



TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

Introduction to big data management and processing

Instructor: Dr. Thanh-Chung Dao
Slides by Dr. Viet-Trung Tran
School of Information and Communication Technology

1

About this course

Tên học phần:	Lưu trữ và xử lý dữ liệu lớn (Big data storage and processing)
Mã số học phần:	IT4043E
Khối lượng:	3(3-1-0-6) – Lý thuyết: 45 tiết – BTL: 15 tiết – Thí nghiệm: 0 tiết
Slide và nhóm	– Trên Microsoft Teams

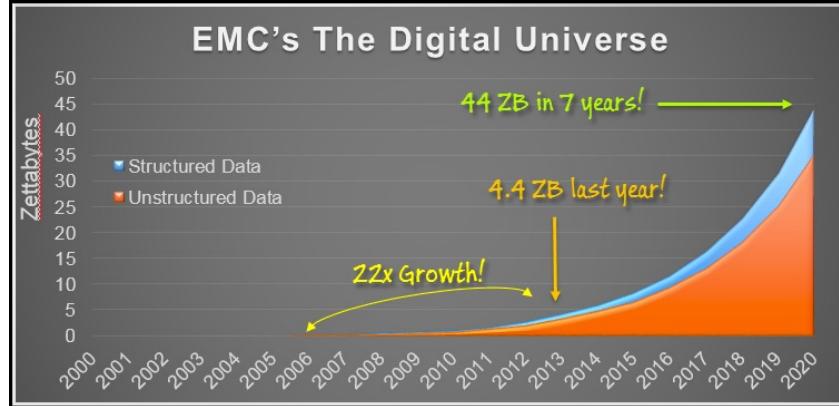
2

Syllabus

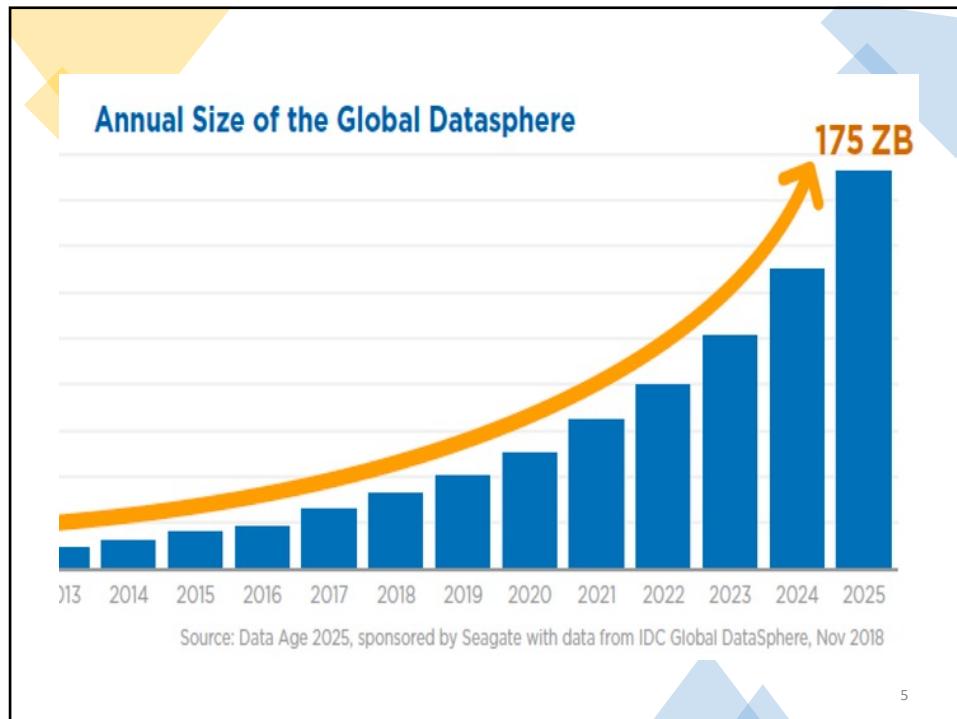
STT	Lecture
1, 2	Tổng quan về lưu trữ và xử lý dữ liệu lớn
	Hệ sinh thái Hadoop (Hadoop ecosystem)
	Hệ thống tập tin phân tán Hadoop HDFS
3, 4	Cơ sở dữ liệu phi quan hệ NoSQL - phần 1 Tổng quan
	Cơ sở dữ liệu phi quan hệ NoSQL - phần 2 Kiến trúc phân tán phổ biến
	Cơ sở dữ liệu phi quan hệ NoSQL - phần 3 Truy vấn SQL trên NoSQL, Elasticsearch
5	Hệ thống truyền thông điệp phân tán
6	Các kỹ thuật xử lý dữ liệu lớn theo khối – Hadoop Mapreduce Map Reduce
7, 8, 9	Các kỹ thuật xử lý dữ liệu lớn theo khối – Apache Spark Apache Spark
10	Các kỹ thuật xử lý luồng dữ liệu lớn Spark Streaming
12	Phân tích dữ liệu lớn và đồ thị hóa dữ liệu Spark ML và Kibana

3

How big is big data?



