

# BIG DATA AND HADOOP

**Venkatesh Vinayakarao**

[venkateshv@cmi.ac.in](mailto:venkateshv@cmi.ac.in)

<http://vvtesh.co.in>

---

Chennai Mathematical Institute

---

Data is the new oil. - Clive Humby, 2006.

# Know Your Instructor

**BE (Computer Science and Engineering)**

Java/J2EE  
Developer

**MS (Information Technology)**

SDE, Search  
Technologies  
Group, Bing,  
Microsoft

Principal  
Engineer, Cloud  
Platforms Group,  
Yahoo

**PhD (Computer Science)**

Intern, Porting ML  
Models to Azure,  
Microsoft  
Research

# Agenda

- Introduction to Big Data
- Course Dynamics
- Evolution of Systems and Technologies
  - Data Storage
  - Data Processing

# What Comes Next?

byte

kilobyte

megabyte

gigabyte

??

???

????

?????

# Sizes

Name	Size
Byte	8 bits
Kilobyte	1024 bytes
Megabyte	1024 kilobytes
Gigabyte	1024 megabytes
Terabyte	1024 gigabytes
Petabyte	1024 terabytes
Exabyte	1024 petabytes
Zettabyte	1024 exabytes
Yottabyte	1024 zettabytes

# The Impact of Big Data



## Your train is on time thanks to **big data**

TNW - 31-Dec-2019

Thanks to thousands of sensors and **big data** analytics, train ... It's this data that keeps the Dutch rail network moving, and helps NS deliver a ...



## The power of **data** in smart city developments

Independent Australia - 03-Jan-2020

Other fascinating **big data** developments that were presented included ... led to the production of the Australian **Cancer Atlas** — an interactive, ...



## At HCA Healthcare, Real-Time **Data Saves Lives**

RTInsights (press release) (blog) - 01-Jun-2019

At HCA Healthcare, Real-Time **Data Saves Lives** ... “Our existing **data** infrastructure was designed for **large-scale** business intelligence and ...

# Big Data is Ubiquitous

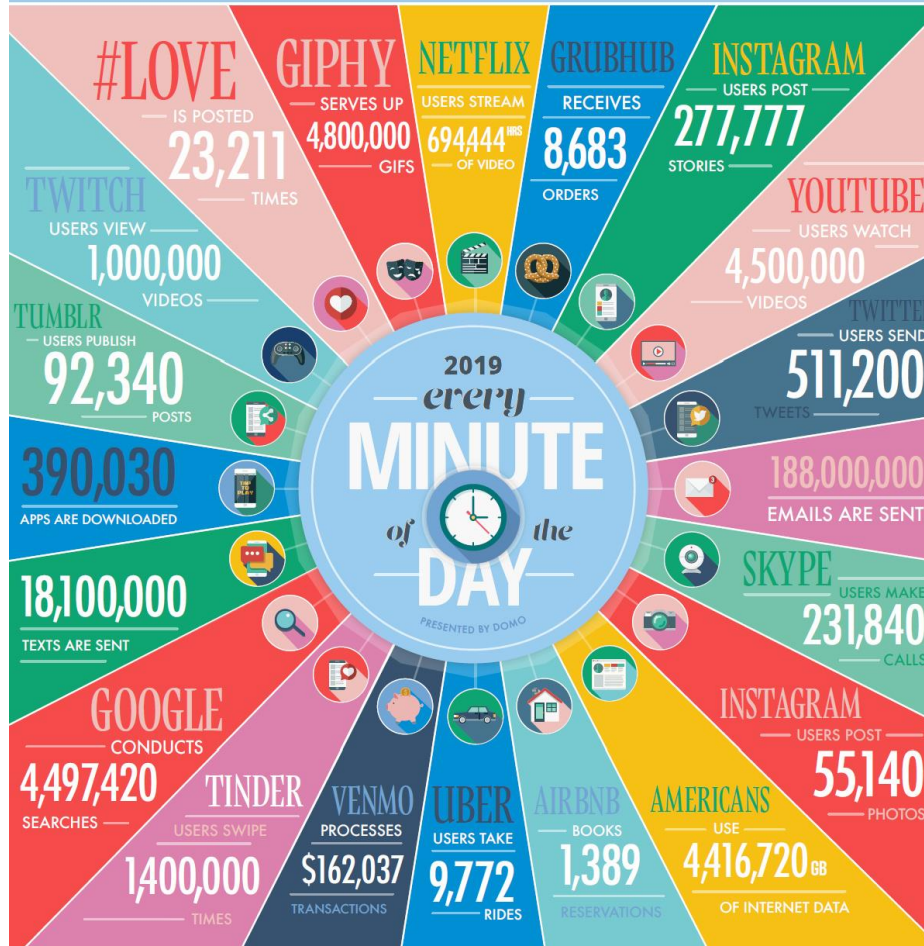
- Facebook (**per day** statistics)
  - 1.5 billion people are active on Facebook **daily!**
  - More than 300 million photos get uploaded **per day!**
  - Totally, more than 2.5 Trillion posts!
- Facebook (per minute statistics)
  - **Every minute** there are 510,000 comments posted and 293,000 statuses updated!
- Youtube (**per minute** statistics)
  - Users watch 4,146,600 YouTube videos!
- **Guess... How many emails are sent per minute?**



# DATA NEVER SLEEPS 7.0

How much data is generated *every minute*?

