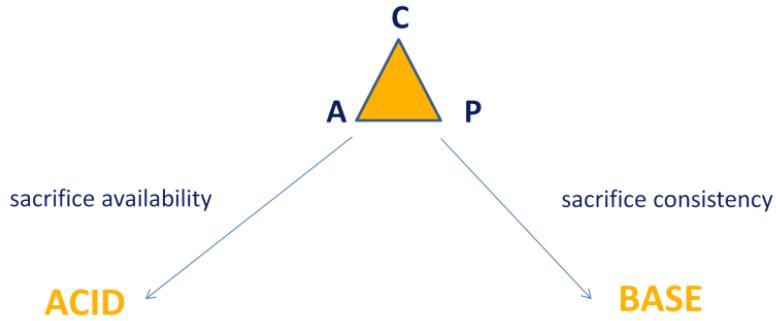
CA/CP/AP systems



In practice, CAP takes place during a timeout. Then a decision should be made:

- Cancel the operation and thus decrease availability
- Continue the operation and be prone to inconsistency

Database transactions: ACID vs BASE



Atomicity: operation is performed on **all** replicas or is not performed on any of them
Consistency: after each operation all replicas reach the same state
Isolation: no operation can see the data from another operation in an intermediate state
Durability: successful writes persist infinitely

Basic Availability: there will be a response to any request, but data may be inconsistent
Soft state: state of system changes over time, even when there is no input
Eventually consistent: data will propagate to all replicas sooner or later

Any realistic data scaling system will have to scale, either horizontally (sharding, replicating same data over multiple nodes) or vertically (splitting data according to their function, e.g. storing user data on one node, product information on another, etc.). Of course, horizontal and vertical scaling can (will) be combined, but in any case we end up with a distributed database system that is subject to the CAP theorem. This means that when a partition occurs, you have to choose between availability and consistency when starting a new transaction with the database.

If we choose consistency, then *database transactions* follow the ACID pattern:

Atomicity. All of the operations in the transaction will complete, or none will.

Consistency. The database will be in a consistent state when the transaction begins and ends.

Isolation. The transaction will behave as if it is the only operation being performed upon the database.

Durability. Upon completion of the transaction, the operation will not be reversed.

Database vendors introduced a technique known as 2PC (two-phase commit) for providing ACID guarantees across multiple database instances. The protocol is broken into two phases:

- First, the transaction coordinator asks each database involved to precommit the operation and indicate whether commit is possible. If all databases agree the commit can proceed, then phase 2 begins.

- The transaction coordinator asks each database to commit the data.

If any database vetoes the commit, then all databases are asked to roll back their portions of the transaction. This way, we are getting consistency across partitions. But if one of the nodes goes down, any transaction will fail. Essentially, a transaction involving two database nodes in a 2PC commit will have the availability of the product of the availability of each database. For example, if we assume each database has 99.9 percent availability, then the availability of the transaction becomes 99.8 percent, or an additional downtime of 43 minutes per month.

On the other hand, if we sacrifice consistency and prefer availability, the processing of transactions follows the BASE acronym:

Basically Available: This constraint states that the system does guarantee the availability of the data as specified by the CAP Theorem -- there will be a response to any request. But, that response could still be a 'failure' to obtain the requested data or the data may be in an inconsistent or changing state.

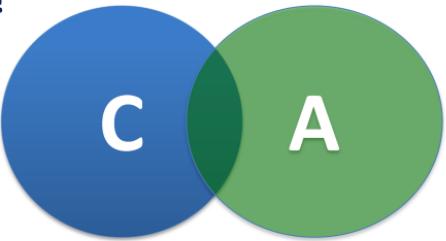
Soft state: The state of the system could change over time, so even during times without input there may be changes going on due to 'eventual consistency,' thus the state of the system is always 'soft.'

Eventual consistency: The system will *eventually* become consistent once it stops receiving input. The data will propagate to everywhere it should sooner or later, but the system will continue to receive input and is not checking the consistency of every transaction before it moves onto the next one.

COPING WITH CAP IN REAL-WORLD SYSTEMS

A popular misconception: 2 out 3

- How about distributed CA?
- Can a distributed system (with unreliable network) really be not tolerant of partitions?



CAP Theorem 12 year later

- Prof. Eric Brewer: father of CAP theorem
 - “The “2 of 3” formulation was always **misleading** because it tended to oversimplify the tensions among properties. ...
 - **CAP prohibits only a tiny part of the design space:** *perfect availability and consistency in the presence of partitions*, which are rare.”



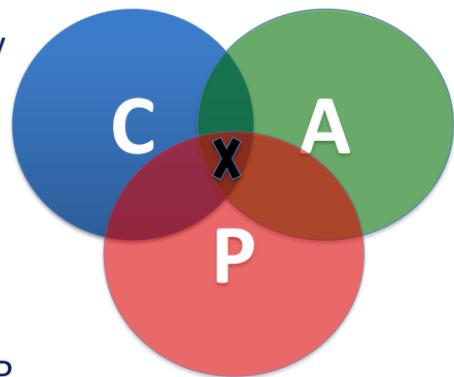
<http://www.infoq.com/articles/cap-twelve-years-later-how-the-rules-have-changed>

The CAP theorem has been victim of its own success and is subject to some misconceptions:

- The popular belief is that CAP means *pick any two*. The confusion is about the existence of CA systems, which pretend that partition tolerance is optional, or claim that partitions don't happen. In reality, you can't sacrifice partition tolerance because partitions happen in real large-scale systems all the time.
- The second misconception is that the CAP theorem means *you can't be consistent and available during partitions*. That's not true. Specifically, the CAP theorem only prevents *everybody* from being consistent and available, not *anybody* (some literature calls this *always available*). It doesn't prevent clients and replicas on the majority side of simple partitions from making progress, and experiencing both consistency and availability. There are other restrictions on this, but they aren't CAP restrictions.
- The third misconception is that the *consistency* in CAP is all or nothing, and that you can't offer any consistency guarantees at all during partitions. In reality, many very useful consistency models can be offered on all sides of a partition. Implementation tricks like session stickiness and client-side caching can allow systems to offer useful models like *read your writes*, *monotonic reads* and even *causal consistency*.

Consistency or Availability

- Consistency and Availability is not “binary” decision
- AP systems relax consistency in favor of availability – but are not inconsistent
- CP systems sacrifice availability for consistency- but are not unavailable
- This suggests both AP and CP systems can offer a degree of consistency, and availability, as well as partition tolerance



The practical implications of the CAP theorem are not as stringent as they appear, because in practical situations you are often allowed to relax (but not give up completely) consistency or availability.

Instead of actually having to choose between a CP or AP system, the important message of the CAP theorem is that you have to balance between Consistency with Availability.

CP: Best Effort Availability

- guarantees consistency, regardless of network behavior
- when communication is typically reliable
 - E.g. servers in same datacenter, partitions are rare (but not impossible)
- Example:
 - Majority protocols
 - Distributed Locking (Google Chubby Lock service)
- Trait:
 - Pessimistic locking
 - Make minority partition unavailable

AP: Best Effort Consistency

- Sometimes being unavailable is not an option
- Inconsistency is not a major problem
 - Best effort for up-to-date data
 - No assurance that all users get the same content
- Example:
 - Web Caching
 - DNS
- Trait:
 - Optimistic
 - Expiration/Time-to-live
 - Conflict resolution

Types of Consistency

- Strong Consistency
 - After the update completes, **any subsequent access** will return the **same** updated value.
- Weak Consistency
 - It is **not guaranteed** that subsequent accesses will return the updated value.
- **Eventual Consistency**
 - Specific form of weak consistency
 - It is guaranteed that if **no new updates** are made to object, **eventually** all accesses will return the same updated value (e.g., *propagate updates to replicas in a lazy fashion*)
 - in principle, does not impose global ordering

Eventual consistency is a consistency model, which is used in many large distributed databases. Such databases require that all changes to a replicated piece of data eventually reach all affected replicas. The storage system guarantees that if no new updates are made to the object, eventually all accesses will return the last updated value.

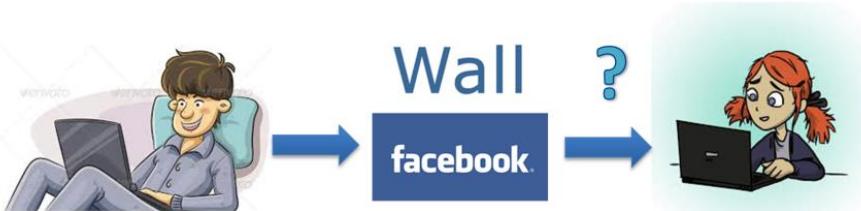
Consider a case where data item R=0 on all three nodes. Assume that we have the following sequence of writes and commits: W(R=3) C W(R=5) C W(R=7) C in node 0. Now read on node 1 could return R=5 and read from node 2 could return R=7. This is eventually consistent as long as eventually read from all nodes return the same value. Note that this final value could be R=5. **Eventual consistency does not restrict the order in which the writes must be executed.**

Definition of Eventual consistency:

- * **Eventual delivery:** An update executed at one node eventually executes at all nodes.
- * **Termination:** All update executions terminate.
- * **Convergence:** Nodes that have executed the same updates eventually reach an equivalent state (and stay).

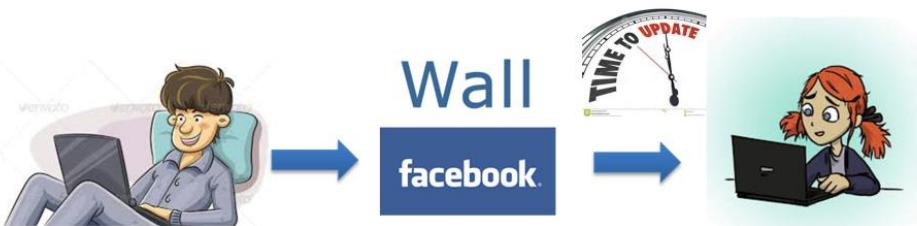
Eventual Consistency: Facebook

- Bob finds an interesting story and shares with Alice by posting on her Facebook wall
- Bob asks Alice to check it out
- Alice logs in her account, checks her Facebook wall but finds:
 - **Nothing is there!**



Eventual Consistency: Facebook

- Bob tells Alice to wait a bit and check out later
- Alice waits for a minute or so and checks back:
 - She finds the story Bob shared with her!



The reason that Alice doesn't see the post immediately is because Facebook applies an eventual consistent model. Facebook has more than 1 billion active users, so it's non-trivial to efficiently and reliably store the huge amount of data generated at any given time. An eventual consistent model offers the option to reduce the load and improve availability.

Eventual Consistency: Facebook

- Reason: it is possible because Facebook uses an **eventual consistent model**
- Why Facebook chooses eventual consistent model over the strong consistent one?
 - Facebook has more than 1 billion active users
 - It is non-trivial to efficiently and reliably store the huge amount of data generated at any given time
 - Eventual consistent model offers the option to **reduce the load and improve availability**

Eventual Consistency: Dropbox

- Dropbox enabled immediate consistency via synchronization in many cases.
- However, what happens in case of a network partition?



Eventual Consistency: Dropbox

- Let's do a simple thought experiment here:
 - Open a file in your Dropbox
 - Disable your network connection (e.g., Wi-Fi, 4G)
 - Try to edit the file in the Dropbox: can you do that?
 - Re-enable your network connection: what happens to your Dropbox folder?

Also Dropbox embraces eventual consistency. Immediate consistency is impossible in case of a network partition. But even when online, it is not desirable to have immediate (strong) consistency: users will feel bad if their Word documents freeze each time they save it, simply due to the large latency to update all devices across WAN.

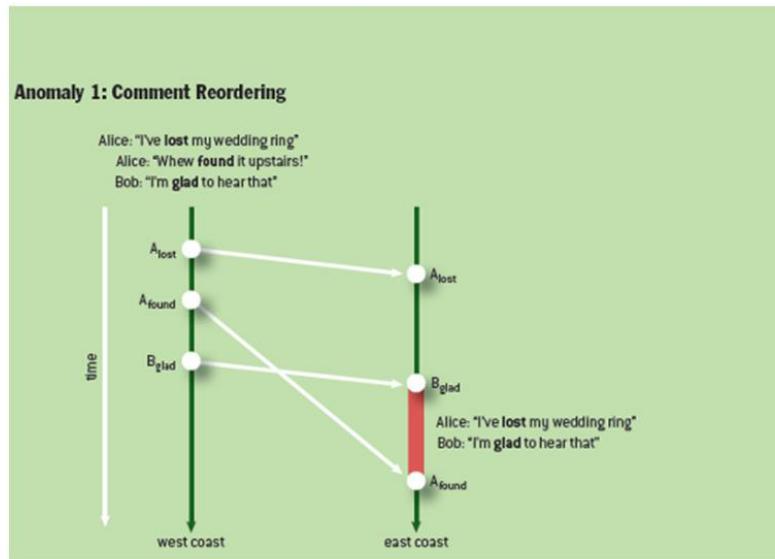
Dropbox is oriented to **personal syncing** and not on collaboration. So having eventual consistency is not a real limitation.

Eventual Consistency: ATM

- In design of automated teller machine (ATM):
 - Strong consistency appear to be a nature choice
 - However, in practice, **A beats C**
 - Higher availability means **higher revenue**
 - ATM will allow you to withdraw money *even if the machine is partitioned from the network*
 - However, it puts **a limit** on the amount of withdraw (e.g., \$200)
 - The bank might also charge you a fee when a overdraft happens



Eventual consistency goes without strict ordering



Eventual Consistency Variations

In practice, clients of the distributed data storage system get better guarantees than *pure* EC

- Read-your-write consistency
 - A process always accesses the data item after it's update operation and never sees an older value
- Session consistency
 - As long as session exists, system guarantees read-your-write consistency
 - Guarantees do not overlap sessions

Read Your Writes Consistency (RYWC) guarantees that a client that has written a version n will thereafter always be able to read a version that is at least as new as n . This helps, for example, to avoid user irritation when a person checks his bank account statement after he has wired some money to another person, but does not see his account credited. If there is no RYWC, he might think that the transfer was unsuccessful and wire the same amount of money again. Generally, RYWC avoids situations where a user or application issues the same request several times because it gets the impression that the request failed the first time. For idempotent operations reissuing requests causes only additional load on the system, while reissuing other requests create severe inconsistencies.

Eventual Consistency Variations

- Monotonic read consistency
 - If a process has seen a particular value of data item, any subsequent access by that process will never return any previous values
- Monotonic write consistency
 - The system guarantees to serialize the writes by the *same* process
- In practice
 - A number of these properties can be combined
 - Monotonic reads and read-your-writes are most desirable

Monotonic Read Consistency guarantees that a client that has read a version n will thereafter always read versions $\geq n$. This is helpful as from an application perspective data visibility might not be instantaneous but versions come at least in chronological order: the system never “goes backward” in time.

Monotonic Write consistency guarantees that two updates by the same client will be serialized in the order that they arrive at the storage system. This is useful to avoid seemingly lost updates when an application first writes and then updates a datum (data entry) but the update is executed before the initial write and is, thus, overwritten.

Dynamic Tradeoff between C and A



Many applications require neither strong consistency nor continual availability.

Applications specify level of continuous consistency: e.g. airline reservation system

- Many free seats – sacrifice consistency
- A few places left – sacrifice availability

When designing your system, CAP only describes that you have to balance (and not choose) between consistency and availability. The optimal trade-off depends on the application at hand and even inside a single application the desired trade-off might vary during execution.

Consider the example of an airline reservation. When most of the seats of a specific flight are available, it is ok to rely on somewhat out-of-date data. Availability is more critical, since being unavailability might mean that customers shift to another airline. However, when the plane is close to be filled, the reservation system needs more accurate data to ensure the plane is not overbooked. In this case consistency is more critical, because the fees that you have to pay to passengers with denied boarding is too high to compensate for the risk of losing a few customers when being unavailable.

Heterogeneity: Segmenting C and A

- No single uniform requirement for entire system
 - Some aspects require strong consistency
 - Others require high availability
- Segment your system into different components
 - Each provides different types of guarantees
- Overall guarantees neither consistency nor availability
 - But each part of the service gets exactly what it needs
- Can be partitioned along different dimensions

The trade-off between C and A should (and most likely will) not be the same for all the components in your system. When designing your system, you should segment it into different components that each receive exactly those guarantees they need.

Discussion

- In an e-commercial system (e.g., Amazon, e-Bay, etc), what are the trade-offs between consistency and availability you can think of? What is your strategy?
- Hint -> Things you might want to consider:
 - Different types of data (e.g., shopping cart, billing, product, etc.)
 - Different types of operations (e.g., query, purchase, etc.)
 - Different types of services (e.g., distributed lock, DNS, etc.)
 - Different groups of users (e.g., users in different geographic areas, etc.)

Examples of partitioning C/A

- Data Partitioning
- Operational Partitioning
- Functional Partitioning
- User Partitioning
- Hierarchical Partitioning

Partitioning Examples

Data Partitioning

- Different data may require different consistency and availability
- Example:
 - Shopping cart: high availability, responsive, can sometimes suffer anomalies
 - Product information need to be available, slight variation in inventory is sufferable
 - Checkout, billing, shipping records must be consistent

Partitioning Examples

Operational Partitioning

- Each operation may require different balance between consistency and availability
- Example:
 - Reads: high availability; e.g., “query”
 - Writes: high consistency, lock when writing; e.g., “purchase”

Partitioning Examples

Functional Partitioning

- System consists of sub-services
- Different sub-services provide different balances
- Example: A comprehensive distributed system
 - Distributed lock service (e.g., Chubby) :
 - Strong consistency
 - DNS service:
 - High availability

Partitioning Examples

User Partitioning

- Try to keep related data close together to assure better performance
- Example: Craigslist
 - Might want to divide its service into several data centers, e.g., east coast and west coast
 - Users get high performance (e.g., high availability and good consistency) if they query servers closest to them
 - Poorer performance if a New York user query Craigslist in San Francisco

Partitioning Examples

Hierarchical Partitioning

- Large global service with local “extensions”
- Different location in hierarchy may use different consistency
- Example:
 - Local servers (better connected) guarantee more consistency and availability
 - Global servers has more partition and relax one of the requirement

What if there are no partitions?

The **occurrence** of failure causes CAP tradeoffs

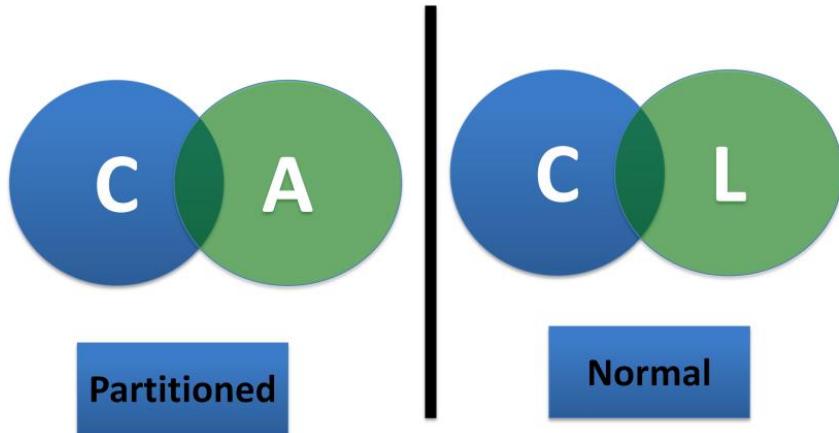
The **possibility** of failure results in
Consistency/Latency tradeoffs

- Availability ~ Latency
 - intuition: unavailable system provides extreme latency
- High availability → need to replicate → consistency problem

Availability and latency are arguably the same thing: an unavailable system provides extreme latency. Latency exists even without network partitions. But on the other hand, when a system runs long enough, at least one component will fail. Highly available systems will thus need to replicate data, which in turns causes challenges to guarantee consistency.

CAP → PACELC

If there is a **partition (P)**, how does the system trades off **availability and consistency (A and C); else (E)**, when the system is running normally in the absence of partitions, how does the system trades off **latency (L) and consistency (C)?**



The acronym PACELC makes clear that even in the **absence** of partitions there's a tradeoff between consistency and latency. In the general case, though not all cases, consistency requires a level of coordination which prevents systems from being always available, and increases latency when no partition is present. The matter of latency, which is of great practical importance in real-world systems, isn't captured at all in CAP.

A more complete description of the space of potential tradeoffs for distributed system is thus given by the PACELC acronym. PACELC stands for the following sentence: “If there is a **partition (P)**, how does the system trade off **availability and consistency (A and C); else (E)**, when the system is running normally in the absence of partitions, how does the system trade off **latency (L) and consistency (C)?**”

This sentence (and the acronym) are written in a seminal article of Daniel Abadi on consistency in modern distributed database system design:

Abadi, Daniel J. "Consistency tradeoffs in modern distributed database system design." Computer-IEEE Computer Magazine 45.2 (2012): 37.

Examples

- **PC/EC Systems:** Refuse to give up consistency and pay the cost of availability and latency
 - BigTable, Hbase, VoltDB/H-Store
- **PA/EL Systems:** Give up both Cs for availability and lower latency
 - Dynamo, Cassandra, Riak
- **PA/EC Systems:** Give up consistency when a partition happens and keep consistency in normal operations
 - MongoDB
- **PC/EL System:** Keep consistency if a partition occurs but gives up consistency for latency in normal operations
 - Yahoo! PNUTS

The first two categories of PACELC are very clear. PC/EC is the most consistent class of systems, which never give up consistency. PA/EL systems don't try hard to be consistent, and rather take the opportunity to reduce latency and gain availability by reducing coordination.

The trickier ground starts with PA/EC. These types of systems give up consistency when there is a partition, and are consistent when there isn't. That's more subtle than it looks. When is there a partition? How long does the network need to be down before there is a partition? Is a single dropped connection or lost packet a partition? That may seem like nit picking, but there's an important line to be drawn between *partition* and *not partition*. PACELC doesn't help there.

If PA/EC is tricky, PC/EL is madness. What does it mean to be more consistent during a partition? Daniel Abadi (who coined PACELC) says that is the wrong question: Yahoo! PNUTS is a PC/EL system. In normal operation, it gives up consistency for latency; however, if a partition occurs, it trades availability for consistency. This is admittedly somewhat confusing: according to PACELC, PNUTS appears to get more consistent upon a network partition. However, PC/EL should not be interpreted in this way. PC does not indicate that the system is fully consistent; rather it indicates that the system does not reduce consistency beyond the baseline consistency level when a network partition occurs—instead, it reduces availability.

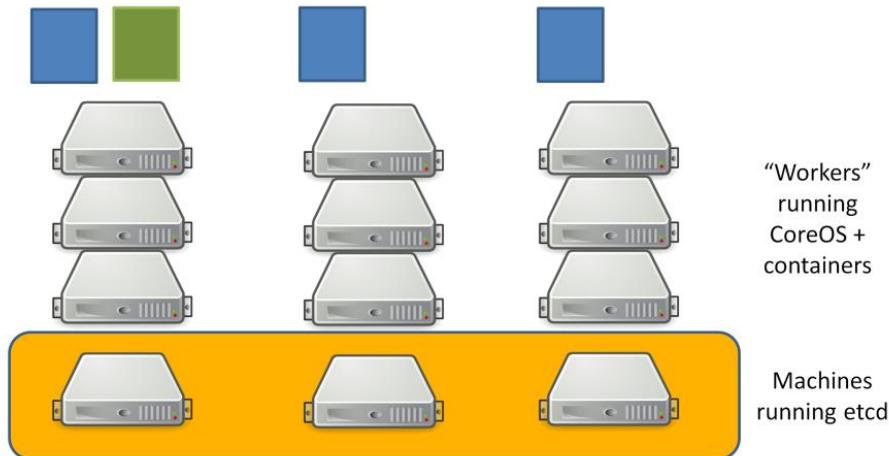
Bibliography: distributed computing

- D. Wang, Cloud Computing, CSE 40822,
<http://www3.nd.edu/~dthain/>
- M. Brooker, CAP and PACELC: Thinking More Clearly About Consistency,
<http://brooker.co.za/blog/2014/07/16/pacelc.html>
- I. Tsalouchidou, The Cap theorem in depth,
http://www.slideshare.net/ioanna_tsalouchidou/cap-in-depth
- D. Bermbach, Benchmarking, Consistency, Distributed Database Management Systems
- D. J. Abadi, Consistency Tradeoffs in Modern Distributed Database System Design
- D. Power, What is ACID and BASE in database theory?
<http://dssresources.com/faq/index.php?action=artikel&id=281>

Outline

- Parallel computing
- Distributed computing
 - what is distributed computing
 - CAP theorem and its extensions
 - RAFT consensus algorithm
 - definition
 - challenges
 - practical example: etcd

Consensus: practical example



Distributed consensus key-value store (etcd) used for

- Cluster management: current state of all machines and containers
- Application details: database connection details, feature flags...

Consensus algorithms allow a collection of machines to work as a coherent group that can survive the failures of some of its members. Because of this, they play a key role in building reliable large-scale software systems. A typical example where consensus algorithms are used is in a server cluster running containers. The set-up consists of worker nodes, complemented with 3-5 machines running central services (including the distributed key-value store). Each of the workers will use the distributed key-value store on the central machines via local proxies.

An example central service is the fleet manager. This fleet manager handles the scheduling of containers across the cluster machines and keeps them running even if the original host they were running on is terminated. The fleet manager needs a shared view of the current state of all the machines and containers running in the cluster. This functionality is provided by e.g. etcd, a distributed, consistent key-value store used for storing shared configuration and information in a cluster.

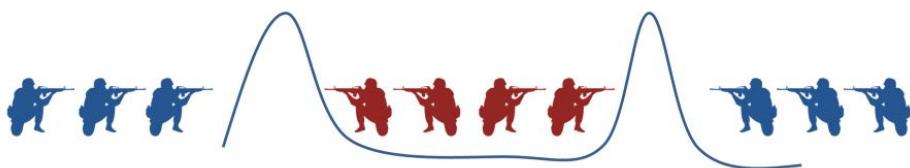
Also applications can read and write data into the key-value store. A simple use-case is to store database connection details or feature flags in etcd as key value pairs. These values can be watched, allowing your app to reconfigure itself when they change. Advanced users take advantage of the consistency guarantees to implement database master elections or do distributed locking across a cluster of workers.

Requirements for consensus

- **Agreement**
 - All correct processes must agree on the same value
- **Termination**
 - All processes must eventually decide on an output value
- **Validity**
 - If all correct processes propose the same value v , then all correct processes decide v
- **Integrity**
 - If a correct process decides v , then v must have been processed by some correct process

The idea behind integrity is that we want to exclude trivial solutions that just decide ‘No’ whatever the initial set of values is. Such an algorithm would satisfy termination and agreement, but would be completely vacuous, and no use to use at all.

Consensus



- generals of left and right blue army must decide on value of attack (= 0 or 1)
 - blue army only wins if both sides attack at the same time
 - to reach consensus, they must send messages
 - but these messages may get lost during transmission
- impossible to reach consensus!
 - my message may be intercepted/lost
 - is the other army still there?
- main reason: asynchronous character of the messaging

Consensus algorithms allow a collection of machines to work as a coherent group that can survive the failures of some of its members. Because of this, they play a key role in building reliable large-scale software systems. Some of the challenges in reaching consensus between distributed entities are exemplified in the well known “Two army problem”. Suppose we have two allied (blue) armies that have encircled the red army. The red army can only be beaten when both blue armies attack simultaneously. So, to reach consensus, the general of the left blue army will have to communicate with the general of the right blue army.

However, using messaging, it is impossible to reach consensus on the decision to attack or not. The major problem is that the armies use asynchronous messaging: there is no upper boundary on the time in which one army may expect an answer from the other army. In other words: for the left army it is impossible to know why the other army has not responded: its own message could be lost, the general of the other army might still be thinking about his strategy, or the reply of the other army can be lost during transmission.

Intuitive proof of the above: assume protocol P is the shortest protocol (i.e. with a minimum number of messages) that solves the attack decision problem. Suppose now that the last message of protocol P does not reach its destination. Since protocol P is correct, consensus must be reached in any case. This means, the last message was useless, and then P could not be the shortest!

FLP theorem: “impossibility result”

No consensus can be guaranteed in an asynchronous communication system in the presence of any failures.

Fischer-Lynch-Patterson (1985)

- The real-world is asynchronous
 - Variations in response time of websites
 - Link failures (\sim infinite response time)
- But... consensus is possible in (some) synchronous settings!
- Real-world distributed consensus systems often assume *partial synchrony*
 - “Timing-based distributed algorithms”
 - Individual processes have some information about time, e.g.
 - Clocks are synchronized within some bound
 - Approximate bounds on message-deliver time
 - Use of timeouts

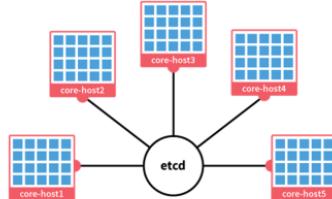
In 1985, Fischer, Lynch and Patterson formulated their “Impossibility Result” theorem which states that no consensus can be guaranteed in an asynchronous communication system in the presence of any failures. In an asynchronous timing assumption, there is no time bound on how long it takes for a message to be delivered.

Although the theorem has been proven mathematically, we provide here only an intuitive explanation. The asynchronous assumption makes it impossible to differentiate between failed and slow processes. A “failed” process may just be slow, or can rise from the dead at exactly the wrong time. Therefore, *termination* (liveness) cannot be guaranteed. On the other hand, an slow process may decide differently than other processes, thus violating the agreement property of consensus.

Because the real-world operates asynchronous, the FLP theorem is very important. It states that we must introduce at least a partial level of synchrony to be able to reach consensus. With partial synchrony, individual processes have some information about time, e.g. using clock synchronization within some bound, by setting approximate bounds on message-deliver time, by the use of timeouts...

RAFT algorithm

- New protocol (2014)
- Designed for *understandability*
- Used in *etcd*
- Before, PAXOS was the default standard
 - exceptionally difficult to understand



"The dirty little secret of the NSDI community is that at most five people really, truly understand every part of PAXOS ☺"

anonymous NSDI reviewer

- very difficult to implement

"There are significant gaps between the description of the PAXOS algorithm and the needs of a real-world system... The final system will be based on an unproven protocol"

authors of Chubby (distributed lock service from Google)

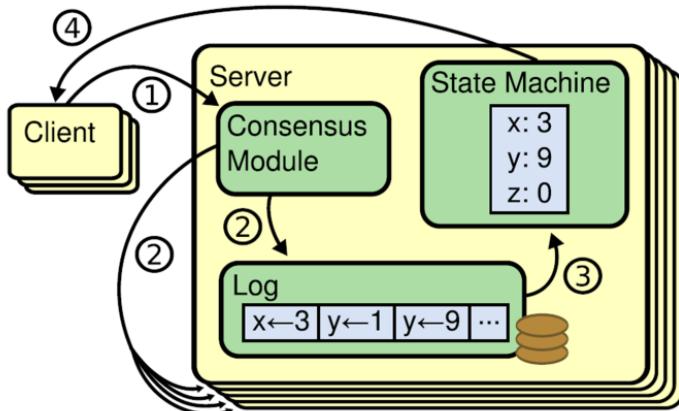
The Paxos consensus algorithm dominated the discussion of consensus algorithms over the last decade. Paxos first defines a protocol capable of reaching agreement on a single decision, such as a single replicated log entry. We refer to this subset as *single-decree Paxos*. Paxos then combines multiple instances of this protocol to facilitate a series of decisions such as a log (*multi-Paxos*). Unfortunately, PAXOS has two significant drawbacks.

The first drawback is that Paxos is exceptionally difficult to understand. The full (original) explanation is notoriously opaque; few people succeed in understanding it, and only with great effort. Paxos' opaqueness derives from its choice of the single-decree subset as its foundation. Single-decree Paxos is dense and subtle: it is divided into two stages that do not have simple intuitive explanations and cannot be understood independently. Because of this, it is difficult to develop intuitions about why the single degree protocol works. The composition rules for multi-Paxos add significant additional complexity and subtlety.

The second problem with Paxos is that it does not provide a good foundation for building practical implementations. One reason is that there is no widely agreed upon algorithm for multi-Paxos. Furthermore, the Paxos architecture is a poor one for building practical systems; this is another consequence of the single-degree decomposition. For example, there is little benefit to choosing a collection of log entries independently and then melding them into a sequential log; this just adds

complexity. It is simpler and more efficient to design a system around a log, where new entries are appended sequentially in a constrained order. Another problem is that Paxos uses a symmetric peer-to-peer approach at its core. This makes sense in a simplified world where only one decision will be made, but few practical systems use this approach. If a series of decisions must be made, it is simpler and faster to first elect a leader, then have the leader coordinate the decisions. As a result, practical systems bear little resemblance to Paxos. Each implementation begins with Paxos, discovers the difficulties in implementing it, and then develops a significantly different architecture.

State machines & the distributed log



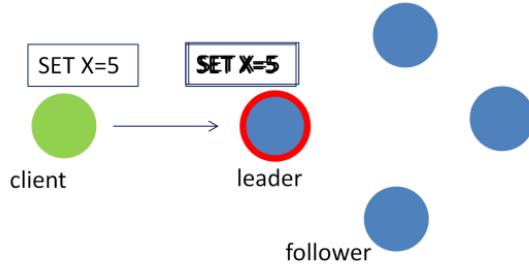
- Consensus algorithms ensures that all logs contain the same commands in the same order
- State machines remain consistent if commands have deterministic results

Consensus algorithms are typically needed in the context of *replicated state machines*. State machines on a collection of servers compute identical copies of the same state and can continue operating even if some of the servers are down. Replicated state machines are used to solve a variety of fault tolerance problems in distributed systems. For example, large-scale systems that have a single cluster leader, such as HDFS, typically use a separate replicated state machine to manage leader election and store configuration information that must survive crashes. Examples of replicated state machines include Chubby, ZooKeeper and Etcd.

Replicated state machines are typically implemented using a replicated log. Each server stores a log containing a series of commands, which its state machine executes in order. Each log contains the same commands in the same order, so each state machine processes the same sequence. Since the state machines are deterministic, each machine computes the same state and the same sequence of outputs.

Keeping the replicated log consistent is the job of the consensus algorithm. The consensus module receives commands from clients and adds them to its log. It communicates with the consensus modules on other servers to ensure that every log eventually contains the same requests in the same order, even if some servers fail. Once commands are properly replicated, each server's state machine processes them in log order, and the outputs are returned to clients. As a result, the servers appear to form a single, highly reliable state machine.

Leaders and followers

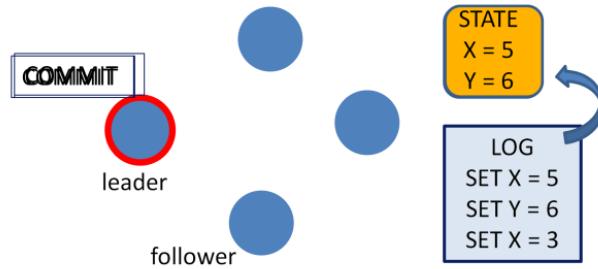


- Cluster nodes elect a single **leader**, others become followers
 - Sole server habilitated to accept commands from clients
 - Will enter them in its log and forward them to the followers
 - Tells followers when it's safe to apply log entries to their state machine

A RAFT cluster consists of several servers (typically five). The cluster can tolerate the failure of any two servers.

Raft implements consensus by first electing one distinguished *leader*, then giving the leader complete responsibility for managing the replicated log. The leader accepts log entries from clients, replicates them on other servers, and tells servers when it is safe to apply log entries to their state machines. Having a leader simplifies the management of the replicated log. For example, the leader can decide where to place new entries in the log without consulting other servers, and data flows in a simple fashion from the leader to other servers. A leader can fail or become disconnected from the other servers, in which case a new leader is elected.

Committing state

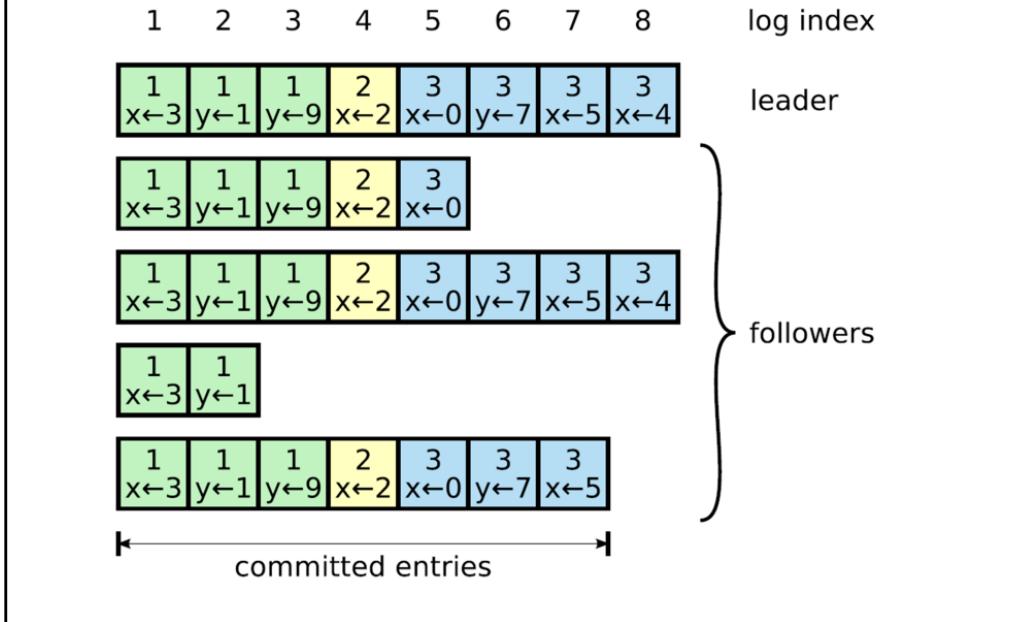


- Leader determines when a log entry is *committed*
 - when it receives ACK from a majority of the nodes
- Followers only apply log update to their state machine after notification of commitment by leader

The leader decides when it is safe to apply a log entry to state machines, such an entry is called *committed*. Raft guarantees that committed entries are durable and will eventually be executed by all of the available state machines. A log entry is committed once the leader that created the entry has replicated it on a majority of the servers. This also commits all preceding entries in the leader's log, including entries created by previous leaders.

The leader keeps track of the highest index it knows to be committed, and it includes that index in future AppendEntries RPCs (including heartbeats) so that the other servers eventually find out. Once a follower learns that a log entry is committed, it applies the entry to its local state machine (in log order).

Log replication



Each follower has a subset of the committed entries of the leader.

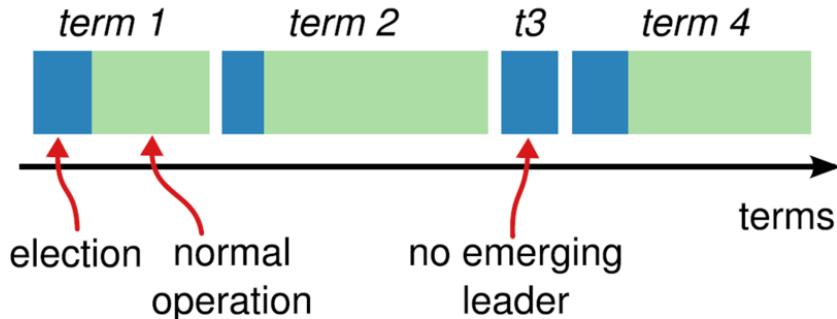
3 subproblems

- Leader election
- Log replication
- Safety

Given the leader approach, Raft decomposes the consensus problem into three relatively independent subproblems:

- Leader election: a new leader must be chosen when an existing leader fails
- Log replication: the leader must accept log entries from clients and replicate them across the cluster, forcing the other logs to agree with its own
- Safety: if a server has applied a log entry at given index to its state machine, no other server will ever apply a different log entry for the same index

Terms

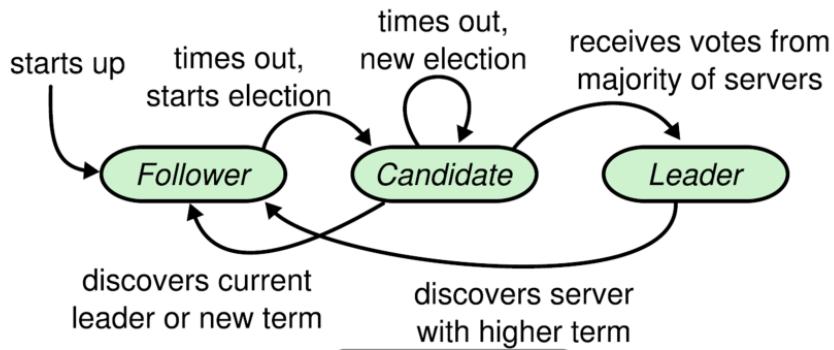


- Terms act as logical clock
 - Each server stores its own *current term*
- Each term begins with a leader election
- A single leader manages the cluster for the entire term
 - Leaders typically operate until they fail
- Different servers may observe the transitions between terms at different times

Raft divides time into *terms* of arbitrary length. Terms are numbered with consecutive integers. Each term begins with an *election*, in which one or more candidates attempt to become leader. If a candidate wins the election, then it serves as leader for the rest of the term. In some situations an election will result in a split vote. In this case the term will end with no leader; a new term (with a new election) will begin shortly. Raft ensures that there is at most one leader in a given term.

Different servers may observe the transitions between terms at different times, and in some situations a server may not observe an election or even entire terms. Terms act as a logical clock in Raft, and they allow servers to detect obsolete information such as stale leaders. Each server stores a *current term* number, which increases monotonically over time. Current terms are exchanged whenever servers communicate; if one server's current term is smaller than the other's, then it updates its current term to the larger value. If a candidate or leader discovers that its term is out of date, it immediately reverts to follower state. If a server receives a request with a stale term number, it rejects the request.

Server state transitions



At any given time each server in a RAFT cluster is in one of three states *leader*, *follower*, or *candidate*. In normal operation there is exactly one leader and all of the other servers are followers. Followers are passive: they issue no requests on their own but simply respond to requests from leaders and candidates. The leader handles all client requests (if a client contacts a follower, the follower redirects it to the leader). If a follower receives no communication, it becomes a candidate and initiates an election. A candidate that receives votes from a majority of the full cluster becomes the new leader. Leaders typically operate until they fail.