4



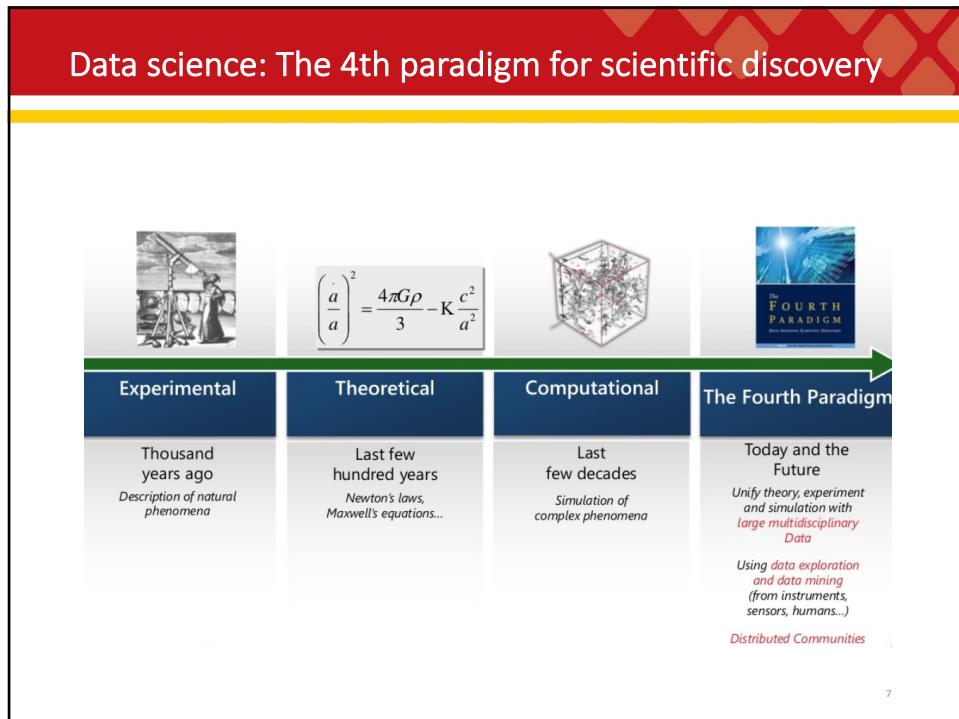
5

5



6

6



7

Big data in 2008

<http://www.wired.com/wired/issue/16-07> September 2008

The End of Science

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data. Welcome to the Petabyte Age.

8

Big data in 2014



THE AVERAGE PERSON TODAY PROCESSES MORE DATA IN A SINGLE DAY THAN A PERSON IN THE 1500'S DID IN AN ENTIRE LIFETIME ▾

LOOK TO THE LEFT, and you see Times Square at dusk. Look to the right, and you see the same location at midmorning. Internationally acclaimed photographer Stephen Wilkes's time-altering image of New York's Times Square is part of his body of work titled *Day to Night*.

The image was created by blending more than 1,400 separate photos taken over the course of 15 hours—a meticulous process that took him nearly three months.

PHOTO: STEPHEN WILKES

9

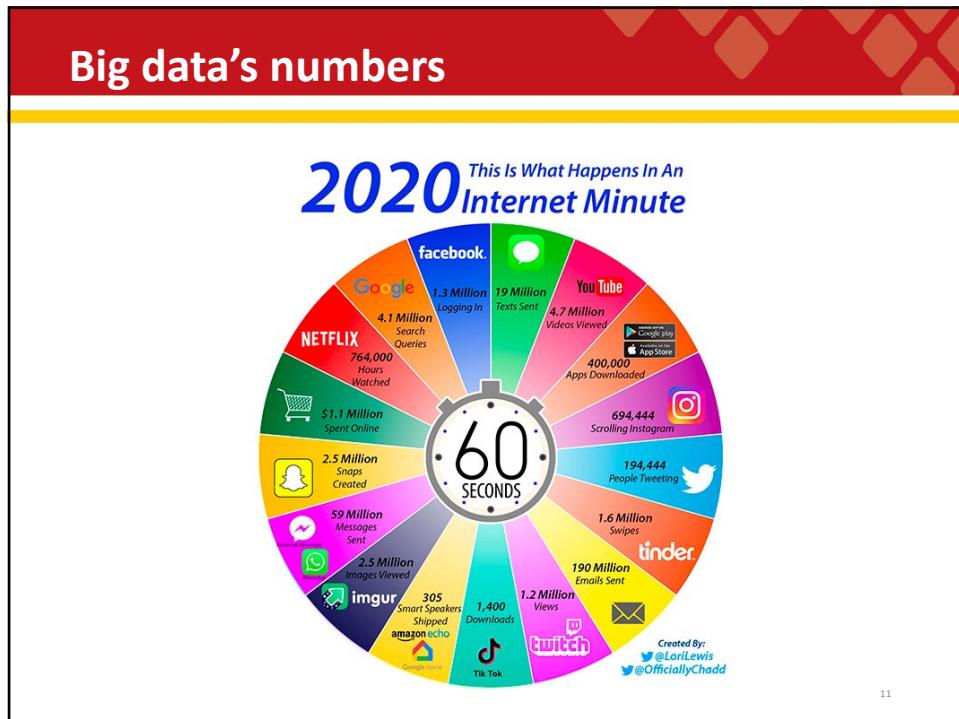
Big data today



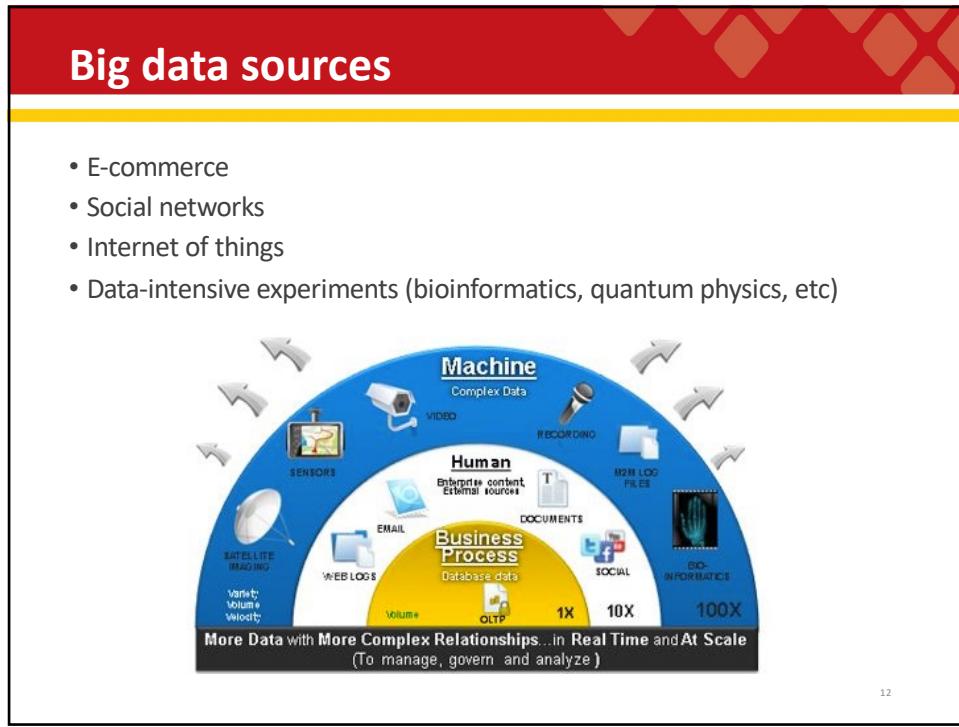
The amount of information generated during the first day of a baby's life today is equivalent to 70 times the information contained in the Library of Congress