There's no way around it: big data just keeps getting bigger. The numbers are staggering, and they're not slowing down. By 2020, there will be 40x more bytes of data than there are stars in the observable universe. In our 7th edition of Data Never Sleeps, we bring you the latest stats on how much data is being created in every digital minute.



SOURCES: STATISTA, INTERNET LIVE STATS, EXPANDED RAMBLINGS, NATIONAL ASSOCIATION OF CITY TRANSPORTATION OFFICIALS, WIRED





# And, It is Growing!

JAN  
2021

## GLOBAL DIGITAL GROWTH

THE YEAR-ON-YEAR CHANGE IN DIGITAL ADOPTION

INTERNET USER NUMBERS NO LONGER INCLUDE DATA SOURCED FROM SOCIAL MEDIA PLATFORMS, SO VALUES ARE **NOT COMPARABLE** WITH PREVIOUS REPORTS

TOTAL  
POPULATION



we  
are  
social

**+1.0%**

JAN 2021 vs. JAN 2020

**+81 MILLION**

UNIQUE MOBILE  
PHONE USERS



we  
are  
social

**+1.8%**

JAN 2021 vs. JAN 2020

**+93 MILLION**

INTERNET  
USERS\*



KEPIOS

**+7.3%**

JAN 2021 vs. JAN 2020

**+316 MILLION**

ACTIVE SOCIAL  
MEDIA USERS\*



we  
are  
social



Hootsuite®

**+13.2%**

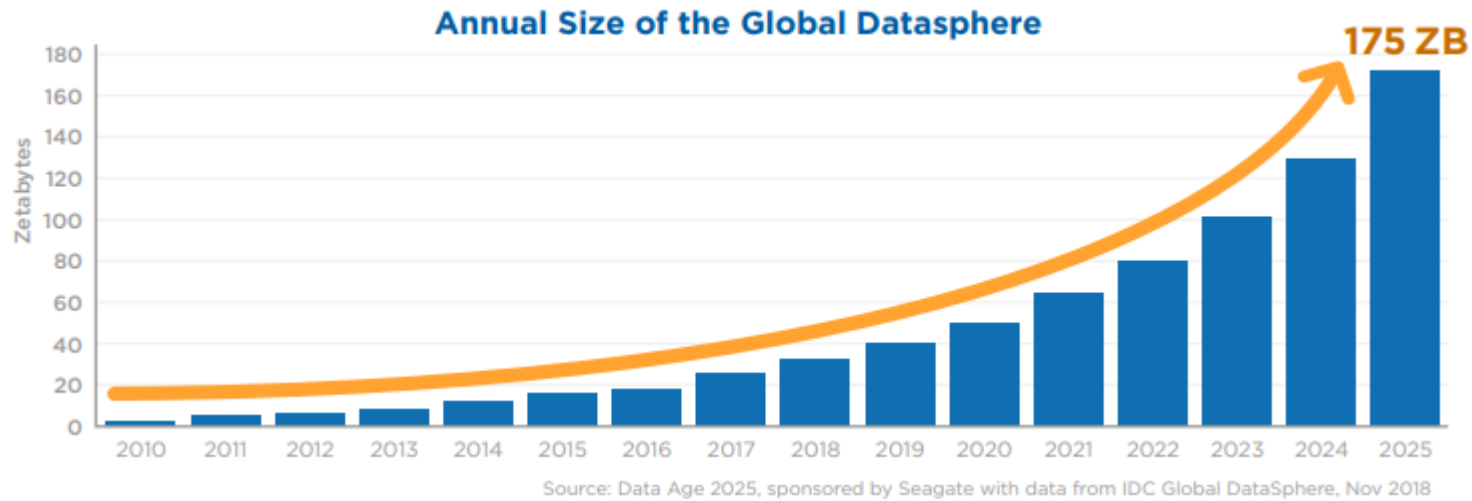
JAN 2021 vs. JAN 2020

**+490 MILLION**

9

**SOURCES:** THE U.N.; LOCAL GOVERNMENT BODIES; GSMA INTELLIGENCE; ITU; GWI; EUROSTAT; CNNIC; APJII; SOCIAL MEDIA PLATFORMS' SELF-SERVICE ADVERTISING TOOLS; COMPANY EARNINGS REPORTS; MEDIASCOPE. **\*ADVISORIES:** INTERNET USER NUMBERS NO LONGER INCLUDE DATA SOURCED FROM SOCIAL MEDIA PLATFORMS, SO VALUES ARE **NOT COMPARABLE** TO DATA PUBLISHED IN PREVIOUS REPORTS. SOCIAL MEDIA USER NUMBERS MAY NOT REPRESENT UNIQUE INDIVIDUALS. **◆ COMPARABILITY ADVISORY:** SOURCE AND BASE CHANGES.

# Data Growth



Mankind's quest to digitize the world!  
33 ZB (2018) → 175 ZB (2025)  
size of global datasphere\*

\*Source: <https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf>

**Global datasphere is growing!**

How have the computers evolved to capture,  
process and analyze these data?

