

# CC5212-1

PROCESAMIENTO MASIVO DE DATOS

OTOÑO 2021

## Lecture 9

### NoSQL: Overview

Aidan Hogan

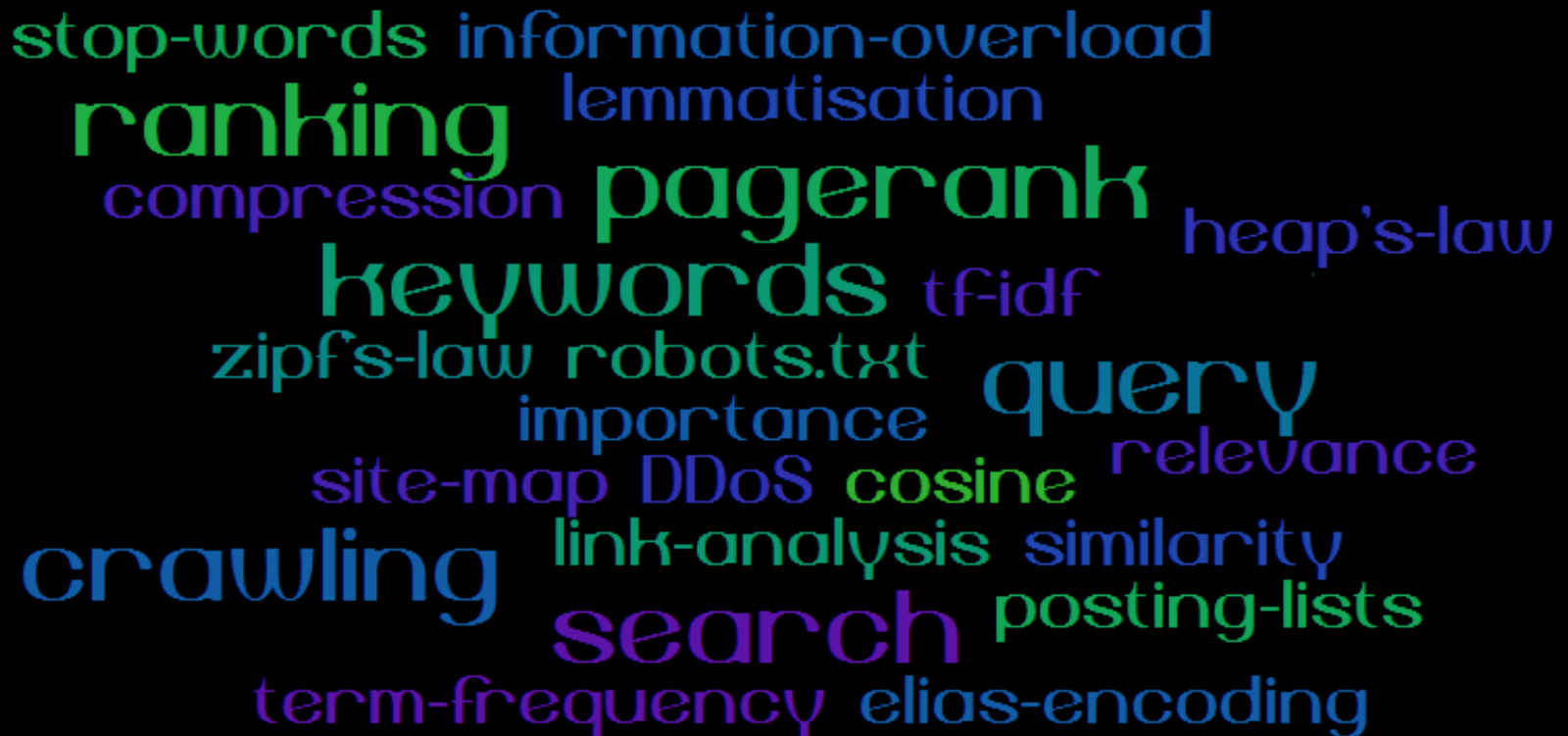
[aidhog@gmail.com](mailto:aidhog@gmail.com)

# Hadoop/MapReduce/Pig/Spark: Processing Un/Structured Information



Information Retrieval:

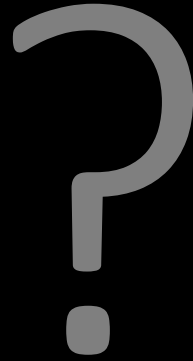
Storing Unstructured Information



A word cloud of information retrieval concepts. The words are arranged in a roughly circular pattern, with some words being significantly larger than others. The colors of the words range from light blue to dark purple. The words include: stop-words, information-overload, ranking, lemmatisation, compression, pagerank, heap's-law, keywords, tf-idf, zipf's-law, robots.txt, query, importance, relevance, site-map, DDoS, cosine, crawling, link-analysis, similarity, search, posting-lists, term-frequency, and elias-encoding.

stop-words information-overload  
ranking lemmatisation  
compression pagerank heap's-law  
keywords tf-idf  
zipf's-law robots.txt query  
importance relevance  
site-map DDoS cosine  
crawling link-analysis similarity  
search posting-lists  
term-frequency elias-encoding

Storing Structured Information??



BIG DATA:

STORING STRUCTURED INFORMATION

# Relational Databases







# Relational Databases: One Size Fits All?

## “One Size Fits All”: An Idea Whose Time Has Come and Gone

Michael Stonebraker  
*Computer Science and Artificial  
Intelligence Laboratory, M.I.T., and  
StreamBase Systems, Inc.*  
stonebraker@csail.mit.edu

Uğur Çetintemel  
*Department of Computer Science  
Brown University, and  
StreamBase Systems, Inc.*  
ugur@cs.brown.edu

### Abstract

*The last 25 years of commercial DBMS development can be summed up in a single phrase: “One size fits all”. This phrase refers to the fact that the traditional DBMS architecture (originally designed and optimized for business data processing) has been used to support many data-centric applications with widely varying characteristics and requirements.*

*In this paper, we argue that this concept is no longer applicable to the database market, and that the commercial world will fracture into a collection of independent database engines, some of which may be unified by a common front-end parser. We use examples from the stream-processing market and the data-warehouse market to bolster our claims. We also briefly discuss other markets for which the traditional architecture is a poor fit and argue for a critical rethinking of the current factoring of systems services into products.*

of multiple code lines causes various practical problems, including:

- *a cost problem*, because maintenance costs increase at least linearly with the number of code lines;
- *a compatibility problem*, because all applications have to run against every code line;
- *a sales problem*, because salespeople get confused about which product to try to sell to a customer; and
- *a marketing problem*, because multiple code lines need to be positioned correctly in the marketplace.

To avoid these problems, all the major DBMS vendors have followed the adage “put all wood behind one arrowhead”. In this paper we argue that this strategy has failed already, and will fail more dramatically off into the future.

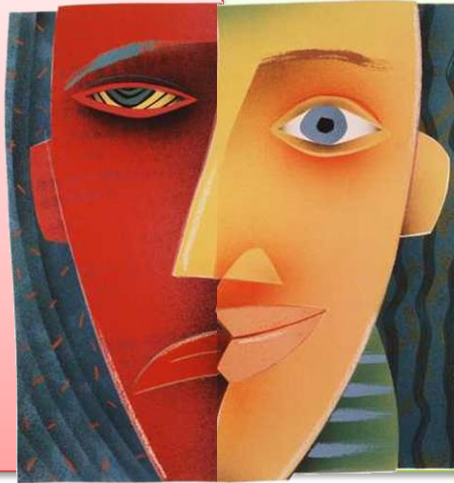
The rest of the paper is structured as follows. In Section 2, we briefly indicate why the single code-line strategy has failed already by citing some of the key characteristics of the data warehouse market. In Section



# SQL

Difficult to optimise

Difficult to distribute



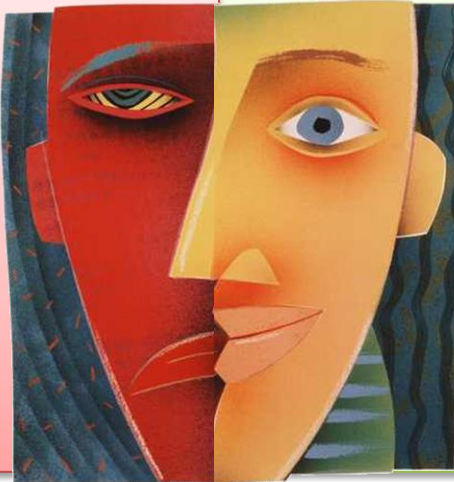
Declarative language

Expressive

# ACID

Costly to implement

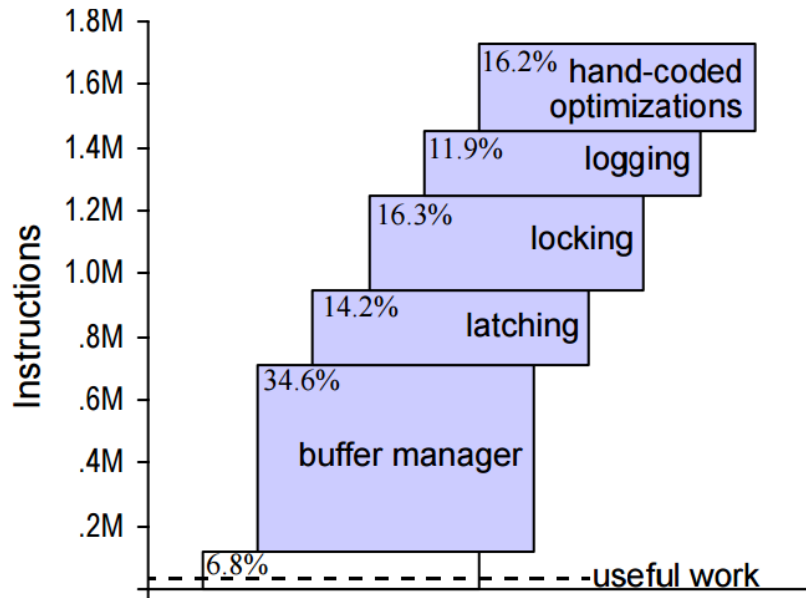
Difficult to distribute



Guarantees correct behaviour

Support transactions

# Transactional overhead: the cost of ACID



- 640 transactions per second for system with full transactional support (ACID)
- 12,700 transactions per second for system without logs, transactions or lock scheduling

## OLTP Through the Looking Glass, and What We Found There

Stavros Harizopoulos  
HP Labs  
Palo Alto, CA  
stavros@hp.com

Daniel J. Abadi  
Yale University  
New Haven, CT  
dna@cs.yale.edu

Samuel Madden      Michael Stonebraker  
Massachusetts Institute of Technology  
Cambridge, MA  
{madden, stonebraker}@csail.mit.edu

### ABSTRACT

Online Transaction Processing (OLTP) databases include a suite of features — disk-resident B-trees and heap files, locking-based concurrency control, support for multi-threading — that were optimized for computer technology of the late 1970's. Advances in modern processors, memories, and networks mean that today's computers are vastly different from those of 30 years ago, such that many OLTP databases will now fit in main memory, and most OLTP transactions can be processed in milliseconds or less. Yet database architecture has changed little.

### 1. INTRODUCTION

Modern general purpose online transaction processing (OLTP) database systems include a standard suite of features: a collection of on-disk data structures for table storage, including heap files and B-trees, support for multiple concurrent queries via locking-based concurrency control, log-based recovery, and an efficient buffer manager. These features were developed to support transaction processing in the 1970's and 1980's, when an OLTP database was many times larger than the main memory, and when the computers that ran these databases cost hundreds of thousands to

# ALTERNATIVES TO RELATIONAL DATABASES FOR BIG DATA?

# NoSQL

Anybody know anything about NoSQL?

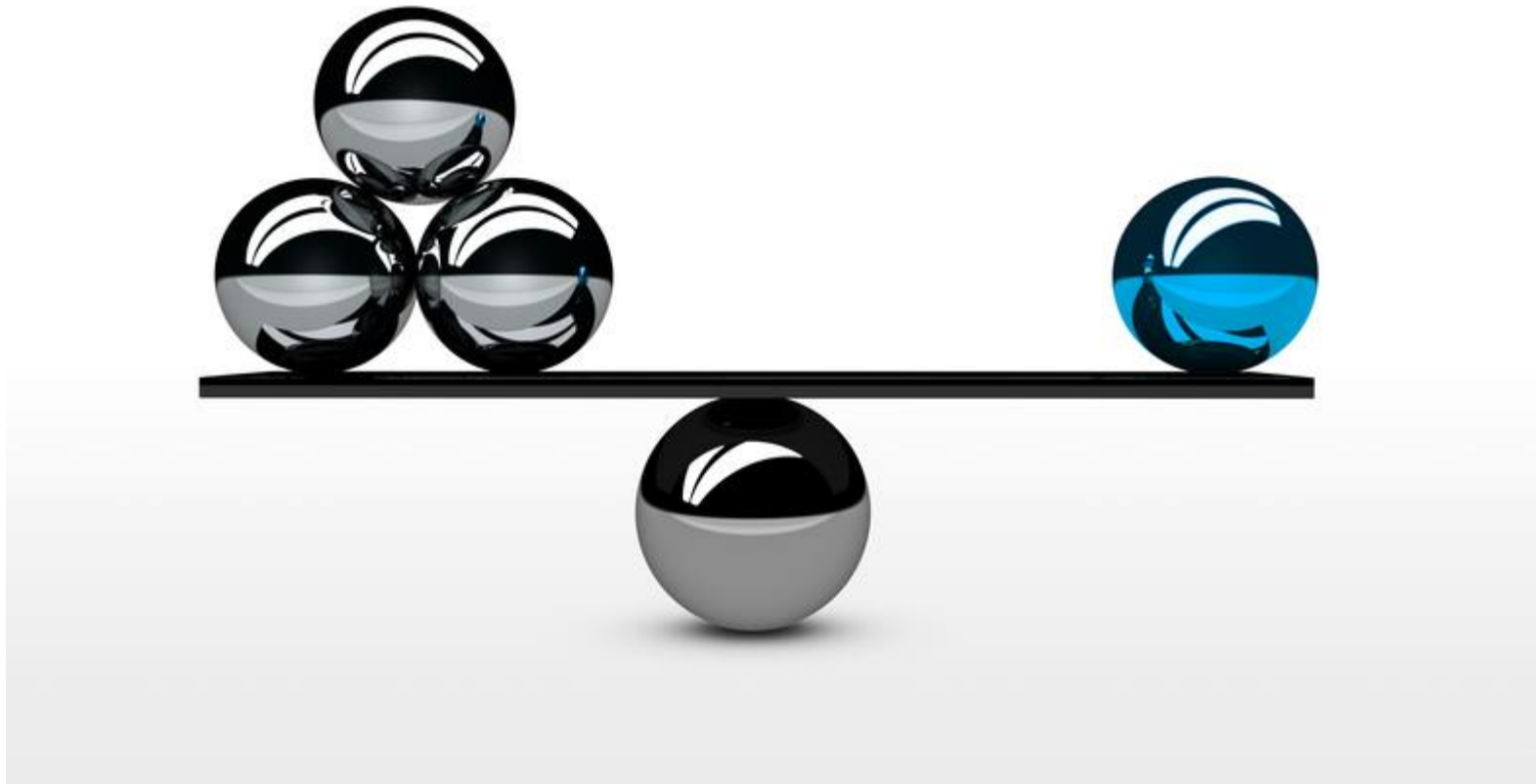


**N**ot  
**O**nly **SQL**



May 2021	Rank		DBMS	Database Model	Score		
	Apr 2021	May 2020			May 2021	Apr 2021	May 2020
1.	1.	1.	Oracle	Relational, Multi-model	1269.94	-4.98	-75.50
2.	2.	2.	MySQL	Relational, Multi-model	1236.38	+15.69	-46.26
3.	3.	3.	Microsoft SQL Server	Relational, Multi-model	992.66	-15.30	-85.64
4.	4.	4.	PostgreSQL	Relational, Multi-model	559.25	+5.73	+44.45
5.	5.	5.	MongoDB	Document, Multi-model	481.01	+11.04	+42.02
6.	6.	6.	IBM Db2	Relational, Multi-model	166.66	+8.88	+4.02
7.	7.	8.	Redis	Key-value, Multi-model	162.17	+6.28	+18.69
8.	8.	7.	Elasticsearch	Search engine, Multi-model	155.35	+3.18	+6.23
9.	9.	9.	SQLite	Relational	126.69	+1.64	+3.66
10.	10.	10.	Microsoft Access	Relational	115.40	-1.33	-4.50
11.	11.	11.	Cassandra	Wide column	110.93	-3.92	-8.22
12.	12.	12.	MariaDB	Relational, Multi-model	96.69	+0.32	+6.61
13.	13.	13.	Splunk	Search engine	92.11	+3.62	+4.36
14.	14.	14.	Hive	Relational	76.19	-2.31	-5.35
15.	15.	23.	Microsoft Azure SQL Database	Relational, Multi-model	70.46	-1.39	+27.70
16.	16.	16.	Amazon DynamoDB	Multi-model	70.07	-0.66	+5.35
17.	17.	15.	Teradata	Relational, Multi-model	69.98	-0.57	-3.91
18.	18.	20.	SAP HANA	Relational, Multi-model	52.75	-0.69	+2.22
19.	20.	21.	Neo4j	Graph	52.23	+1.19	+2.47
20.	21.	18.	Solr	Search engine, Multi-model	51.19	+0.59	-1.39
21.	19.	17.	SAP Adaptive Server	Relational, Multi-model	49.96	-1.70	-4.03
22.	22.	19.	FileMaker	Relational	46.73	+0.32	-4.23
23.	23.	22.	HBase	Wide column	43.24	-0.92	-6.48
24.	24.	26.	Google BigQuery	Relational	37.63	+2.05	+10.04
25.	25.	24.	Microsoft Azure Cosmos DB	Multi-model	34.71	+1.19	+4.03

# NoSQL: features vs. scale/performance

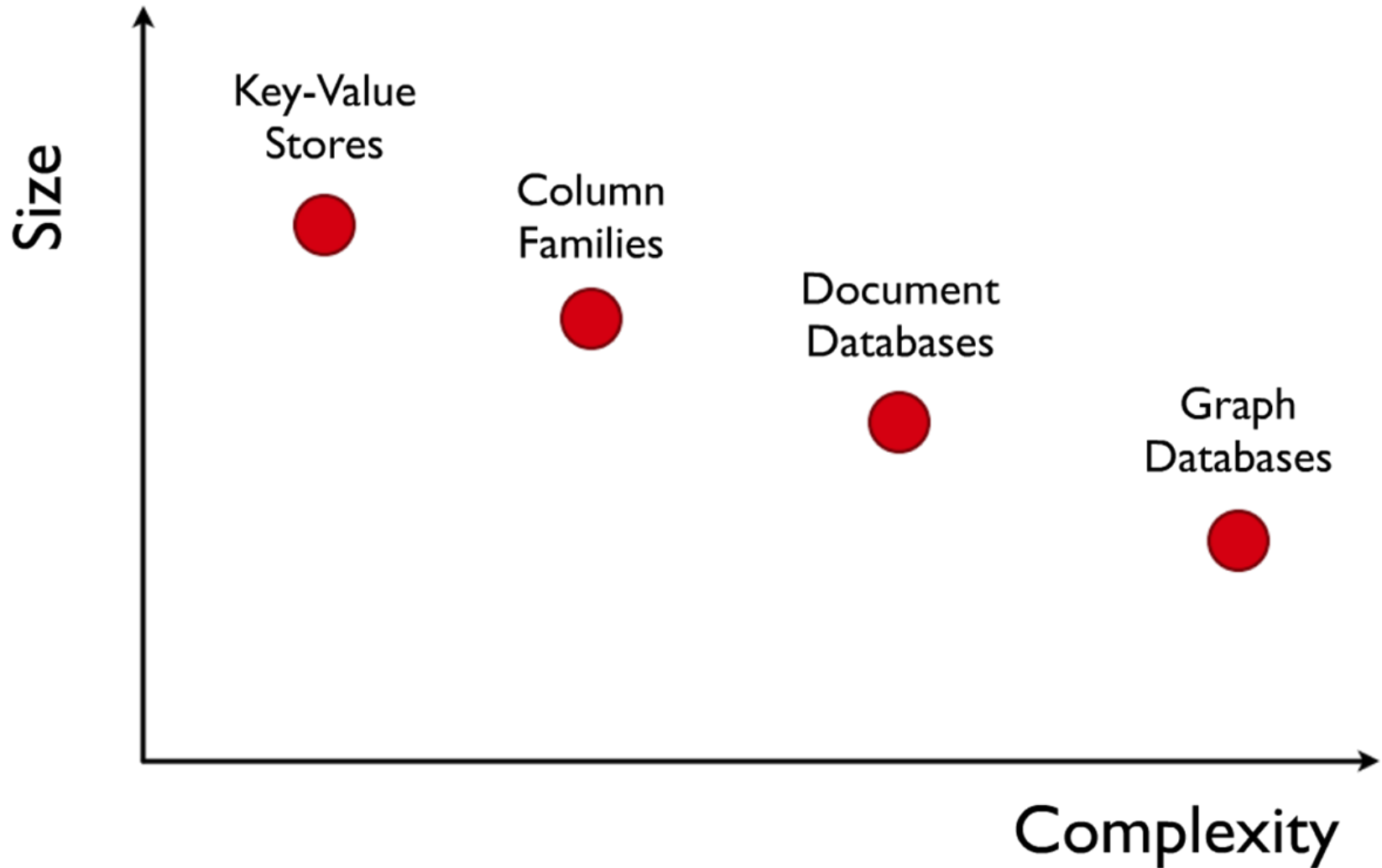




# NoSQL: common characteristics

- Often **distributed**
- Often **simpler** languages than SQL
- **Different flavours** (for different scenarios)