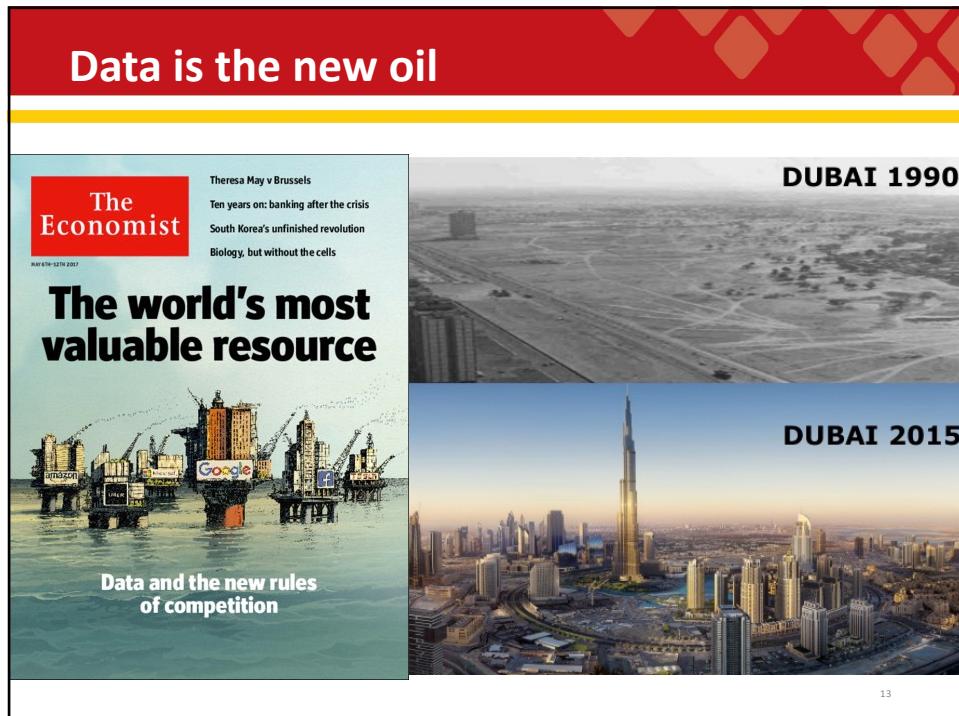
10

10

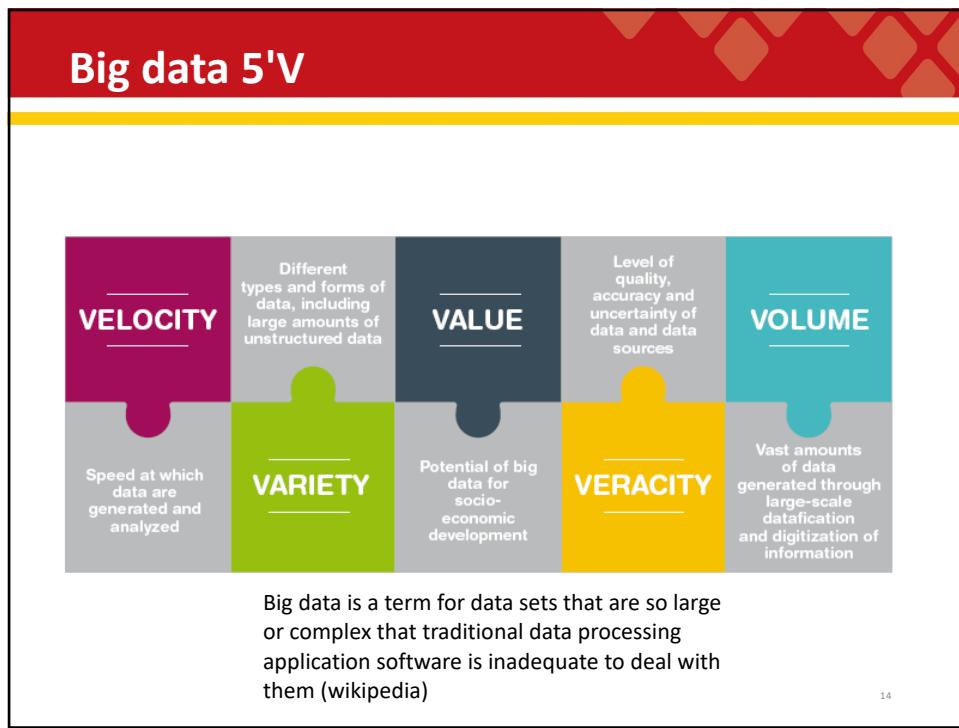


11

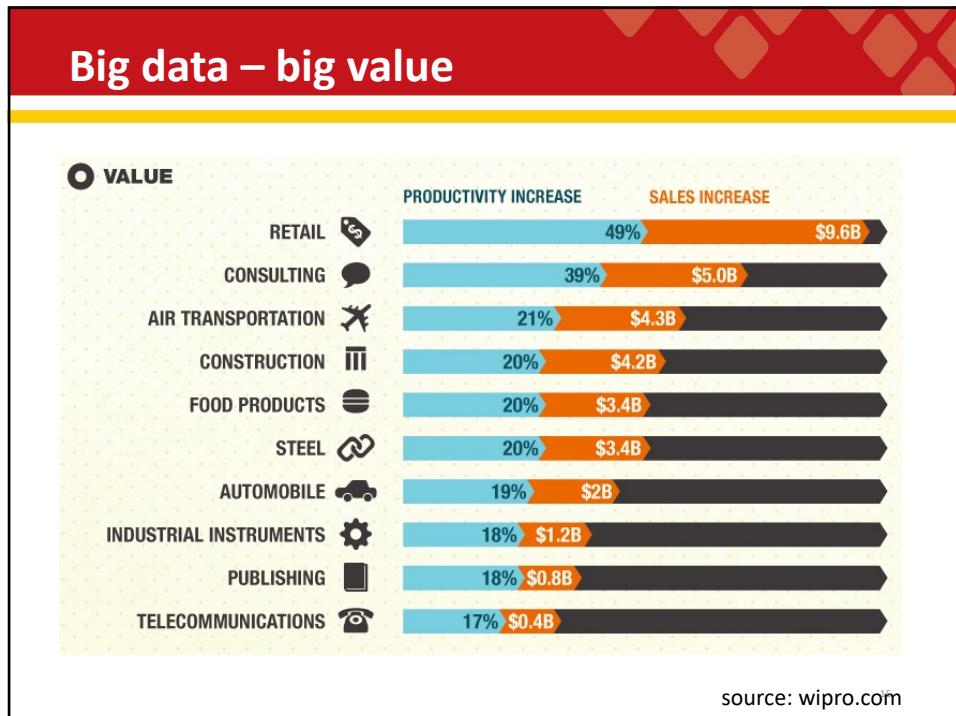




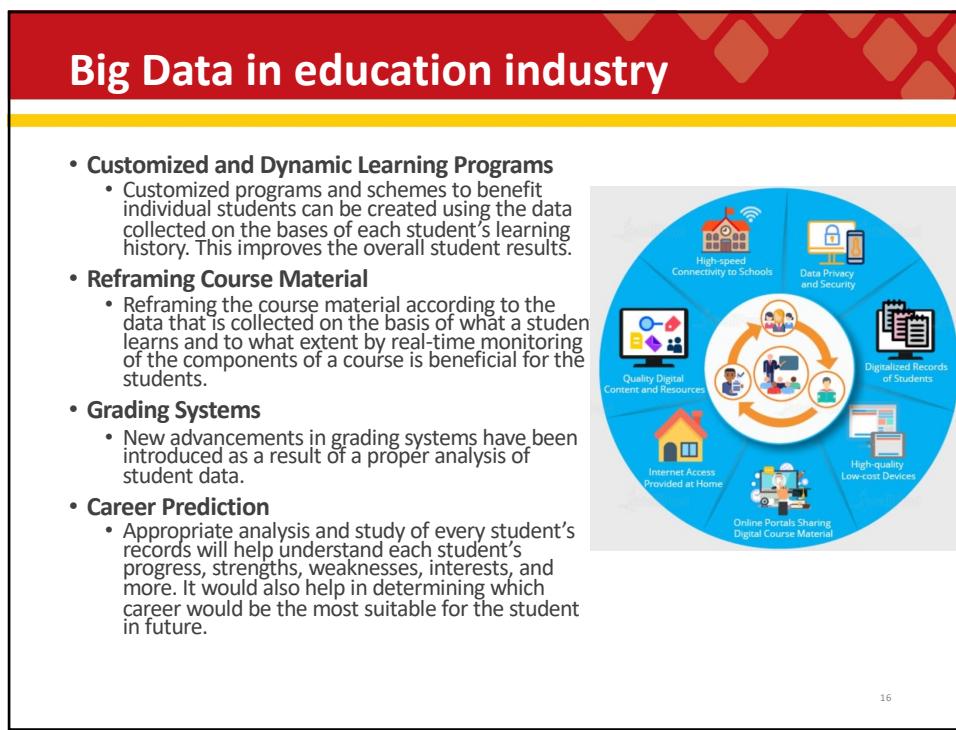
13



14



15



16

Edtech

- Coursera
- VioEdu
- <https://byjus.com/>
 - Engaging Video Lessons
 - Personalized Learning Journeys
 - Mapped to the Syllabus
 - In-depth Analysis
 - Engaging Interactive Questions

17

Big Data in healthcare industry

- Big data reduces costs of treatment since there is less chances of having to perform unnecessary diagnosis.
- It helps in predicting outbreaks of epidemics and also in deciding what preventive measures could be taken to minimize the effects of the same.
- It helps avoid preventable diseases by detecting them in early stages. It prevents them from getting any worse which in turn makes their treatment easy and effective.
- Patients can be provided with evidence-based medicine which is identified and prescribed after doing research on past medical results.

18

Big Data in government sector

• Welfare Schemes

- In making faster and informed decisions regarding various political programs
- To identify areas that are in immediate need of attention
- To stay up to date in the field of agriculture by keeping track of all existing land and livestock.
- To overcome national challenges such as unemployment, terrorism, energy resources exploration, and much more.

• Cyber Security

- Big Data is hugely used for deceit recognition.
- It is also used in catching tax evaders.



19

19

Big Data in media and entertainment industry

- Predicting the interests of audiences
- Optimized or on-demand scheduling of media streams in digital media distribution platforms
- Getting insights from customer reviews
- Effective targeting of the advertisements
- Example
 - Spotify, an on-demand music providing platform, uses Big Data Analytics, collects data from all its users around the globe, and then uses the analyzed data to give informed music recommendations and suggestions to every individual user.
 - Amazon Prime that offers, videos, music, and Kindle books in a one-stop shop is also big on using big data.



