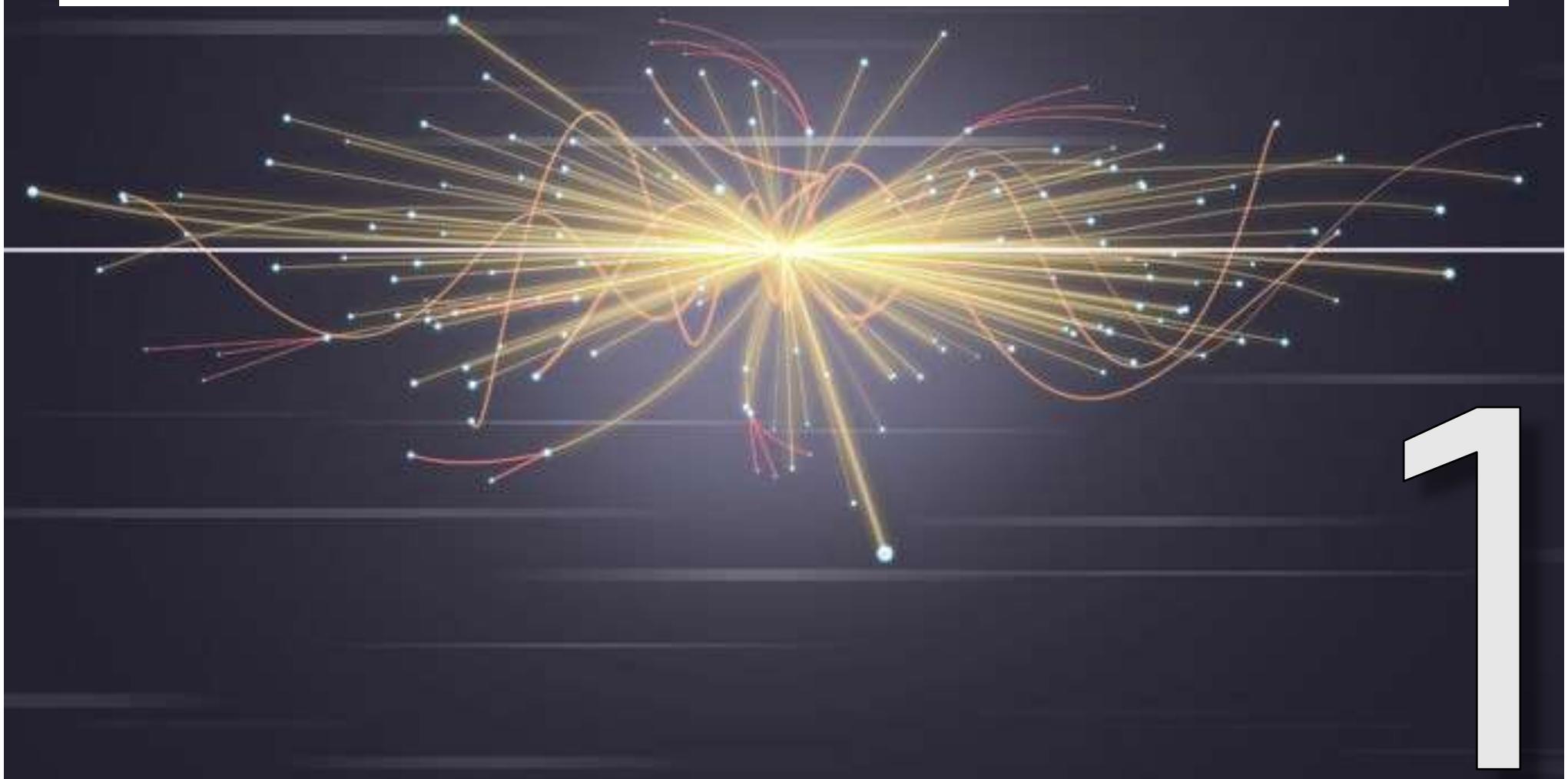


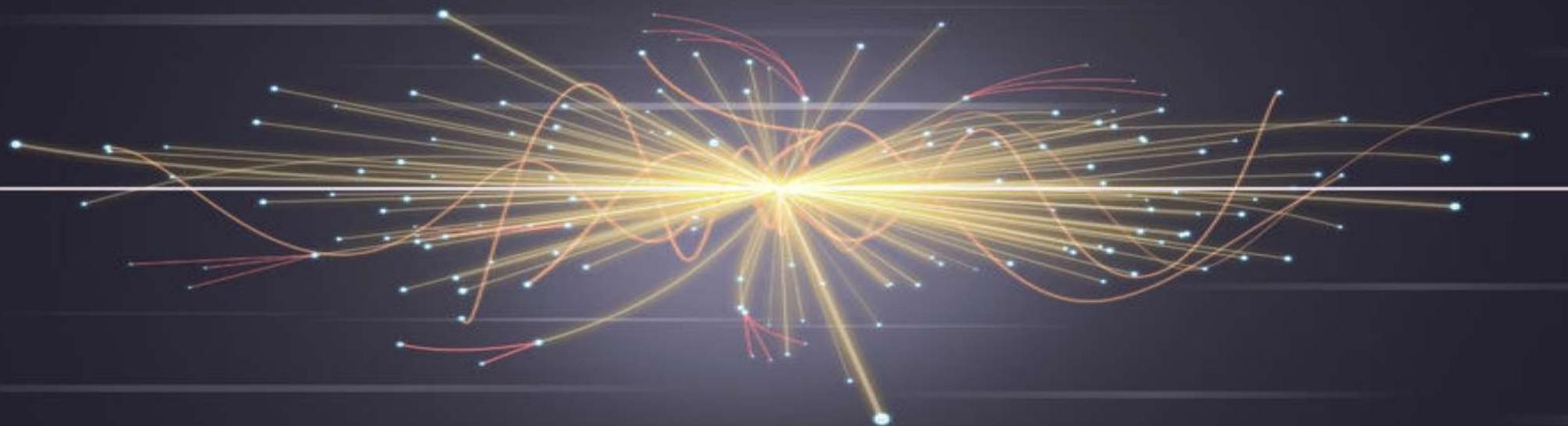
Ghislain Fourny

Big Data for Engineers Spring 2020

1. Introduction

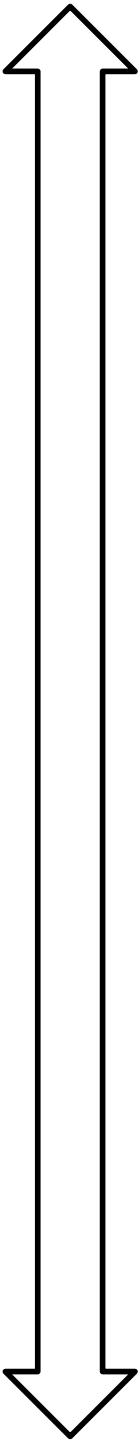


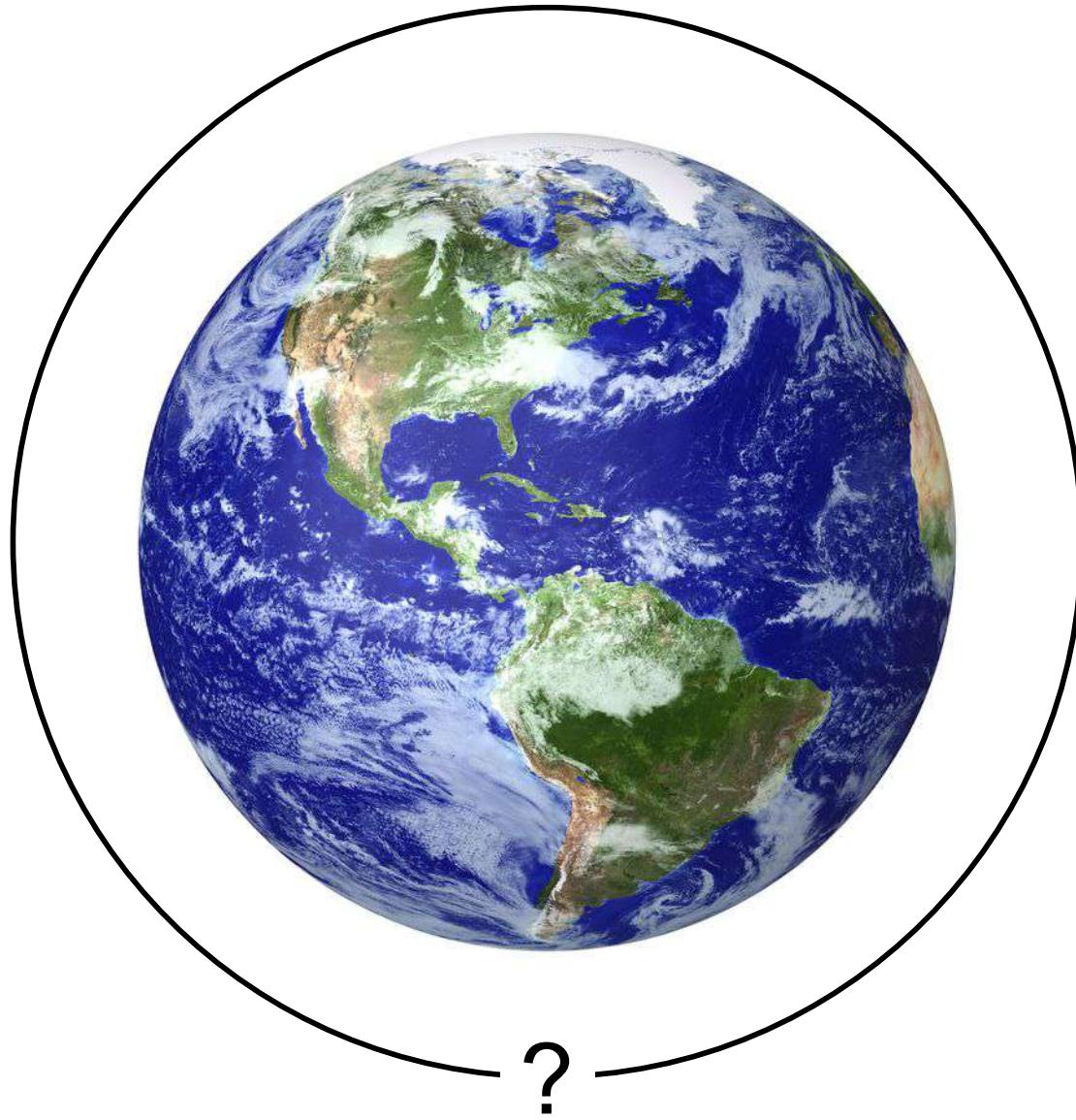
Scale

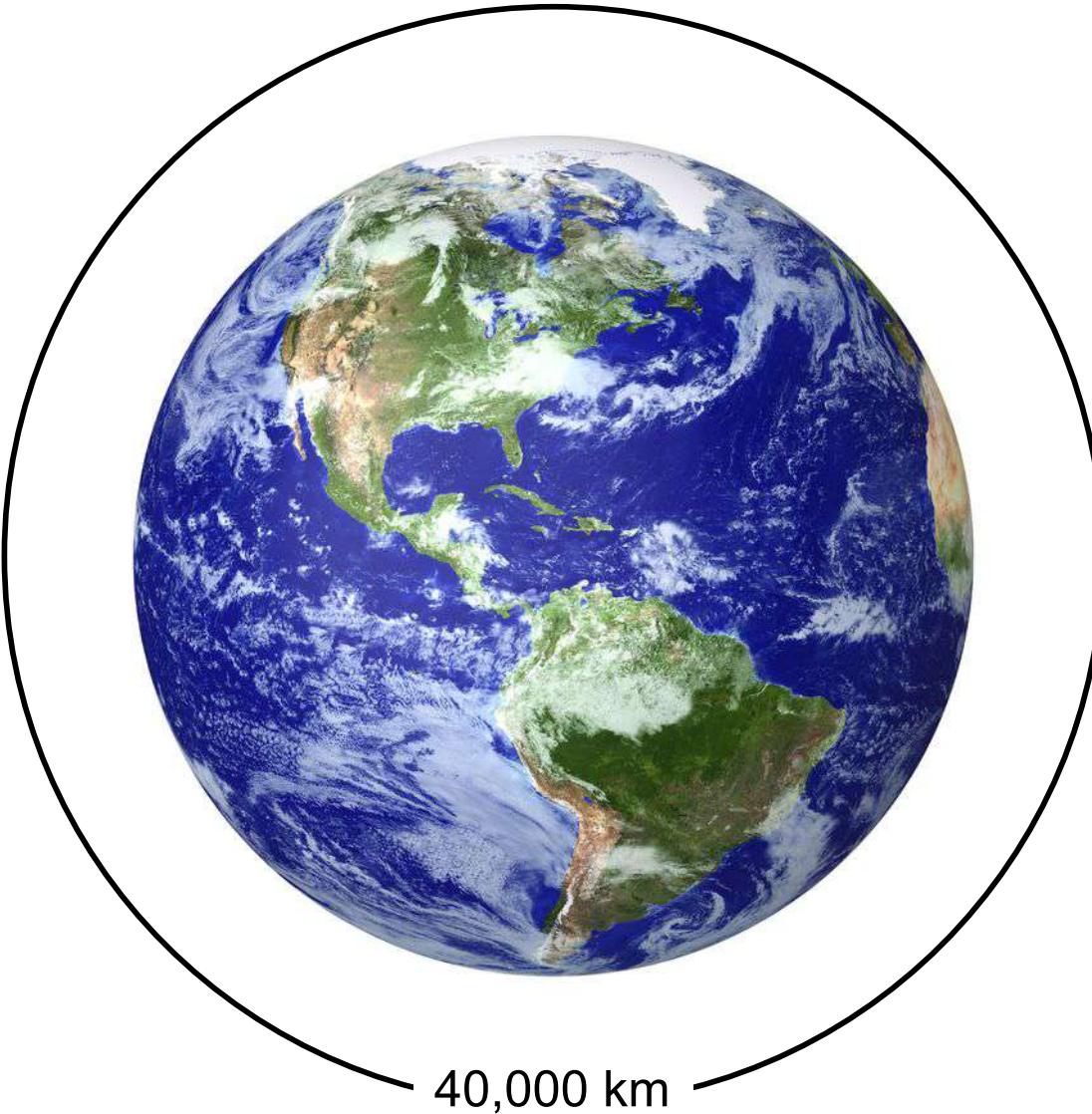


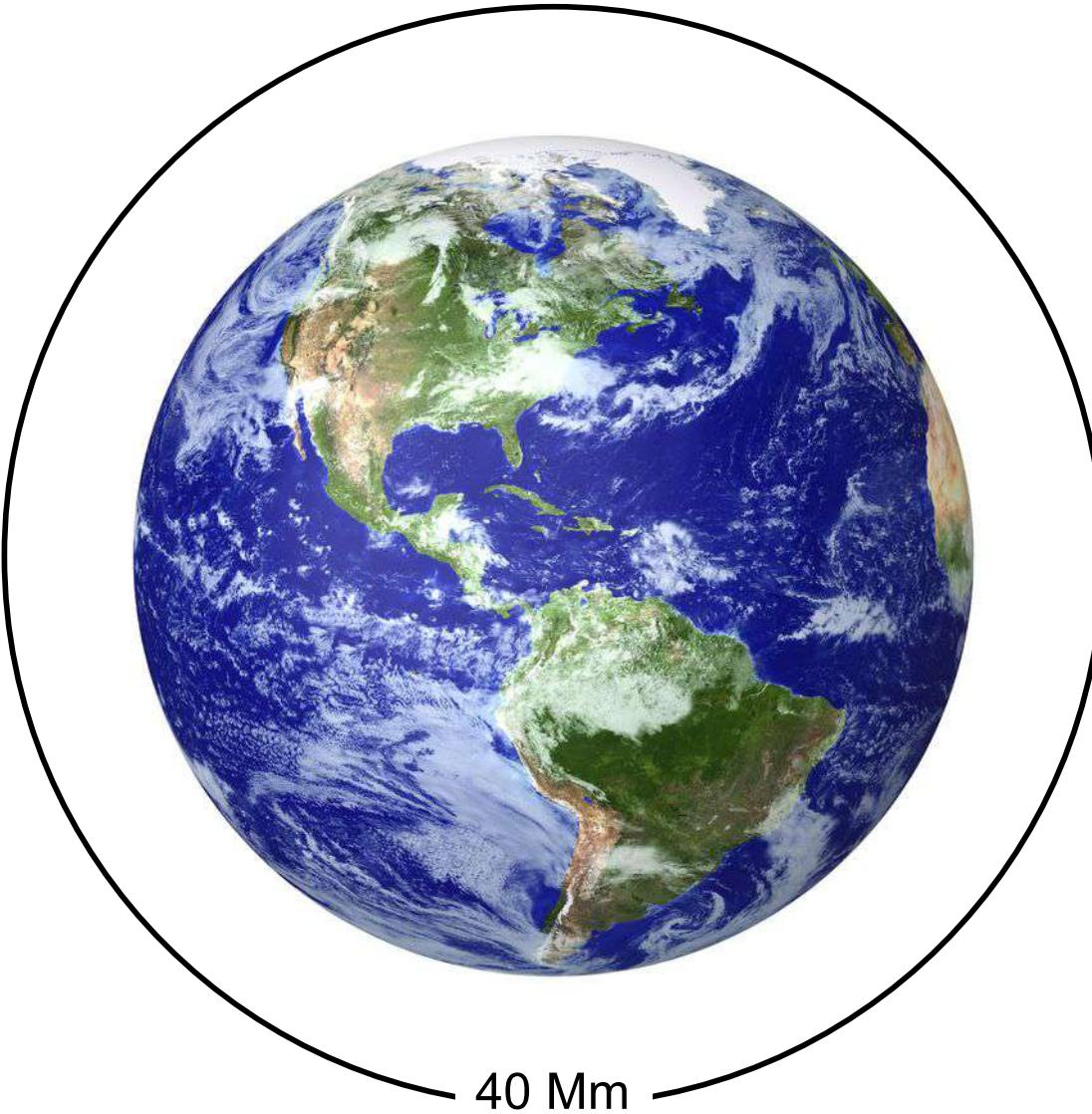


1.85m

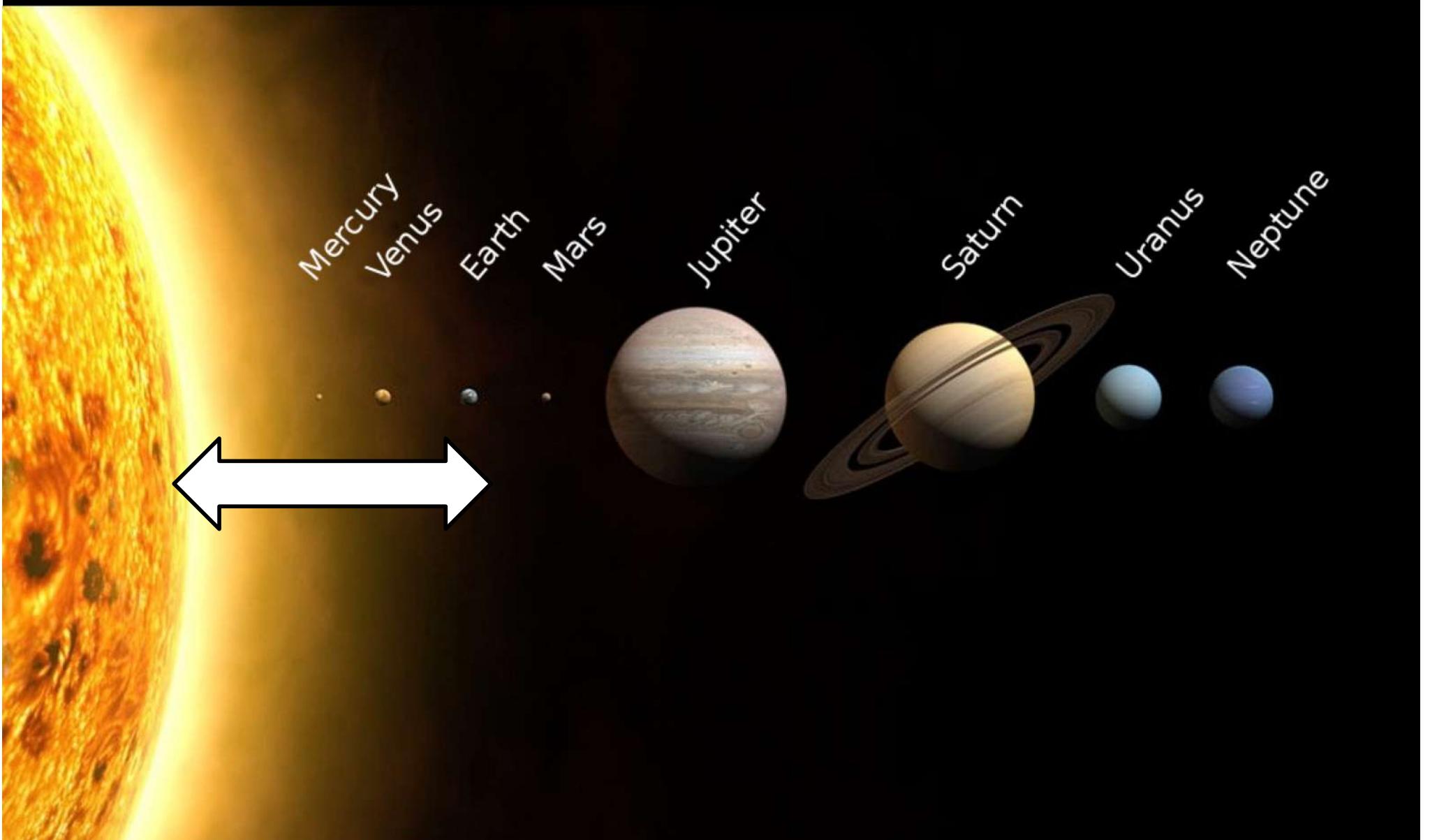


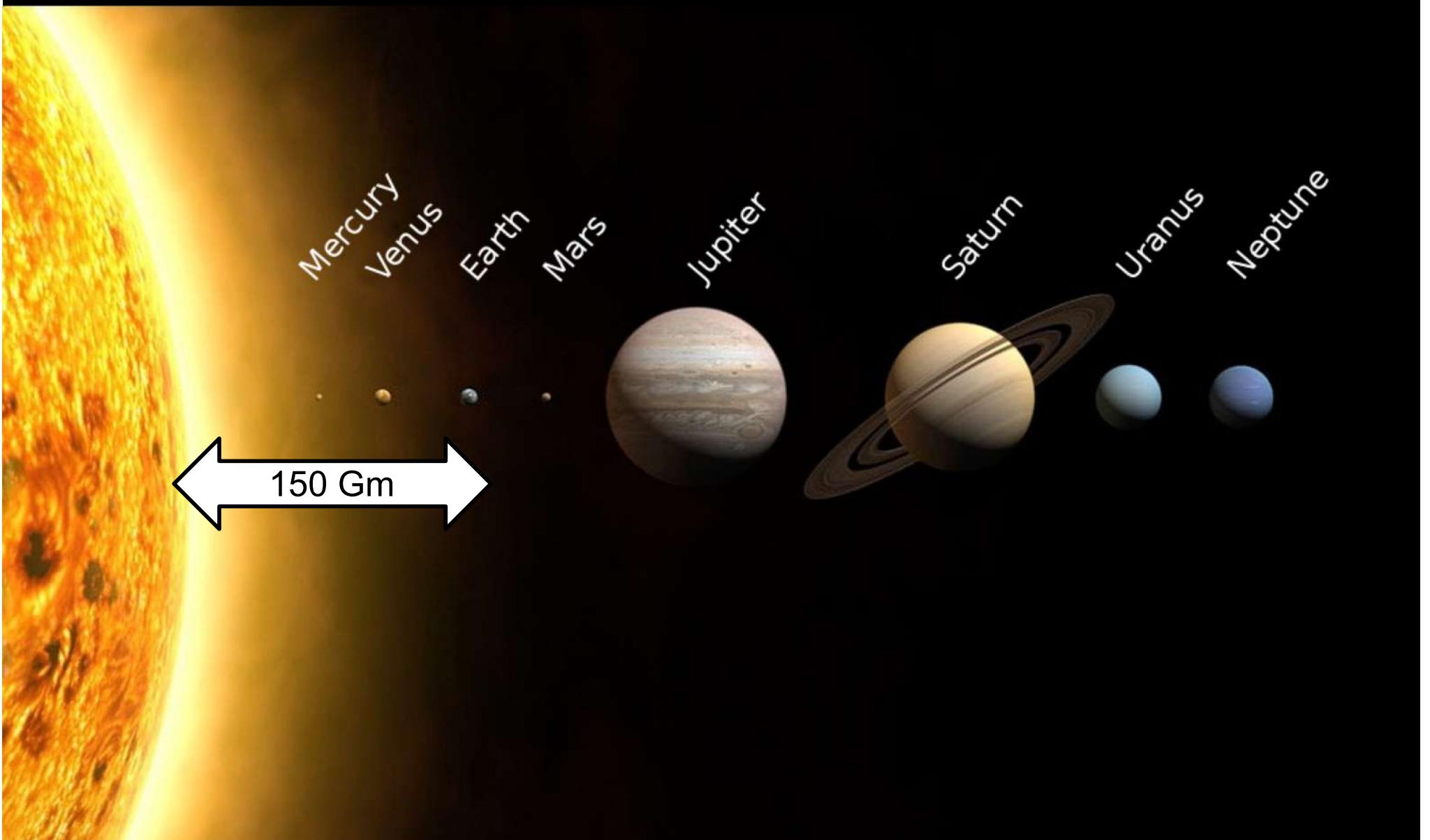




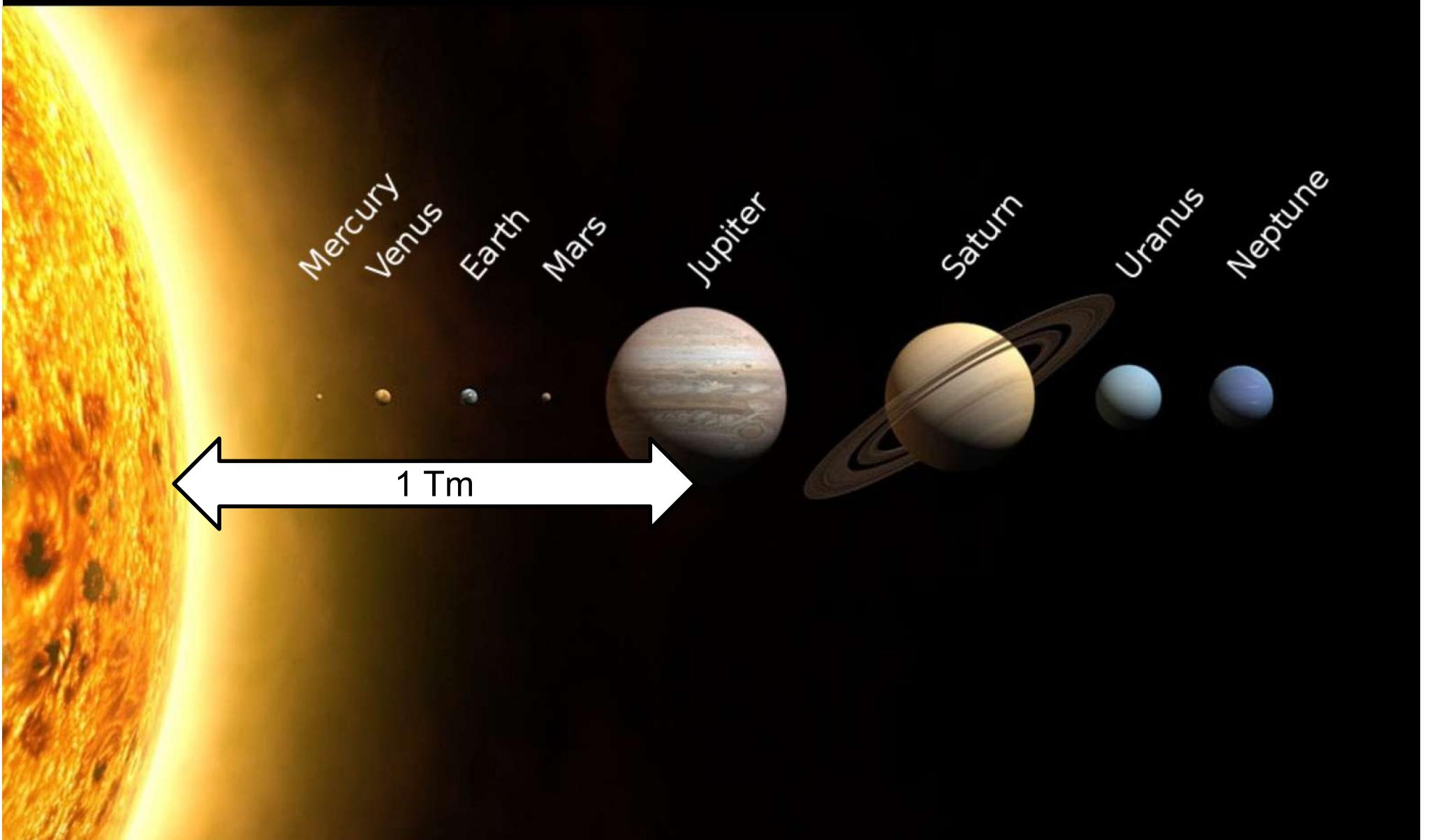


40 Mm

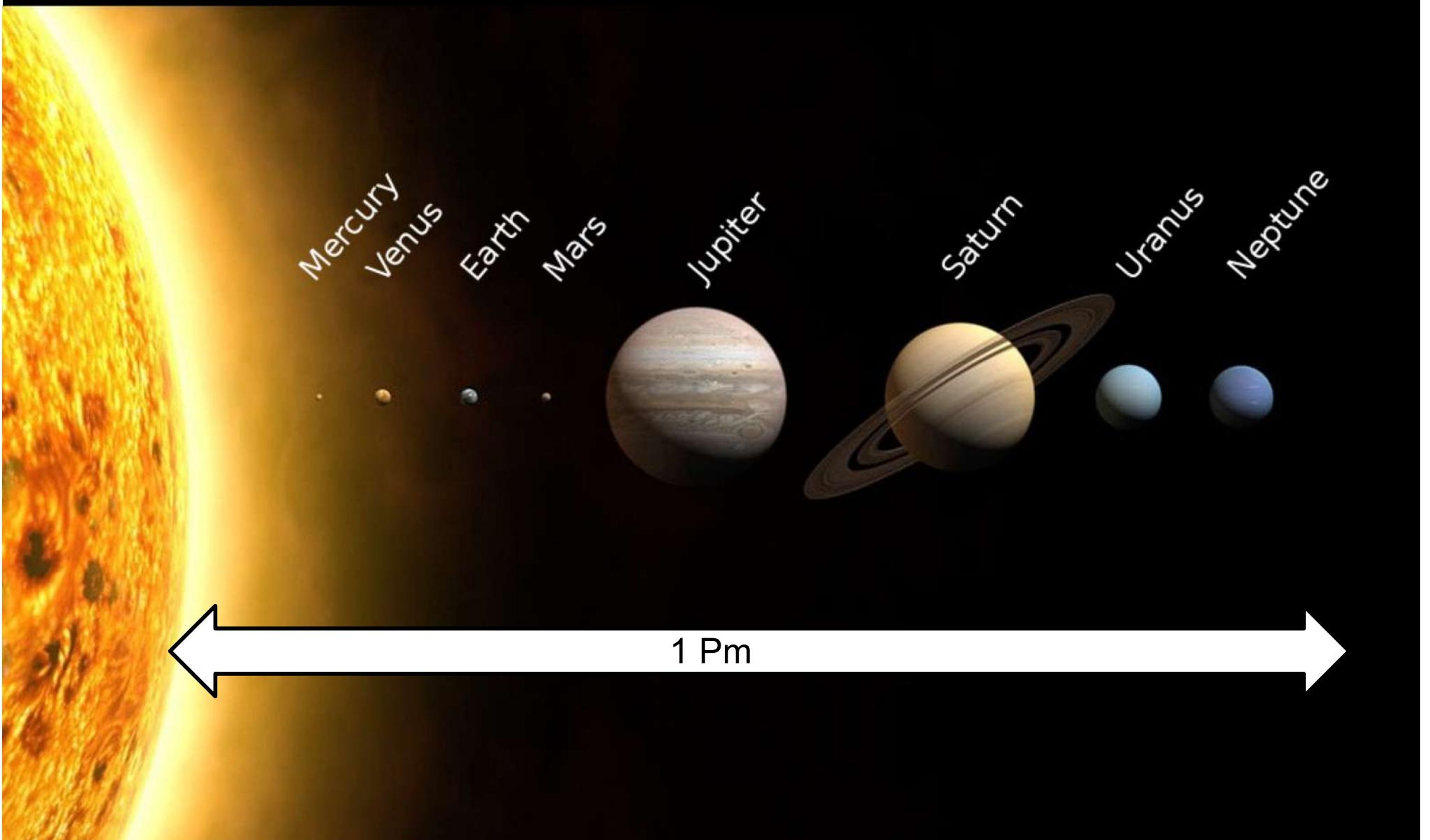




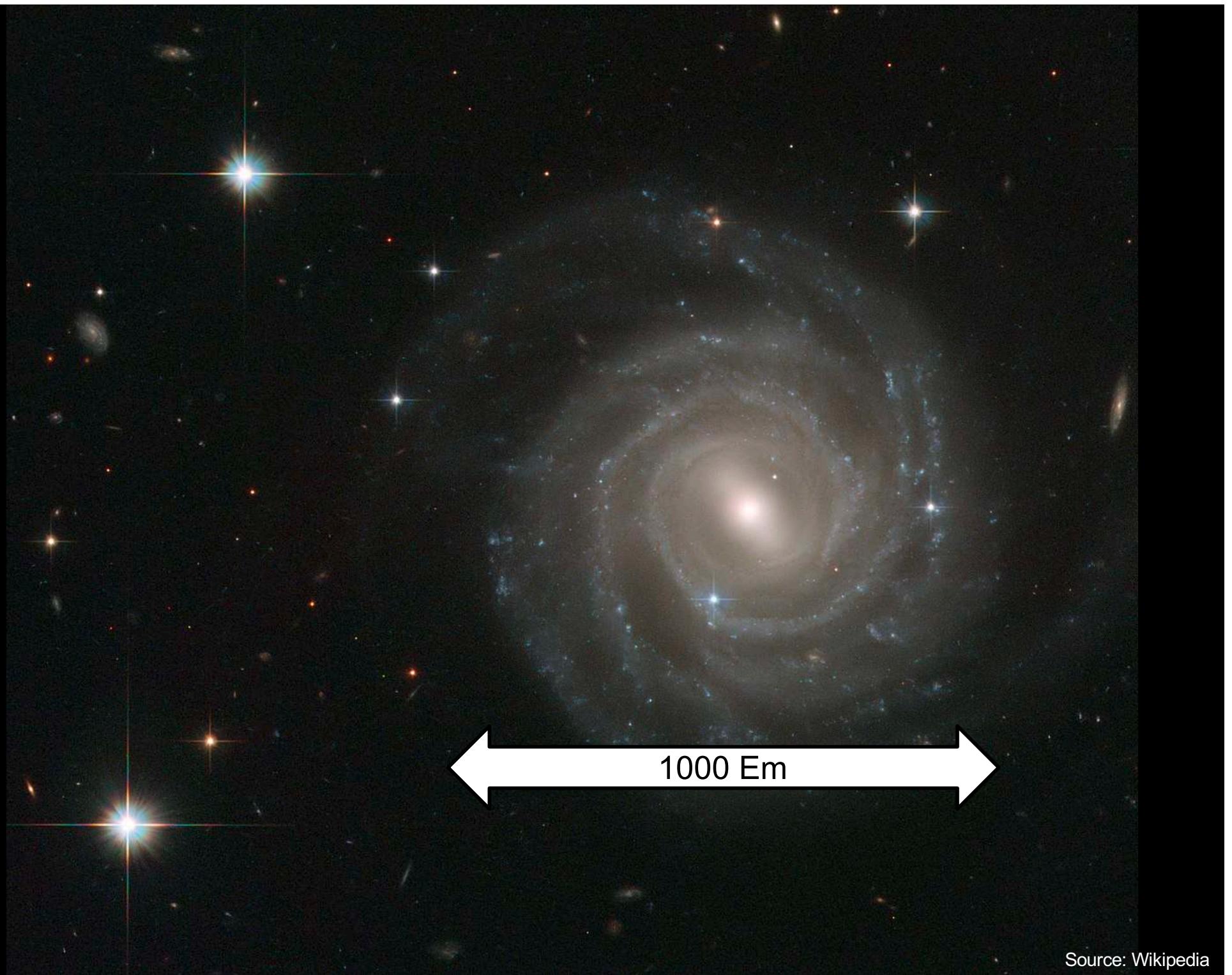
Source: Wikipedia



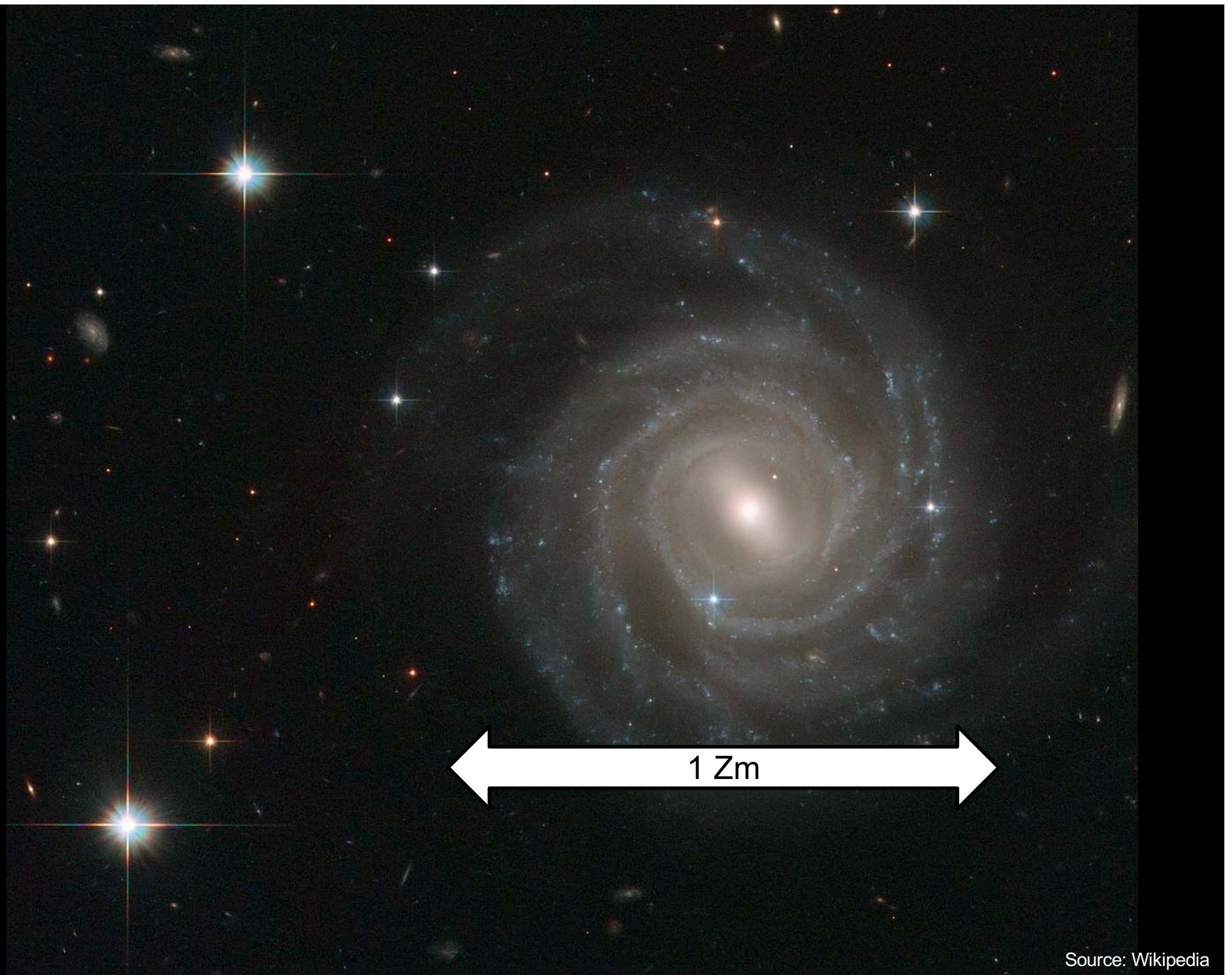
Source: Wikipedia



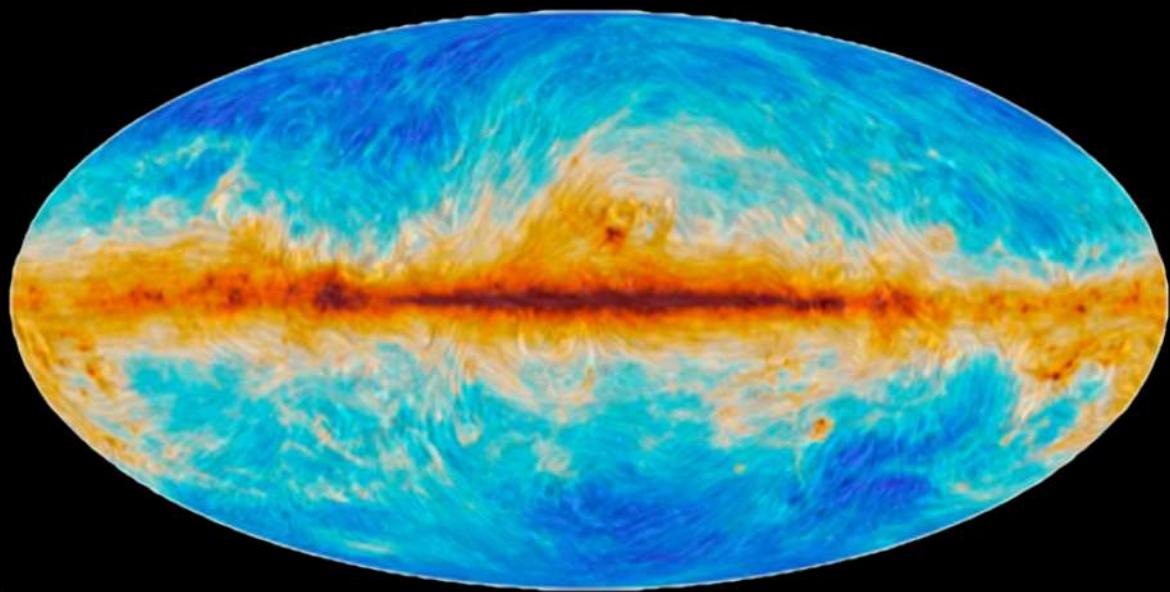
Source: Wikipedia



Source: Wikipedia

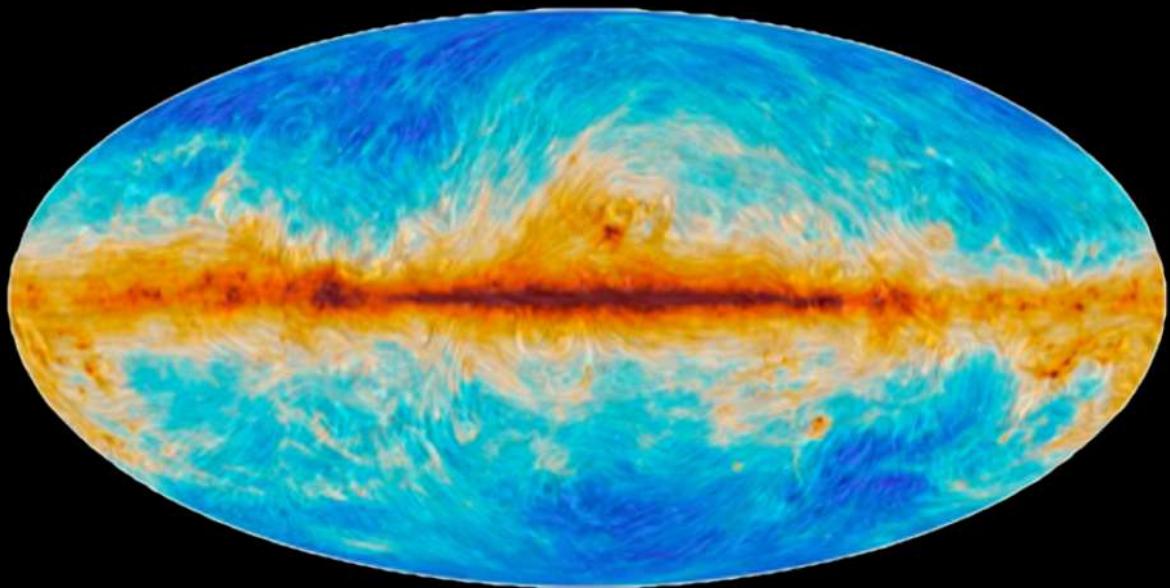


Source: Wikipedia

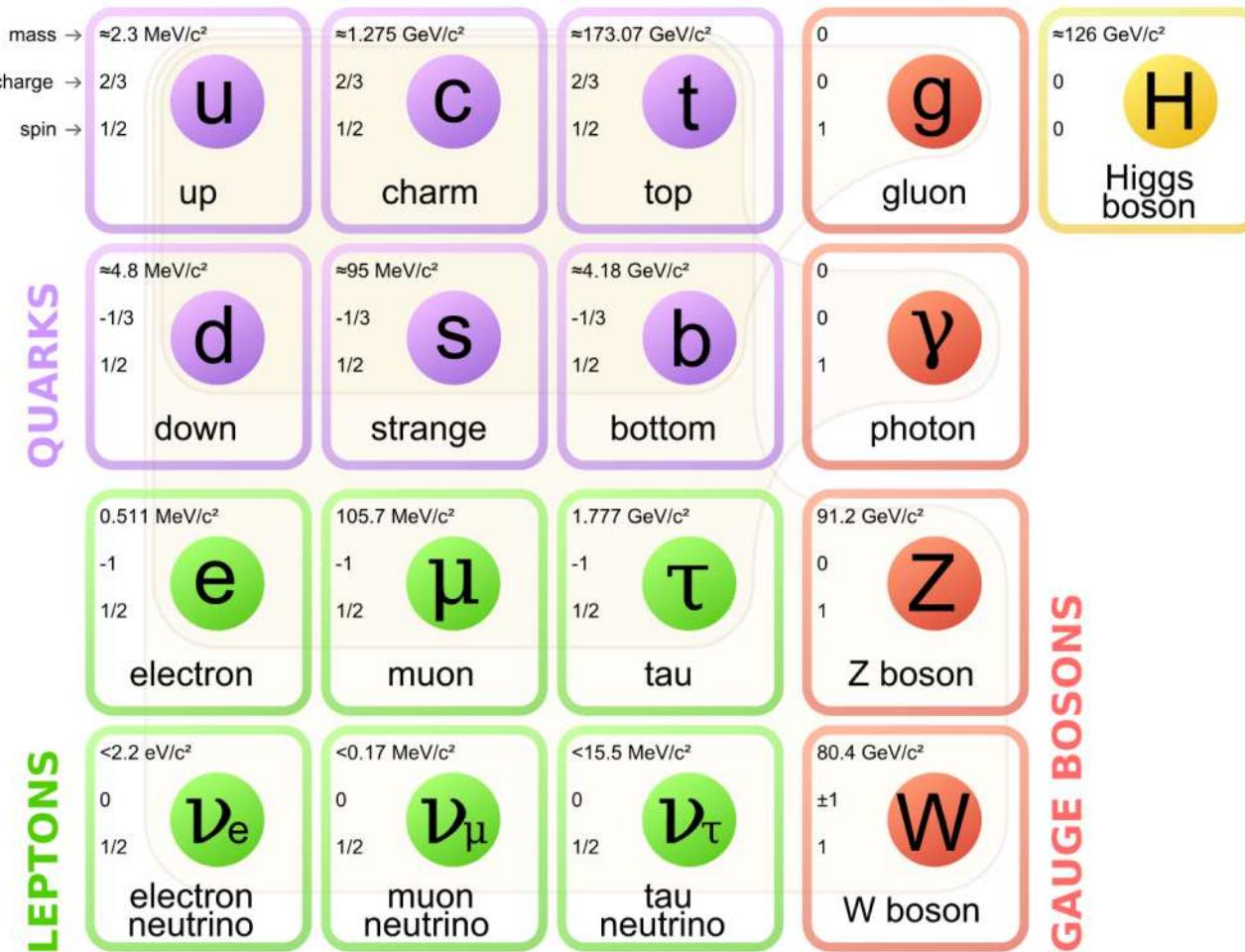


← 138 Ym →