Leader election: candidates

<https://ramcloud.stanford.edu/~ongaro/raftscope/>

- After election timeout, follower transitions to candidate
 - increases its term
 - sends “RequestVote” to other cluster members
- Election outcome
 - majority (incl itself) votes for this candidate
 - servers vote on first-come-first-served and vote only once per term
 - candidate becomes leader for this term
 - candidate receives notification of other leader in at least the same term
 - candidate resigns
 - split vote
 - new election procedure
 - randomized time-outs make this very unlikely; but not impossible

Raft uses a heartbeat mechanism to trigger leader election. When servers start up, they begin as followers. A server remains in follower state as long as it receives valid RPCs from a leader or candidate. Leaders send periodic heartbeats (AppendEntriesRPCs that carry no log entries) to all followers in order to maintain their authority. If a follower receives no communication over a period of time called the *election timeout*, then it assumes there is no viable leader and begins an election to choose a new leader.

To begin an election, a follower increments its current term and transitions to candidate state. It then votes for itself and issues RequestVote RPCs in parallel to each of the other servers in the cluster. A candidate continues in this state until one of three things happens:

- (a) it wins the election,
- (b) another server establishes itself as leader, or
- (c) a period of time goes by with no winner.

A candidate wins an election if it receives votes from a majority of the servers in the full cluster for the same term. Each server will vote for at most one candidate in a given term, on a first-come-first-served basis. The majority rule ensures that at most one candidate can win the election for a particular term. Once a candidate wins an election, it becomes leader. It then sends heartbeat messages to all of the other servers to establish its authority and prevent new elections.

While waiting for votes, a candidate may receive an AppendEntries RPC from another server claiming to be leader. If the leader's term (included in its RPC) is at least as large as the candidate's current term, then the candidate recognizes the leader as legitimate and returns to follower state. If the term in the RPC is smaller than the candidate's current term, then the candidate rejects the RPC and continues in candidate state.

The third possible outcome is that a candidate neither wins nor loses the election: if many followers become candidates at the same time, votes could be split so that no candidate obtains a majority. When this happens, each candidate will time out and start a new election by incrementing its term and initiating another round of RequestVote RPCs. However, without extra measures split votes could repeat indefinitely.

Raft uses randomized election timeouts to ensure that split votes are rare and that they are resolved quickly. To prevent split votes in the first place, election timeouts are chosen randomly from a fixed interval (e.g., 150–300ms). This spreads out the servers so that in most cases only a single server will time out; it wins the election and sends heartbeats before any other servers time out. The same mechanism is used to handle split votes. Each candidate restarts its randomized election timeout at the start of an election, and it waits for that timeout to elapse before starting the next election; this reduces the likelihood of another split vote in the new election.

3 subproblems

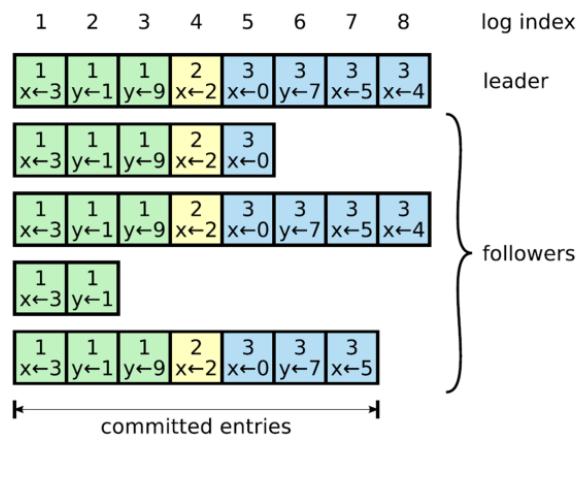
- Leader election
- Log replication
- Safety

Given the leader approach, Raft decomposes the consensus problem into three relatively independent subproblems:

- Leader election: a new leader must be chosen when an existing leader fails
- Log replication: the leader must accept log entries from clients and replicate them across the cluster, forcing the other logs to agree with its own
- Safety: if a server has applied a log entry at given index to its state machine, no other server will ever apply a different log entry for the same index

Log replication

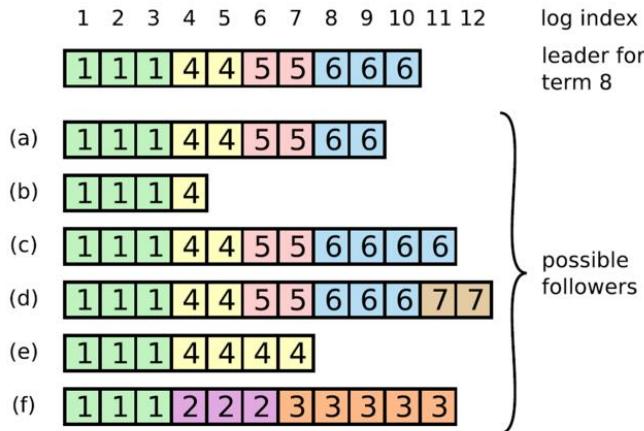
- Two entries in different logs having the same index and term store the same command
 - a leader creates at most one entry with a given log index in a given term
- Two entries in different logs having the same index and term are identical in all preceding entries
 - leader includes index/term of previous entry in its log to append request



Once a leader has been elected, it begins servicing client requests. Each client request contains a command to be executed by the replicated state machines. The leader appends the command to its log as a new entry, then issues AppendEntries RPCs in parallel to each of the other servers to replicate the entry. When the entry has been safely replicated, the leader applies the entry to its state machine and returns the result of that execution to the client. If followers crash or run slowly, or if network packets are lost, the leader retries AppendEntries RPCs indefinitely (even after it has responded to the client) until all followers eventually store all log entries.

Logs are organized as shown. Each log entry stores a state machine command along with the term number when the entry was received by the leader. The term numbers in log entries are used to detect inconsistencies between logs. Each log entry also has an integer index identifying its position in the log.

Bringing logs to consistency



- Conflicting entries in followers logs are overwritten with entries from the leader's log
- Finding latest log entry where leader/follower log agree is done by a consistency check after each AppendEntries RPC

During normal operation, the logs of the leader and followers stay consistent, so the AppendEntries consistency check never fails. However, leader crashes can leave the logs inconsistent (the old leader may not have fully replicated all of the entries in its log). These inconsistencies can compound over a series of leader and follower crashes. The above figure illustrates the ways in which followers' logs may differ from that of a new leader. A follower may be missing entries that are present on the leader, it may have extra entries that are not present on the leader, or both. Missing and extraneous entries in a log may span multiple terms.

When the leader at the top in the figure comes to power, it is possible that any of scenarios (a–f) could occur in follower logs. Each box represents one log entry; the number in the box is its term. A follower may be missing entries (a–b), may have extra uncommitted entries (c–d), or both (e–f). For example, scenario (f) could occur if that server was the leader for term 2, added several entries to its log, then crashed before committing any of them; it restarted quickly, became leader for term 3, and added a few more entries to its log; before any of the entries in either term 2 or term 3 were committed, the server crashed again and remained down for several terms.

In Raft, the leader handles inconsistencies by forcing the followers' logs to duplicate its own. This means that conflicting entries in follower logs will be overwritten with entries from the leader's log. To bring a follower's log into consistency with its own, the leader must find the latest log entry where the two logs agree, delete any entries

in the follower's log after that point, and send the follower all of the leader's entries after that point. All of these actions happen in response to the consistency check performed by AppendEntries RPCs. The leader maintains a *nextIndex* for each follower, which is the index of the next log entry the leader will send to that follower. When a leader first comes to power, it initializes all *nextIndex* values to the index just after the last one in its log (11 in the figure above). If a follower's log is inconsistent with the leader's, the AppendEntries consistency check will fail in the next AppendEntries RPC. After a rejection, the leader decrements *nextIndex* and retries the AppendEntries RPC. Eventually *nextIndex* will reach a point where the leader and follower logs match. When this happens, AppendEntries will succeed, which removes any conflicting entries in the follower's log and appends entries from the leader's log (if any). Once AppendEntries succeeds, the follower's log is consistent with the leader's, and it will remain that way for the rest of the term.

With this mechanism, a leader does not need to take any special actions to restore log consistency when it comes to power. It just begins normal operation, and the logs automatically converge in response to failures of the AppendEntries consistency check. A leader never overwrites or deletes entries in its own log.

3 subproblems

- Leader election
- Log replication
- Safety

The mechanisms described so far are not quite sufficient to ensure that each state machine executes exactly the same commands in the same order. For example, a follower might be unavailable while the leader commits several log entries, then it could be elected leader and overwrite these entries with new ones; as a result, different state machines might execute different command sequences.

Election restriction

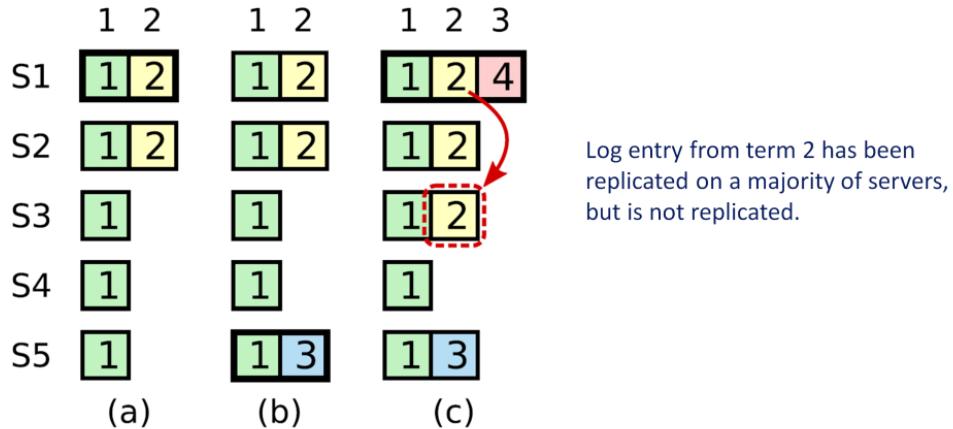
- Candidate cannot win election unless its log contains all committed entries
 - avoids transfer of entries between old and new leader
- Intuitive proof:
 - candidate must contact a majority of the cluster to be elected
 - every committed update is present in at least one of those servers
 - voter denies its vote if own log is more up-to-date

Raft uses a simple approach that guarantees that all the committed entries from previous terms are present on each new leader from the moment of its election, without the need to transfer those entries to the leader. This means that log entries only flow in one direction, from leaders to followers, and leaders never overwrite existing entries in their logs.

Raft uses the voting process to prevent a candidate from winning an election unless its log contains all committed entries. A candidate must contact a majority of the cluster in order to be elected, which means that every committed entry must be present in at least one of those servers. If the candidate's log is at least as up-to-date as any other log in that majority (where "up-to-date" is defined precisely below), then it will hold all the committed entries. The RequestVote RPC implements this restriction: the RPC includes information about the candidate's log, and the voter denies its vote if its own log is more up-to-date than that of the candidate. Raft determines which of two logs is more up-to-date by comparing the index and term of the last entries in the logs. If the logs have last entries with different terms, then the log with the later term is more up-to-date. If the logs end with the same term, then whichever log is longer is more up-to-date.

Committing entries from previous terms

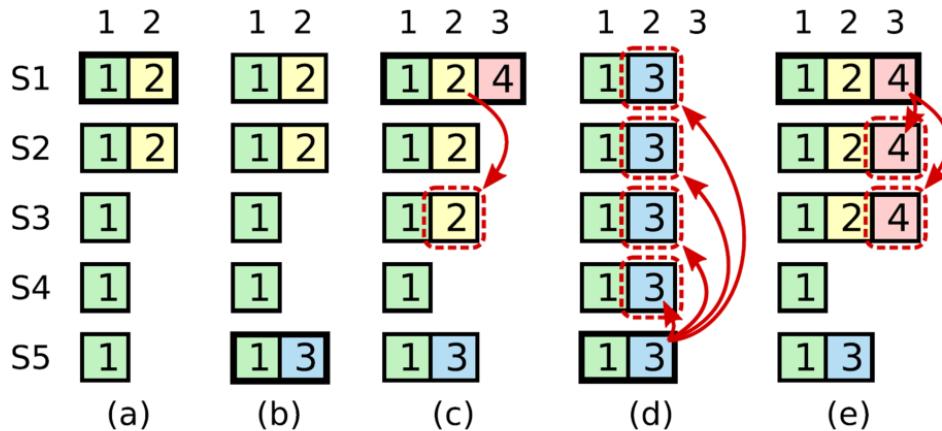
- A leader knows that an entry from its **current term** is committed once stored on a majority of servers
- A leader cannot determine commitment of entries from **previous terms** based on this majority rule



A leader knows that an entry from its current term is committed once that entry is stored on a majority of the servers. If a leader crashes before committing an entry, future leaders will attempt to finish replicating the entry. However, a leader cannot immediately conclude that an entry from a previous term is committed once it is stored on a majority of servers.

To illustrate this problem, we use the above time sequence showing why a leader cannot determine commitment using log entries from older terms. In (a) S1 is leader and partially replicates the log entry at index 2. In (b) S1 crashes; S5 is elected leader for term 3 with votes from S3, S4, and itself, and accepts a different entry at log index 2. In (c) S5 crashes; S1 restarts, is elected leader, and continues replication. At this point, the log entry from term 2 has been replicated on a majority of the servers, but it is not committed. It can be overwritten by a future leader (as we show on the next slide)

Committing entries from previous terms



If S1 crashes, as in (d), S5 could be elected leader (with votes from S2, S3 and S4) and overwrite the entry with its own entry from term 3.

However, if S1 replicates an entry from its **current** term on a majority of the servers before crashing, as in (e), then this entry is committed (S5 cannot win an election). At this point all preceding entries in the log are committed as well.

Thus, to eliminate commitment problems with previous terms, Raft never commits log entries from previous terms by counting replicas. Only log entries from the leader's current term are committed by counting replicas; once an entry from the current term has been committed in this way, then all prior entries are committed indirectly (because servers will first synchronize their complete log with the leader, as described before).

Bibliography: distributed systems

- D. Thain, Foundations of Distributed Systems
- The paper trail, A brief tour of FLP impossibility,
<http://the-paper-trail.org/blog/a-brief-tour-of-flp-impossibility/>
- A. Brown, Distributed agreement,
<http://www.cs.toronto.edu/~demke/469F.07/Lectures/Lecture18.ppt>
- R. Wattenhofer, Principles of Distributed Computing (lecture collection),
http://disco.ethz.ch/lectures/podc_allstars/

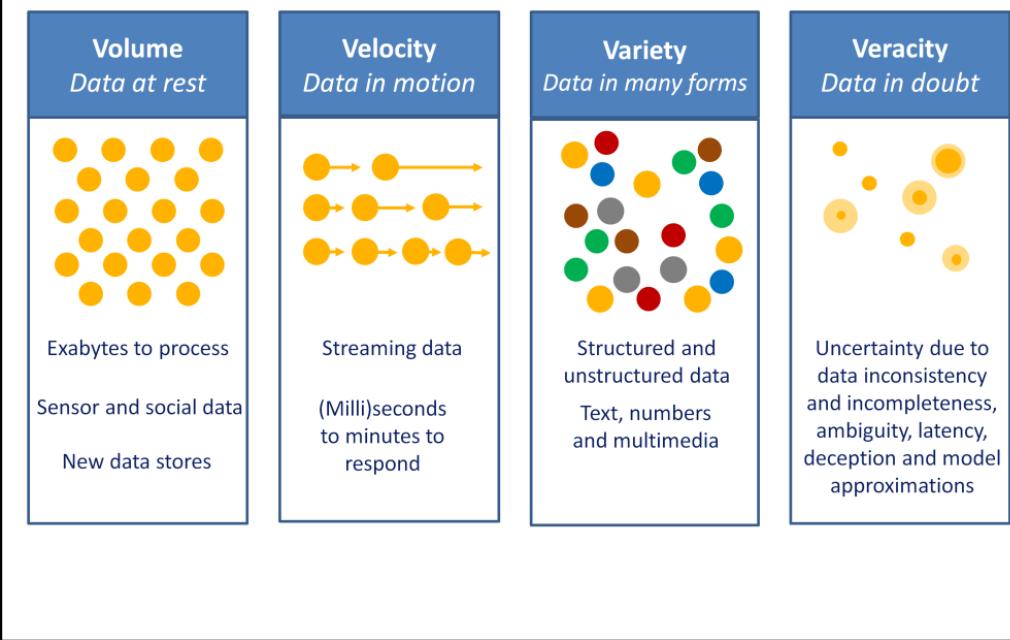


We are living in the era of big data, where exponential growth of phenomena such as web, social networking, smartphones, and so on are producing petabytes of data on a daily basis. Gaining insights from analyzing these very large amounts of data has become a *must-have* competitive advantage for many industries. However, the size and the possibly unstructured nature of these data sources make it impossible to use traditional solutions such as relational databases to store and analyze these datasets.

Storage, processing, and analyzing petabytes of data in a meaningful and timely manner require many compute nodes with thousands of disks and thousands of processors together with the ability to efficiently communicate massive amounts of data among them. Such a scale makes failures such as disk failures, compute node failures, network failures, and so on a common occurrence making fault tolerance a very important aspect of such systems. Other common challenges that arise include the significant cost of resources, handling communication latencies, handling heterogeneous compute resources, synchronization across nodes, and load balancing. As you can infer, developing and maintaining distributed parallel applications to process massive amounts of data while handling all these issues is not an easy task.



Big Data definition



The term “big data” remains difficult to understand because it can mean so many different things to different people. Your understanding will be different if you look at big data through a technology lens, versus a business lens or industry lens.

Essentially, big data refers to two major phenomena:

- The breathtaking speed at which we are now generating new data
- Our improving ability to store, process and analyze that data

To describe the phenomenon that is big data, people have been using the four Vs: Volume, Velocity, Variety and Veracity.

Volume refers to the vast amount of data generated every second. Just think of all the emails, Twitter messages, photos, video clips and sensor data that we produce and share every second. We are not talking terabytes, but zettabytes or brontobytes of data. On Facebook alone we send 10 billion messages per day, click the like button 4.5 billion times and upload 350 million new pictures each and every day. If we take all the data generated in the world between the beginning of time and the year 2000, it is the same amount we now generate every minute! This increasingly makes data sets too large to store and analyze using traditional database technology. With big data technology we can now store and use these data sets with the help of distributed systems, where parts of the data is stored in different locations, connected by networks and brought together by software.

Velocity refers to the speed at which new data is generated and the speed at which

data moves around. Just think of social media messages going viral in minutes, the speed at which credit card transactions are checked for fraudulent activities or the milliseconds it takes trading systems to analyze social media networks to pick up signals that trigger decisions to buy or sell shares. Big data technology now allows us to analyze the data while it is being generated without ever putting it into databases.

Variety refers to the different types of data we can now use. In the past we focused on structured data that neatly fits into tables or relational databases such as financial data (for example, sales by product or region). In fact, 80 percent of the world's data is now unstructured and therefore can't easily be put into tables or relational databases—think of photos, video sequences or social media updates. With big data technology we can now harness differed types of data including messages, social media conversations, photos, sensor data, video or voice recordings and bring them together with more traditional, structured data.

Veracity refers to the messiness or trustworthiness of the data. With many forms of big data, quality and accuracy are less controllable, for example Twitter posts with hashtags, abbreviations, typos and colloquial speech. Big data and analytics technology now allows us to work with these types of data. The volumes often make up for the lack of quality or accuracy.

What is a data system?

Answers questions based on information that was acquired in the past up to the present

Query = function (*all* data)



What is this person's name?
How many friends does this person have?

Several pieces combined to produce information



What is my current balance?
What recent transactions have occurred on my account ?

Not all information is equal

Data is the raw source from which information is derived

At the most fundamental level, a data system answers questions based on information that was acquired in the past up to the present. So a social network profile answers questions like “What is this person’s name?” and “How many friends does this person have?”. A bank account web page answers questions like “What is my current balance?” and “What transactions have occurred on my account recently”?

Data systems don’t just memorize and regurgitate information. They combine bits and pieces together to produce their answers. A bank account balance, for example, is based on combining the information about all transactions on the account.

Another crucial observation is that not all bits of information are equal. Some information is derived from other pieces of information. A bank account balance is derived from a transaction history. A friend count is derived from a friend list, and a friend list is derived from all the times a user added and removed friends from their profile.

When you keep tracing back where information is derived from, you eventually end up at information that is not derived from anything. This is the rawest information you have: information you hold to be true simply because it exists. We call this information *data*.

STORING BIG DATA

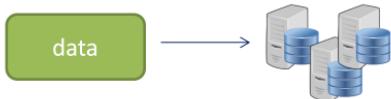
Outline

- Types of distributed storage solutions
 - Distributed file system
 - case study: Hadoop DFS2
 - Distributed data stores
 - key-value; columnar; document; graph
- Reasons for distributing
 - sharding
 - replication

Scaling data storage

Distributing data across nodes is a necessity when:

- **volume**: dataset too large for a single node
- **velocity**: single node cannot cope with required read/write rates



Common requirements for each storage solution:

- scalability: preferably horizontal scaling using commodity hardware
- fault tolerance: data is not lost when one or more node(s) fail
[failure is the norm, rather than the exception]

Distributed file systems



write-once-read-many

Distributed data stores



mongoDB

amazon web services S3



cassandra

random read/write

Storage solutions for big data have to be innately distributed. First, as the volume of big data grows, in the near or far feature you will reach the capacity limits of a single node. Second, in some cases a fast processing of big data is required. This means that there can be important constraints on the tolerable latency for a single read or write operation, even under heavy load with many concurrent users.

There is no one-size-fits-all solution for big data. At a high level, we can however distinguish two categories: file systems and data stores.

Distributed file systems, like the Hadoop Distributed File System (HDFS), are a perfect fit when data is written only once (e.g. raw sensor readings, logs, transaction data): the only write operation is to add *new* data (and not to modify existing data). Conversely, you want to perform regular calculations on this data, so reading operations are frequent.

Distributed data stores often prefer availability over consistency (cfr. CAP) and have advanced indexing mechanisms for rapid execution of random read/write messages. These kind of storage solutions are geared to update (overwrite) existing data. Distributed data stores are often grouped under the umbrella “NoSQL”, but also relational databases can be distributed. We will discuss this later.

Both distributed file systems and distributed data stores have their unique strengths and compromises and even within a single category each available solution provides

different performance guarantees. It is up to the developer to pick the right tool – and to correctly engineer it. In the next slides, we will cover the major principles that should arm you to choose the best solution for the problem at hand.

Hadoop Distributed File System



- Suitable for applications with large data sets
 - Highly fault-tolerant
 - High throughput
 - Runs on commodity hardware clusters
- “Moving computation is cheaper than moving data”
- Data characteristics
 - write-once-read-many
 - streaming data access (append-only write)
 - large files (gigabytes to terabytes)
 - batch processing rather than interactive

The Hadoop Distributed File System (HDFS) is a distributed and scalable file system that manages how data is stored across a cluster of commodity storage nodes. It has built-in capacities for fault-tolerance and favors high throughput reading operations, even on commodity hardware.

A computation requested by an application is much more efficient if it is executed near the data it operates on. This is especially true when the size of the data set is huge. This minimizes network congestion and increases the overall throughput of the system. The assumption is that it is often better to migrate the computation closer to where the data is located rather than moving the data to where the application is running. HDFS provides interfaces for applications to move themselves closer to where the data is located.

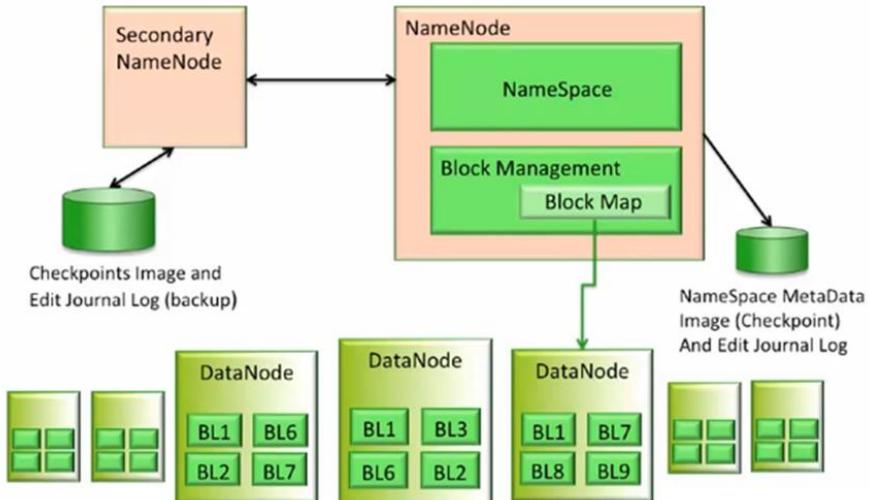
Applications built on top of HDFS need a write-once-read-many access model for files. A file once created, written and closed need not be changed. This assumption simplifies data coherency issues and enables high throughput data access. A Map/Reduce application or a web crawler application fits perfectly with this model.

HDFS provides streaming data access: it is not possible to start reading from a random point in the file. Write operations are limited to appending to the end of an existing file, but this is very complex and should be avoided. If you want to add new data, then simply add an additional file to the filesystem.

The data should be organized in files of sufficient size (GBs per file), otherwise the efficiency of the distributed filesystem will degrade. As we will see later in more detail, application code is brought to the data nodes. Having many small files would mean that the application code has to be deployed on many data nodes, which introduces bookkeeping overhead.

HDFS architecture

- Files are broken into blocks of equal size
- Blocks are **spread** over multiple nodes for scalability and to enable parallel processing
- Blocks are **replicated** across multiple nodes for fault tolerance



HDFS is deployed across multiple servers, typically called a *cluster* and HDFS manages how data is stored across the cluster.

HDFS consists of **NameNode** and **DataNode** services providing the basis for the distributed filesystem. When you upload a file to HDFS, the file is first chunked into blocks of a fixed size, typically between 64 MB and 256 MB. Each block is then replicated across multiple DataNodes (typically three).

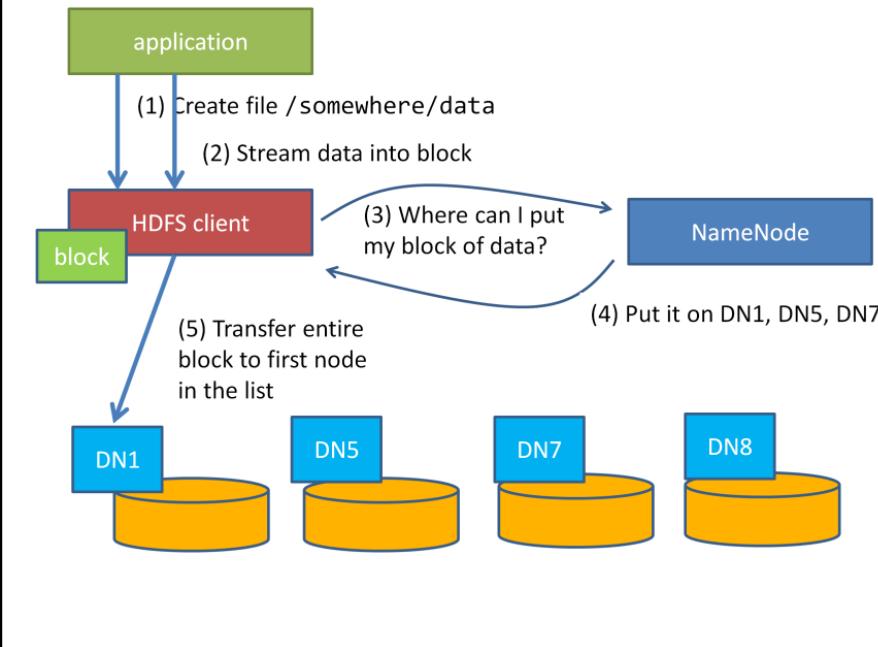
The NameNode stores, manages, and serves the metadata of the filesystem, but does not store any real data blocks. When retrieving data, client applications first contact the NameNode to get the list of locations the requested data resides in and then contact the DataNodes directly to retrieve the actual data.

Hadoop brings in several performance, scalability, and reliability improvements.

- **High Availability (HA)** support for the HDFS NameNode: manual and automatic failover capabilities for the HDFS NameNode service, avoiding the NameNode single point of failure weakness
- **HDFS Federation** enables the usage of multiple independent HDFS namespaces in a single HDFS cluster. These namespaces would be managed by independent NameNodes, but share the DataNodes of the cluster to store the data. The HDFS federation feature improves the horizontal scalability of HDFS by allowing us to distribute the workload of NameNodes.

- Other important improvements of include the support for HDFS snapshots, heterogeneous storage hierarchy support and in-memory data caching support

Writing phase 1: staging on client

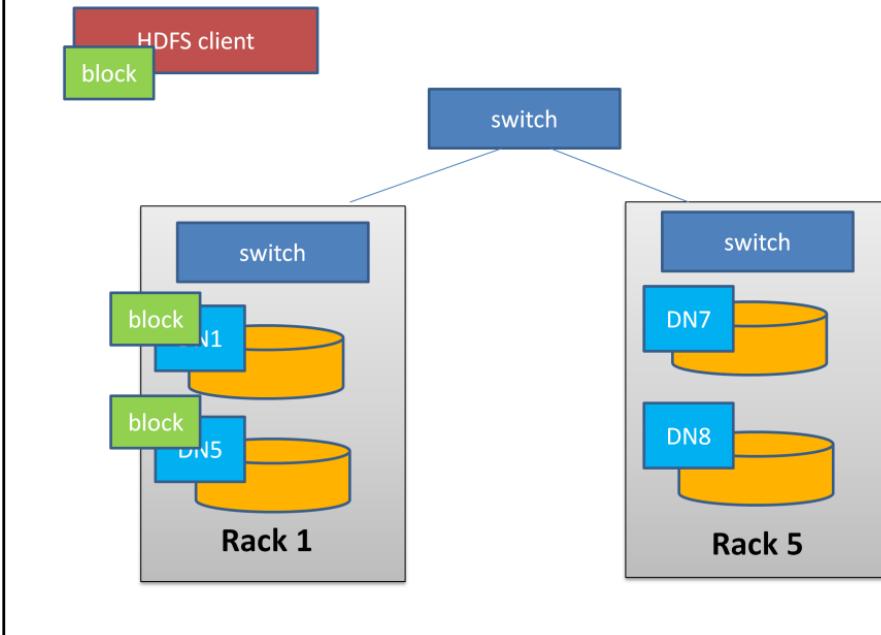


A client request to create a file does not reach the NameNode immediately. In fact, initially the HDFS client caches the file data into a temporary local file. Application writes are transparently redirected to this temporary local file. When the local file accumulates data worth over one HDFS block size, the client contacts the NameNode. The NameNode inserts the file name into the file system hierarchy and allocates a data block for it. The NameNode responds to the client request with the identity of the DataNode(s) and the destination data block. Then the client flushes the block of data from the local temporary file to the first of the specified DataNodes by communicating with the DataNode daemon (DNx boxes in the figure).

When a file is closed, the remaining un-flushed data in the temporary local file is transferred to the DataNode. The client then tells the NameNode that the file is closed. At this point, the NameNode commits the file creation operation into a persistent store. If the NameNode dies before the file is closed, the file is lost.

This design approach has been adopted after careful consideration of target applications that run on HDFS. These applications need streaming writes to files. If a client writes to a remote file directly without any client side buffering, the network speed and the congestion in the network impacts throughput considerably. This approach is not without precedent.

Writing phase 2: pipelined replication



When a client is writing data to an HDFS file, its data is first written to a local file as explained in the previous slide. Suppose the HDFS file has a replication factor of three. When the local file accumulates a full block of user data, the client retrieves a list of DataNodes from the NameNode. This list contains the DataNodes that will host a replica of that block. The client then flushes the data block to the first DataNode. The first DataNode starts receiving the data in small portions, writes each portion to its local repository and transfers that portion to the second DataNode in the list. The second DataNode, in turn starts receiving each portion of the data block, writes that portion to its repository and then flushes that portion to the third DataNode. Finally, the third DataNode writes the data to its local repository. Thus, a DataNode can be receiving data from the previous one in the pipeline and at the same time forwarding data to the next one in the pipeline. Thus, the data is pipelined from one DataNode to the next.

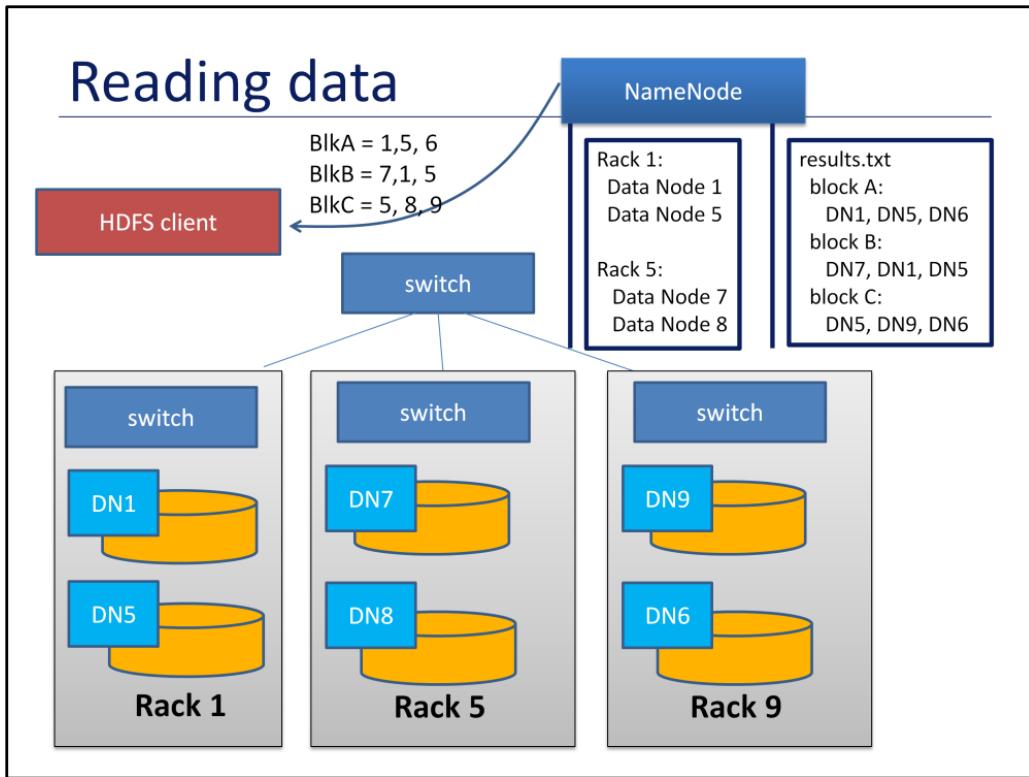
The placement of replicas is critical to HDFS reliability and performance. Optimizing replica placement is a feature that needs lots of tuning and experience. HDFS has the concept of “rack awareness”. The purpose of a rack-aware replica placement policy is to improve data reliability, availability, and network bandwidth utilization. The current default implementation for the replica placement policy is only a first effort and more sophisticated placement strategies may be introduced in the next versions.

Large HDFS instances run on a cluster of computers that commonly spread across

many racks. Communication between two nodes in different racks has to go through switches. In most cases, network bandwidth between machines in the same rack is greater than network bandwidth between machines in different racks. The NameNode knows the rack id each DataNode belongs to.

A simple but non-optimal policy is to place replicas on unique racks. This prevents losing data when an entire rack fails and allows use of bandwidth from multiple racks when reading data. This policy evenly distributes replicas in the cluster which makes it easy to balance load on component failure. However, this policy increases the cost of writes because a write needs to transfer blocks to multiple racks.

For the common case, when the replication factor is three, HDFS's placement policy is to put one replica on one node in the local rack, another on a different node in the local rack, and the last on a different node in a different rack. This policy cuts the inter-rack write traffic which generally improves write performance. The chance of rack failure is far less than that of node failure; this policy does not impact data reliability and availability guarantees. However, it does reduce the aggregate network bandwidth used when reading data since a block is placed in only two unique racks rather than three. With this policy, the replicas of a file do not evenly distribute across the racks. One third of replicas are on one node, two thirds of replicas are on one rack, and the other third are evenly distributed across the remaining racks. This policy improves write performance without compromising data reliability or read performance.

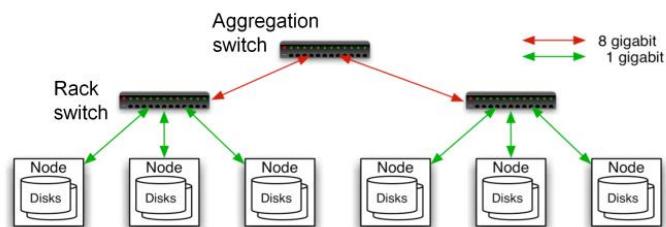


When a Client wants to retrieve a file from HDFS, perhaps the output of a job, it again consults the Name Node and asks for the block locations of the file. The Name Node returns a list of each Data Node holding a block, for each block. The client then picks one of the block location for each block.

To minimize global bandwidth consumption and read latency, the NameNode tries to satisfy a read request from a replica that is closest to the reader. The HDFS client (the reader) is not always located outside of the cluster. In practice, read requests will often come from applications deployed on datanodes (see later, when we discuss YARN/MapReduce). If there exists a replica on the same rack as the reader node, then that replica is preferred to satisfy the read request. For this reason, the order in which the NameNode returns the list of nodes per block is important: the first node in the list is preferred from the viewpoint of cluster load. Unless the first node has crashed and the namenode has not noticed this yet, (legitimate) clients have no reason for not choosing the first replica in the list.

Data model for reading

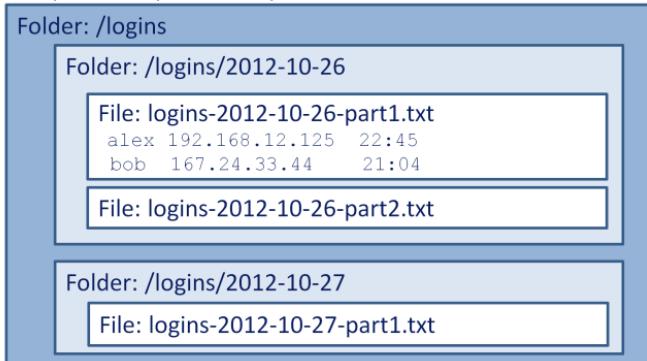
- Client operation
 - query NameNode for block location
 - client accesses data directly from DataNode
- Rack-aware replica selection and placement
 - heartbeat: NameNode detects DataNode failures
 - NameNode balances disk usage and communication traffic



Tweaking HDFS performance

Conflicting requirements:

- Avoid many small files
 - NameNode keeps metadata per file and per block
 - for same amount of data, large sequential read outperforms several random reads
 - Each block is processed by individual worker
- Vertical partitioning in directories
 - operate only on data you actually need



Small file problem

HDFS does not work well with lots of small files and instead wants fewer large files. A small file can be defined as any file that is significantly smaller than the Hadoop block size. The Hadoop block size is usually set to 64,128, or 256 MB, trending toward increasingly larger block sizes. However, the small file problem does not just affect small files. If a large number of files in your Hadoop cluster are marginally larger than an increment of your block size you will encounter the same challenges as small files. For example if your block size is 128MB but all of the files you load into Hadoop are 136MB you will have a significant number of small 8MB blocks. The good news is that this can be easily solved by choosing an appropriate (larger) block size.

There are two primary reasons that HDFS has a small file problem:

- **Namenode memory** - Every directory, file, and block in Hadoop is represented as an object in memory on the NameNode. As a rule of thumb, each object requires 150 bytes of memory. If you have a billion files each requiring just one block, this will require 300GB of memory and that is assuming every file is in the same folder! In addition, the NameNode must constantly track and check where every block of data is stored in the cluster. This is done by listening for data nodes to report on all of their blocks of data. The more blocks a data node must report, the more network bandwidth it will consume. Even with high-speed interconnects between the nodes, simple block reporting at this scale could become disruptive.
- **Processing delays** – As we will see later on, HDFS is typically used in conjunction

with batch processing tools like MapReduce. A large number of small files means a large number of random disk IO. Disk IO is often one of the biggest limiting factors in MapReduce performance. One large sequential read will always outperform reading the same amount of data via several random reads. A second processing delay comes from how batch processing frameworks work. Each block will be processed by an individual worker. In the example of MapReduce, each worker (“map task”) is run in its own Java Virtual Machine. If you have 10 000 files each containing 10 MB of data, you will have the overhead of spinning up and tearing down just 10 000 Java Virtual Machines. If instead you have 800 files of 128 MB each, you only need 800 workers.

Vertical partitioning

Data is of course stored to be used by applications. You can make your processing application much more efficient if you partition your data so that the application can easily access data relevant to its computation. This process is called *vertical partitioning*, and avoid reading in and filtering out data that you don’t need.

Vertically partitioning data on a distributed file system can be done by sorting your data into separate folders. For example, suppose you are storing login information on a distributed file system. Each login contains a username, IP address and timestamp. To vertically partition by day, you can create a separate folder for each day of data. Now if you only want to look at a particular subset of your dataset, you can just look at the files in those particular folders and ignore the other files.

As a file can only reside in one directory, this vertical partitioning means splitting the data in different files. Hence, you should trade-off the benefits of vertical partitioning with the performance penalty of smaller file sizes.

Outline

- Types of distributed storage solutions
 - Distributed file system
 - case study: Hadoop
 - Distributed data stores
 - key-value; columnar; document; graph
- Reasons for distributing
 - sharding
 - replication

Distributed data stores

- Write-once-read-many does not fit all needs
 - fast random reads
 - real-time updating of data views
- Data stores provide different data model
 - data is indexed
 - data organization allows for fast random writes
- We will discuss the following aspects
 - matching big data with relational databases
 - types of distributed data stores and their trade-offs
 - techniques for distributing indexed data

The write-once-read-many paradigm of distributed file systems does not match all applications needs. Other applications may be better off with random read/write operations.

Databases are efficient in random reads because they index data. Moreover, many databases have advanced ways of organizing data, which favors fast random write operations.

As we will see, relational databases do not match well the requirements of many big data applications. We will study a new breed of data storage solutions – often referred to as “NoSQL”. Our primary focus is on how these techniques ensure scalability and fault tolerance. Because these new technologies often lack support for SQL queries, we often refer to them as data *stores* rather than databases.