# NoSQL: four main flavours



# LIMITATIONS OF DISTRIBUTED COMPUTING: CAP THEOREM

# What is CAP?

Three *guarantees* a distributed sys. could make

## 1. Consistency:

- All nodes have a consistent view of the system

## 2. Availability:

- Every read/write is acted upon

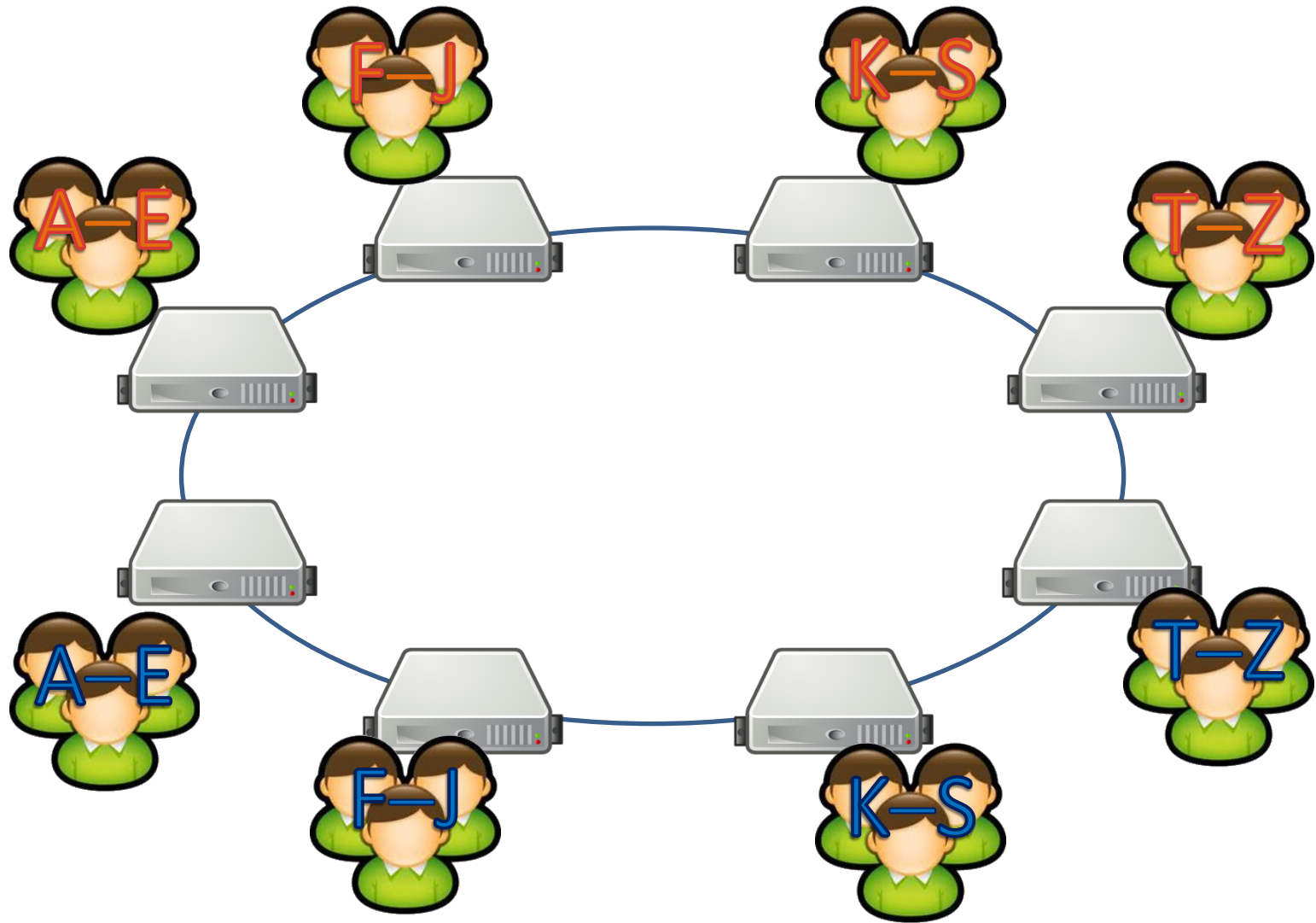
## 3. Partition-tolerance:

- The system works even if messages are lost

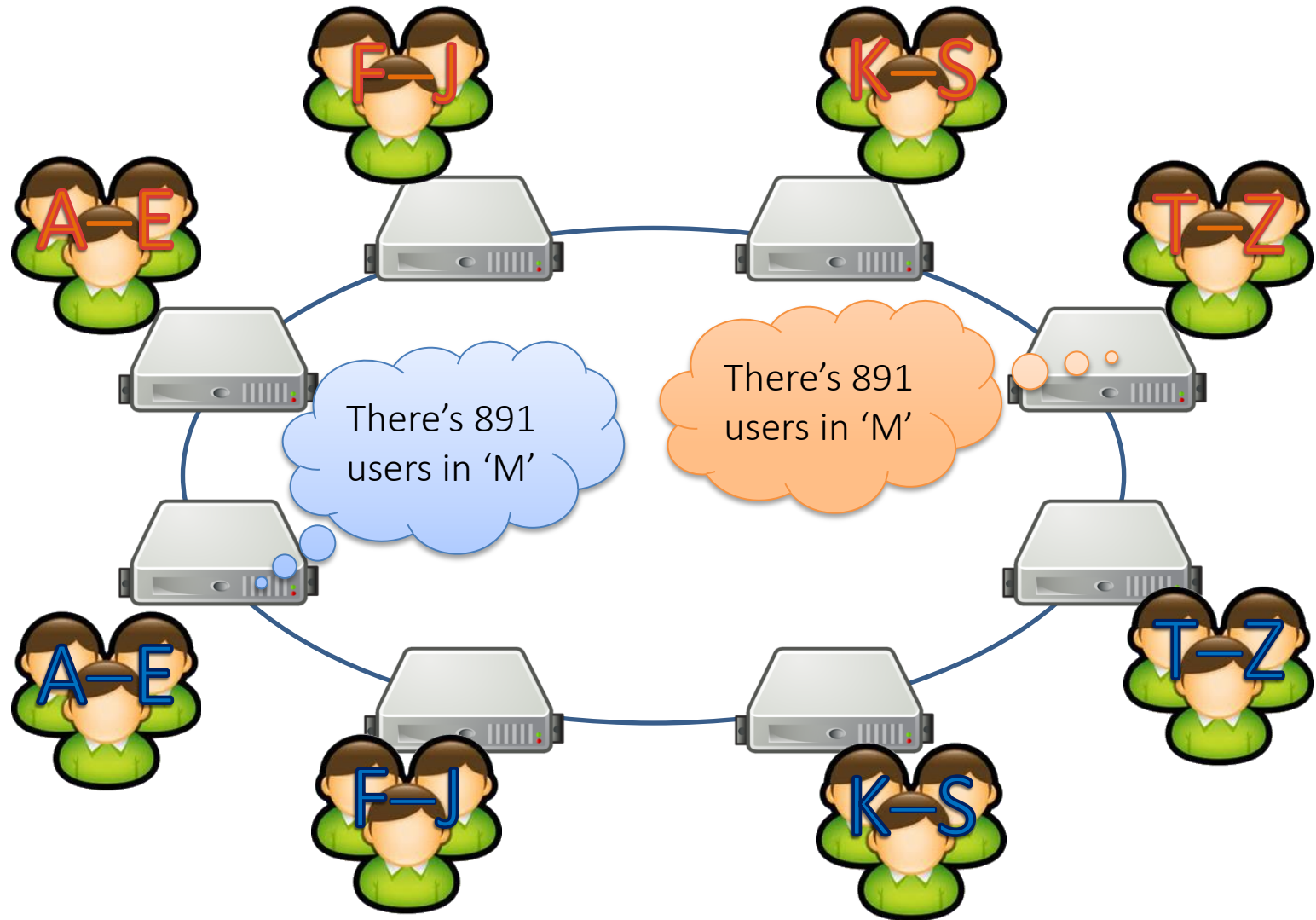
CA in CAP not the same as CA in ACID!!



# A Distributed System (with Replication)

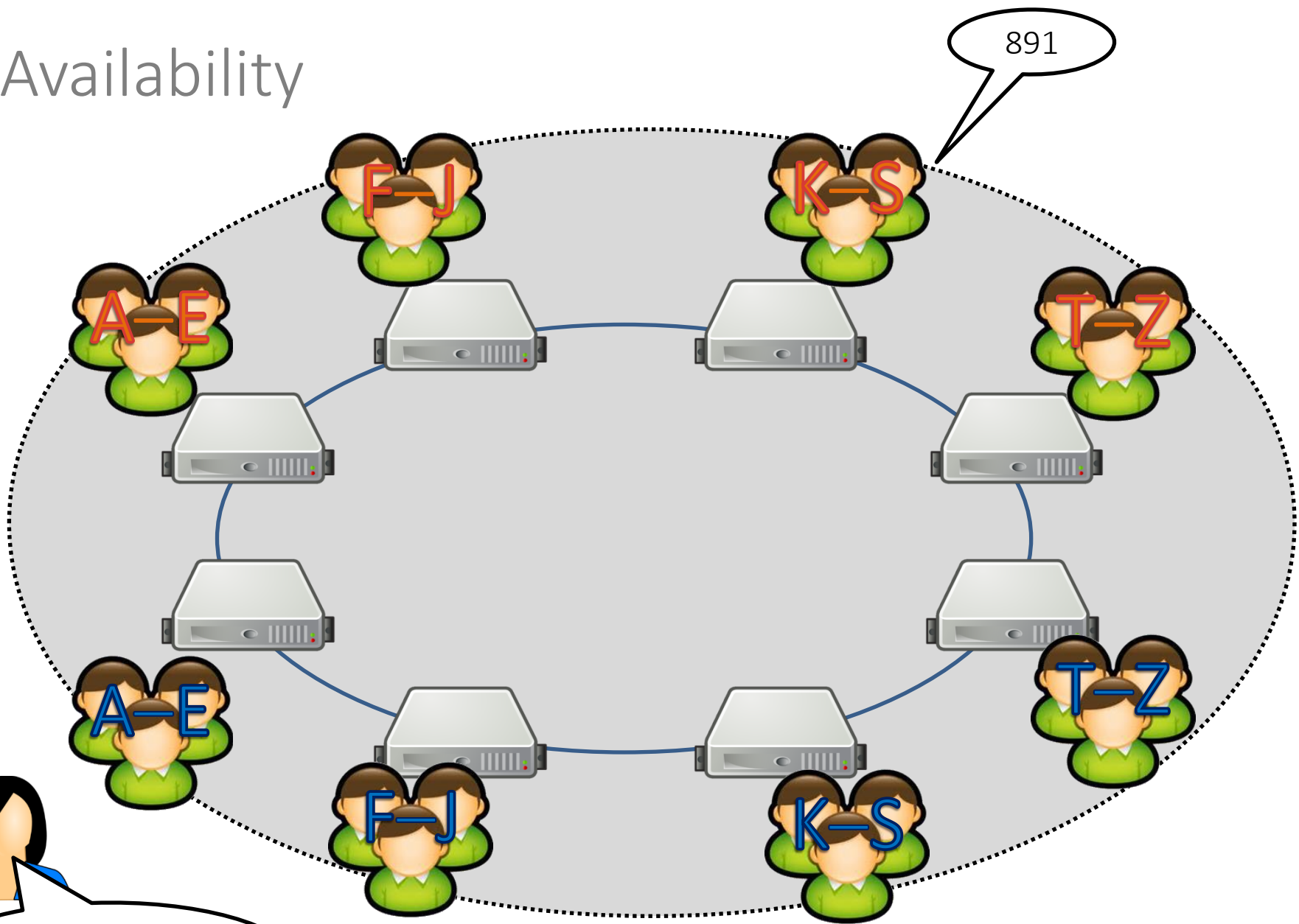


# Consistency



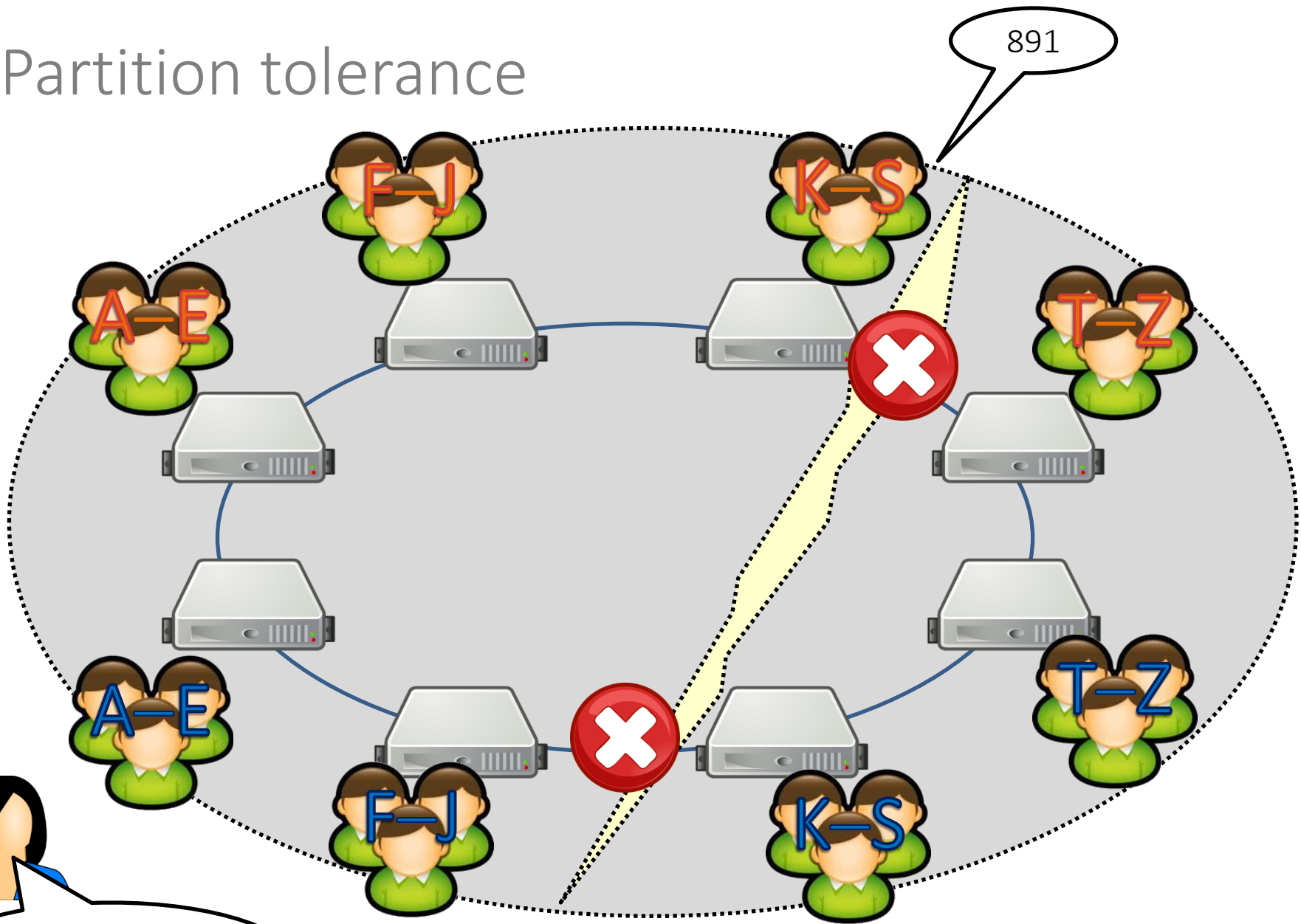


# Availability



How many users start with 'M'

# Partition tolerance



891

How many users start with 'M'

# The CAP Question

Can a distributed system guarantee

**consistency** (all nodes have the same up-to-date view),

**availability** (every read/write is acted upon) **and**

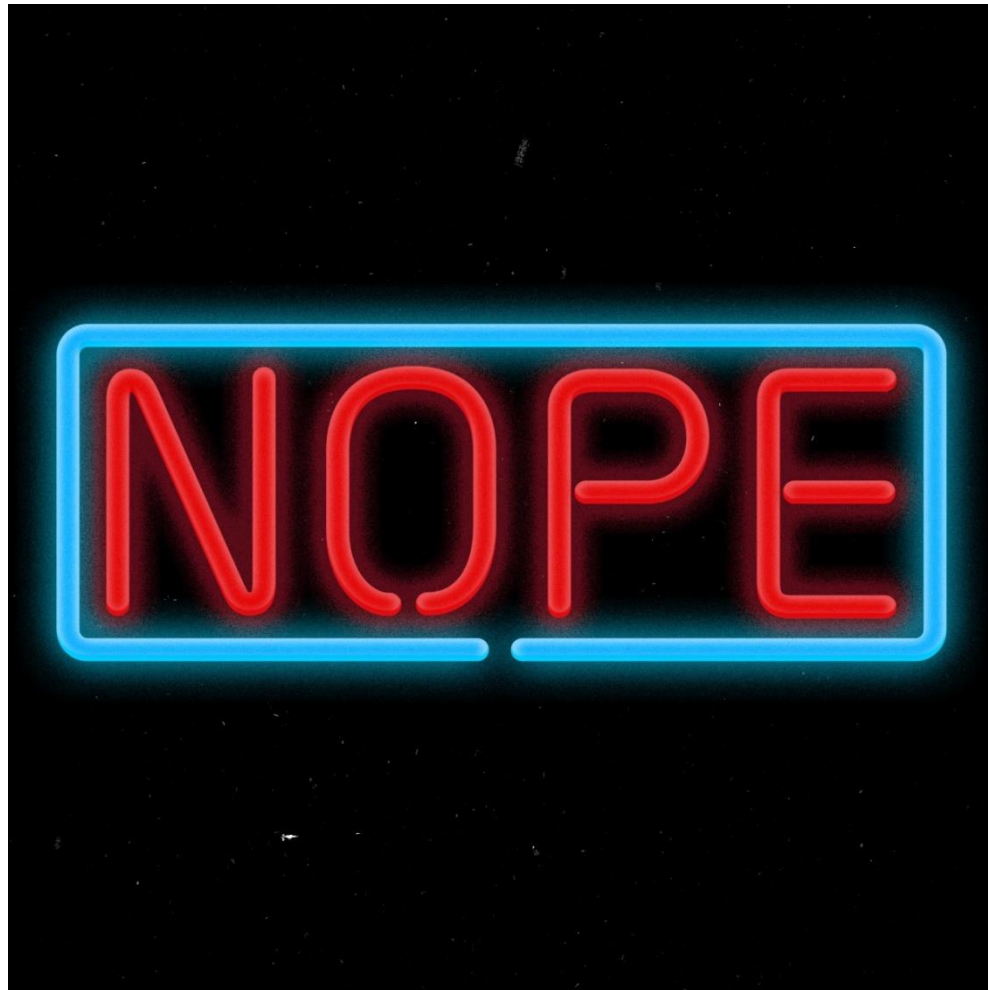
**partition-tolerance** (the system works if messages are lost)

at the same time?

What do you think?



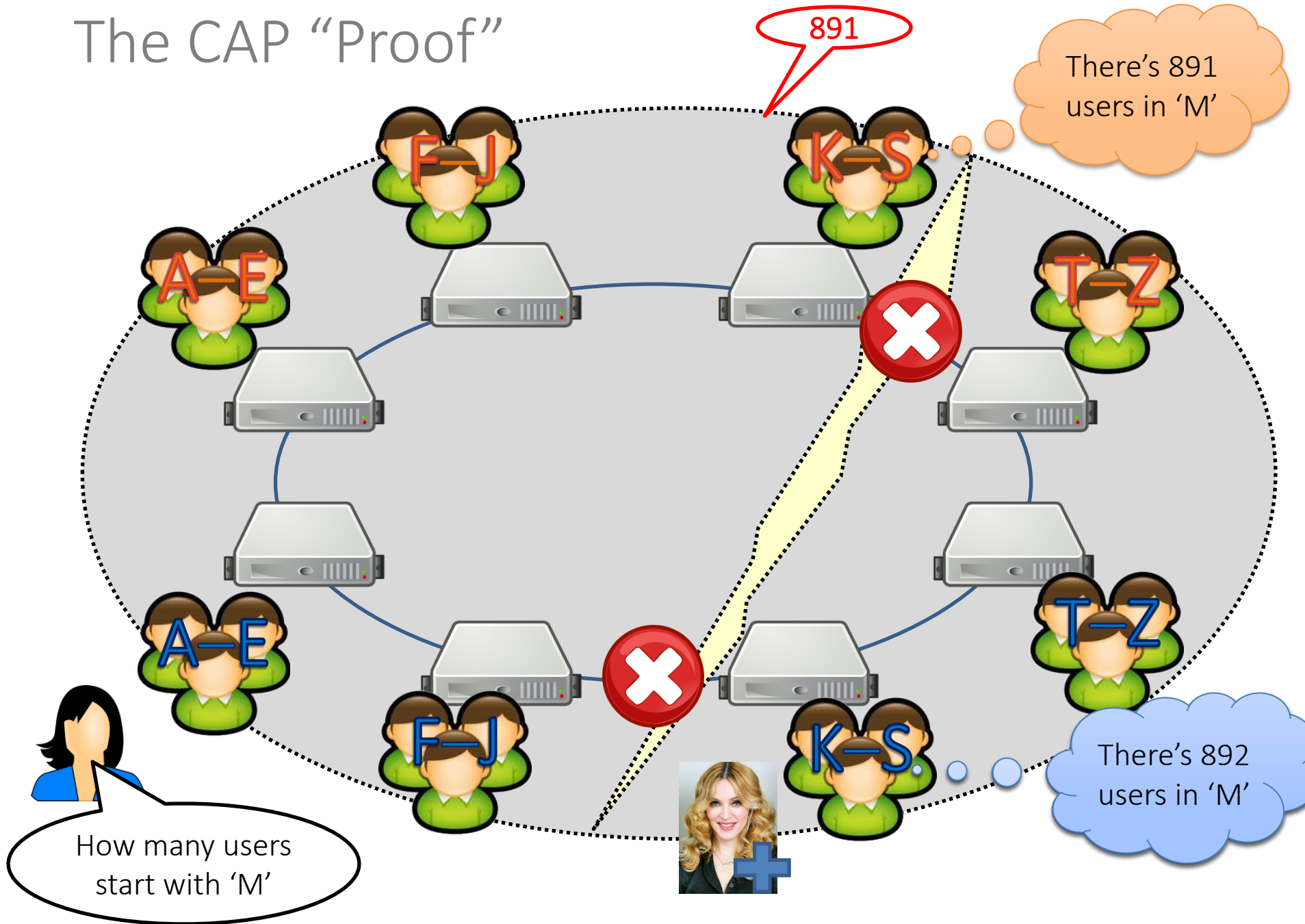
# The CAP Answer



# The CAP Theorem

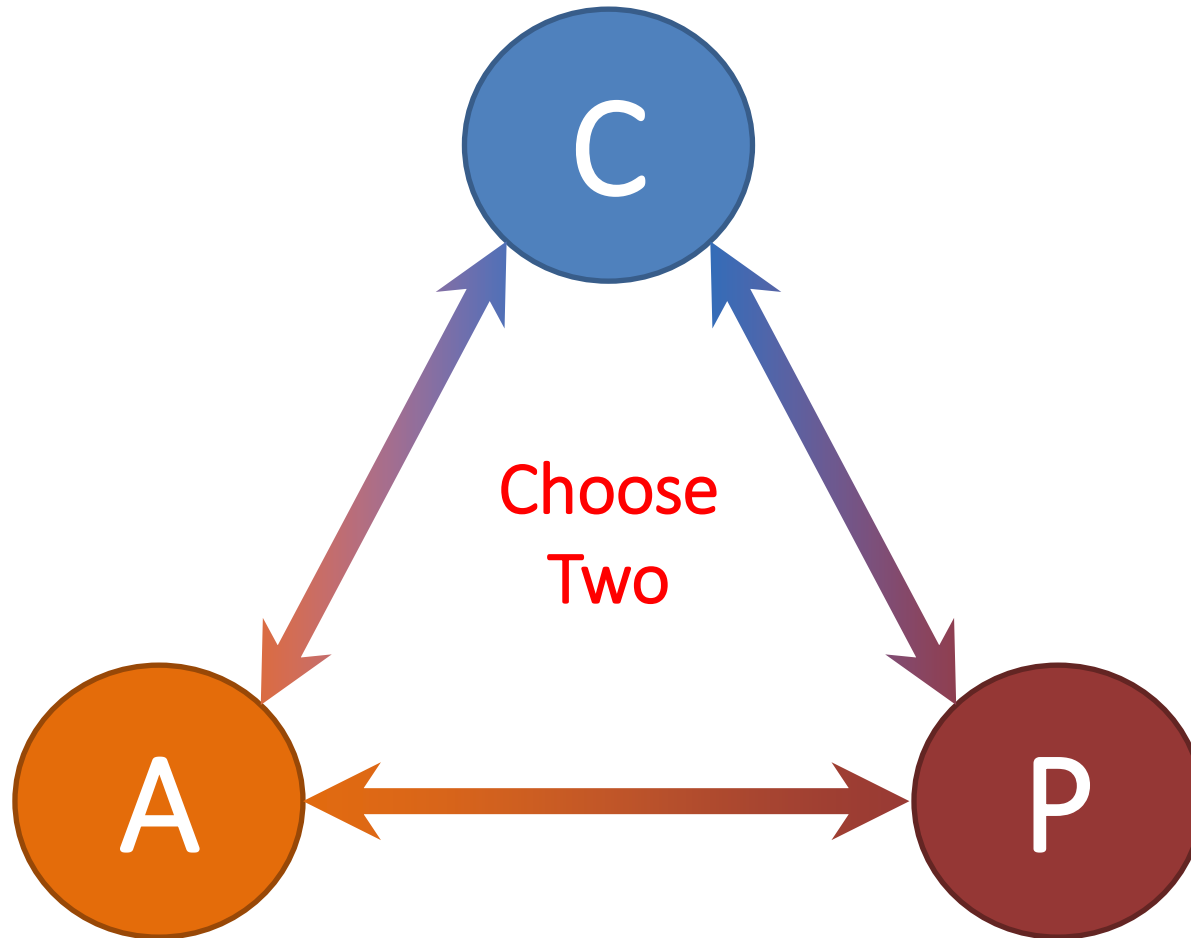
A distributed system cannot guarantee  
**consistency** (all nodes have the same up-to-date view),  
**availability** (every read/write is acted upon) **and**  
**partition-tolerance** (the system works if messages are lost)  
at the same time!