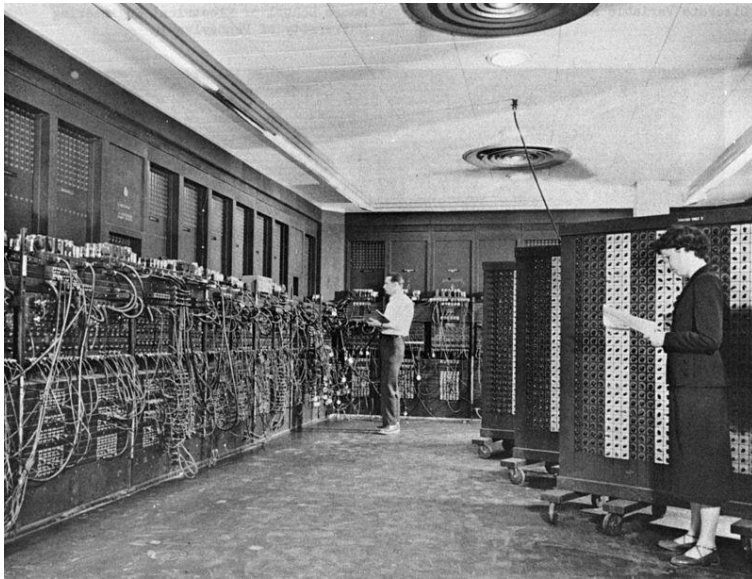
# Course Dynamics

<https://vvtesh.github.io/teaching/bdh-2022.html>

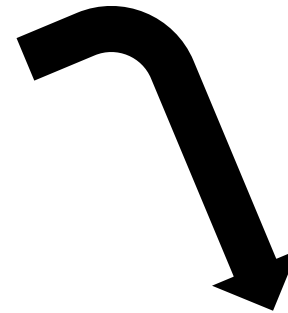
# Student Presentations

- Please register
  - your choice of presentation topic,
  - your team details (Team size: 3 or 4)
- Deadlines
  - (1 Mark) Register your topic before Feb 10<sup>th</sup> .  
Registration link will be available on moodle.
  - (2 Marks) Send a one-page abstract of the talk before mid-term. Clearly mention the team members.
  - (12 Marks) Complete the presentation one week before the final exam.