Relational databases and big data

- Relational databases have unique strengths:
 - Consistency
 - Advanced querying
- But this doesn't always fit big data's needs

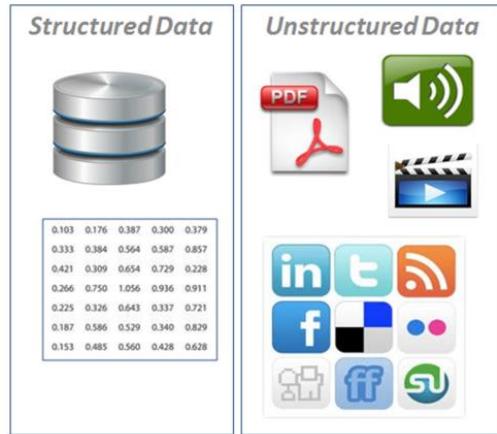


Relational databases provide solid, mature services according to the ACID properties. We get transaction-handling, efficient logging to enable recovery etc. These are core services of relational databases, and the ones that they are good at. They are hard to customize, and might be considered as a bottleneck, especially if you don't need them in a given application (e.g. serving website content with low importance).

Lots of "big data" problems don't require these strict constraints, for example web analytics, web search or processing moving object trajectories, as they already include uncertainty by nature.

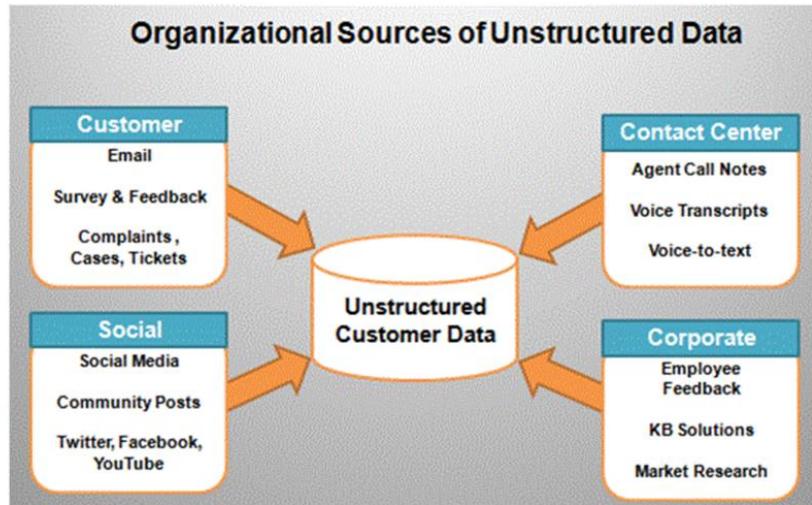
Variety: unstructured data

- big data is semi-structured or unstructured
 - blogs, e-mails, videos, tweets...
 - structure differs between records
 - cannot be mapped to tables (think Excel sheet) in a meaningful way



Big data is typically semi-structured. This means that the individual data records do not conform with a formal data model (as is the case in relational databases), but nonetheless they contain tags or other markers to separate semantic elements and enforce hierarchies of records and fields within the data. These tags and markers make the data records a self-describing structure.

Variety example: customer data



More insight is gathered via text analysis than via SQL queries

Unstructured data is often text-heavy. The useful knowledge from big data is gathered from applying text analysis on this text (e.g. using tools like MapReduce), rather than from launching advanced SQL queries.

Variety example: storing e-mails

Sender	Date	Body
serge.brin@abc.xyz	12-08-15 01:15 PM
rector@ugent.be	14-08-15 08:05 AM
...

- E-mails can be saved in semi-structured way
- Only very limited knowledge about customers can be derived using queries
 - What would be a useful query for the body column?

Continuing the example of the previous slide, let's look at the e-mails received. Although all e-mails have a body, it is hard to imagine a useful query for the body column.

Variety example: webshop

Product information of a webshop

- Every product has a name, unit price and vendor
- Attributes differ per product:
 - CPU has clock rate, cache size, # of cores
 - monitor has size, resolution
 - RAM has capacity, technology type

How to store this information in a relational database?

?? Very wide table with hundreds of fields for any possible product attribute ??
?? Separate table per product category ??
?? ...



Store product description as a complete document
with varying internal structure

Consider the example of a webshop. Each product (category) has many different attributes. It is not straightforward to map this efficiently in the table structure of a relational database.

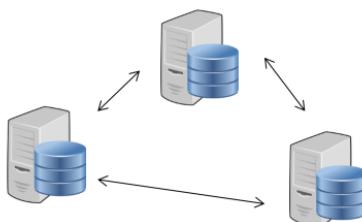
You could create a very wide table with hundreds of field for any possible attribute some product could have, but for most product most of these fields will be NULL. Adding new columns to a relational DB requires the system be shut down and ALTER TABLE commands to be run. When a database is large, this process can impact system availability, costing time and money.

You could have a separate table for each product category, thus introducing a lot of tables and relationships.

Although all options are valid, none of them is really satisfying. The fundamental problem is the rigid database schema structures imposed by relational DBs that do not match our needs.

Volume/Velocity: Distributing relational databases

- JOINS require data from multiple nodes
- ACID guarantees hard to maintain
 - Atomicity/Consistency/Isolation/Durability
 - enforced via transactions
 - intrinsically slower because of locking
 - CAP → system is unavailable from time to time



When the data becomes too large for a single node, it will be split over multiple nodes (shards). At a high level, this can be done in two ways for a relational DB: either you split the rows of a large table, either you distribute on a per-table basis (in practice, both approaches are combined). But in both ways multiple nodes are involved when we perform a JOIN operation: JOINS require the full data of each table being joined.

Another strength of relational DBs is their ACID-compliance. In a relational database, the standards (and most of the existing products) are built around the concept of *transactions*, which are maintained in a definite order by a *transaction log*, which is by its nature a global scope data structure. The transaction log allows a database to guarantee that the current state of the database can be provably reached from some earlier state of the database by the application, in sequence, of a series of transactions. This is the basis of various ACID guarantees. ACID-compliant transaction means the database is designed so it absolutely will not lose data:

- ✓ Each operation moves the database from one valid state to another (Atomic).
- ✓ Everyone has the same view of the data at any point in time (Consistent).
- ✓ Operations on the database don't interfere with each other (Isolation).
- ✓ When a database says it has saved data, you know the data is safe (Durable).

Regarding Consistency (the "C" in ACID), real-time consistency is particularly difficult to implement in a partitioned system, because transactions that occur across (or

draw data from) multiple partitions need to be globally ordered ("serializable" in transaction terminology) in order to ensure that operations are drawing from a consistent view of the system.

NoSQL systems that use partitioning tend to solve this by relaxing the Consistency requirement such that it applies at a per-partition level, and not at a system level. This is a very reasonable engineering trade-off for most systems, and removes the greatest barrier for scalability that exists in distributed systems.

NoSQL



- often favor “A” over “C” in CAP
- key features :
 - schema agnostic (schema-on-read)
 - non-relational, aggregate data models
 - scaling out on commodity hardware
 - highly distributable

Four core features apply to most NoSQL databases.

Schema Agnosticism

A database schema is the description of all possible data and data structures in a relational database. With a NoSQL database, a schema isn’t required, giving you the freedom to store information without doing up-front schema design. You are not required to do a lot of up-front design work before you can store data in NoSQL databases.

An alternative interpretation of schema agnosticism is *schema on read*. You need to know how the data is stored only when constructing a *query* (a coded question that retrieves information from the database), so for practical purposes, this feature is exactly what it says: You need to know the schema on read.

Nonrelational

Relations in a database establish connections between tables of data. For example, a list of transaction details can be connected to a separate list of delivery details. With a NoSQL database, this information is stored as an aggregate — a single record with everything about the transaction, including the delivery address.

Commodity hardware

Some databases are designed to operate best (or only) with specialized storage and

processing hardware. With a NoSQL database, cheap off-the-shelf servers can be used. Adding more of these cheap servers allows NoSQL databases to scale to handle more data.

Highly distributable

Distributed databases can store and process a set of information on more than one device. With a NoSQL database, a cluster of servers can be used to hold a single large database.

NoSQL features

- aggregate-oriented
 - mismatch between relational database and data structures of application developers
 - collection of data that we interact with as a unit
 - key-value, document, column family, graph
 - inter-aggregate relationships difficult to handle
- distribution models
 - sharding
 - different subset of data across servers
 - replication
 - copy data across servers
 - master-slave or peer-to-peer

Application developers have been frustrated with the impedance mismatch between the relational data structures and the in-memory data structures of the application. Using NoSQL databases allows developers to develop without having to convert in-memory structures to relational structures.

Aggregate Data Models:

Relational database modelling is vastly different than the types of data structures that application developers use. Using the data structures as modelled by the developers to solve different problem domains has given rise to movement away from relational modelling and towards aggregate models. An aggregate is a collection of data that we interact with as a unit. These units of data or aggregates form the boundaries for ACID operations with the database. Key-value, Document, and Column-family data stores (all discussed later) can all be seen as forms of aggregate-oriented database.

Aggregates make it easier for the database to manage data storage over clusters, since the unit of data now could reside on any machine and when retrieved from the database gets all the related data along with it. Aggregate-oriented databases work best when most data interaction is done with the same aggregate, for example when there is need to get an order and all its details, it better to store order as an aggregate object but dealing with these aggregates to get item details on all the orders is not elegant.

Aggregate-oriented databases make inter-aggregate relationships more difficult to handle than intra-aggregate relationships. Aggregate-ignorant databases are better when interactions use data organized in many different formations.

Distribution Models:

Aggregate oriented databases make distribution of data easier, since the distribution mechanism has to move the aggregate and not have to worry about related data, as all the related data is contained in the aggregate. There are two styles of distributing data:

- **Sharding:** Sharding distributes different data across multiple servers, so each server acts as the single source for a subset of data.
- **Replication:** Replication copies data across multiple servers, so each bit of data can be found in multiple places. Replication comes in two forms,
 - Master-slave replication makes one node the authoritative copy that handles writes while slaves synchronize with the master and may handle reads.
 - Peer-to-peer replication allows writes to any node; the nodes coordinate to synchronize their copies of the data.

Master-slave replication reduces the chance of update conflicts but peer-to-peer replication avoids loading all writes onto a single server creating a single point of failure. A system may use either or both techniques.

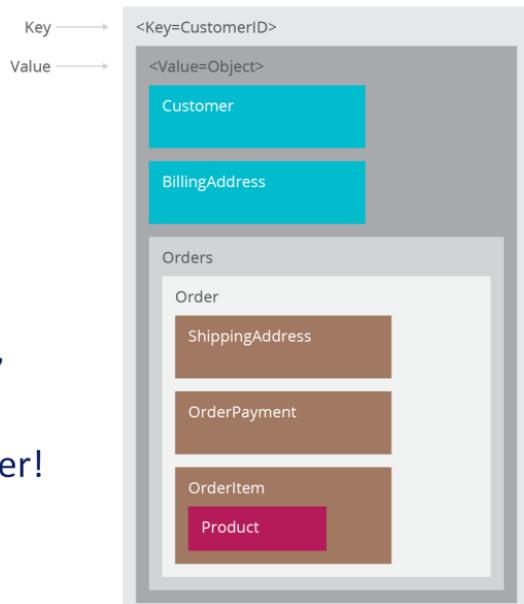
Common NoSQL aggregation types

Type	Description	Typical uses
Key-value	Associate large data file with a simple text string	Dictionary, image/document/file store, query cache, lookup tables
Graph	Store nodes and arcs of a graph	Social network queries, friend-of-friends queries, inference, rules system, pattern matching
Column family	Store a sparse matrix data using a row and column index	Web crawling, large sparsely populated tables, highly-adaptable systems, systems with high variance
Document	Store tree-structured hierarchical information in a single unit	Any data that has a natural container structure: office documents, sales orders, invoices, product descriptions, forms, web pages

One of the challenges for users of NoSQL systems is there are many different architectural patterns from which to choose. The table above lists the significant data architecture patterns associated with the NoSQL movement.

Key-value databases

- very simple API
- value is a “blob” for database
- easy to scale
- examples: Riak, Redis, Memcached, Amazon S3
 - properties may differ!



A *key-value store* is a simple database that when presented with a simple string (the key) returns an arbitrary large BLOB of data (the value). Key-value stores have no query language; they provide a way to add and remove key-value pairs (a combination of key and value where the key is bound to the value until a new value is assigned) into/from a database.

The dictionary is a simple key-value store where word entries represent keys and definitions represent values. Like the dictionary, a key-value store is also indexed by the key; the key points directly to the value, resulting in rapid retrieval, regardless of the number of items in your store.

One of the benefits of not specifying a data type for the value of a key-value store is that you can store any data type that you want in the value. The system will store the information as a BLOB and return the same BLOB when a GET (retrieval) request is made. It's up to the application to determine what type of data is being used, such as a string, XML file, or binary image. All key-value databases are not the same, there are major differences between these products, for example: Memcached data is not persistent while in Riak it is, these features are important when implementing certain solutions. Let's consider we need to implement caching of user preferences, implementing them in memcached means when the node goes down all the data is lost and needs to be refreshed from source system, if we store the same data in Riak we may not need to worry about losing data but we must also consider how to

update stale data. It's important to not only choose *a* key-value database based on your requirements, it's also important to choose *which* key-value database.

Example key-value: RIAK

- Key features:
 - Resilience
 - Link walking and support for map-reduce
- Keys are organized in buckets



- Values can be anything: XML, JSON, images...

Riak is a distributed key-value database where values can be anything—from plain text, JSON, or XML to images or video clips—all accessible through a simple HTTP interface. Riak is also fault-tolerant. Servers can go up or down at any moment with no single point of failure. But this flexibility has some trade-offs. Riak lacks robust support for ad hoc queries, and key-value stores, by design, have trouble linking values together

(in other words, they have no foreign keys). Riak attacks these problems by providing mechanisms like link walking and map reduce.

Riak breaks up classes of keys into *buckets* to avoid key collisions—for example, a key for java the *language* will not collide with java the *drink*.

PUT - GET

http://SERVER:PORT/riak/BUCKET/KEY

```
curl -i -X PUT http://SERVER:PORT/riak/albums/born-this-way \
-H "Content-Type: application_json" \
-d '{"release-date": "2011-05-23"}'
```

```
curl http://SERVER:PORT/riak/albums/born-this-way
```

```
HTTP/1.1 200 OK
Content-Type: application/json
Content-Length: ...

{"release-date": "2011-05-23"}
```

Interfacing with Riak happens through HTTP requests. You query via URLs, headers, and verbs, and Riak returns assets and standard HTTP response codes.

In the example on the slide, we *put* a JSON file with metadata with key “born-this-way” in the bucket “albums”. If you use a POST request to /riak/albums, then RIAK will generate the key itself. A GET request to the same location will retrieve the value.

RIAK – links to add metadata

One-way, multiple links per key possible

```
curl -i -X PUT http://SERVER:PORT/riak/albums/born-this-way \
-H "Content-Type: application/json" \
-H "Link: <riak/songs/judas>; riaktag=\"lists_track\""
-d '{"release-date": "2011-05-23"}'
```

```
curl -i -X PUT http://SERVER:PORT/riak/songs/judas \
-H "Content-Type: application/json" \
-H "Link: <riak/mp3s/62542>; riaktag=\"is_audio_file\""
-d '{"release-date": "2011-05-23"}'
```

```
curl -i -X PUT http://SERVER:PORT/riak/mp3/62542 \
-H "Content-Type: audio/mpeg3" \
-H "Link: <riak/mp3s/62542>; riaktag=\"is_audio_file\""
-d ...
```

One of the ways that we are able to extend the fairly-limited data model provided by a key/value store is with the notion “links” and a type of query known as “link walking.”

Links are metadata that establish one-way relationships between objects by associating one key to other keys. The basic structure is this:

Link: </riak/bucket/key>; riaktag="whatever"

The key to where this value links is in pointy brackets (<...>), followed by a semicolon and then a tag describing how the link relates to this value (any string). The links are attached by adding a “Link” header in the PUT (or POST) HTTP request. In the example we are attaching a link with tag “lists_track” to the key “Judas” in the bucket “songs”. We are also adding the key “Judas”, which points to a JSON with metadata about the song and has a link “is_audio_file” to the key corresponding with the correct mp3 in the bucket “mp3”. You can of course add multiple links to one object. The object that the link is pointing to should not even be in the database at the moment you register the link.

Links are unidirectional: an object is not aware of the links pointing to it.

RIAK – link walking

bucket, tag, keep

Get all objects in the bucket “songs” that are linked with tag “lists_track” on the album “born-this-way”:

```
curl http://SERVER:PORT/riak/albums/born-this-way/songs,lists_track,1
```

Get all objects in any bucket that are linked with tag “lists_track” on the album “born-this-way”:

```
curl http://SERVER:PORT/riak/albums/born-this-way/_lists_track,1
```

Get all mp3s that are on the album “born-this-way”:

```
curl http://SERVER:PORT/riak/albums/born-this-way/_lists_track,_mp3,is_audio_file,1
```

Once you have tagged objects in Riak with links, you can then traverse them with an operation called “Link Walking.” With links, you create lightweight pointers between your data, for example, from ‘projects’ to ‘milestones’ to ‘tasks’, and then select data along that hierarchy using simple client API commands. This can substitute as a lightweight graph database, as long as the number of links attached to a given key are kept reasonably low. Links are an incredibly powerful feature of Riak.

Getting the linked data is achieved by appending a *link spec* to the end of the URL that is structured like this: /bucket,tag,keep

- Bucket – a bucket name to limit the links to
- Tag – the “riaktag” to limit the links
- Keep – 0 or 1, whether to return results from this step or phase

Each of these three elements can be replaced by an underscore (_) representing wildcards.

You can walk any number of links with one request by chaining multiple *link specs*. In the example above, we request from the database all objects in the bucket “mp3” that are pointed to by a link tagged “is_audio_file” from any object (in any bucket) that is inked with a tag “lists_track” from the object with key “born-this-way” in the bucket “albums”. By default, Riak will only include the objects found by the last step. If we would have put a “1” in the first link spec, then Riak would also return the intermediate objects in the GET response.

Document databases

- Documents encapsulate and encode data in some standard format (XML, JSON...)
- Document is automatically indexed into tree-like structure
- Rich query language
- Example: MongoDB, CouchDB

```
<Key=CustomerID>
{
  "customerid": "fc986e48ca6" ←
  "customer":
  {
    "firstname": "Pramod",
    "lastname": "Sadalage",
    "company": "ThoughtWorks",
    "likes": [ "Biking", "Photography" ]
  }
  "billingaddress":
  {
    "state": "AK",
    "city": "DILLINGHAM",
    "type": "R"
  }
}
```

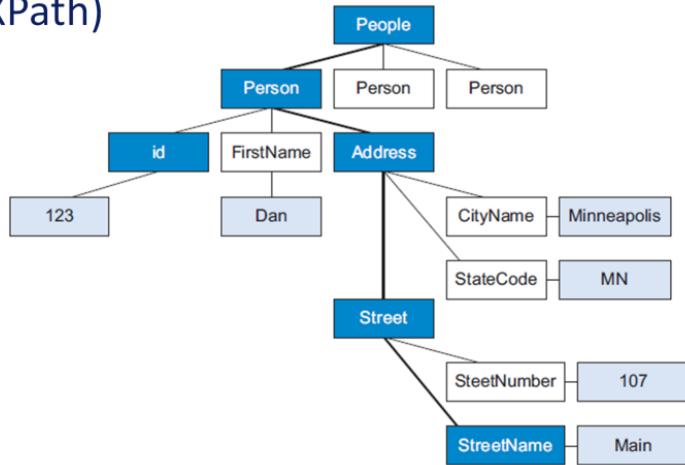
Documents are the main concept in document databases. The database stores and retrieves documents, which can be XML, JSON, BSON, and so on. These documents are self-describing, hierarchical tree data structures which can consist of maps, collections, and scalar values. **The documents stored are similar to each other but do not have to be exactly the same.** Compared to relational databases, for example, collections could be considered analogous to tables and documents analogous to records. But they are different: every record in a table has the same sequence of fields, while documents in a collection may have fields that are completely different.

Document databases store documents in the value part of the key-value store. While key-value stores, when presented with a key, return a blob of data that lacks a formal structure and is not indexed or searchable, document stores work in the opposite manner: the key may be a simple ID which is never used or seen.

But you can get almost any item out of a document store by querying any value or content within the document. Document databases such as MongoDB provide a rich query language and constructs such as database, indexes etc allowing for easier transition from relational databases.

Document DBs: tree structure

Retrieve documents based on their content
(cfr. XPath)



`People/Person[id='123']/Address/Street/StreetName/text()`

Think of a document store as a tree-like structure. Document trees have a single root element (or sometimes multiple root elements). Beneath the root element there is a sequence of branches, sub-branches, and values. Each branch has a related path expression that shows you how to navigate from the root of the tree to any given branch, sub-branch, or value. Each branch may have a value associated with that branch. Sometimes the existence of a branch in the tree has specific meaning, and sometimes a branch must have a given value to be interpreted correctly.

Each document store has an API or query language that specifies the path or path expression to any node or group of nodes in the tree. With the query on the slide you begin by selecting a subset of all people records that have the identifier 123. Often this points to a single person. Next you look in the Address section of the record and select the text from the Address street name. The full path name to the street name is the following: `People/Person[id='123']/Address/Street/StreetName/text()`.

Example document store: Mongo

- document-oriented
 - schemaless
 - JSON (BSON) format
- advanced server-side JavaScript queryability
 - ad hoc queries by field, range, regular expression
 - queries can include user-defined JS functions
 - map-reduce
 - indexing
- scalability and redundancy
 - load balancing
 - replication

Mongo is designed as a scalable database—the name Mongo comes from “humongous”—with performance and easy data access as core design goals. Mongo hits a sweet spot between the powerful **queryability** of a relational database and the **distributed nature** of other datastores like Riak or HBase. Mongo is a JSON document database (though technically data is stored in a binary form of JSON known as BSON). A Mongo document can be likened to a relational table row without a schema, whose values can nest to an arbitrary depth. It enforces however no schema (similar to Riak), so documents can optionally contain fields or types that no other document in the collection contains.

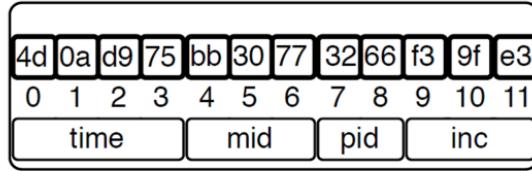
Despite being schemaless, Mongo offers extensive options for querying, including querying by user-defined functions, map-reduce functionality and indexing. These queries are written in JavaScript and executed at the server: the client only receives the result of the query (like with traditional SQL queries). We will discuss some of these features in the next slides.

Mongo’s other strength lies in its ability to handle huge amounts of data (and huge amounts of requests) by replication and horizontal scaling.

Mongo: collections

- Interaction through JavaScript commands in Mongo shell
- Documents are inserted into collections (~ buckets in Riak)
- Each document is automatically given an ID
 - Different nodes in cluster generate non-colliding IDs

```
> db.towns.insert({  
  name: "New York",  
  population: 22200000,  
  last_census: ISODate("2009-07-31"),  
  famous_for: [ "statue of liberty", "food" ],  
  mayor : {  
    name : "Michael Bloomberg",  
    party : "I"  
  }  
})
```



A mongo database comprises multiple *collections*. The concept of collections is similar to a bucket in Riak nomenclature. Interaction happens via JavaScript in a mongo shell. Once you have opened a mongo shell, the variable “db” is a JavaScript object that contains information about the current database. “db.x” is a JavaScript object representing a collection (named x). Commands on these functions are just JavaScript functions.

The command above shows how a new document is added to the collection “towns”. Since Mongo is schemaless, there is no need to define anything up front; merely using it is enough. In the example above, the collection towns didn’t exist before the first document was added to it.

Mongo automatically adds an `_id` field of type `ObjectId` to each JSON document stored. This is akin to a numeric key. The `ObjectId` is always 12 bytes, composed of a timestamp, client machine ID, client process ID, and a 3-byte incremented counter. The crux about this autonumbering scheme is that each process on every machine can handle its own ID generation without colliding. This design choice gives a hint of Mongo’s distributed nature.

Mongo is schemaless

```
> db.towns.insert({  
  name: "Punxsutaway",  
  population: 6200,  
  last_census: ISODate("2008-31-01"),  
  famous_for: [ "phil the groundhog"],  
  mayor : {  
    name : "Jim Wehrle"  
  }  
})  
  
> db.towns.insert({  
  name: "Portland",  
  population: 582000,  
  last_census: ISODate("2007-20-09"),  
  famous_for: [ "beer", "food"],  
  mayor : {  
    name : "Sam Adams",  
    party: "D"  
  }  
})
```

The documents added in a single bucket can have varying structure. Not all documents need to have all fields, fields can be nested with arbitrary depth, etc...

In the example above, the political party of the mayor of Punxsutaway is unknown. In a relational database, we should enter NULL in the corresponding column. In a schemaless document store, we simply omit that value.

While the overlap in the current example is still relatively high; this needs not be the case: data items in the same collection can have a completely different structure.

Mongo: querying

```
> db.towns.find(
  { _id : ObjectId("4d0ada1fbb30773266f39fe4") }, { name : 1 })  
  
{  
  "_id" : ObjectId("4d0ada1fbb30773266f39fe4"),  
  "name" : "Punxsutawney"  
}
```

```
> db.towns.find(
  { _id : ObjectId("4d0ada1fbb30773266f39fe4") }, { name : 0 })  
  
{  
  "_id" : ObjectId("4d0ada1fbb30773266f39fe4"),  
  "population" : 6200,  
  "last_census" : "Thu Jan 31 2008 00:00:00 GMT-0800 (PST)",  
  "famous_for" : [ "phil the groundhog" ]  
}
```

To access a specific document, you call the `find()` function. The first argument is a the `_id` of the document you want to retrieve. The `find()` function also accepts an optional second parameter: a `fields` object you can use to filter which fields are retrieved. If you want only the town name (along with `_id`), pass in “name” with a value resolving to 1 (or true).

Mongo: advanced querying

Use of regex and range operators:

```
> db.towns.find(  
    { name : /^P/, population : { $lt : 10000 } },  
    { _id : 0, name : 1, population : 1 }  
)  
  
{ "name" : "Punxsutawney", "population" : 6200 }
```

Query by matching nested array data:

```
> db.towns.find(  
    { famous_for : /statue/ },  
    { _id : 0, name : 1, famous_for : 1 }  
)  
  
{ "name" : "New York", "famous_for" : ["statue of liberty", "food" ]}
```

In Mongo you can construct ad hoc queries by field values, ranges, or a combination of criteria. The first example on the slide uses a regular expression and a range operator to find all towns that begin with the letter P and have a population less than 10 000. The second example demonstrates the ability to match nested array data.

Mongo: matching subdocuments

```
db.towns.find(  
  { 'mayor.party' : 'I' },  
  { _id : 0, name : 1, mayor : 1 }  
)  
  
{  
  "name" : "New York",  
  "mayor" : {  
    "name" : "Michael Bloomberg",  
    "party" : "I"  
}
```

```
> db.towns.find(  
  { 'mayor.party' : { $exists : false } },  
  { _id : 0, name : 1, mayor : 1 }  
)  
  
{ "name" : "Punxsutawney", "mayor" : { "name" : "Jim Wehrle" } }
```

The true power of Mongo stems from its ability to dig down into a document and return the results of deeply nested subdocuments. To query a subdocument, your field name is a string separating nested layers with a dot.

The examples on the slide shows how to find all towns with independent mayors, or those with mayors who don't have a party.

Mongo: indexing

Collection of phone number documents:

```
> db.phones.find().limit(2)
{ "_id" : 1800555000, "components" : { "country" : 1, "area" : 800,
  "prefix" : 555, "number" : 555000 }, "display" : "+1 800-555000"
}
{ "_id" : 8800555001, "components" : { "country" : 8, "area" : 800,
  "prefix" : 555, "number" : 555001 }, "display" : "+8 800-555001"
}
```

Create your own index on the `display` field

```
> db.phones.ensureIndex(
  { display : 1 },
  { unique : true, dropDups : true }
)
```

One of Mongo's useful built-in features is indexing to increase query performance — something that's not available on all NoSQL databases. MongoDB provides several of data structures for indexing, such as the classic B-tree, and other additions such as two-dimensional and spherical GeoSpatial indexes.

To demonstrate the feature, we are going to do a little experiment of MongoDB's B-tree index on a collection with a series of phone number documents. A B-tree is a generalization of a binary search tree in that a node can have more than two children.

Whenever a new collection is created, Mongo automatically creates an index by the `_id` field. Most queries will include more fields than just the `_id`, so we need to make indexes on those fields. In our example, we create a B-tree index on the `display` field by calling `ensureIndex(fields, options)`. The `fields` parameter is an object containing the fields to be indexed against. The `options` parameter describes the type of index to make. In this case, we are building a unique index on `display` that should just drop duplicate entries.

Speed of indexing

Before indexing:

```
> db.phones.find({display: "+1 800-5650001"}).explain()
{
  "cursor" : "BasicCursor",
  "nscanned" : 109999,
  "nscannedObjects" : 109999,
  "n" : 1,
  "millis" : 52,
  "indexBounds" : { }
}
```

After indexing:

```
> db.phones.find({display: "+1 800-5650001"}).explain()
{
  "cursor" : "BtreeCursor display_1",
  "nscanned" : 1,
  "nscannedObjects" : 1,
  "n" : 1,
  "millis" : 3,
  "indexBounds" : { "display" : [ [ "+1 800-5650001", "+1 800-5650001" ] ] }
}
```

The `explain()` method can be used to output details of a given operation. In the output, the `millies` field shows the time needed to complete the query.

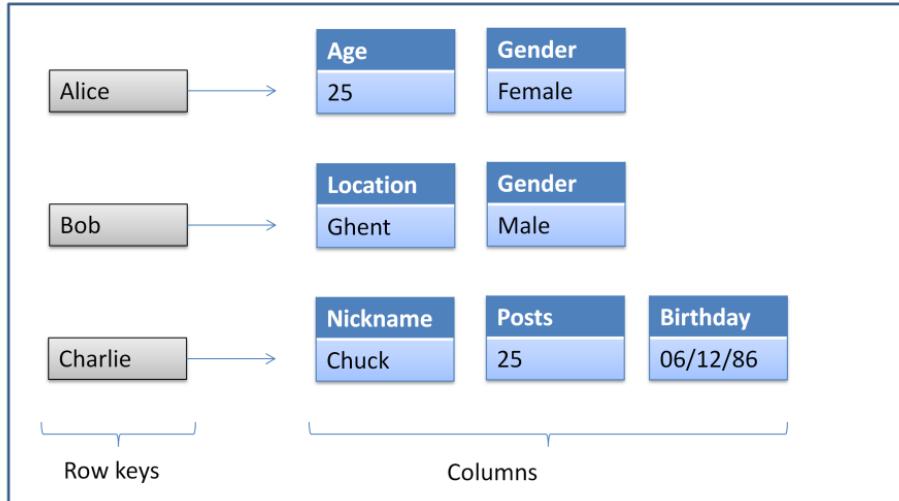
When we rerun `explain` after indexing, the `millies` values will be dropped by at least an order of magnitude. Note how the cursor has changed from a Basic to a B-tree cursor. Mongo is no longer doing a full collection scan but instead walking the B-tree to retrieve the value. Importantly, the number of scanned objects dropped from 109 999 to 1 since it has become just a unique lookup (based on the index).

Just like queries can be nested, you can build your index on nested values. If you wanted to index on all area codes, use the dot-notated field representation. In production, you should always build indexes in the background. This gives the following command:

```
> db.phones.ensureIndex({ "components.area": 1 }, { background : 1 })
```

Creating an index on a large collection can be slow and resource-intensive. You should always consider these impacts when building an index by creating indexes off-peak times, running index creation in the background, and running them manually rather than using automated index creation.

Column family data stores



Column family: users

Examples: Cassandra, BigTable

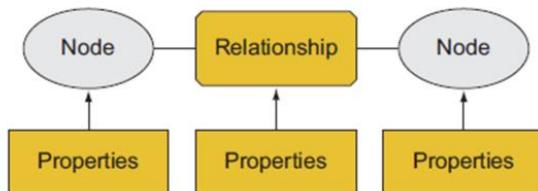
Column-family databases store data in column families as rows that have many columns associated with a row key. The row and column identifiers are used as general purpose keys for data lookup. Various rows do not have to have the same columns, and columns can be added to any row at any time without having to add it to other rows.

This type is often referred to as *data stores* rather than *databases*, since this kind of storage solutions lacks features you may expect to find in traditional databases. For example, they lack typed columns, secondary indexes, triggers, and query languages.

An alternative view on the datamodel is to consider it as a map with sorted maps as values (or optionally, a sorted map with sorted maps as values). To expand upon the terms:

- Column families – column families are analogous to tables in relational databases, and each column family stores a completely independent set of information.
- Keys – if you consider a column family as a giant map, keys are the top-level entries in the map. Keys are used to partition a column family across a cluster
- Columns – each key points to another map of name/value pairs called columns. All columns for a key are physically stored together, making it inexpensive to access ranges of columns. Different keys can have different sets of columns, and it's possible to have millions of columns for a given key.

Graph databases



- Store three datafields: nodes, relationships and properties
- Difficult to shard (unless you have an unconnected graph)
- Analyze relationships between objects or visit all nodes in a graph in a particular manner (graph traversal)
- Use cases:
 - link analysis, e.g. friends-of-friends in social network, fraud detection
 - rules and inference (semantic technologies)
 - integration of public datasets
- Example: Neo4J, Titan

A *graph store* is a system that contains a sequence of nodes and relationships that, when combined, create a graph. A graph store has three data fields: *nodes*, *relationships*, and *properties*. Some types of graph stores are referred to as *triple stores* because of their node-relationship-node structure.

Unlike other NoSQL patterns, graph stores are difficult to scale out on multiple servers due to the close connectedness of each node in a graph. Data can be replicated on multiple servers to enhance read and query performance, but writes to multiple servers and graph queries that span multiple nodes can be complex to implement.

We list three uses cases where a graph store can be used to effectively solve a particular business problem:

- *link analysis*: for business problems that are solved by traversing graph data to search and look for patterns and relationships. The canonical example is to find friends of your friends in a social network, but graph stores are appropriate for identifying distinct patterns of connections between nodes. For example, creating a graph of all incoming and outgoing phone calls between people in a prison might show a concentration of calls (patterns) associated with organized crime. Analyzing the movement of funds between bank accounts might show patterns of money laundering or credit card fraud.
- *rules and inference*: relationships form the basis of the *semantic web stack* (e.g.

Using ontologies and querying languages like SPARQL).

- *dataset integration*: graph stores can be used to automatically join datasets that were created by different organizations (e.g. linked open data). This requires that nodes in the individual datasets can be unambiguously correlated, so that they appear as one node in the merged graph.

Example



This is an example of how data is modelled in graph data stores like Neo4J. The black blobs indicate nodes. The upper left node has a single property ("name" – not shown) with the value "Wine Expert Monthly". This magazine reviewed the wine "Prancing Wolf Ice Wine 2007", which is represented as a node with one property ("name"). The edge between these two nodes is annotated with the type of relationship: "reported_on".

This particular ice wine is created from the *riesling* grape. We could add this as a property directly to the wine node, but riesling is a general category that could apply to other wines. For this reason, we better create a new node and set its property to [name: "riesling"]. We also introduce a new relationship as *grape_type* and give it the property [style: "ice wine"].

Lastly, we do similar operations for the wine manufacturer that produces various wines.

If you want to see some more examples using Neo4J, you can refer to this slide deck:
<http://www.slideshare.net/peterneubauer/neo4j-5-cool-graph-examples-4473985>

NoSQL = end of SQL?

NO ☺

- NoSQL is not a panacea/silver bullet
 - it answers fundamentally different data problems than relational DBs
 - it has different trade-offs (CAP, ACID)
- “SQL doesn’t scale” is a myth!
 - Facebook does it...
 - but requires good engineering
- NewSQL databases: Google Spanner, ...
 - borrow some ideas from NoSQL
 - retain support for SQL queries and/or ACID

Outline

- Types of distributed storage solutions
 - Distributed file system
 - case study: Hadoop
 - Distributed data stores
 - key-value; columnar; document; graph
- Reasons for distributing
 - replication
 - sharding

The need for distributing data



You **will** run against the limitations of a single node

- CPU, memory, disk speed, data size, network bandwidth



This is even more likely in the cloud

- mostly commodity hardware
- multi-tenancy causes resource contention

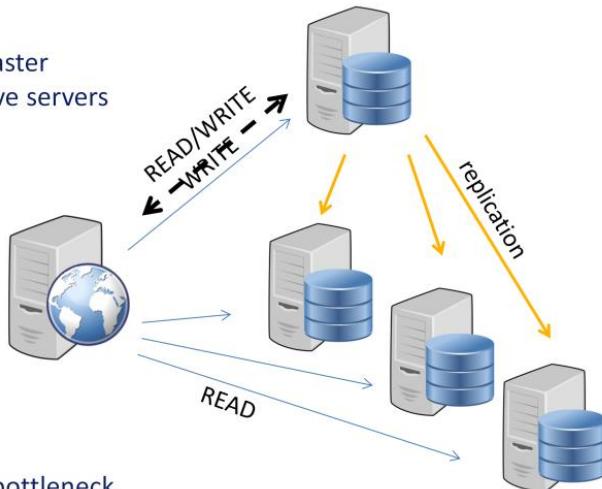
Database systems with large data sets and high throughput applications can challenge the capacity of a single server. High query rates can exhaust the CPU capacity of the server. Larger data sets exceed the storage capacity of a single machine. Finally, working set sizes larger than the system's RAM stress the I/O capacity of disk drives. Also the network bandwidth from/to the database node can be bottleneck.

This limited capacity is exacerbated with cloud database services because the database is running on commodity hardware and the database server is multitenant.

Scaling through replication

Read-write split

- Writes go to master
- Reads go to slave servers



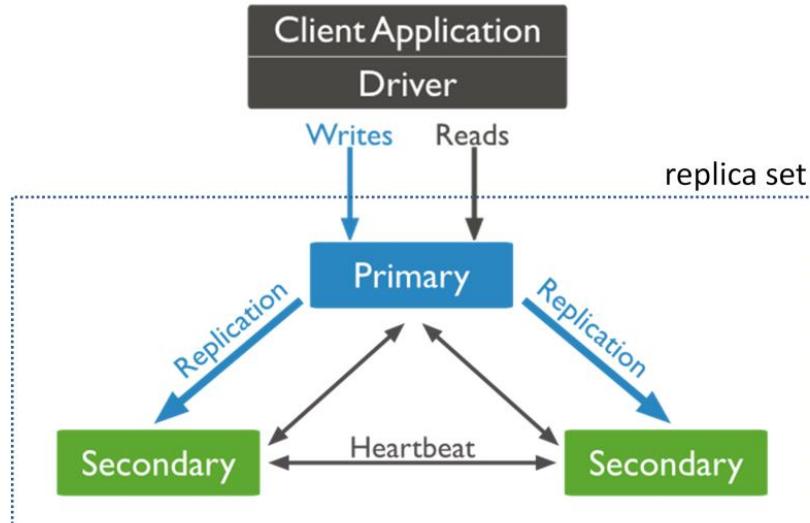
- master server is bottleneck
- eventually consistent slaves

Replication provides redundancy and increases data availability. With multiple copies of data on different database servers, replication protects a database from the loss of a single server. Replication also allows you to recover from hardware failure and service interruptions. With additional copies of the data, you can dedicate one to disaster recovery, reporting, or backup. In some cases, you can use replication to increase read capacity. Clients have the ability to send read and write operations to different servers. You can also maintain copies in different data centers to increase the locality and availability of data for distributed applications.

Master-Slave partitioning is the simplest option, with a single Master server for all write (Create Update or Delete) operations, and one or many additional Slave servers that provide read-only operations. The Master uses standard, near-real-time database replication to each of the Slave servers. The Master/Slave model can speed overall performance to a point, allowing read-intensive processing to be offloaded to the Slave servers, but there are several limitations with this approach:

- The single Master server for writes is a clear limit to scalability, and can quickly create a bottleneck.
- Slaves are *eventually consistent*, meaning that the Slave servers are not guaranteed to have a current picture of the data that is in the Master (but they will later). While this is fine for some applications, if your applications always require an up-to-date view, this approach is unacceptable.

Case study master-slave: MongoDB



A *replica set* in MongoDB is a group of mongod processes that maintain the same data set. Replica sets provide redundancy and high availability, and are the basis for all production deployments. Mongod is the primary daemon process for the MongoDB system. It handles data requests, manages data access, and performs background management operations.

One mongod, the **primary**, receives all writes operations. A replica set can have only one primary. To support replication, the primary records all changes to its data sets in its log (MongoDB calls this the oplog). All other instances, **secondaries**, replicate the primary's oplog and apply the operations to their data sets such that the secondaries' data sets reflect the primary's data set.

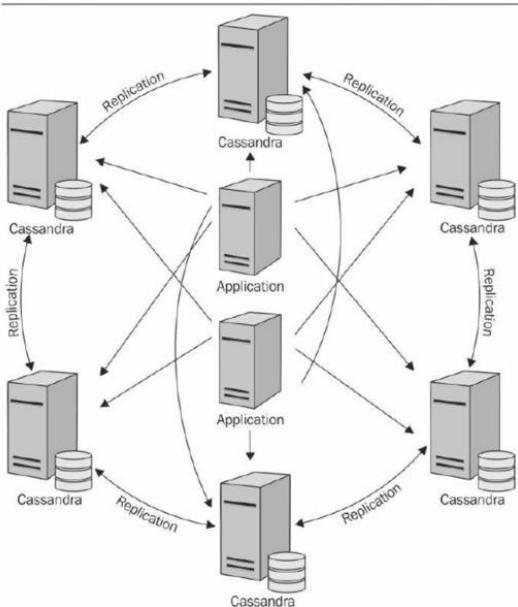
When a primary does not communicate with the other members of the set for more than 10 seconds, the replica set will attempt to select another member to become the new primary. The first secondary that receives a majority of the votes becomes primary.