# The CAP “Proof”





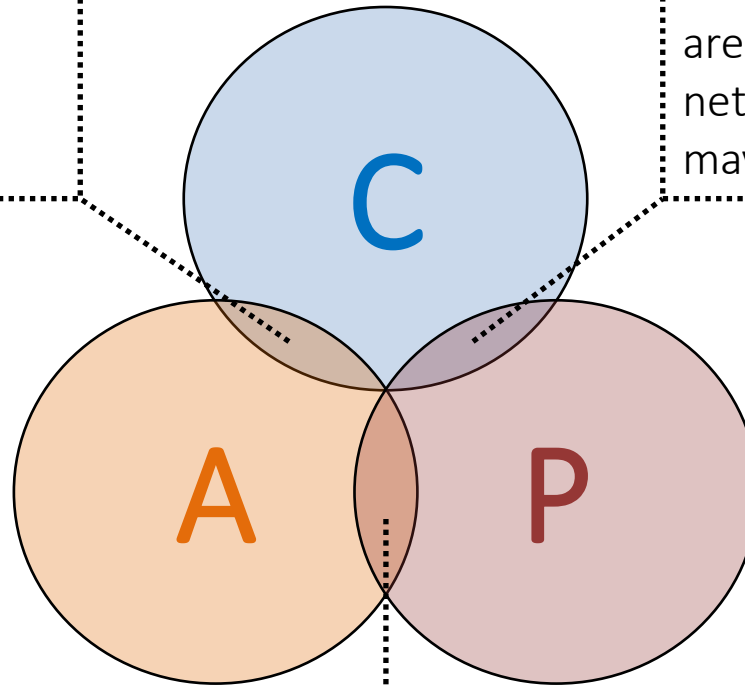
# The CAP Triangle



# CAP Systems

**CA**: Guarantees to give a correct response but only while network works fine  
*(Centralised / Traditional)*

**CP**: Guarantees responses are correct even if there are network failures, but response may fail  
*(Weak availability)*



(No intersection)

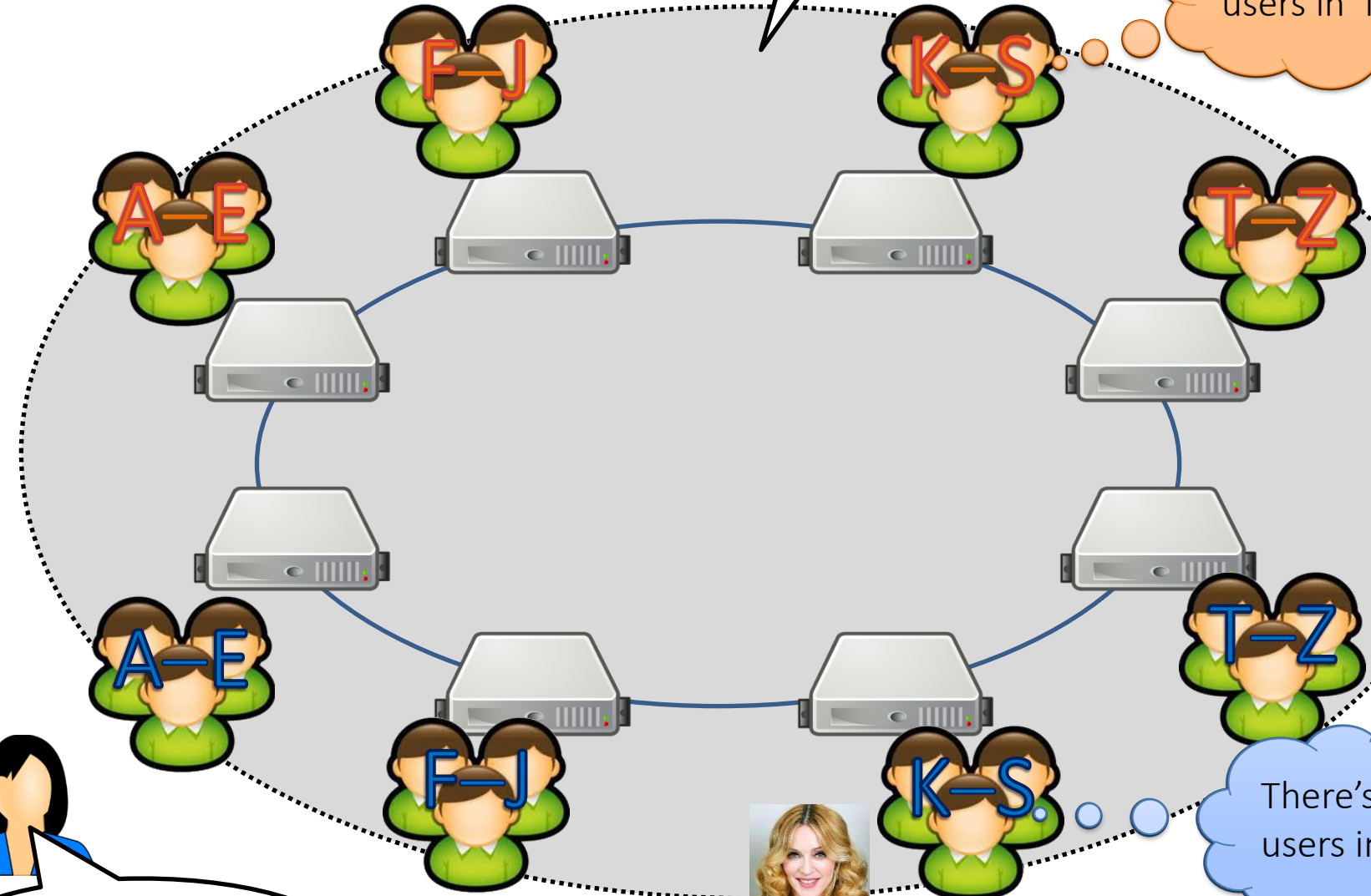
**AP**: Always provides a “best-effort” response even in presence of network failures  
*(Eventual consistency)*

# CA System

892

There's 892 users in 'M'

There's 892 users in 'M'



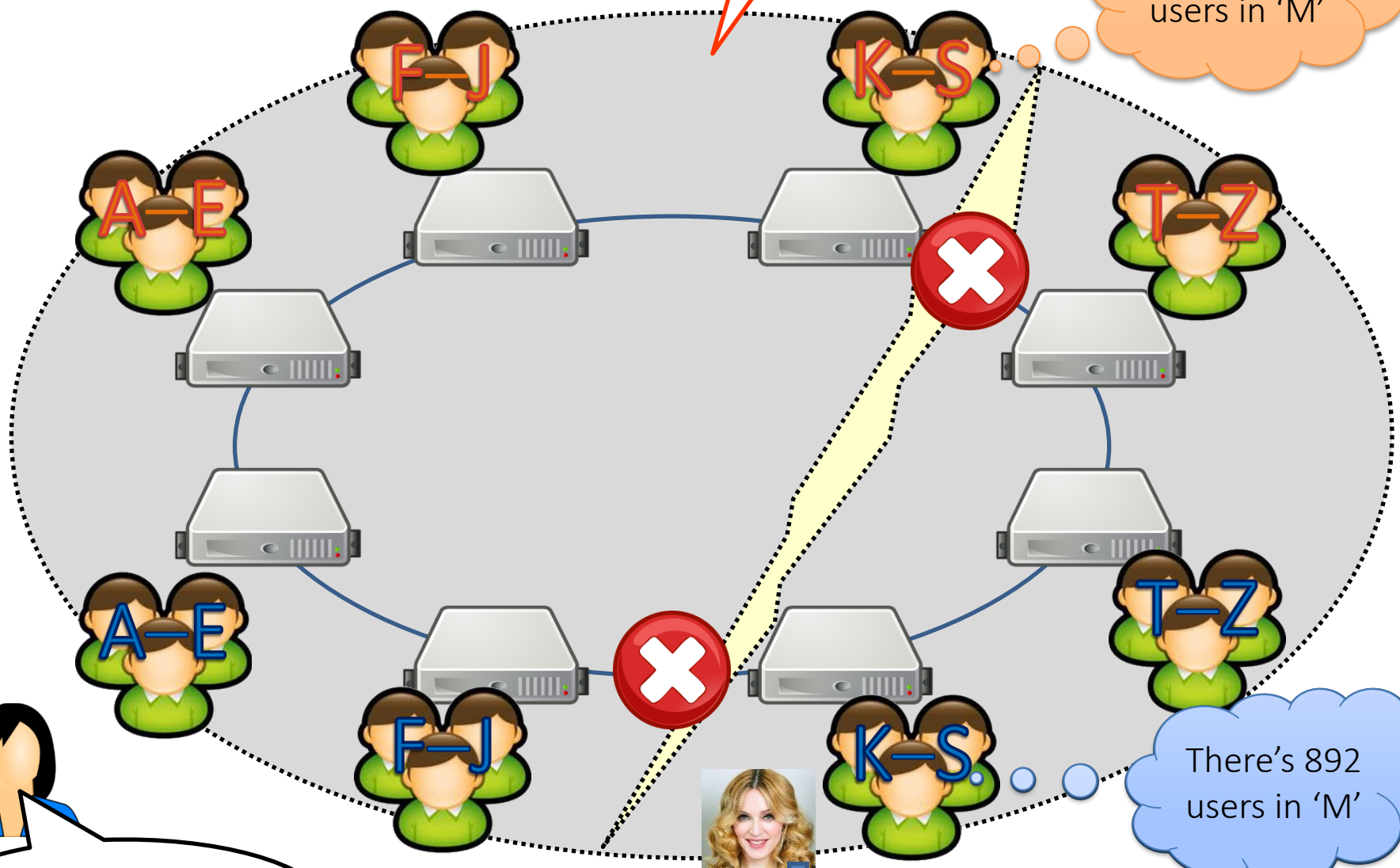
How many users start with 'M'



# CP System

Error

There's 891 users in 'M'



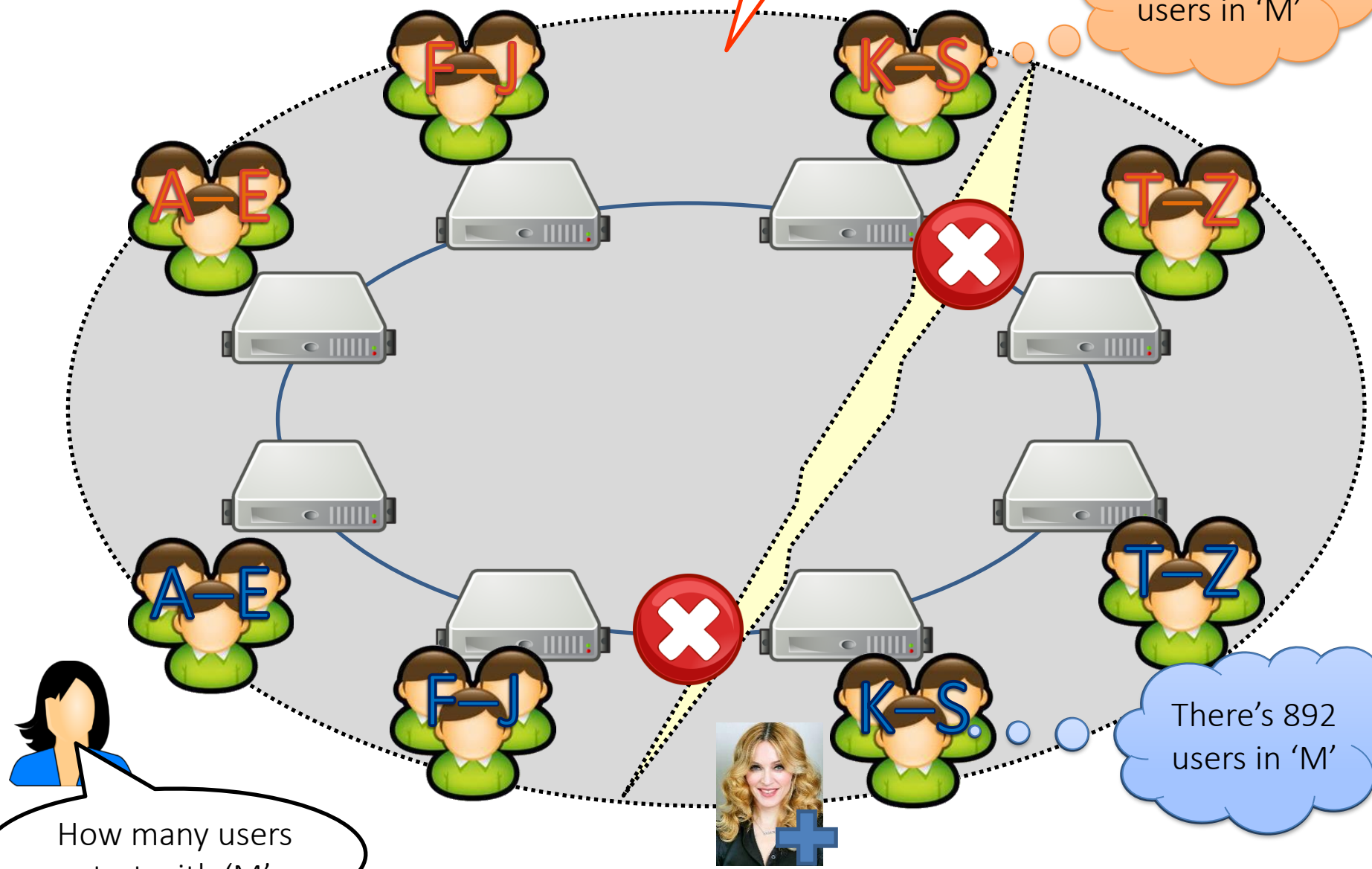
How many users start with 'M'

There's 892 users in 'M'

# AP System

891

There's 891 users in 'M'



How many users start with 'M'

There's 892 users in 'M'



# BASE (AP)

- **B**asically **A**vailable
  - Almost always “up”
- **S**oft State
  - Replicated, cached data
- **E**ventual Consistency
  - Stale data tolerated, for a while

In what way does Twitter act as a BASE (AP) system?



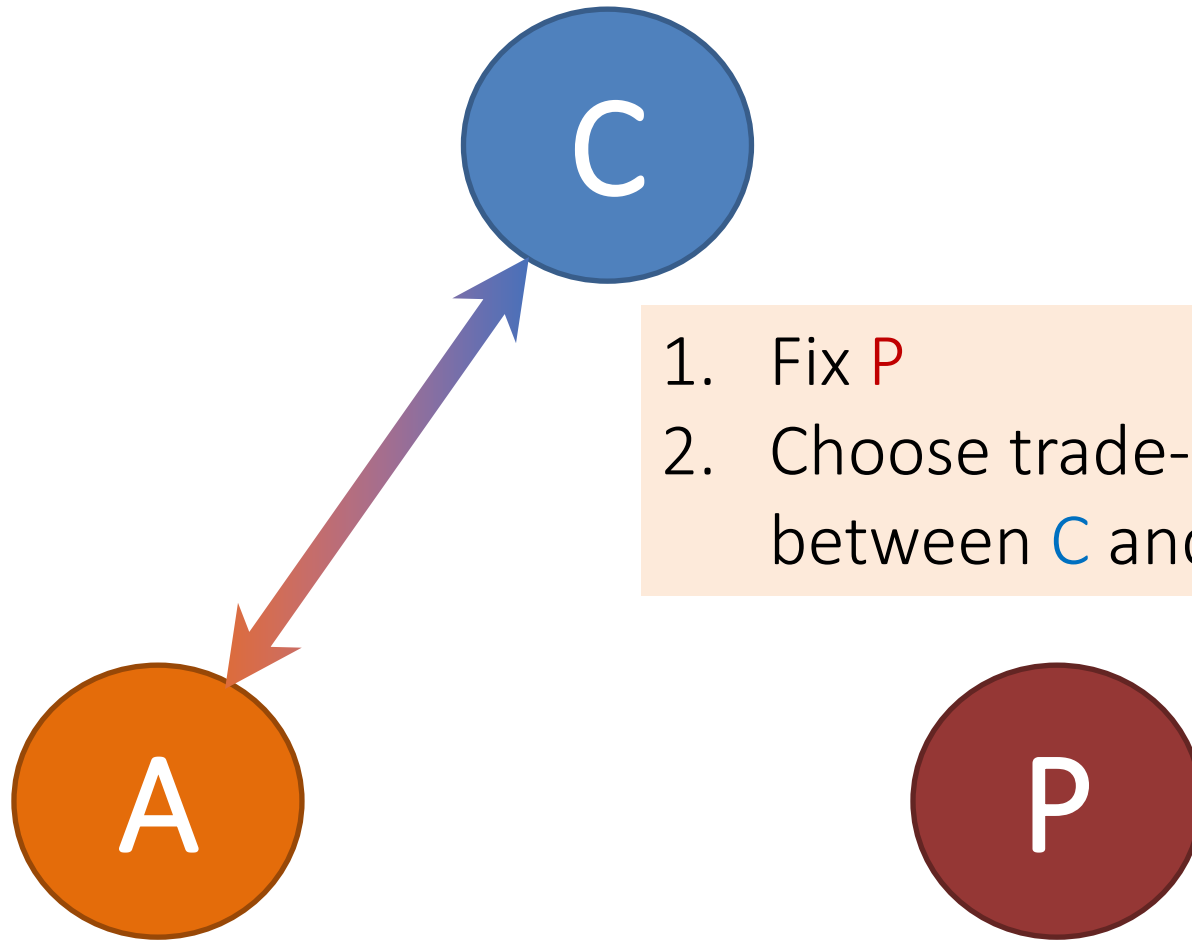
# High-fanout creates a “partition”

Rank ↕	Change (monthly)	Account name ↕	Owner ↕	Followers (millions) ↕	Occupation ↕	Country ↕
1	—	@BarackObama	Barack Obama	129.9	44th President of the United States	United States
2	—	@justinbieber	Justin Bieber	114.1	Musician	Canada
3	—	@katyperry <sup>[a]</sup>	Katy Perry	109.1	Musician	United States
4	—	@rihanna	Rihanna	102.5	Musician and businesswoman	Barbados
5	—	@Cristiano	Cristiano Ronaldo	92.1	Footballer	Portugal
6	—	@taylorswift13	Taylor Swift	88.6	Musician	United States
7	—	@ladygaga	Lady Gaga	83.9	Musician and actress	United States
8	—	@ArianaGrande	Ariana Grande	83.2	Musician and actress	United States
9	—	@TheEllenShow	Ellen DeGeneres	78.5	Comedian and television hostess	United States
10	—	@YouTube	YouTube	73	Online video platform	United States

Users may see retweets of celebrity tweets before the original tweet.

Later when the original tweet arrives the timeline will be reordered and made consistent.