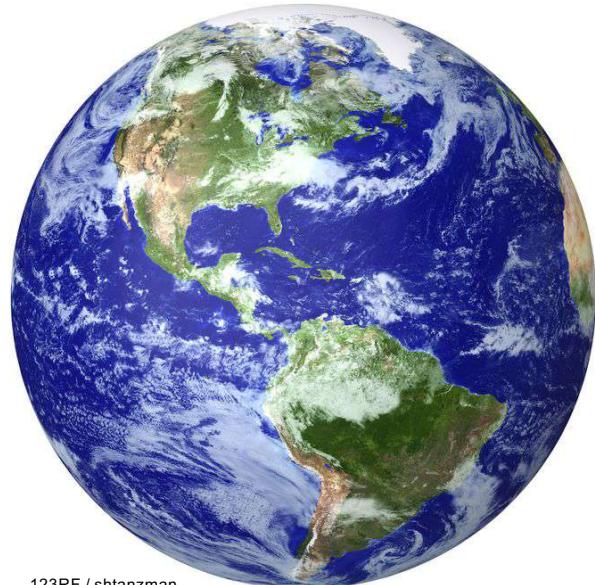
Exploring the infinitely big...



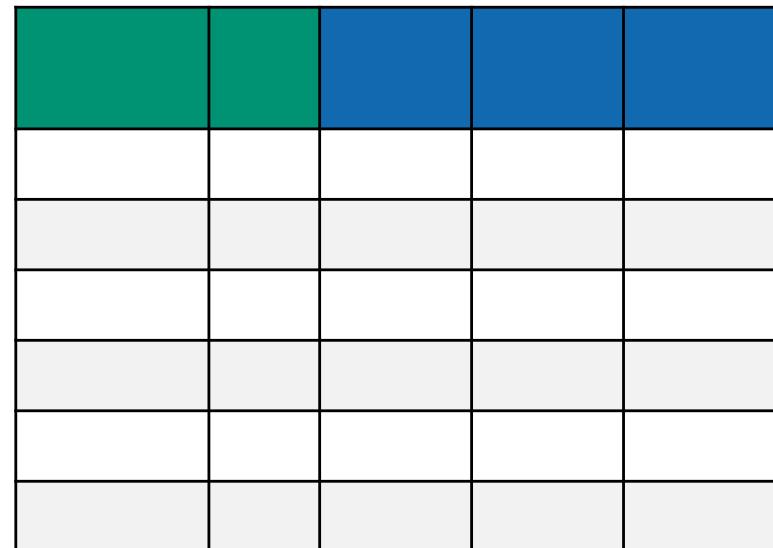
... means exploring the infinitely small



Data is like matter



123RF / shtanzman



Study of the real world

Physics

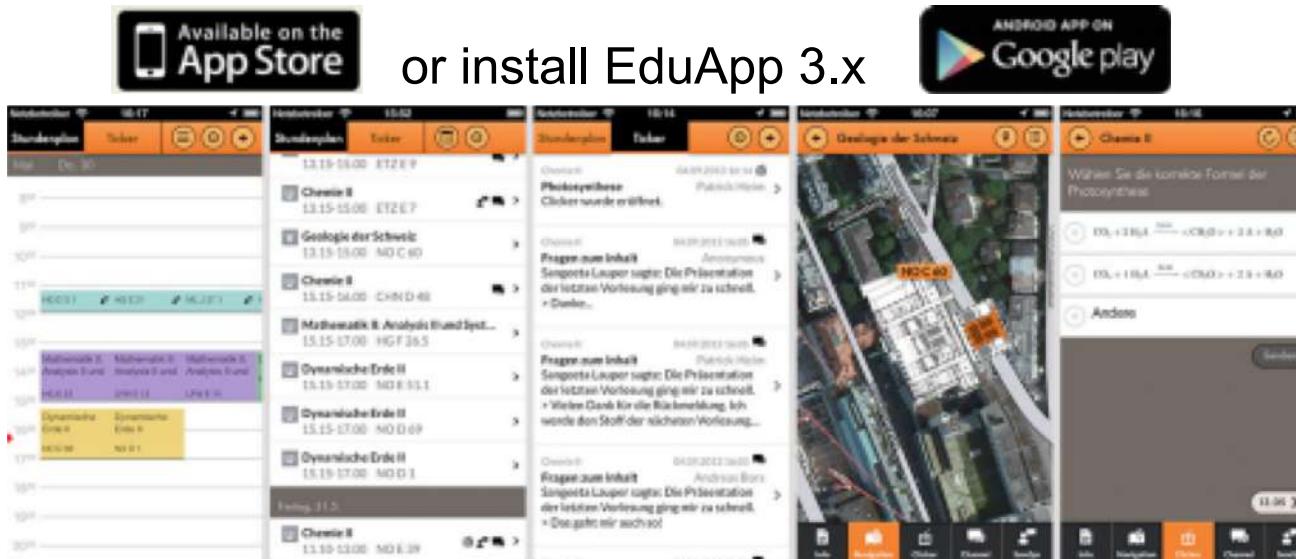
Study of the data world

Data Science

Poll

Go *now* to:

<https://eduapp-app1.ethz.ch/>



My How-We-Do-Science Matrix

Mathematics



Physics



My How-We-Do-Science Matrix

Thinking

With our brain
(natural)



Ontological

The world as it must be
(necessary)

Mathematics

theoretical



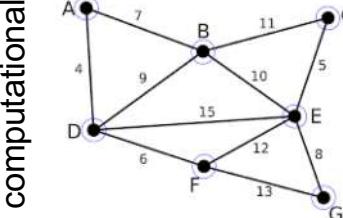
Epistemic

The world as it is
(contingent)

Physics

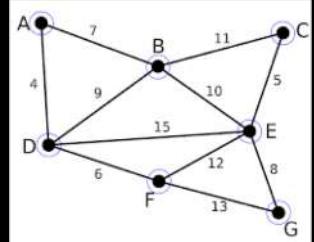


My How-We-Do-Science Matrix

		<i>Ontological</i>	<i>Epistemic</i>
		The world as it must be (necessary)	The world as it is (contingent)
