# 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

# Things to Remember

- Lectures usually start on time and end on time.

- All exams are closed book exams.

- Visit your course page for slides, readings and more.

- Discussion during lecture sessions are encouraged. Do not wait till the end if you have questions/comments.

- Exams will be in-person (not online).

- All assignment submissions will be on moodle. Assignment deadlines penalties are very strict.
    - 1 min delay: 0.5 marks
    - 1 hour delay: 1 marks
    - 12 hours delay: 2 marks
    - 24 hours delay: 3 marks
    - Beyond 24 hours, a new make-up assignment might be given.

- Please keep your video "on" if possible so that I can assign faces to names.

- I provide recommendation letters only when you have worked with me beyond the course work (some research project, internships, etc).

# 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 |

# Trivia

- How many websites are there in the world wide web?

- If each website has 100 pages on average, how many pages would we have?

- If each page consumes 1 MB (on average), how much space is required to store all the pages? (Use 1000 instead of 1024 for ease of calculation)

# Trivia

- How many websites are there in the world wide web?
  - Approximately 1.1 Billion
- If each website has 100 pages on average, how many pages would we have?
  - 110 Billion
- If each page consumes 1 MB (on average), how much space is consumed?
  - $110 * 10^9 * 1 MB = 110 * 10^9 * 10^6$ bytes = 110 PB

# Activity

- Visit https://news.google.com/
- Search for "Big data"
- Paste the most interesting news title that you saw which got published in the recent one year.

# 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, ...

Uber

Challenges and Opportunities to Dramatically Reduce the Cost of Uber's Big Data

# The Impact of Big Data



**Welcome to the United Nations**
2022 UN Big Data Hackathon

26 Sept 2022

**UnivDatos Market Insights**
Unleashing the Power of Big Data: Revolutionizing Sports Analytics

24 Jun 2023

**GE Aerospace**
All Systems Go! FlightPulse Unlocks the Power of Data and Analytics for Airline Pilots
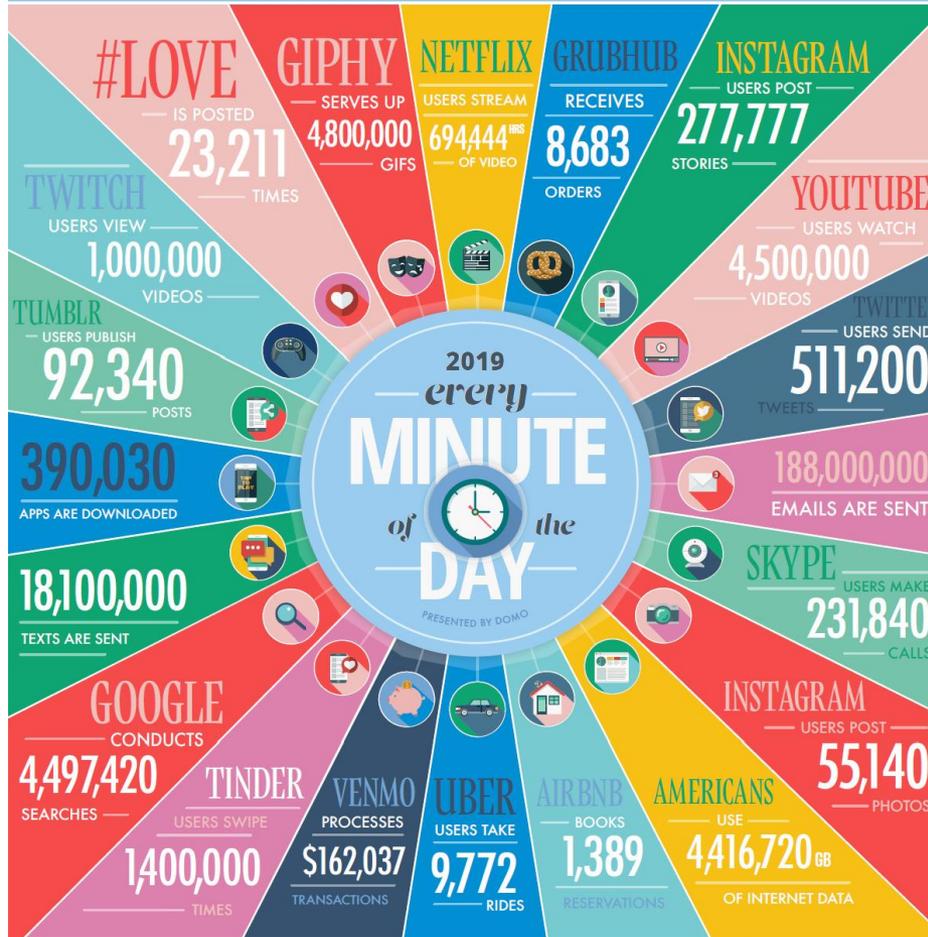
16 May 2024

# 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/

# Data Never Sleeps 11.0

Domo has been keeping tabs on the world's data usage—in a minute—for over a decade now. What the numbers consistently show is that how we use data is always evolving—and that data isn't slowing down. We're also seeing some big changes. The rise of Artificial Intelligence (AI) is reshaping the way we communicate, work, and create. Digital payments continue to replace traditional transactions. Taylor Swift streams in countless headphones. And a rash of cybercrime grows alongside these digital experiences.