# CAP in practical distributed systems

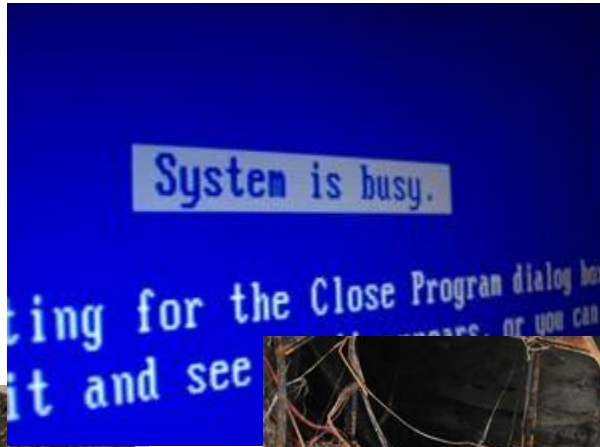


1. Fix **P**
2. Choose trade-off point between **C** and **A**



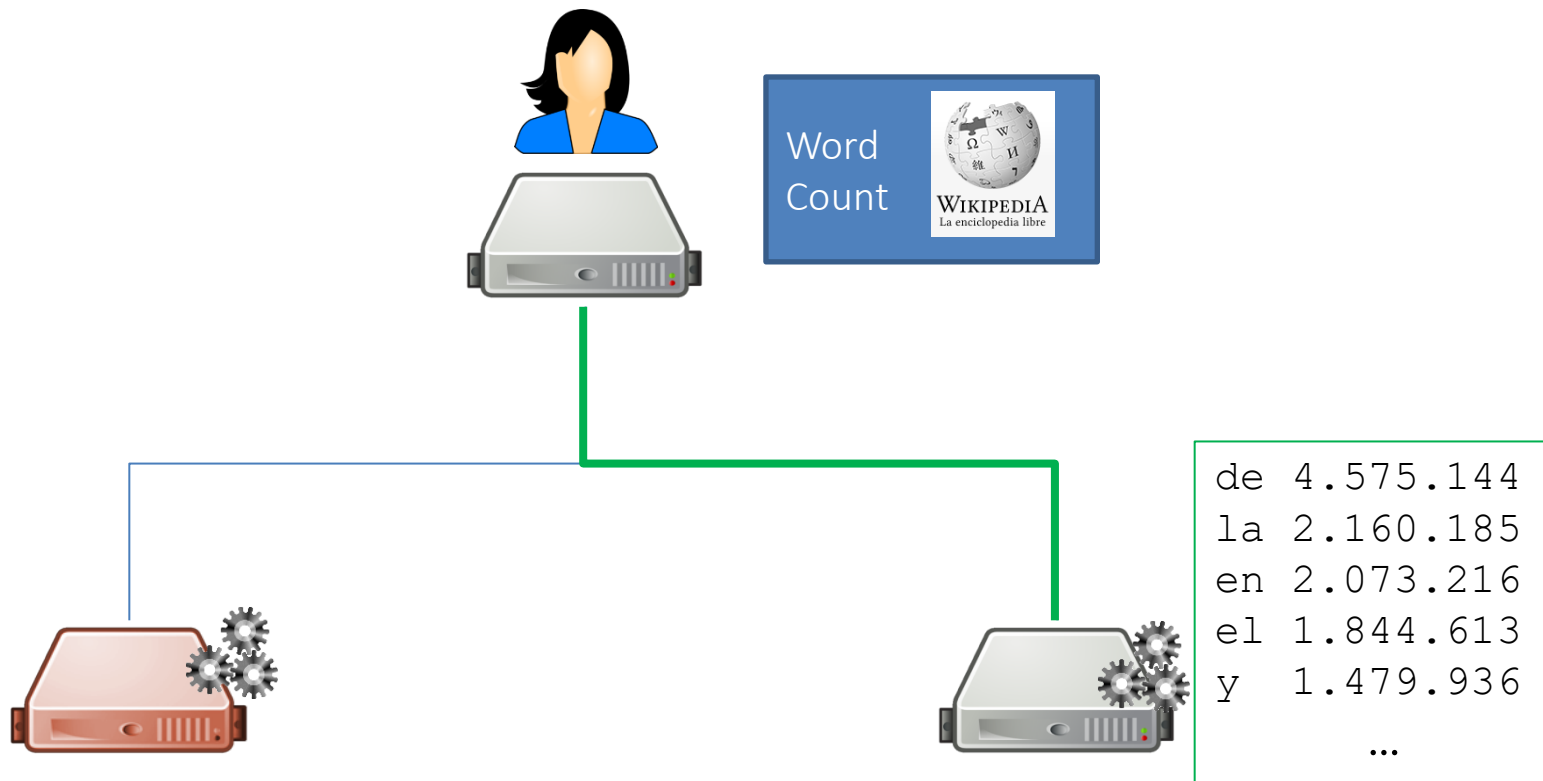
# PARTITION TOLERANCE

# Faults



# Fail–Stop Fault

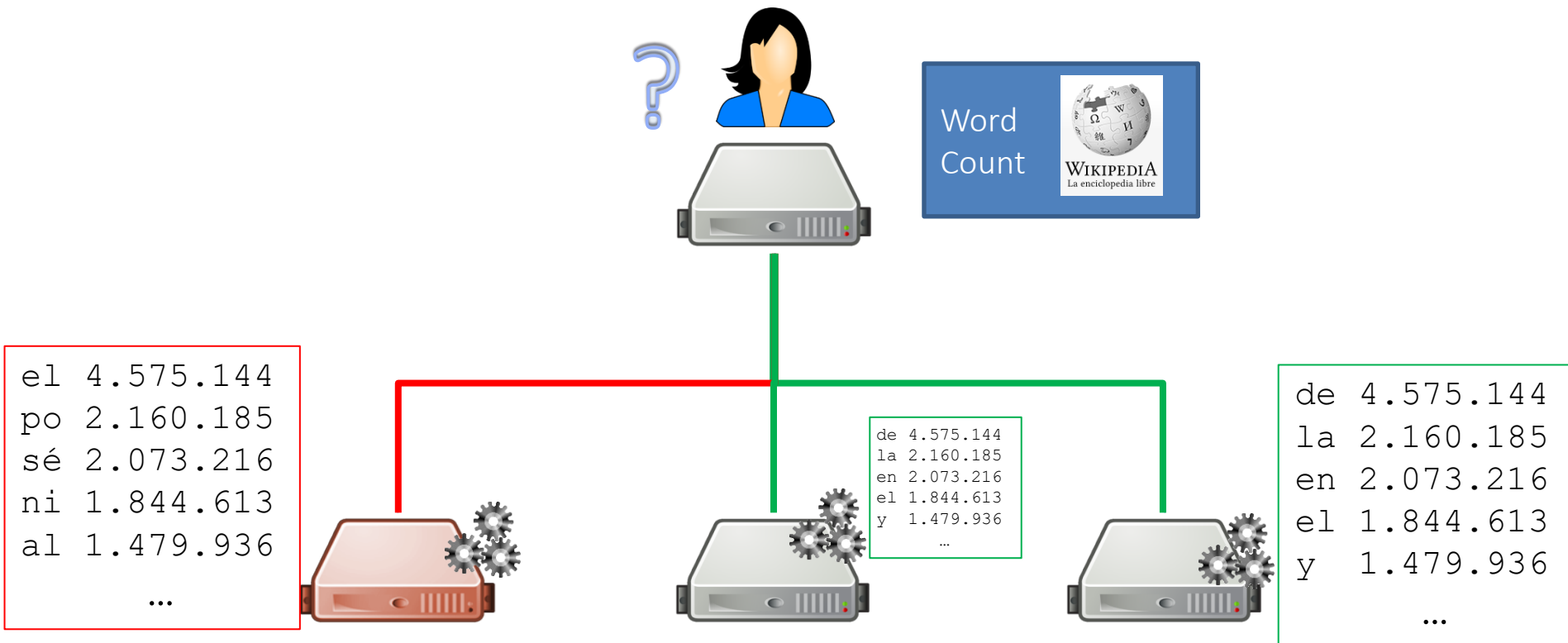
- A machine fails to respond or times-out
  - often hardware or load
  - need at least  $f + 1$  replicated machines
    - $f$  = number of fail-stop failures



# Byzantine Fault

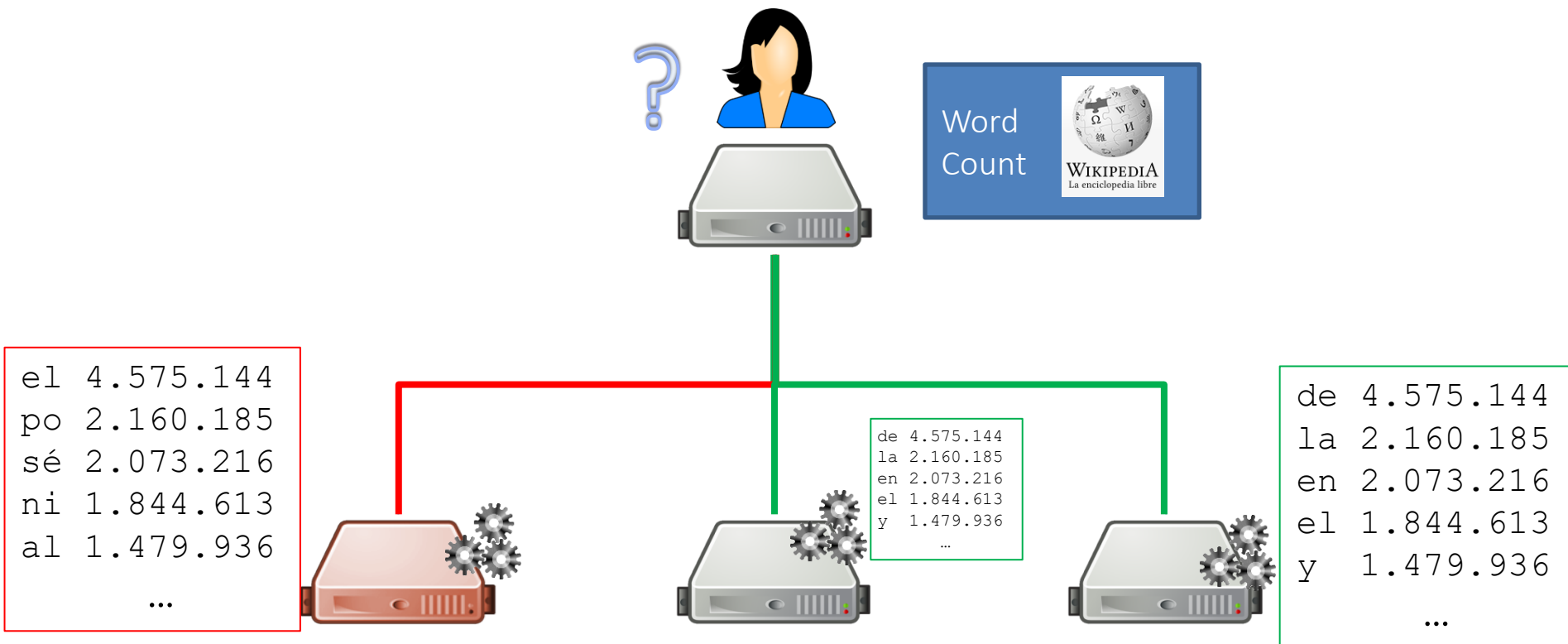
- A machine responds incorrectly/maliciously

How many working machines do we need in the general case to be robust against Byzantine faults?



# Byzantine Fault

- A machine responds incorrectly/maliciously
  - Need *at least*  $2f+1$  replicated machines
    - $f$  = number of (possibly Byzantine) failures



# DISTRIBUTED CONSENSUS

# Distributed Consensus

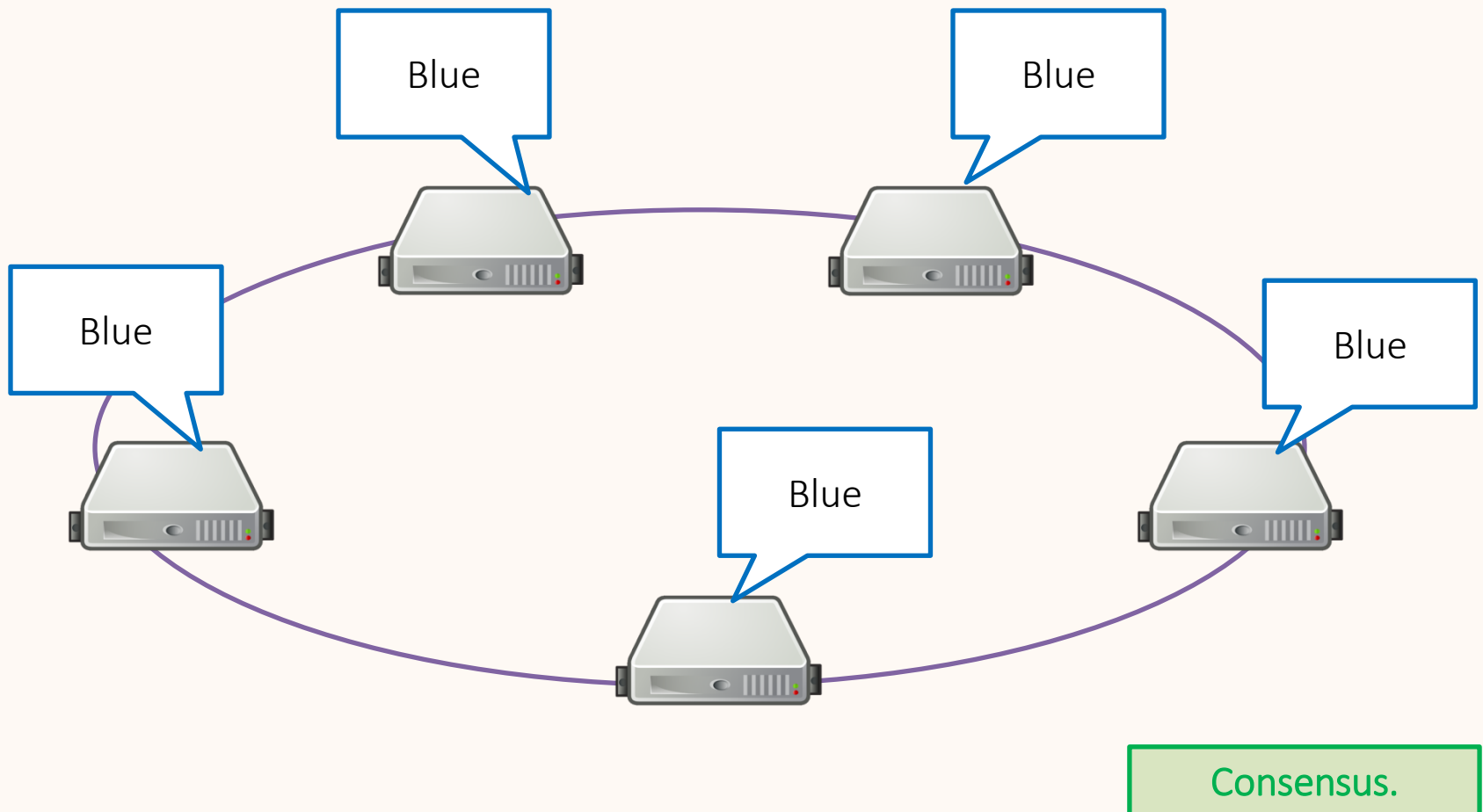


Colour of the dress?



# Distributed Consensus

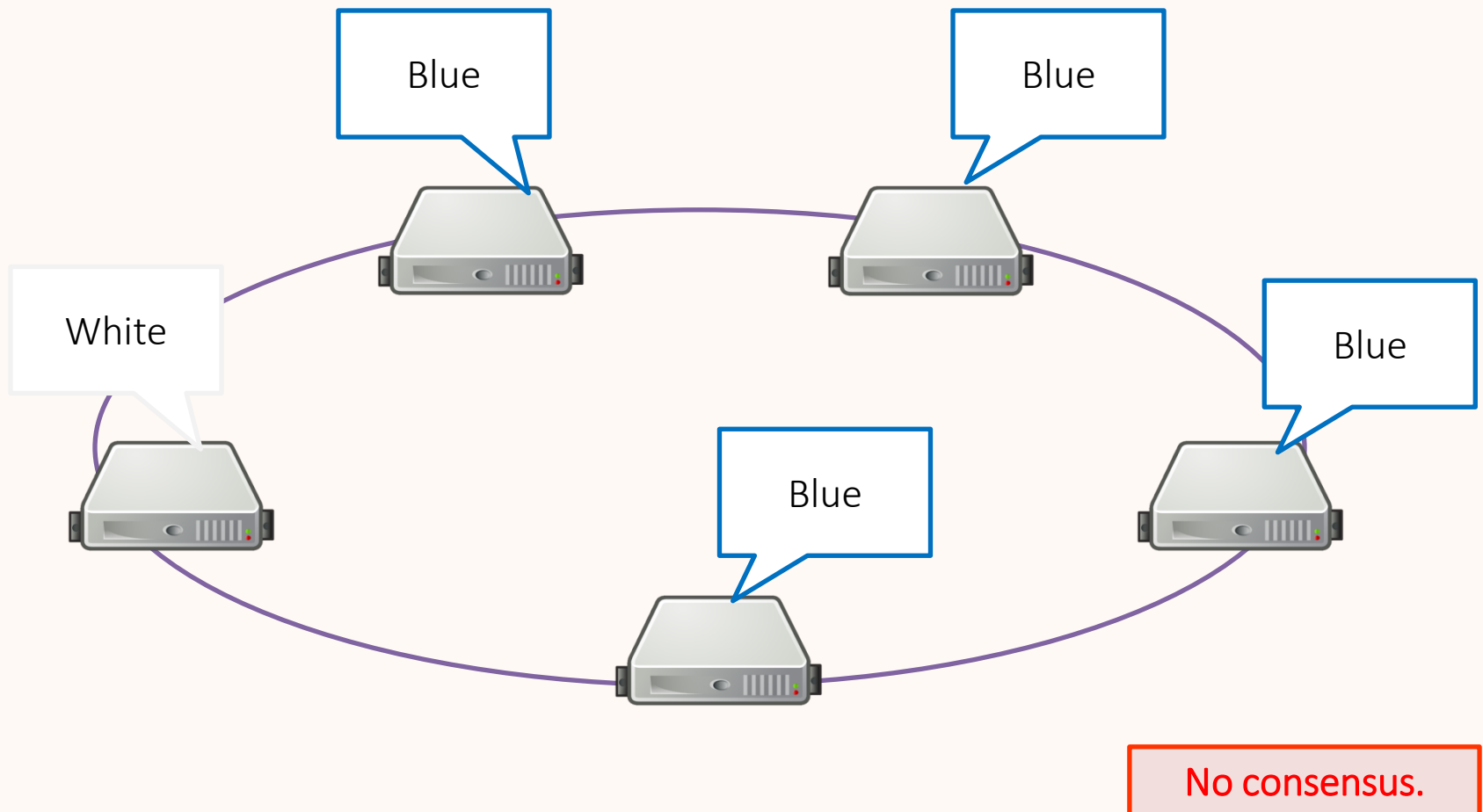
Strong consensus: All nodes need to agree