<i>Thinking</i>	With our brain (natural) 	Mathematics 	Physics 
<i>Computing</i>	With a machine (artificial) 	Computer Science 	

My How-We-Do-Science Matrix

The four paradigms

		<i>Ontological</i>	<i>Epistemic</i>
<i>Thinking</i>	With our brain (natural)	The world as it must be (necessary) Mathematics 	The world as it is (contingent) Physics 
<i>Computing</i>	With a machine (artificial)	theoretical Computer Science 	empirical Data Science 

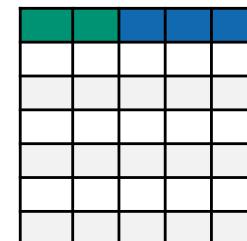
**A good decision is based on
knowledge, not on numbers.**

- Plato

Data Management Lectures at ETH

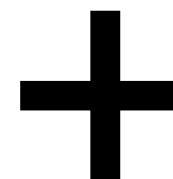
Computer Science
Data Science
CBB MSc with CS background

Other departments



Information Systems for Engineers

Spring 2020



Big Data
(Fall)

This lecture



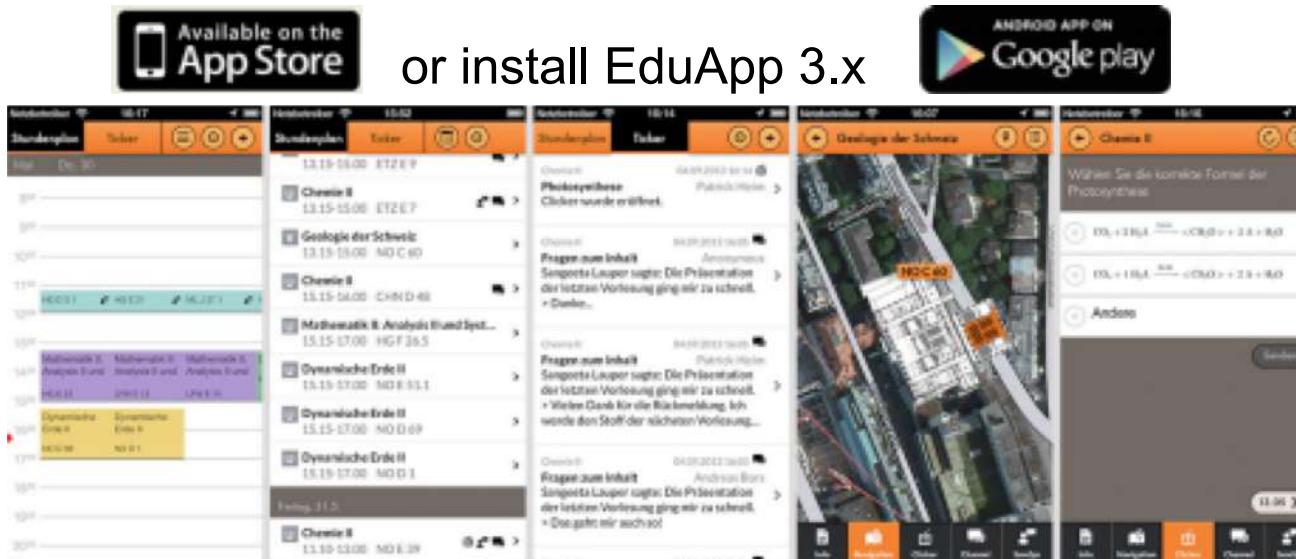
Big Data for Engineers

Spring 2020

Poll

Go *now* to:

<https://eduapp-app1.ethz.ch/>





123RF / Samantha Craddock

A Short History of Databases

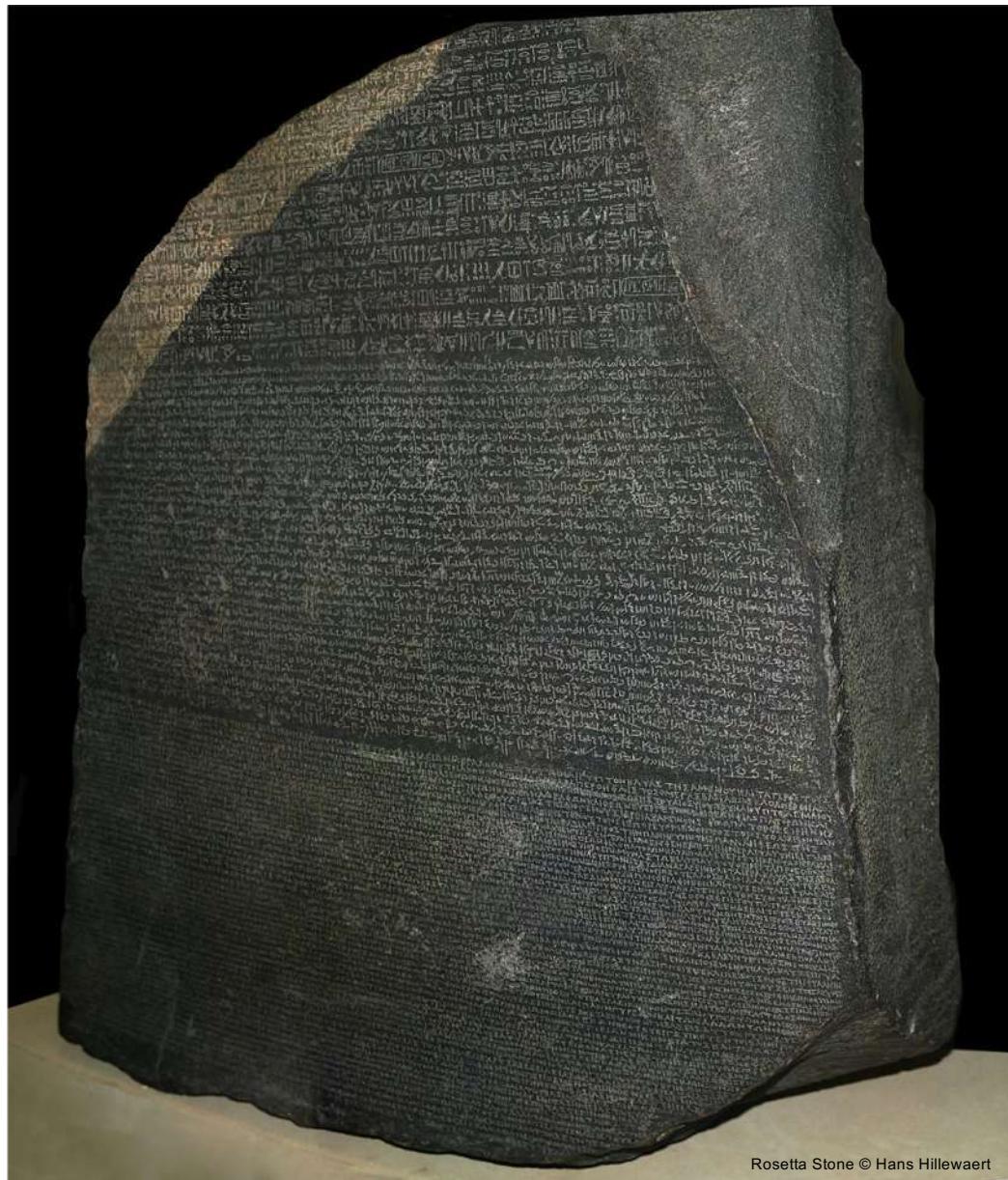


Database Prehistory

123RF / andreykuzmin



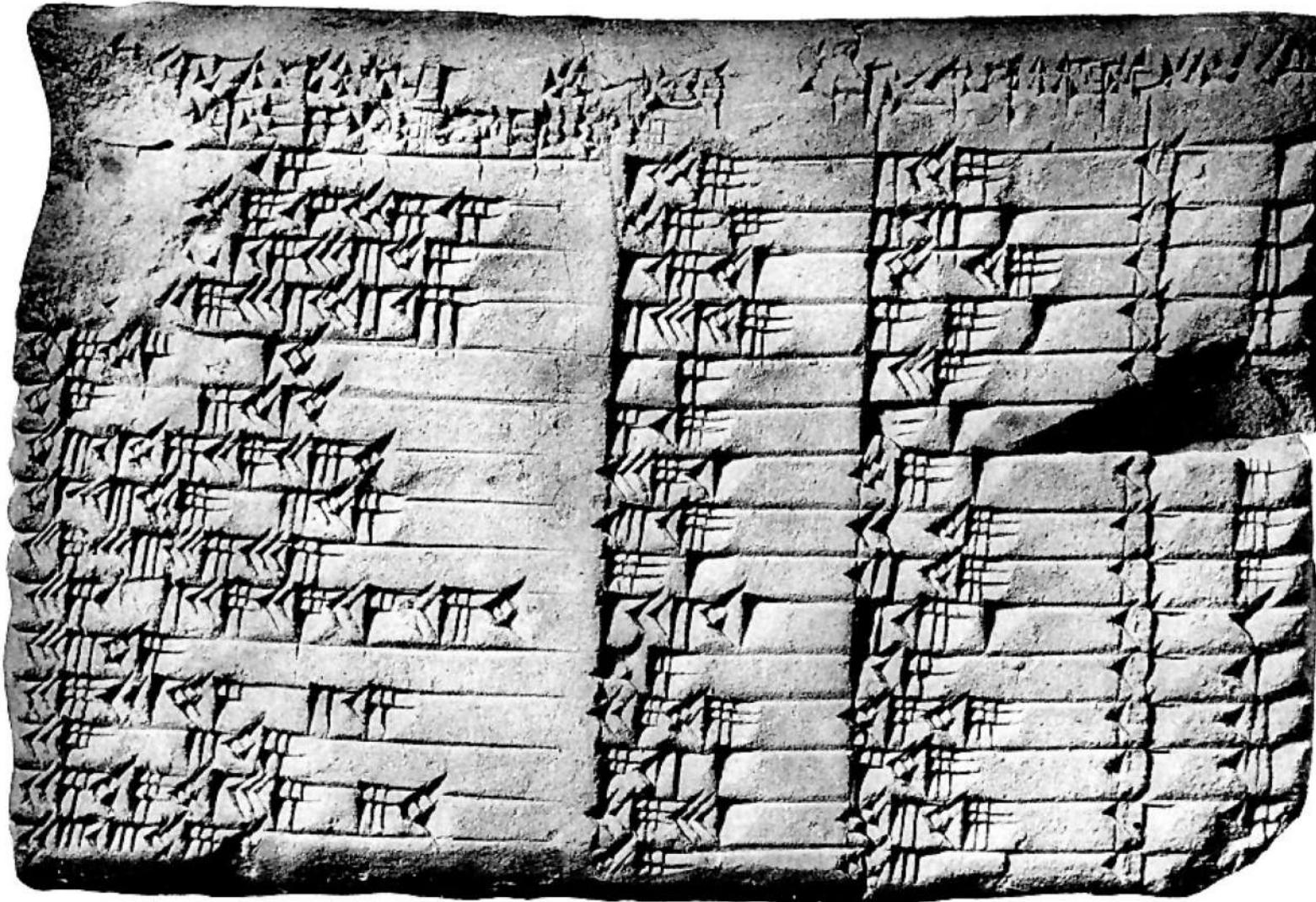
Speaking/Singing



Rosetta Stone © Hans Hillewaert

Writing

Accounting

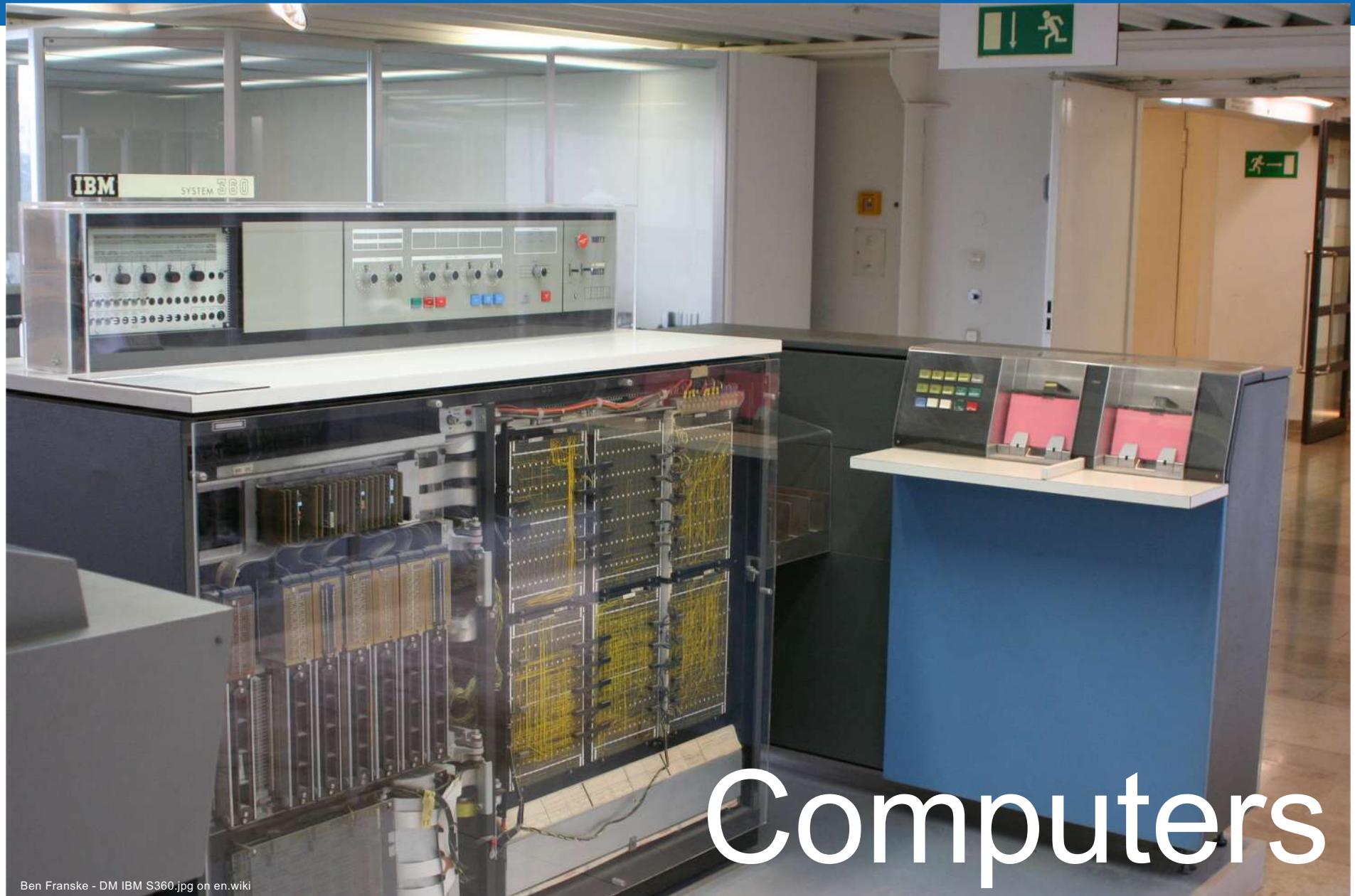


Plimpton 322 (Public Domain)



Willi Heidelbach

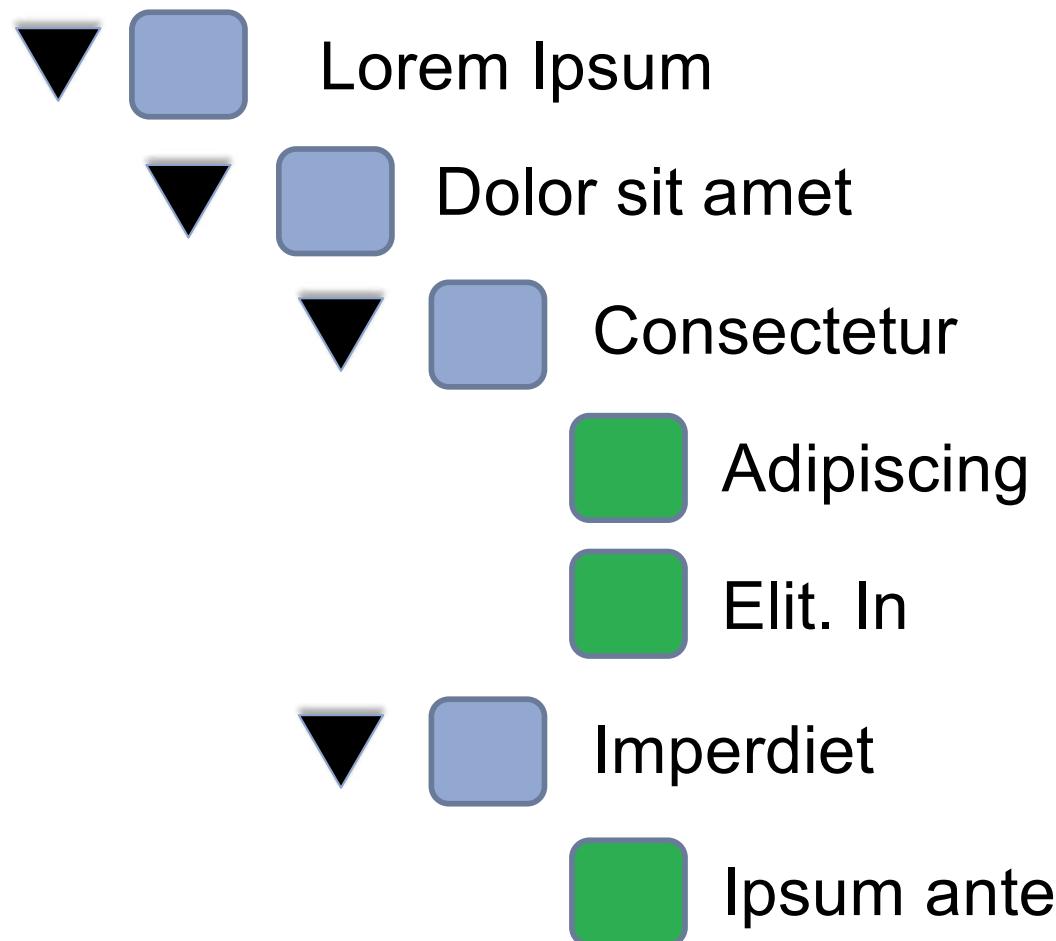
Printing



Ben Franske - DM IBM S360.jpg on en.wiki

Computers

1960s: File Systems





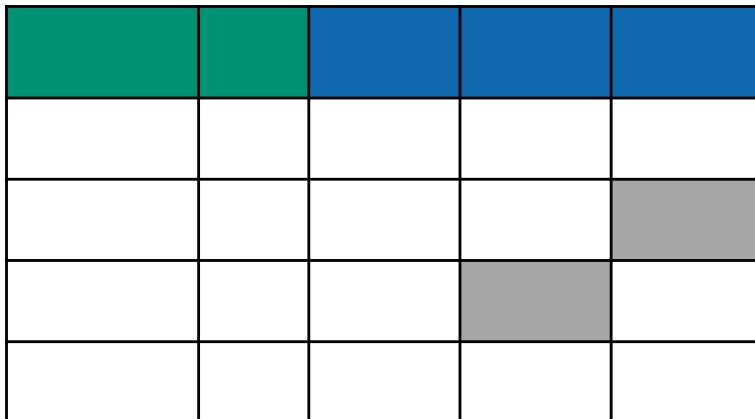
Database History

1970s: The Relational Era

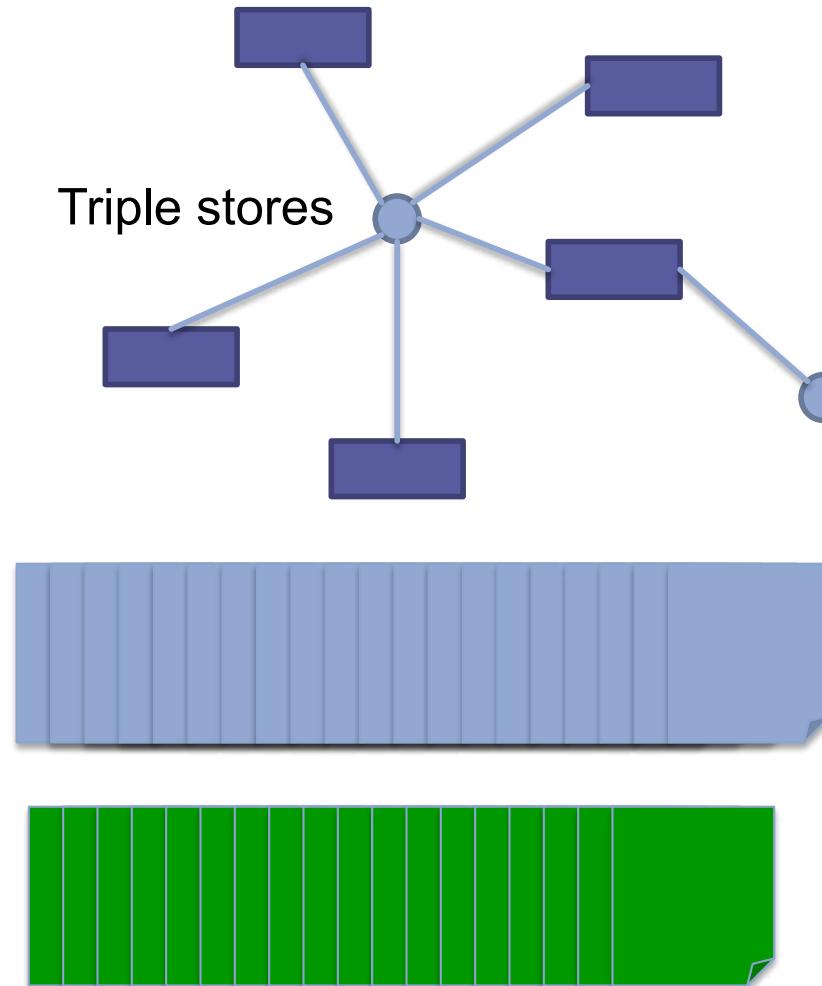
2000s: The NoSQL Era



Key-value stores



Column stores



Document stores

In short?



We threw data
at computers.

1970

In short?



We threw computers
at computers.

1990

In short?



We threw computers
at data.

2000

In short?



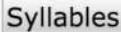
We are throwing
data at data.
now



Big Data

It's a buzzword!

 **buzzword** 

[CITE](#) [buhz-wurd]  

 [Examples](#) [Word Origin](#)
 [See more synonyms on Thesaurus.com](#)



noun

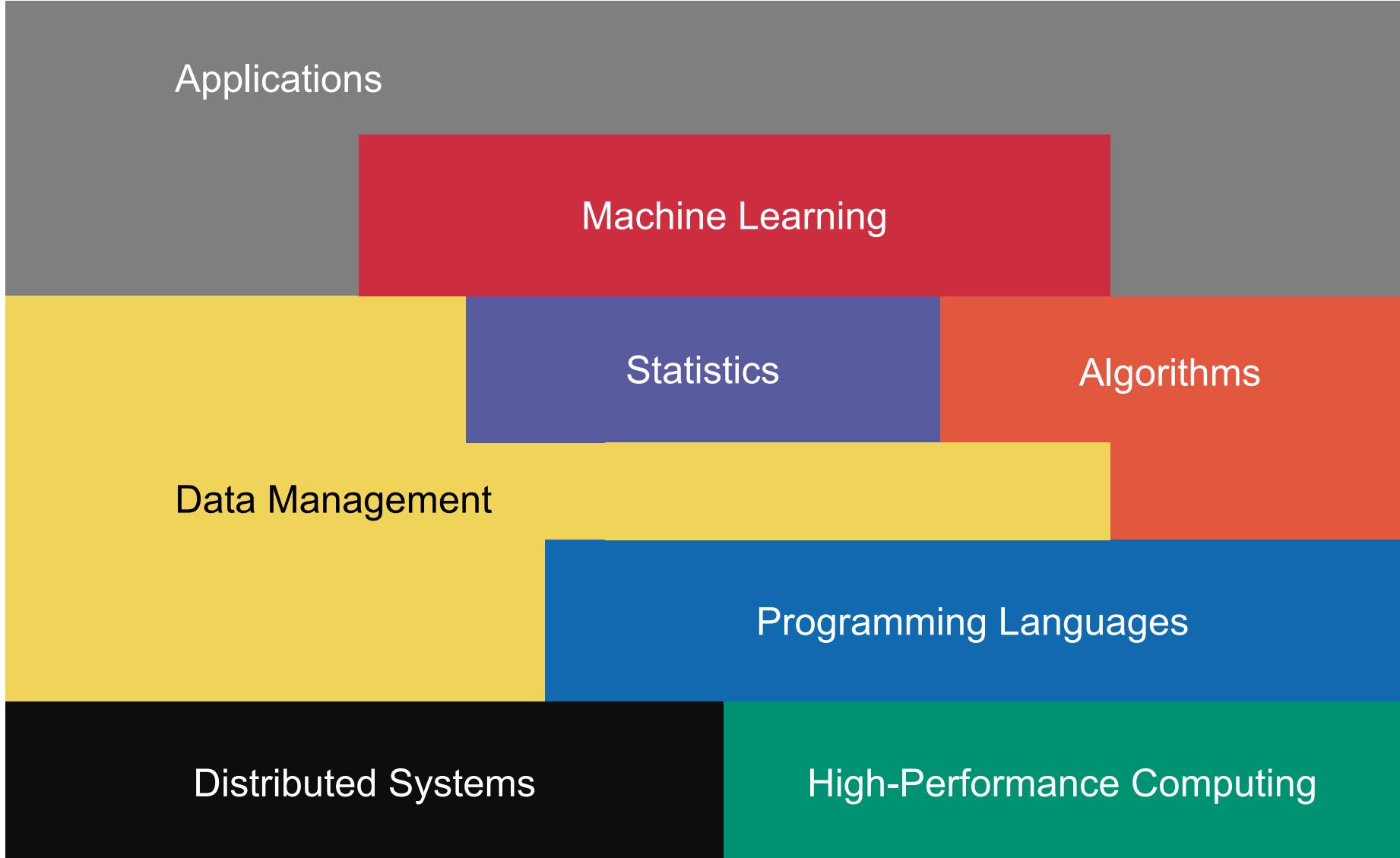
1. a word or phrase, often sounding authoritative or technical, that is a vogue term in a particular profession, field of study, popular culture, etc.

Origin of buzzword 

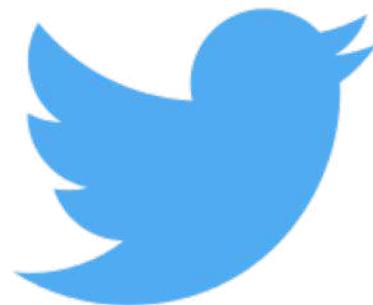
1965-1970 

1965-70; [buzz¹](#)+ [word](#)

Big Data goes across disciplines



Big Data involves a lot of proprietary technology

The Facebook logo is the word "facebook" in white lowercase letters on a solid blue rectangular background.The Yahoo! logo is the word "YAHOO!" in a large, bold, purple, three-dimensional font.

Google Dataset Search

Beta



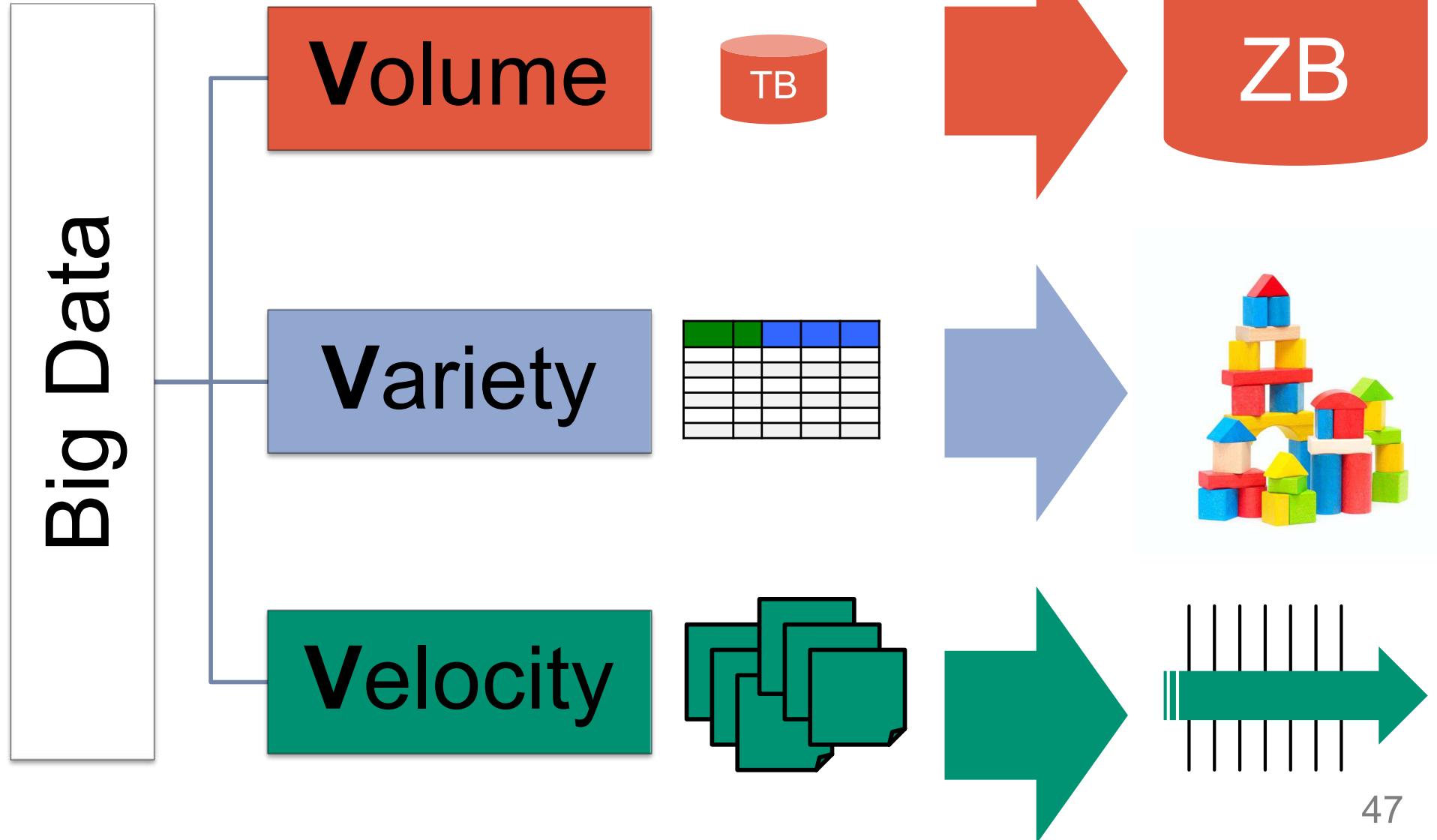
Try [boston education data](#) or [weather site:noaa.gov](#)



123RF / Patricia Hofmeester

The Big in Big Data

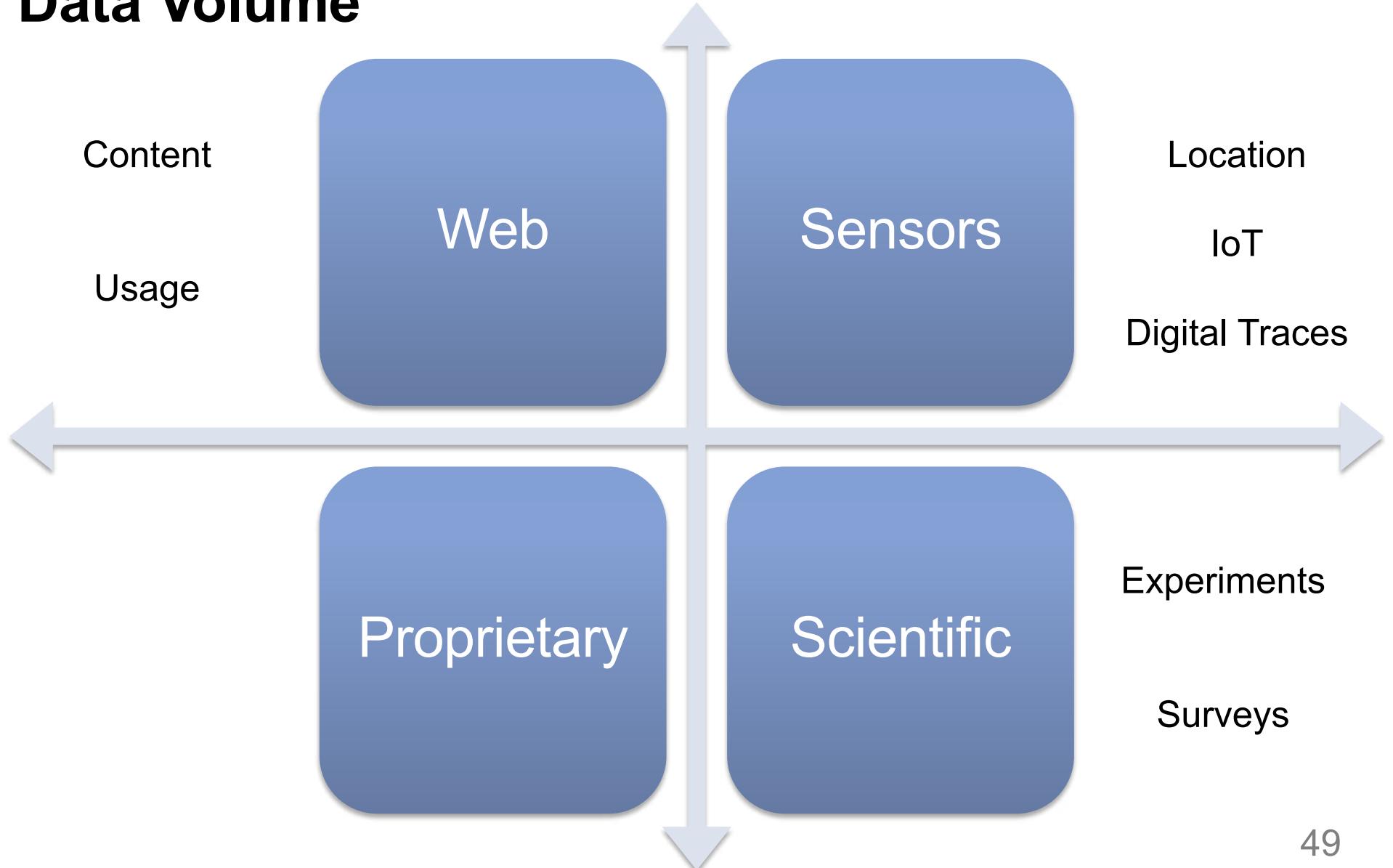
The Three Vs



MORE
MORE
MORE

Data Volume

Data Volume



Data Volume

... because
we
can!