By default, an application directs its read operations to the primary member in the replica set. Because write operations are issued to the primary, reading from the primary always returns the latest version of a document. Although it is possible to send read requests to secondary nodes, this is in general not recommended, because:

- All members of a replica have roughly equivalent write traffic; as a result, secondaries will service reads at roughly the same rate as the primary.
- Replication is asynchronous and there is some amount of delay between a successful write operation and its replication to secondaries; reading from a secondary can return out-of-date data.
- Distributing read operations to secondaries can compromise availability if *any* members of the set become unavailable because the remaining members of the set will need to be able to handle all application requests.
- Mongo runs a balancer in the background, who shifts chunks of data between nodes to balance node usage. For clusters with the balancer active, secondaries may return stale results with missing or duplicated data because of incomplete or terminated chunk migrations.

Sharding increases read and write capacity by distributing read and write operations across a group of machines, and is often a better strategy for adding capacity.

Case study ring: Cassandra



no special nodes

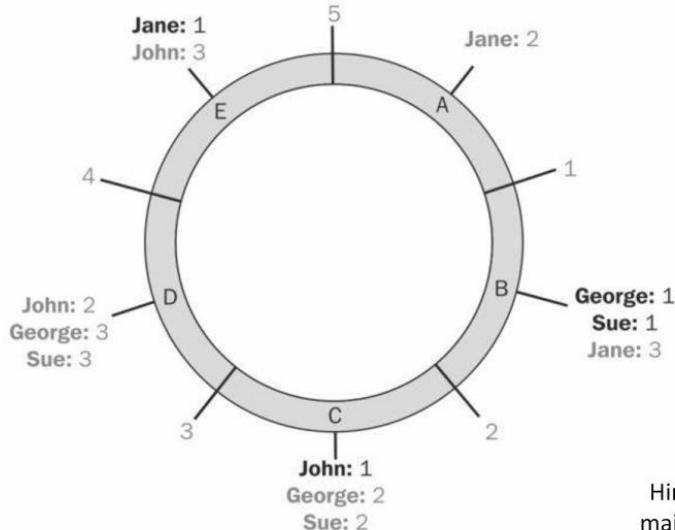
applications can query
any node

supports rack- and
datacenter-awareness

Unlike either monolithic or master-slave designs, Cassandra makes use of an entirely peer-to-peer architecture. All nodes in a Cassandra cluster can accept reads and writes, no matter where the data being written or requested actually belongs in the cluster. Internode communication takes place by means of a gossip protocol, which allows all nodes to quickly receive updates without the need for a master coordinator. Note that in contrast to the monolithic and master-slave architectures, there are no special nodes. In fact, all nodes are essentially identical, and as a result Cassandra has no single point of failure.

Cassandra employs a sophisticated replication system that allows fine-grained control over replica placement and consistency guarantees. There are two strategies available, one for single data center deployments and one when you have multiple datacenters.

Cassandra: replication in single data center



Hinted handoffs allow to maintain replication factor even when a replica node fails

We will only discuss the replication strategy in a single data center, assuming the default replication factor of three. As we will discuss later in more detail, data is assigned to the node in the cluster via a hash algorithm. Each node owns its own range of the hashed key space. The primary replica for each key is assigned to the node owning the hashed key value. Cassandra will then walk the ring in a clockwise direction to place each additional replica.

In the example on the slide, we have a cluster of 5 nodes (A-E), each covering some part of the hash keyspace. In the diagram, the keys in bold represent the primary replicas. Taking the hash of "Jane", we find that the first replica must be placed on node E. Then Cassandra takes the next two nodes (A and B) to place a replica.

Maintaining the replication factor when a node fails

You can specify a consistency level when you write to a node (see next slide). One key way in which Cassandra maintains fault tolerance even during node failure is through a mechanism called *hinted handoff*. If one of the replica nodes is unreachable during a write, then the system will store a hint on the coordinator node (the node that receives the write). This hint contains the data itself along with information about where it belongs in the cluster. Hints are replayed to the replica node once the coordinator learns via gossip that the replica node is back online.

Cassandra: tunable consistency

Consistency Level	Reads	Writes
ANY	Not supported	Data must be written to at least one node. Hinted handoffs tolerated.
ONE	Replica from closest node	Idem ANY; but no hinted handoffs tolerated
TWO	Replicas from two closest node	Idem ONE, but two replicas must be written
QUORUM	Replicas from a quorum will be compared and replica with latest timestamp is returned	Data must be written to a quorum of replica nodes.
ALL	Idem QUORUM, but for all nodes	Data must be written to all replica nodes.

Closely related to replication is the idea of consistency between replicas. Cassandra is often described as an eventually consistent system, but it is more accurate to describe it as having tunable consistency. The precise degree of consistency guarantee can be specified on a per-statement level. This gives the application developer strong control over the trade-offs between consistency, availability, and performance at call level – rather than forcing a one-size-fits-all strategy.

On every read and write operation, the caller must specify a consistency level, which lets Cassandra know what level of consistency to guarantee for that one call. For any operation, it is possible to achieve either strong consistency or eventual consistency. In the former case, we can know for certain that the copy of the data that Cassandra returns will be the latest. In the case of eventual consistency, the data returned may or may not be the latest, or there may be no data returned at all if the node is unaware of newly inserted data. Under eventual consistency, it is also possible to see deleted data if the node you're reading from has not yet received the delete request.

There are numerous combinations of read and write consistency levels, all with different consistency guarantees. To illustrate this point, let's assume that you would like to guarantee absolute consistency for all read operations. On the surface, it might seem as if you would have to read with a consistency level of ALL, thus sacrificing availability in the case of node failure. But there are alternatives depending on your use case.

There are actually two additional ways to achieve strong read consistency:

Write with consistency level of ALL:

This has the advantage of allowing the read operation to be performed using ONE, which lowers the latency for that operation. On the other hand, it means the write operation will result in `UnavailableException` if one of the replica nodes goes offline.

Read and write with QUORUM or LOCAL_QUORUM:

Since QUORUM requires a majority of nodes, using this level for both the write and the read will result in a full consistency guarantee, while still maintaining availability during a node failure.

You should carefully consider each use case to determine what guarantees you actually require. For example, there might be cases where a lost write is acceptable, or occasions where a read need not be absolutely current. At times, it might be sufficient to write with a level of QUORUM, then read with ONE to achieve maximum read performance, knowing you might occasionally and temporarily return stale data. Cassandra gives you this flexibility, but it's up to you to determine how to best employ it for your specific data requirements.

Balancing replication factor with consistency

Assume a cluster of **10 nodes**, with various replication factors (RF) and consistency levels (CL)

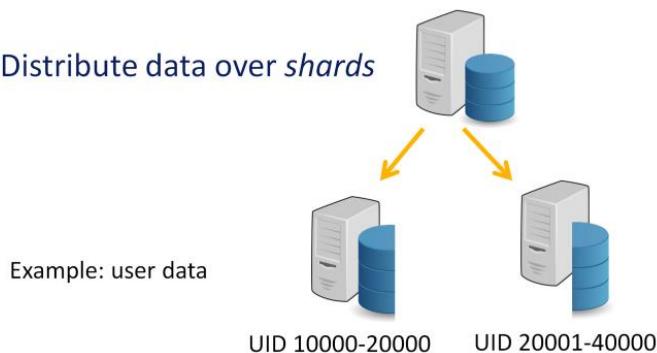
RF	Write CL	Read CL	Consistency	Availability
1	ONE QUORUM ALL	ONE QUORUM ALL	Consistent	No replica loss tolerated
2	ONE	ONE	Eventual	Tolerates loss of one replica
2	ONE	QUORUM ALL	Consistent	Tolerates loss of one replica on writes, but none on reads
3	ONE	ONE	Eventual	Tolerates loss of two replicas
3	ONE	QUORUM	Eventual	Tolerates loss of two replicas on write and one on reads

Achieving the desired availability, consistency, and performance targets requires coordinating your replication factor with your application's consistency level configurations.

In the table above, we consider a single data center cluster of 10 nodes and examine the impact of various configuration combinations on consistency and availability. We leave it as an exercise to the reader to study other combinations.

Scaling through sharding

Distribute data over *shards*



Shared-nothing

- no disk, memory, CPU sharing between nodes

Horizontal partitioning of data

Shard/Partition key

- determines on which shard to store a data unit

Database sharding provides a method for scalability across independent servers, each with their own CPU, memory and disk. The basic concept of database sharding is very straightforward: take a large database, and break it into a number of smaller databases across servers.

Sharding is a *horizontal* scaling strategy in which resources from each shard (or node) contribute to the overall capacity of the sharded database. Database shards are said to implement a *shared nothing* architecture that simply means that nodes do not share with other nodes; they do not share disk, memory, or other resources. Sharding thus supports high throughput and large data sets:

- Sharding reduces the number of operations each shard handles. Each shard processes fewer operations as the cluster grows. As a result, a cluster can increase capacity and throughput *horizontally*. For example, to insert data, the application only needs to access the shard responsible for that record.
- Sharding reduces the amount of data that each server needs to store. Each shard stores less data as the cluster grows. For example, if a database has a 1 terabyte data set, and there are 4 shards, then each shard might hold only 256GB of data. If there are 40 shards, then each shard might hold only 25GB of data.

A specific database column designated as the *shard key* determines which shard node stores any particular database row. The shard key is needed to access data. As a naïve but easily understood example, the shard key is the *username* column and the first

letter is used to determine the shard. Any usernames starting with A-J are in the first shard, and K-Z in the second shard. When your customer logs in with their username, you can immediately access their data because you have a valid shard key.

How to shard?

There is no single
secret sauce



- Some basic building blocks
- More about what *not* to do rather than a specific recipe
- Optimal scheme is highly application specific

It is important to note that Database Sharding is effective because it offers an application specific technique for massive scalability and performance improvements. The degree of effectiveness is directly related to how well the sharding algorithms themselves are tailored to the application problem at hand. There are numerous methods for deciding how to shard your data, and its important to understand your transaction rates, table volumes, key distribution, and other characteristics of your application.

There are multiple shard schemes possible, each designed to address a specific type of application problem. Each scheme has inherent performance and/or application characteristics and advantages when applied to a specific problem domain. In fact, using the wrong shard scheme can actually inhibit performance and the very results you are trying to obtain. It is also not uncommon for a single application to use more than one shard scheme, each applied to a specific portion of the application to achieve optimum.

General rule: avoid hotspots

Avoid uneven distribution of data
and/or operations across the shards



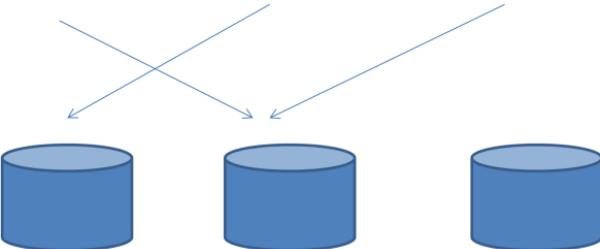
Examples of cross-shard operations: sorting, counting, joining

Although data tier frameworks may allow this kind of operations, the response time of such operations will dramatically impact your application

You should organize your shards according to data access patterns. So it is important to pick the right sharding key.

Data modelling

Hash([Alice]) Hash([Bob]) Hash([Charlie])



Two high-level (and conflicting) goals for your data model:

- Spread data evenly around the cluster by picking a good shard key
- Minimize the number of reads across multiple shards

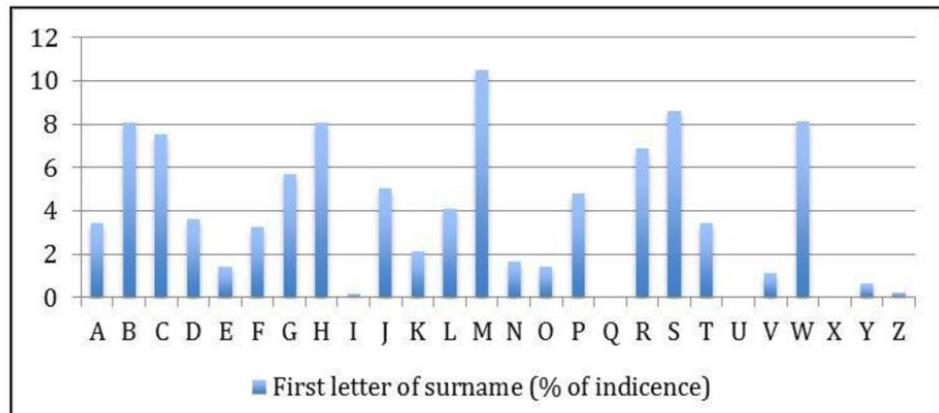
These are the two high-level goals for your data model:

- 1) Spread data evenly around the cluster
- 2) Minimize the number of partitions read

Data items (rows, documents...) are spread around the cluster based on a hash of the shard key. When you issue a read query, you want to read from as few partitions as possible, because each shard may reside on a different node. The node that receives a read request will generally need to issue separate commands to separate nodes for each partition you request. This adds a lot of overhead and increases the variation in latency. Furthermore, even on a single node, it's more expensive to read from multiple partitions than from a single one due to the way rows are stored.

These two goals often conflict and you need to balance these. We discuss this topic in more detail in the next slides.

Hotspots: uneven distribution of data

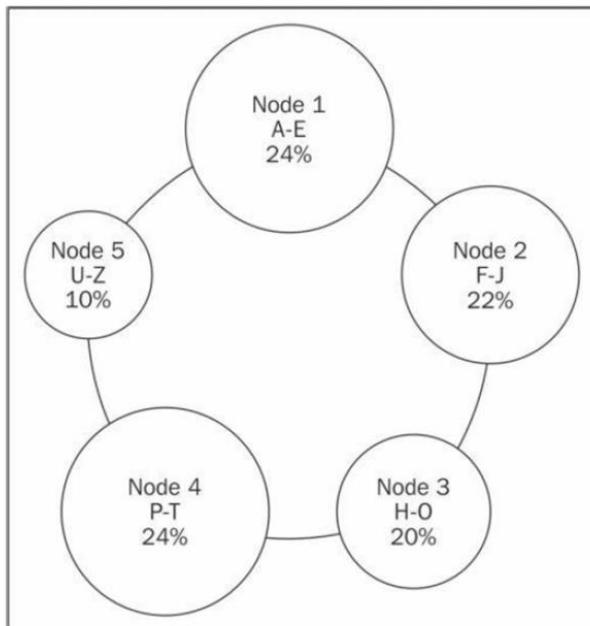


US Census Data, 2000

Let's assume, for example, that you're storing an address book, where the keys represent the last name of the contact. You use the last name of the contact to represent the keys. The graph above shows the distribution among the 26 letters, using 2000 United States Census Data.

As one would expect, last names in the United States are not evenly distributed by the first letter. In fact, the distribution is quite uneven. If we presume that each node owns a subset of the keys alphabetically, the result will resemble the diagram on the following slide.

Hotspots: uneven distribution of data



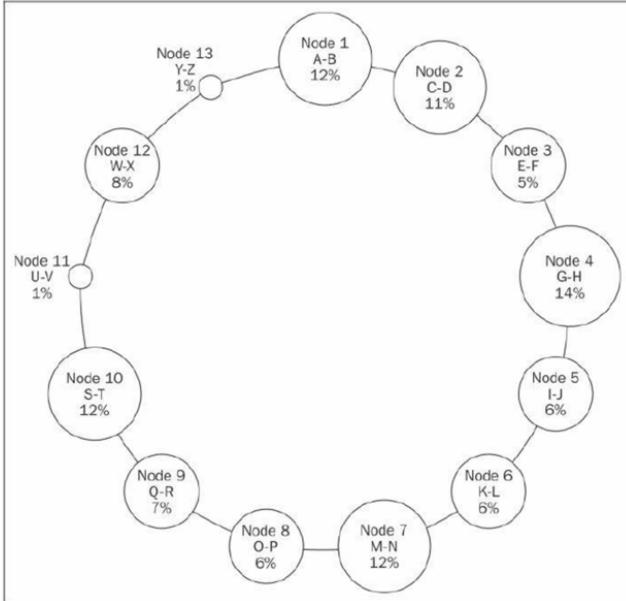
Imbalance on data placement

If queries follow popularity of data, then this also creates an impact on read/write operations

Using the last name as the shard key is likely to result in uneven distribution. Using the data from the previous slide, we see that we have created hotspots in node 1 and node 4, while node 5 is significantly underutilized.

One perhaps less obvious side effect of this imbalance is the impact on reads and writes. If we presume that both reads and writes follow the same distribution as the data itself (which is a logical assumption in this case), the heavier data nodes will also be required to handle more operations than the lighter data nodes.

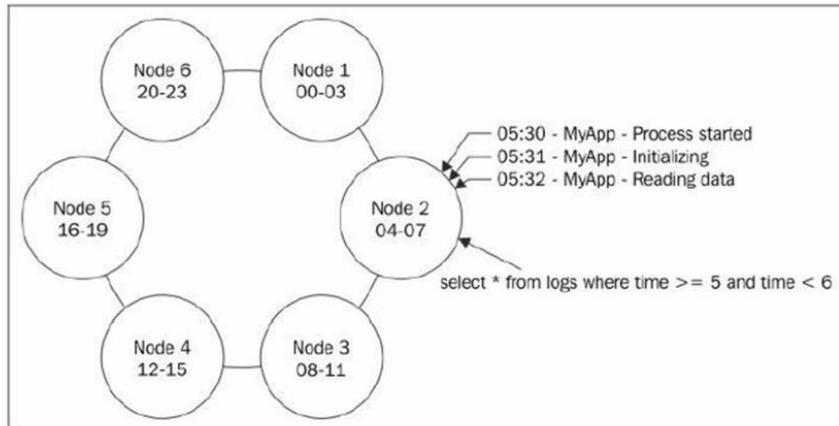
Hotspots: uneven distribution of data



Problem becomes worse in larger clusters

As is often the case in large systems, scaling out does not help to address this problem. In fact, the imbalance only gets worse when nodes are added. Still using the same data distribution from the previous example, the slide now shows the imbalance in a cluster of size 13. While in the five-node cluster, only one node was significantly underutilized, the larger cluster has eight out of 13 nodes doing half or less than half of the work as compared to the other nodes. In fact, two of the nodes own almost no data at all.

Hotspots: imbalanced queries



- Data placement is more balanced (assuming application is equally busy each hour)
- Read/write operations are focused on a single node

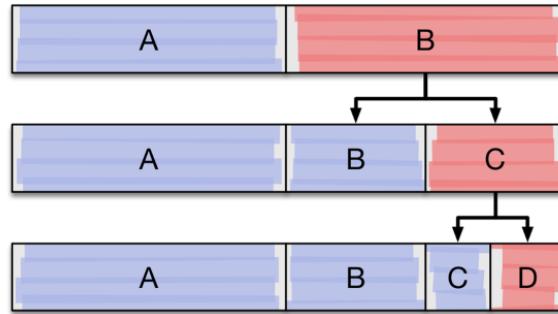
A common use case for big data storage systems is storing time-series data. Let's assume our use case involves writing log-style data, where we are always writing current timestamps and reading from relatively recent ranges of time.

Let's assume we have a six node cluster where the shard key corresponds to the time of day. If you are always writing current time, your writes will always go to a single node. Even worse, presuming you are reading recent ranges, your reads will also go to that same node.

Using the timestamp as shard key, time-series reads and writes will concentrate on a small subset of nodes. In the figure on the slide, node 2 is the only node doing any work. Each time the hour shifts, the workload will move to the next node in the ring. While the distribution of data in this model might be balanced (or it might not, depending on whether the application is busier at certain times), the workload will always experience hotspots.

Note: please do not conclude that using the timestamp as shard key is always a bad choice. This really depends on your query pattern. For example, if you calculate queries over the complete day, then data is equally retrieved from all nodes in the cluster.

Hotspots: imbalanced queries



Only last range will receive inserts!

This slide provides an alternative view on the uneven distribution of data write operations when using timestamps as shard key. Only the last shard will receive all write operations. Adding additional shards won't work, all requests keep going to the last shard.

A better approach is to use a combined partitioning (or shard) key. We will discuss this further in the Cassandra case study.

Bad sharding: “by application”

- Each service gets its own node
- Result:
 - Data distribution is non-uniform, massive hot spots
 - Every data access pattern is unique
 - Very little efficiency of scale



Each service might have different access patterns to the data it needs. If we put all these databases on different nodes, then we organize a non-uniform access pattern across the nodes. Some nodes might be overwhelmed, while others might have spare capacity.

Sharding approaches

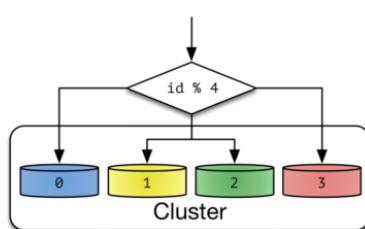
- NoSQL (key-value operations)
 - algorithmic sharding
 - dynamic sharding
 - case study: MongoDB
- SQL
 - entity grouping

This slide provides an overview of the different approaches to sharding we will study in the next slides.

NoSQL databases are often based on key-value operations. Mongo, Cassandra, RIAK all have some sort of key (or indexing). For this kind of databases, we have two typical approaches to sharding. In algorithmic sharding, the client can determine a given partition's database without any help. In dynamic sharding, a separate locator service tracks the partitions amongst the nodes.

Many (SQL) databases have more expressive querying and manipulation capabilities. They provide features such as joins, indexes and transactions that reduce complexity for an application. For this style of databases, we apply entity grouping: we store related entities in the same partition to provide additional capabilities within a single partition.

Algorithmic sharding



Sharding function maps partition key to node ID

Can be calculated by application itself (and provide as argument to query)

Data distribution does not consider payload size or space utilization

(over-simplistic) example:
`user_ID % NUM_NODES`

- Suitable for key-value databases with homogeneous values
- Resharding data is challenging:
 - Update sharding function (possibly in applications)
 - Move data around the cluster
- Example: memcached

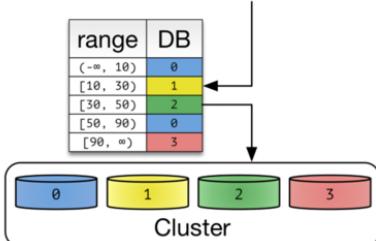
Algorithmically sharded databases use a sharding function (*partition_key*) -> *database_id* to locate data. A simple sharding function may be "*hash(key) % NUM_DB*".

Reads are performed within a single database as long as a partition key is given. Queries without a partition key require searching every database node. Non-partitioned queries do not scale with respect to the size of cluster, thus they are discouraged. Algorithmic sharding distributes data by its sharding function only. It doesn't consider the payload size or space utilization. To uniformly distribute data, each partition should be similarly sized. Fine grained partitions reduce hotspots — a single database will contain many partitions, and the sum of data between databases is statistically likely to be similar. For this reason, algorithmic sharding is suitable for key-value databases with homogeneous values.

Resharding data can be challenging. It requires updating the sharding function and moving data around the cluster. Doing both at the same time while maintaining consistency and availability is hard. Clever choice of sharding function can reduce the amount of transferred data. Consistent hashing (discussed later) is such an algorithm.

Examples of such system include Memcached. Memcached is not sharded on its own, but expects client libraries to distribute data within a cluster. Such logic is fairly easy to implement at the application level.

Dynamic sharding



External *locator service* determines location of entries

Locators can be created, split and reassigned to redistribute data

To read and write data, clients need to consult the locator service first

- More resilient to non-uniform distribution of data and load
- Locator service becomes single point of contention and failure
 - not simple to cache or replicate locators
- Auto-(re)sharding is possible, although challenging
- Used in many popular storage solutions
 - HDFS: Name Node
 - MongoDB: ConfigServer

In dynamic sharding, an external locator service determines the location of entries. It can be implemented in multiple ways. If the cardinality of partition keys is relatively low, the locator can be assigned per individual key. Otherwise, a single locator can address a range of partition keys.

To read and write data, clients need to consult the locator service first. Operation by primary key becomes fairly trivial. Other queries also become efficient depending on the structure of locators. In the example of range-based partition keys, range queries are efficient because the locator service reduces the number of candidate databases. Queries without a partition key will need to search all databases.

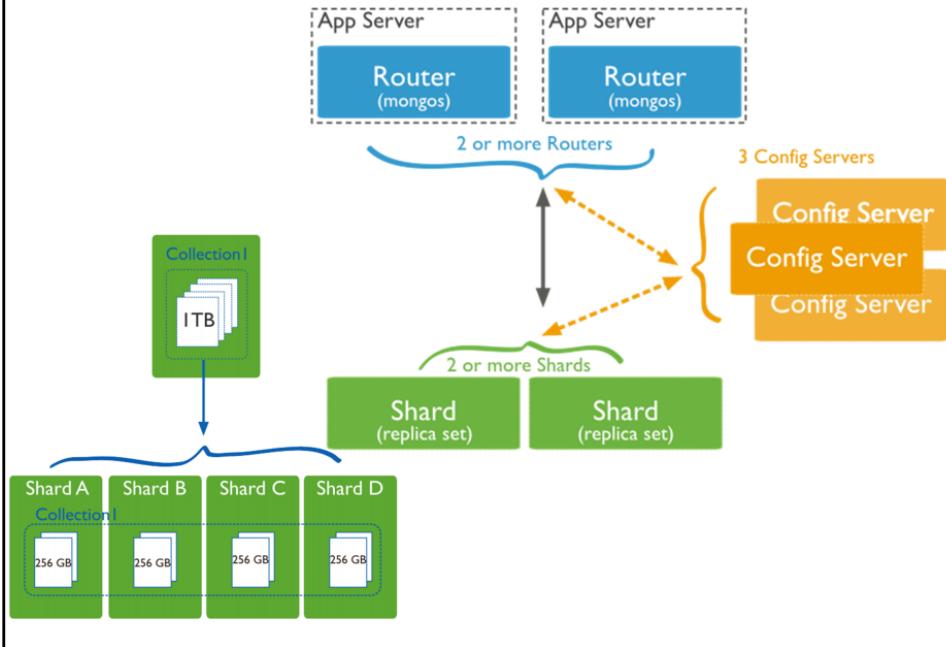
Dynamic sharding is more resilient to non-uniform distribution of data. Locators can be created, split, and reassigned to redistribute data. However, relocation of data and update of locators need to be done in unison. This process has many corner cases with a lot of interesting theoretical, operational, and implementational challenges.

The locator service becomes a single point of contention and failure. Every database operation needs to access it, thus performance and availability are a must. However, locators cannot be cached or replicated simply. Out of date locators will route operations to incorrect databases. Misrouted writes are especially bad — they become undiscoverable after the routing issue is resolved.

Since the effect of misrouted traffic is so devastating, many systems opt for a high consistency solution. Consensus algorithms and synchronous replications are used to store this data. Fortunately, locator data tends to be small, so computational costs associated with such a heavyweight solution tends to be low.

Due to its robustness, dynamic sharding is used in many popular databases. **HDFS** uses a NameNode to store filesystem metadata. In **MongoDB**, the ConfigServer stores the sharding information, and mongos performs the query routing. ConfigServer uses synchronous replication to ensure consistency. When a config server loses redundancy, it goes into read-only mode for safety. Normal database operations are unaffected, but shards cannot be created or moved. (We'll explain all these Mongo terms in more detail later).

Case study: MongoDB



Sharding is an important part of how MongoDB achieves its scalability. MongoDB supports sharding through the configuration of *sharded clusters*.

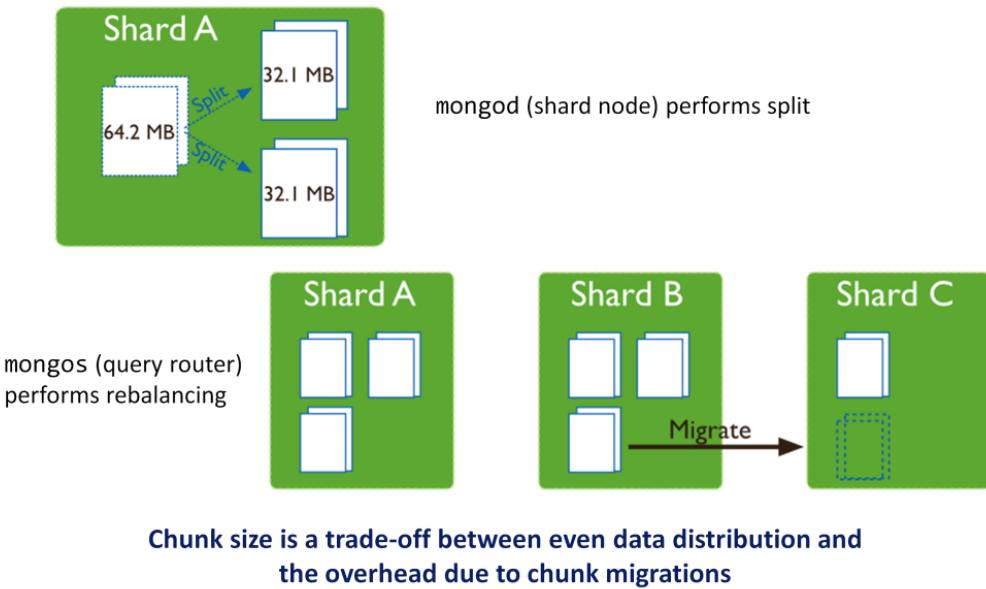
A sharded cluster has three component types: shards, query routers and config servers.

Shards store the data. To provide high availability and data consistency, in a production sharded cluster, each shard is a replica set.

Query Routers, or `mongos` instances (`mongos` = mongo server daemon), interface with client applications and direct operations to the appropriate shard or shards. The query router processes and targets operations to shards and then returns results to the clients. A sharded cluster can contain more than one query router to divide the client request load. A client sends requests to one query router. Most sharded clusters have many query routers.

Config servers store the cluster's metadata. This data contains a mapping of the cluster's data set to the shards. The query router uses this metadata to target operations to specific shards. Production sharded clusters have exactly 3 config servers.

Splitting and balancing



In MongoDB, data is split between shards using a shard key. The shard key is either an indexed field or an indexed compound field that exists in every document in the collection. MongoDB partitions data in the collection using ranges of shard key values. Each range, or **chunk**, defines a non-overlapping range of shard key values. MongoDB distributes the chunks, and their documents, evenly among the shards in the cluster.

The addition of new data or the addition of new servers can result in data distribution imbalances within the cluster, such as a particular shard containing significantly more chunks than another shard or a size of a chunk is significantly greater than other chunk sizes. MongoDB ensures a balanced cluster using two background processes: splitting and the balancer.

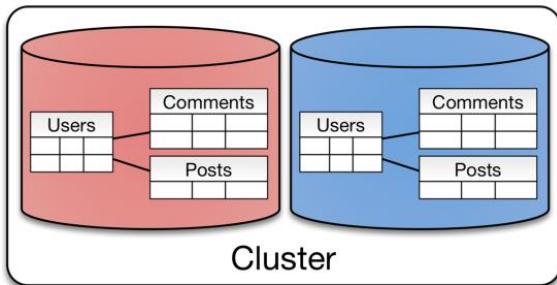
Splitting is a background process that keeps chunks from growing too large. When a chunk grows beyond a specified chunk size, MongoDB splits the chunk in half. Inserts and updates of data trigger splits. Splits are an efficient meta-data change. To create splits, MongoDB does not migrate any data or affect the shards.

The balancer is a background process that manages chunk migrations. The balancer can run from any of the query routers in a cluster. When the distribution of a sharded collection in a cluster is uneven, the balancer process migrates chunks from the shard that has the largest number of chunks to the shard with the least number of chunks.

until the collection balances. For example: if collection *users* has 100 chunks on shard 1 and 50 chunks on shard 2, the balancer will migrate chunks from shard 1 to shard 2 until the collection achieves balance. The shards manage chunk migrations as a background operation between an origin shard and a destination shard. During a chunk migration, the destination shard is sent all the current documents in the chunk from the origin shard. Next, the destination shard captures and applies all changes made to the data during the migration process. Finally, the metadata regarding the location of the chunk on config server is updated.

The default chunk size in MongoDB is 64 megabytes. Small chunks lead to a more even distribution of data at the expense of more frequent migrations. This creates expense at the query routing (*mongos*) layer. Large chunks lead to fewer migrations. This is more efficient both from the networking perspective *and* in terms of internal overhead at the query routing layer. But, these efficiencies come at the expense of a potentially more uneven distribution of data. For many deployments, it makes sense to avoid frequent and potentially spurious migrations at the expense of a slightly less evenly distributed data set.

Entity groups



Store related entities in the same partition

- Queries within a single physical shard are efficient
- Stronger consistency semantics within a shard
- Store data across multiple partitions to support efficient reads
 - E.g. chat messages between two users
- Replicate *reference data* to maintain shard autonomy
 - complete table is duplicated (<-> sharded)
 - e.g. list of ZIP codes
- Example: Google App Engine NDB Datastore

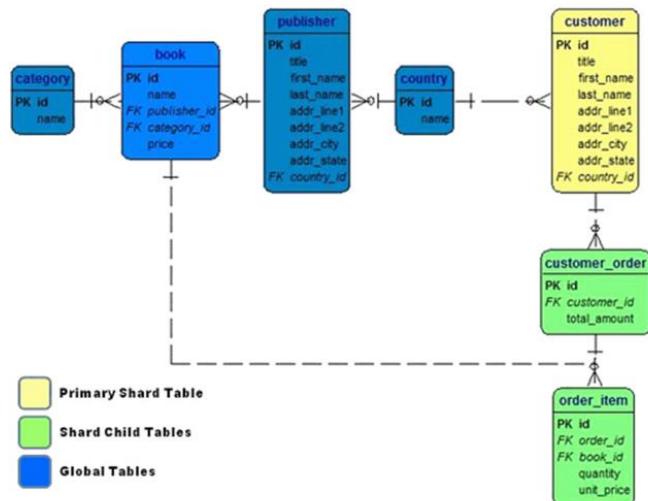
The previous discussion on algorithmic and dynamic sharding was geared towards key-value operations. But sharding can also be realized in relational databases through the concept of entity grouping.

The concept is very simple: store related entities in the same partition to provide additional capabilities within a single partition. These additional capabilities are:

- 1) Queries within a single physical shard are more efficient
- 2) Stronger consistency semantics can be achieved within a shard

Queries spanning multiple partitions typically have looser consistency guarantees than a single partition query. They also tend to be inefficient, so such queries should be done sparingly. However, a particular cross-partition query may be required frequently and efficiently. In this case, data needs to be stored in multiple partitions to support efficient reads. For example, chat messages between two users may be stored twice – partitioned by both senders and recipients. All messages sent or received by a given user are stored in a single partition. Another technique is the replication of so-called global tables: the relatively static lookup tables that are common utilized when joining to much larger primary tables. Tables containing values as status codes, countries, types, and even products fall into this category.

Example: bookstore



- Data shared by **customer.id** attribute
 - All related rows in two child tables are sharded as well
- Global Tables: common lookup tables, relatively low activity, replicated to all shards to avoid cross-shard joins

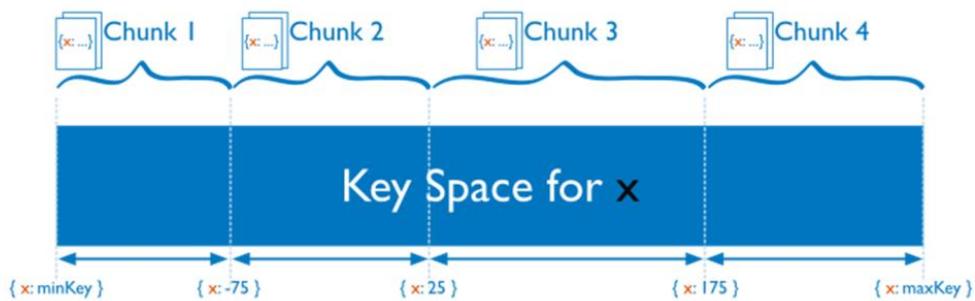
The example shows a simple schema for a Bookstore. The primary shard table that is used to shard the data is the ‘customer’ entity. The ‘customer’ table is the parent of the shard hierarchy, with the ‘customer_order’ and ‘order_item’ entities as child tables. The data is sharded by the ‘customer.id’ attribute, and all related rows in the child tables associated with a given ‘customer.id’ are sharded as well.

The global tables are the common lookup tables, which have relatively low activity, and these tables are replicated to all shards to avoid cross-shard joins. While this example is very basic, it does provide the basic considerations for determining how to shard a given database application.

Picking the right partition key

- Range-based partitioning
- Hash-based partitioning
 - Consistent hashing
- Combined keys
 - Case study: time-series in Cassandra

Range-based partitioning



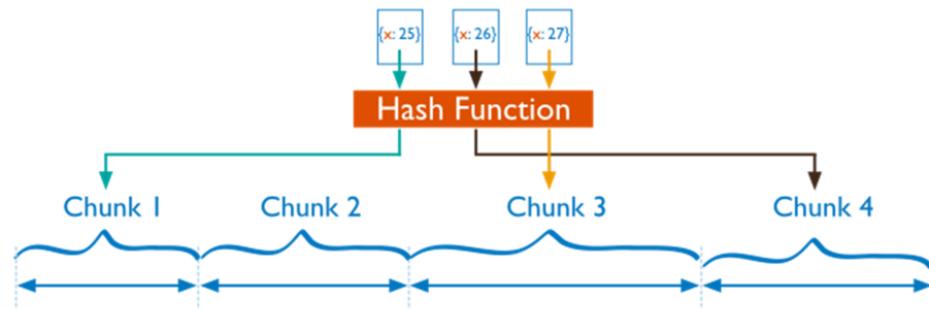
- “close” keys are likely to be on the same node
- ideal for range-based queries

For *range-based sharding*, you divide the data set into ranges determined by the shard key values to provide **range based partitioning**. Consider a numeric shard key: If you visualize a number line that goes from negative infinity to positive infinity, each value of the shard key falls at some point on that line. The key space is partitioned into smaller, non-overlapping ranges. In MongoDB, these ranges are called **chunks** where a chunk is range of values from some minimum value to some maximum value.

Given a range based partitioning system, data units with “close” shard key values are likely to be in the same chunk, and therefore on the same shard. Range based partitioning supports more efficient range queries. Given a range query on the shard key, the query router can easily determine which chunks overlap that range and route the query to only those shards that contain these chunks.

However, range based partitioning can result in an uneven distribution of data, which may negate some of the benefits of sharding. For example, if the shard key is a linearly increasing field, such as time, then all requests for a given time range will map to the same chunk, and thus the same shard. In this situation, a small set of shards may receive the majority of requests and the system would not scale very well.

Hash-based partitioning

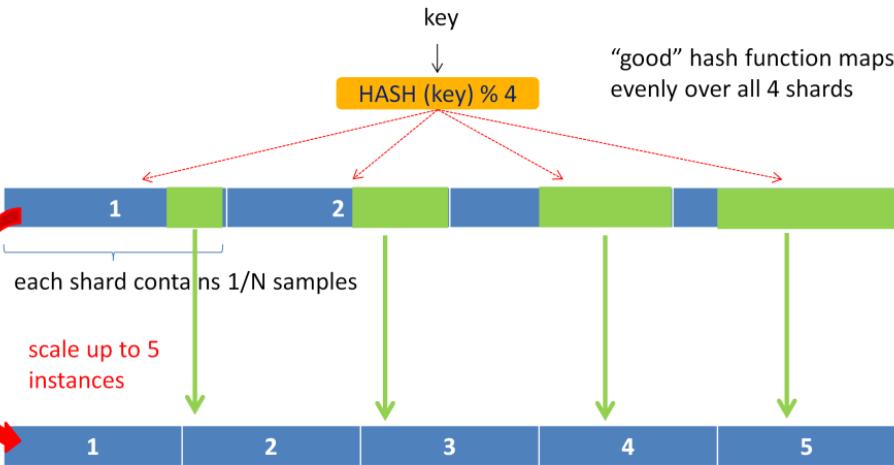


- more random spread of data
- but... range queries likely to touch many nodes

For *hash based partitioning*, you compute a hash of a field's value, and then use these hashes to create chunks. With hash based partitioning, two documents with "close" shard key values are *unlikely* to be part of the same chunk.

Hash based partitioning ensures a more random (and thus even) distribution of data at the expense of efficient range queries. Hashed key values results in random distribution of data across chunks and therefore shards. But random distribution makes it more likely that a range query on the shard key will not be able to target a few shards but would more likely query every shard in order to return a result.

Bad sharding: “fixed hashing”



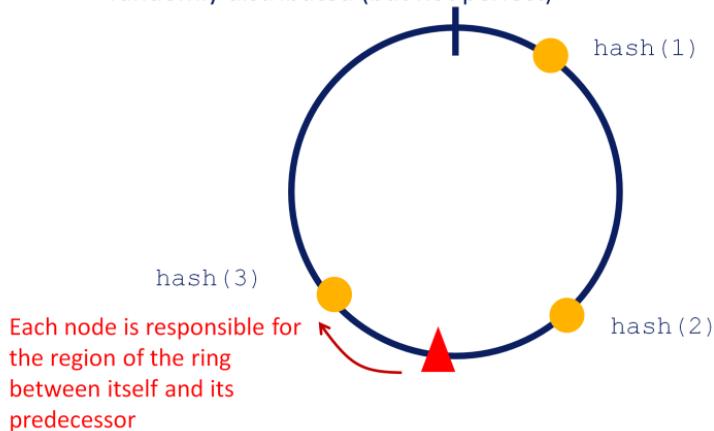
- recalculating all shard keys ($\text{mod } (N+1)$ instead of $\text{mod } (N)$)
- moving a lot of data to other shards

Sharding key schemes should also be robust to scaling up and/or node failures. Since good hashing functions generate output uniformly over their key space, modding this hash output over the number of nodes uniformly spreads the data units over all shards. Otherwise stated: the remainder of the division of the hash of the original key by the number of nodes is used to determine on which node the corresponding value is stored.

However, such a fixed hashing scheme performs badly when we scale horizontally, or when a node fails. Since the number of nodes has changed, we must recalculate for all data items the remainder of the division of the keys by the new number of nodes. For many data items, this means that they have to be shifted to another node in the cluster.