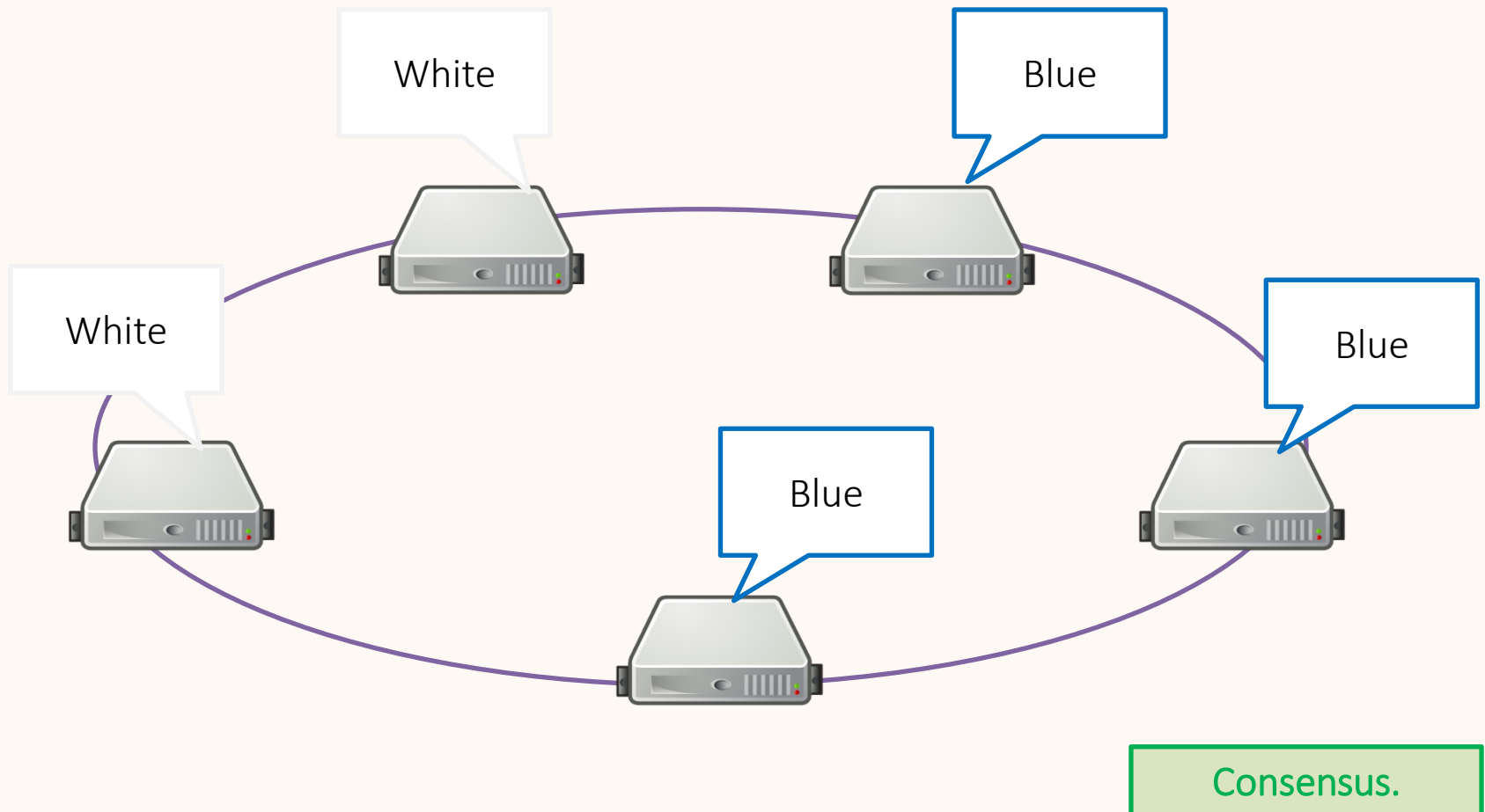
# Distributed Consensus

Strong consensus: All nodes need to agree



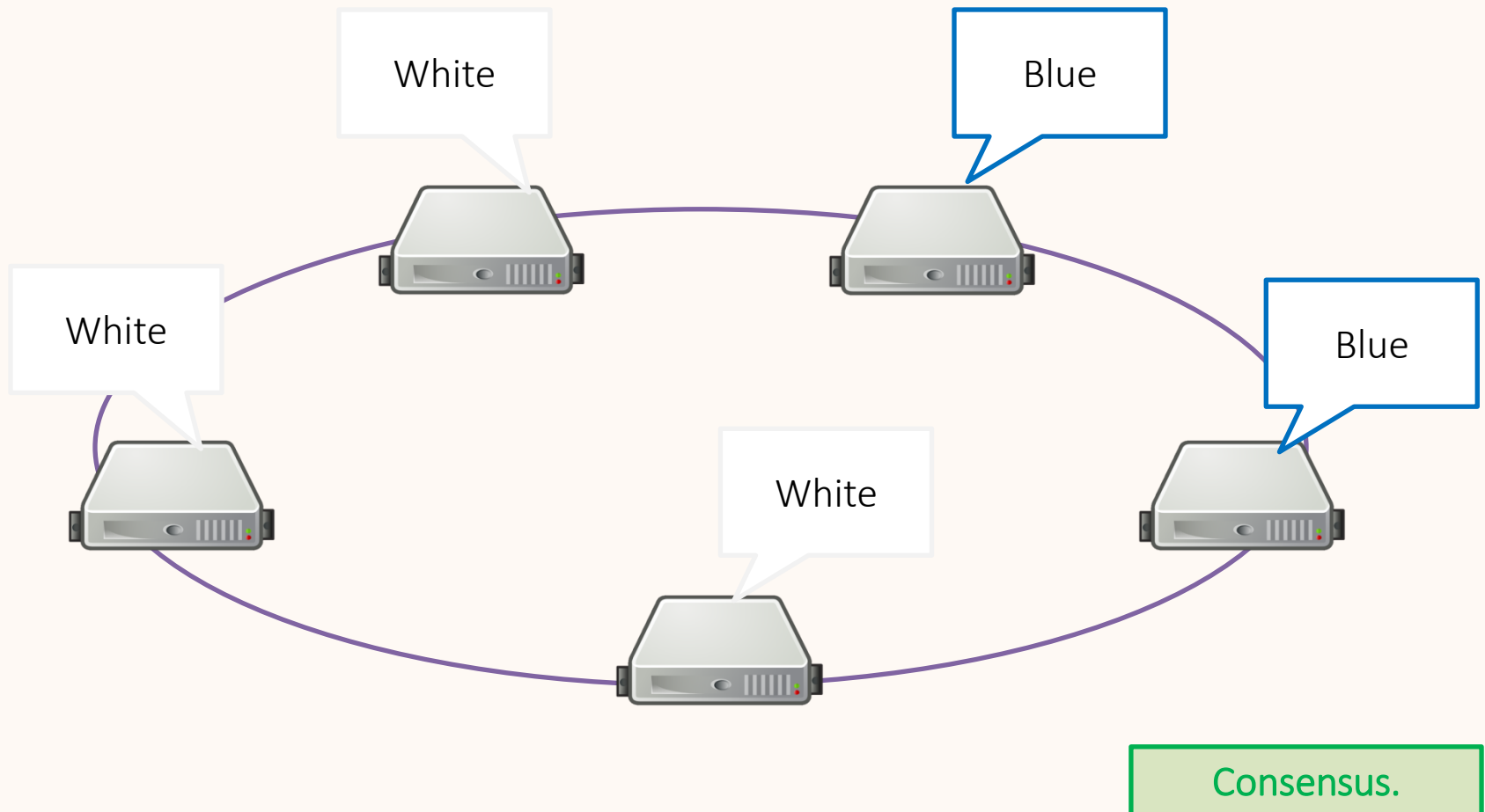
# Distributed Consensus

Majority consensus: A majority of nodes need to agree



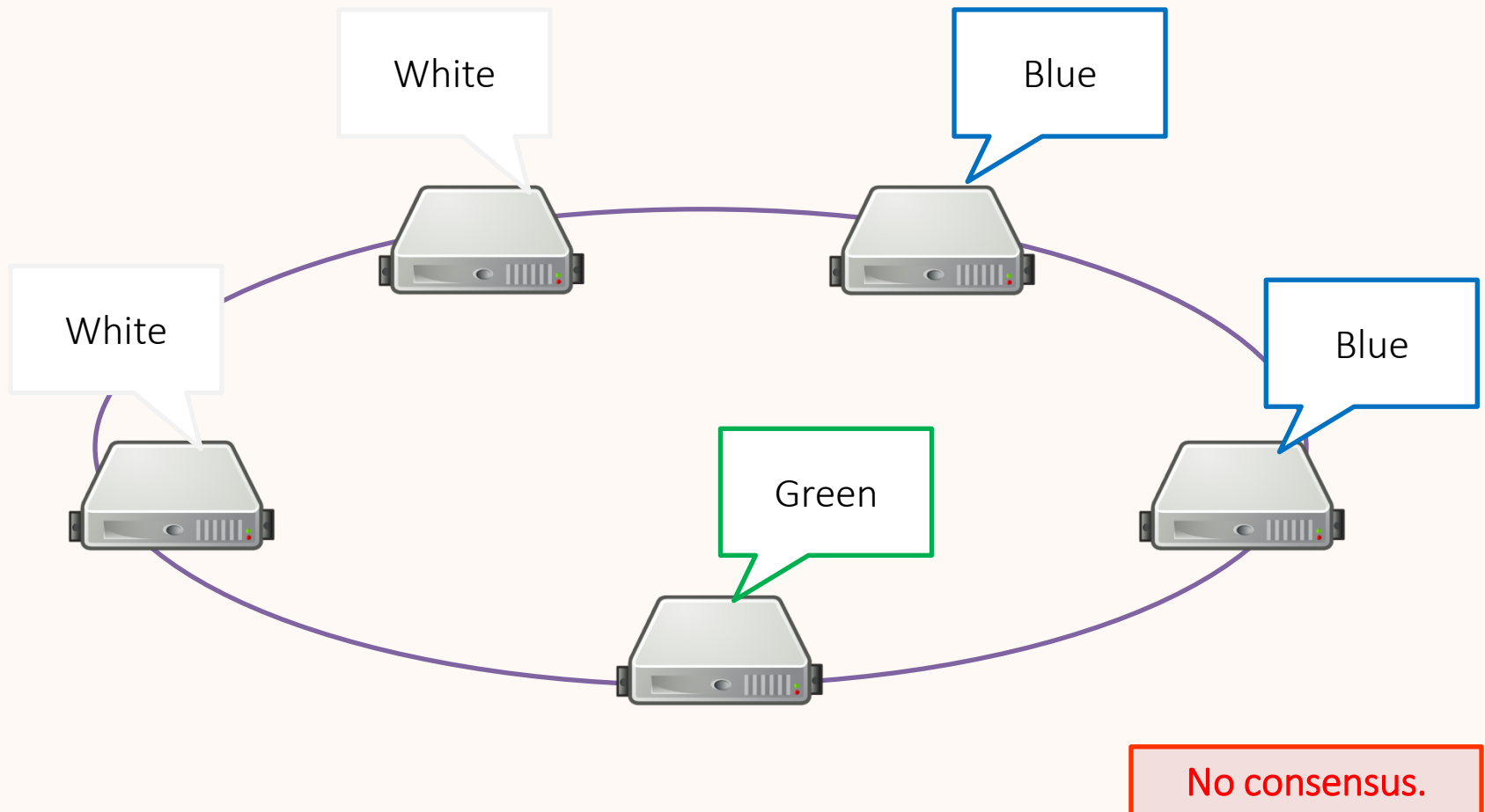
# Distributed Consensus

Majority consensus: A majority of nodes need to agree



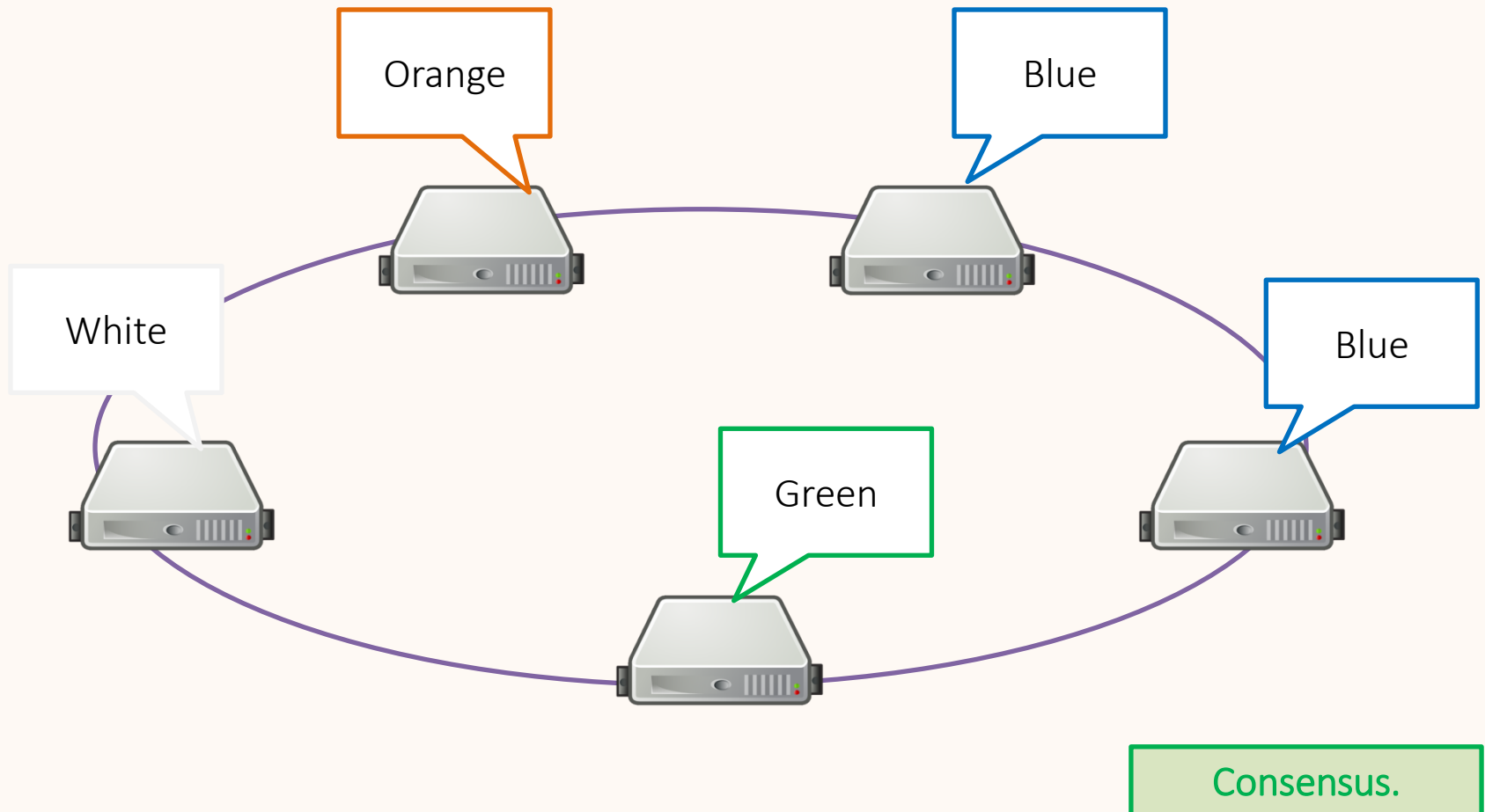
# Distributed Consensus

Majority consensus: A majority of nodes need to agree



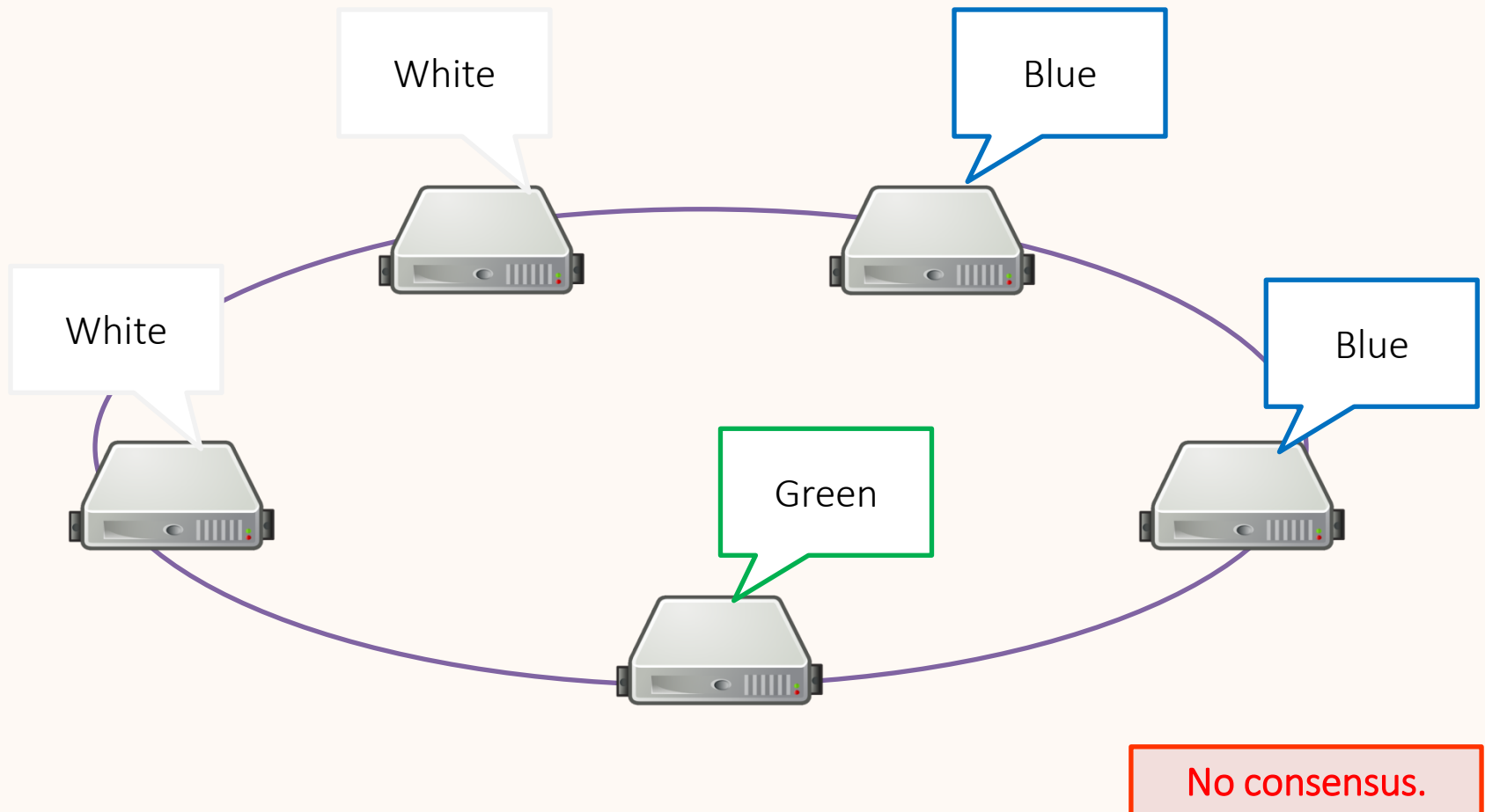
# Distributed Consensus

Plurality consensus: A plurality of nodes need to agree



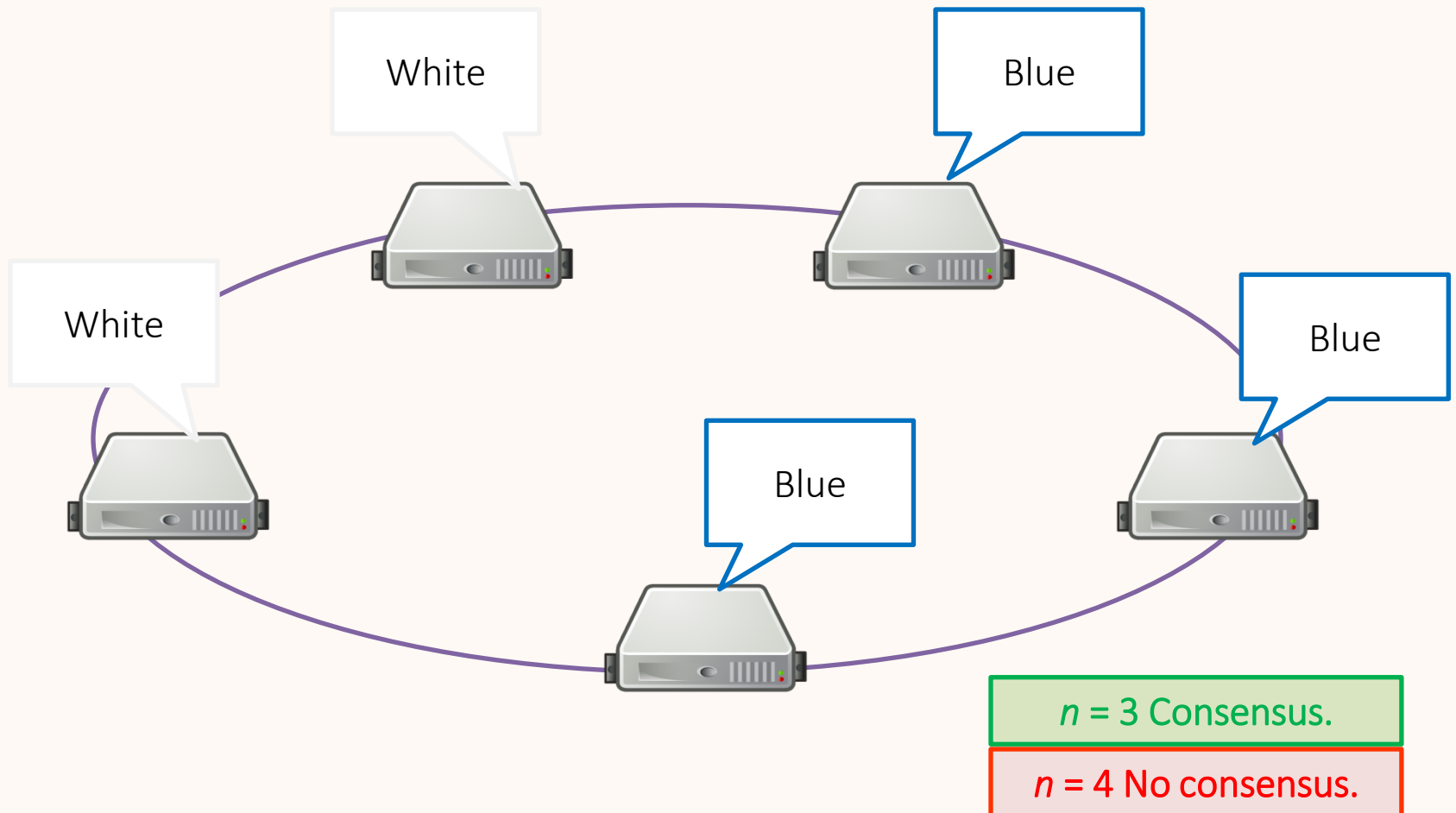
# Distributed Consensus

Plurality consensus: A plurality of nodes need to agree



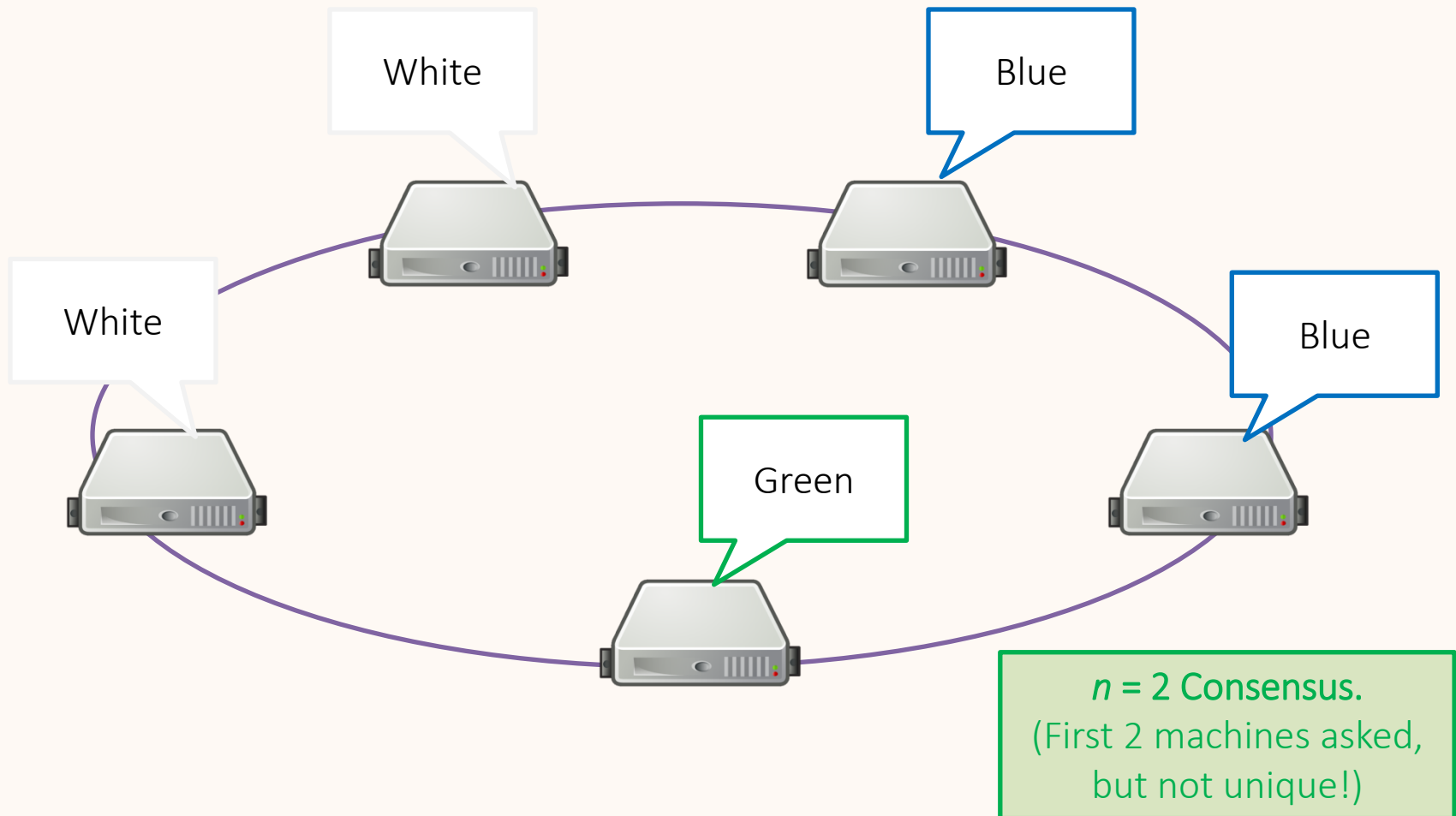
# Distributed Consensus

Quorum consensus:  $n$  nodes need to agree



# Distributed Consensus

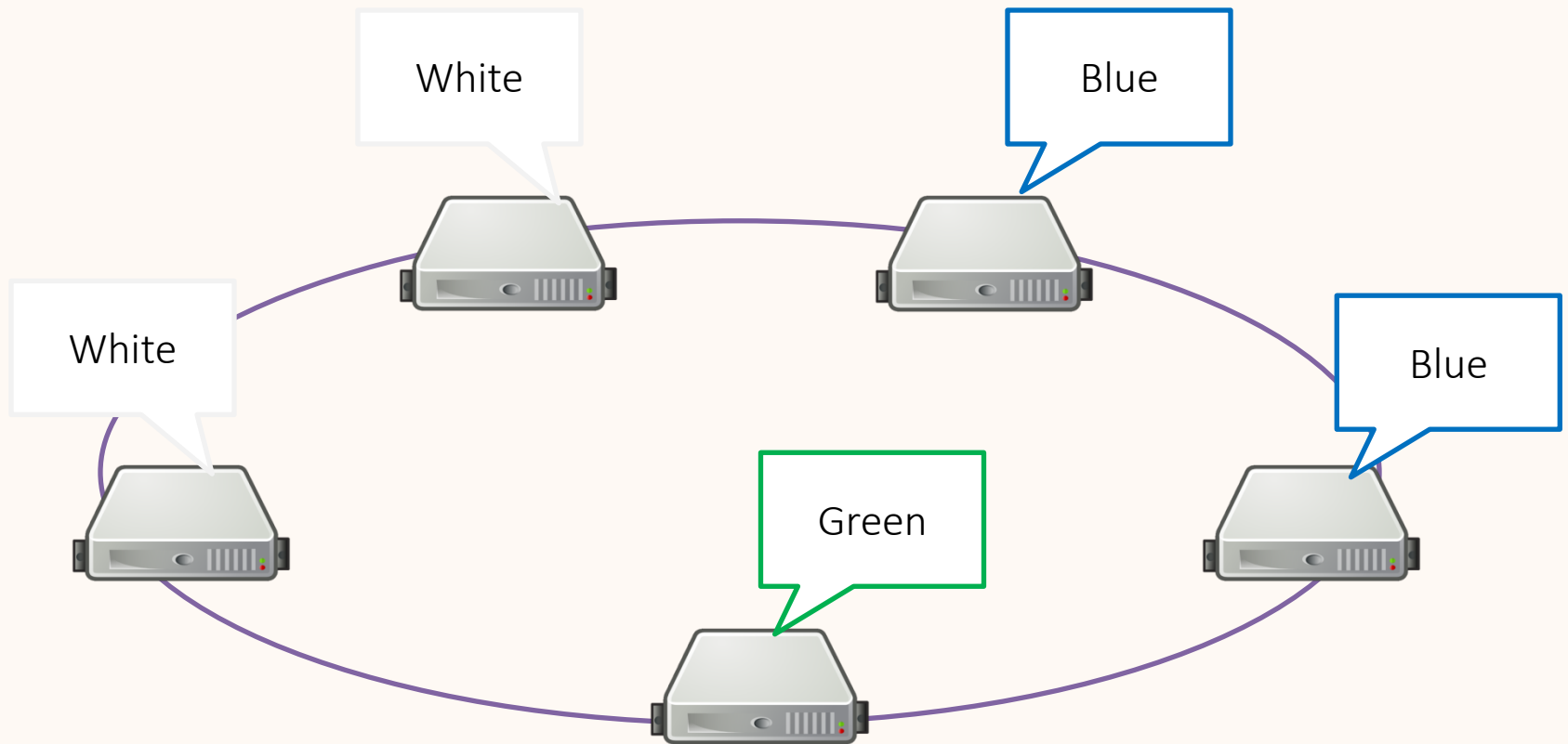
Quorum consensus:  $n$  nodes need to agree