In Domo's 11th edition of Data Never Sleeps, we take the pulse of our digital age, where every click, swipe, and stream fuels an ever-expanding digital universe. These are not just numbers; they are the heartbeat of a world where data reigns supreme.

**EVERY MINUTE 01:00 OF THE DAY**

PRESENTED BY DOMO

**AIRBNB** GUESTS BOOK **747** STAYS

VIEWERS WATCH **43 YEARS** OF STREAMING **CONTENT**

**AMAZON** SHOPPERS SPEND **$455K**

**X** USERS SEND **360K** TWEETS

**6.3M** SEARCHES HAPPEN ON **GOOGLE**

**WHATSAPP** USERS SEND **41.6M** MESSAGES

**LINKEDIN** USERS SUBMIT **6,060** RESUMES

**3,720** USERS DOWNLOAD **INSTAGRAM** THREADS

**CHATGPT** USERS SEND **6,944** PROMPTS

**FACEBOOK** USERS LIKE **4M** POSTS

FANS STREAM A **TAYLOR SWIFT** SONG **69.4K** TIMES

**CYBER-CRIMINALS** LAUNCH **30** DDOS ATTACKS

**INSTAGRAM** USERS SEND **694K** REELS VIA DM

**DOORDASH** DINERS PLACE **$122K** IN ORDERS

**VENMO** USERS SEND **$463K** IN PAYMENTS

GLOBAL **INTERNET** USERS SPEND **25.1M** HOURS ONLINE

THE AVERAGE **PERSON** PRODUCES **102 MB** OF DATA

PEOPLE TRADE **$398M** IN TREASURY **BONDS**

PEOPLE SEND **241M** EMAILS

**TWITCH** USERS WATCH **48K** HOURS OF CONTENT

The world's internet population continues to grow significantly year-over-year. As of November 2023, the internet represents 5.2 billion people—approximately 64.6% of the global population. According to Statista, the total amount of data predicted to be created, captured, copied, and consumed globally in 2023 is 120 zettabytes, a number projected to grow to 181 zettabytes by 2025.
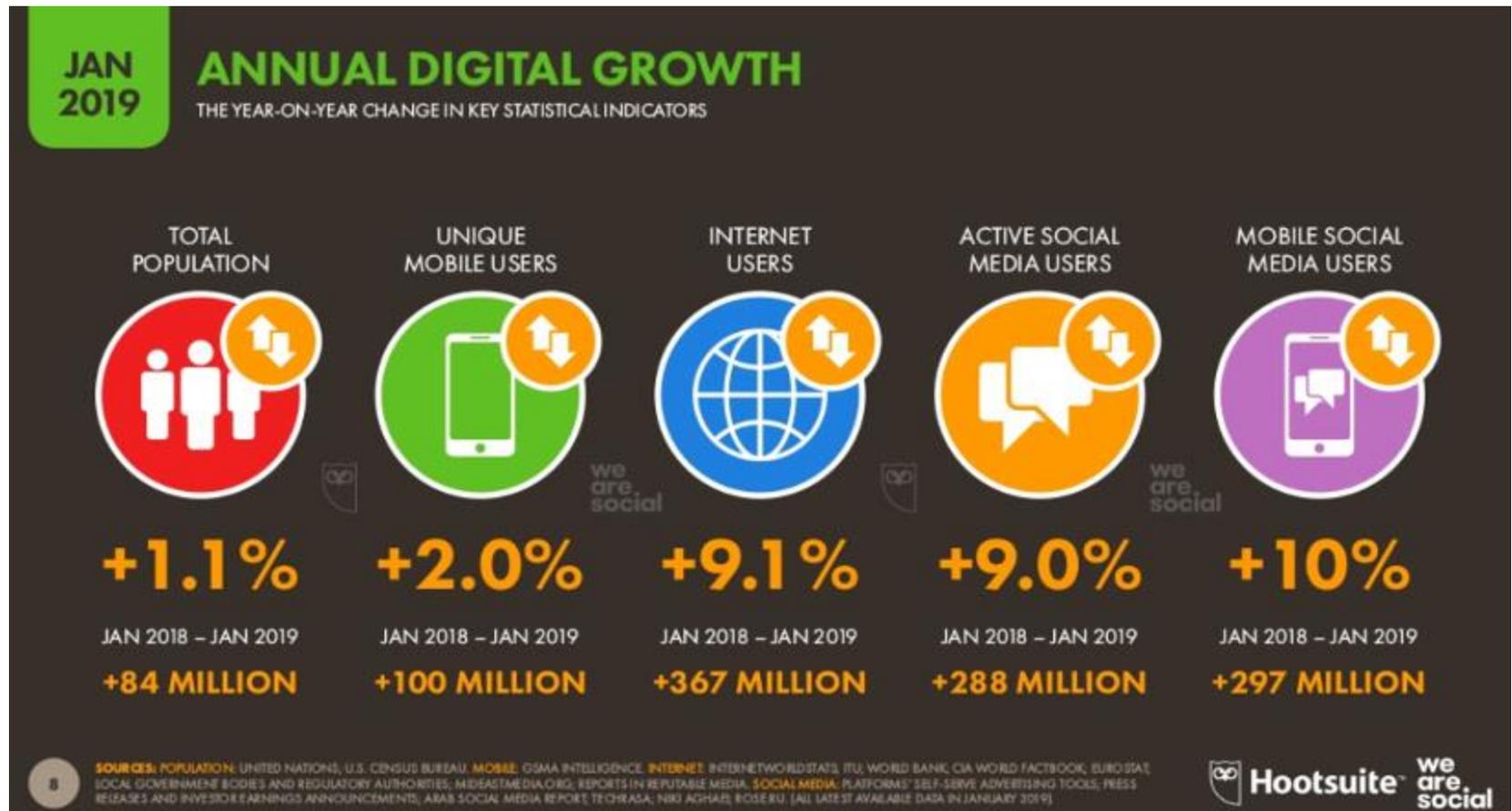
As data grows and evolves, businesses need to grow and evolve, too. Domo helps you harness the power of data so you can change as quickly as the world changes and make data-driven decisions that set you apart from the crowd. Let Domo help you make sense of all the clicks, swipes, and shares so you can see the big picture that a lot of small decisions make.