20

20



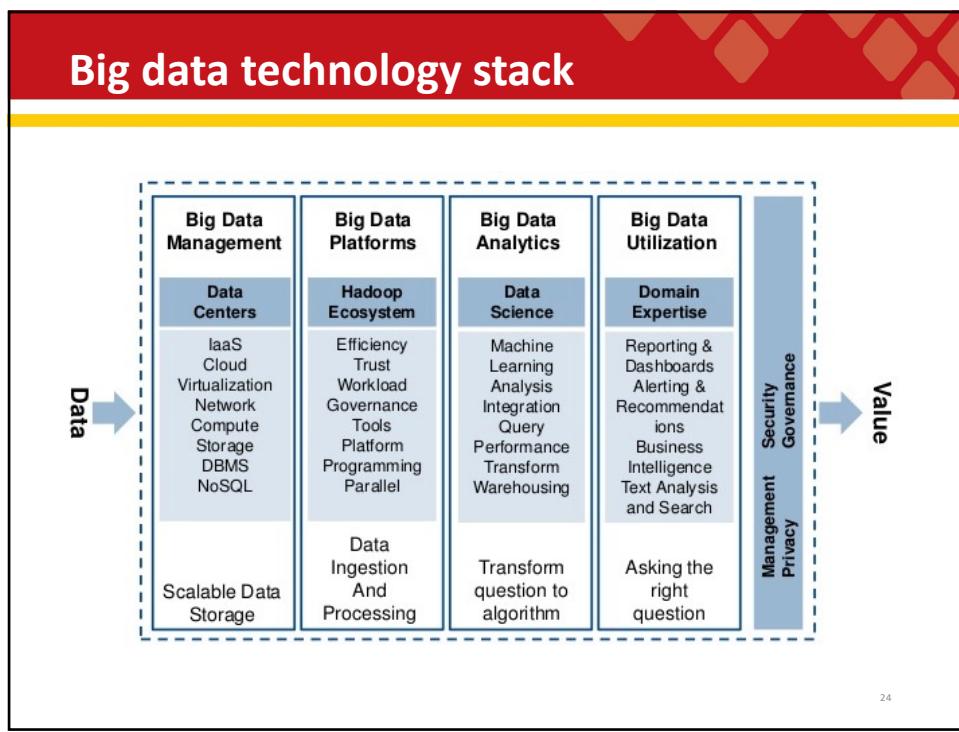
21



22



23



24

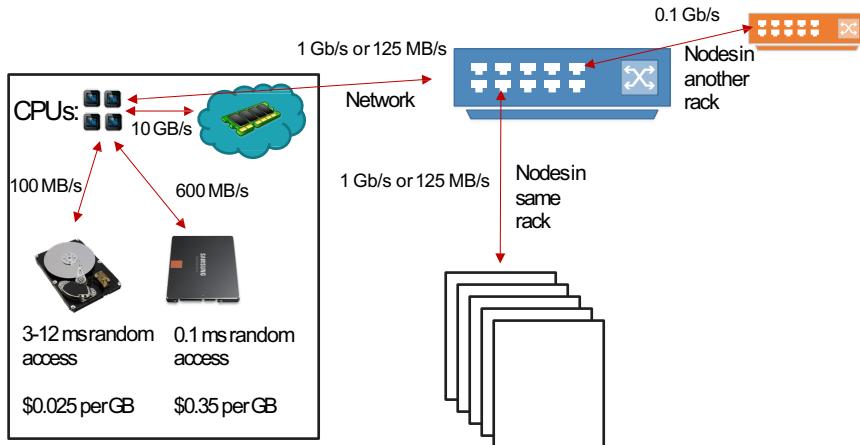
Scalable data management

- Scalability
 - Able to manage increasingly big volume of data
- Accessibility
 - Able to maintain efficiency in reading and writing data (I/O) into data storage systems
- Transparency
 - In distributed environment, users should be able to access data over the network as easily as if the data were stored locally.
 - Users should not have to know the physical location of data to access it.
- Availability
 - Fault tolerance
 - The number of users, system failures, or other consequences of distribution shouldn't compromise the availability.

25

25

Data I/O landscape



26

26

Scalable data ingestion and processing

- Data ingestion
 - Data from different complementing information systems is to be combined to gain a more comprehensive basis to satisfy the need
 - How to ingest data efficiently from various, distributed heterogeneous sources?
 - Different data formats
 - Different data models and schemas
 - Security and privacy
- Data processing
 - How to process massive volume of data in a timely fashion?
 - How to process massive stream of data in a real-time fashion?
 - Traditional parallel, distributed processing (OpenMP, MPI)
 - Big learning curve
 - Scalability is limited
 - Fault tolerance is hard to achieve
 - Expensive, high performance computing infrastructure
 - Novel realtime processing architecture
 - Eg. Mini-batch in Spark streaming
 - Eg. Complex event processing in Apache Flink

27

27

Scalable analytic algorithms

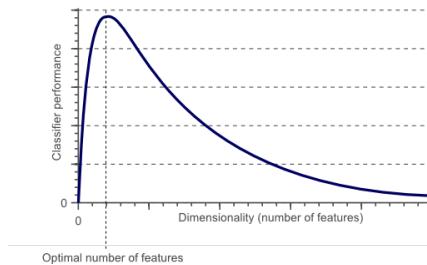
- Challenges
 - Big volume
 - Big dimensionality
 - Realtime processing
- Scaling-up Machine Learning algorithms
 - Adapting the algorithm to handle Big Data in a single machine.
 - Eg. Sub-sampling
 - Eg. Principal component analysis
 - Eg. feature extraction and feature selection
 - Scaling-up algorithms by parallelism
 - Eg. k-nn classification based on MapReduce
 - Eg. scaling-up support vector machines (SVM) by a divide and-conquer approach

28

28

Eg. Curse of dimensionality

- The required number of samples (to achieve the same accuracy) grows exponentially with the number of variables!
- In practice: number of training examples is fixed!
=> the classifier's performance usually will degrade for a large number of features!