# Distributed Consensus

Quorum consensus:  $n$  nodes need to agree



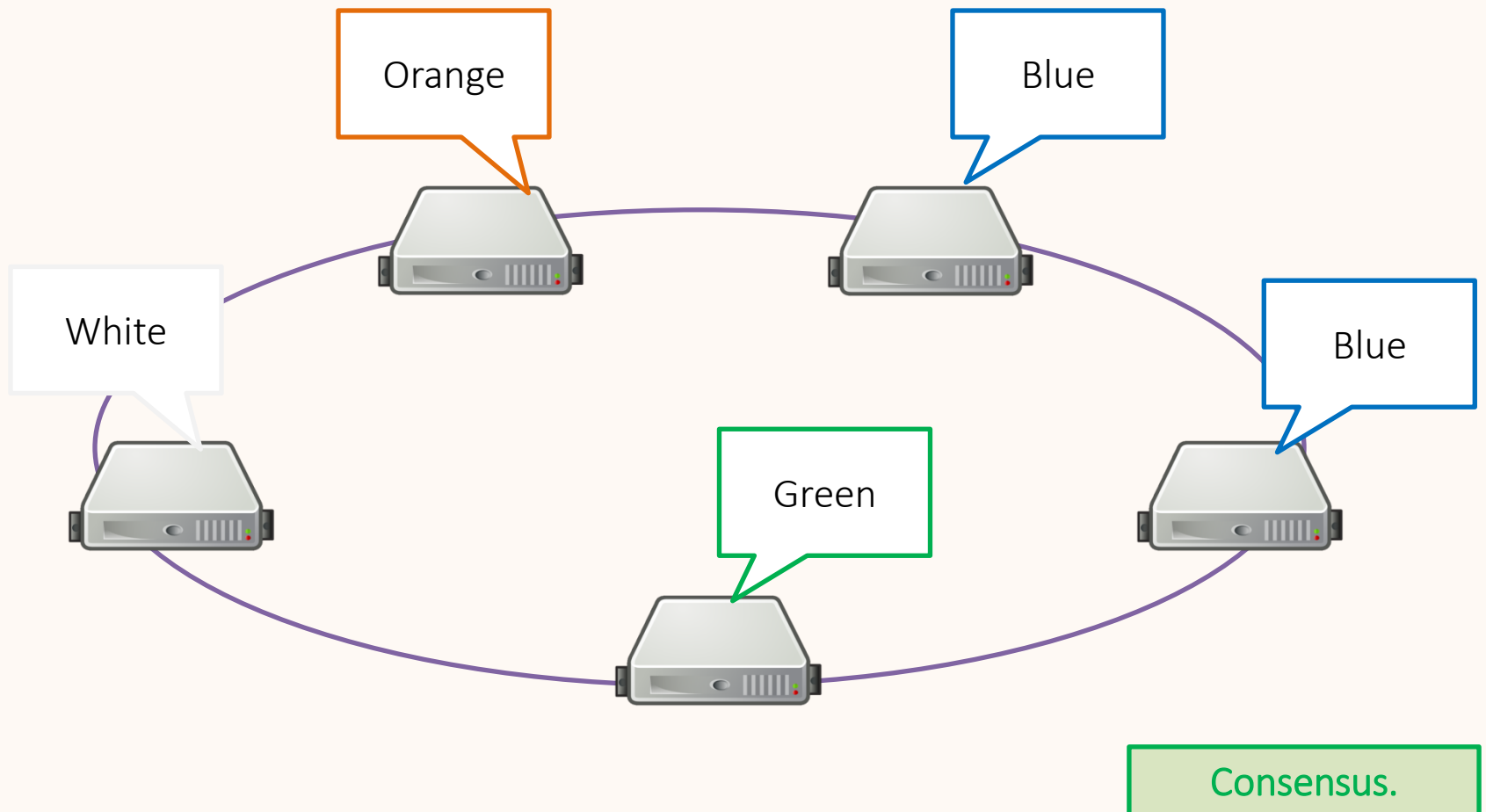
Value of  $n$  needed for unique consensus with  $N$  nodes?



$n > N/2$

# Distributed Consensus

Consensus off: Take first answer



# Distributed Consensus

CP vs. AP?



Strong consensus: All nodes need to agree

CP

Majority consensus: A majority of nodes need to agree

Plurality consensus: A plurality of nodes need to agree

Quorum consensus: “Fixed”  $n$  nodes need to agree

Consensus off: Take first answer

AP

# Distributed Consensus

Scale?



Strong consensus: All nodes need to agree

More replication

Majority consensus: A majority of nodes need to agree

Plurality consensus: A plurality of nodes need to agree

Quorum consensus: “Fixed”  $n$  nodes need to agree

Consensus off: Take first answer

Less replication

# Distributed Consensus

Strong consensus: All nodes need to agree

Majority consensus: A majority of nodes need to agree

**Choice is application dependent:**

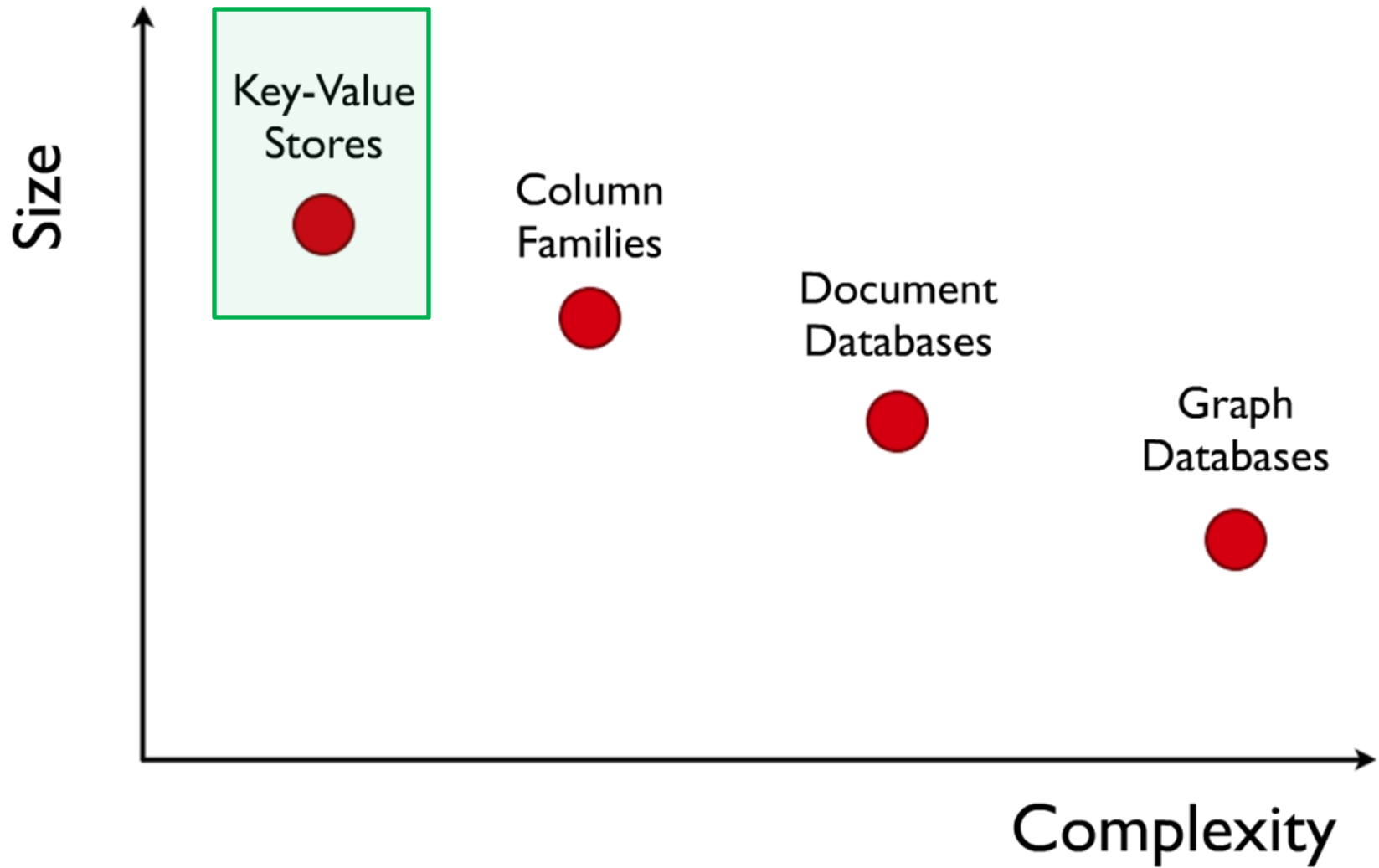
Plurality Many NoSQL stores allow you to choose level of consensus/replication

Quorum consensus: “Fixed”  $n$  nodes need to agree

Consensus off: Take first answer

# NoSQL: KEY-VALUE STORES

# NoSQL: Key-Value Stores



# Key–Value Store Model

It's just a Map / Associate Array / Dictionary 😊

- `put (key, value)`
- `get (key)`
- `delete (key)`

Key	Value
Afghanistan	Kabul
Albania	Tirana
Algeria	Algiers
Andorra la Vella	Andorra la Vella
Angola	Luanda
Antigua and Barbuda	St. John's
...	...



# But You Can Do a Lot With a Map

Key	Value
country:Afghanistan	capital@city:Kabul,continent:Asia,pop:31108077#2011
country:Albania	capital@city:Tirana,continent:Europe,pop:3011405#2013
...	...
city:Kabul	country:Afghanistan,pop:3476000#2013
city:Tirana	country:Albania,pop:3011405#2013
...	...
user:10239	basedIn@city:Tirana,post:{103,10430,201}
...	...

*... actually you can model any data in a map (but possibly with a lot of redundancy and inefficient lookups if unsorted).*

# THE CASE OF AMAZON

# The Amazon Scenario

Products Listings: prices, details, stock

The screenshot shows an Amazon search results page for the keyword "presenter". The page layout includes a top navigation bar with the Amazon logo, a search bar containing "presenter", and a secondary navigation bar with links for "Shop by Department", "Aidan's Amazon.com", "Today's Deals", "Gift Cards", "Sell", and "Help". Below the navigation bar, it indicates "1-16 of 19,088 results for 'presenter'".

On the left side, there is a "Show results for" section with a list of categories: "Office Products" (with sub-items: Office Presentation Remotes, Office Presentation Pointers), "Computers & Accessories" (with sub-items: Tablet Accessories, Computer Mice), "Cell Phones & Accessories" (with sub-item: Cell Phone Accessories), and "Software" (with sub-item: Presentations). A link to "See All 29 Departments" is also present. Below this is a "Refine by" section with filters for "International Shipping" (Ship to Ireland), "Eligible for Free Shipping" (Free Shipping by Amazon), and "Brand" (Kensington, Logitech, Targus, Satechi, Infiniter, August).


The main content area displays three product listings:

- Point and zoom in presentations with Myo Armband**: Includes an image of a black armband and a "Shop now" link.
- Logitech Wireless Presenter R400**: Includes an image of the presenter, the price "\$44.29" (crossed out from "\$49.99"), the Prime logo, and the delivery date "Get it by Wednesday, May 27".
- Logitech Professional Presenter R800 with Green Laser Pointer**: Includes an image of the presenter, the price "\$57.49" (crossed out from "\$79.99"), the Prime logo, and the delivery date "Get it by Wednesday, May 27".
- Kensington Wireless Presenter with Laser Pointer**: Includes a partial image of the presenter and the text "by Kensington".

Related Searches: logitech presenter, mpow presenter, wireless presenter.





# The Amazon Scenario

Customer info: shopping cart, account, etc.

 **Shopping Cart** Already a customer?  
[Sign in](#)

[See more items like those in your Cart](#)

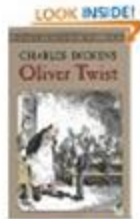
**subtotal = \$88.77**  
Make any changes below? [Update](#)

Shopping Cart Items--To Buy Now		Price:	Qty:
<small>Item added on May 22, 2009</small>	<a href="#">The Principles of Beautiful Web Design</a> - Jason Beard; <b>Paperback</b> Condition: New In Stock	<b>\$26.37</b> You Save: \$13.58 (34%)	<input type="text" value="1"/>
<a href="#">Save for later</a>	<a href="#">Delete</a>	 Eligible for FREE Super Saver Shipping	
	<a href="#">Add gift-wrap/note</a>  <a href="#">(Learn more)</a>		
<small>Item added on May 22, 2009</small>	<a href="#">Don't Make Me Think: A Common Sense Approach to Web Usability, 2nd Edition</a> - Steve Krug; <b>Paperback</b> Condition: New In Stock	<b>\$26.40</b> You Save: \$13.60 (34%)	<input type="text" value="1"/>
<a href="#">Save for later</a>	<a href="#">Delete</a>	 Eligible for FREE Super Saver Shipping	
	<a href="#">Add gift-wrap/note</a>  <a href="#">(Learn more)</a>		

# The Amazon Scenario

Recommendations, etc.:

## Customers Who Bought This Item Also Bought



Oliver Twist (Dover Thrift Editions)

> Charles Dickens

★★★★☆ (213)

Paperback

\$3.50



David Copperfield (Dover Thrift Editions)

> Charles Dickens

★★★★☆ (196)

Paperback

\$5.00



JANE EYRE

> Charlotte Bronte

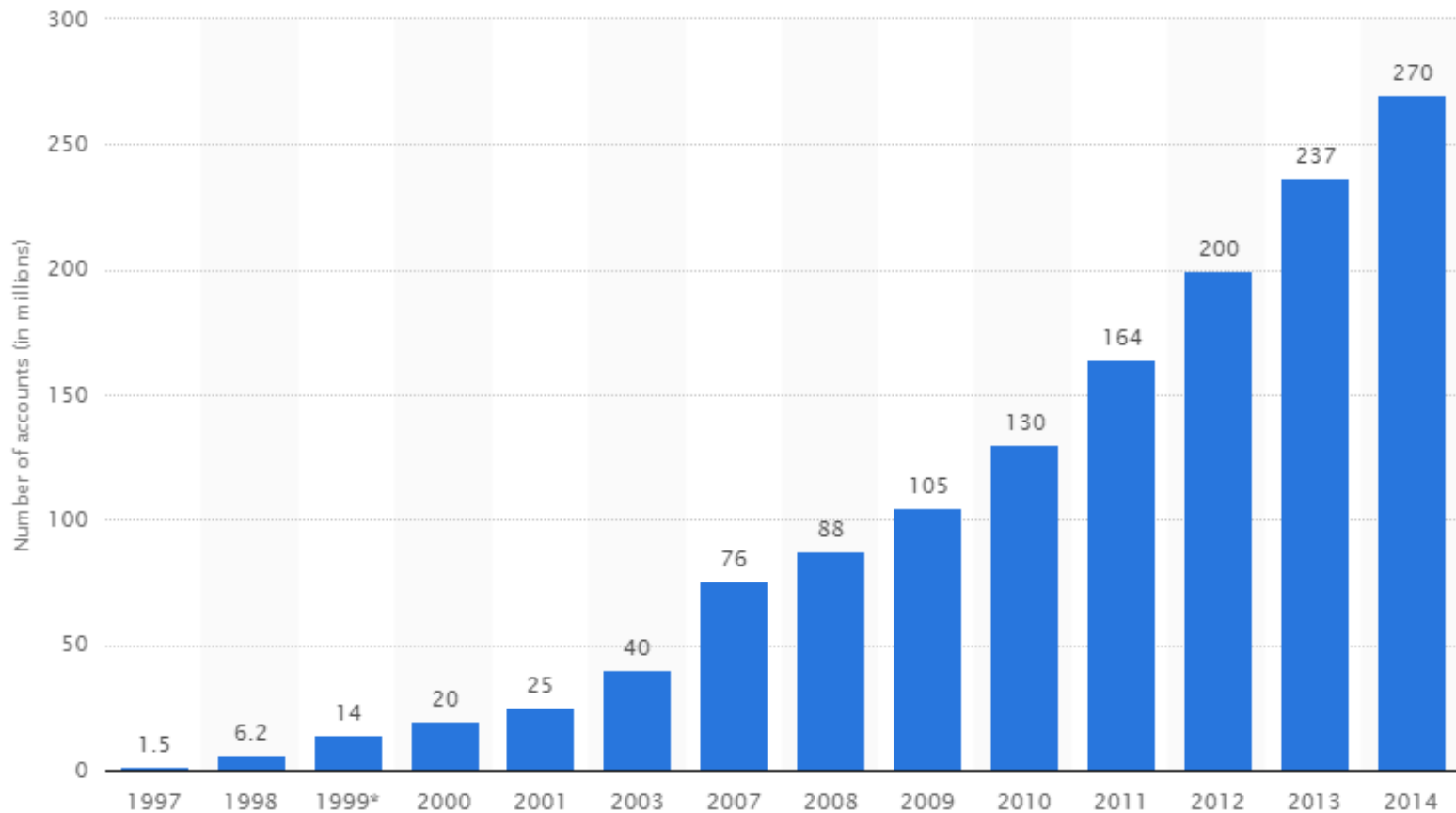
★★★★☆ (1,045)

Paperback

\$2.99

# The Amazon Scenario

- Amazon customers:

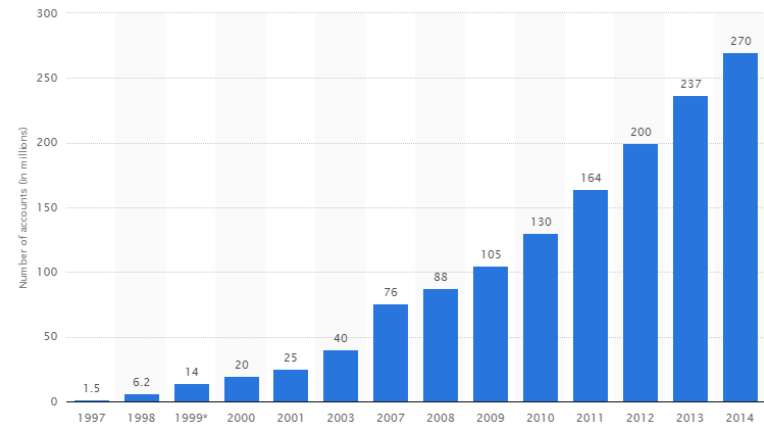


# The Amazon Scenario



# The Amazon Scenario

Databases struggling ...



But many Amazon services don't need:

- **SQL** (a simple map often enough)

or even:

- **transactions, strong consistency, etc.**



# Key–Value Store: Amazon Dynamo(DB)

## Dynamo: Amazon's Highly Available Key-value Store

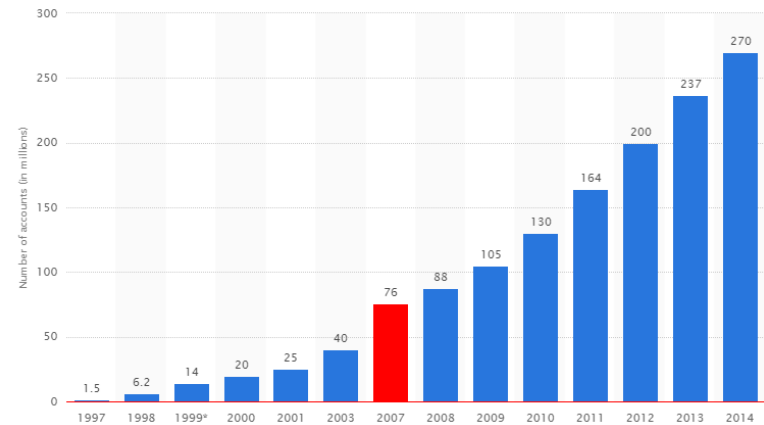
Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall and Werner Vogels

Amazon.com

### ABSTRACT

Reliability at massive scale is one of the biggest challenges we face at Amazon.com, one of the largest e-commerce operations in the world; even the slightest outage has significant financial consequences and impacts customer trust. The Amazon.com platform, which provides services for many web sites worldwide, is implemented on top of an infrastructure of tens of thousands of servers and network components located in many datacenters

One of the lessons our organization has learned from operating Amazon's platform is that the reliability and scalability of a system is dependent on how its application state is managed. Amazon uses a highly decentralized, loosely coupled, service oriented architecture consisting of hundreds of services. In this environment there is a particular need for storage technologies that are always available. For example, customers should be able to view and add items to their shopping cart even if disks are failing, network routes are changing, or data centers are being



## Goals:

- Scalability (able to grow)
- High availability (reliable)
- Performance (fast)

Don't need full SQL, don't need full ACID

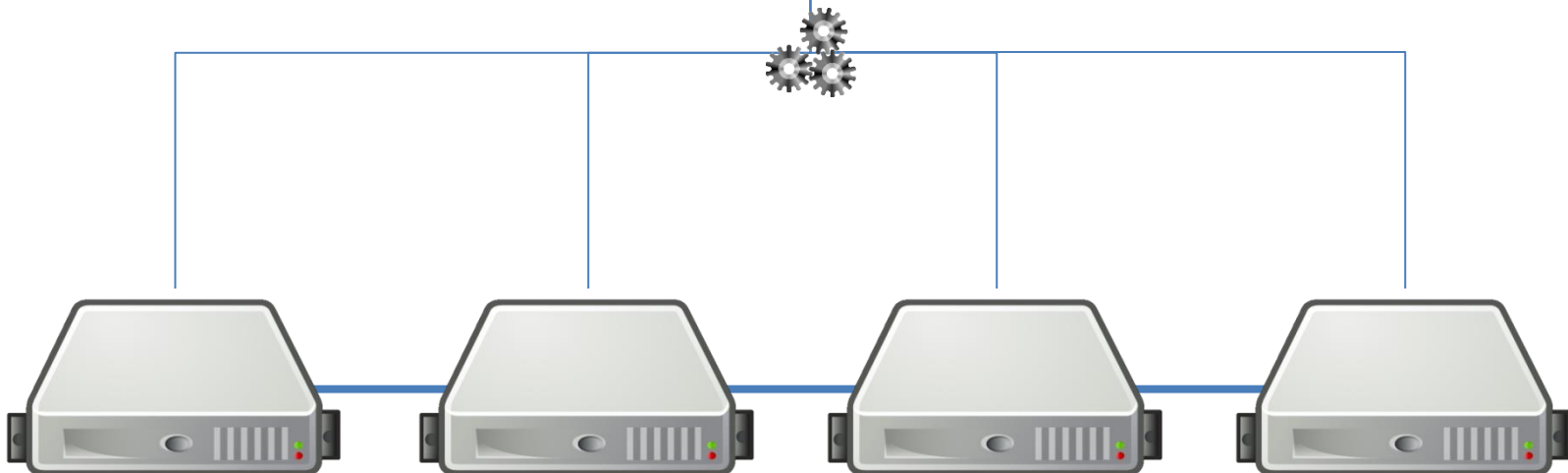
# Key–Value Store: Distribution

How might we distribute a key–value store over multiple machines?