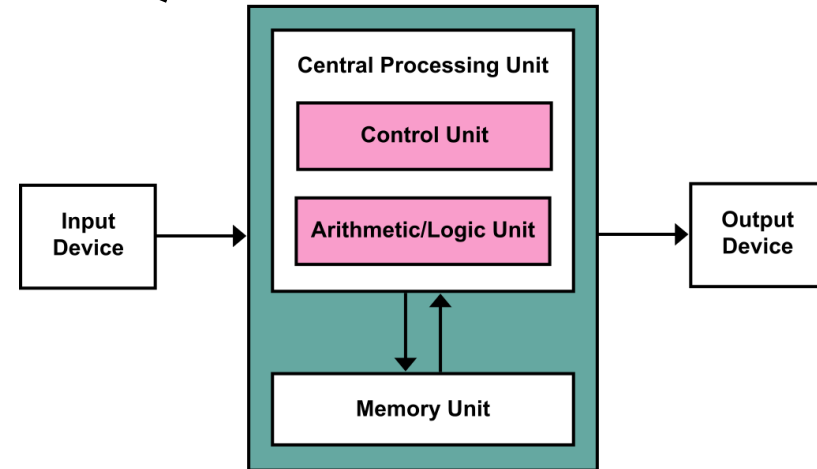
# Evolution of Computers



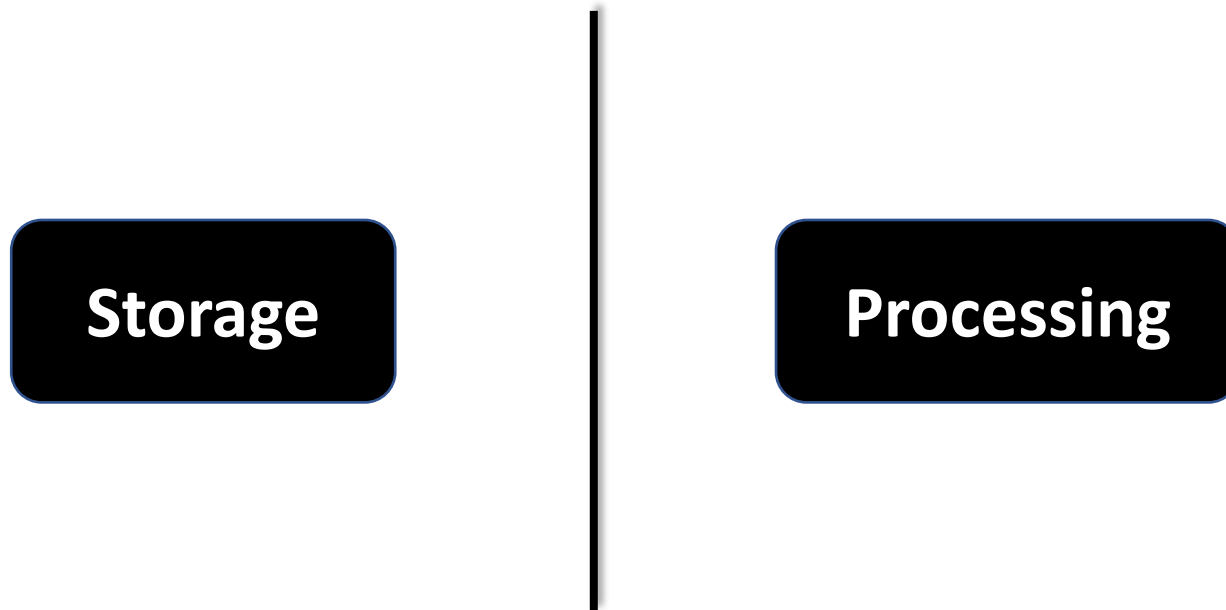
**ENIAC**  
Early 1900s



**Stored-program  
Von Neumann  
Architecture  
1940**



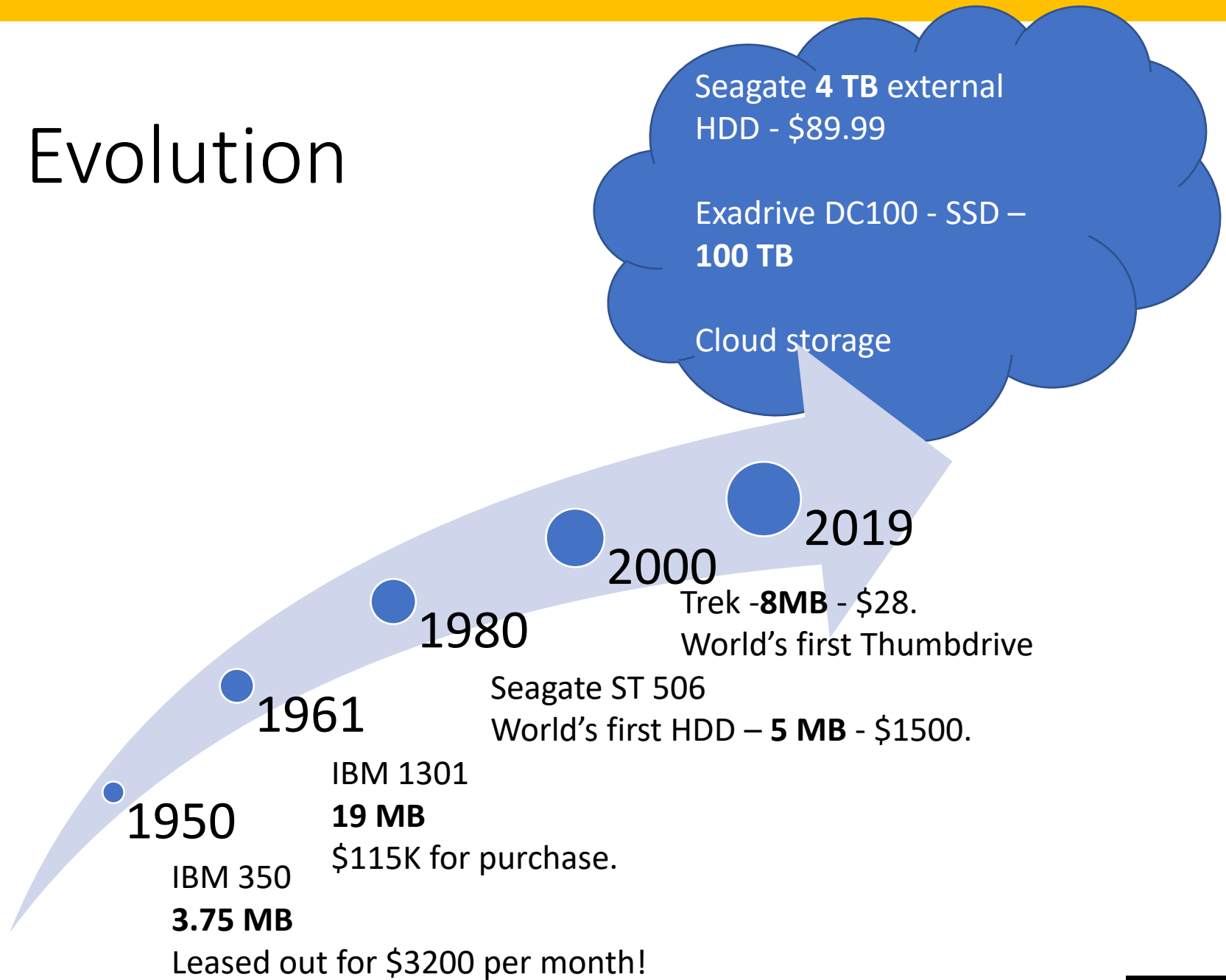
# Two Kinds of Problems



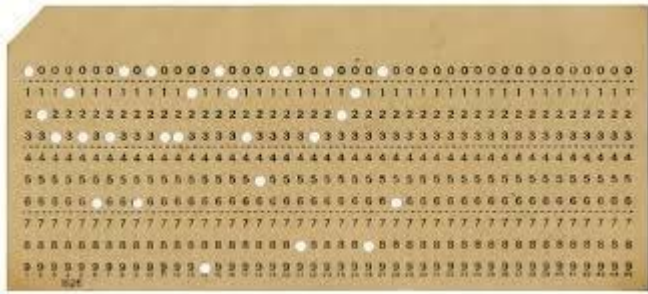
# Data Storage



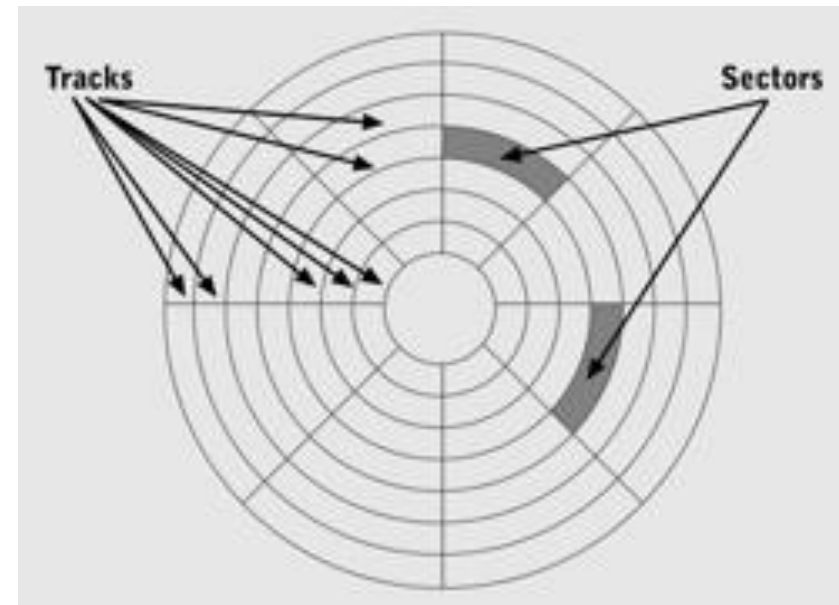
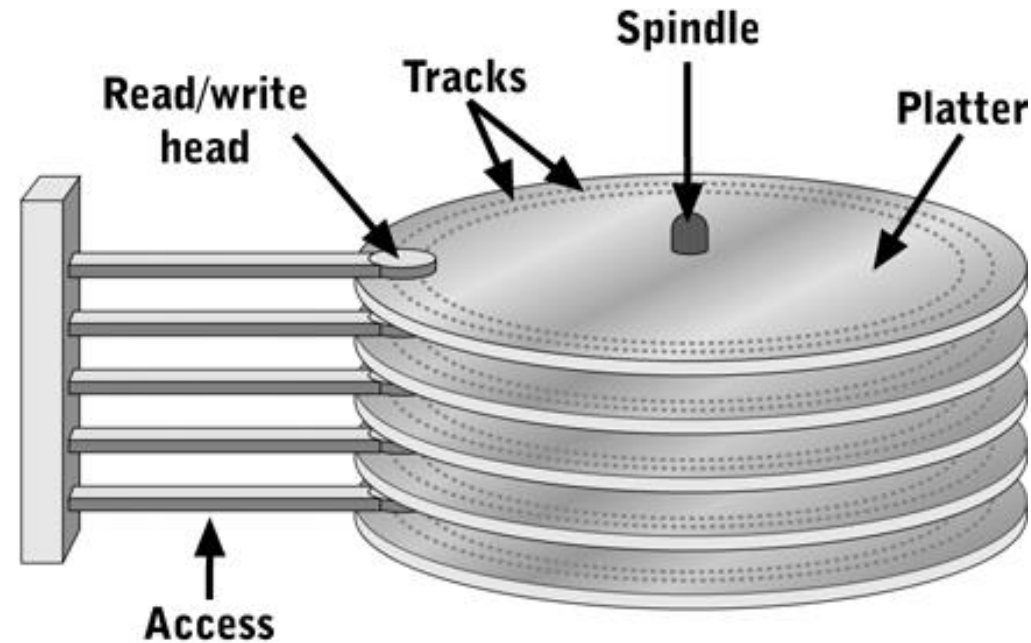
# Evolution



# (Secondary) Storage Technologies



# Disk Drive and Access Time



Source: Systems Architecture, Fifth Edition



Roll over image to zoom in

## Seagate 500GB SATA Laptop Hard Disk

by [Seagate](#)

★★★★☆ 279 ratings | 493 answered questions

M.R.P.: ₹2,999.00

Price: **₹ 1,433.00** + ₹ 77.00 Delivery charge [Details](#)

You Save: **₹ 1,566.00 (52%)**

Inclusive of all taxes



Pay on  
Delivery



10 Days  
Replacement



Amazon  
Delivered



1 Year  
Warranty

**In stock.**

Delivery by: **Jan 8 - 10** [Details](#)

[📍 Deliver to Venkatesh - Chennai 600014](#)

Sold by [KCM\\_STORE](#) (3.9 out of 5 stars | 29 ratings).

**New (20)** from **₹ 1,510.00** + FREE Shipping

- 500 GB capacity
- 5400 RPM spin speed, 16 MB cache buffer
- Designed for durability and low-power consumption
- SATA 3GB interface with native command queuing
- Perpendicular recording technology for increased storage capacity
- Fast performance and whisper quiet acoustics

# Average Access Time

- Head switching time is considered negligible (H)
- Head seek time (S)
- Rotational delay = Time taken for  $\frac{1}{2}$  a rotation (average) (R)
- Read time = time to spin an entire sector (T)
- Average Access Time =  $H + S + R + T$

\*Sector is a minimum storage unit

# Quiz

- If disk spins at 6000 RPM, compute the rotational delay.

# Quiz

- If disk spins at 6000 RPM, compute the rotational delay.
  - One turn takes  $1/6000$  min or  $1/100$  sec = 10ms
  - $\frac{1}{2}$  a turn takes 5ms.

# Read Time

- If the drive spins at 6000RPM and the disk has 20 sectors per track, what is the read time?

- Time for 1 full spin is  $\frac{1}{6000} \text{ min} = \frac{1}{100} \text{ sec} = 10\text{ms}$

- Time for 1/20 of a spin is  $10\text{ms} \times \frac{1}{20} = 0.5\text{ms}$



# Average Access Time

- Drive spins at 7200RPM and has average seek time of 8ms. The disk has 24 sectors per track. What is the average access time?

