# DISTRIBUTED COMPUTING AND BIG DATA

**Venkatesh Vinayakarao**

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)**

Principal Engineer, Search, Here Technologies

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, ...

**rw** Fagen wasanni

The Future of European Healthcare: Embracing Big Data Analytics for Improved Patient Outcomes

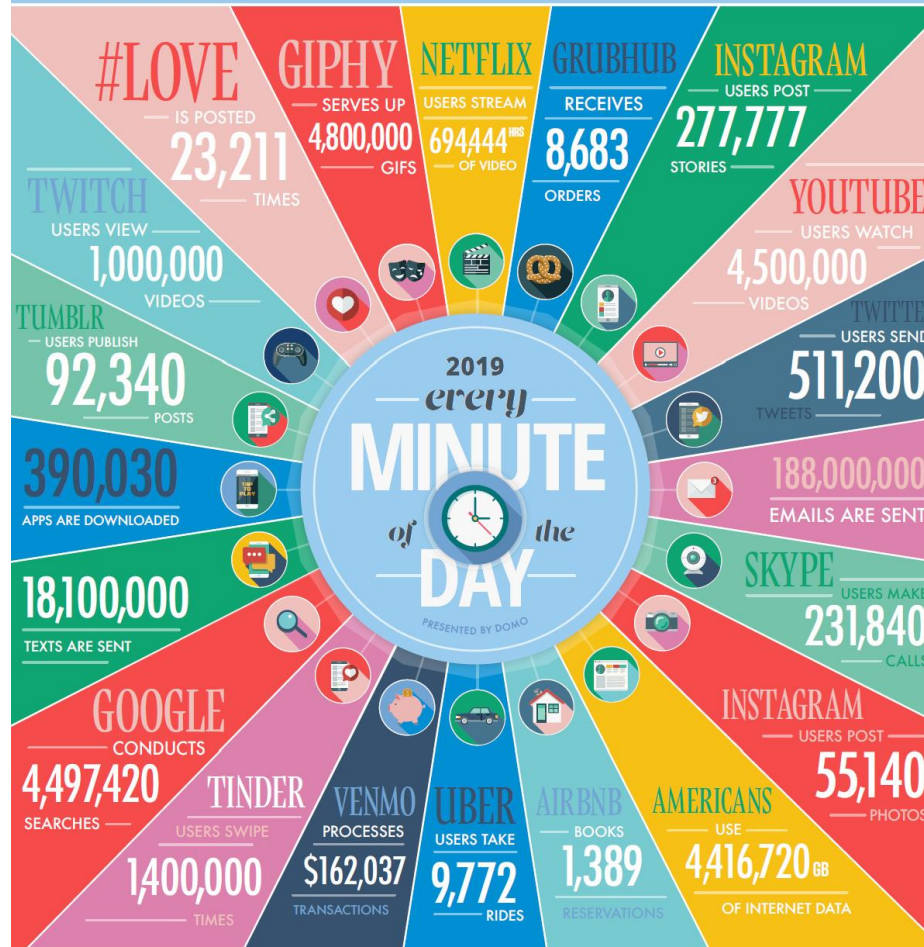3 days ago

# 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!

Source: Forbes

Source: https://www.visualcapitalist.com/big-data-keeps-getting-bigger/

# And, It is Growing!



**ANNUAL DIGITAL GROWTH**
JAN 2019
THE YEAR-ON-YEAR CHANGE IN KEY STATISTICAL INDICATORS

| TOTAL POPULATION | UNIQUE MOBILE USERS | INTERNET USERS | ACTIVE SOCIAL MEDIA USERS | MOBILE SOCIAL MEDIA USERS |
|---|---|---|---|---|
| +1.1% | +2.0% | +9.1% | +9.0% | +10% |
| JAN 2018 – JAN 2019 | JAN 2018 – JAN 2019 | JAN 2018 – JAN 2019 | JAN 2018 – JAN 2019 | JAN 2018 – JAN 2019 |
| +84 MILLION | +100 MILLION | +367 MILLION | +288 MILLION | +297 MILLION |