Key	Value
country:Afghanistan	capital@city:Kabul,continent:Asia,pop:31108077#2011
country:Albania	capital@city:Tirana,continent:Europe,pop:3011405#2013
...	...
city:Kabul	country:Afghanistan,pop:3476000#2013
city:Tirana	country:Albania,pop:3011405#2013
...	...
user:10239	basedIn@city:Tirana,post:{103,10430,201}
...	...

$$\text{mod}(\text{hash}(\textit{key}), m)$$



# Key–Value Store: Distribution

What happens if a machine leaves or joins afterwards?

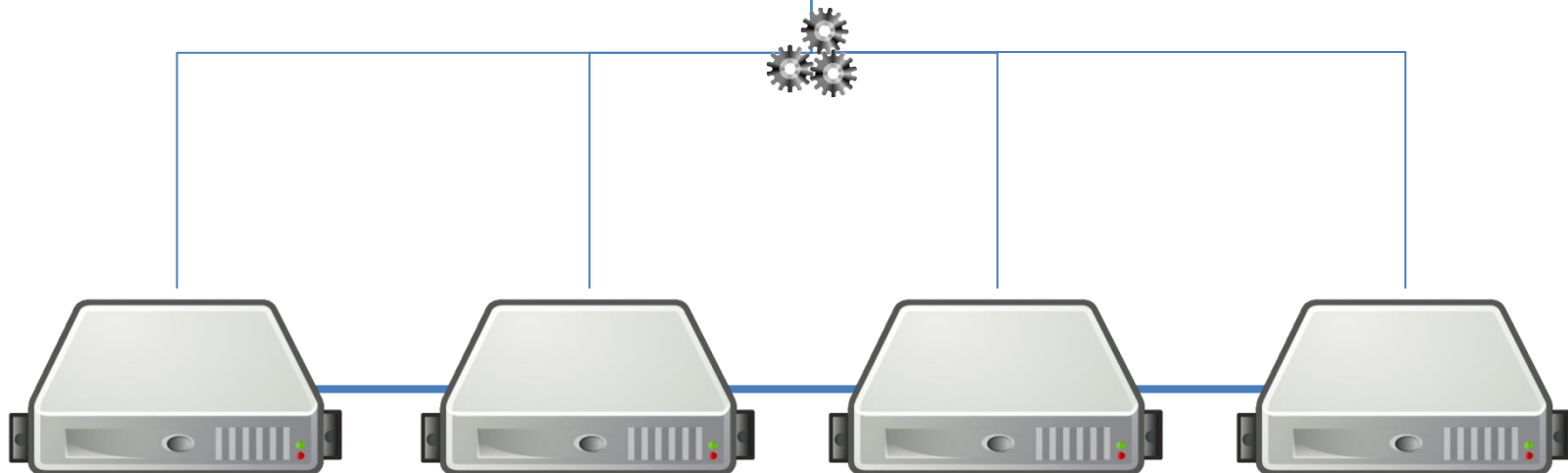


How can we avoid rehashing everything?



Key	Value
country:Afghanistan	capital@city:Kabul,continent:Asia,pop:31108077#2011
country:Albania	capital@city:Tirana,continent:Europe,pop:3011405#2013
...	...
city:Kabul	country:Afghanistan,pop:3476000#2013
city:Tirana	country:Albania,pop:3011405#2013
...	...
user:10239	basedIn@city:Tirana,post:{103,10430,201}
...	...

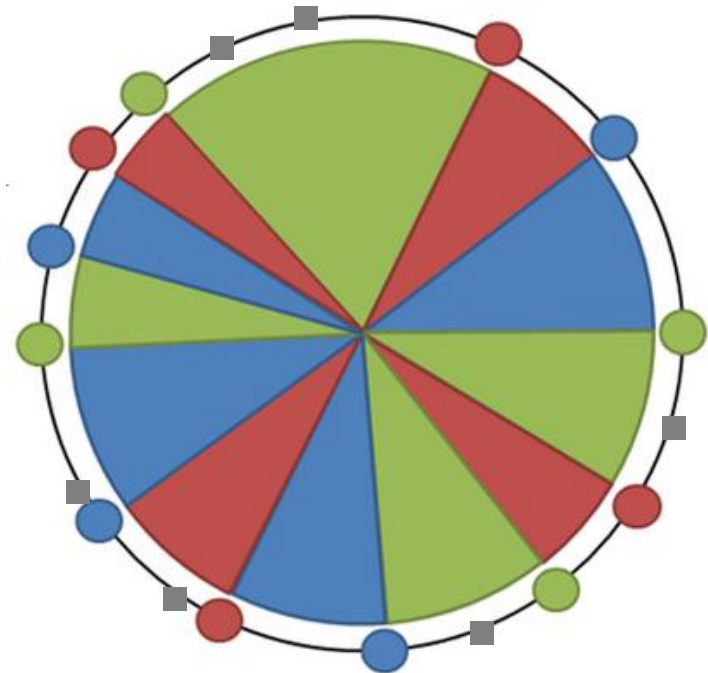
$$\text{mod}(\text{hash}(\textit{key}), m)$$



# Consistent Hashing

## Avoid re-hashing everything

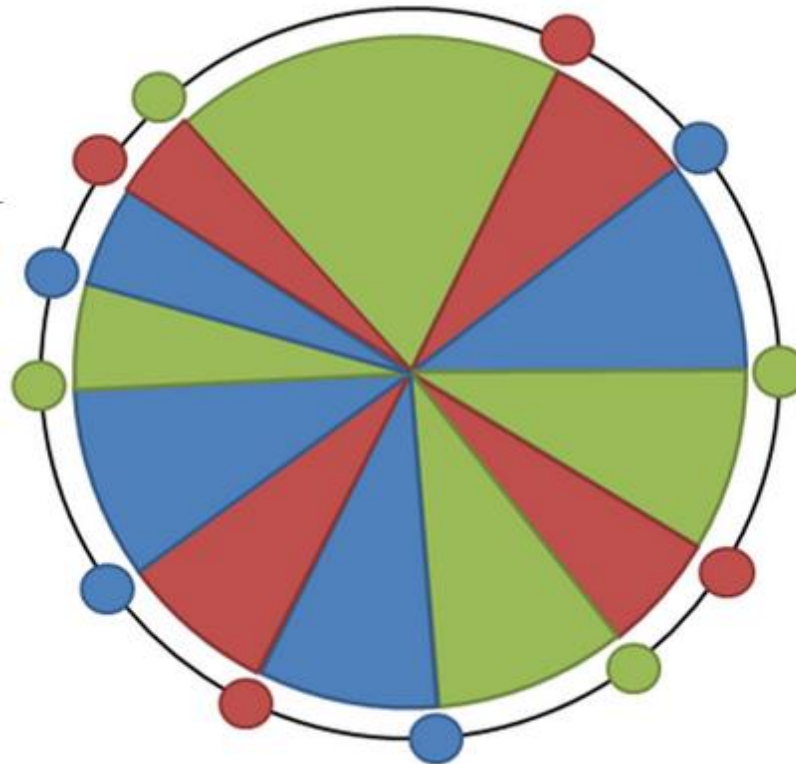
- Hash using a ring
- Each machine picks  $n$  pseudo-random points on the ring
- Machine responsible for arc after its point
- Objects mapped to ring
- If a machine leaves, its range moves to previous machine
- If a machine joins, it picks new points



# Amazon Dynamo: Hashing

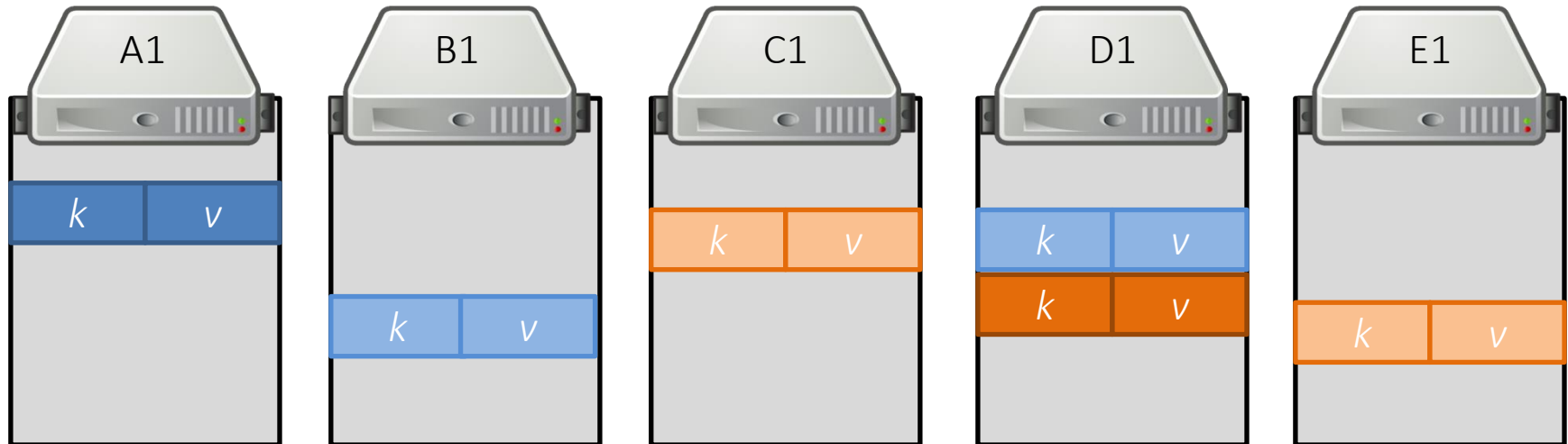


- Consistent Hashing (128-bit MD5)



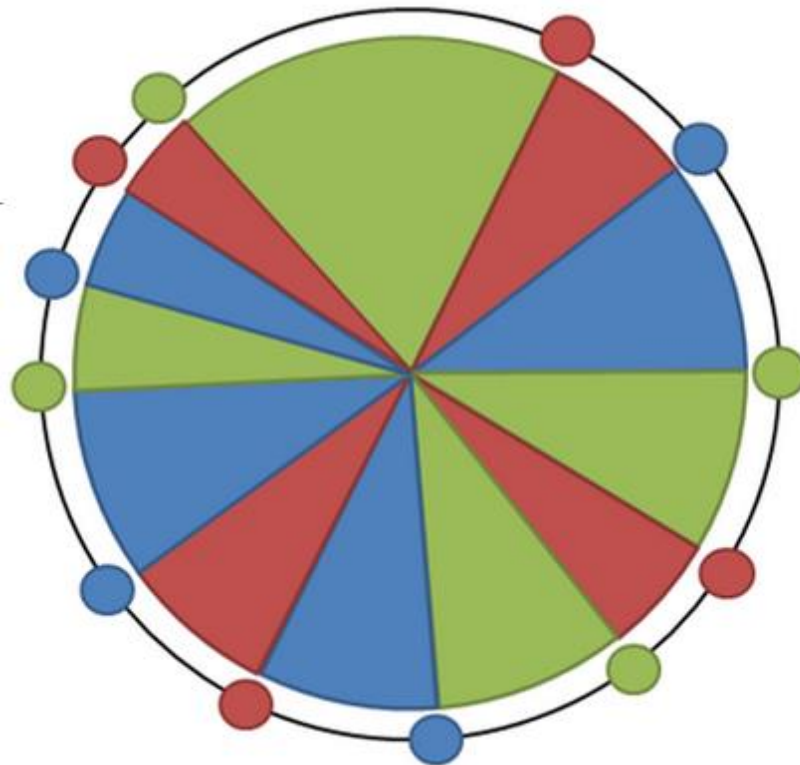
# Amazon Dynamo: Replication

- A set replication factor (e.g., 3)
- Commonly primary / secondary replicas
  - Primary replica elected from secondary replicas in the case of failure of primary



# Amazon Dynamo: Replication

- Replication factor of  $n$ ?
  - Easy: pick  $n$  next buckets (different machines!)



# Amazon Dynamo: Model

- Named table with primary key and a value
- Primary key is hashed / unordered

Countries	
Primary Key	Value
Afghanistan	capital:Kabul,continent:Asia,pop:31108077#2011
Albania	capital:Tirana,continent:Europe,pop:3011405#2013
...	...

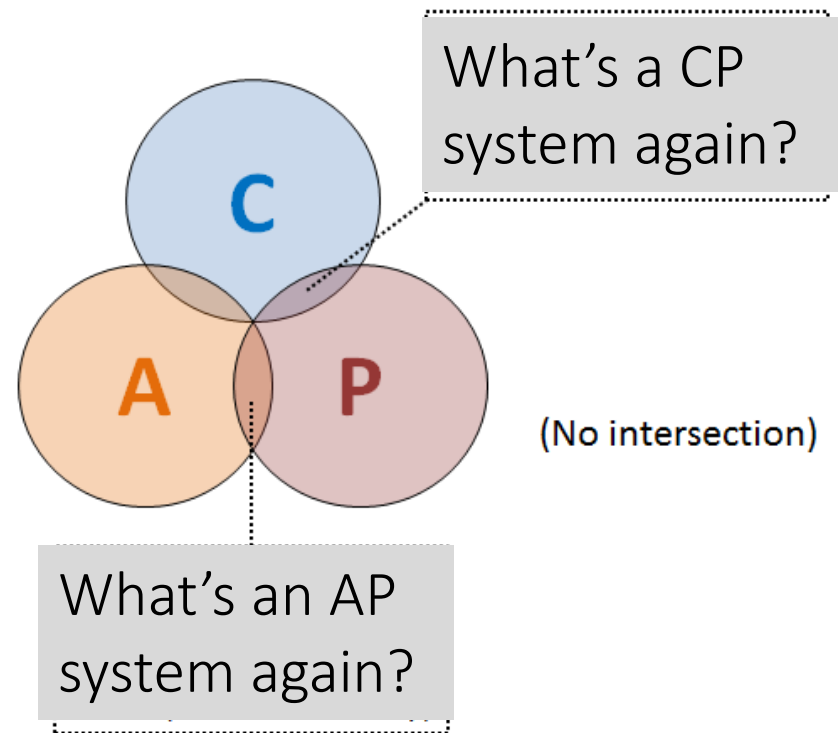
Cities	
Primary Key	Value
Kabul	country:Afghanistan,pop:3476000#2013
Tirana	country:Albania,pop:3011405#2013
...	...



# Amazon Dynamo: CAP

Two options for each table:

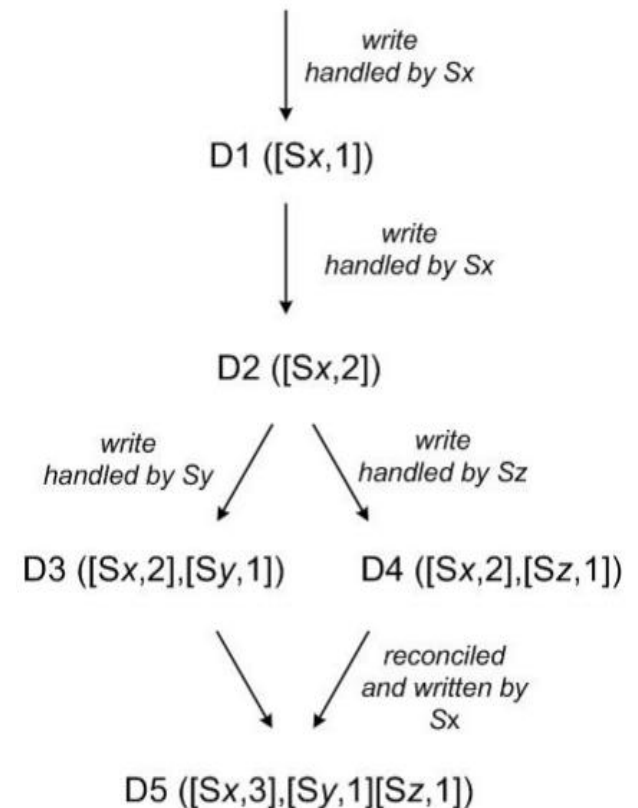
- **AP**: Eventual consistency,  
High availability
- **CP**: Strong consistency,  
Lower availability



# Amazon Dynamo: Consistency

- **Vector Clock:**


- A list of pairs indicating a node and operation count
- Used to track branches of revisions




# Amazon Dynamo: Consistency

- Two versions of one shopping cart:

## Shopping Cart

	Price	Quantity
	<b>\$90.99</b> You save: \$49.00 (35%)	1
<small>In Stock Eligible for FREE Shipping <input type="checkbox"/> This is a gift <a href="#">Learn more</a> <a href="#">Delete</a> <a href="#">Save for later</a></small>		
	<b>\$44.29</b> You save: \$5.70 (11%)	1
<small>Logitech In Stock Eligible for FREE Shipping <input type="checkbox"/> This is a gift <a href="#">Learn more</a> <a href="#">Delete</a> <a href="#">Save for later</a></small>		
<b>Subtotal (2 items): \$135.28</b>		
Total savings: <b>\$54.70</b>		

## Shopping Cart

	Price	Quantity
	<b>\$99.00</b> You save: \$30.00 (23%)	1
<small>Only 2 left in stock. Shipped from: Sam Ash Gift options not available. <a href="#">Learn more</a> <a href="#">Delete</a> <a href="#">Save for later</a></small>		
	<b>\$44.29</b> You save: \$5.70 (11%)	1
<small>Logitech In Stock Eligible for FREE Shipping <input type="checkbox"/> This is a gift <a href="#">Learn more</a> <a href="#">Delete</a> <a href="#">Save for later</a></small>		
<b>Subtotal (2 items): \$143.29</b>		
Total savings: <b>\$35.70</b>		


How best to merge multiple conflicting versions of a value  
(known as reconciliation)?




Application knows best

(... and must support multiple versions being returned)

# Amazon Dynamo: Consistency



Key	Value
country:Afghanistan	capital@city:Kabul,continent:Asia,pop:31108077#2011
country:Albania	capital@city:Tirana,continent:Europe,pop:3011405#2013
...	...
city:Kabul	country:Afghanistan,pop:3476000#2013
city:Tirana	country:Albania,pop:3011405#2013
...	...
user:10239	basedIn@city:Tirana,post:{103,10430,201,408}
...	...



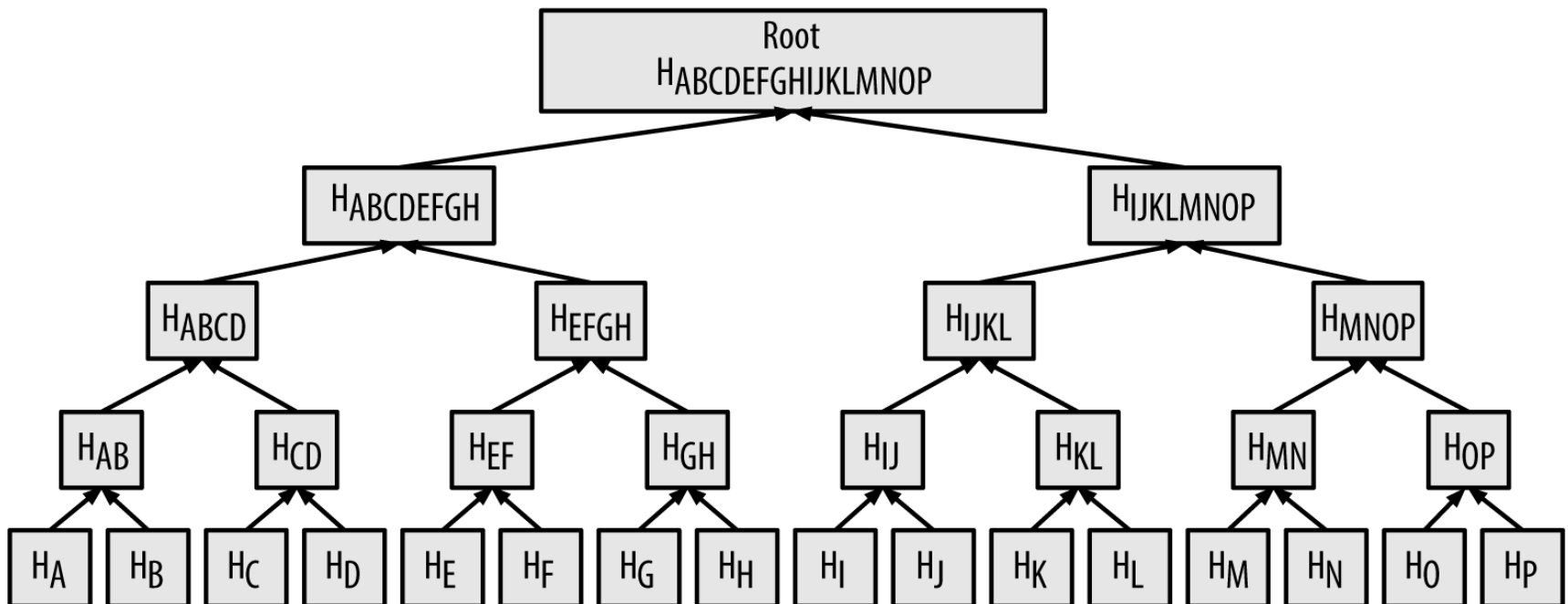
Key	Value
country:Afghanistan	capital@city:Kabul,continent:Asia,pop:31108077#2011
country:Albania	capital@city:Tirana,continent:Europe,pop:3011405#2013
...	...
city:Kabul	country:Afghanistan,pop:3476000#2013
city:Tirana	country:Albania,pop:3011405#2013
...	...
user:10239	basedIn@city:Tirana,post:{103,10430,201}
...	...

How can we efficiently verify that two copies of a block of data are the same (and find where the differences are)?



# Amazon Dynamo: Merkle Trees

- **Merkle tree:** A hash tree
  - Leaf node compute hashes from data
  - Non-leaf nodes have hashes of their children
  - Find differences between two trees level-by-level



Read More ...