## Global Internet Population Growth
(IN BILLIONS)

| 2013 | 2016 | 2020 | 2023 |
|------|------|------|------|
| 2.1 | 3.4 | 4.5 | 5.2 |

**Learn more at domo.com**

# And, It is Growing!

# Data Growth



**Annual Size of the Global Datasphere**

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*

# 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
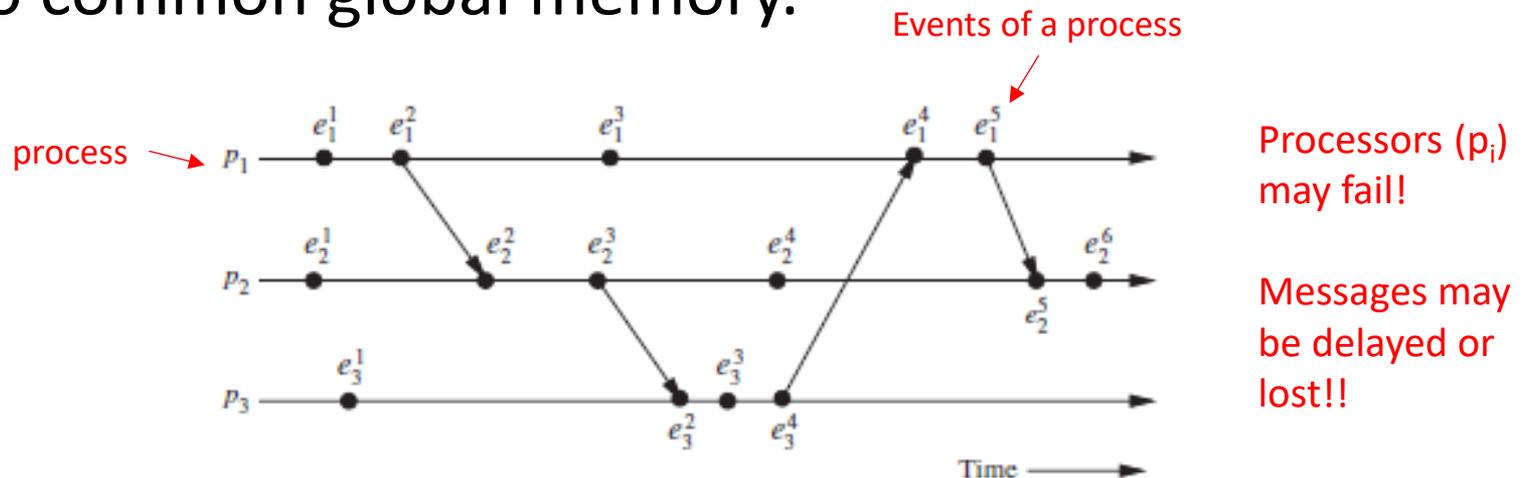
**Global datasphere is growing!**

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

Processors ($p_i$) may fail!

Messages may be delayed or lost!!

# Course Dynamics

Visit the course website. Bookmark it.