Head seek time	0.008 sec (Given)
Rotational delay	$1/120 * (1/2) = 0.0042$ sec
Read time	$0.0084$ (full spin) / 24 sectors = $0.00035$ sec
<b>Avg Seek Time</b>	<b>= <math>0.008 + 0.0042 + 0.00035</math></b> <b>= <math>0.01255</math> sec or 12.55 ms</b>

# Characteristics

Attribute	Description
Speed	Time to read/write
Volatility	Data persistence even when powered off
Access Method	Serial, Parallel
Portability	Internal, External
Capacity	Volume of data storage

# Storing/Managing/Processing Data

- RDBMS
- ETL
- OLTP
- Data Warehouse
- Data Lake
- Cloud Storage
- STaaS

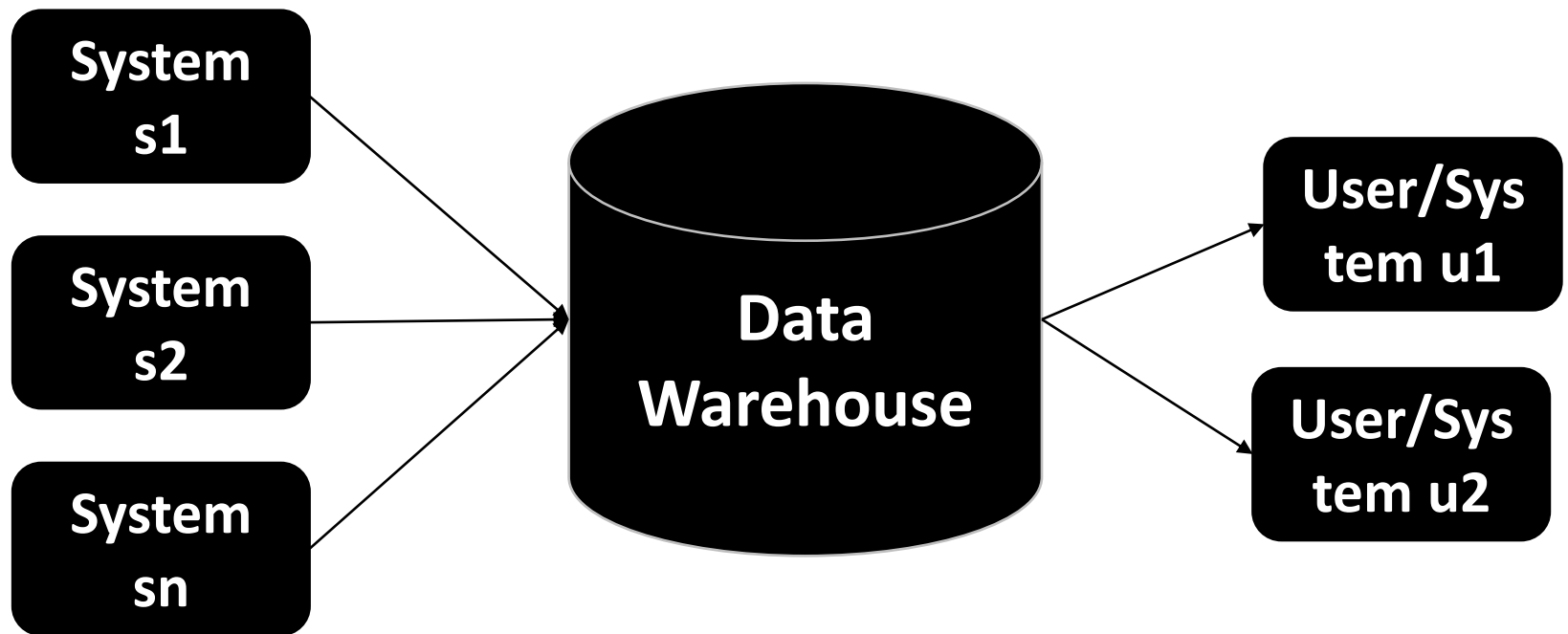
# ETL

- Extract, Transform and Load (ETL)
  - general procedure followed to address data variety
- Variety of data sources
  - Tabs, Sensors, Desktops, Bots, Multiple databases, Files,...
- Variety of data formats
  - Text, PDFs, XML, JSON, Images, Videos, ...

# Fast OLTP

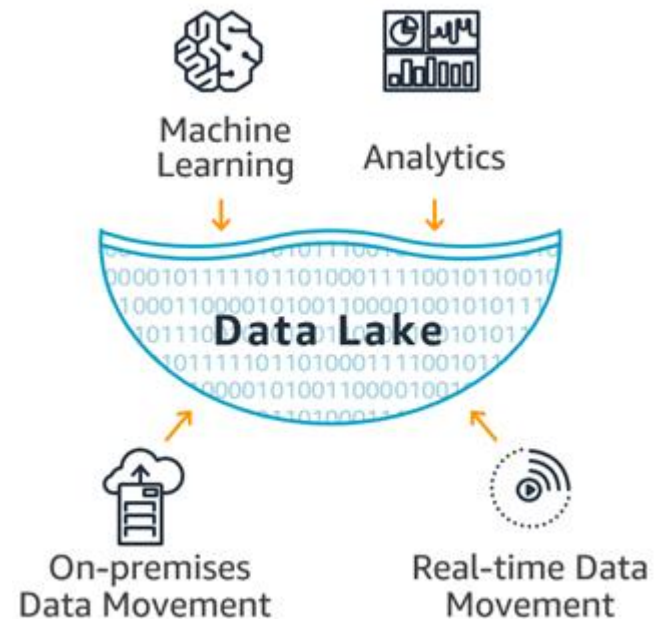
- Online Transaction Processing (OLTP)
- Real-time/Near Real-time Performance. Finds application in:
  - Banking
  - Railway Reservations
  - Stock Market Trading
    - Handle transactions in milliseconds.
      - VoltDB, MemSQL, ...