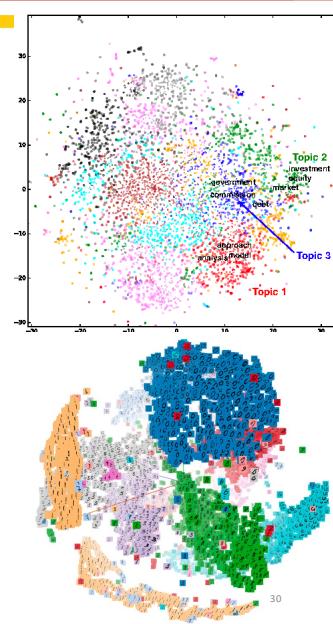
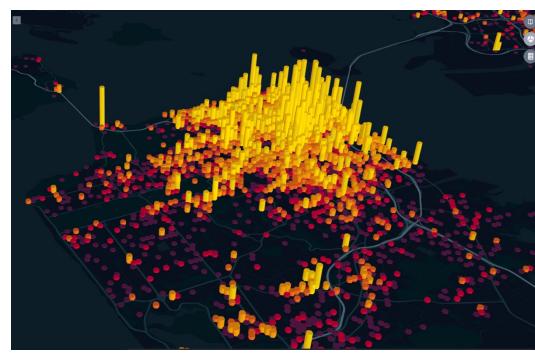
In fact, after a certain point, increasing the dimensionality of the problem by adding new features would actually degrade the performance of classifier.

29

29

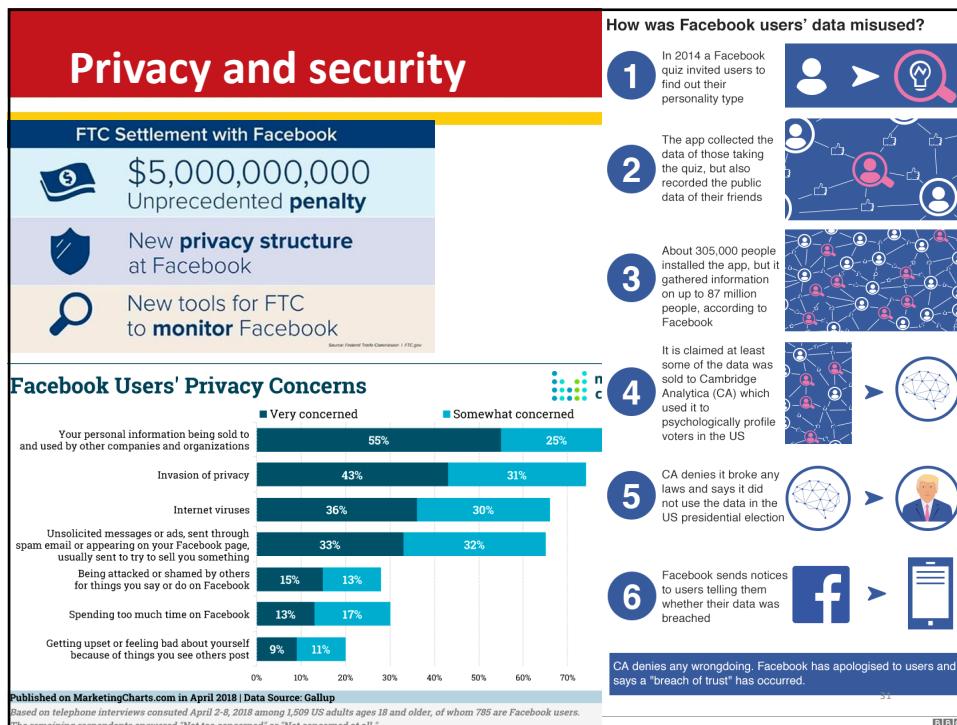
Utilization and interpretability of big data

- Domain expertise to findout problems and interpret analytics results
- Scalable visualization and interpretability of million data points
 - to facilitate their interpretability and understanding