Remedy: Consistent hashing

- Allow to determine the location of an object in spite of constant shifting of nodes in and out of the cluster
- predefined range of hash keys
 - take hash of machine node number
 - randomly distributed (but not perfect)

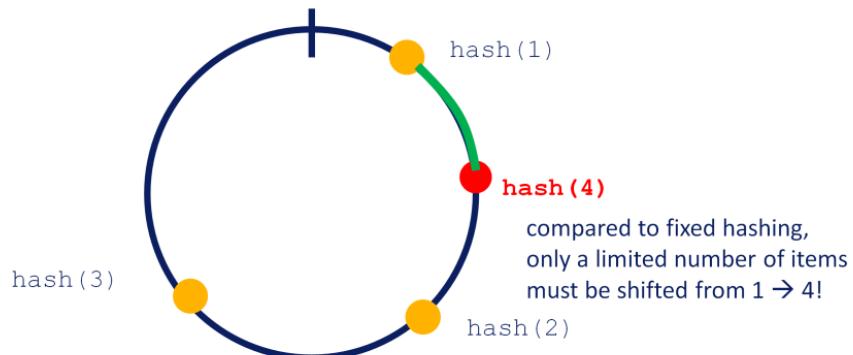


Consistent hashing provides a remedy for the issue presented on the previous slide.

The circle represents a predefined range of hash keys, organized in a ring. Keys are then hashed to produce a value that lies somewhere along the ring. If we take the hash of machine node numbers, than the resulting hashes will be more or less spread uniformly over the key space. We then state that each node is responsible for the region of the ring between itself and its predecessor. If we then calculate the hash of a data unit (the triangle on the slide), then we look at the first machine number hash that is strictly larger than the hash of the data unit to determine where to store the data unit.

(note: the order of hashes is not necessarily identical to the order of the elements hashed: $\text{hash}(N+1)$ is not necessarily larger than $\text{hash}(N)$)

Adding a node



Remaining problems:

- Irregular distribution of keys possible (better with higher number of nodes)
- Only items from a single node are shifted to the newly added node

If we add a fourth node, then only a limited number of items must be shifted compared to fixed hashing.

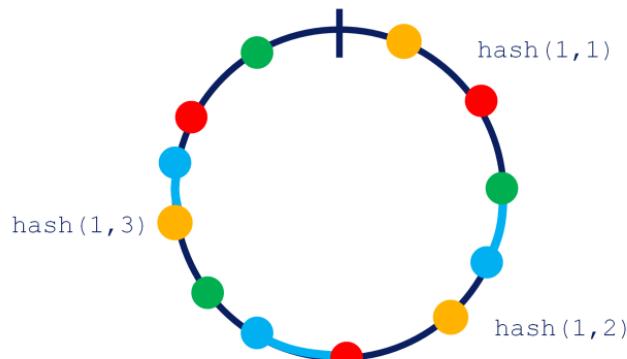
However, two problems remain:

- With a limited number of nodes, the distribution of the key range among the nodes might be irregular
- Only items of a single node are shifted to the newly added node, meaning that one node will be involved in an intensive process of copying data to the new node, while the others are left untouched

Adding a node

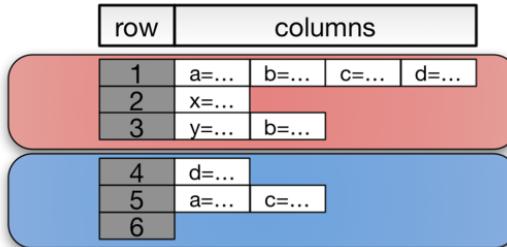
Map multiple “replicas” per machine

- More nodes on the circle = more uniform spreading
- Smaller number of nodes are redistributed to new machine, from almost each of the other machines in the cluster



A better solution is to map multiple replicas per machine. In the example above, each machine is mapped to three nodes on the “key circle”. If we add an additional node (e.g. the red dots), then the circle will be further divided. As you can see, items will be shifted from almost each of the other machines.

Combined keys in Cassandra



- Model column families around query patterns
 - what queries do I want to support?
- De-normalize and duplicate for read performance
 - no JOINs in Cassandra
 - +/- one table (column family) per query pattern
 - repeating data is not a shame...

Picking the right partition (or shard) key is not an easy task. Most often, using a single field in your data will not be sufficient. Instead, you will have to resort to combined (or compound) keys.

We will study this for the Cassandra database that was discussed earlier. To make the most out of Cassandra, you need to follow two rules:

- Cassandra is optimized for high write throughput, so write operations are relatively cheap. Read operations tend to be more expensive and are much more difficult to tune. So, if you can perform extra writes to improve the efficiency of your read queries, it's almost always a good tradeoff. For this, you must carefully study which queries are typically performed on your data.
- Disk space is generally the cheapest resource (compared to CPU, memory, disk IO or network), and Cassandra is architected around that fact. In order to get the most efficient reads, you often need to de-normalize and duplicate data. In the relational world, the pros of normalization are well understood: less data duplication, fewer data modification anomalies, conceptually cleaner, easier to maintain, and so on. The cons are also understood: that queries might perform slowly if many tables are joined, etc. The same holds true in Cassandra, but the cons are magnified since it's distributed and Cassandra supports no joins.

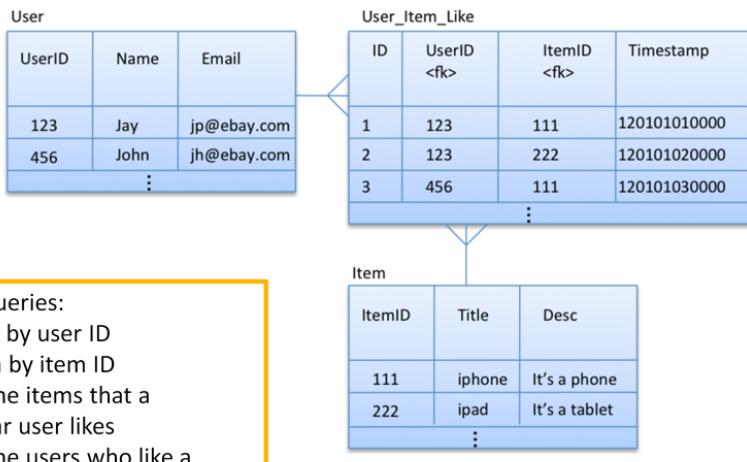
How to do this?

First, try to determine *exactly* what queries you need to support. This can include a lot of considerations that you may not think of at first. For example, think about grouping by an attribute, ordering by an attribute, filtering, enforcing uniqueness, etc... Changes to just one of these query requirements will frequently warrant a data model change for maximum efficiency.

Second, try to create a table (= column family) where you can satisfy your query by reading (roughly) one partition. In practice, this means you will use roughly one table per query pattern. If you need to support multiple query patterns, you usually need more than one table. To put this another way, each table should pre-build the “answer” to a high-level query that you need to support. If you need different types of answers, you usually need different tables. This is how you optimize for reads. Many of your tables may (/will) repeat the same data.

Example: “Like” relationship

Relational view:



The example concerns the functionality of an e-commerce system where users can like one or more items. One user can like multiple items and one item can be liked by multiple users, leading to a many-to-many relationship as shown in the relational model on the slide.

For this example, let's say we would like to query data as follows:

- *Get user by user id*
- *Get item by item id*
- *Get all the items that a particular user likes*
- *Get all the users who like a particular item*

Option 1: Exact replica of relational model

User			Item		
	Name	Email		Title	Desc
123	Jay	jp@ebay.com	111	iphone	It's a phone
	:			:	
User_Item_Like					
	UserID	ItemID			
1	123	111			
	:			:	



No easy way to query:
- all the items that a particular user likes
- all the users who like a particular item

This model supports querying user data by user id and item data by item id. But there is no easy way to query all the items that a particular user likes or all the users who like a particular item. This is the worst way of modeling for this use case. Basically, User_Item_Like is not modeled correctly here.

Note that the ‘timestamp’ column (storing when the user liked the item) is dropped from User_Item_Like for simplicity.

Option 2: normalized entities with custom indexes

User		Item	
123	Name	Email	
	Jay	jp@ebay.com	
	:		:
User_By_Item		Item_By_User	
111	123	456	
	null	null	...
	:		
123	111	222	
	null	null	...
	:		

Problem

many additional queries when we look up usernames who like a given item and vice versa

This model has fairly normalized entities, except that user id and item id mapping is stored twice, first by item id and second by user id. Here, we can easily query all the items that a particular user likes using `Item_By_User`, and all the users who like a particular item using `User_By_Item`. We refer to these column families as custom secondary indexes, but they're just other column families.

Let's say we always want to get the item title in addition to the item id when we query items liked by a particular user. In the current model, we first need to query `Item_By_User` to get all the item ids that a given user likes; and then for each item id, we need to query `Item` to get the title.

Similarly, let's say we always want to get all the usernames in addition to user ids when we query users who like a particular item. With the current model, we first need to query `User_By_Item` to get the ids for all users who like a given item; and then for each user id, we need to query `User` to get the username.

It's possible that one item is liked by a couple hundred users, or an active user has liked many items — which will cause many additional queries when we look up usernames who like a given item and vice versa. So, it's better to optimize by de-normalizing item title in `Item_by_User`, and username in `User_by_Item`, as shown in option 3 on the next slide.

Option 3: normalized entities with de-normalization into custom indexes

User		Item	
	Name		Desc
123	Jay	111	iphone
	jp@ebay.com		It's a phone
	:		:

User_By_Item		Item_By_User	
	123	456	
111	Jay	John	111
	222
	:	:	...

- Efficient querying of all item titles liked by a given user (and vice versa)
- What if we want all information (title, desc...)?
 - only needed when user asks for it (by clicking on a title)
- Also fairly efficient (only two read operations) for following query patterns:
 - For a given item id, get all item data and names of users who liked that item
 - For a given user id, get all user data along with item titles liked by that user

In this model, title and username are de-normalized in User_By_Item and Item_By_User respectively. This allows us to efficiently query all the item titles liked by a given user, and all the user names who like a given item. This is a fair amount of de-normalization for this use case.

What if we want to get all the information (title, desc, price, etc.) about the items liked by a given user? But we need to ask ourselves whether we really need this query, particularly for this use case. We can show all the item titles that a user likes and pull additional information only when the user asks for it (by clicking on a title). So, it's better not to do extreme de-normalization for this use case. (However, it's common to show both title and price up front. It's easy to do)

Let's consider the following two query patterns:

- For a given item id, get all of the item data (title, desc, etc.) along with the names of the users who liked that item.
- For a given user id, get all of the user data along with the item titles liked by that user.

These are reasonable queries for item detail and user detail pages in an application. Both will perform well with this model. Both will cause two lookups, one to query item data (or user data) and another to query user names (or item titles). As the user becomes more active (starts liking thousands of items?) or the item becomes hotter

(liked by a few million users?), the number of lookups will not grow; it will remain constant at two. That's not bad, and de-normalization may not yield as much benefit as we had when moving from option 2 to option 3.

Cassandra: clustering

```
CREATE TABLE playlists (
    id uuid,
    song_order int,
    song_id uuid,
    title text,
    album text,
    artist text,
    PRIMARY KEY (id, song_order )
);
```

Compound primary key consists of:

- Partition key → determines which node stores which row
- Additional columns → determine clustering, how is data sorted on disk

A compound primary key consists of the partition key and one or more additional columns that determine clustering. The partition key determines which node stores the data. It is responsible for data distribution across the nodes. The additional columns determine per-partition clustering. Clustering is a storage engine process that sorts data within the partition.

The data for each partition is clustered by the remaining column or columns of the primary key definition. On a physical node, when rows for a partition key are stored in order based on the clustering columns, retrieval of rows is very efficient.

Time series data modelling

Example: weather station creating temperature reading every minute

Option 1: add column every minute



Figure 1

```
CREATE TABLE temperature (
    weatherstation_id text,
    event_time timestamp,
    temperature text,
    PRIMARY KEY (weatherstation_id,event_time) );
```

1234ABCD	2015-04-03 07:01:00 11	2015-04-03 07:02:00 11	2015-04-03 07:03:00 12	2015-04-03 07:04:00 12	2015-04-03 07:05:00 12
----------	------------------------------	------------------------------	------------------------------	------------------------------	------------------------------

Cassandra's data model is an excellent fit for handling data in sequence regardless of datatype or size. When writing data to Cassandra, data is sorted and written sequentially to disk. When retrieving data by row key and then by range, you get a fast and efficient access pattern due to minimal disk seeks – time series data is an excellent fit for this type of pattern.

The simplest model for storing time series data is creating a wide row of data for each source. In this first example, we will use the weather station ID as the row key. The timestamp of the reading will be the column name and the temperature the column value. Since each column is dynamic, our row will grow as needed to accommodate the data. The event_time is used to sort the data on disk.

Time series data modelling

```
SELECT ( weatherstation_id, event_time, temperature
FROM temperature
WHERE wheaterstation_id='1234ABCD'
AND event_time >= '2015-04-03 07:01:00'
AND event_time <= '2015-04-03 07:05:00'
```

sequential read on disk!

1234ABCD	2015-04-03 07:01:00	2015-04-03 07:02:00	2015-04-03 07:03:00	2015-04-03 07:04:00	2015-04-03 07:05:00
	11	11	12	12	12

If we then query the temperature measured for a specific weatherstation in a specific period, then the query can be executed very efficiently:

- The node with the row containing the requested data is easily found by the partition key (all data of a single weather station is stored on a single node)
- We then only perform a sequential read on the disk, since Cassandra orders the columns by the timestamp value

Compound partition key

In the previous lay-out, you would ultimately hit the 2 billion column limit

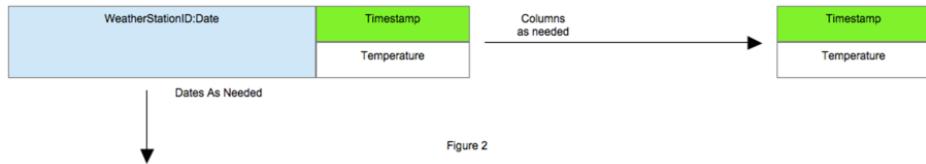


Figure 2

```
CREATE TABLE temperature_by_day (
    weatherstation_id text,
    date text,
    event_time timestamp,
    temperature text,
    PRIMARY KEY ((weatherstation_id,date),event_time) );
```

Compound partition key will group all weather data for a single day on a single row (and thus node)

In some cases, the amount of data gathered for a single device isn't practical to fit onto a single row. Cassandra can store up to 2 billion columns per row, but if we're storing data every millisecond you wouldn't even get a month's worth of data. The solution is to use a pattern called row partitioning by adding data to the row key to limit the amount of columns you get per device. Using data already available in the event, we can use the date portion of the timestamp and add that to the weather station id. This will give us a row per day, per weather station, and an easy way to find the data.

Note the `(weatherstation_id,date)` portion. When we do that in the PRIMARY KEY definition, the key will be compounded with the two elements. Now when we insert data, the key will group all weather data for a single day on a single row.

This comes however at the cost of having to perform reads from (possibly) two partitions if you want to read the data of multiple days of the same weather station.

Further reading

- D. McCreary, Making Sense of NoSQL: A guide for managers and the rest of us
- A. Fowler, NoSQL for Dummies
- R. Strickland, Cassandra High Availability

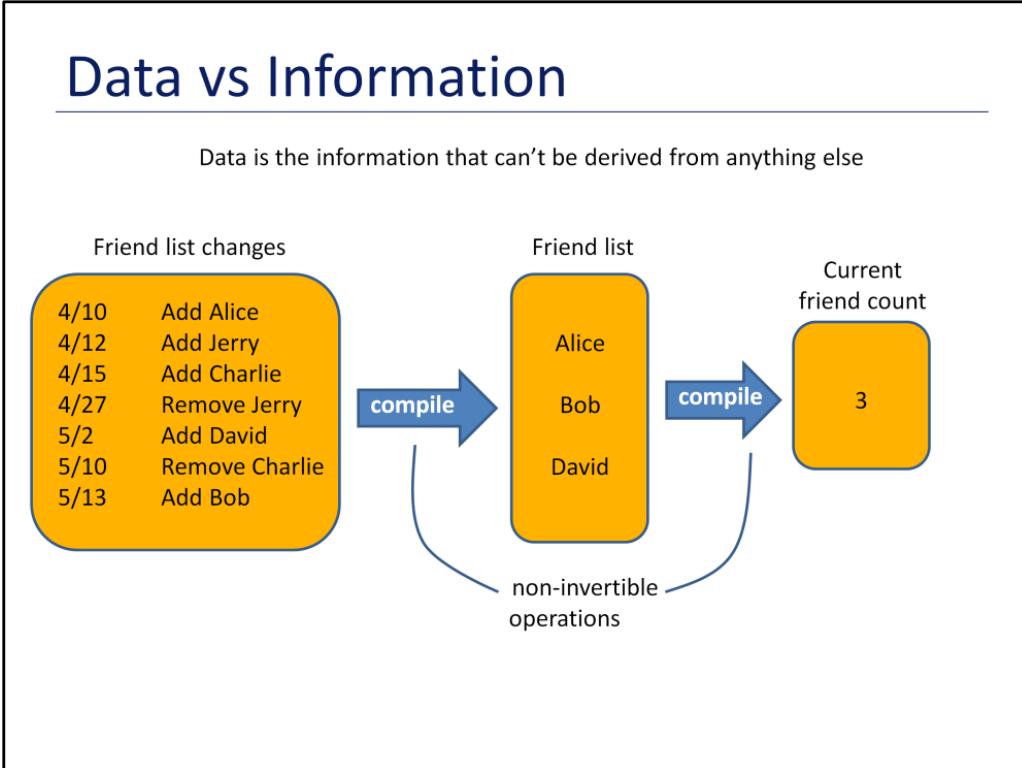
THE PROPERTIES OF DATA

Raw
Immutable
Perpetual



Data vs Information

Data is the information that can't be derived from anything else



The example illustrates information dependency. Each layer of information can be derived from the one to its left, but it's a one-way process. From the sequence of friend and unfriend events, you can determine the other quantities. But if you only have the number of friends, it is impossible to determine exactly who they are.

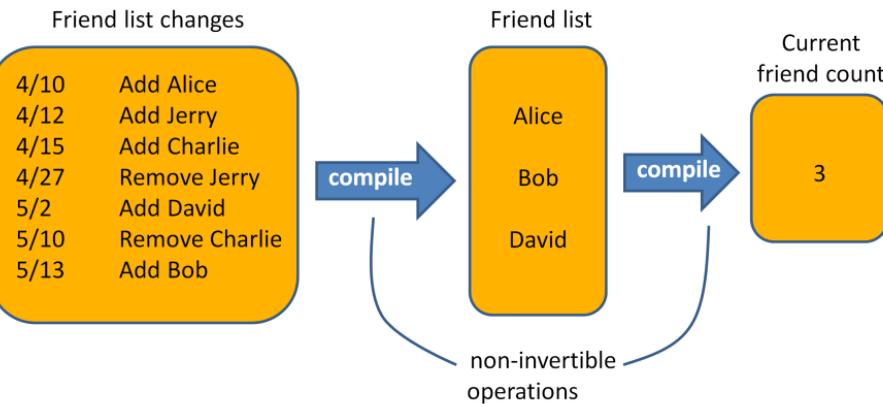
This shapes the following definitions:

- *Information* is the general collection of knowledge relevant to your big data system. It's synonymous with the colloquial usage of the word *data*.
- *Data* refers to the information that can't be derived from anything else. Data serves as the axioms from which everything else derives.
- *Queries* are questions you ask of your data. For example, you query your financial transaction history to determine your current bank account balance.
- *Views* are information that has been derived from your base data. They are built to assist with answering specific types of queries.

Data is *raw*



Query = function (all data)



- Answer as many questions as possible
- The rawer your data, the more questions you can ask of it
- You rarely know in advance all the questions

A big data system must be able to answer as many questions as possible. We'll colloquially call the property *rawness*. The rawer your data, the more questions you can ask of it. If you can, you want to store the rawest data you can get your hands on.

The example shows the data you might keep when designing a new social network. Each layer of information can be derived from the one to its left, but it's a one-way process. From the sequence of friend and unfriend events, you can determine the other quantities. But if you only have the number of friends, it is impossible to determine exactly who they are.

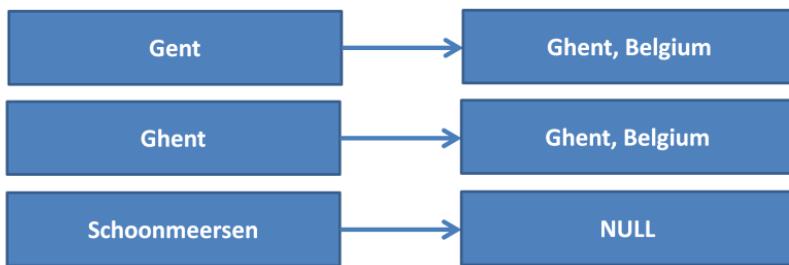
This shapes the following definitions:

- *Information* is the general collection of knowledge relevant to your big data system. It's synonymous with the colloquial usage of the word *data*.
- *Data* refers to the information that can't be derived from anything else. Data serves as the axioms from which everything else derives.

Storing raw data is hugely valuable because you rarely know in advance all the questions you want answered. By keeping the rawest data possible, you maximize your ability to obtain new insights, whereas summarizing, overwriting, or deleting information limits what your data can tell you. The trade-off is that rawer data typically entails more of it—sometimes much more. But Big Data technologies are designed to manage exabytes of data.

Guidelines: normalization?

Unstructured data is rawer than normalized data



- Semantic normalization may improve over time
- Only store normalized data if extraction is simple and accurate (e.g. age, sex)

Although the concept of *rawness* is straightforward, it is not always clear what information you should store as your raw data. A common hazy area is the line between *parsing* and *semantic normalization*. Semantic normalization is the process of reshaping free-form information into a structured form of data or pre-defined vocabulary (e.g. a dictionary of medical terms).

In our social network example, users may input anything for their location. A semantic normalization would try to match the input with a known place. We argue that it is better to store the unstructured string, because your semantic normalization algorithm may improve over time. If you store the unstructured string, you can renormalize that data at a later time when you have improved your algorithms. In the example, you may adapt the algorithm to recognize Schoonmeersen as a campus of Ghent University in Ghent.

As a rule of thumb, you should store the results of an algorithm for data extraction when that algorithm is simple and accurate, like extracting an age from an HTML page. If the algorithm is subject to change, due to improvements or broadening the requirements, store the unstructured form of the data.

Guidelines: more data = rawer data?



Bigger
is not always
better

- More data does not always equate to rawer data
- Example: URL of blog in social media profile
 - some HTML tags provide useful information on content
 - Javascript, CSS do not

It's easy to presume that more data equates to rawer data, but that's not always the case. Suppose that a social media user posts the URL of his new blog post on his profile. What exactly should you store? Storing the pure text of the blog entries is certainly a possibility. But any phrases in italics, boldface, or titles were deliberately emphasized by the user and could prove useful in text analysis. For example, you could use this additional information for an index to make your social network searchable. We'd thus argue that the annotated text entries (HTML tags) are a rawer form of data than the simple ASCII text strings. On the other hand, storing the color scheme, stylesheets and Javascript code of the blog website can't be used to derive any further information about the user. They serve only as the container for the contents of the site and shouldn't be part of the raw data. You will strip this information of the HTML page before storing it as raw data.

Data is *immutable*



Query = function (all data)

Friend list changes

4/10	Add Alice
4/12	Add Jerry
4/15	Add Charlie
4/27	Remove Jerry
5/2	Add David
5/10	Remove Charlie
5/13	Add Bob

Friend list

Alice
Bob
David

compile

Current friend count

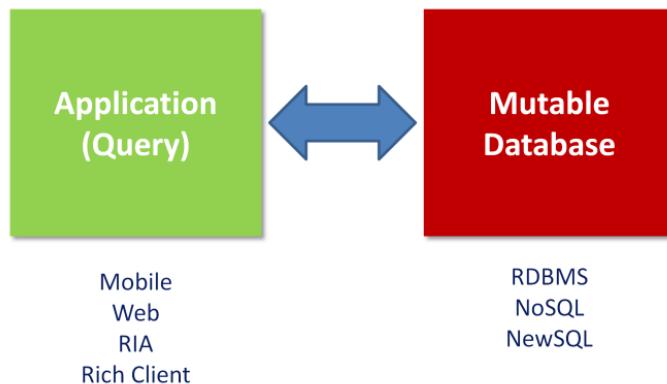
3

non-invertible operations

- Never delete anything...

You are unlikely to modify a data item already written to the storage. Instead; you will always append data.

Today's **incremental** architectures



Source of Truth is **mutable**: CRUD pattern

- Lack of Human Fault Tolerance
- Potential loss of information/data
- Difficult to achieve eventual consistency

What characterizes traditional application architectures is the use of read/write databases and maintaining the state in those databases *incrementally* as new data is seen. Applications continuously update items in the database. For example, an incremental approach to counting pageviews would be to process a new pageview by adding one to the counter for its URL. The problems with incremental architectures are significant: intolerance for human errors and the operational complexity as well as the constraints of the CAP theorem may lead to loss of information/data.

Note that this characterization of incremental architectures is a lot more fundamental than the discussion of relation vs. non-relational databases – the vast majority of both relation and non-relation databases deployments are done as fully incremental architectures.

Lack of human fault tolerance

- Bugs will be deployed to production over the lifetime of a data system
 - Operational mistakes will be made
 - Humans are part of the overall system
 - Just like hard disks, CPUs, memory, software
 - Design for human error like you design for any other fault
 - Examples of human error:
 - Deploy a bug that increments counters by two
 - Accidentally delete data from database
 - Accidental DOS on important internal service
- As long as an error does not lose or corrupt good data,
you can fix what went wrong.

Operational complexity

- Parts on the disk become unused as items are modified or deleted
- Compaction is the process of reclaiming space
- Compaction is expensive
 - lowers performance of machine while running
 - can even cause cascading failures

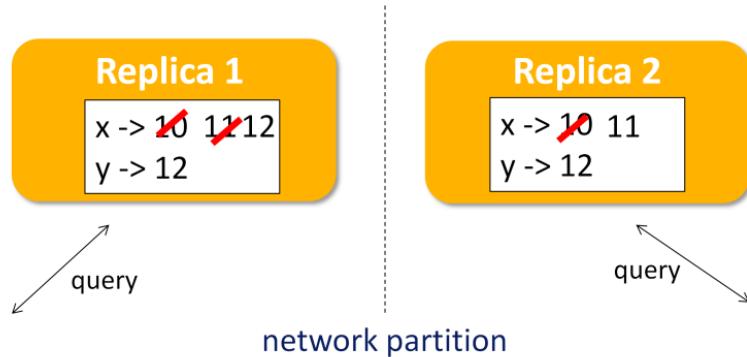
One of the difficulties in operating a production infrastructure hosting an incremental architecture is the need for read/write databases to perform online compaction. In a read/write database, as a disk index is incrementally added to and modified, parts of the index become unused. These unused parts take up space and eventually need to be reclaimed to prevent the disk from filling up. Reclaiming space as soon as it becomes unused is too expensive, so the space is occasionally reclaimed in bulk in a process called *compaction*.

Compaction is an intensive operation. The server places substantially higher demand on the CPU and disk during compaction, which dramatically lowers the performance of that machine during that period. The performance loss can even cause cascading failure – if too many machines compact at the same time, the load they were supporting will have to be handled by other machines in the cluster. This can potentially overload the rest of your cluster, causing total failure.

A competent operational staff can correctly schedule compactions on each node so that not too many nodes are affected at one. But this requires knowledge on how long a compaction takes (as well as the variance) and enough disk capacity on the nodes to last between compactions. A better approach is to avoid online compaction as much as possible.

Achieving eventual consistency

Example: highly available system
(CAP → no consistency in case of network partition)



Another complexity of incremental architectures emerges when trying to make systems highly available. A highly available system allows for queries and updates even in the presence of machine or partial network failure. As stated by the CAP theorem, it is impossible to realize consistency in a highly available system when a network partition occurs. So a highly available system sometimes returns stale results during a network partition.

In order for a highly available system to return to consistency once a network partition ends (*eventual consistency*), a lot of help is required from your application. Take for example the basic use case of maintaining a count in a database. The obvious way to go is to increment a number in the database whenever an event is received that requires the count to go up. To achieve high availability, distributed databases will keep multiple replicas. The information remains available even if a machine goes down or the network gets partitioned. During a network partition, a highly available system has clients update whatever replicas are reachable to them. This causes replicas to diverge and receive different sets of updates. Only when the partition goes away can the replicas be merged together into a common value. In the example above, the network partition is resolved when x has the value of 12 in replica 1 and a value of 11 in replica 2. What should the merged value be? Although the correct answer is 13 (since x had the value 10 when the network partition started), there is no way to know just by looking at the numbers 12 and 11. The replicas could have diverged at 11 (in this case the answer would be 12), or at 0 (in this case the answer

would be 23.

To do highly available counting correctly, you need a data structure that is amenable to merging when values diverge and need to implement repairing code. That is an amazing amount of complexity just to maintain a simple count.

Mutability vs immutability

Name	Location
Benteke	Birmingham
Denayer	Manchester
...	...

ID	Name	timestamp
1	Benteke	03-12-1990
2	Denayer	28-06-1995
...



Name	Location
Benteke	Liverpool
Denayer	Manchester
...	...

Update the current state of the world

ID	Location	timestamp
1	Birmingham	04-06-2012
2	Glasgow	17-07-2014
1	Liverpool	05-05-2015
2	Manchester	14-05-2015

Incrementally capture historical records of events (log)

In traditional database, you update existing records. In the example, we keep track of the location of football players. When a player updates his location because he is transferred to another team, in a mutable system we will update the location field in the table. The mutable system thus stores a current snapshot of the world.

In an immutable system, you create a separate record every time the data changes. Accomplishing this requires two changes. First, you track each field in a separate table. Second, you tie each unit of data to a moment in time when the information is known to be true.

By keeping each field in a separate table, you only record the information that changed. This requires less space for storage and guarantees that each record is new information and is not simply carried over from the last record.

Advantages of data immutability

- Human-fault tolerance
 - No data can be lost
 - Delete bad data units and recompute if necessary
- Simplicity
 - No data index needed
 - Only need to append new data units
- Trade-offs with data storage

Using an immutable schema for big data systems means that there are no updates or deletes of data. Instead, you only add more data, which gains you two vital advantages:

- **Human-fault tolerance** – People will make mistakes, and you must limit the impact of such mistakes and have mechanisms for recovering from them. With a mutable data model, a mistake can cause data to be lost, because values are actually overridden in the database. With an immutable data model, *no data can be lost*. If bad data is written, earlier (good) data units still exist. Fixing the data system is just a matter of deleting the bad data units and recomputing the views built from the master dataset.
- **Simplicity** – Mutable data models imply that the data must be indexed in some way so that specific data objects can be retrieved and updated. In contrast, with an immutable data model you only need the ability to append new data units to the master dataset. This doesn't require an index for your data, which is a huge simplification. Storing a master dataset can be as simple as storing a flat file.

One of the trade-offs of the immutable approach is that it uses more storage than a mutable schema. First, the user ID is specified for every property, rather than just once per row, as with a mutable approach. Additionally, the entire history of events is stored rather than just the current view of the world. You should take advantage of the ability to store large amounts of data using Big Data technologies to get the benefits of immutability. The importance of having a simple and strongly human-fault tolerant master dataset can't be overstated.

Data is perpetual



- Each piece of data is true in perpetuity
 - Once true, always true
- Tagging each piece of data with a timestamp is a practical way to achieve this

Deleting data is a statement about the *value* and not about truthfulness

- Garbage collection
- Regulations

The key consequence of data immutability is that each piece of data is true in perpetuity. A piece of data, once true, must always be true; Immutability wouldn't make sense without this property, and tagging each piece of data with a timestamp is a practical way to make data eternally true.

In general, the master dataset consistently grows by adding new immutable and eternally true pieces of data. There are some special cases in which you do delete data, but these cases are not incompatible with data being eternally true.

- Garbage collection - you delete all data units that have low value. You can use garbage collection to implement data retention policies that control the growth of the dataset. For example, you may decide to keep only one location per person per year instead of the full history of each time a user changes locations.
- Regulations – government regulations may require you to purge data from your databases under certain conditions

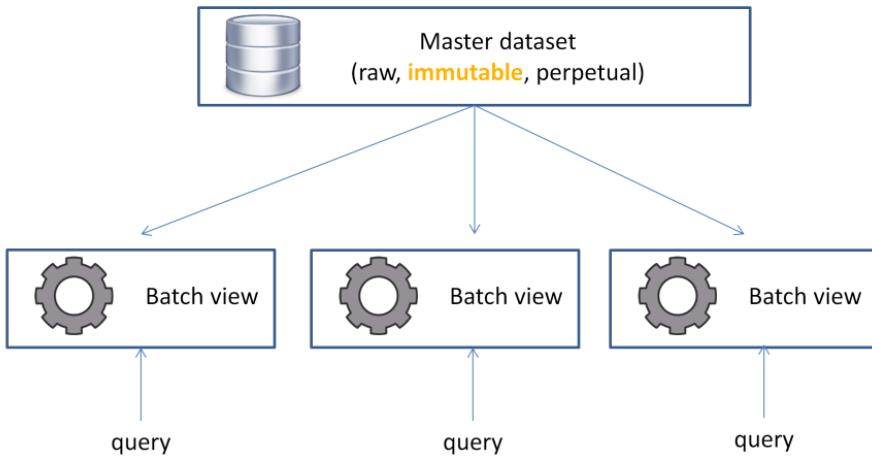
In both cases, deleting the data is not a statement about the truthfulness of the data. Instead, it is a statement about the value of the data. Although the data is eternally true, you may prefer to forget the information either because you must or because it doesn't provide enough value for the storage cost.

BATCH PROCESSING

Generating views on data

Batch view = function (all data)

Query = function (batch view)



Because of the many problems with today's incremental architectures, we will study another approach to cope with data. The foundation of the approach lies in the construction of the master dataset. This master dataset is structured along the RIP principles discussed in the previous section: it contains raw data that is no longer touched and kept forever.

Views are information that has been derived from the master dataset. They are pre-computed to assist with answering specific types of queries. The precomputed view is indexed so that it can be accessed with random reads. In this system, you run a function on all the data to get the view. Then, when you want to know the value for a query, you run a function on that view. In the master dataset, no data is stored redundantly. For efficiency reasons, the result of the batch views may contain duplicate data: one piece of data from the master dataset may get indexed into many batch views.

Of course, the master dataset is continuously expanding with new data. By the time the result of the view calculation is available, it is very likely that the result is already out-of-date. We will discuss how to cope with this later in this course. For now, we focus on *batch* processing: processing a snapshot of the master dataset.

Example

Master dataset

Friend list changes

4/10	Add Alice
4/12	Add Jerry
4/15	Add Charlie
4/27	Remove Jerry
5/2	Add David
5/10	Remove Charlie
5/13	Add Bob



Batch views



Queries

"Are Tom and Jerry friends?"

"How many friends does Tom have?"

- Batch views are indexed for fast querying
- Duplicate information between different batch views possible

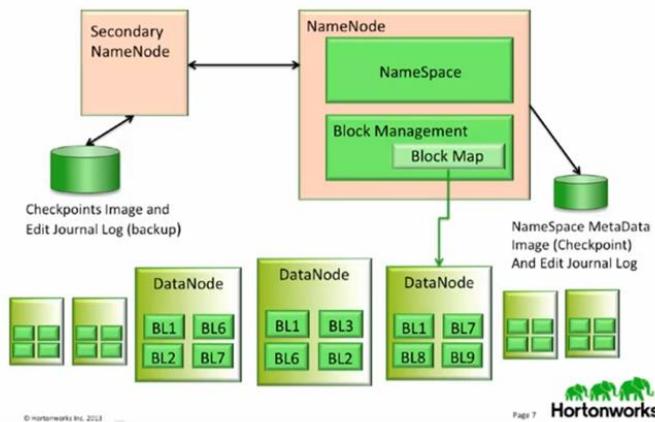
MapReduce

- Execution framework and programming model
- large-scale data processing
- builds on a distributed file system
- Apache Hadoop MapReduce is the most widely known and widely used open source implementation of the Google MapReduce paradigm
- Unlike traditional HPC clusters, Hadoop uses the same set of compute nodes for data storage as well as to perform the computations

Hadoop Distributed File System

recap

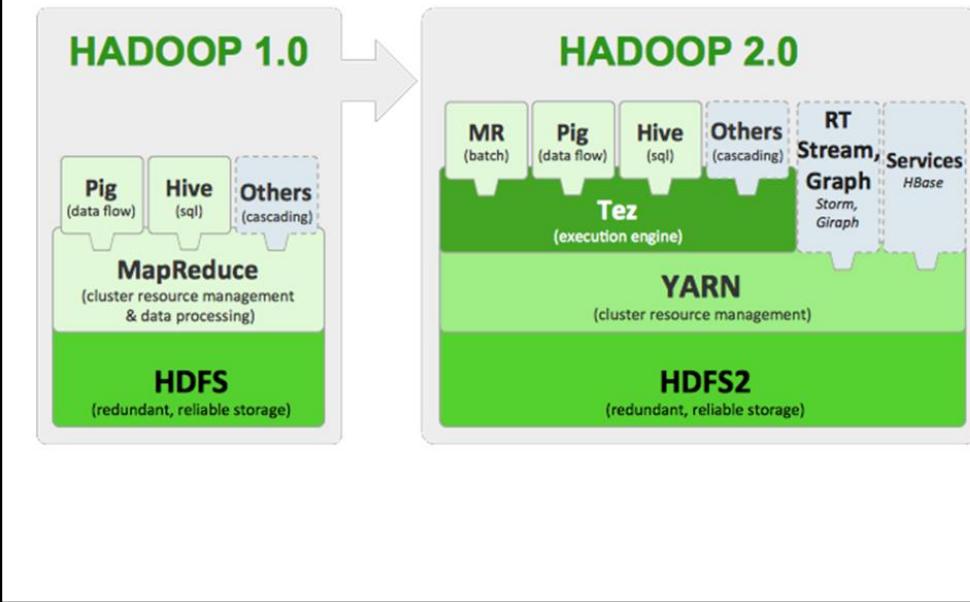
- High throughput parallel file system
- Massive amounts of data stored on commodity computers



The lowest layer of the MapReduce framework is the Hadoop Distributed File System. Since we discussed this earlier, we briefly recapitulate the concepts most relevant for batch processing.

HDFS consists of **NameNode** and **DataNode** services providing the basis for the distributed filesystem. NameNode stores, manages, and serves the metadata of the filesystem. NameNode does not store any real data blocks. DataNode is a per node service that manages the actual data block storage in the DataNodes. When retrieving data, client applications first contact the NameNode to get the list of locations the requested data resides in and then contact the DataNodes directly to retrieve the actual data.

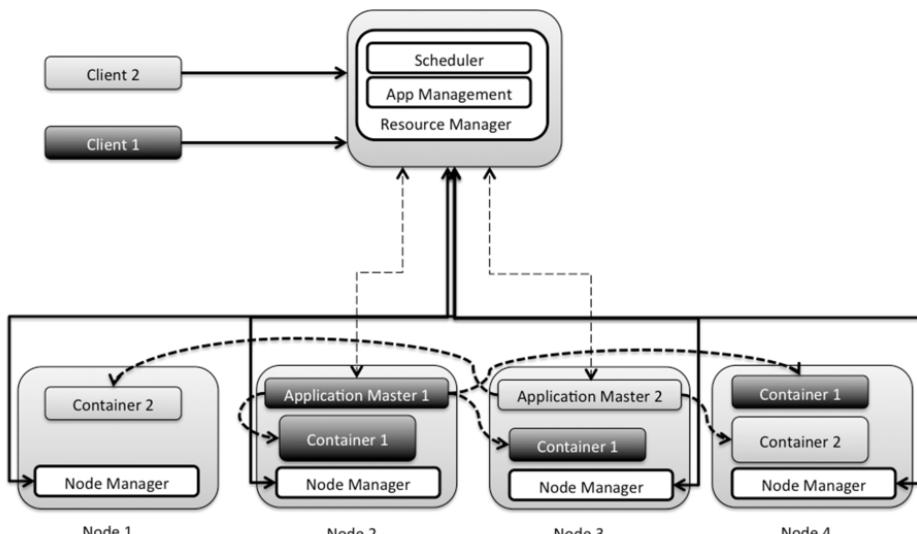
Yet Another Resource Negotiator



YARN (Yet Another Resource Negotiator) is a resource management system that allows multiple distributed processing frameworks to effectively share the compute resources of a Hadoop cluster and to utilize the data stored in HDFS. The primary goal of YARN is to separate concerns relating to resource management and application execution. By separating these functions, it provides a common platform for many different types of distributed applications. Users can thus utilize multiple distributed application frameworks side by side sharing a single cluster and the HDFS filesystem.

The batch processing based MapReduce framework was the only natively supported data processing framework in Hadoop v1. While MapReduce works well for analyzing large amounts of data, MapReduce by itself is not sufficient enough to support the growing number of other distributed processing use cases such as real-time data computations, graph computations, iterative computations, and real-time data queries.

YARN ApplicationMaster



This slide lays out the architecture of a YARN-based cluster. YARN has a concept called containers, which is the unit of resource allocation. Each allocated container has the rights to a certain amount of CPU and memory in a particular compute node. Applications can request resources from YARN by specifying the required number of containers and the CPU and memory required by each container.

YARN abstracts out resource management functions to a platform layer called **ResourceManager (RM)**. The RM is a per-cluster daemon that solely manages and allocates resources to the different applications (also known as jobs) submitted to the cluster. It has two main components: the Scheduler and the ApplicationsManager.

The Scheduler is responsible for allocating resources to the various applications that are running in the cluster. It does not have any insight into the status of the application: it does not guarantee restarts on application or hardware failures. It uses queues and capacity parameters during the allocation process. The ApplicationsManager is the component responsible for handling application submissions made by clients. It also bootstraps applications by negotiating the container on behalf of the application for the Application Master. The ApplicationsManager also provides the services of restarting the Application Master in case of failures.