# Data Warehouse



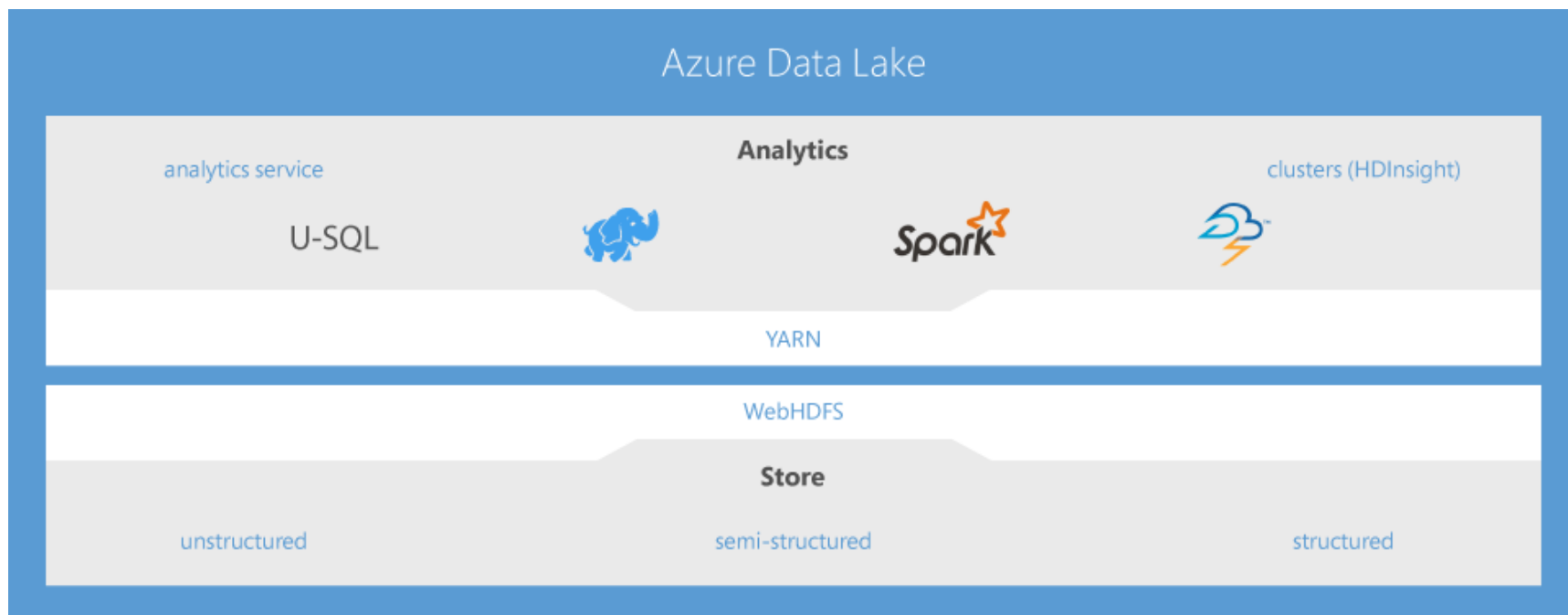
# Data Lakes

- No schema definition.
- Store everything
  - often without or with very little pre-processing, /cleaning.
- Use ML, analytics to query, or gather insights.



Source: <https://aws.amazon.com/big-data/datalakes-and-analytics/what-is-a-data-lake/>

# Microsoft's Azure Data Lake



More details at <https://azure.microsoft.com/en-us/resources/videos/azure-data-lake-making-big-data-easy/>



# Data Warehouse Vs. Data Lake

Characteristics	Data Warehouse	Data Lake
Data	<b>Relational</b> from transactional systems, operational databases, and line of business applications	<b>Non-relational and relational</b> from IoT devices, web sites, mobile apps, social media, and corporate applications
Schema	Designed <b>prior</b> to the DW implementation (schema-on-write)	Written <b>at the time of analysis</b> (schema-on-read)
Price/Performance	Fastest query results using <b>higher cost storage</b>	Query results getting faster using <b>low-cost storage</b>
Data Quality	Highly <b>curated data</b> that serves as the central version of the truth	Any data that may or may not be curated (ie. <b>raw data</b> )
Users	<b>Business analysts</b>	<b>Data scientists</b> , Data developers, and Business analysts (using curated data)
Analytics	<b>Batch</b> reporting, BI and visualizations	Machine <b>Learning</b> , Predictive analytics, data discovery and profiling

Source: <https://aws.amazon.com/big-data/datalakes-and-analytics/what-is-a-data-lake/>