Technology

Software

Hardware

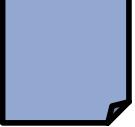
Infrastructure

Data Volume

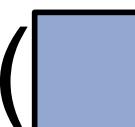


... because
data
carries
value

Data is worth more than the sum of its parts

Utility( **+**  **)**

>

Utility( **) + Utility(**  **)**

Data totality: one must have complete data

**All flights
All hotels
All shops**

...



Prefixes (International System of Units)

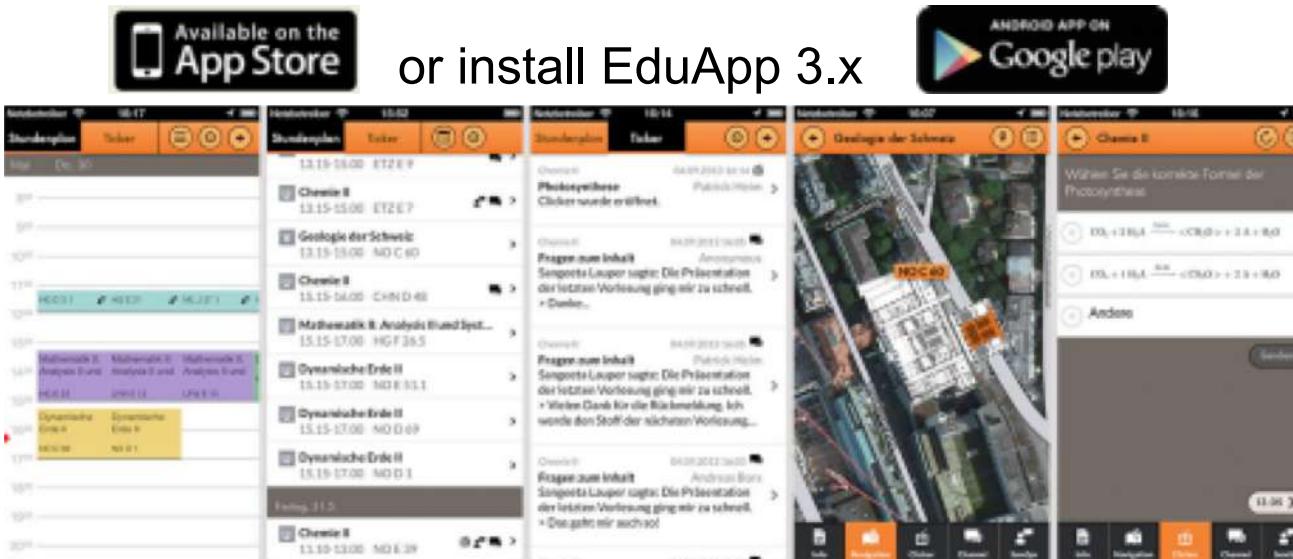
kilo (k)	1,000 (3 zeros)
Mega (M)	1,000,000 (6 zeros)
Giga (G)	1,000,000,000 (9 zeros)
Tera (T)	1,000,000,000,000 (12 zeros)
Peta (P)	1,000,000,000,000,000 (15 zeros)
Exa (E)	1,000,000,000,000,000,000 (18 zeros)
Zetta (Z)	1,000,000,000,000,000,000,000 (21 zeros)
Yotta (Y)	1,000,000,000,000,000,000,000,000 (24 zeros)

You must know this by ❤️!

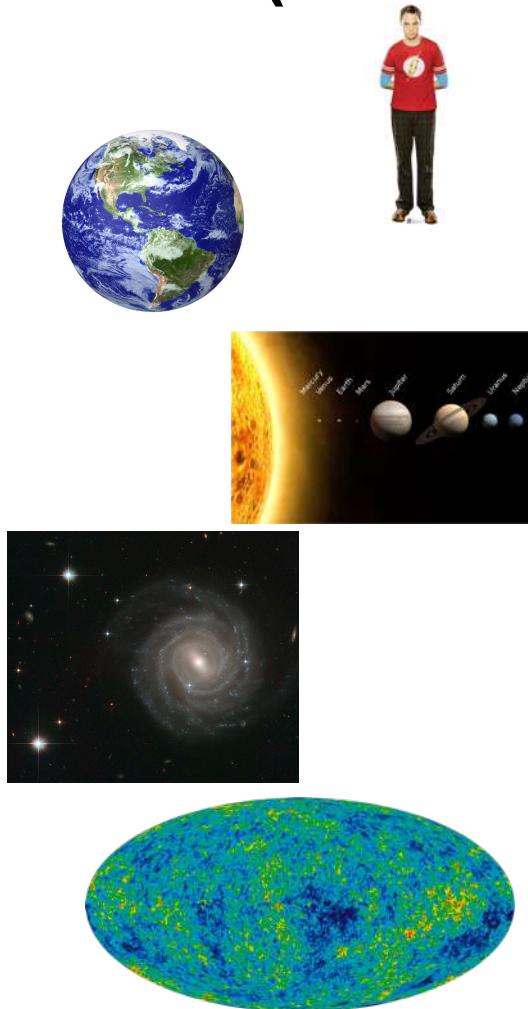
Clicker question

Go *now* to:

<https://eduapp-app1.ethz.ch/>



Prefixes (International System of Units)



kilo (k)
Mega (M)
Giga (G)
Tera (T)
Peta (P)
Exa (E)
Zetta (Z)
Yotta (Y)



Prefixes (International System of Units)

kibi (ki)	1,024 (2^{10})
Mebi (Mi)	1,048,576 (2^{20})
Gibi (Gi)	1,073,741,824 (2^{30})
Tebi (Ti)	1,099,511,627,776 (2^{40})
Pebi (Pi)	1,125,899,906,842,624 (2^{50})
Exbi (Ei)	1,152,921,504,606,846,976 (2^{60})
Zebi (Zi)	1,180,591,620,717,411,303,424 (2^{70})
Yobi (Yi)	1,208,925,819,614,629,174,706,176 (2^{80})

You must NOT know this by ❤ !



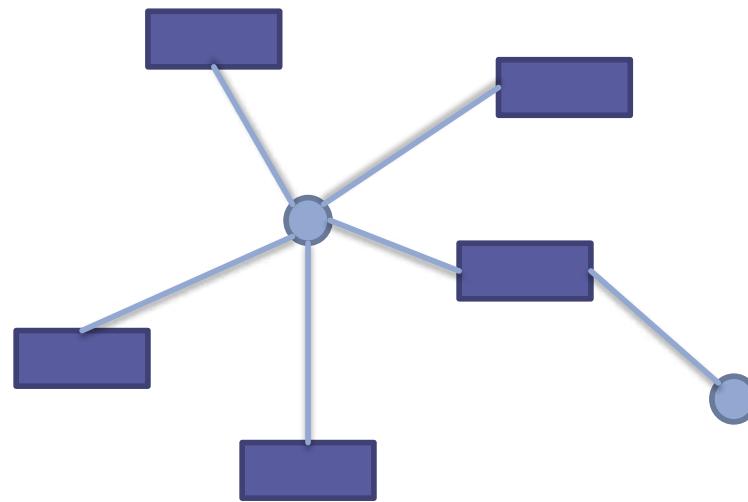
Data Variety

Data Shapes: Tables

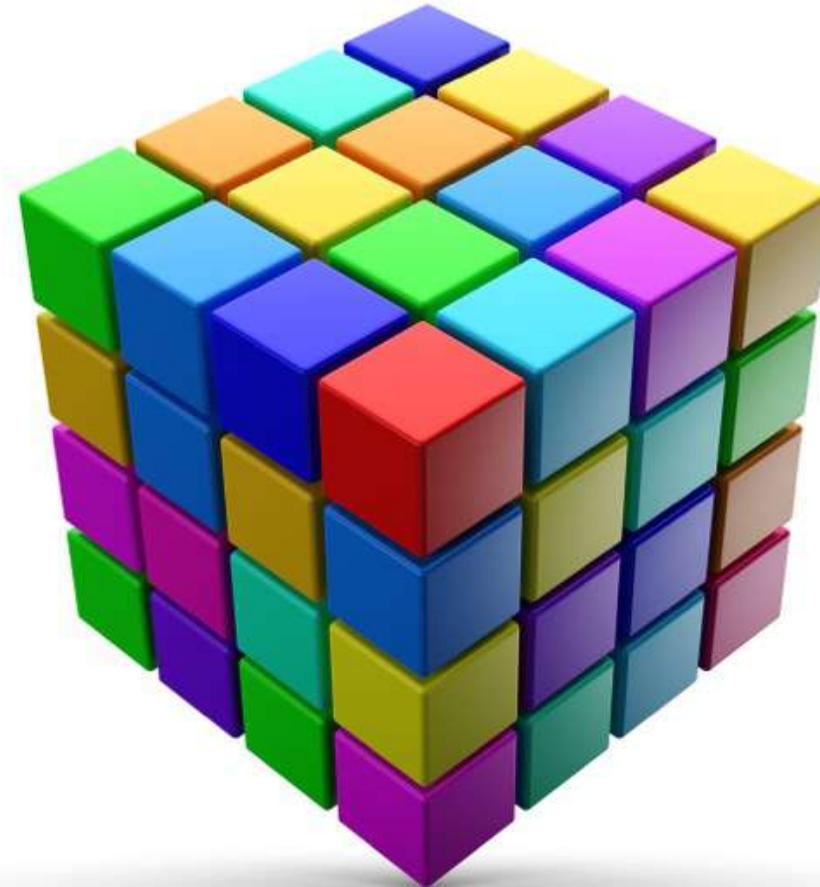
Data Shapes: Trees



Data Shapes: Graphs



Data Shapes: Cubes



Data Shapes: Text

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam vel erat nec dui aliquet vulputate sed quis nulla. Donec eget ultricies magna, eu dignissim elit. Nullam sed urna nec nisl rhoncus ullamcorper placerat et enim. Integer varius ornare libero quis consequat. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aenean eu efficitur orci. Aenean ac posuere tellus. Ut id commodo turpis.

Praesent nec libero metus. Praesent at turpis placerat, congue ipsum eget, scelerisque justo. Ut volutpat, massa ac lacinia cursus, nisl dui volutpat arcu, quis interdum sapien turpis in tellus. Suspendisse potenti. Vestibulum pharetra justo massa, ac venenatis mi condimentum nec. Proin viverra tortor non orci suscipit rutrum. Phasellus sit amet euismod diam. Nullam convallis nunc sit amet diam suscipit dapibus. Integer porta hendrerit nunc. Quisque pharetra congue porta. Suspendisse vestibulum sed mi in euismod. Etiam a purus suscipit, accumsan nibh vel, posuere ipsum. Nulla nec tempor nibh, id venenatis lectus. Duis lobortis id urna eget tincidunt.



Data Velocity

Data is generated automatically



Data is a realtime byproduct of human activity



Three paramount factors

Capacity

Throughput

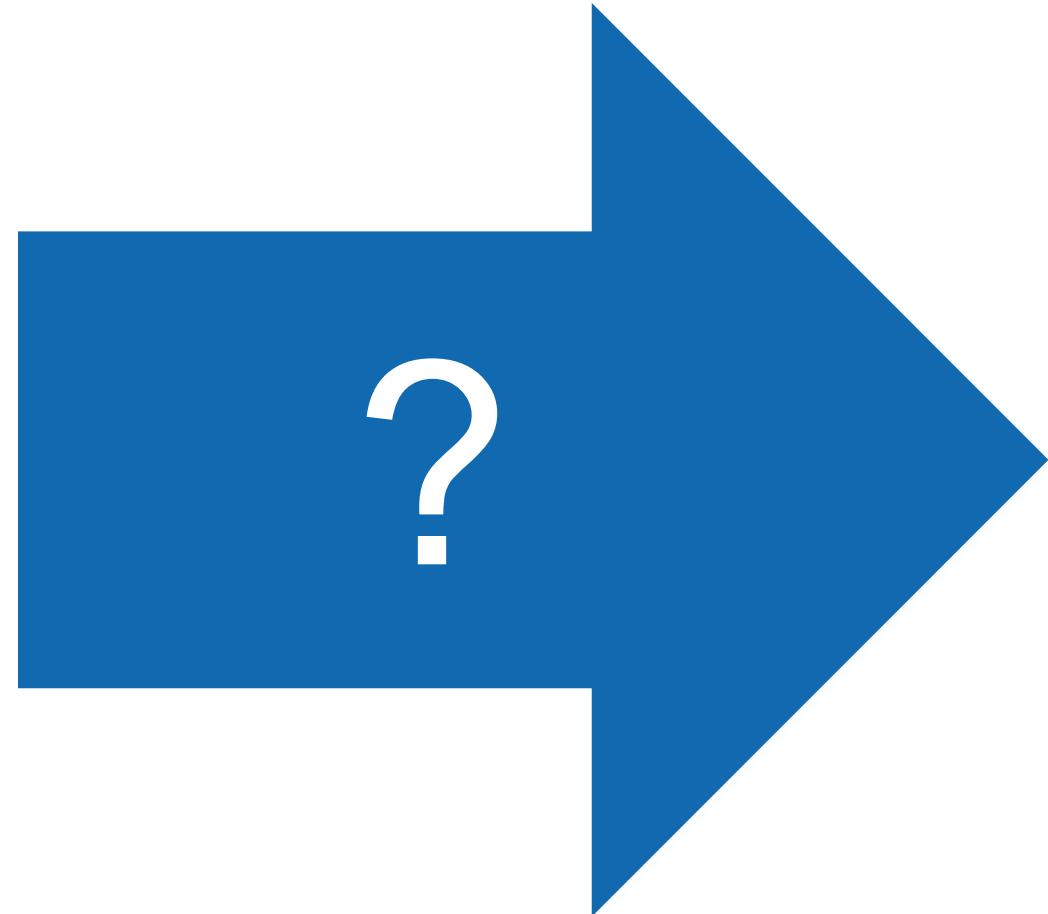
Latency

Capacity



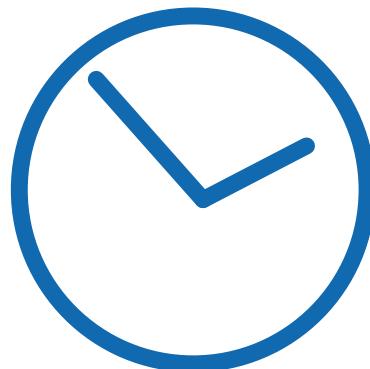
"How much data can we store?"

Throughput



"How fast can we transmit data?"

Latency



"When do I start receiving data?"

1956: IBM RAMAC 350



1956: IBM RAMAC 350



2020: Western Digital Ultrastar DC HC650



2020: Western Digital Ultrastar DC HC650



2020: Western Digital Ultrastar DC HC650



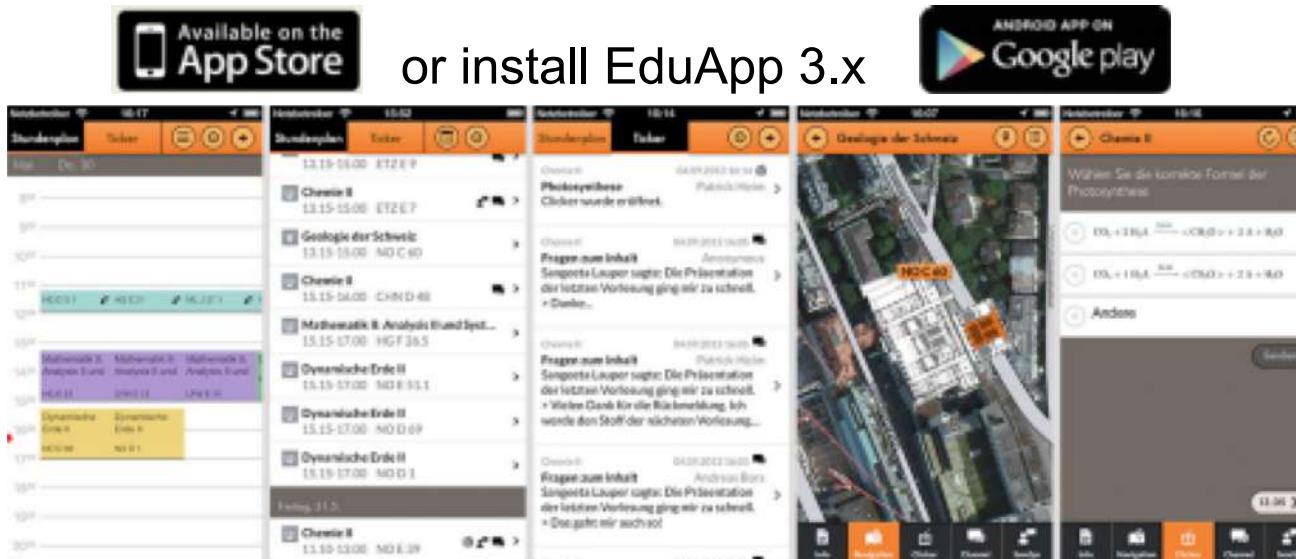
2020: Western Digital Ultrastar DC HC650



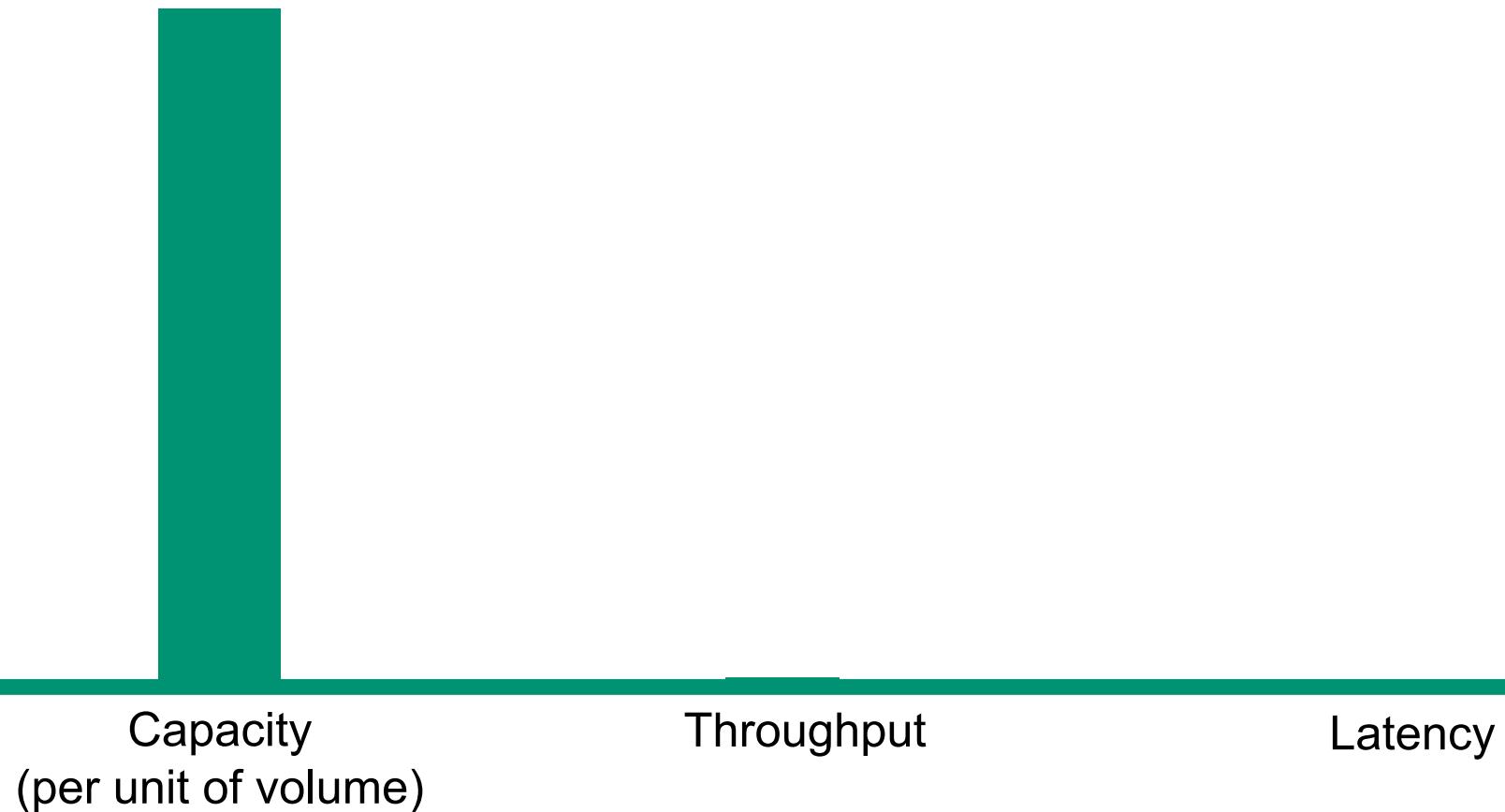
Clicker question

Go *now* to:

<https://eduapp-app1.ethz.ch/>



The progress made (1956-2020)



The progress made (1956-2020)