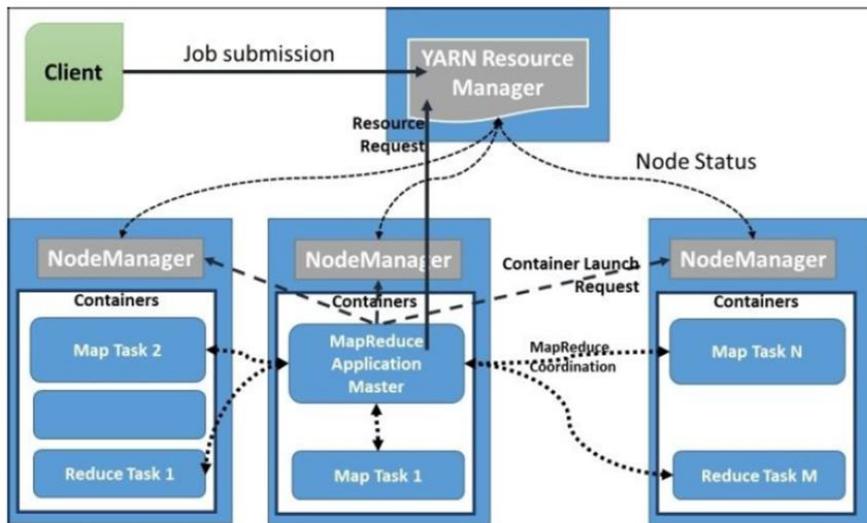
The **NodeManager (NM)** is a per-node daemon that does local container

management, ranging from authentication to resource monitoring (CPU, memory, disk health). The NM reports monitoring parameters to the RM. The RM scheduler can take decisions on container scheduling based on the load or health of the node.

The **ApplicationMaster (AM)** is a per-application process that coordinates the computations for a single application. The first step of executing a YARN application is to deploy the AM. After an application is submitted by a YARN client, the RM allocates a container and deploys the AM for that application. Once deployed, the AM is responsible for requesting and negotiating the necessary resource containers from the RM. Once the resources are allocated by the RM, AM coordinates with the NM to launch and monitor the application containers in the allocated resources.

Having separate AMs for each submitted application improves the scalability of the cluster as opposed to having a single process bottleneck to coordinate all the application instances.

YARN ApplicationMaster

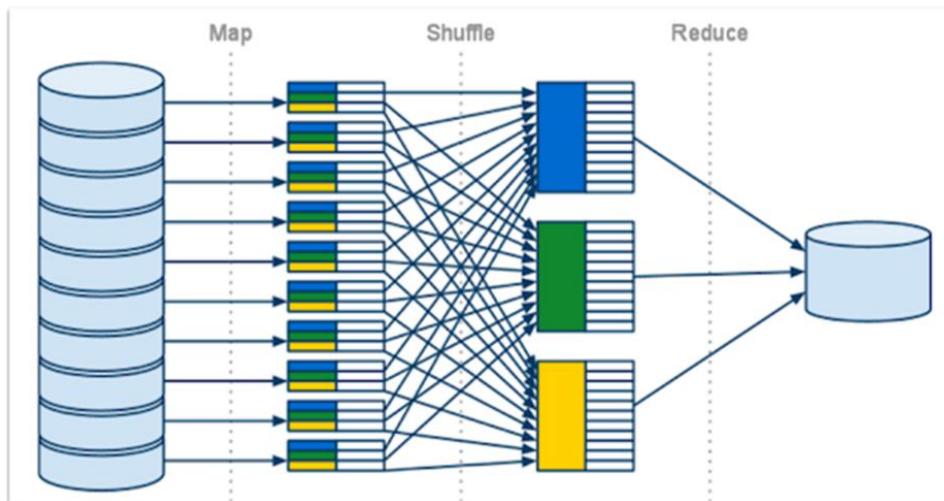


The diagram depicts the interactions between various YARN components, when a MapReduce application is submitted to the cluster. When the MapReduce ApplicationManager requests the Resource Manager for containers, it assigns a Map task to the container whose data is local or close to the allocated container node. The decision of which Map task is executed in the container thus happens *after* the AM receives the containers.

This concept is called late binding: the container spawned might not directly be related to the AM's request. The state at which the AM requests resources might change by the time the resource is allocated.

MapReduce programming model

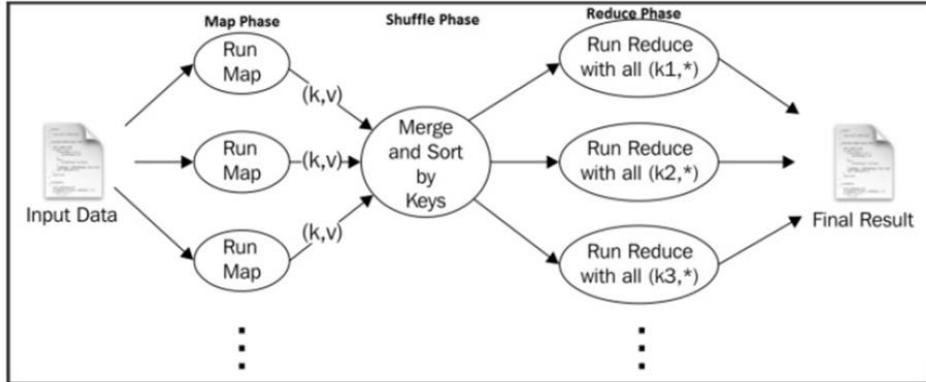
“Divide and conquer”



After having discussed the execution framework of MapReduce, we now turn our attention to the programming model.

The MapReduce programming model adopts a divide and conquer strategy to handle large data sets. First, the dataset is partitioned into smaller, independent data chunks to be processed in parallel (“map”). Then, the results from the previous step are combined, merged or otherwise aggregated (“reduce”).

MapReduce programming model



User specifies two functions:

- map: $(k_1, v_1) \rightarrow \text{list } [k_2, v_2]$
- reduce: $(k_2, \text{list } [v_2]) \rightarrow \text{list } [k_3, v_3]$

Sorting of intermediate keys between map and reduce phase.

The MapReduce programming model consists of Map and Reduce functions. The Map function receives each record of the input data (lines of a file, rows of a database, and so on) as key-value pairs and outputs key-value pairs as the result. By design, each Map function invocation is independent of each other allowing the framework to use divide and conquer to execute the computation in parallel. This also allows duplicate executions or re-executions of the Map tasks in case of failures or load imbalances without affecting the results of the computation. Typically, Hadoop creates a single Map task instance for each HDFS data block of the input data. The number of Map function invocations inside a Map task instance is equal to the number of data records in the input data block of the particular Map task instance.

Hadoop MapReduce groups the output key-value records of all the Map tasks of a computation by the **key** and distributes them to the Reduce tasks. This distribution and transmission of data to the Reduce tasks is called the Shuffle phase of the MapReduce computation. Input data to each Reduce task would also be sorted and grouped by the key. The Reduce function gets invoked for each key and the group of values of that key (*reduce <key, list_of_values>*) in the sorted order of the keys. In a typical MapReduce program, users only have to implement the Map and Reduce functions and Hadoop takes care of scheduling and executing them in parallel. Hadoop will rerun any failed tasks and also provide measures to mitigate any unbalanced computations.

Example: word count

Count how many times each word occurs in a given (huge) text

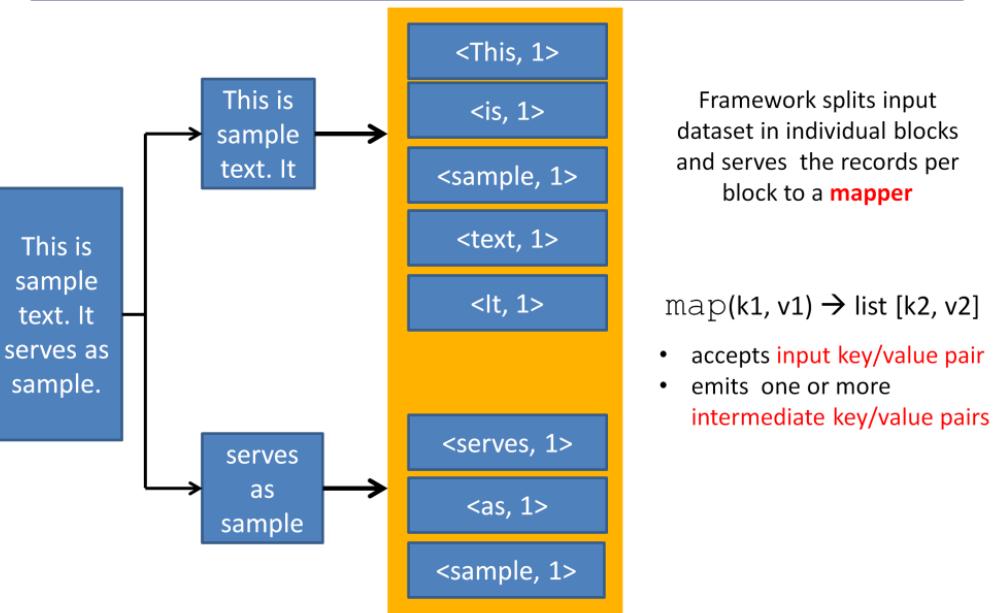
input

This is sample text
serving as a sample.

output

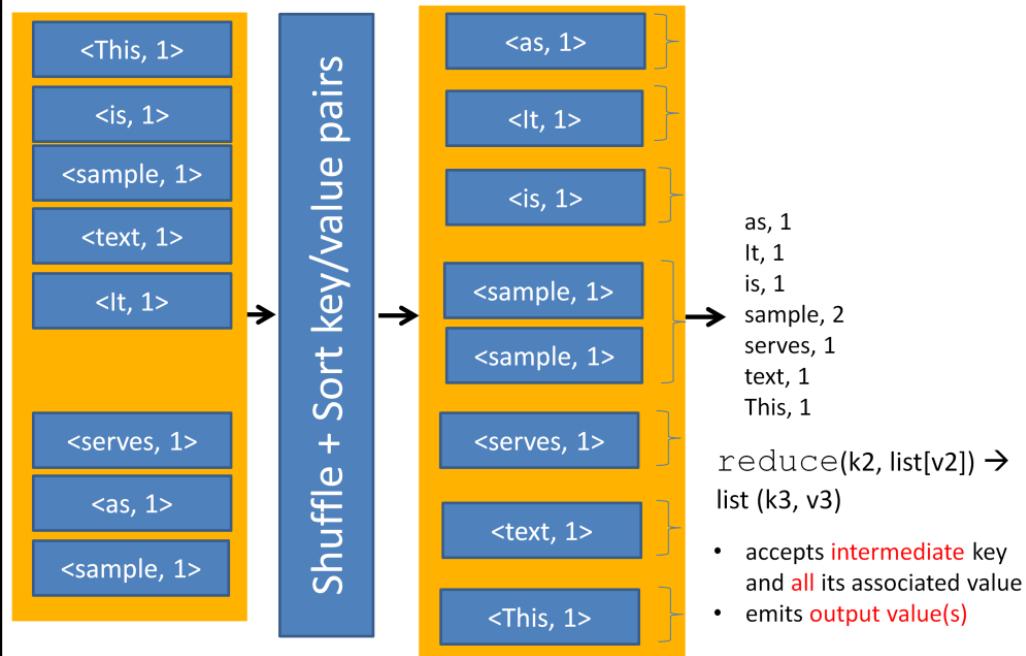
```
<this, 1>  
<is, 1>  
<sample, 2>  
...
```

Implementation: mapping phase



In this example, each word serves as key in the intermediate $\langle \text{key}, \text{value} \rangle$ pairs. The value is always 1 (unsigned integer).

Implementation: reduce phase



The intermediate pairs are shuffled and sorted across the different nodes of the cluster (based on the key value). Each reducer then handles a number of intermediate keys (and the associated value lists). In this simple example, the reducer simply sums all values in the list and emits this value.

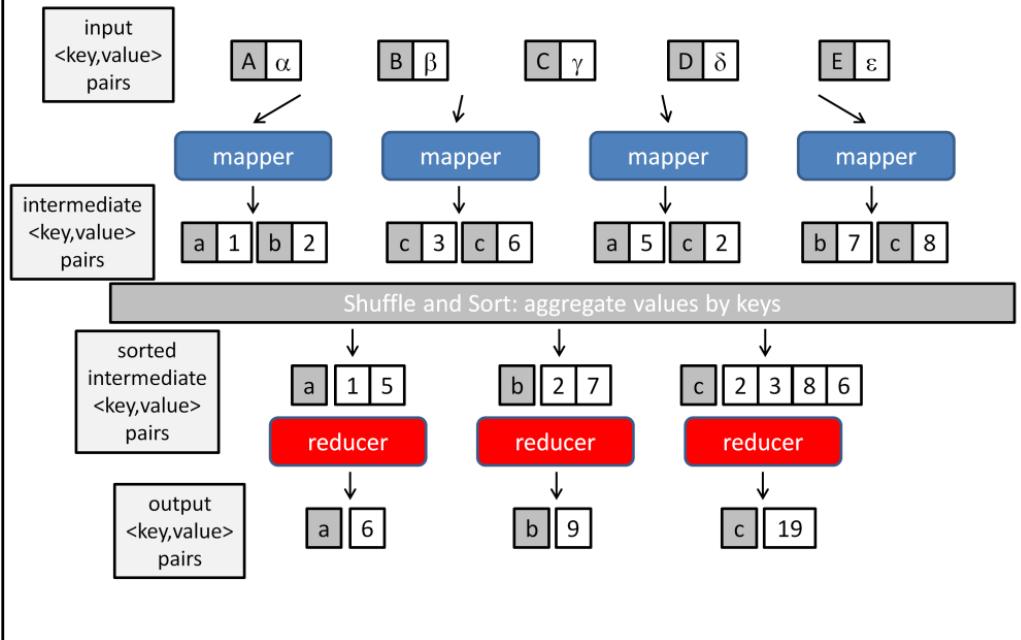
Note that in general the type of intermediate keys/values should not be the same as the type of output keys/values.

Pseudocode

```
function map(record) {
    for word in record:
        emit (word, 1)
}
```

```
function reduce(key, values) {
    emit (key, sum(values))
}
```

Mappers and reducers are executed **in parallel**



Mappers run on the worker nodes and execute map tasks by applying the map function to a part of the input data. Different map tasks can be executed in parallel. This means that there can be no dependencies between different map tasks. In other words, it should be possible to process part of the input data without having access to the remaining parts.

Reducers run on the worker nodes and execute reduce tasks by applying the reduce function on a single intermediate key and its corresponding values. Different reduce tasks can be executed in parallel. This means that there can be no dependencies between different reduce tasks. In other words, it should be possible to process a certain intermediate key and all of its corresponding values without having access to other intermediate keys and values.

MapReduce Design pattern

- User responsibilities
 - prepare the data
 - Implement mapper/reducer (+ combiner/partitioner – see further)
- All the rest is handled by the framework
- Possible extras by the user:
 - Complex, user defined data types in <key, value> pairs
 - User-specific initialization code in mapper/reducer
 - Ability to preserve state across multiple inputs
 - Determine the sort order of intermediate terms
 - Control the partitioning of the intermediate key space

MapReduce design patterns

- Complex algorithms
 - Require sequence or hierarchy of jobs
 - External (sequential) program (“driver”) can control and launch MapReduce jobs
- Design patterns help control the scalability
 - data X 2 → time X 2
 - # nodes X 2 → time/2

Pattern: local aggregation

```
function map(record) {
    for word in record:
        emit (word, 1)
}
```

Total number of emitted <k,v> pairs is equal to the **total number of words** in all documents

Local aggregation reduces the amount of intermediate <k, v> pairs:

- reduce amount of network traffic and disk I/O (pairs are streamed to disk)
- alleviates effect of “stranglers” (tasks that take excessively long)

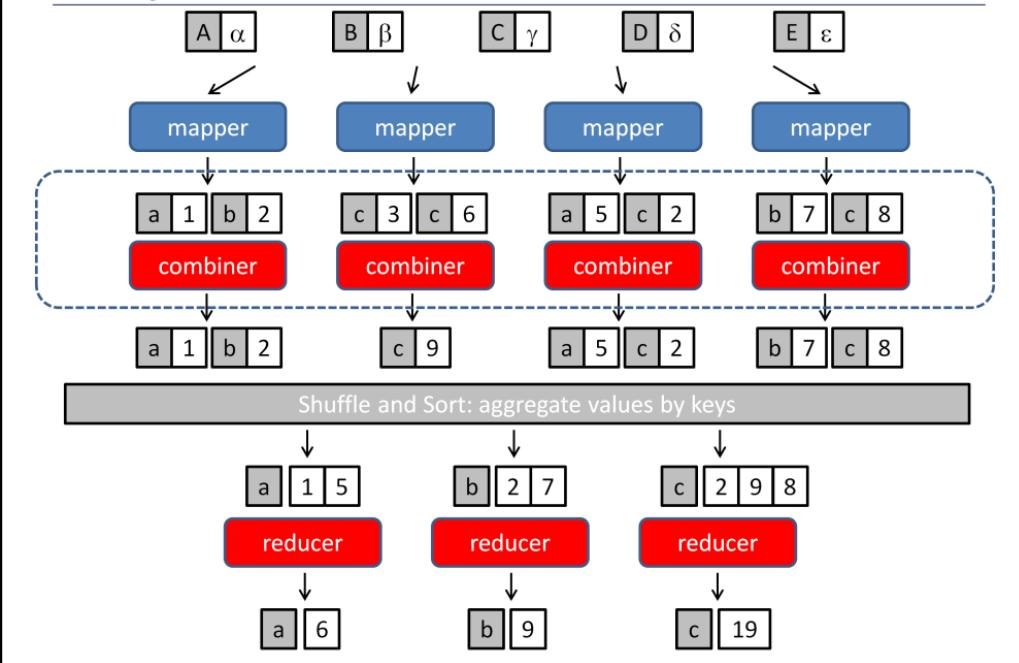
MapReduce combiners

- Optional method implemented by user
- Does not change the output of the program!

Even though this implementation is correct, it is not really efficient. The number of emitted <key, value> pairs equals the total number of words in a document (key = a word, value = always the number 1). All these data have to be written to disk and sorted, which can be time-consuming for extremely large datasets.

To optimize such scenarios, Hadoop supports a special function called **combiner**, which performs local aggregation of the Map task output key-value pairs. When provided, Hadoop calls the combiner function on the Map task outputs before persisting the data on the disk to shuffle the Reduce tasks. This can significantly reduce the amount of data shuffled from the Map tasks to the Reduce tasks. It should be noted that the combiner is an optional step of the MapReduce flow. Even when you provide a combiner implementation, Hadoop may decide to invoke it only for a subset of the Map output data or may decide to not invoke it at all.

MapReduce Combiners



Combiners do **not** reduce the number of intermediate $\langle \text{key}, \text{value} \rangle$ pairs emitted by the mappers. They do, however, reduce the number of intermediate $\langle \text{key}, \text{value} \rangle$ pairs that have to be **sorted**.

Pseudocode

```
function map(record){  
    for word in record:  
        emit (word, 1)  
}  
  
function combine(key, [int1, int2, int3]){  
    emit (key, sum([int1, int2, int3]))  
}  
  
function reduce(key, [intA, intB]){  
    emit (key, sum([intA, intB]))  
}
```

This slide shows pseudocode for a combiner in the word counting example. The combiner takes as input a key and a list of values, generated by one or more mappers. In this example, the key will correspond to a specific word. The combiner then emits a pair that has the same key, and the sum of its value list.

Note that the code still works correct when the combiner is not executed. This is an important check you have to do as programmer. The reasoning is quite simply: the combiner simply makes a partial sum of the number of word occurrences in a part of the dataset.

Combiners

- operates like a reducer but only on a subset of the key/values output
- input and output key-value data types must match the mapper output type
- called at the discretion of the framework
 - program must give same result with/without combiners
 - combiners may only be called for subset of key/values
- Only for commutative and associative operations
 - $a \times b = b \times a$
 - $a \times (b \times c) = (a \times b) \times c$

Local aggregation correctness

Compute the **mean of values** associated with the same key

```
function map(string t, integer r){  
    emit (t, r)  
}  
  
function reduce (string t, integer [r1, r2...]) {  
    s = sum(r1, r2 ...)  
    cnt = [r1, r2 ...].length()  
    emit (t, s/cnt)  
}
```

Cannot use **reducer** as **combiner**:

$\text{mean}(1,2,3,4,5) \neq \text{mean}(\text{mean}(1,2), \text{mean}(3,4,5))$

In this case, we cannot use a combiner: the mean of the numbers [1, 2, 3, 4, 5] is 3. If we would use combiners, which each calculate the mean on a subset of the data (say [1,2] and [3,4,5]), then we get a different and wrong results.

The reducer cannot distinguish between output coming from a mapper and output coming from a combiner.

What is wrong with this code?

```
function map(string t, integer r){  
    emit(t, r)  
}  
  
function combine(string t, integer [r1, r2...]) {  
    s = sum(r1, r2 ...)  
    cnt = [r1, r2 ...].length()  
    emit (t, pair(sum, cnt))  
}  
  
function reduce(string t, pairs [(s1,c1), (s2,c2)...]) {  
    sum = sum_1st([(s1,c1), (s2,c2)...]) //first element  
    cnt = sum_2nd([(s1,c1), (s2,c2)...]) //second element  
    r_avg = sum/cnt  
    emit(t, r_avg)  
}
```

To address the problem of mean calculation introduced on the previous slide using combiners, one could argue that we let combiners only calculate the partial sum and the count of how many numbers are in that sum. The combiner thus emits a value that consists of a pair. The elements of this pair are the partial sum, and the number of elements summed. This way, the reducer would be able to correctly calculate the mean: it sums the partial sums, and it knows how many numbers there were in total in the dataset.

The code above is however incorrect in one aspect...

Must match with output of mapper!

```
function map(string t){
    emit(t, pair(r, 1))
}

function combine(string t, pairs [(s1,c1), (s2,c2)...]){
    s = sum(r1, r2 ...)
    cnt = [r1, r2 ...].length()
    emit (t, pair(sum, cnt))
}

function reduce(string t, pairs [(s1,c1), (s2,c2)...]){
    sum = sum_1st([(s1,c1), (s2,c2)...]) //first element
    cnt = sum_2nd([(s1,c1), (s2,c2)...]) //second element
    r_avg = sum/cnt
    emit(t, r_avg)
}
```

The code must keep giving the correct result, even if the combiner is only applied to a subset of the key-value pairs generated by the mapper! Otherwise stated, the reducer cannot know which of the `<key, value>` pairs it receives are coming from a mapper and which ones are coming from a reducer.

This also means that the output of the mapper and the output of the combiner must have the same format. Hence, we must adjust the mapper to also emit a value pair. The count element in the emitted pair is always set to 1.

Pattern: constructing complex keys/values

- package data in more complex key/value types
 - e.g. pair, array
- example: co-occurrence matrix $[m_{ij}]$
 - m_{ij} = # of co-occurrences of w_i and w_j
 - space requirement $O(n^2)$
 - many use cases
 - to identify correlated product purchases
 - to identify words likely to occur nearby a specific word
- two patterns
 - “pairs”: emit $(\langle w_i, w_j \rangle, 1)$
 - “stripes”: emit $(w_i, H(w_j))$ (H = associative array)

One common approach for synchronization in MapReduce is to construct complex keys and values in such a way that data necessary for a computation are naturally brought together by the execution framework. We first touched on this technique in the previous slides, in the context of “packaging” partial sums and counts in a complex value (i.e., pair) that is passed from mapper to combiner to reducer.

In the next slides, we will study another common MapReduce example: the calculation of co-occurrences between item i and item j. Examples are e.g. to identify correlated product purchases (how many times are item i and item j bought together), or to identify common word patterns (how many times is word i found next to word j in all documents in the dataset)?

We will study two approaches to solve this pattern. Viewed abstractly, the pairs and stripes algorithms represent two different approaches to counting co-occurring events from a large number of observations. This general description captures the gist of many algorithms in fields as diverse as text processing, data mining, and bioinformatics. For this reason, these two design patterns are broadly useful and frequently observed in a variety of applications.

Pairs approach

```
function map(docid a, doc d){  
    for all word in d do  
        for all u in neighbours(word, d) do  
            emit(pair(word, u), 1)  
        done  
    done
```

```
function reduce(pair p, integer [r1, r2...]) {  
    emit(p, sum([r1, r2]))  
}
```

- each emitted pair by the reduce corresponds to a **cell** of the co-occurrence matrix
- generates enormous amount of key-value pairs (even after combining)
- no memory issues even for very large datasets (only pairs are emitted)

Pseudo-code for the “pairs” approach is shown here. Document ids and the corresponding document content make up the input key-value pairs. The mapper processes each input document and emits intermediate key-value pairs with each co-occurring word pair as the key and the integer one (i.e., the count) as the value. This is straightforwardly accomplished by two nested loops: the outer loop iterates over all words (the left element in the pair), and the inner loop iterates over all neighbors of the first word (the right element in the pair). The neighbors of a word can either be defined in terms of a sliding window (e.g. no further away than N words) or some other contextual unit such as a sentence.

The MapReduce execution framework guarantees that all values associated with the same key are brought together in the reducer. Thus, in this case the reducer simply sums up all the values associated with the same co-occurring word pair to arrive at the absolute count observed in the dataset, which is then emitted as the final key-value pair. Each pair corresponds to a cell in the word co-occurrence matrix. This algorithm illustrates the use of complex keys in order to coordinate distributed computations.

Stripes approach

```
function map(docid a, doc d){  
    for all word in d do  
        H = new AssociativeArray()  
        for all u in neighbours(word, d) do  
            H(u) = H(u) + 1  
        done  
        emit(word, H)  
    done
```

Tally co-occurrence of words
u with *word*

```
function reduce(string word, AssociativeArray [H1, H2...]) {  
    Hf = new AssociativeArray()  
    for all H in [H1, H2 ...] do  
        element_sum(Hf, H)  
    done  
    emit(word, Hf)  
}
```

Element-wise sum

- each emitted pair by the reduce corresponds to a **row** of the co-occurrence matrix
- key-space is smaller
 - less sorting/shuffling needed
 - more opportunities for local aggregation
- **possible memory issues:** associative array may grow very large

An alternative approach, dubbed the “stripes” approach, is presented here. Like the pairs approach, co-occurring word pairs are generated by two nested loops. However, the major difference is that instead of emitting intermediate key-value pairs for each co-occurring word pair, co-occurrence information is first stored in an associative array, denoted *H*. The mapper emits key-value pairs with words as keys and corresponding associative arrays as values, where each associative array encodes the co-occurrence counts of the neighbors of a particular word.

The MapReduce execution framework guarantees that all associative arrays with the same key will be brought together in the reduce phase of processing. The reducer performs an element-wise sum of all associative arrays with the same key, accumulating counts that correspond to the same cell in the co-occurrence matrix. The final associative array is emitted with the same word as the key. In contrast to the pairs approach, each final key-value pair encodes a row in the co-occurrence matrix.

It is immediately obvious that the pairs algorithm generates an immense number of key-value pairs compared to the stripes approach. The stripes representation is much more compact, since with pairs the left element is repeated for every co-occurring word pair. The stripes approach also generates fewer and shorter intermediate keys, and therefore the execution framework has less sorting to perform. However, values in the stripes approach are more complex, and come with more serialization and deserialization overhead than with the pairs approach.

Both algorithms can benefit from the use of combiners, since the respective operations in their reducers (addition and element-wise sum of associative arrays) are both commutative and associative. However, combiners with the stripes approach have more opportunities to perform local aggregation because the key space is the vocabulary— associative arrays can be merged whenever a word is encountered multiple times by a mapper. In contrast, the key space in the pairs approach is the cross of the vocabulary with itself, which is far larger—counts can be aggregated only when the same co-occurring word pair is observed multiple times by an individual mapper (which is less likely than observing multiple occurrences of a word, as in the stripes case).

Pairs vs Stripes

- In general, stripes is faster
 - but has potential memory bottleneck
- approaches are endpoints of a continuum
 - divide entire vocabulary in b buckets (e.g. by hashing)
 - words co-occurring with w spread over b assoc. arrays
 - keys emitted are form $(w, 1), (w, 2) \dots (w, b)$

It is important to consider potential scalability bottlenecks of either algorithm. The stripes approach makes the assumption that, at any point in time, each associative array is small enough to fit into memory—otherwise, memory paging will significantly impact performance. The size of the associative array is bounded by the vocabulary size, which is itself unbounded with respect to the total size of the dataset (in worst case, each word in the dataset is different!). Therefore, as the size increase, this will become an increasingly pressing issue—perhaps not for gigabyte-sized datasets, but certainly for terabyte-sized and petabyte-sized sets that will be commonplace tomorrow. The pairs approach, on the other hand, does not suffer from this limitation, since it does not need to hold intermediate data in memory.

To conclude, it is worth noting that the pairs and stripes approaches represent endpoints along a continuum of possibilities. The pairs approach individually records each co-occurring event, while the stripes approach records all co-occurring events with respect to a conditioning event. A middle ground might be to record a subset of the co-occurring events with respect to a conditioning event. We might divide up the entire vocabulary into b buckets (e.g., via hashing), so that words co-occurring with word i (w_i) would be divided into b smaller “sub-stripes”, associated with ten separate keys, $(w_i, 1), (w_i, 2), \dots, (w_i, b)$. This would be a reasonable solution to the memory limitations of the stripes approach, since each of the sub-stripes would be smaller. In the case of $b = |V|$, where $|V|$ is the vocabulary size, this is equivalent to the pairs approach. In the case of $b = 1$, this is equivalent to the standard stripes approach.

Pattern: order inversion

- occurrence matrix has absolute counts
 - some words appear more frequently than others
 - $[m_{ij}]$ can be high simply because w_i is very common
- convert to relative frequencies
 - stripes
 - All data is present in reducer == trivial
 - Memory bottleneck still present
 - pairs
 - pairs $\langle w_i, w_j \rangle$ are sent to the reducer
 - denominator is not available in the reducer

$$f(w_j | w_i) = \frac{N(w_i, w_j)}{\sum_{w'} N(w_i, w')}$$

Let us build on the pairs and stripes algorithms and continue with our running example of constructing the word co-occurrence matrix M for a large text dataset. Recall that in this large square $n \times n$ matrix, where $n = |V|$ (the vocabulary size), cell m_{ij} contains the number of times word w_i co-occurs with word w_j within a specific context. The drawback of absolute counts is that it doesn't take into account the fact that some words appear more frequently than others. Word w_i may co-occur frequently with w_j simply because one of the words is very common. A simple remedy is to convert absolute counts into relative frequencies, $f(w_j | w_i)$. That is, what proportion of the time does w_j appear in the context of w_i ?

In the formula above, $N(\cdot, \cdot)$ indicates the number of times a particular co-occurring word pair is observed in the corpus. We need the count of the joint event (word co-occurrence), divided by what is known as the marginal (the sum of the counts of the conditioning variable co-occurring with anything else).

Computing relative frequencies with the stripes approach is straightforward. In the reducer, counts of all words that co-occur with the conditioning variable (w_i in the above example) are available in the associative array. Therefore, it suffices to sum all those counts to arrive at the marginal and then divide all the joint counts by the marginal to arrive at the relative frequency for all words. This implementation requires minimal modification to the original stripes algorithm. Through appropriate structuring of keys and values, one can use the MapReduce execution framework to

bring together all the pieces of data required to perform a computation. Note that, as with before, this algorithm also assumes that each associative array fits into memory

In the pairs approach, the reducer receives (w_i, w_j) as the key and the count as the value. From this alone it is not possible to compute $f(w_j | w_i)$ since we do not have the marginal.

Order inversion

- General idea:

- compute denominator *before* the pair count
- mapper emits extra $\langle \langle w_i, * \rangle, 1 \rangle$ pair for each $\langle w_i, w_j \rangle$ pair

```
function map(docid a, doc d) {
    for all word in d do
        for all u in neighbours(word, d) do
            emit(pair(word, u), 1)
            emit(pair(word, *), 1)
        done
    done
```

- Requirements for the reducer

- sum $\langle \langle w_i, * \rangle, N \rangle$ pairs
- same reducer must receive all $\langle w_i, * \rangle$ and $\langle w_i, w_j \rangle$ keys!

Recall that in the basic pairs algorithm, each mapper emits a key-value pair with the co-occurring word pair as the key. To compute relative frequencies, we modify the mapper so that it additionally emits a “special” key of the form $(w_i, *)$, with a value of one, that represents the contribution of the word pair to the marginal. Through use of combiners, these intermediate pairs can be aggregated into partial marginal counts $(\langle w_i, * \rangle, N)$ before being sent to the reducers. In the reducer, we can then sum the values of all $(\langle w_i, * \rangle, N)$ pairs to get the number of occurrences of w_i .

To accomplish this, we must keep “state” between individual calls to the reduce method. This is possible in principle, using an initialize method in your reducer class. We initialize an array in which we store all intermediate pairs. After all have received, we iterate over the array and emit the calculated relative frequencies. But there is a better way...

Sorting and shuffling

MapReduce will sort the output of the mapping phase based on the mapping key:

```
< <dog, *>, [104, 23, 1, ...] >  
< <dog, aardvark>, [10, 7, 1] >           (mapping output after possible local  
    < <dog, cat>, [2, 1, 1, 1] >          aggregation)  
    < <doge, *>, [682, ...] >
```

same reducer must receive all $\langle w_i, * \rangle$ and $\langle w_i, w_j \rangle$ keys!

- not guaranteed by default
- intermediate key is sent to reducer i :
 $i = \text{hash}(\text{intermediate key}) \bmod \# \text{reducers}$
- we have to define our own partitioner, so that hash value only depends on first part of the key!

$\langle\langle w_i, * \rangle, 1 \rangle$ pairs must arrive before $\langle\langle w_i, w_j \rangle, 1 \rangle$

- define “sorting” for $\langle w_i, * \rangle$ keys (default is OK: concatenation)

Although the algorithm in the previous slide will work, it suffers from the same drawback as the stripes approach: at some point we might run out of memory because we keep all intermediate pairs in arrays in the reducer.

If we could however influence the order in which the intermediate pairs are delivered to the reducers, we can reduce the memory consumption. MapReduce allows programmers to define the sorting order of keys so that intermediate pairs needed earlier is presented to the reducer before intermediate pairs that are needed later.

Note that in this example, we work with keys composed of two words. If we make sure that the key $\langle w_i, * \rangle$ comes “lexicographically” before any other key $\langle w_i, w_j \rangle$, then we can simply sum the values of $\langle w_i, * \rangle$ in the reducer to get the total count of w_i . There is no need for storing intermediate pairs in an array. As the intermediate pairs with key $\langle w_i, w_j \rangle$ arrive, we can simply divide the sum of their values by this total count and emit the relative frequency.

First, the reducer is presented with the special key $\langle \text{dog}, * \rangle$ and a number of values, each of which represent a partial marginal contribution from the map phase. The reducer accumulates these values to arrive at the total number of occurrences of the word “dog”. The reducer holds this value in its internal state as it processes subsequent keys, like $\langle \text{dog}, \text{aardvark} \rangle$ and $\langle \text{dog}, \text{cat} \rangle$. When it encounters a new special key ($\langle \text{doge}, * \rangle$), it resets its internal state and the process begins again.

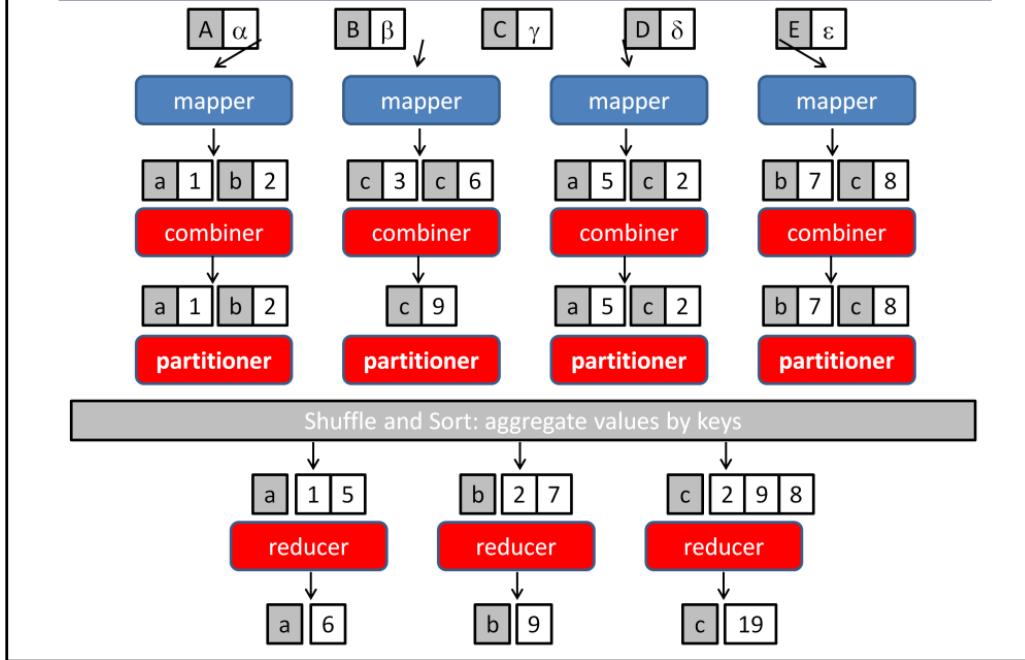
However, even with correct sorting, we must ensure that all pairs with the same left word (w_i) in the key are sent to the **same reducer**. This, unfortunately, does not happen automatically: the default partitioner is based on the hash value of the intermediate key, modulo the number of reducers. For a complex key, the raw byte representation is used to compute the hash value. As a result, there is no guarantee that, for example, intermediate pairs with key (dog, aardvark) and with key (dog, zebra) are assigned to the same reducer. To produce the desired behavior, we must define a custom partitioner that only pays attention to the left word. That is, the partitioner should partition based on the hash of the left word only.

This design pattern, which we call “order inversion”, occurs surprisingly often and across applications in many domains. It is so named because through proper coordination, we can access the result of a computation in the reducer (for example, an aggregate statistic) before processing the data needed for that computation. The key insight is to convert the sequencing of computations into a sorting problem. In most cases, an algorithm requires data in some fixed order: by controlling how keys are sorted and how the key space is partitioned, we can present data to the reducer in the order necessary to perform the proper computations. This greatly cuts down on the amount of partial results that the reducer needs to hold in memory.

To summarize, the specific application of the order inversion design pattern for computing relative frequencies requires the following:

- Emitting a special key-value pair for each co-occurring word pair in the mapper to capture its contribution to the marginal.
- Controlling the sort order of the intermediate key so that the key-value pairs representing the marginal contributions are processed by the reducer before any of the pairs representing the joint word co-occurrence counts.
- Defining a custom partitioner to ensure that all pairs with the same left word are shuffled to the same reducer.
- Preserving state across multiple keys in the reducer to first compute the marginal based on the special key-value pairs and then dividing the counts by the marginals to arrive at the relative frequencies

Complete MapReduce framework

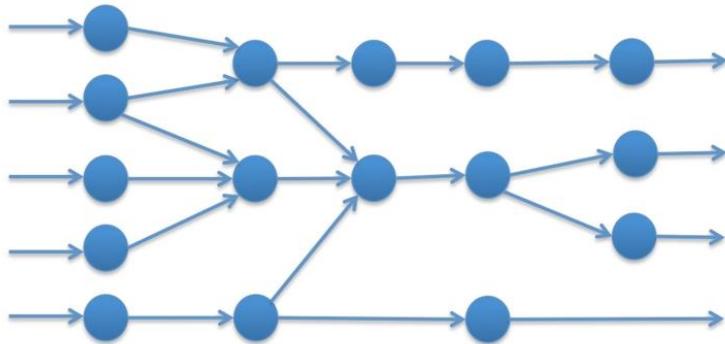


Further reading

- N. Marz and J. Warren, “Big Data”
- Hadoop MapReduce v2 Cookbook
- J. Lin and C. Dyer, “Data-Intensive Text Processing with MapReduce”
- S. Karanth, “Mastering Hadoop”

STREAM PROCESSING

Stream (real-time) processing



The Velocity in Big Data

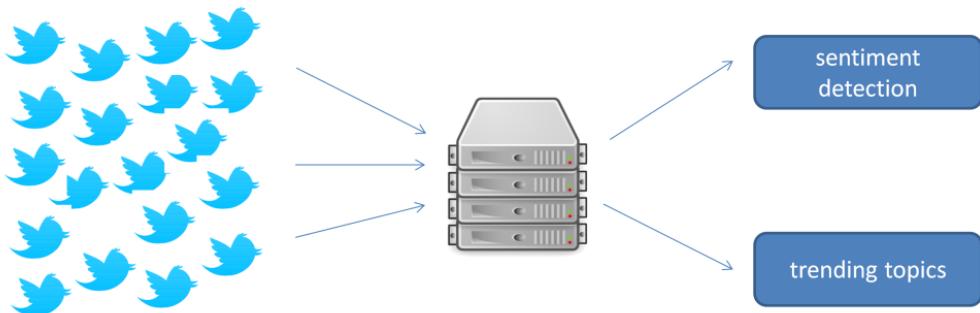
- Batch processing is store-first, process-second
 - Unable to scale for real-time Big Data applications
- Stream processing is an alternative
 - Network of concurrently executing but independent continuous queries
 - Operate on data streams as they flow through
 - Scale with number of incoming streams, not the size of the data

Batch processing is store-first, process-second. However, in many cases we need to be able to analyze data as it comes in (e.g. from sensors, bank card transactions, etc.). Real time (or stream) data processing and analytics allows to take immediate action for those times when acting within seconds or minutes is significant. The goal is to obtain the insight required to act prudently at the right time - which increasingly means immediately.

The high level architecture of a stream processor is a network of processing nodes, where each node performs some action or transformation on the data as it flows through. Each node in the system is a continuously executing and independent query that performs operations on data streams such as filtering, aggregation and analytics.

- Each node is an independent continuous query, that is, a query that never ends.
- All nodes execute concurrently, subscribing or consuming one or more input data streams, and generating one or more output streams.
- The stream processing platform is responsible for the scheduling, query optimization, and runtime execution management, including the movement (or clocking) of data through the system.