# Storing on the Cloud

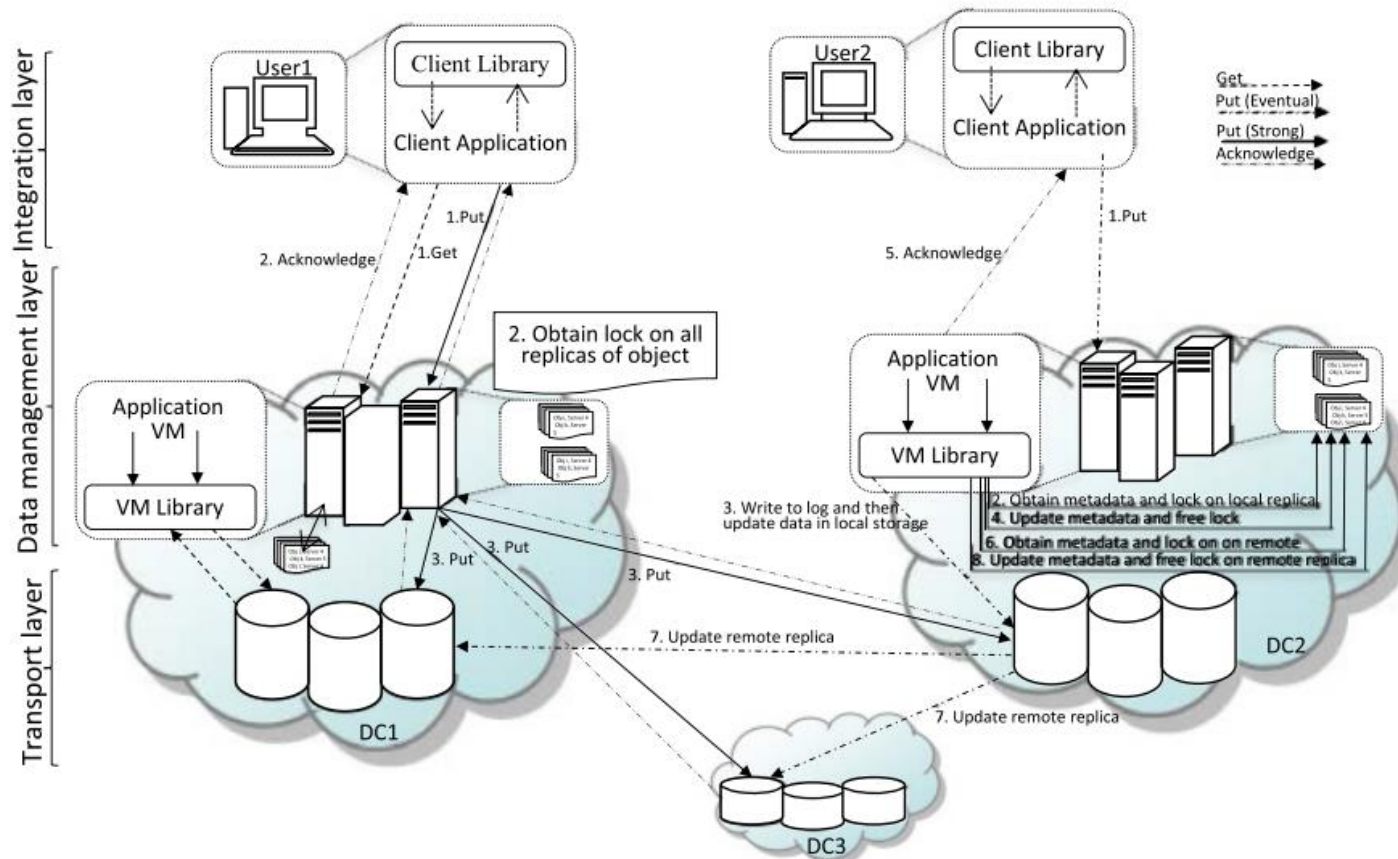
- Gmail: Gives 15 GB of free storage (as of 2020)
- Several online sites for storing images, apps, files, ...
  - Security
  - Ease of sharing
  - Backups
  - Availability



# Storage as a Service (STaaS)

- What is it?
  - A business model in which a company rents space in their storage infrastructure to another company or individual.
- How does it work?
  - STaaS provider rents space
  - cost-per-gigabyte-stored and cost-per-data-transfer basis.
- Benefits
  - Shifting from Capital Expenditure to Operational Expenditure
  - Scale up/down at will (temporarily)

# Cloud Storage

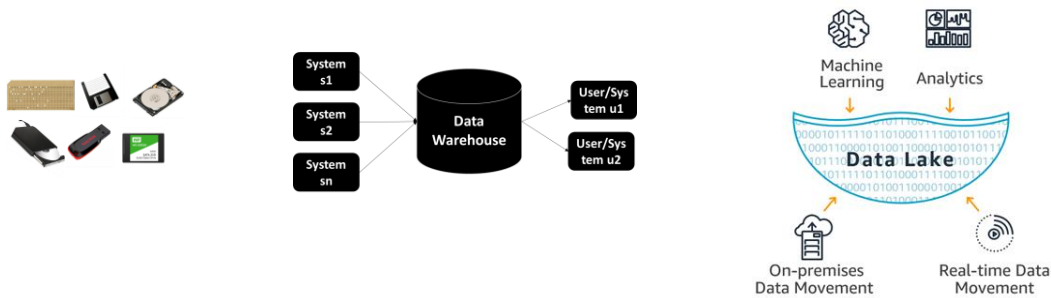


# Cloud Computing



A data center

# Summary



Data Storage - Summary