200,000,000,000x



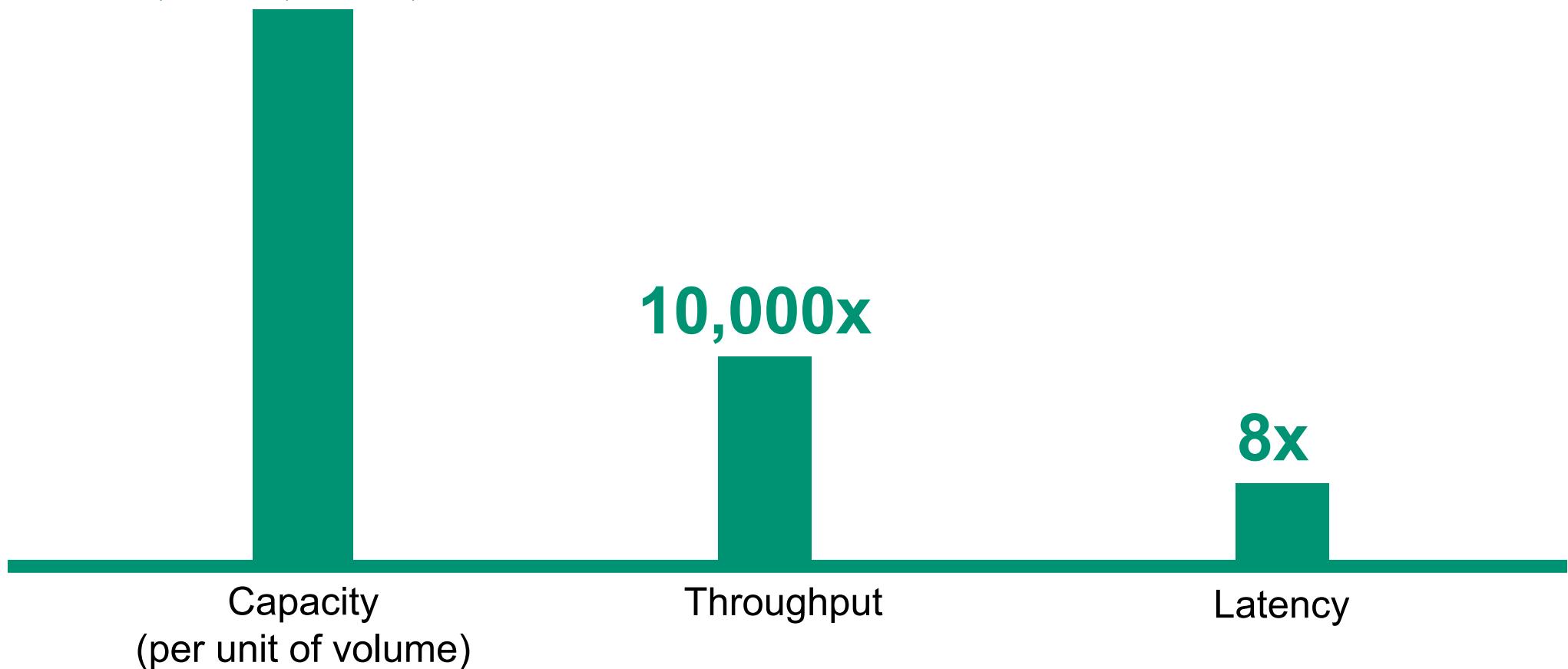
Capacity
(per unit of volume)

10,000x

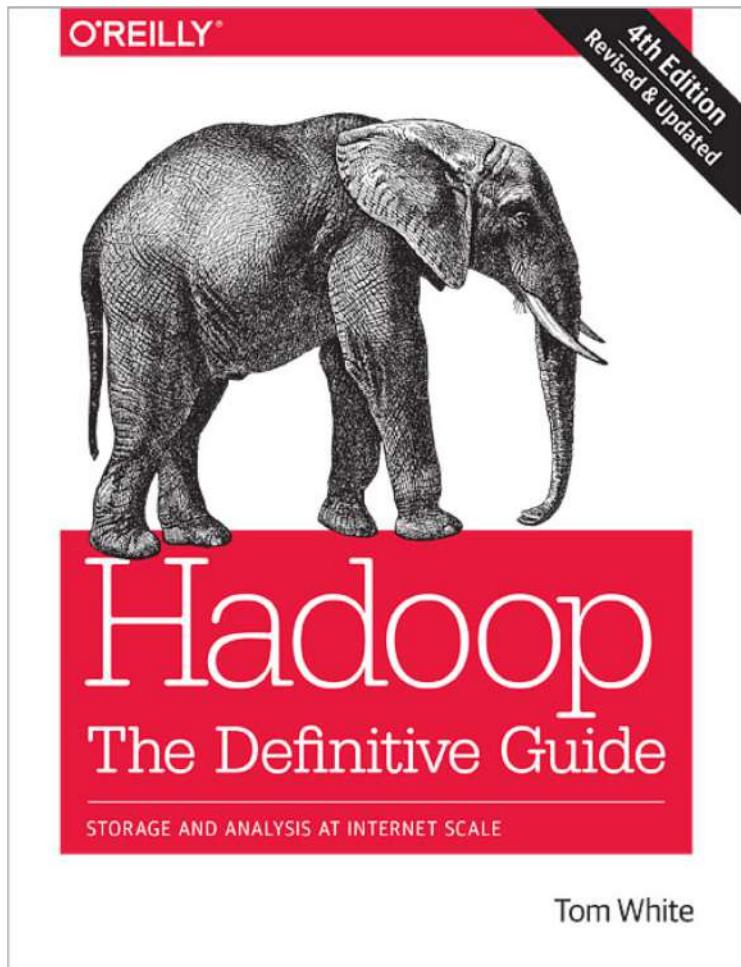
8x

Latency

The progress made (1956-2020): Logarithmic **200,000,000,000x**

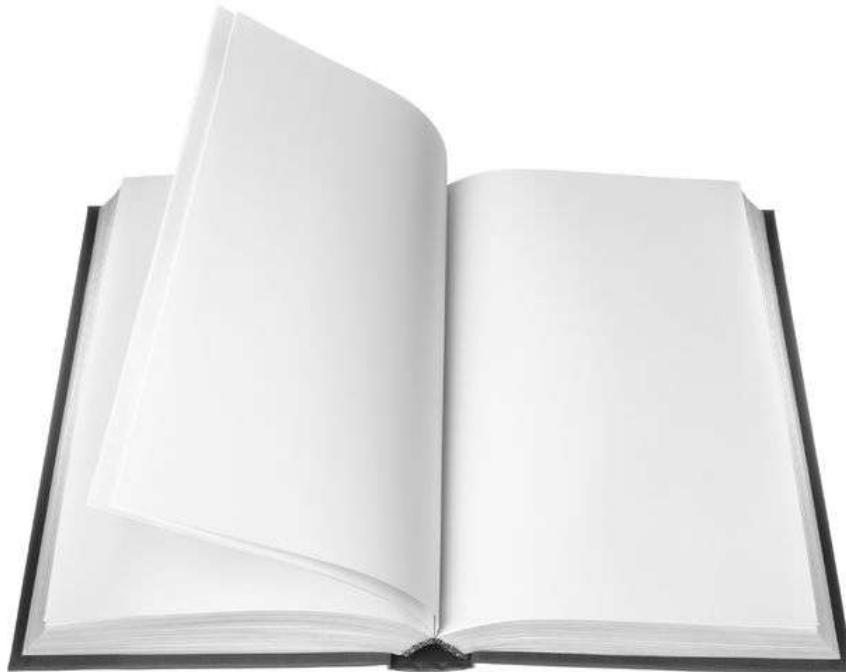


Capacity: Example



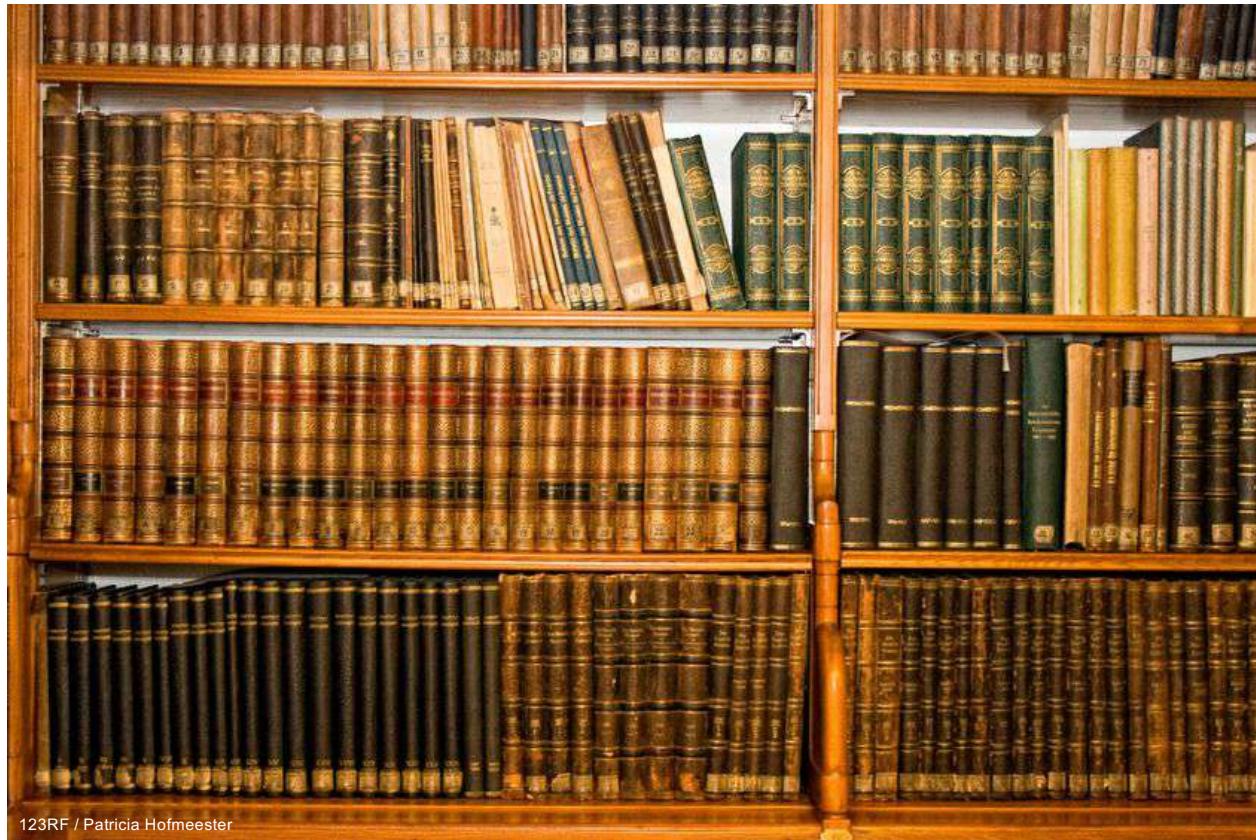
~ 600,000 words.

Throughput: Example



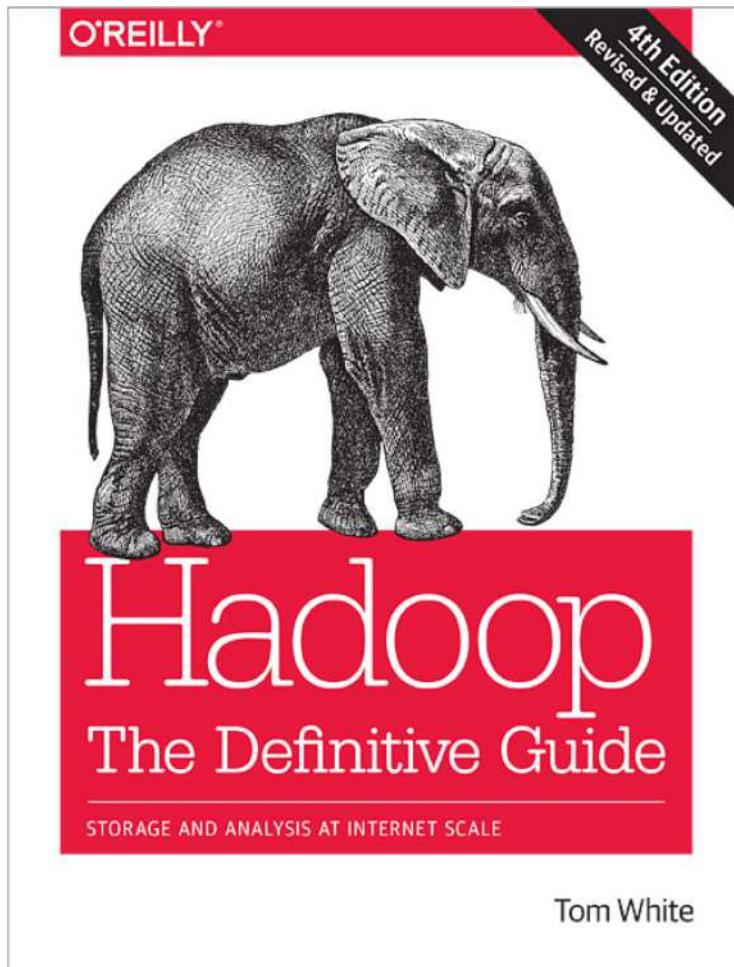
~1,000 word
per minute

Latency: Example



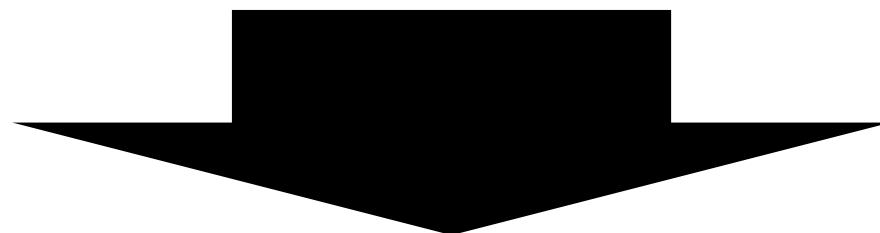
~ 1 minute to
stand up,
go to the shelf,
pick the book,
find the page.