Use case: Twitter



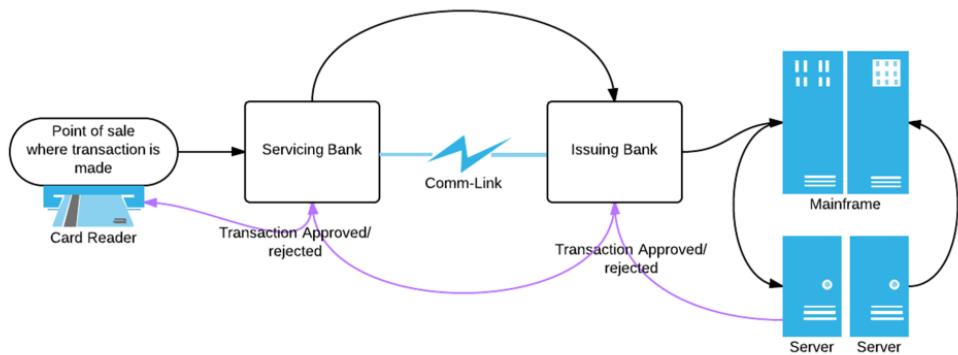
- canonical case for real-time streaming processing
- 6 000 tweets per second
- sentiment detection: advertising, political campaigns...

Twitter allows users to post tweets, messages of up to 140 characters, on its social network. The number of tweets sent per day is sometimes daunting. Every second, on average, around 6,000 tweets are tweeted on Twitter [1], which corresponds to over 350,000 tweets sent per minute, **500 million tweets per day** and around 200 billion tweets per year.

Most work to date has focused on post-facto analysis of tweets, with results coming days or even months after the collection time. However, because tweets are short and easy to send, they lend themselves to quick and dynamic expression of instant reactions to current events. Automated real-time sentiment analysis of this user-generated data can provide fast indications of changes in opinion, showing for example how an audience reacts to particular candidate's statements during a political debate.

[1] www.internetlivestats.com/twitter-statistics

Use case: credit card fraud detection



- System must validate transaction in less than 5 seconds
- Outlier detection via sequence mining (e.g. many transactions within small time window)
- How to handle billions of transactions **in parallel?**

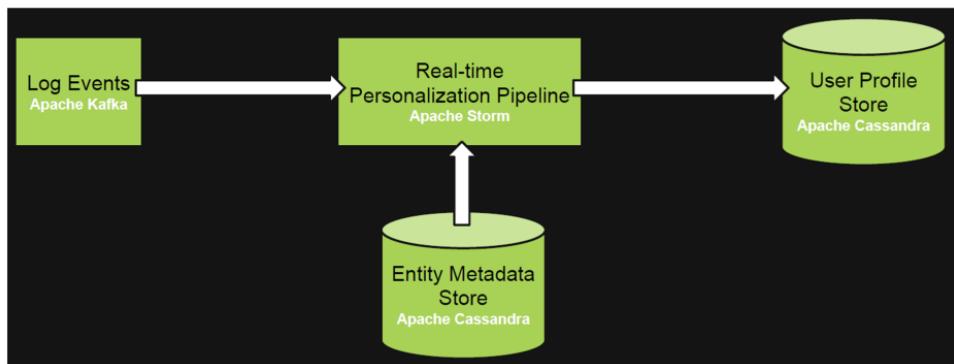
When we make any transaction and swipe our debit or credit card for payment, the duration within which the bank has to validate or reject the transaction in less than five seconds. In less than five seconds, data or transaction details have to be encrypted, travel over secure network from servicing back bank to the issuing bank, then at the issuing back bank the entire fuzzy logic for acceptance or decline of the transaction has to be computed, and the result has to travel back over the secure network.

The challenges such as network latency and delay can be optimized to some extent, but to achieve the preceding featuring transaction in less than 5 seconds, one has to design an application that is able to churn a considerable amount of data and generate results within 1 to 2 seconds.

A sample of how fraud detection is implemented can be found here:
<https://pkghosh.wordpress.com/2013/10/21/real-time-fraud-detection-with-sequence-mining/>

Use case: Spotify

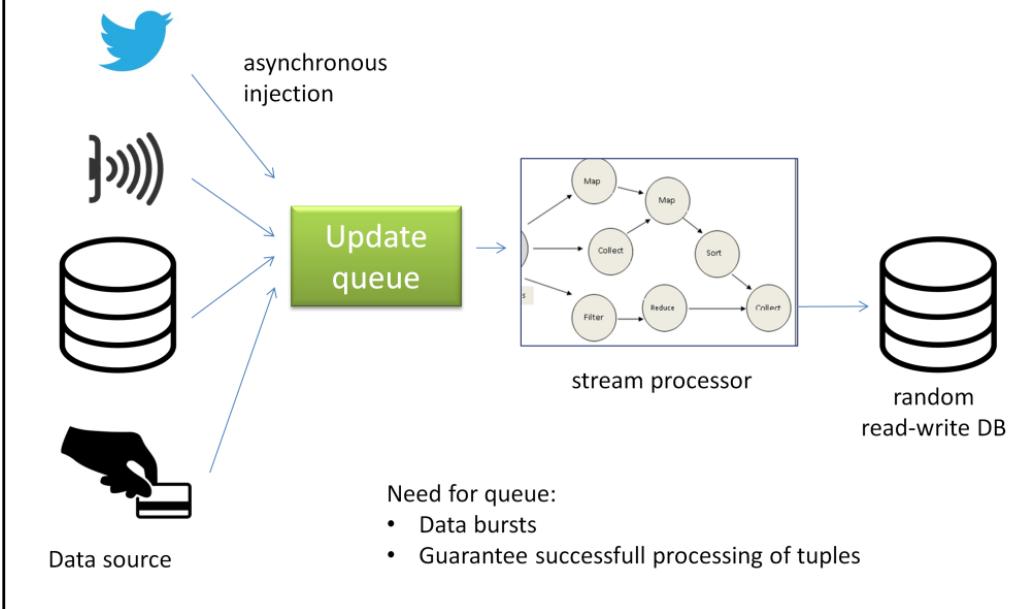
- 200 000 tuples per second
- recommendation, ads, monitoring, analytics
 - real-time changing contexts and conditions



Also Spotify uses real-time streaming for recommendation, advertisement, monitoring and analytics. Recalculating recommendations over their entire batch of data (playlists, play history, user profile) would be infeasible: users can quickly skip a number of recommended tracks, which would quickly exhaust the calculated list of suggestions.

Moreover, it is important to cope quickly with changing contexts and conditions. For example, a metal genre listener might not enjoy an announcement for a metal genre album when they are trying to put their kid to sleep and playing kid's music at night.

Generic architecture

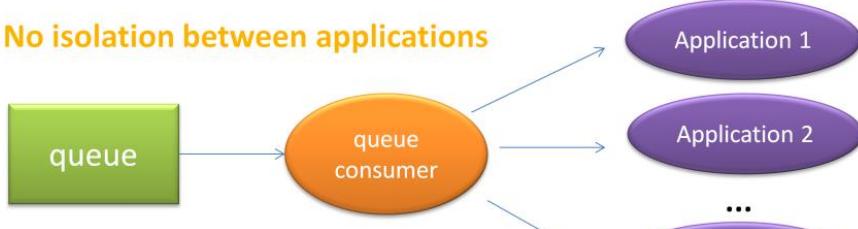


The high-level architecture of any stream processing framework is depicted on the slide. The architecture is asynchronous: data sources simply inject their data into the stream processing framework without waiting for confirmation that the data has been persisted in the DB. Data sources are often external; or serve many consumers, so working synchronously would block them.

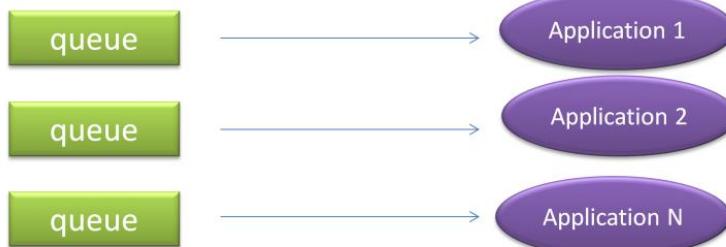
Data sources are injecting their data into an update queue. Data sources can generate data from specific events (tweets, sensor data, credit card transactions), but may also emit data from a database. To understand the need for a queue, let's hypothesize on a system without a queue. In such a system, events would be handed directly to worker nodes in the stream processor that would process each event independently. This fire-and-forget system cannot guarantee that all the data is successfully processed. A worker can die before completing its assigned task, but there is no mechanism to detect or correct the error. The architecture is also susceptible to bursts in traffic that exceed the resources of the processing cluster. In such a scenario, the client would be overwhelmed and messages could be lost.

Need for multi-consumer queue

No isolation between applications



Load proportional to #appl times # events



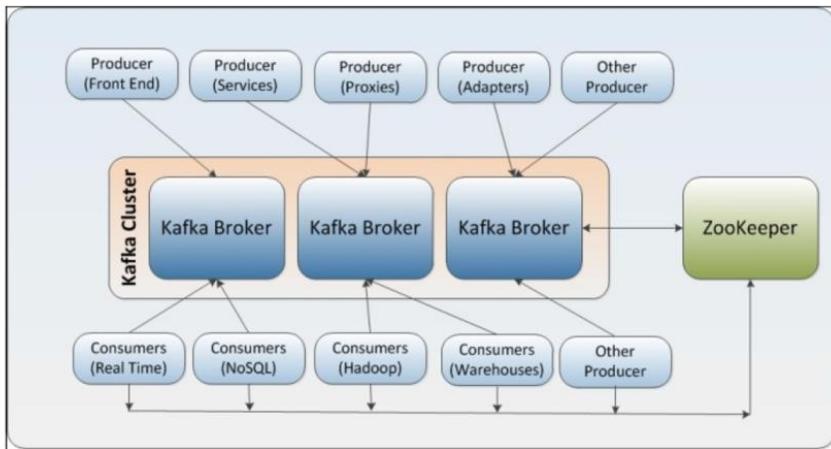
While it is now clear that an asynchronous architecture needs a queue to persist an event stream, the semantics of good queue design require further discussion. Many implementations, including the standard Java Queue and RabbitMQ are single-consumer queues. This design is based on the idea that when you read an event from the queue, the event is not immediately removed. Instead, the item taken from the queue contains an identifier that you later use to acknowledge success or report failure for the processing of the event. Only when an event is acked will it be removed from the queue. If the event processing fails or a timeout occurs, the queue server will allow another client to retrieve the same event via another get call. An event may therefore be processed multiple times with this approach (for example, when a client processes an event but dies before acknowledging it), but each event is guaranteed to be processed at least once.

There is a deep flaw in this queue design: what if multiple applications want to consume the same stream? This is quite common. For example, given a stream of pageviews, one application could build a view of pageviews over time while another could build a view of unique visitors over time. One possible solution would be to wrap all the applications within the same consumer. This is a bad design as it eliminates any isolation among independent applications. If one application has a bug, it could potentially affect all the other applications running within the same consumer. With a single-consumer queue, the only way to achieve independence between applications is to maintain a separate queue for each consumer application.

If you have three applications, you maintain three separate copies of the queue. The load is now proportional to the number of applications multiplied by the number of incoming events, rather than just to the number of incoming events.

What we really need is a single queue that can be used by many consumers, where adding a consumer is simple and introduces a minimal increase in load. The fundamental issue with a single-consumer queue is that the queue is responsible for keeping track of what's consumed. Because of the restrictive condition that an item is either 'consumed' or 'not consumed', the queue is unable to gracefully handle multiple clients wanting to consume the same item.

Apache Kafka



Open source, distributed, partitioned and replicated commit-log publish-subscribe messaging system

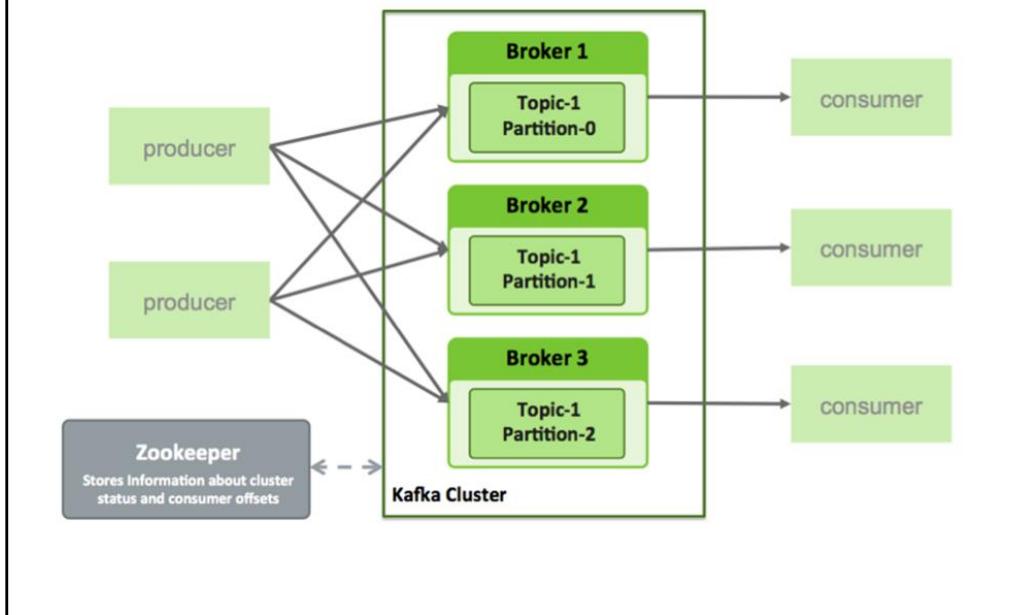
- Persistent messaging – no information loss
- High throughput – even on commodity hardware
- Real time – messages are immediately visible to consumer threads

Apache Kafka is an open source, distributed, partitioned, and replicated commit-log-based publish-subscribe messaging system, mainly designed with the following characteristics:

- **Persistent messaging:** To derive the real value from big data, any kind of information loss cannot be afforded. Apache Kafka is designed with O(1) disk structures that provide constant-time performance even with very large volumes of stored messages that are in the order of TBs. With Kafka, messages are persisted on disk as well as replicated within the cluster to prevent data loss.
- **High throughput:** Keeping big data in mind, Kafka is designed to work on commodity hardware and to handle hundreds of MBs of reads and writes per second from large number of clients.
- **Distributed:** Apache Kafka with its cluster-centric design explicitly supports message partitioning over Kafka servers and distributing consumption over a cluster of consumer machines while maintaining per-partition ordering semantics. Kafka cluster can grow elastically and transparently without any downtime.
- **Real time:** Messages produced by the producer threads should be immediately visible to consumer threads; this feature is critical to event-based systems.

The diagram on the slide shows a typical big data aggregation-and-analysis scenario supported by the Apache Kafka messaging system, with different kinds of producers and different kinds of consumers.

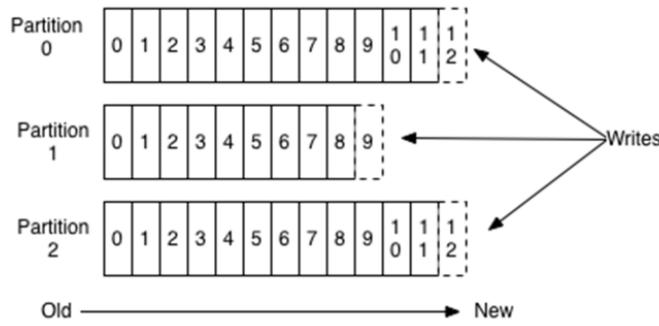
Basic terminology



A Kafka cluster primarily has five main components:

- Kafka maintains feeds of messages in categories called **topics**. A topic is a category or feed name to which messages are published by the message producers. In Kafka, topics are partitioned.
- Kafka cluster consists of one or more servers where each one may have one or more server processes running and is called the **broker**. Topics are created within the context of broker processes.
- **ZooKeeper** serves as the coordination interface between the Kafka broker and consumers. ZooKeeper is akin to etcd: it allows distributed processes to coordinate with each other. It stores coordination data: status information, configuration, location information, and so on.
- Process that publish messages to a Kafka topic are referred to as **producers**.
- Processes that subscribe to topics and process the feed of published messages are referred to as **consumers**.

Topics: anatomy and publishing

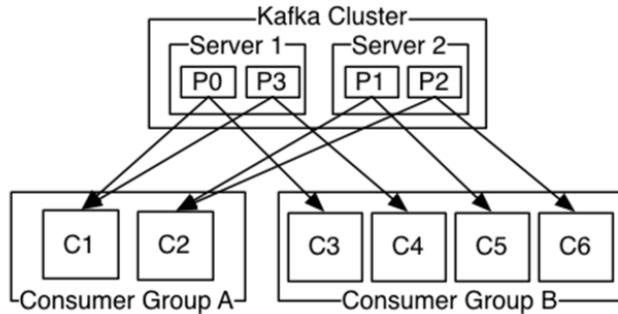


- Topic comprises a partitioned log
- Each log is ordered, immutable sequence of messages
- Append-only
- Producers choose appropriate partition within the topic to write to

For each topic, Kafka maintains a partitioned log as shown on the slide. Each partition is an ordered, immutable sequence of messages that is continually appended to – a commit log. The messages in the partitions are each assigned a sequential id number called the *offset* that uniquely identifies each message within the partition.

The partitions in the log serve several purposes. First, they allow the log to scale beyond a size that will fit on a single server. Each individual partition must fit on the servers that host it, but a topic may have many partitions so it can handle an arbitrary amount of data. Second they act as the unit of parallelism. Producers publish data to the topics by choosing the appropriate partition within the topic. For load balancing, the allocation of messages to the topic partition can be done in a round-robin fashion or using a custom defined function.

Consuming a topic



- Consumer processes ordered in **consumer groups**
- Message within a topic is consumed by a single process per consumer group
- Ordering guarantees *within* a partition: one partition → one consumer process
- Broker maintains offset per consumer, but offset is controlled entirely by consumer
- Cluster has retention period – older messages are automatically discarded

The Kafka platform has the concept of consumer groups. Here, each consumer is represented as a process and these processes are organized within groups called **consumer groups**.

A message within a topic is consumed by a single process (consumer) within the consumer group and, if the requirement is such that a single message is to be consumed by multiple processes, all these consumer processes need to be kept in different consumer groups.

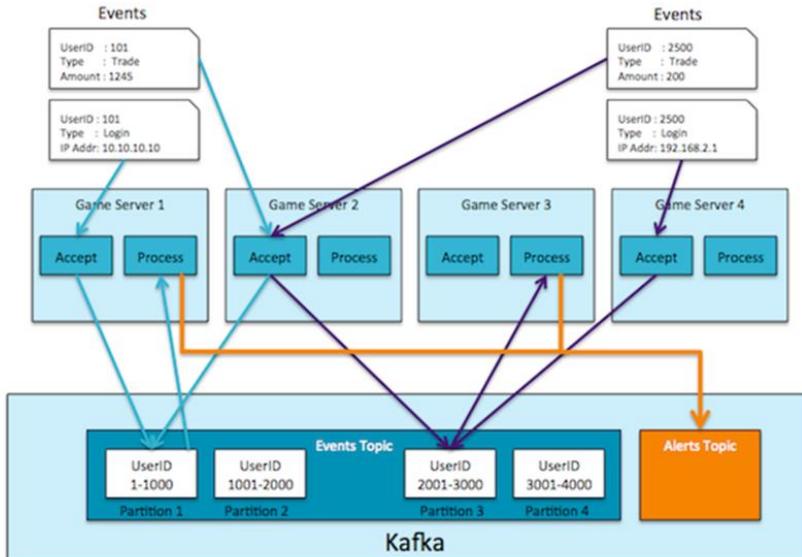
Kafka only provides a total order over messages *within* a partition. To achieve this, Kafka assigns the partitions in the topic to the consumers in the consumer group: each partition is consumed by exactly one consumer in the group. This ensures that a single process is the only reader of that partition and consumes the data in order. Since there are many partitions this still balances the load over many consumer instances. Note however that there cannot be more consumer processes than partitions.

Consumers always consume messages from a particular partition sequentially and also acknowledge the message offset. This acknowledgement implies that the consumer has consumed all prior messages. Consumers issue an asynchronous pull request containing the offset of the message to be consumed to the broker and get the buffer of bytes.

Brokers are stateless, which means the message state (messages was successfully processed or not) of any consumed message is maintained within the message consumer, and the Kafka broker does not maintain a record of what is consumed by whom. The only metadata retained on a per-consumer basis is the position of the consumer in the log (“the offset”). This offset is controlled by the consumer: normally the consumer will advance its offset as it reads data, but it can consume messages in any order it likes. For example a consumer can reset to an older offset to reprocess.

The Kafka cluster retains all published for a configurable period of time. Messages older than the retention period will be deleted -- whether or not they have been consumed.

Kafka at work: MMOG



We will illustrate Kafka with the use case of a massive multiplayer online game (MMOG). In these games, players cooperate and compete with each other in a virtual world. Often players trade with each other, exchanging game items and money, so as game developers it is important to make sure players don't cheat: trades will be flagged if the trade amount is significantly larger than normal for the player and if the IP the player is logged in with is different than the IP used for the last 20 games. In addition to flagging trades in real-time, we also want to load the data to Apache Hadoop, where our data scientists can use it to train and test new algorithms.

For the real-time event flagging, it will be best if we can reach the decision quickly based on data that is cached on the game server memory, at least for our most active players. Our system has multiple game servers and the data set that includes the last 20 logins and last 20 trades for each player can fit in the memory we have, if we partition it between our game servers.

Our game servers have to perform two distinct roles: The first is to accept and propagate user actions and the second to process trade information in real time and flag suspicious events. To perform the second role effectively, we want the whole history of trade events for each user to reside in memory of a single server. This means we have to pass messages between the servers, since the server that accepts the user action may not have his trade history. To keep the roles loosely coupled, we use Kafka to pass messages between the servers, as you'll see below.

We have configured Kafka with a single topic for logins and trades. The reason we need a single topic is to make sure that trades arrive to our system after we already have information about the login (so we can make sure the gamer logged in from his usual IP). Kafka maintains order within a topic, but not between topics.

When a user logs in or makes a trade, the accepting server immediately sends the event into Kafka. We send messages with the user id as the key, and the event as the value. This guarantees that all trades and logins from the same user arrive to the same Kafka partition. Each event processing server runs a Kafka consumer, each of which is configured to be part of the same group—this way, each server reads data from few Kafka partitions, and all the data about a particular user arrives to the same event processing server (which can be different from the accepting server). When the event-processing server reads a user trade from Kafka, it adds the event to the user's event history it caches in local memory. Then it can access the user's event history from the local cache and flag suspicious events without additional network or disk overhead.

It's important to note that we create a partition per event-processing server, or per core on the event-processing servers for a multi-threaded approach. This may sound like a circuitous way to handle an event: send it from the game server to Kafka, read it from another game server and only then process it. However, this design decouples the two roles and allows us to manage capacity for each role as required. In addition, the approach does not add significantly to the timeline as Kafka is designed for high throughput and low latency; even a small three-node cluster can process close to a million events per second with an average latency of 3 ms.

When the server flags an event as suspicious, it sends the flagged event into a new Kafka topic—for example, Alerts—where alert servers and dashboards pick it up. Meanwhile, a separate process reads data from the Events and Alerts topics and writes them to Hadoop for further analysis.

Models for stream processing



	One-at-a-time	Micro-batched
Lower latency	✓	
Higher throughput		✓
At-least-once semantics		✓
Exactly-once semantics	in some cases	✓
Simpler programming model	✓	

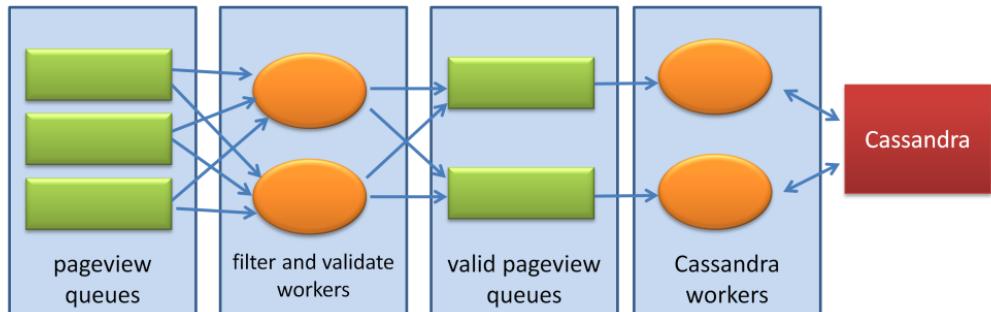
One-at-a-time and micro-batched are two recent models of stream processing. There are trade-offs to consider, since each has its strengths and weaknesses. They are very much complementary – some applications are better suited for one-at-a-time stream processing, and micro-batch stream processing is a better choice for others.

A big advantage of one-at-a-time stream processing is that it can process streams with lower latency than micro-batched processing. Example applications that greatly benefit from this attribute include alerting and financial trading.

The latency and throughput characteristics are different for micro-batch processing. For any individual tuple, the latency from when it's added to the source queue to when it's fully processed is much higher in micro-batch processing. There is a small but significant amount of overhead to coordinating batches that increases latency, and instead of waiting for just one tuple to complete, processing needs to be done on many more tuples. In practice, this turns out to be latency on the order of hundreds of milliseconds to seconds.

But micro-batch processing can have higher throughput than one-at-a-time processing (number of tuples processed per second). Whereas one-at-a-time processing must do tracking on an individual tuple level, micro-batch processing only has to track at a batch level. This means fewer resources are needed on average per tuple, allowing micro-batch processing to have higher throughput than one-at-a-time processing.

Queues and workers paradigm



- first set of workers partition outgoing stream by URL to avoid write race conditions
 - poor fault tolerance if one of the Cassandra workers goes down
- operational burden due to queues between every set of workers
- queues add latency and decrease the system throughput
- each intermediate queue needs to be managed, monitored and scaled
- tedious to build

The slide illustrates a one-at-a-time stream processing model for a pageview (per URL) counting problem. The first set of workers reads pageview events from a set of queues, validates each pageview to filter out invalid URLs, and then passes the events to a second set of workers. The second set of workers then updates the pageview counts of the valid URLs.

The queues-and-workers paradigm is straightforward but not necessarily simple. One subtlety is the need to ensure that multiple workers don't attempt to update (increment) the pageview count of the same URL at the same time to avoid race conditions. To meet this guarantee, the first set of workers partitions its outgoing stream by the URL. With this partitioning, the entire set of URLs will still be spread among the queues, but pageview events for any given URL will always go to the same queue (e.g. by taking the modulo of the hash of the URL). Unfortunately, a consequence of partitioning over queues is poor fault tolerance. If a worker that updates the pageview counts in the database goes down, no other workers will update the database for that portion of the stream. You'll have to manually start the failed worker somewhere else, or build a custom system to automatically do so.

Another problem is that having queues between every set of workers adds to the operational burden of your system. If you need to change the topology of your processing, you will need to coordinate your actions so that the intermediate queues are cleared before your redeploy.

Queues also add latency and decrease the throughput of your system because each event passed from worker to worker is forced to go through a third party, where it must be persisted to disk.

On top of everything else, each intermediate queue needs to be managed and monitored and adds yet another layer that needs to be scaled.

Perhaps the biggest problem with this approach is how tedious it is to build. Much code is required for (de)serialization, to pass objects through queues, routing logic to connect worker pools, and instructions for deploying workers over a cluster of servers.

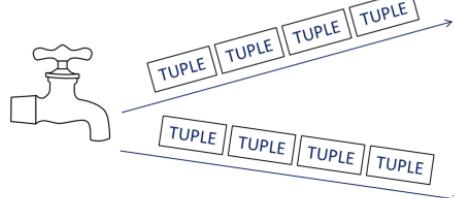
Clearly, there is need for a higher-level abstraction framework.

Apache Storm model

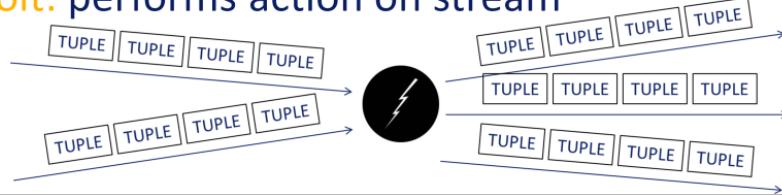
- **stream:** infinite sequence of **tuples**



- **spout:** source of streams



- **bolt:** performs action on stream

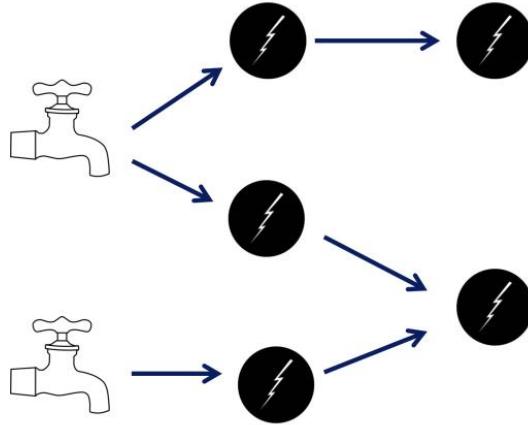


At the core of the Storm model are *streams*. A stream is an infinite sequence of *tuples*, where a tuple is simply a named list of values. In essence, the Storm model is about transforming streams into new streams, potentially updating databases along the way.

The next abstraction in the Storm model is the *spout*. A spout is a source of streams in a topology. A spout could read from some data source (e.g. a sensor, a Kafka queue) and turn the data into a tuple stream, or a timer spout could emit a tuple into its output stream every 10 seconds.

While spouts are sources of streams, the *bolt* abstraction performs actions on streams. A bolt takes any number of streams as input and produces any number of streams as output. Bolts implement most of the logic in a topology, they run functions, filter data, compute aggregations, do streaming joins, update databases, and so forth.

Topology

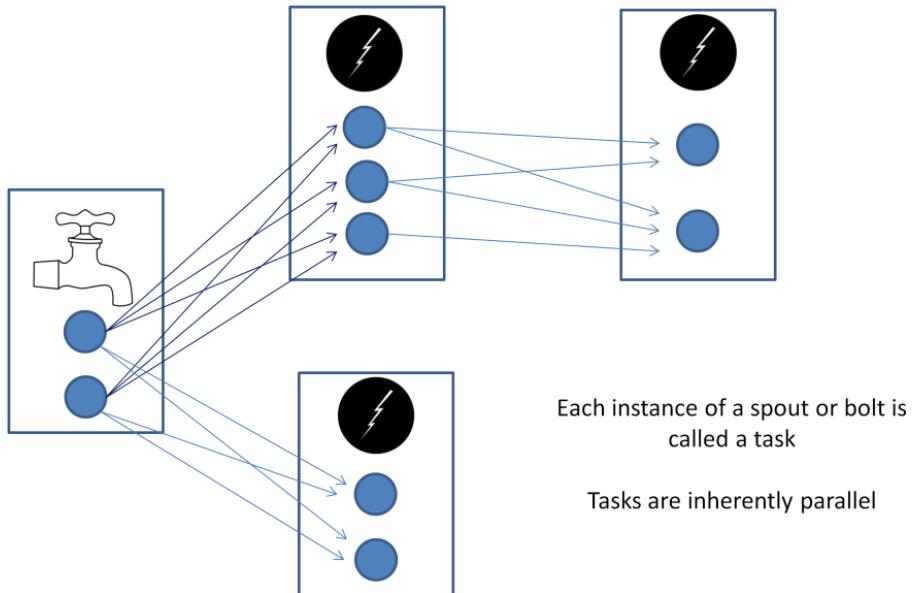


- Network of spouts and bolts
- Each edge represents a bolt that processes the output stream of another spout or bolt.
- Defines how tuples flow through a Storm application

The Storm model represents the entire stream-processing pipeline as a graph of computation called a topology. Rather than write separate programs for each node of the topology and connect them manually, as required in the queues-and-workers scheme, the Storm model involves a single program that's deployed across a cluster. This flexible approach allows a single executable to filter data in one node, compute aggregates with a second node, and update realtime view databases with a third. Serialization, message passing, task discovery, and fault tolerance can be handled for you by the abstractions, and this can all be done while achieving very low latency.

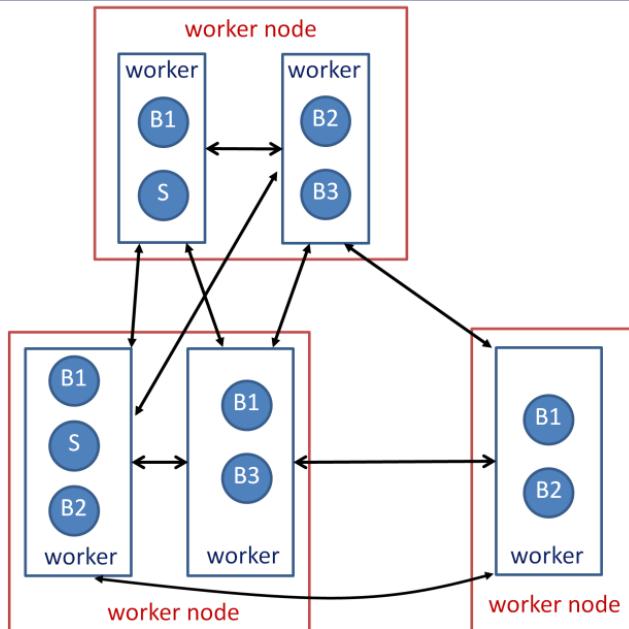
A topology is therefore a network of spouts and bolts with each edge representing a bolt that processes the output stream of another spout or bolt. The topology defines how tuples flow through a Storm application.

Tasks



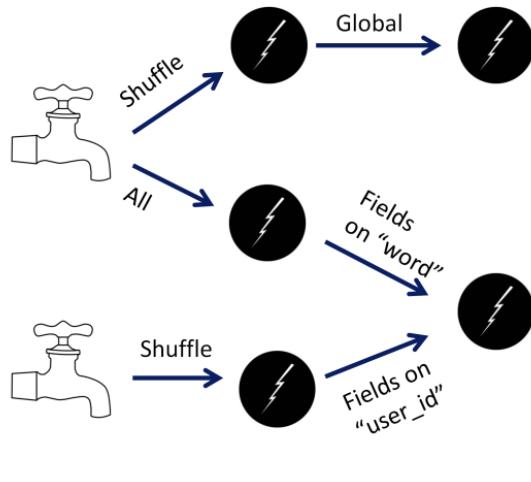
Each instance of a spout or bolt is called a *task*. The key to the Storm model is that tasks are inherently parallel – exactly like how map and reduce task are inherently parallel in MapReduce. Spouts and bolts consist of multiple tasks that are executed in parallel. A bolt task receives tuples from all tasks that generate the bolt's input stream.

Physical view



All the tasks for a given spout or bolt will not necessarily run on the same machine. Instead, they are spread among the different workers of the cluster. The above figure depicts a topology grouped by physical machines.

Stream groupings



Stream groupings

- shuffle
- fields
- all
- global

The fact that spouts and bolts run in parallel brings up a key question: when a task emits a tuple, which of the consuming tasks should receive it? The Storm model requires *stream groupings* to specify how tuples should be partitioned among consuming tasks.

The simplest kind of stream grouping is a *shuffle grouping* that distributes tuples using a random round-robin algorithm. This grouping evenly splits the processing load by distributing the tuples randomly but equally to all consumers. Another common grouping is the *fields grouping* that distributes tuples by hashing a subset of the tuple fields and modding the result by the number of consuming tasks.

All grouping replicates the stream across all the bolt's tasks. This grouping must be used with care.

Global grouping makes the entire stream go to a single one of the bolt's tasks. Specifically, it goes to the task with the lowest id.

Example: word counting



```
TopologyBuilder builder = new TopologyBuilder();

builder.setSpout("spout", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentence(),8)
    .shuffleGrouping("spout");
builder.setBolt("count", new WordCount(),12)
    .fieldsGrouping("split", new Fields("word"));
```

Just as word count is the de facto introductory MapReduce example, let's see what the streaming version of word count looks like in the Storm model.

The splitter bolt transforms a stream of sentences into a stream of words, and the word-count bolt consumes the words to compute the word counts. The key here is the fields grouping between the splitter bolt and the word-count bolt. That ensures that each word-count task sees every instance of every word they receive; making it possible for them to compute the correct count.

Counting bolt

```
public static class WordCount extends BaseBasicBolt {
    Map<String, Integer> counts = new HashMap<String, Integer>();

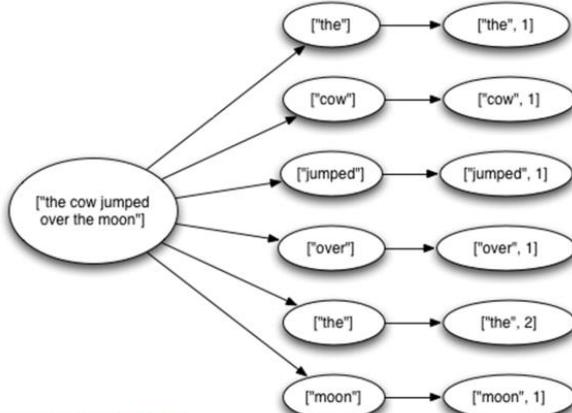
    @Override
    public void execute(Tuple tuple, BasicOutputCollector collector) {
        String word = tuple.getString(0);
        Integer count = counts.get(word);
        if (count == null)
            count = 0;
        count++;
        counts.put(word, count);
        collector.emit(new Values(word, count));
    }

    @Override
    public void declareOutputFields(OutputFieldsDeclarer declarer) {
        declarer.declare(new Fields("word", "count"));
    }
}
```

This is the code implementing the WordCount bolt in Storm. The execute method receives a tuple and a collector to emit output to. As you can see, you can keep state per Bolt instance. To make sure that your state is consistent between different Bolts, you need to use the appropriate stream groupings. In our word count example, this means that each word should be in the Map object of only one Bolt instance.

You can also see that the Storm model requires no logic around where to send tuples or how to serialize tuples. That is all handled underneath the Storm abstractions.

At-least-once guarantee



At-least-once guarantee:

- Storm tracks “Tuple DAG”
- Tuples are retried from the spout upon downstream failure
- Requires *anchoring* and *acking* of tuples

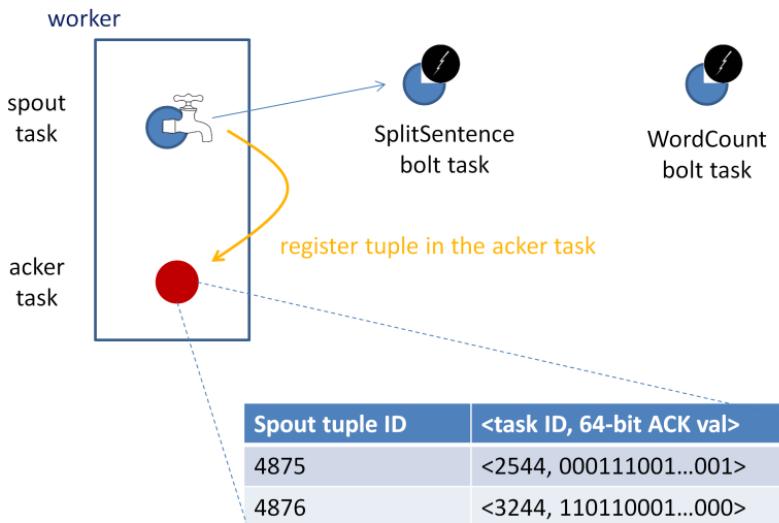
Worker nodes may go down during the processing. Storm guarantees that each tuple coming from a spout is processed *at least* once. To understand the mechanism, let's take a look at what the processing of a tuple looks like in the word-count topology.

When a sentence tuple is generated by the spout, it's sent to whatever bolts subscribe to that spout. In this case, the word-splitter bolt creates six new tuples based on that spout tuple. Those word tuples go on to the word-count bolt, which creates a single tuple for every one of those word tuples. You can visualize all the tuples created during the processing of a single spout tuple as a directed acyclic graph (DAG). Let's call this the *tuple DAG* (shown in the slide). You could imagine much larger tuple DAGs for more involved topologies.

Storm uses an efficient and scalable algorithm for tracking tuple DAGs and retrying tuples from the spout if there's a failure somewhere downstream. Retrying tuples from the spout will cause the entire tuple DAG to be regenerated. Retrying from the spout may seem a step backward, since one failure somewhere deep in the DAG means redoing all the calculations, even those intermediate stages that had completed successfully. But upon further inspection, this model is no different than the queues-and-workers model. With queues and workers, a stage could succeed in processing, fail right before acknowledging the message and letting it be removed from the queue, and then be tried again. In both scenarios, the processing guarantee is still an at-least-once guarantee.

There's two things you have to do as a user to benefit from Storm's reliability capabilities. First, you need to tell Storm whenever you're creating a new link in the tree of tuples. Second, you need to tell Storm when you have finished processing an individual tuple. By doing both these things, Storm can detect when the tree of tuples is fully processed and can ack or fail the spout tuple appropriately. Storm's API provides a concise way of doing both of these tasks. Specifying a link in the tuple tree is called *anchoring*, and each bolt in the DAG will acknowledge the processing of the tuple.

General principle



A Storm topology has a set of special "acker" tasks that track the DAG of tuples for every spout tuple. When an acker sees that a DAG is complete, it sends a message to the spout task that created the spout tuple to ack the message. Storm defaults the number of "acker tasks" to be equal to the number of workers configured in the topology.

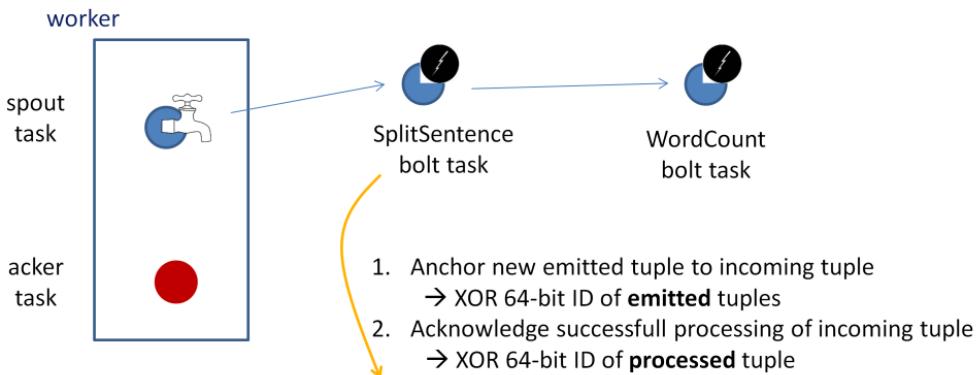
When a tuple is created in a topology (be it in a spout or in a bolt), it is given a random 64 bit id. These ids are used by ackers to track the tuple DAG for every spout tuple.