## **Dynamo: Amazon's Highly Available Key-value Store**

Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati,  
Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall  
and Werner Vogels

Amazon.com

### **ABSTRACT**

Reliability at massive scale is one of the biggest challenges we face at Amazon.com, one of the largest e-commerce operations in the world; even the slightest outage has significant financial consequences and impacts customer trust. The Amazon.com platform, which provides services for many web sites worldwide, is implemented on top of an infrastructure of tens of thousands of servers and network components located in many datacenters

One of the lessons our organization has learned from operating Amazon's platform is that the reliability and scalability of a system is dependent on how its application state is managed. Amazon uses a highly decentralized, loosely coupled, service oriented architecture consisting of hundreds of services. In this environment there is a particular need for storage technologies that are always available. For example, customers should be able to view and add items to their shopping cart even if disks are failing, network routes are flapping, or data centers are being

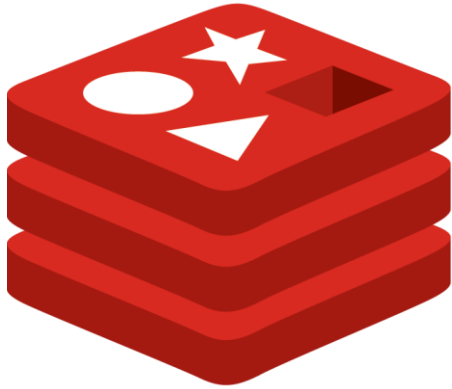
# OTHER KEY-VALUE STORES

# Other Key–Value Stores





# Other Key–Value Stores



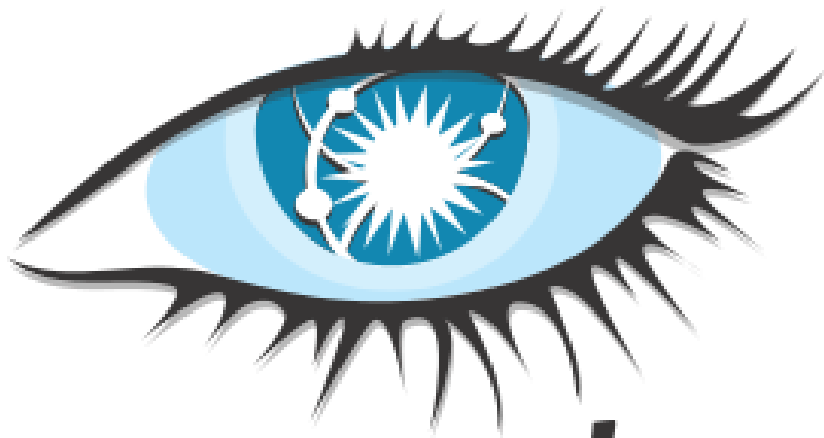
redis



StackExchange 



# Other Key-Value Stores



***cassandra***



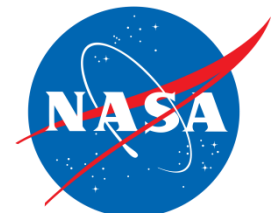
**Instagram**  
Fast beautiful photo sharing

**accenture**  
High performance. Delivered.

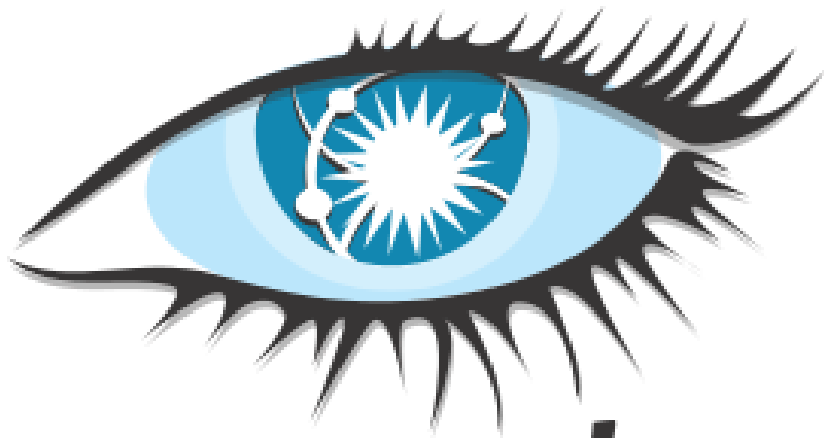
**Answers.com**



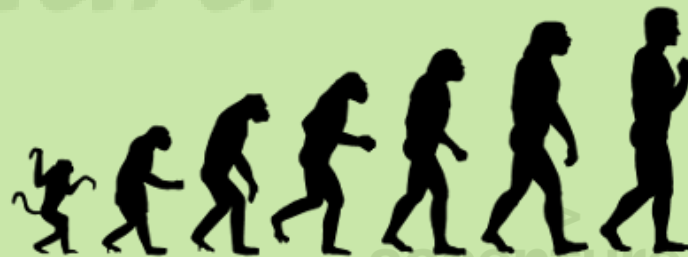
**Disney**



# Other Key-Value Stores



cassandra



Instagram  
Fast beautiful

accenture  
High performance. Delivered.

Answers.com

Evolved into a  
tabular store ...

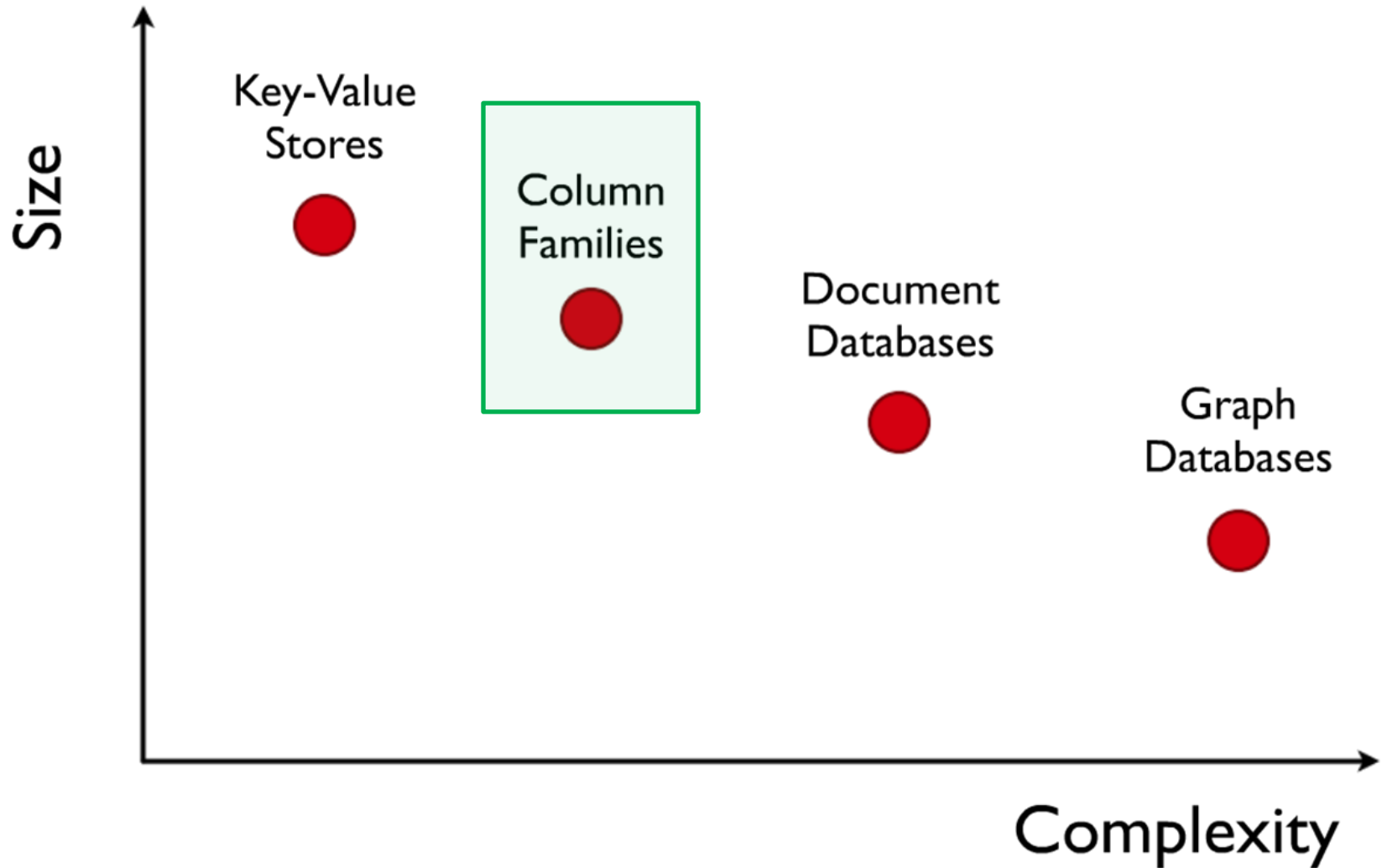


DISNEY



TABULAR / COLUMN FAMILY

# NoSQL: Column Family Stores



# Key–Value = a Distributed Map

Countries	
Primary Key	Value
Afghanistan	capital:Kabul,continent:Asia,pop:31108077#2011
Albania	capital:Tirana,continent:Europe,pop:3011405#2013
...	...

# Tabular = Multi-dimensional Maps

Countries				
Primary Key	capital	continent	pop-value	pop-year
Afghanistan	Kabul	Asia	31108077	2011
Albania	Tirana	Europe	3011405	2013
...	...	...	...	...

# Bigtable: The Original Whitepaper

MapReduce  
authors

## Bigtable: A Distributed Storage System for Structured Data

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach,  
Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber

{fay,jeff,sanjay,wilsonh,kerr,m3b,tushar,fikes,gruber}@google.com

*Google, Inc.*

### Abstract

Bigtable is a distributed storage system for managing structured data that is designed to scale to a very large size: petabytes of data across thousands of commodity servers. Many projects at Google store data in Bigtable, including web indexing, Google Earth, and Google Finance. These applications place very different demands on Bigtable, both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data service). Despite these varied demands, Bigtable has successfully provided a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over data layout and format, and we describe the design and implementation of Bigtable.

achieved scalability and high performance, but Bigtable provides a different interface than such systems. Bigtable does not support a full relational data model; instead, it provides clients with a simple data model that supports dynamic control over data layout and format, and allows clients to reason about the locality properties of the data represented in the underlying storage. Data is indexed using row and column names that can be arbitrary strings. Bigtable also treats data as uninterpreted strings, although clients often serialise various forms of structured and semi-structured data into these strings. Clients can control the locality of their data through careful choices in their schemas. Finally, Bigtable schema parameters let clients dynamically control whether to serve data out of memory or from disk.

Section 2 describes the data model in more detail, and Section 3 provides an overview of the client API. Sec-

# Bigtable used for ...



Google Analytics



+Amit Web Images Videos Maps News Shopping Gmail More

Amit Singhal 2 Share

Google chikoo

Search 50 personal results and 419,000 other results (0.61 seconds)

Manilkara zapota - Wikipedia, the free encyclopedia  
en.wikipedia.org/wiki/Manilkara\_zapota  
Sapodilla is known as **chikoo** ("चिकू," or chiku, "चोफ़,") in India and Pakistan and sapota in some parts of India (Tamil Nadu, Kerala, Karnataka, Andhra ...  
Description - Other names - See also - References  
You've visited this page 3 times. Last visit 12/4/11

Images for chikoo - Report images

You	Shilpa Singhal	You	Shilpa Singhal

Chikoo - a simple file organizer for the Mac  
codingturtle.com/chikoo/  
Chikoo, a simple file organizer for the Mac. Download (30-day trial). Version 0.9.1 ...



orkut by Google™





# Bigtable: Sorted Keys

Primary Key	capital		pop-value		pop-year	
Asia:Afghanistan	t <sub>1</sub>	Kabul	t <sub>1</sub>	31143292	t <sub>1</sub>	2009
			t <sub>2</sub>	31120978		
			t <sub>4</sub>	31108077	t <sub>4</sub>	2011
Asia:Azerbaijan	...	...	...	...	...	...
...	...	...	...	...	...	...
Europe:Albania	t <sub>1</sub>	Tirana	t <sub>1</sub>	2912380	t <sub>1</sub>	2010
			t <sub>3</sub>	3011405	t <sub>3</sub>	2013
Europe:Andorra	...	...	...	...	...	...
...	...	...	...	...	...	...

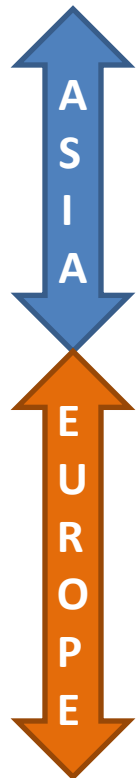
S  
O  
R  
T  
E  
D

Benefits of sorted vs. hashed keys?



Range queries and ...

# Bigtable: Tablets



Primary Key	capital		pop-value		pop-year	
Asia:Afghanistan	t <sub>1</sub>	Kabul	t <sub>1</sub>	31143292	t <sub>1</sub>	2009
			t <sub>2</sub>	31120978		
			t <sub>4</sub>	31108077	t <sub>4</sub>	2011
Asia:Azerbaijan	...	...	...	...	...	...
...	...	...	...	...	...	...
Europe:Albania	t <sub>1</sub>	Tirana	t <sub>1</sub>	2912380	t <sub>1</sub>	2010
			t <sub>3</sub>	3011405	t <sub>3</sub>	2013
Europe:Andorra	...	...	...	...	...	...
...	...	...	...	...	...	...

Benefits of sorted vs. hashed keys?



Range queries and ...

... locality of processing

# A real-world example of locality/sorting



Primary Key	language		title		links	
com.imdb	t <sub>1</sub>	en	t <sub>1</sub>	IMDb Home	t <sub>1</sub>	...
			t <sub>2</sub>	IMDB - Movies		
			t <sub>4</sub>	IMDb	t <sub>4</sub>	...
com.imdb/title/tt2724064/	t <sub>1</sub>	en	t <sub>2</sub>	Sharknado	t <sub>2</sub>	...
com.imdb/title/tt3062074/	t <sub>1</sub>	en	t <sub>2</sub>	Sharknado II	t <sub>2</sub>	...
...	...	...	...	...	...	...
org.wikipedia	t <sub>1</sub>	multi	t <sub>1</sub>	Wikipedia	t <sub>1</sub>	...
			t <sub>3</sub>	Wikipedia Home	t <sub>3</sub>	...
org.wikipedia.ace	t <sub>1</sub>	ace	t <sub>1</sub>	Wikipèdia bahsa Acèh	...	...
...	...	...	...	...	...	...

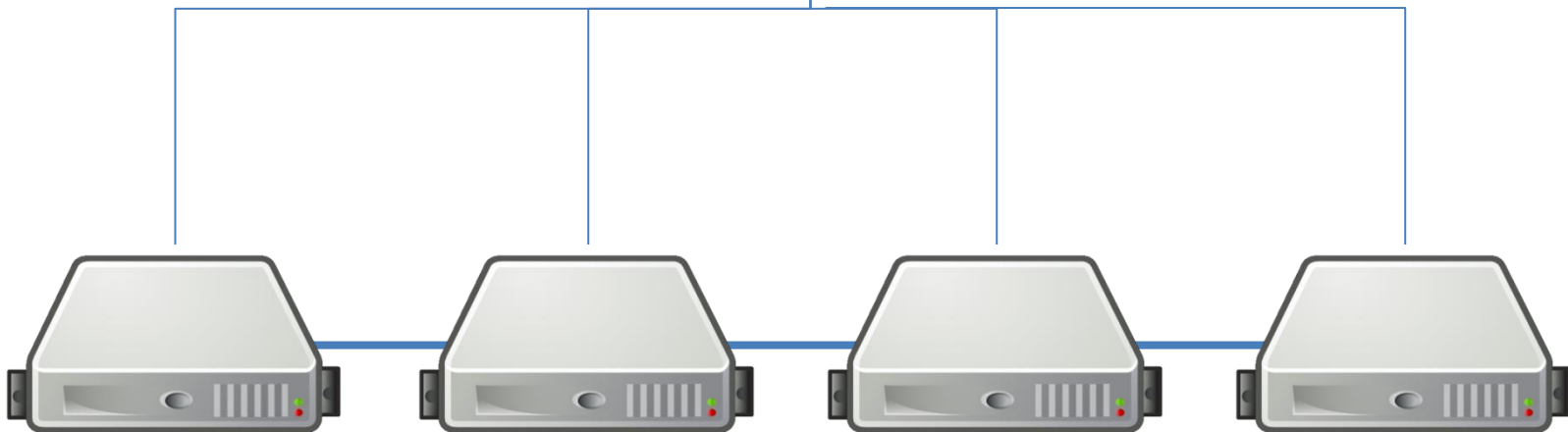


# Bigtable: Distribution

Asia:Afghanistan	t <sub>1</sub>	Kabul	t <sub>1</sub>	31143292	t <sub>1</sub>	2009
			t <sub>2</sub>	31120978		
			t <sub>4</sub>	31108077	t <sub>4</sub>	2011
Asia:Azerbaijan	...	...	...	...	...	...
...	...	...	...	...	...	...

Europe:Albania	t <sub>1</sub>	Tirana	t <sub>1</sub>	2912380	t <sub>1</sub>	2010
			t <sub>3</sub>	3011405	t <sub>3</sub>	2013
Europe:Andorra	...	...	...	...	...	...
...	...	...	...	...	...	...

Split by tablet



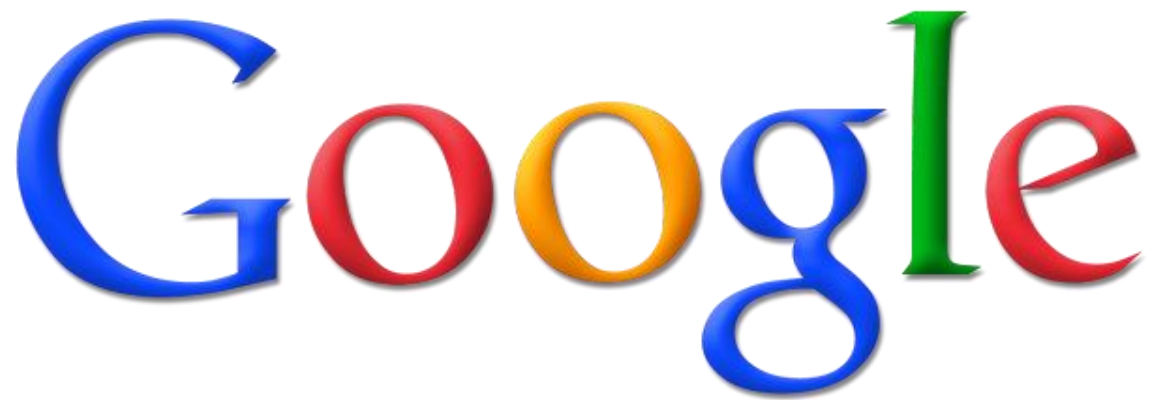
Horizontal range partitioning

# Bigtable: Column Families

Primary Key	pol:capital		demo:pop-value		demo:pop-year	
Asia:Afghanistan	t <sub>1</sub>	Kabul	t <sub>1</sub>	31143292	t <sub>1</sub>	2009
			t <sub>2</sub>	31120978		
			t <sub>4</sub>	31108077	t <sub>4</sub>	2011
Asia:Azerbaijan	...	...	...	...	...	...
...	...	...	...	...	...	...
Europe:Albania	t <sub>1</sub>	Tirana	t <sub>1</sub>	2912380	t <sub>1</sub>	2010
			t <sub>3</sub>	3011405	t <sub>3</sub>	2013
Europe:Andorra	...	...	...	...	...	...
...	...	...	...	...	...	...

- Group logically similar columns together
  - Accessed efficiently together
  - Access-control and storage: column family level
  - If of same type, can be compressed

Read More ...



## **Bigtable: A Distributed Storage System for Structured Data**

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach  
Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber

{fay,jeff,sanjay,wilsonh,kerr,m3b,tushar,fikes,gruber}@google.com

*Google, Inc.*

### **Abstract**

Bigtable is a distributed storage system for managing structured data that is designed to scale to a very large size: petabytes of data across thousands of commodity servers. Many projects at Google store data in Bigtable, including web indexing, Google Earth, and Google Finance. These applications place very different demands on Bigtable, both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data serving). Despite these varied demands, Bigtable has successfully provided a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over data layout and format, and we describe the design and implementation of Bigtable.

achieved scalability and high performance, but Bigtable provides a different interface than such systems. Bigtable does not support a full relational data model; instead, it provides clients with a simple data model that supports dynamic control over data layout and format, and allows clients to reason about the locality properties of the data represented in the underlying storage. Data is indexed using row and column names that can be arbitrary strings. Bigtable also treats data as uninterpreted strings, although clients often serialize various forms of structured and semi-structured data into these strings. Clients can control the locality of their data through careful choices in their schemas. Finally, Bigtable schema parameters let clients dynamically control whether to serve data out of memory or from disk.

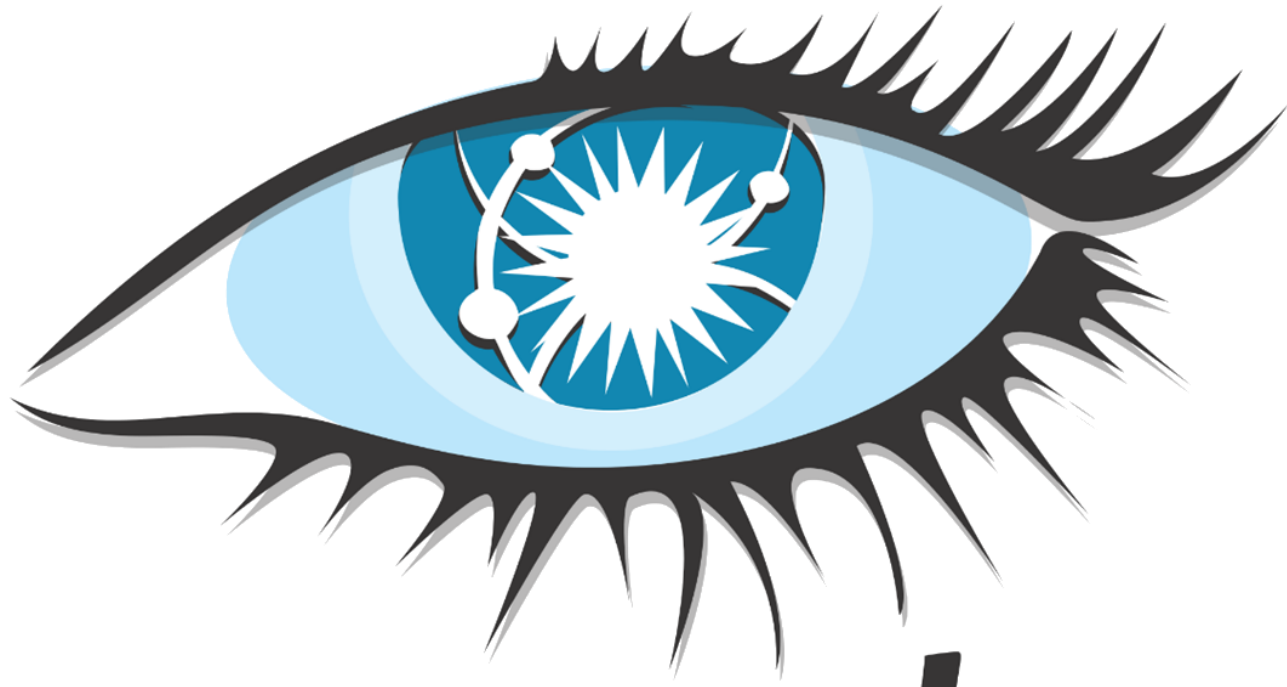
Section 2 describes the data model in more detail, and

Tabular Store: Apache HBase





Tabular Store: Cassandra



***cassandra***



Questions?