2020 – Analogy with a book



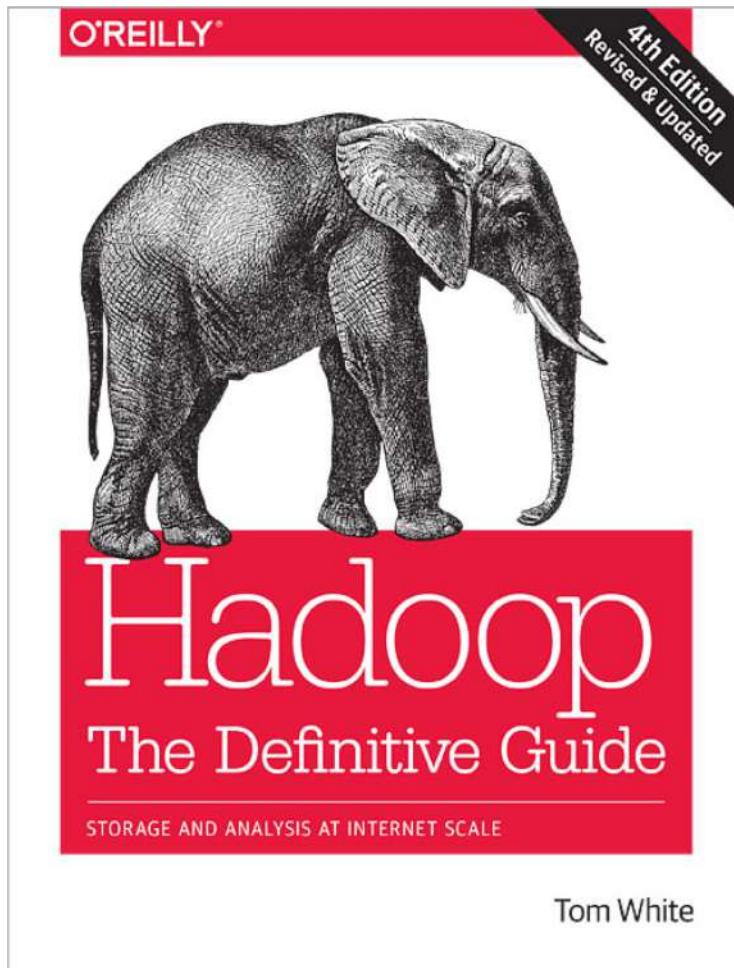
600,000 words

1,000 words per minute



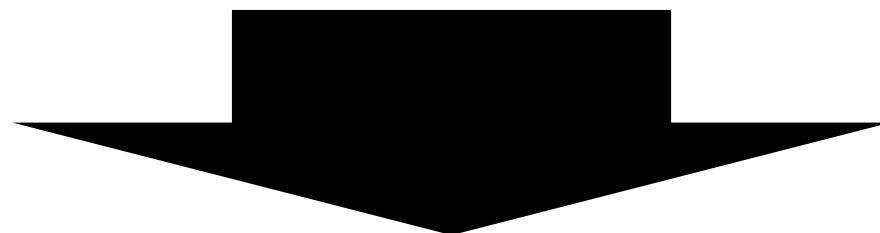
10 hours

2220 – Analogy with a book



120,000,000,000,000 words

10,000,000 words per minute



22,800 years

The progress made (1956-2020): Logarithmic

200,000,000,000x



10,000x



8x



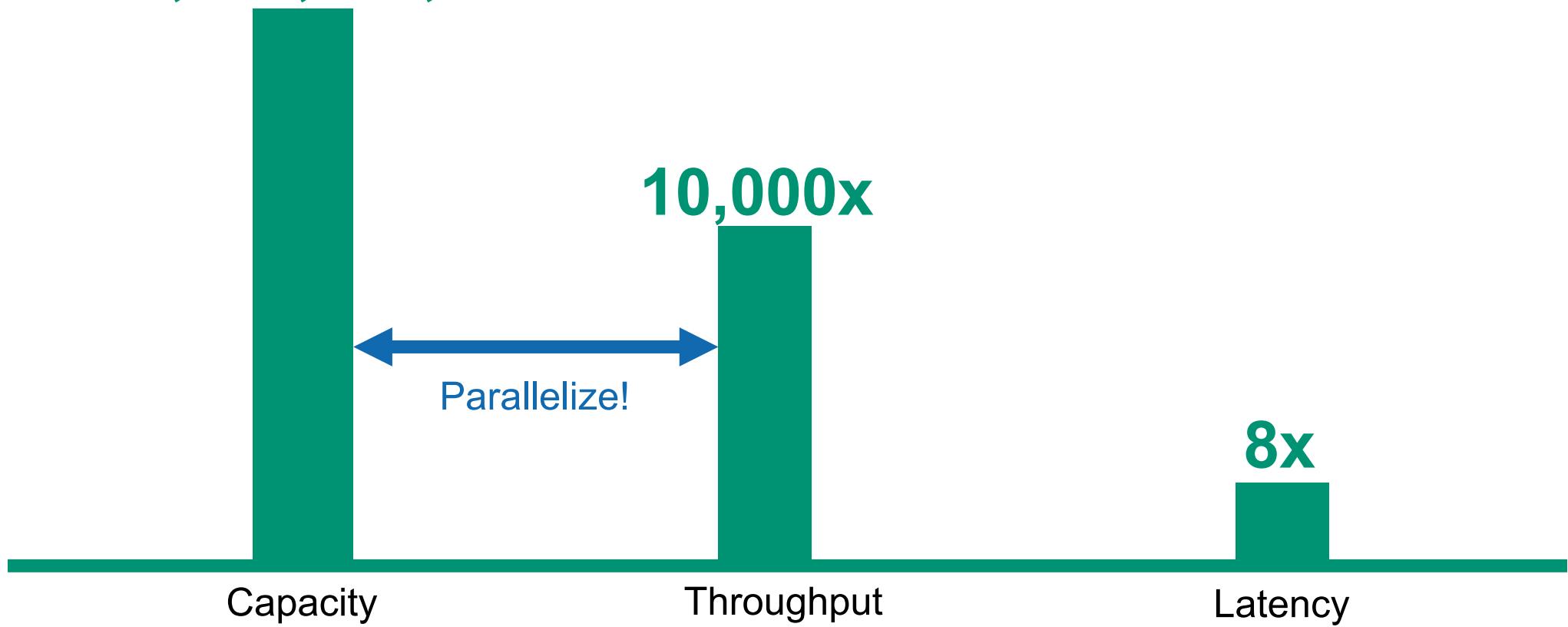
Capacity

Throughput

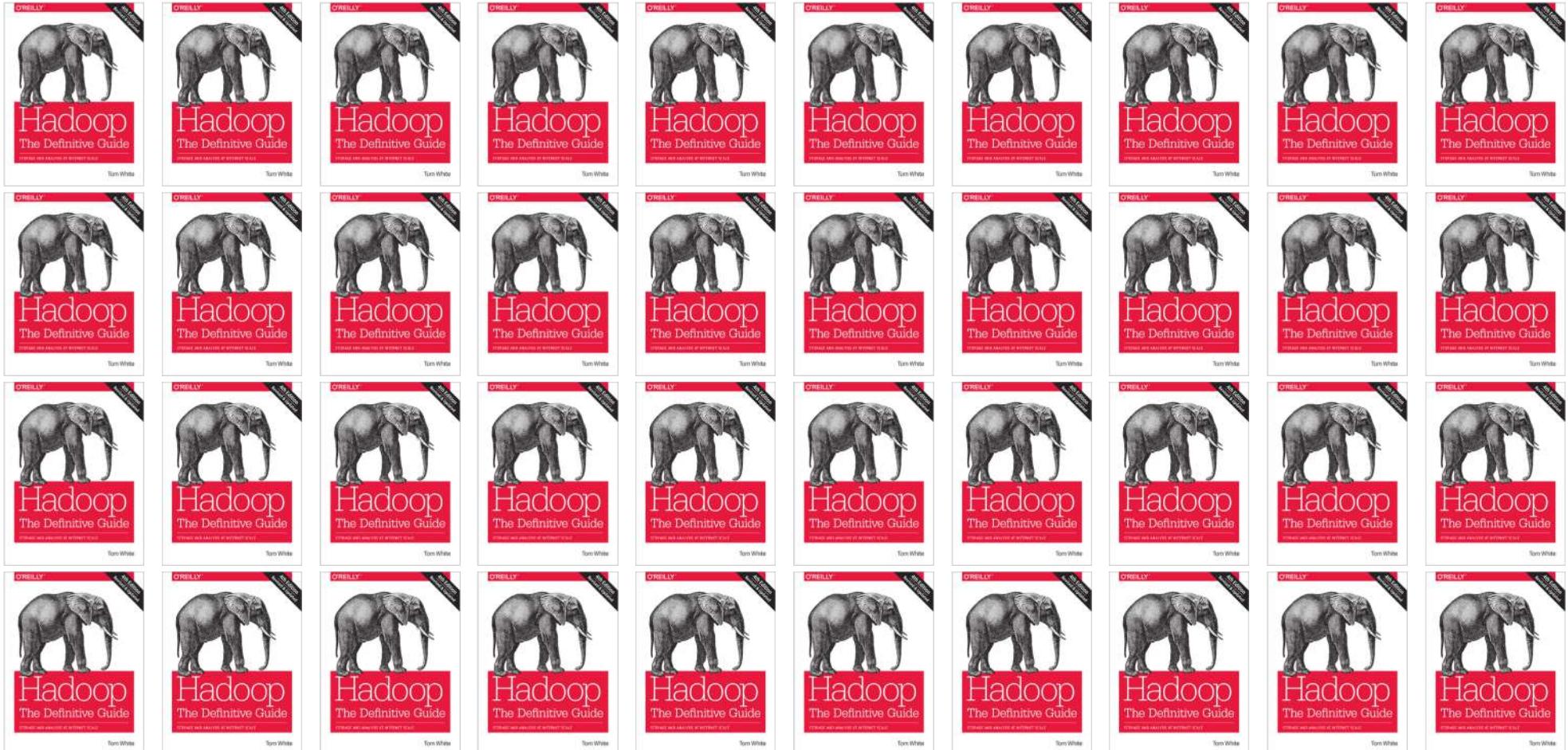
Latency

The progress made (1956-2020): Logarithmic

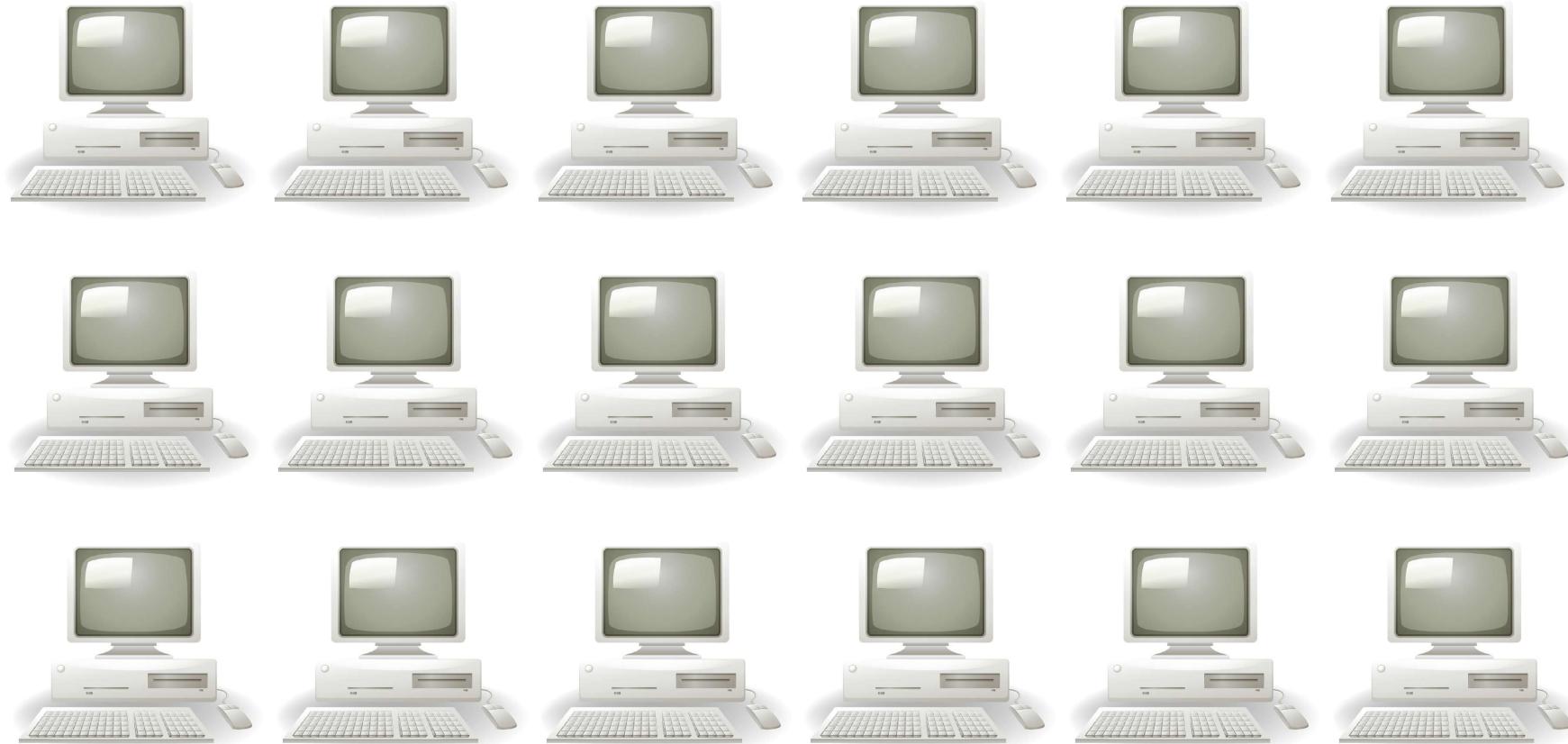
200,000,000,000x



2220: 200,000,000 persons
could read it all in 10 hours.

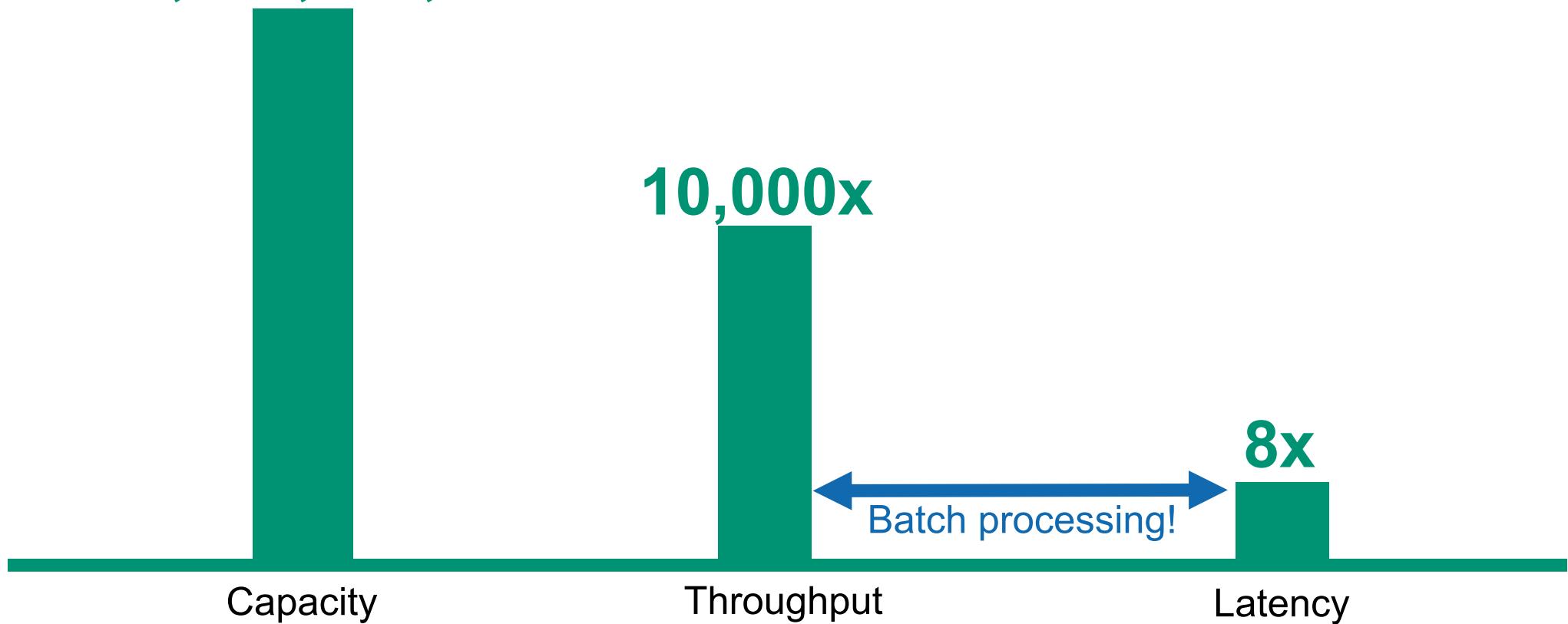


Data centers: clusters of machines (10,000s)



The progress made (1956-2020): Logarithmic

200,000,000,000x



What is Big Data (my definition)?

Big Data is a portfolio of technologies
that were designed to

store, manage and analyze data that is
too large to fit on a single machine

while accommodating for the issue of

growing discrepancy between
capacity, throughput and latency.

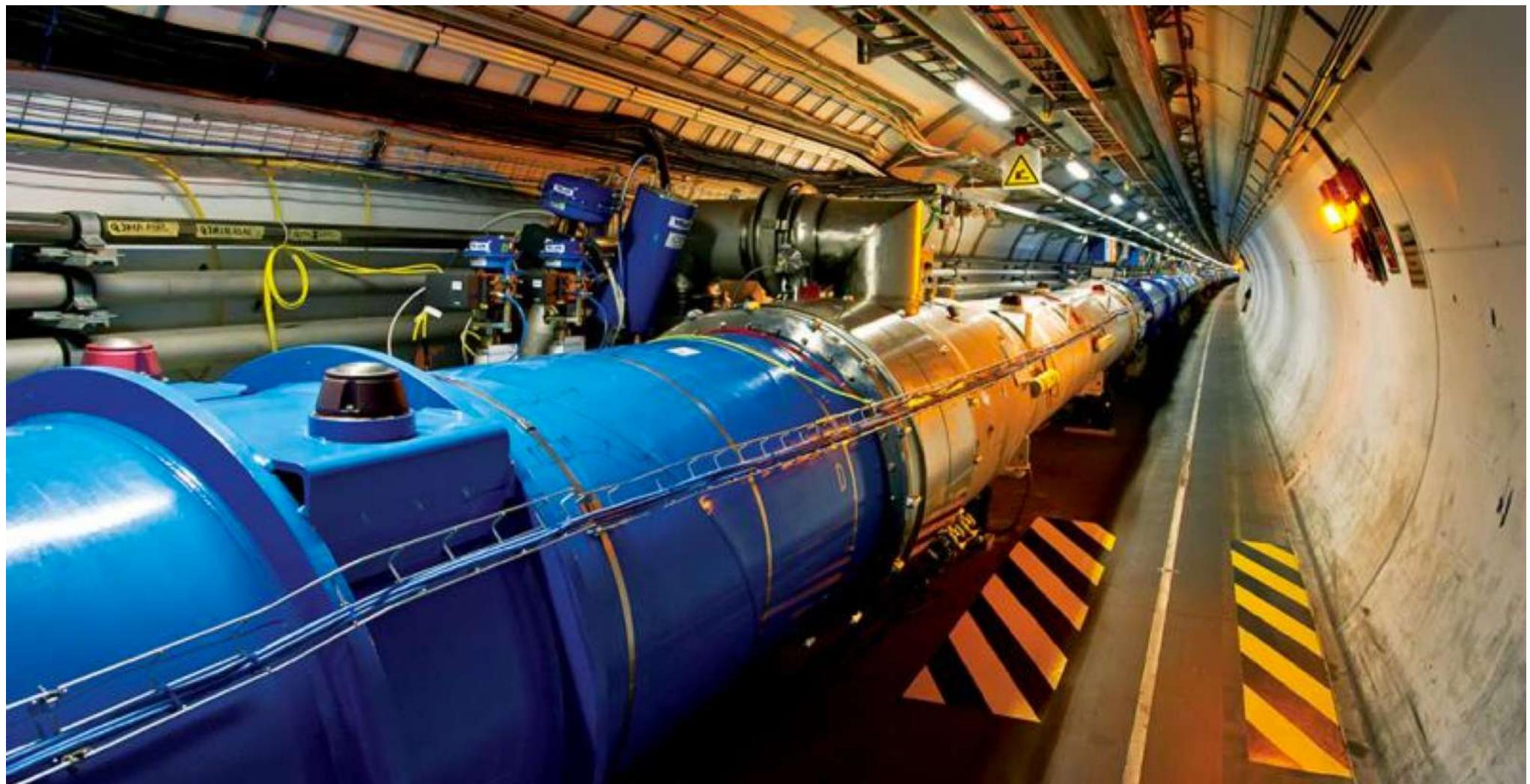


Picture: pcanzo/123RF

Big Data in the Sciences

Physics: CERN pioneers, produces 30 PB/year

Picture: CERN



Wait!
Actually,
that was
three years
ago!

Physics: CERN pioneers, produces 50 PB/year

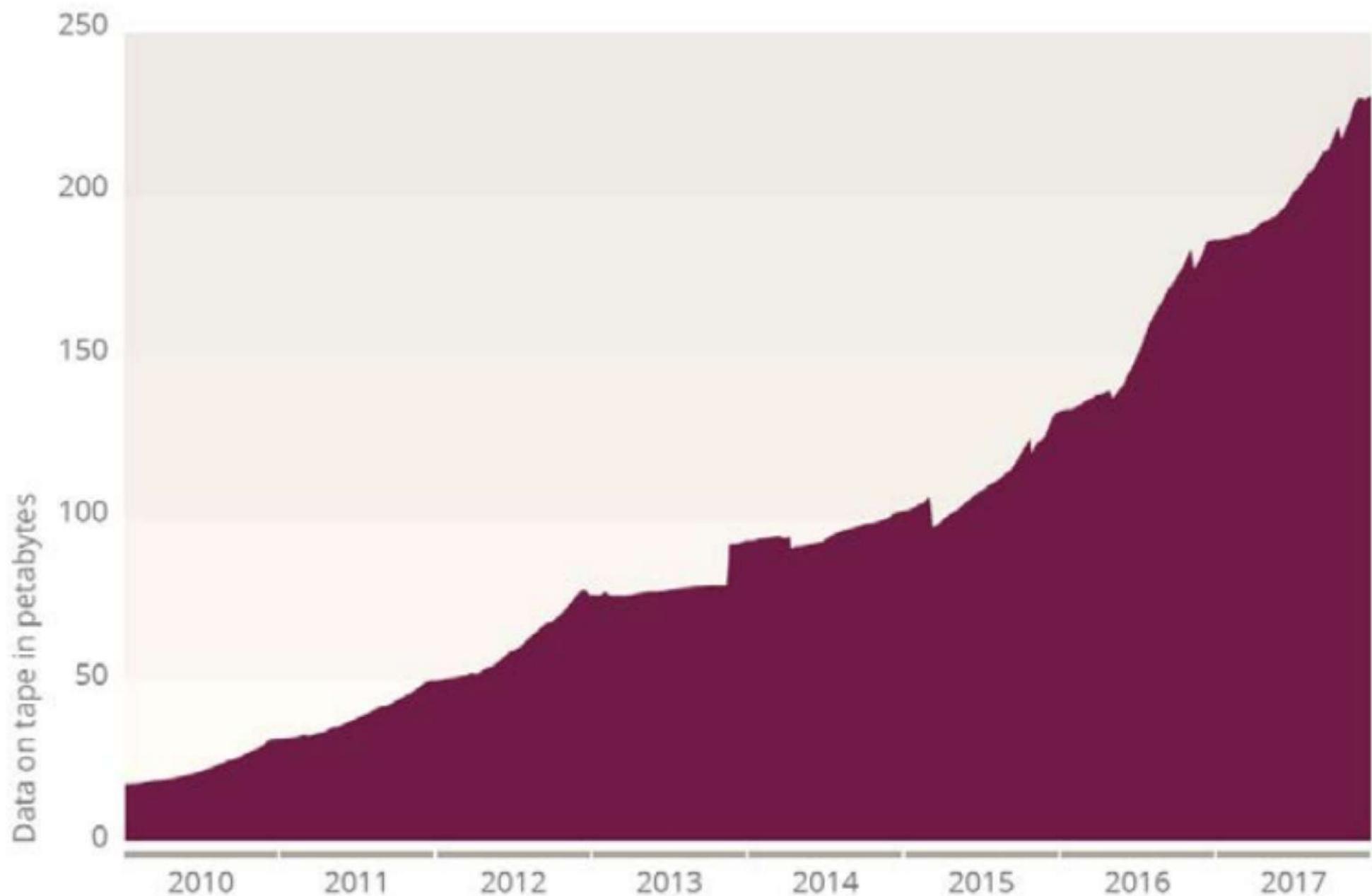
1,000,000,000 collisions/second

15,000 servers

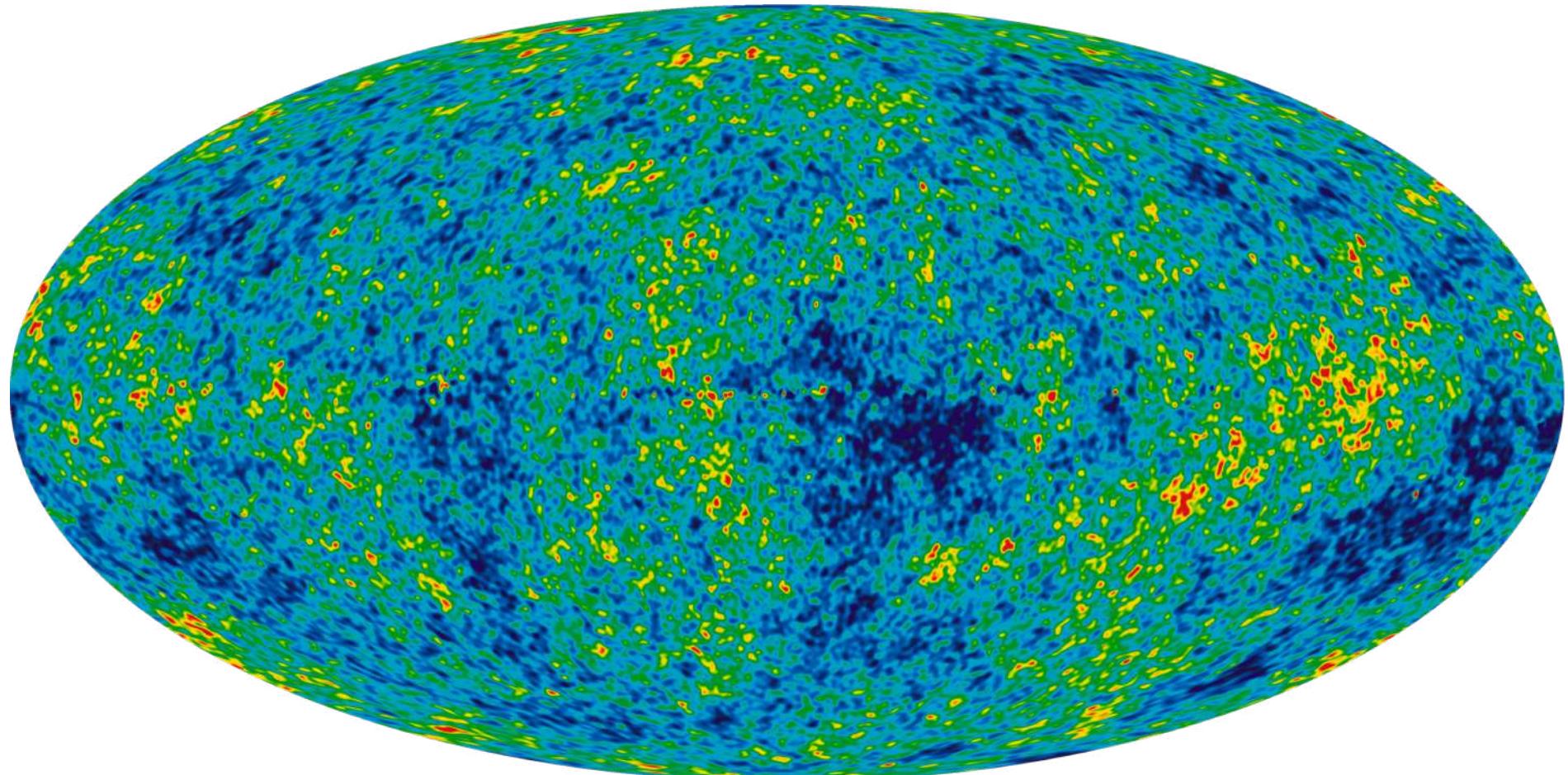
230,000+ cores

More on <http://monit-grafana-open.cern.ch/d/000000884/it-overview?orgId=16>

EVOLUTION OF THE TOTAL AMOUNT OF DATA STORED ON TAPE AT CERN IN PB



Astronomy: Sloan Digital Sky Survey



Astronomy: Sloan Digital Sky Survey

Since 2000, now in phase IV till 2020

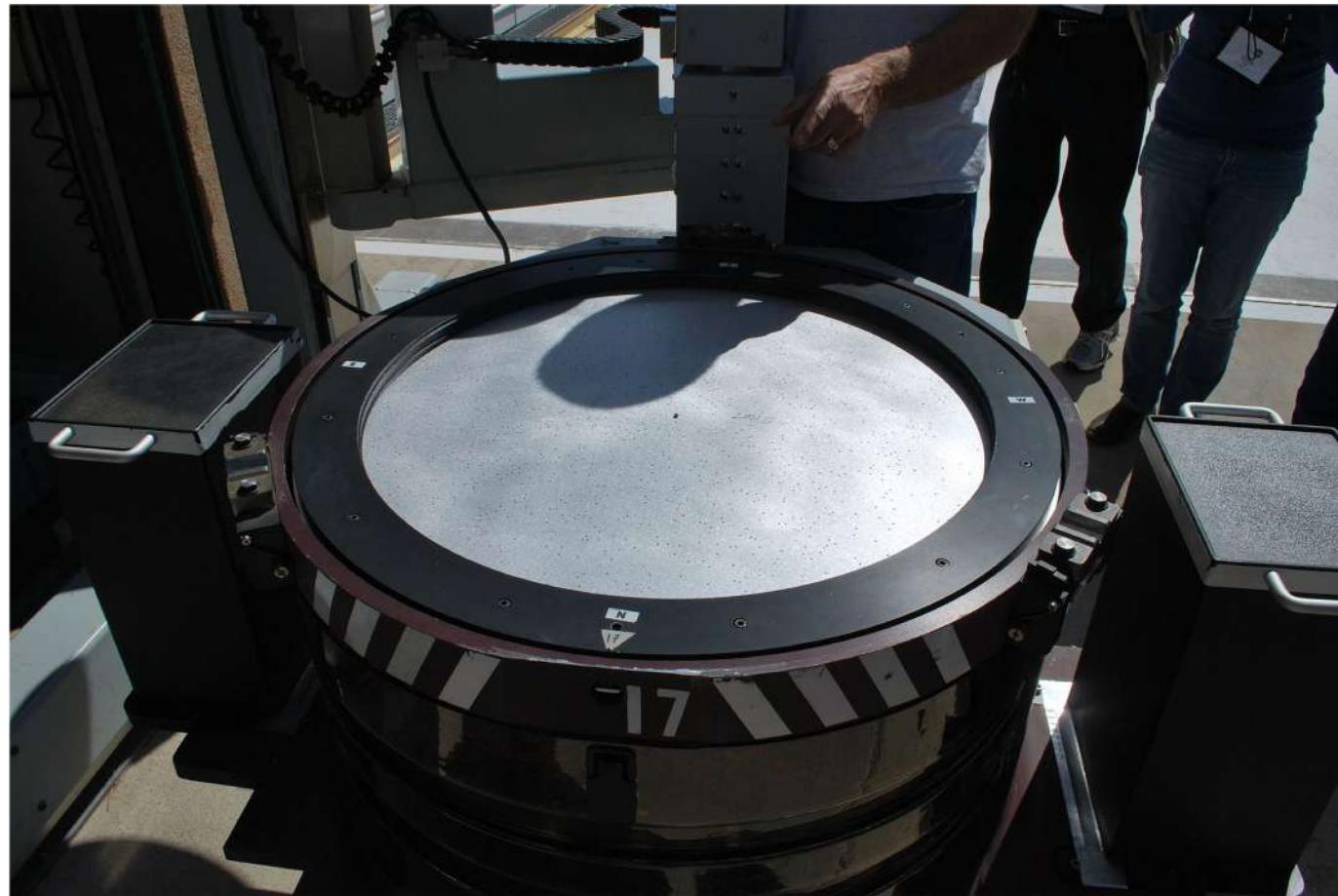
The **most detailed**
3D maps of the Universe
ever made