When a spout task emits a new tuple, it simply sends a message to the appropriate acker telling it that its task id is responsible for that spout tuple. Then when an acker sees a tree has been completed, it knows to which task id to send the completion message. Acker tasks do not track the tree of tuples explicitly. For large tuple trees with tens of thousands of nodes (or more), tracking all the tuple trees could overwhelm the memory used by the ackers. Instead, the ackers take a different strategy that only requires a fixed amount of space per spout tuple (about 20 bytes). This tracking algorithm is the key to how Storm works and is one of its major breakthroughs.

An acker task stores a map from a spout tuple id to a pair of values. When emitting a tuple, the Spout provides a (self-chosen) "message id" that allows it to identify this

tuple. The first value in the map is the task id that created the spout tuple which is used later on to send completion messages. The second value is a 64 bit number called the "ack val". It is initialized to the random 64-bit ID that was given to the tuple by Storm.

General principle (2)



Spout tuple ID	<task ID, 64-bit ACK val>
4875	<2544, 011111000...101>
4876	<3244, 110010101...110>

The tuples emitted by the spout are received by an instance (task) of the SplitSentence bolt. This bolt will process the incoming tuple (containing a complete sentence) and emit a number of new tuples (one per word). Storms automatically assigns a new random 64-bit ID to each of these tuple IDs. These newly emitted tuples represent new links in the tuple DAG that originates in the sentence tuple.

First, the bolt must tell Storm that it has created new links in the tuple DAG. This is called *anchoring* and should be done each time you emit a new tuple. Anchoring means that you contact the appropriate acker and XOR the 64-bit of all the **newly emitted tuples** in the 64-bit ACK val.

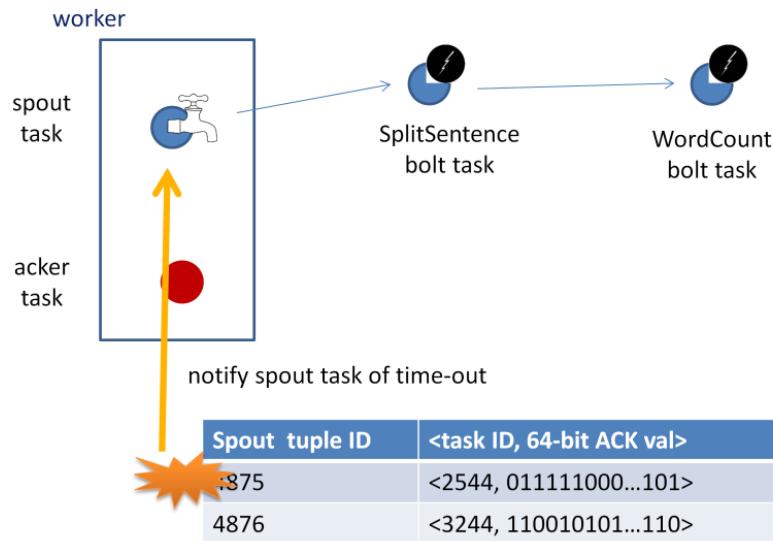
Second, the bolt must tell Storm that it has successfully processed an individual tuple. This is called *acking*. The 64-bit ID of the acked tuple is XORed with the 64-bit ACK val of the original Spout message ID.

A similar acking is done by the WordCount task that receives the word tuples. If this WordCount bolt task is the last one in the topology, it will not emit new tuples and only acknowledge the received tuples. This means that the 64-bit ACK value will become zero.

The ack val is thus a representation of the state of the entire tuple tree, no matter how big or how small. When an acker task sees that an "ack val" has become 0, then it knows that the tuple tree is completed. Since tuple ids are random 64 bit numbers,

the chances of an "ack val" accidentally becoming 0 is extremely small and even then, it will only cause data loss if that tuple happens to fail in the topology.

General principle (3)



Upon time-out of a non-zero 64-bit ACK value, the acker task will notify the spout task that registered this tuple that the tuple was not successfully processed along the entire topology. The spout task will then re-emit this tuple.

Anchoring and acking

```
public static class SplitSentence extends BaseRichBolt {  
    OutputCollector _collector;  
  
    public void execute(Tuple tuple) {  
        String sentence = tuple.getString(0);  
        for(String word: sentence.split(" ")) {  
            _collector.emit(tuple, new Values(word));  
        }  
        _collector.ack(tuple);  
    }  
}
```

ack processing of received tuple

anchoring the newly emitted tuples to the spout tuple ID (contained in the tuple received)

Extend BaseBasicBolt class if you want automatically to:

- anchor all outgoing bolts to the input tuple
- ack the input tuple at the end of the execute method

This bolt splits a tuple containing a sentence into a tuple for each word. Each word tuple is *anchored* by specifying the input tuple as the first argument to emit. Since the word tuple is anchored, the spout tuple at the root of the tree will be replayed later on if the word tuple failed to be processed downstream.

In contrast, let's look at what happens if the word tuple is emitted like this:

```
_collector.emit(new Values(word));
```

Emitting the word tuple this way causes it to be *unanchored*. If the tuple fails to be processed downstream, the root tuple will not be replayed. Depending on the fault-tolerance guarantees you need in your topology, sometimes it can be appropriate to emit an unanchored tuple.

Aggregating/joining streams

incoming tuple not immediately ACK'ed

```
public static class MultiAnchorer extends BaseRichBolt {  
    OutputCollector _collector;  
    List<Tuple> _buffer = new ArrayList<Tuple>();  
    int _sum = 0;  
  
    public void execute(Tuple tuple) {  
        _sum += tuple.getInteger(0);  
        if (_buffer.size() < 100) {  
            _buffer.add(tuple);  
        }  
        else {  
            _collector.emit(_buffer, new Values(_sum)); ←  
            for (Tuple _tuple : buffer)  
                _collector.ack(_tuple);  
            _buffer.clear();  
            _sum = 0;  
        }  
    }  
}
```

Emit tuple with
sum, anchored
to **all tuples** in
the buffer

Ack all tuples in
the buffer

This is an example of advanced anchoring/acking. This bolt emits the sum of 100 tuples received. Once 100 values have been summed, one tuple containing the sum is emitted. The first argument to the emit method is now a list of tuples, and the newly emitted tuple is anchored to the DAG tree of *all* tuples in this list. Subsequently, we individually ack all tuples in the buffer, since they have been sucessfully consumed.

Note that this bolt is implemented as a subclass of the BaseRichBolt class, requiring you to explicitly handle anchoring and acking of tuples. It is a very common pattern for bolts to anchor all outgoing tuples to the input tuple, and then ack that tuple at the end. To automate this behavior, Storm provides a BaseBasicBolt class that takes care of this style of anchoring/acking. This style was illustrated in the source code of the WordCount bolt presented a few slides earlier.

Exactly-one semantics

- one-at-a-time stream processing
 - very low latency
 - simple
 - at-least-once processing guarantee during failure
 - inaccuracy (e.g. counting) sometimes unacceptable
- micro-batch stream processing
 - full accuracy all of the time
 - at the cost of higher latency

One-at-a-time stream processing is very low latency and simple to understand. But it can only provide an at-least-once processing guarantee during failures. Although this doesn't affect accuracy for *idempotent* operations, like adding elements to a set, it does affect accuracy for other operations such as counting.

Sometimes, this level accuracy is not sufficient. In those cases, micro-batch processing can give you the fault-tolerant accuracy you need, at the cost of higher latency in the order of hundreds of milliseconds to seconds.

Strongly-ordered processing

One-at-a-time processing provides no exactly-once semantics:

```
process (tuple) {  
    counter.increment();  
}
```

Key idea:

- Enforce **strong ordering** on the processing of the input stream
- Store result along with ID of latest tuple processed



With one-at-a-time stream processing, tuples are processed independently of each other. Failures are tracked at an individual tuple level, and replays also happen at an individual tuple level. The one-at-a-time stream processing is very low latency and simple, but it can only provide an at-least-once processing guarantee during failures. Although this doesn't affect accuracy for certain operations, like adding elements to a set, it does affect accuracy for other operations such as counting.

In the one-at-a-time code on the slide, tuples will be replayed after failure but when it comes time to increment the count, you have no idea if that tuple was processed already or not. It is possible you incremented the count but then crashed immediately before acking the tuple. The only way to know is to store the ID of every tuple you have processed – but that's not a very viable solution.

The key to achieving exactly-once semantics is to enforce a strong ordering on the processing of the input stream. Assume you only process one tuple at a time and you don't move on to the next tuple until the current one is successfully processed through the entire topology. In addition, we assume that every tuple has a unique ID associated that is always the same no matter how many times it is replayed. The key idea is rather than just store the count, you store the count along with the ID of the latest tuple processed. When you are update the count, you first check the stored ID:

- The stored ID is the same as the current tuple ID. In this case, you know that the count already reflects the current tuple, so you do nothing
- The stored ID is different from the current tuple ID. In this case, you increment the

counter and update the stored ID. This works because tuples are processed in order, and the count and ID are updated atomically.

This update strategy is resilient to all failure scenarios. If the processing fails after updating the count, then the tuple will be replayed and the update will be skipped the second time around. If the processing fails before updating the count, then the update will occur the second time around.

Micro-batch stream processing

- Process tuples in discrete batches
- Batches processed in order
- All tuples of a batch must be processed before moving on to the next batch
 - parallelization of tuple processing possible

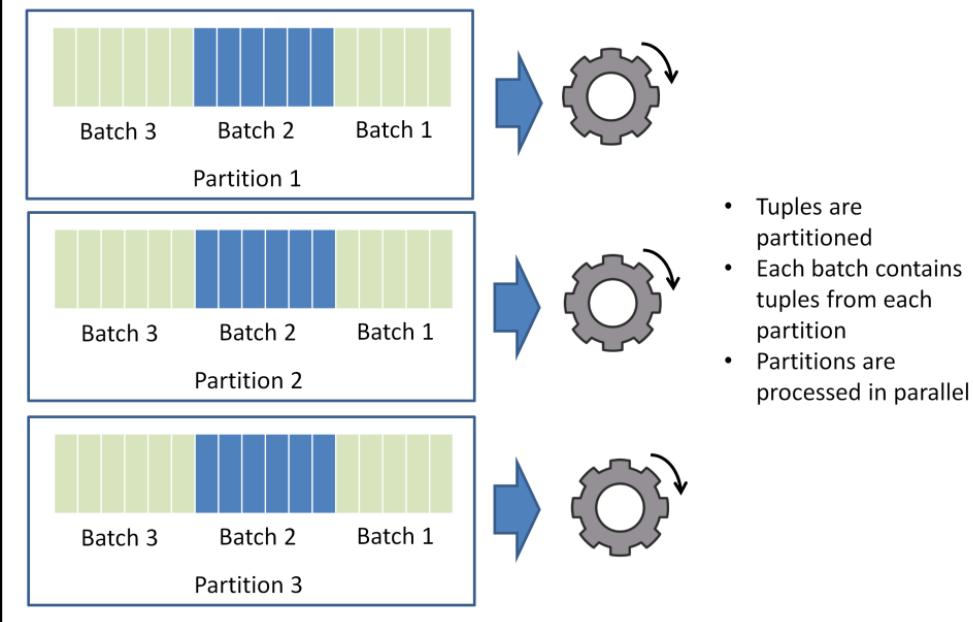


Processing one tuple at a time is of course highly inefficient. A better approach is to process the tuples as discrete batches. This is known as micro-batch stream processing.

The batches are processed in order, and each batch has a unique ID that is always the same on every replay. Because many tuples are processed per iteration rather than just one, the processing can be made scalable by parallelizing it. Batches must be processed to completion before moving on to the next batch.

Retaking our global counting example, we now store the count along with the latest batch ID involved in updating the count. Now suppose that after the state in the database is updated, something fails in the stream processor and the message that the batch was finished is never received. The stream processor will timeout the batch and retry the batch. When it comes time to update the database, it sees that the state has already been updated to include that batch. So rather than increment the count, it does nothing and moves on to the next batch.

Micro-batch processing topologies



Just like how MapReduce and one-at-a-time stream processing partition data and process each partition in parallel, the same is done with micro-batched processing. Processing a batch of words looks like the figure. A single batch includes tuples from all partitions in the incoming stream.

The parallel processing of the different partitions have an impact when you have to save your state. This is further discussed on the next slide.

Stateful computation

- State across batches
 - global count, word count, top-3 words ...
- Two caveats:
 - ensure idempotent processing as failures may occur anytime
 - e.g. using the batch ID
 - avoid race conditions when accessing database
 - e.g. by having only one worker counting a specific word

		Count	Batch ID
Apple	→	15	3
Pear	→	21	18
Banana	→	11	3

When a failure occurs, the entire batch will be replayed. This poses no problem for batch-local computation: computation that occurs solely within a single batch. Examples are repartitioning a stream according to a given field (e.g. word) or simply counting the number of tuples in a batch.

Sometimes you need to keep state across all batches, an example is updating a global count or per-word count over all batches. This is where you have to be really careful about how you do update the state (e.g. in a database).

A first caveat is that you should avoid race conditions between workers when updating a particular field in the database. In the canonical word counting example, this means that you must first repartition the tuples by the field containing the word itself, so that only one worker will update the state for that specific word.

A second caveat is that you should ensure that all state updates are idempotent. The trick of storing the batch ID with the state is a particular way of achieving this.

Let's consider a failure scenario. Suppose a machine dies in the cluster while a batch is being processed, and only some partitions succeeded in updating the database. Some words will have counts reflecting the current batch, and others won't be updated yet. When the batch is replayed, the words that have state including the current batch ID won't be updated, whereas the words that haven't been updated yet will be updated like normal.

Further Reading

- N. Garg, Learning Apache Kafka (2nd edition)
- S. Saxena, Real-time Analytics with Storm and Cassandra

BIG DATA ARCHITECTURES

Data system requirements

Query = function(all data)



latency



timeliness



accuracy

Challenges



Machines will break



Humans will make errors

A data system answers questions based on data you have seen in the past. Or put more formally: a data system computes queries that are functions of all the data you have ever seen (query = function(all data)). There are a number of properties you are concerned about with your queries:

- **Latency** – The time it takes to run a query. In many cases, your latency requirements will be very low – on the order of milliseconds. Other times it is okay for a query to take a few seconds. When doing ad hoc analysis, your latency requirements are often very lax, even on the order of hours.
- **Timeliness** – How up-to-date the query results are. A completely timely query takes into account all data ever seen in the past, whereas a less timely query may not include results from the recent minutes or hours.
- **Accuracy** – In many cases, in order to make queries performant or scalable, you must make approximations in your query implementations.

A huge part of building data systems is making them fault tolerant. You have to plan for how your system will behave when you encounter machine failures. Oftentimes this means making trade-offs with the preceding properties. For example, there is a fundamental tension between latency and timeliness. The CAP theorem shows that under partitions, a system can either be consistent (queries take into account all previous written data) or available (queries are answered at the moment). Consistency is just a form of timeliness, and availability just means the latency of the query is bounded. An eventually consistent system chooses latency over timeliness

(queries are always answered, but may not take into account all prior data during failure scenarios).

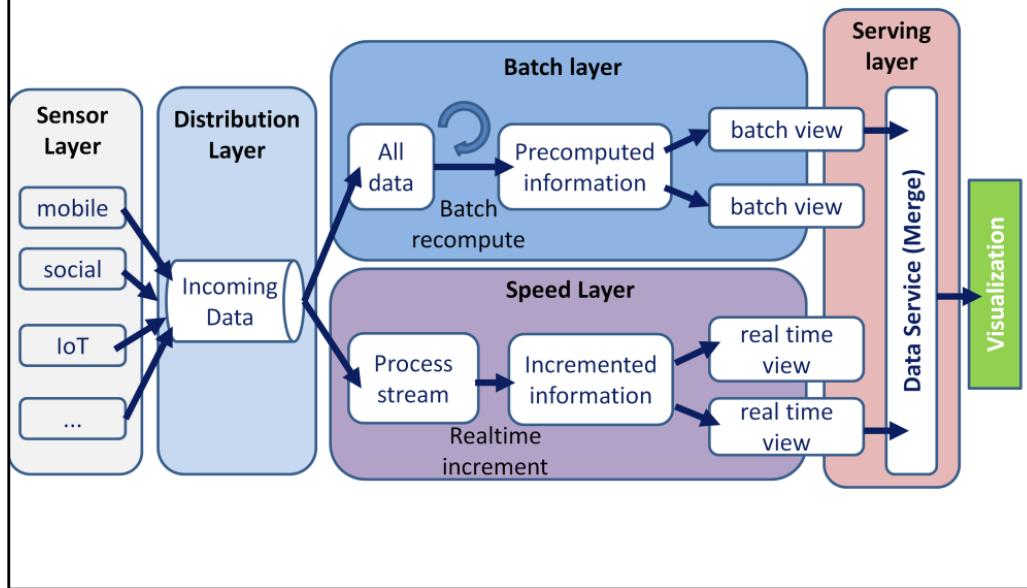
Because data systems are dynamic, changing systems built by humans and with new features and analyses deployed all the time, humans are an integral part of any data system. Humans can and will fail. We saw earlier how mutability is fundamentally not human-fault tolerant. If a human can mutate data, then a mistake can mutate data. The only solution is to make your core data *immutable*, with the only write operation allowed being appending new data to your ever-growing set of data. You can set permissions on your core data to disallow deletes and updates, making your system far more robust.

This leads us to the basic model of data systems:

- A master dataset consisting of an ever-growing set of data
- Queries as functions that take in the entire master dataset as input.

The Lambda architecture is one possible implementation of this basic data model, and is discussed in the next slide.

Lambda architecture



Computing arbitrary functions on an arbitrary dataset in real time is a daunting task. There is no single tool that provides a complete solution. Instead, you have to use a variety of tools and techniques to build a complete Big Data system. The main idea of the Lambda architecture is to build Big Data systems as a series of layers. Each layer satisfies a subset of the properties and builds upon the functionality provided by the layers beneath it.

Everything starts from the $query = function(all\ data)$ equation. Running queries on the fly would take a huge amount of resources and be unreasonably expensive. Imagine having to read a petabyte dataset every time you wanted to answer the query of someone's current location. The most obvious alternative approach is to precompute the query function, which we call batch views. The batch layer needs to be able to do two things: store an immutable, constantly growing master dataset (a very large list of records) and compute arbitrary functions on that dataset. This type of processing is best done using batch-processing systems like Hadoop. Conceptually, the batch layer runs in a `while(true)` loop and continuously recomputes the batch views from scratch, since there is a continuous infeed of data from the sensor layer.

The batch layer emits batch views as the result of its functions. The next step is to load the views somewhere so that they can be queried. This is where the serving layer comes in. The serving layer is a specialized distributed database that loads in a batch view and makes it possible to do random reads on it. When new batch views

are available, the serving layer automatically swaps those in so that more up-to-date results are available. A serving layer database supports batch updates and random reads. Most notably, it doesn't need to support random writes. This is a very important point, as random writes cause most of the complexity in databases. By not supporting random writes, these databases are extremely simple. That simplicity makes them robust, predictable, easy to configure and easy to operate. ElephantDB is one example, built from only a few thousands lines of code.

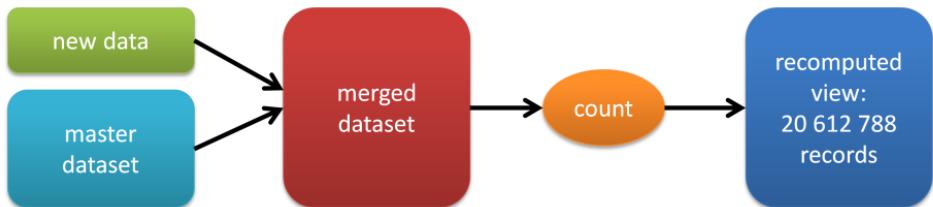
The batch and serving layers support arbitrary queries on an arbitrary dataset with the trade-off that queries will be out of date by a few hours. It takes a new piece of data a few hours to propagate through the batch layer into the serving layer where it can be queried. To have a fully real-time data system that allows to compute arbitrary functions on arbitrary data in real-time, we need to compensate for the data that came in while the batch precomputation was running. This is the purpose of the speed layer.

The speed layer is similar to the batch layer in that it produces views based on the data it receives. One big difference is that the speed layer only looks at recent data, whereas the batch layer looks at all the data at once. Another big difference is that in order to achieve the smallest latency possible, the speed layer doesn't look at all the new data at once. Instead, it updates the realtime views as it receives new data instead of recomputing the views from scratch like the batch layer does. The speed layer does incremental computation instead of the recomputation done in the batch layer.

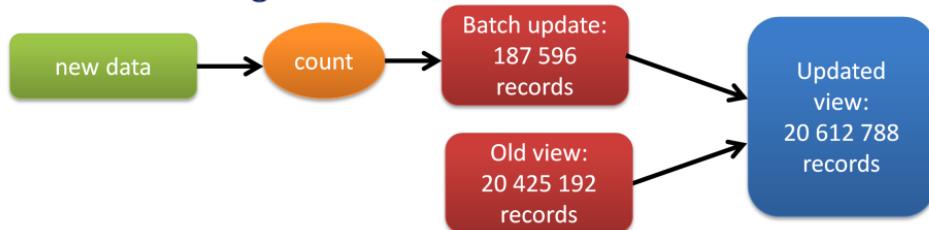
You resolve queries by looking at both the batch and realtime views and merging the results together. The speed layer uses databases that support random reads and random writes. Because these databases support random writes, they are orders of magnitude more complex than the databases you use in the serving layer, both in terms of implementation and operation.

Batch layer

Recomputation algorithm:



Incremental algorithm:



Because your master dataset is continually growing, you must have a strategy for updating your batch views when new data becomes available. You could choose a recomputation algorithm, throwing away the old batch views and recomputing functions over the entire master dataset. Alternatively, an incremental algorithm will update the views directly when new data arrives.

As a basic example, consider a batch view containing the total number of records in your master dataset. A recomputation algorithm would update the count by first appending the new data to the master dataset and then counting all records from scratch. An incremental algorithm, on the other hand, would count the number of new data records and add it to the existing count.

You might be wondering why you would ever use a recomputation algorithm when you can use a vastly more efficient incremental algorithm instead. The key trade-offs between the two approaches are performance, human-fault tolerance and the generality of the algorithm. We'll discuss both types of algorithms in regard to each of these issues.

Performance

Resources required to update batch view with new data



Size of the batch views produced

Example 1: average page views for each URL of a domain

URL	Avg. page views
foo.com	46.22
foo.com/blog	44.18
foo.com/about	2.24
foo.com/faq	7.36

URL	Avg. page views	Total count
foo.com	46.22	1543
foo.com/blog	44.18	1475
foo.com/about	2.24	75
foo.com/faq	7.36	245

Example 2: Number of unique visitors for each URL

URL	# Unique visitors
foo.com	2217
foo.com/blog	1899
foo.com/about	524
foo.com/faq	413

URL	# Unique visitors	Visitor IDs
foo.com	2217	1, 4, 5...
foo.com/blog	1899	2, 3, 5...
foo.com/about	524	3, 6, 7...
foo.com/faq	413	12, 17 ...

There are two aspects to the performance of a batch-layer algorithm: the amount of resources required to update a batch view with new data, and the size of the batch views produced. An incremental algorithm always uses significantly less resources to update a view because it uses new data and the current state of the batch view to perform an update. For a task such as computing URL pageviews over time, the view will be significantly smaller than the master dataset because of the aggregation. A recomputation algorithm looks at the entire master dataset, so the amount of resources needed for an update can be multiple orders of magnitude higher than an incremental algorithm. But the size of the batch view for an incremental algorithm can be significantly larger than the corresponding batch view for a recomputation algorithm. This is because the view needs to be formulated in such a way that it can be incrementally updated. We demonstrate this trade-off through two separate examples.

First, suppose you need to compute the average number of pageviews for each URL within a particular domain. The batch view generated by a recomputation algorithm would contain a map from each URL to its corresponding average. But this isn't suitable for an incremental algorithm, because updating the average incrementally requires that you also know the number of records used for computing the previous average. An incremental view would therefore store both the average and the total count for each URL, increasing the size of the incremental view over the recomputation-based view by a constant factor.

In other scenarios, the increase in the batch view size for an incremental algorithm is much more severe. Consider a query that computes the number of unique visitors for each URL. A recomputation view only requires a map from the URL to the unique count. In contrast, an incremental algorithm only examines the new pageviews, so its view must contain the full set of visitors for each URL so it can determine which records in the new data correspond to return visits. As such, the incremental view could potentially be as large as the master dataset.

The batch view generated by an incremental algorithm isn't always this large, but it can be far larger than the corresponding recomputation-based view.

Fault tolerance and generality

- Human-fault tolerance
 - recomputation: fix algorithm and redeploy code
 - incremental: which records were affected?
- Generality of the algorithms
 - incremental algorithms must often be tailored
 - incremental algorithms shift complexity to on-the-fly computations

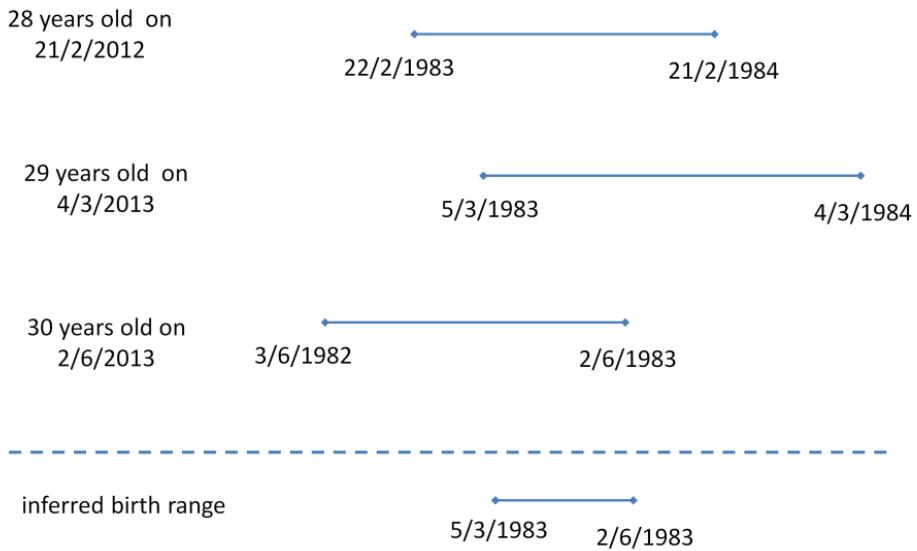
The other two aspects to consider when evaluating incremental vs recomputation algorithms are fault tolerance and generality of the algorithms.

Recomputation algorithms are inherently human-fault tolerant. Consider as an example a batch-layer algorithms that computes a global count of the number of records in the master dataset. Now suppose you make a mistake and deploy an algorithm that increments the global count for each record by two instead of by one. If your algorithm is recomputation-based, all that is required to fix the algorithm and redeploy the code: your batch view will be correct the next time the batch layer runs, since the algorithm recomputes the batch view from scratch. But if your algorithm is incremental, then correcting your view is not so simple. The only option is to identify the records that were overcounted, determine how many times each one was overcounted, and then correct the count for each affected record.

Although incremental algorithms can be faster to run, they must often be tailored to address the problem at hand. Incremental algorithms also shift complexity to on-the-fly computations, which increases their latency. Suppose you have improved a normalization component. You no longer only map “Ghent”, “Zwijnaarde”, “Ledeberg” to “Ghent, Belgium” but now you make a distinction between these city regions. In an incremental algorithm, this means that your batch view has to keep every name that was ever mapped to “Ghent, Belgium” to apply this update. Moreover, you will have to renormalize each city name every time a query is performed. This increases the latency of the on-the-fly component and could very

well take too long for your application's requirements.

Example: coping with messy data



Let us consider another example where the choice between incremental and recomputation algorithms is more difficult: the “birthday inference” problem. Imagine you are writing a web crawler that collects people’s ages from their public profiles. The profile does not contain a birthday, but only what the person’s age is at the moment you crawled that web page. Given this raw data of [age, timestamp] pairs, your goal is to deduce the birthday of each person.

The idea of the algorithm is illustrated in the slide. Imagine you crawl the profile of Tom on January 4, 2012 and see his age is 23. Then you crawl his profile again on January 11, 2012 and see his age is 24. You can deduce that his birthday happened sometime between those dates. Likewise, if you crawl the profile of Jill on October 20, 2013 and see she is 43, and then crawl it again on November 4, 2013 and see she is still 43, you know her birthday is not between those dates. The more age samples you have, the better you can infer that someone’s birthday is within a small range of dates.

In the real world, data can get messy. Someone may have incorrectly entered their birthday and then changed it a later date. This may cause your age inference algorithm to fail because every day of the year has been eliminated as possible birthday. You might modify your algorithm to search for the smallest number of age samples it can ignore to produce the smallest range of possible birthdays. The algorithm might prefer to use recent age samples over older age samples.

If you implement your birthday-inference batch layer using recomputation, it's easy. Your algorithm can look at all age samples for a person at once and do everything necessary to deal with messy data and emit a single range of dates as output. But incrementalizing the algorithm is much trickier: it is hard to see how you can deal with the messy data problem without having access to the full range of age samples.

Partial recomputation

Avoid full recomputation, but still use entire master dataset

1. For the new batch of data, find all people who have a new age sample
2. Retrieve all age samples from the master dataset for all people in step 1
3. Recompute the birthdays for all people in step 1 using the age samples from step 2 and the age samples in the new batch
4. Merge the newly computed birthdays into the existing server layer views



Avoid repartitioning (group-by, join)
Map-only job: iterate over dataset and emit only relevant data

Partial recomputation is an alternative that blurs the line between incrementalization and recomputation and gets you the best of both worlds. In the example of the previous slide: if a person has no new age samples since the last time the batch layer ran, then the inferred birthday for that person will not change at all. The idea is to do the following:

1. For the new batch of data, find all people who have a new age sample
2. Retrieve all age samples from the master dataset for all people in step 1
3. Recompute the birthdays for all people in step 1 using the age samples from step 2 and the age samples in the new batch
4. Merge the newly computed birthdays into the existing server layer views

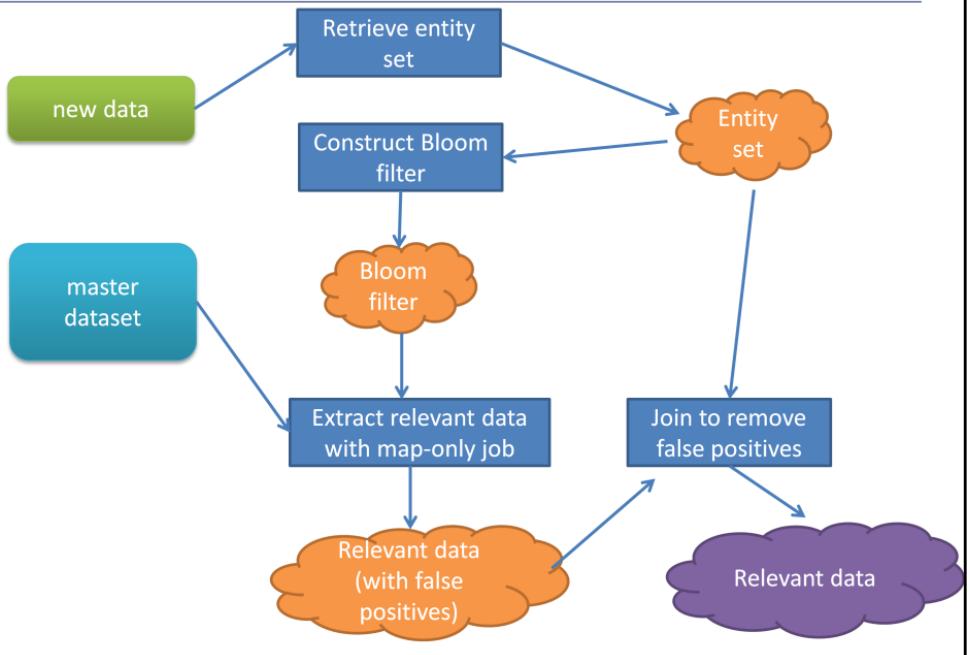
This alternative implementation is not fully incremental because it still makes use of the master dataset. But it avoids most of the cost of a full recomputation by ignoring anyone who hasn't changed in the latest set of data.

The key idea to partial recomputation is to retrieve all the relevant data for the entities that changes, run a normal recompute algorithm on the retrieved data plus the new data, and then merge those results into the existing views. The nice thing about partial recomputes is that they can be implemented very efficiently. The most expensive step – looking over the entire master dataset to find relevant data – can be done relatively cheaply.

The key to making it efficient is to avoid having to repartition the entire master

dataset, as this is the most expensive part of batch algorithms. For example, repartitioning happens whenever you do a group-by operation or a join. Partitioning involves serialization/deserialization, network transfer and possibly buffering on disk. In contrast, operations that don't require partitioning can quickly scan through the data and operate on each piece of data as it's seen. Retrieving relevant data for a partial recompute can be done using the latter method.

Bloom join



The first step to retrieving relevant data is to construct a set of all the entities for which you need relevant data. You then scan over the entire master dataset and only emit data for those entities that exist in the set (each task would have a copy of that set). In a batch-processing system like Hadoop, this would correspond to a map-only job (with a trivial reducer).

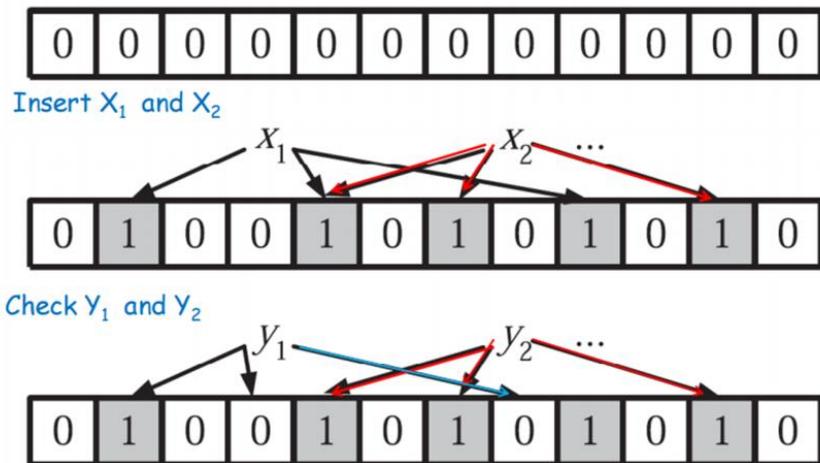
You are limited by memory, so your set can only be so big. But a data structure called a *Bloom filter* can make this work for much larger sets of entities. A Bloom filter is a compact data structure that represents a set of elements and allows you to ask if it contains an element. A Bloom filter is much more compact than a set, but as a trade-off, query operations on it are probabilistic. A Bloom filter will sometimes incorrectly tell you that an element exists in the set, but it will never tell you an element that was added to the set is not in the set. So a Bloom filter has false positives but no false negatives.

If you use a Bloom filter to retrieve relevant data from the master dataset, you will filter out the vast majority of the master dataset. Due to the false positives, though, some data will be emitted that you didn't want to retrieve. You can then do a join between the retrieved data and the list of desired entities to filter out the false positives. A join requires a partitioning, but because the vast majority of the master dataset was already filtered out, getting rid of the false positives is not an expensive operation.

Bloom filter

- array of **m** bits representing a set $S = \{x_1, x_2 \dots x_n\}$ of **n** elements
 - initialized to 0
- **k** independent hash functions $h_1 \dots h_k$ with range {1, 2 ... m}
 - assume each hash function maps each item in the universe to a random number *uniformly* over the range {1... m}
- for each element x in S , the bit $h_i(x)$ in the array is set to 1, for $1 \leq i \leq k$
 - a bit in the array may be set to 1 multiple times for different elements

Bloom filter example



Standard Bloom filter (cont.)

- To check membership of y in S , check whether $h_i(y)$, $1 \leq i \leq k$ are all set to 1
 - If not: y is definitely not in S
 - Else, we conclude that y is in S , but sometimes this conclusion is wrong (false positive)
- For many applications, false positives are acceptable as long as the probability of a false positive is small enough

Dimensioning the filter

- Optimal number of hash functions k

$$k = \frac{m}{n} \ln 2$$

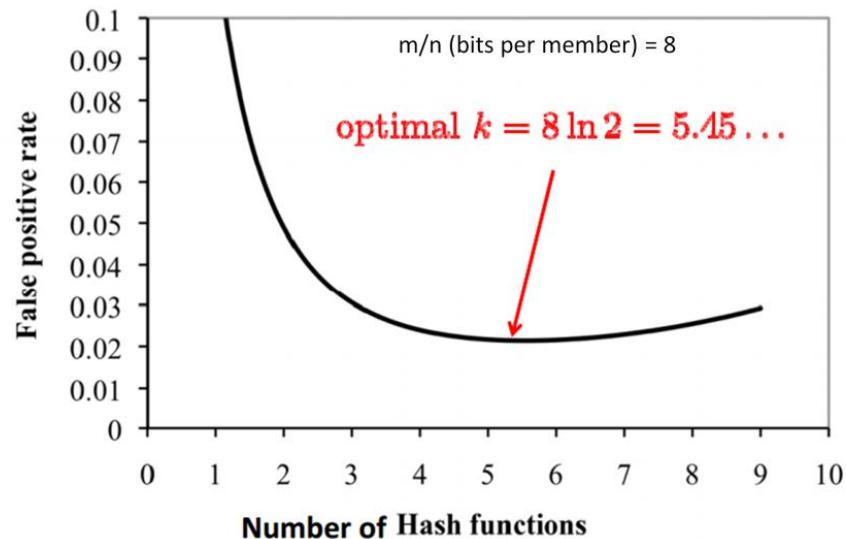
- Required number of bits m , given

- n (number of inserted elements)
- desired false positive probability p
- the optimal value of k (see above) is used

$$m = -\frac{n \ln p}{(\ln 2)^2}$$

proportional to size of
input set n !

False positive rate vs. k



Kappa architecture



Questioning the Lambda Architecture

The Lambda Architecture has its merits, but alternatives are worth exploring.

by Jay Kreps | @jaykreps | +Jay Kreps | Comments: 19 | July 2, 2014



Recognizes **advantages** of Lambda architecture:

- Retain input data unchanged
 - Modeling data transformation as a series of materialized stages from original input
- Highlights problem of reprocessing data
 - application evolves, bug fixes

Cons of Lambda architecture:

- Different and diverging programming paradigms
 - Very different code for MapReduce and Storm
 - also involves debugging and interaction with other products

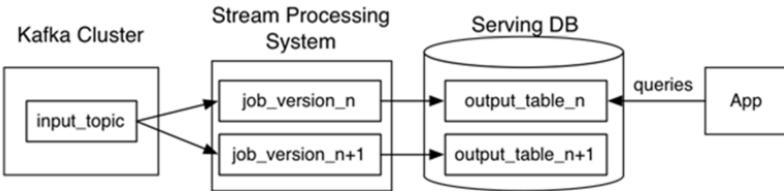
In July 2014, Jay Kreps coined the term “Kappa Architecture” in a blog post (<http://radar.oreilly.com/2014/07/questioning-the-lambda-architecture.html>). At that time, Jay Kreps was working on the big data systems at LinkedIn.

In his blog post, he first acknowledges some merits of the Lambda architecture:

- It highlights the importance of keeping the input data unchanged. Jay Kreps also believes that modeling data transformation as a series of materialized stages from an original input has a lot of merit.
- It highlights the problem of reprocessing data: processing input data over again to re-derive output. Reprocessing is needed because applications evolve (e.g. you want to compute new output fields) or bugs have to be fixed

But he also refers to some problems he sees with the Lambda architecture. Programming in distributed frameworks like Storm and Hadoop is complex. Inevitably, code ends up being specifically engineered toward the framework it runs on. In the Lambda architecture, you need to code in both (types of) frameworks. In addition, the operational burden of running and debugging two systems is going to be very high.

Kappa architecture



- Solely uses stream processing (also for reprocessing)
- Abstraction to DAGs is already used in data warehouses and MapReduce frameworks

- 1) Use a system to retain the full log of the data and that allows for multiple subscribers
- 2) Start a second instance of the stream processing job
 - 1) starts processing from the beginning of the retained data,
 - 2) direct this output data to a new output table
- 3) When the second job has caught up, switch the application to read from the new table
- 4) Stop the old version of the job, and delete the old output table.

The crux of the Kappa architecture is that it uses stream processing systems to handle the reprocessing when code changes. The argument is that stream processing systems already have a notion of parallelism; so why not just handle reprocessing by increasing the parallelism and replaying history very, very fast.

Many people have a notion that stream processing is inherently something that computes results off some ephemeral streams and then throws all the underlying data away. But there is no reason this should be true. The fundamental abstraction in stream processing is data flow DAGs (directed acyclic graphs), which are exactly the same underlying abstraction in traditional data warehouse (a la Volcano) as well as being the fundamental abstraction in the MapReduce successor Tez. Stream processing is just a generalization of this data-flow model that exposes checkpointing of intermediate results and continual output to the user.

So, how can the reprocessing be done directly from the stream processing job?

- 1) Use Kafka (or some other system) to retain the full log of the data you want to be able to reprocess and that allows for multiple subscribers
- 2) When you want to do the reprocessing, start a second instance of your stream processing job that starts processing from the beginning of the retained data, but direct this output data to a new output table
- 3) When the second job has caught up, switch the application to read from the new table
- 4) Stop the old version of the job, and delete the old output table.

Polyglot persistence/processing



The term polyglot is borrowed and redefined for big data as a set of applications that use several core database technologies, and this is the most likely outcome of your implementation planning. The official definition of *polyglot* is “someone who speaks or writes several languages.” It is going to be difficult to choose one persistence style no matter how narrow your approach to big data might be.

Further reading

- S. Lam, Bloom Filters [online - 1]
- G. Schmutz, Big Data and Fast Data – Lambda Architecture in Action [online - 2]
- J. Kreps, Questioning the Lambda Architecture [online - 3]
- N. Marz and J. Warren, Big Data [book]

[1] <http://www.cs.utexas.edu/users/lam/386p/slides/Bloom%20Filters.pdf>

[2] <http://www.slideshare.net/gschmutz/big-data-and-fast-data-lambda-architecture-in-action?related=1>

[3] <http://radar.oreilly.com/2014/07/questioning-the-lambda-architecture.html>