# DATA NEVER SLEEPS 10.0
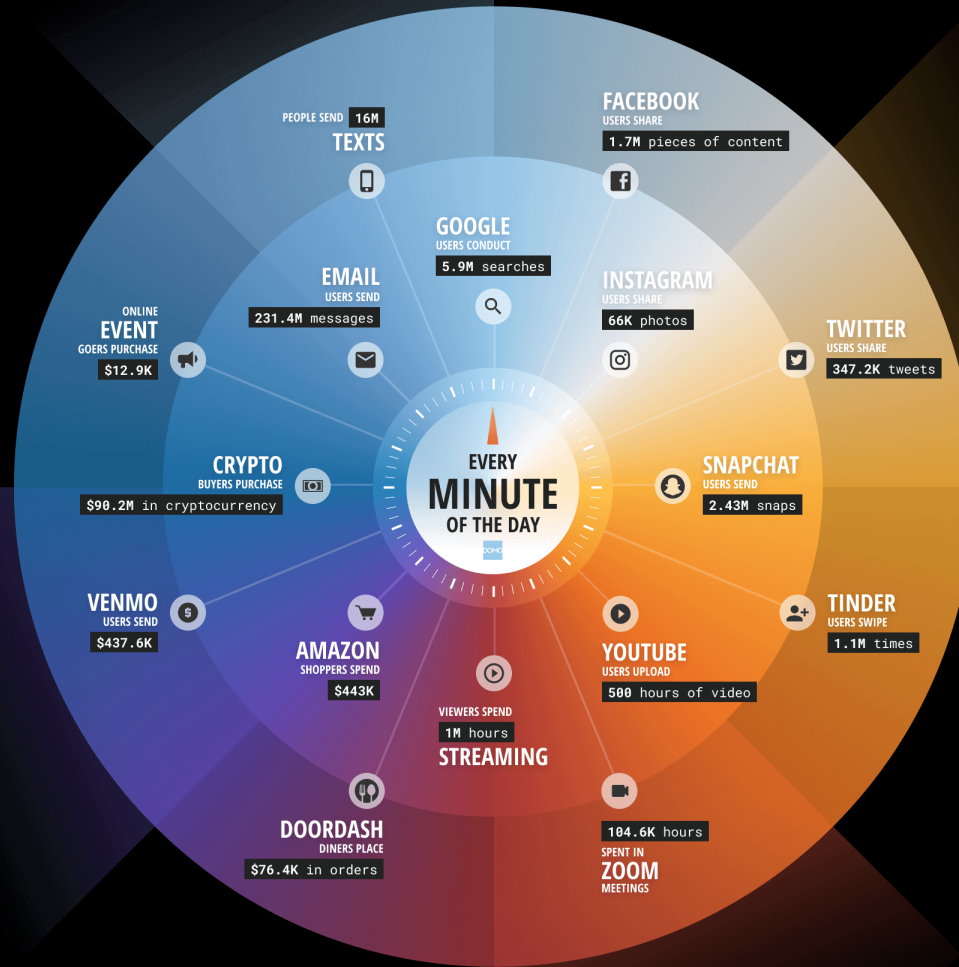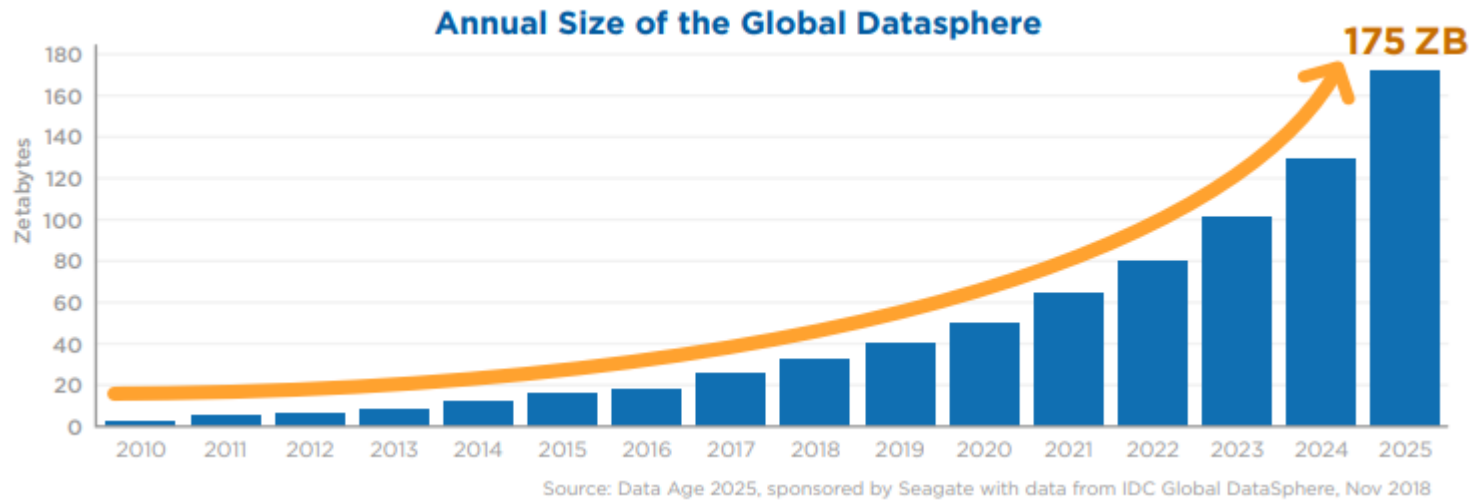
Over the last ten years, digital engagement through social media, streaming content, online purchasing, peer-to-peer payments and other activities has increased hundreds and even thousands of percentage points. While the world has faced a pandemic, economic ups and downs, and global unrest, there has been one constant in society:

our increasing use of new digital tools to support our personal and business needs, from connecting and communicating to conducting transactions and business. In this 10th annual "Data Never Sleeps" infographic, we share a glimpse at just how much data the internet produces each minute from some of this activity, marveling at the volume and variety of information that has been generated.

**EVERY MINUTE OF THE DAY**

PEOPLE SEND **16M**
TEXTS

FACEBOOK
USERS SHARE
**1.7M** pieces of content

GOOGLE
USERS CONDUCT
**5.9M** searches

EMAIL
USERS SEND
**231.4M** messages

INSTAGRAM
USERS SHARE
**66K** photos

ONLINE EVENT
GOERS PURCHASE
**$12.9K**

TWITTER
USERS SHARE
**347.2K** tweets

CRYPTO
BUYERS PURCHASE
**$90.2M** in cryptocurrency

SNAPCHAT
USERS SEND
**2.43M** snaps

VENMO
USERS SEND
**$437.6K**

TINDER
USERS SWIPE
**1.1M** times

AMAZON
SHOPPERS SPEND
**$443K**

YOUTUBE
USERS UPLOAD
**500** hours of video

VIEWERS SPEND
**1M** hours
STREAMING

DOORDASH
DINERS PLACE
**$76.4K** in orders

**104.6K** hours
SPENT IN
ZOOM
MEETINGS

https://www.domo.com/data-never-sleeps#

# Data Growth



**Annual Size of the Global Datasphere**

175 ZB

Zetabytes

Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018
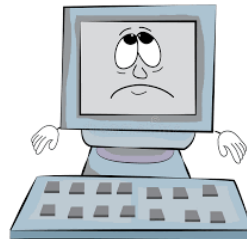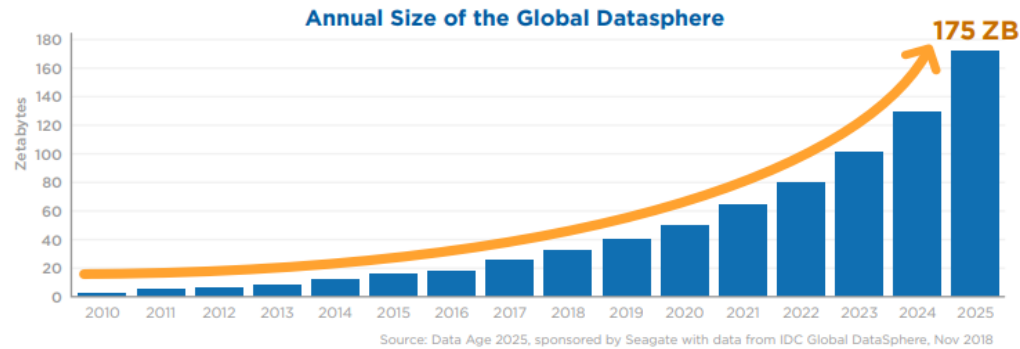
Mankind's quest to digitize the world!
33 ZB (2018) → 175 ZB (2025)
size of global datasphere*