35% covered so far

1G objects, 4 spectra

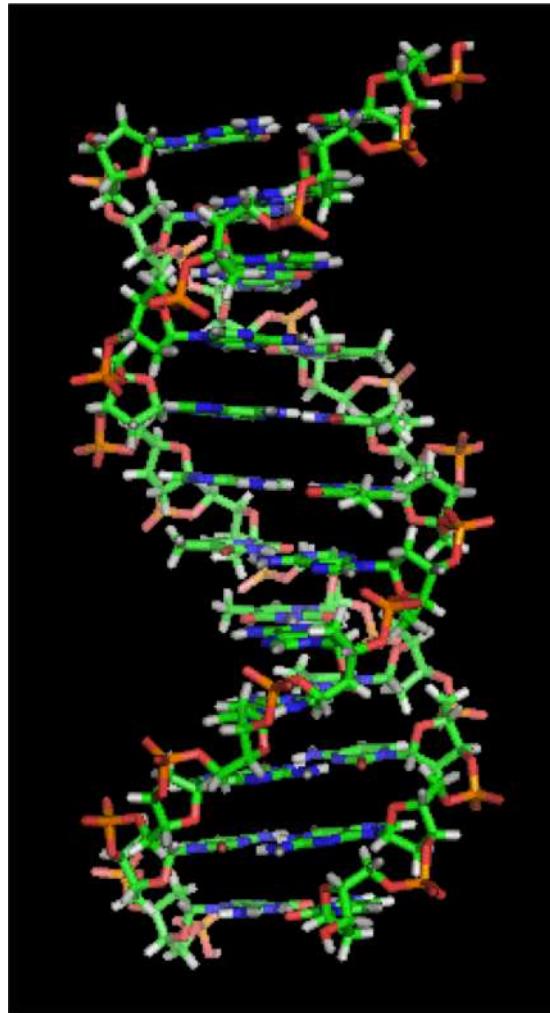
Astronomy: Sloan Digital Sky Survey

200 GB/night



Picture: Wikipedia/EdPost

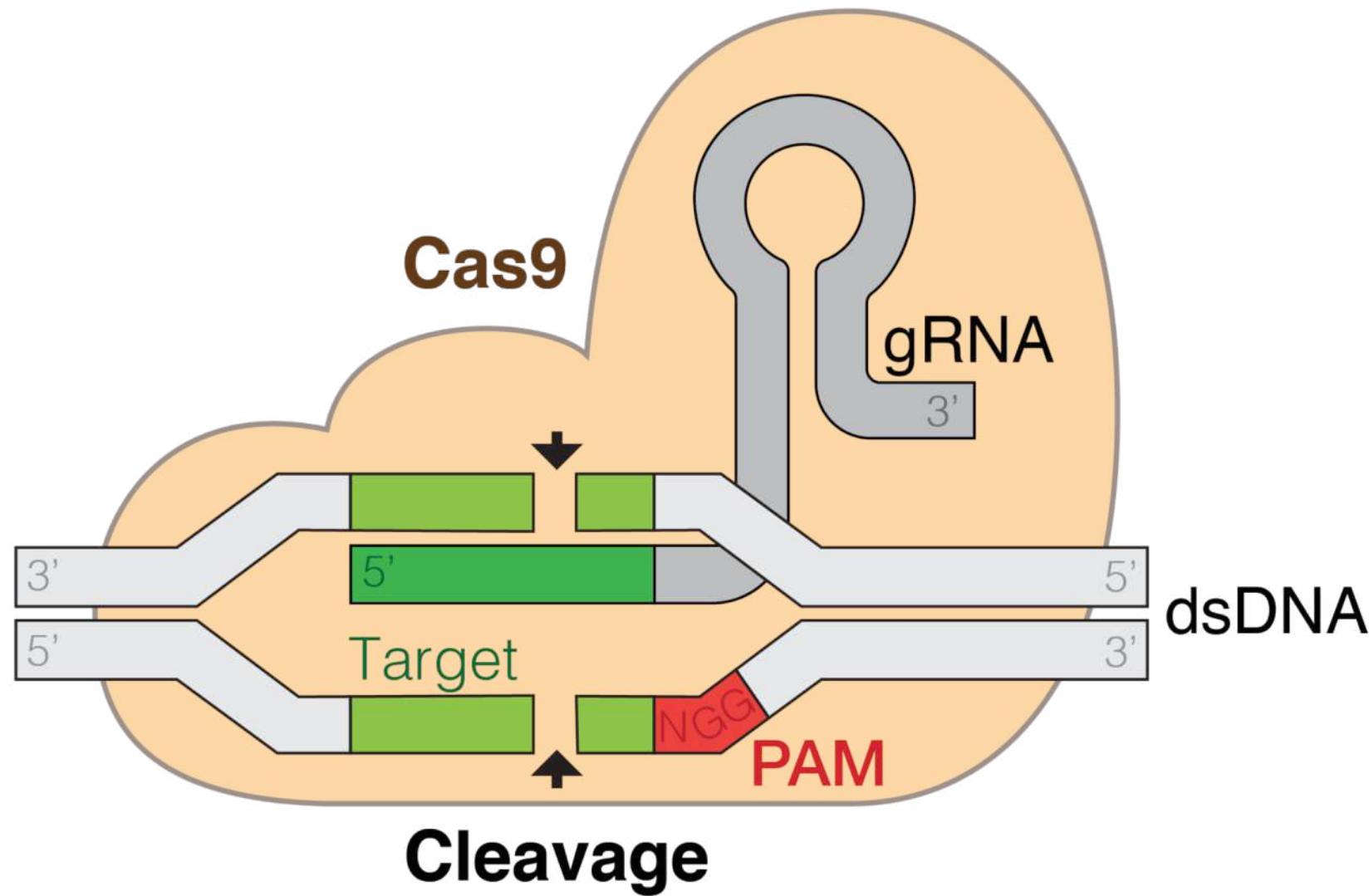
Genomics: the complete human genome



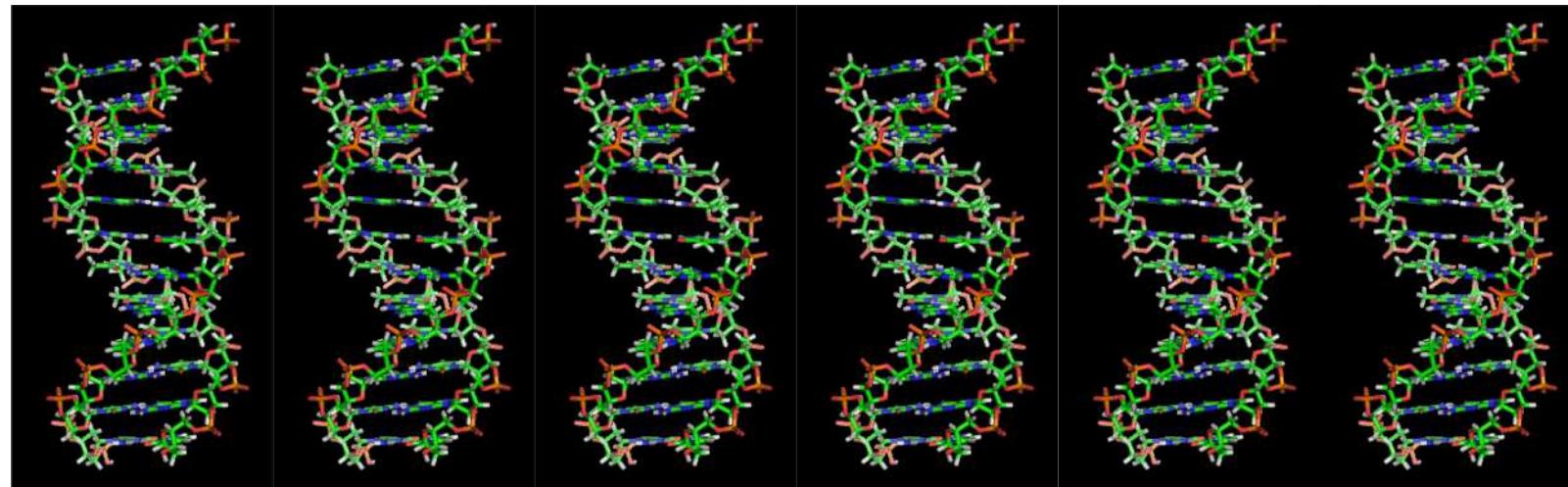
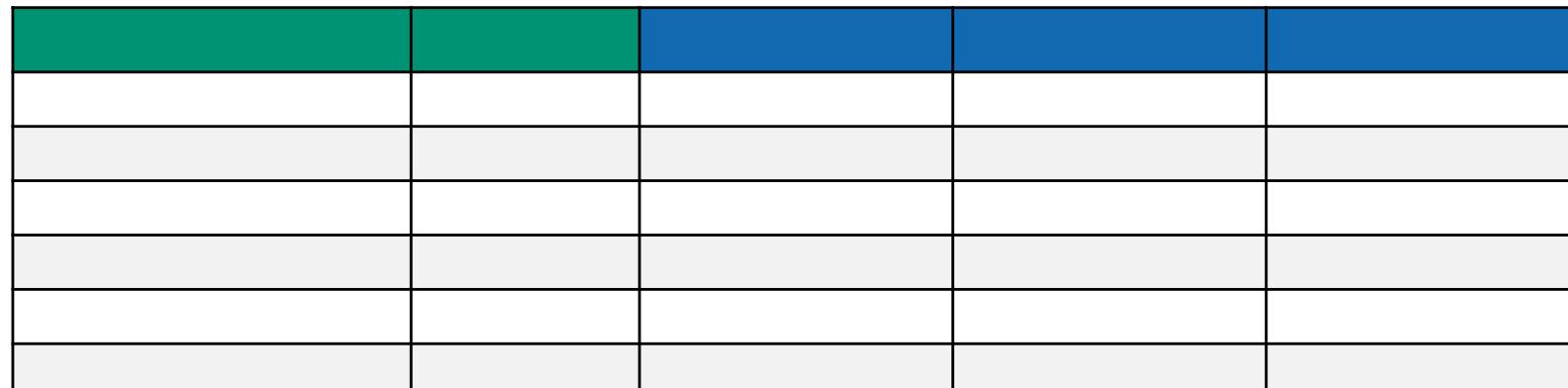
Picture: Wikipedia/Zephyris

3B base pairs

Genomics: CRISPR-Cas9



New (2018): DNA as a storage layer



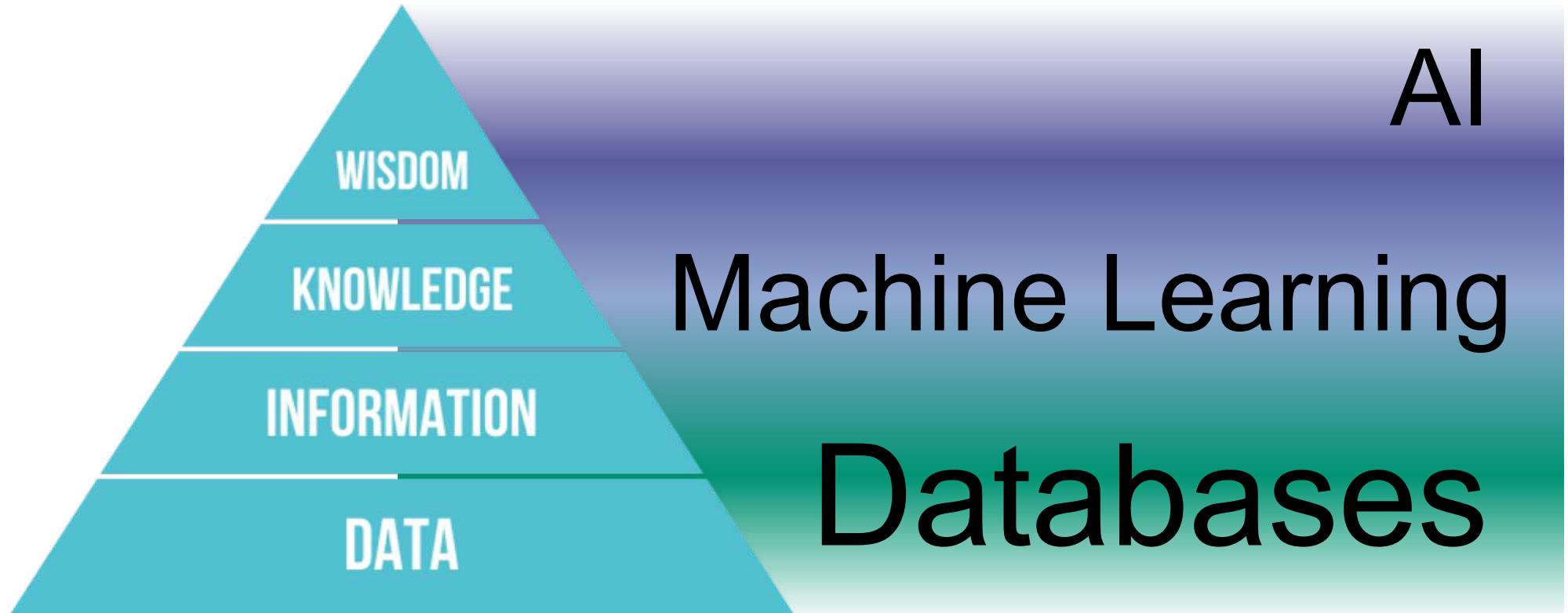


Lecture Scope

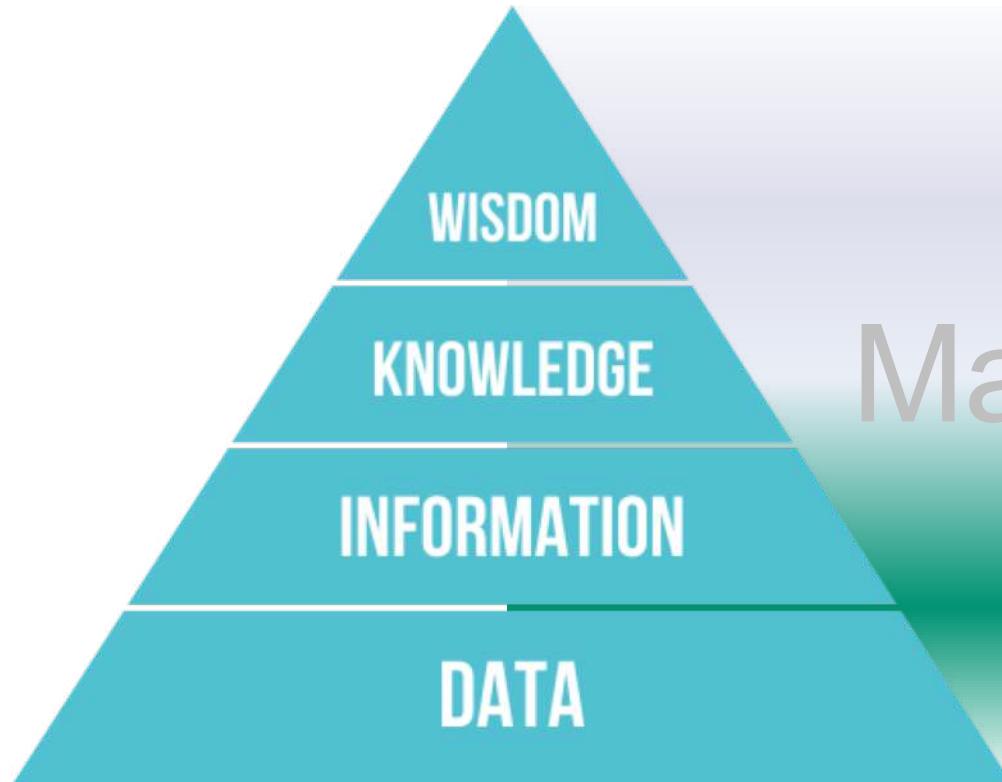
Lecture scope



Lecture scope



Lecture scope: databases only



Machine Learning

Databases

AI

Lecture Team



Ghislain Fourny



Rodrigo Bruno
(Head TA)



Damien Desfontaines
(TA)



Dan Graur
(TA)



Niels Gleinig
(TA)



Giulia Lanzillotta
(TA) 107

Lecture Overview



	Concepts	Technologies	
Storage	Object storage	S3, Azure Blob Storage	
	Distributed file systems	HDFS	
	Syntax	XML, JSON	
Models	Wide column stores	HBase	
	Data models and schemas	XML/JSON Schema	
Processing	2-step distributed query processing	Hadoop MapReduce	
	Resource management	YARN	
	DAG-based distributed query processing	Spark	
Management	Document storage	MongoDB	
	Query languages	JSONiq	



What is expected

Attendance of the **weekly lecture**
(2 hours/w Tuesdays 10-12)

What is expected

Attendance of the **weekly lecture**
(2 hours/w Tuesdays 10-12)

Attendance of the **exercise session**
(2 hours/w Wednesdays/Fridays)

What is expected

Attendance of the **weekly lecture**
(2 hours/w Tuesdays 10-12)

Attendance of the **exercise session**
(2 hours/w Wednesdays/Fridays)

Hands-on self-study, read the books,
play with technology (1-2 hours/w)

What is expected

Attendance of the **weekly lecture**
(2 hours/w Tuesdays 10-12)

Attendance of the **exercise session**
(2 hours/w Wednesdays/Fridays)

Hands-on self-study, read the books,
play with technology (1-2 hours/w)

Passing the **written exam**
(150 minutes, Summer session)

What is expected

Attendance of the **weekly lecture**

(2 hours/w Tuesdays 10-12)

Attendance of the **exercise session**

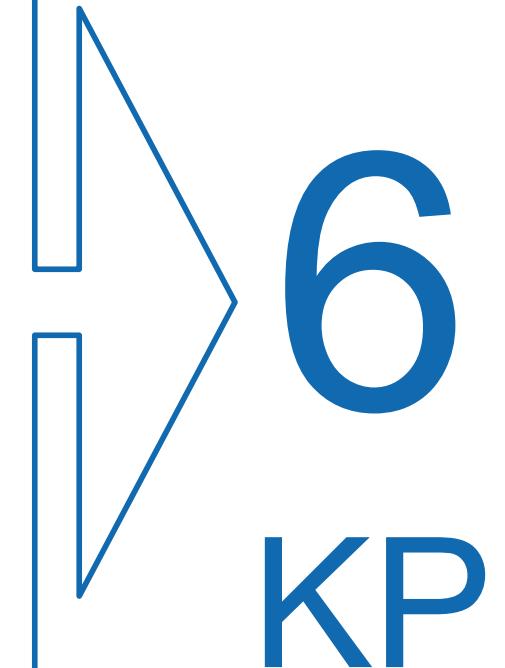
(2 hours/w Wednesdays/Fridays)

Hands-on self-study, read the books,

play with technology (1-2 hours/w)

Passing the **written exam**

(150 minutes, Summer session)



Bonus points!

Bonus points!



0.25

Bonus points!

You SHOULD solve the weekly exercise sheets (theoretical, practical)

Bonus points!

You SHOULD solve the weekly exercise sheets (theoretical, practical)



We will grade the exercises marked as such (25 of them)

Bonus points!

You SHOULD solve the weekly exercise sheets (theoretical, practical)



We will grade the exercises marked as such (25 of them)



You get 0.01 extra point per passed assignment: at most 0.25

Bonus points!

You SHOULD solve the weekly exercise sheets (theoretical, practical)



We will grade the exercises marked as such (25 of them)



You get 0.01 extra point per passed assignment: at most 0.25

"0.n" will thus be added to your exam grade before rounding

Self-study: Azure