*Source: https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf
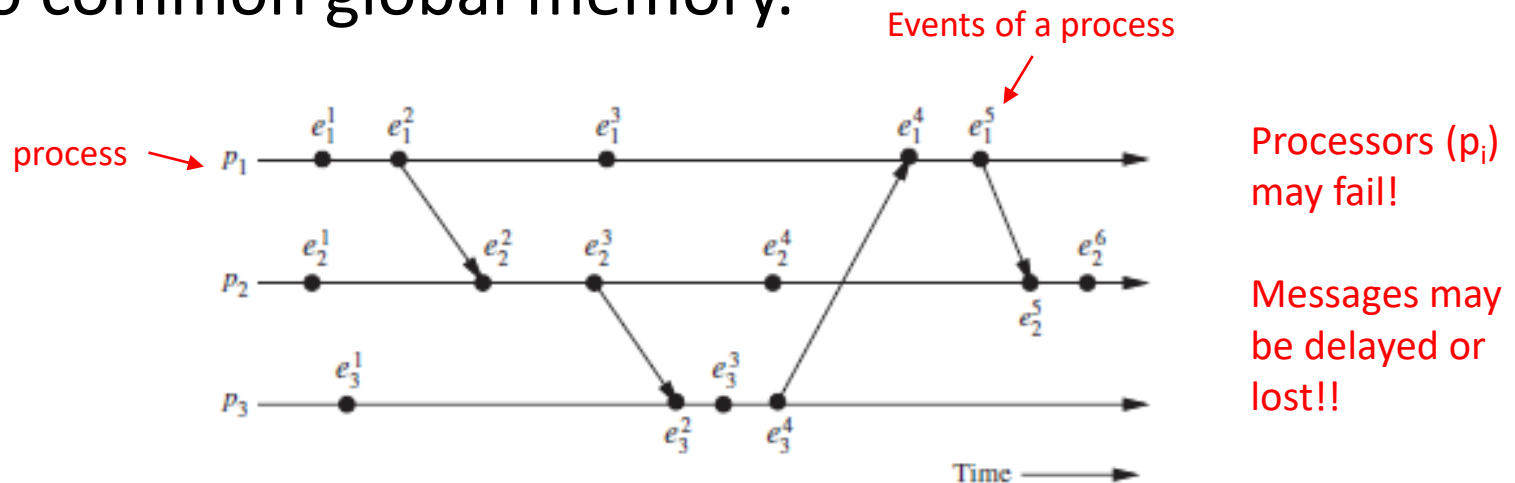
# Beyond a Single Machine



**Annual Size of the Global Datasphere**

175 ZB

Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

**Global datasphere is growing!**

How has computing evolved to capture, process and analyze these data?

# A Model of a Distributed System

- A set of processes connected by a communication network.

- Communication by information exchange.

- No physical global clock.

- No common global memory.

Events of a process

process → $P_1$

$e_1^1$  $e_1^2$  $e_1^3$  $e_1^4$  $e_1^5$

$P_2$  $e_2^1$  $e_2^2$  $e_2^3$  $e_2^4$  $e_2^5$  $e_2^6$

$P_3$  $e_3^1$  $e_3^2$  $e_3^3$  $e_3^4$

Time →

Processors ($p_i$) may fail!

Messages may be delayed or lost!!

# Course Dynamics

https://vvtesh.github.io/teaching/dcbd-saiu-2023.html