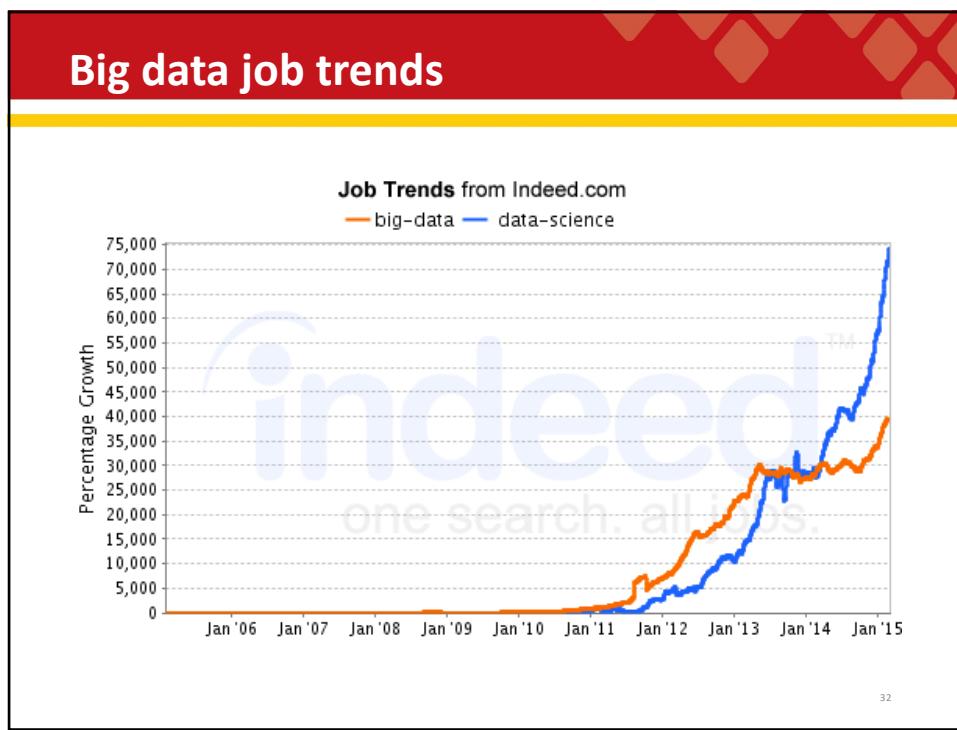


30

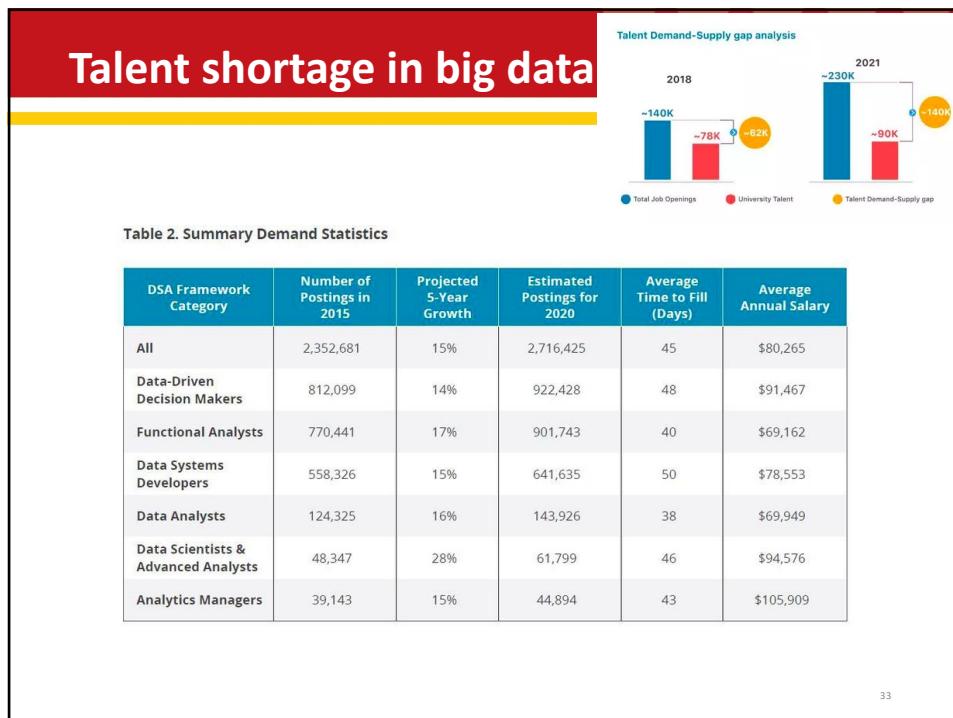
15



31

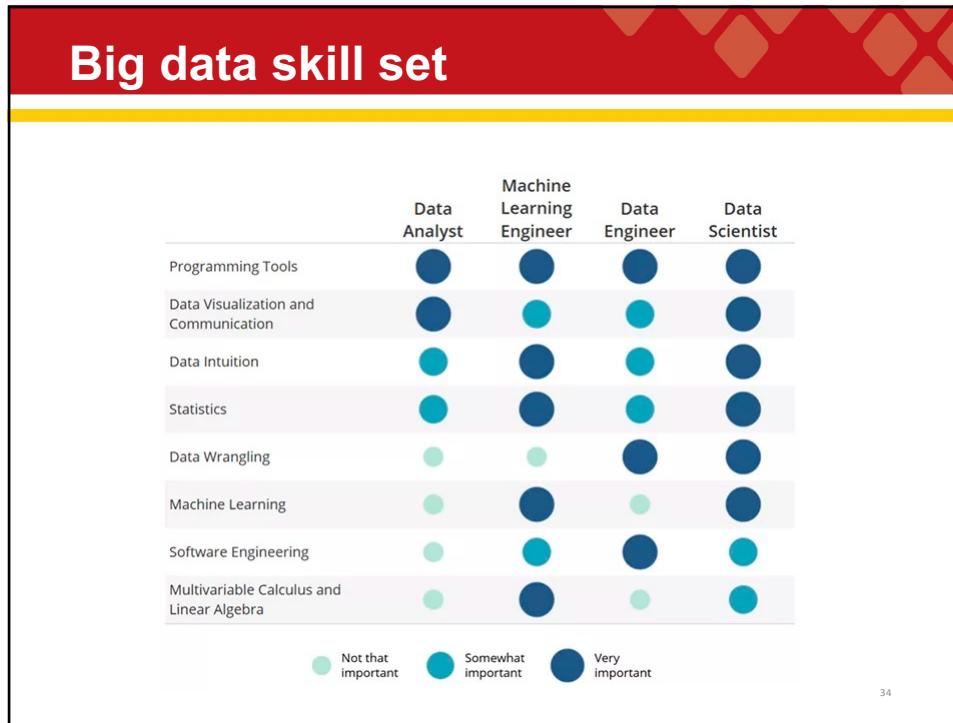


32



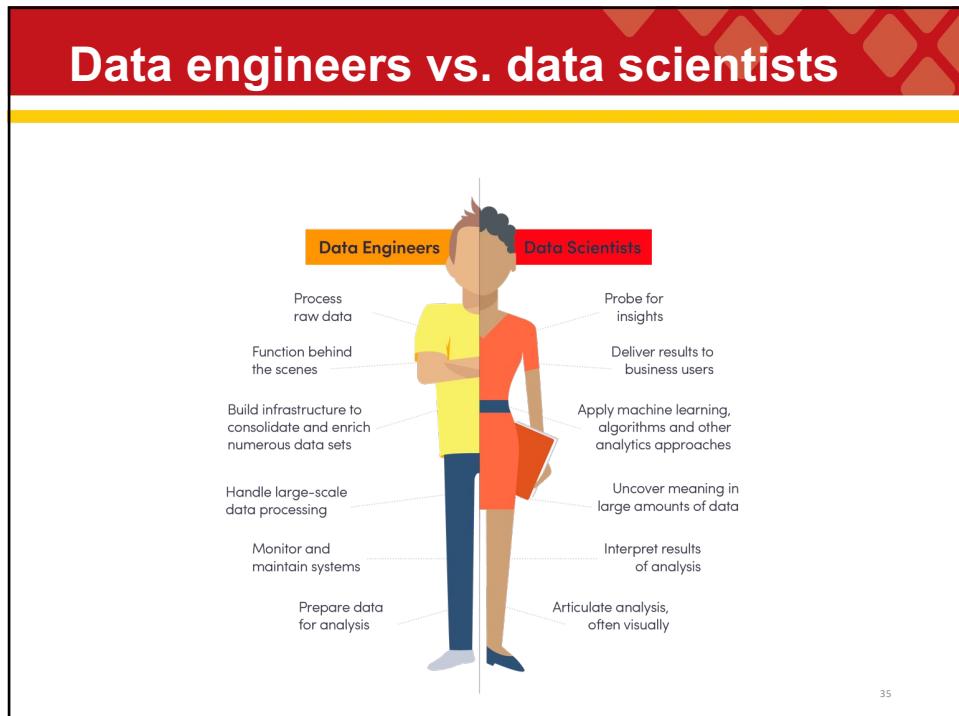
33

33



34

34



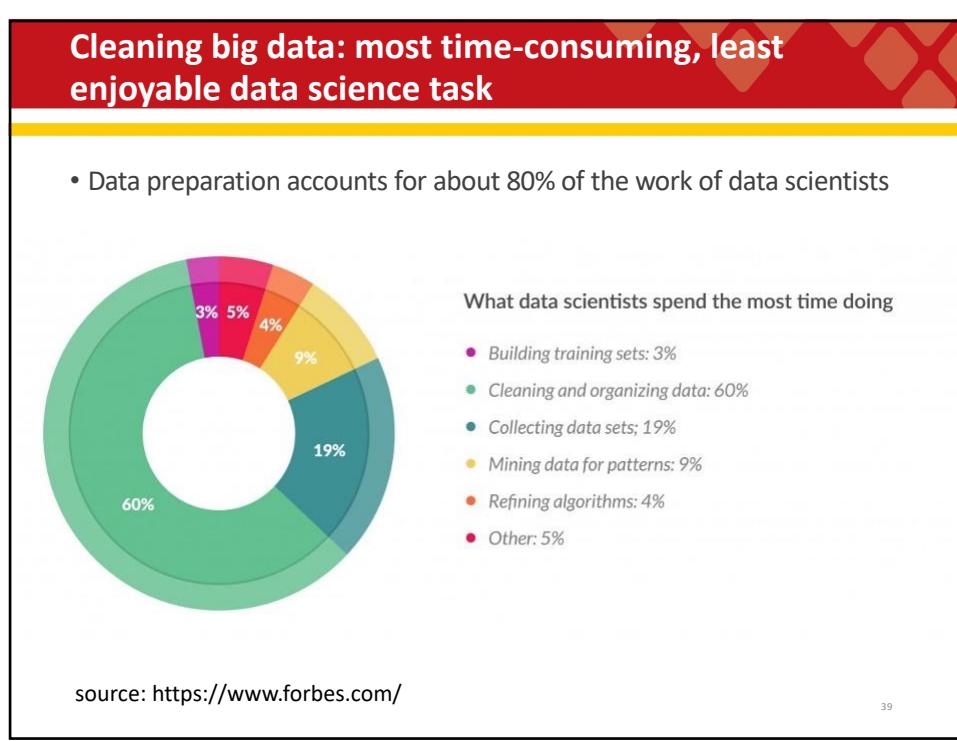
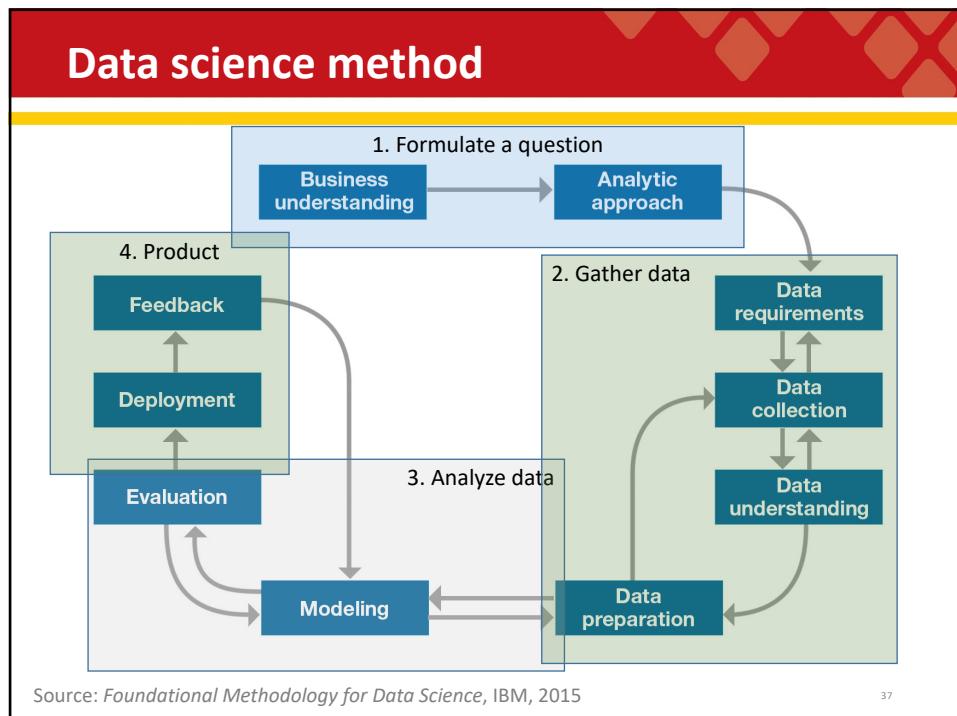
35

How to land big data related jobs

- Learn to code
 - Coursera
 - Udacity
 - Freecodecamp
 - Codecademy
- Math, Stats and machine learning
 - Kaggle
- Hadoop, NoSQL, Spark
- Visualization and Reporting
 - Tableau
 - Pentaho
- Meetup & Share
- Find a mentor
- Internships, projects

36

36



Cleaning big data: most time-consuming, least enjoyable data science task

- 57% of data scientists regard cleaning and organizing data as the least enjoyable part of their work and 19% say this about collecting data sets.

Part of Data Science	Percentage
Cleaning and organizing data	57%
Collecting data sets	21%
Building training sets	10%
Mining data for patterns	3%
Refining algorithms	4%
Other	5%

What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

40

40

References

- [1] Tiwari, Shashank. Professional NoSQL. John Wiley & Sons, 2011.
- [2] Lam, Chuck. Hadoop in action. Manning Publications Co., 2010.
- [3] Miner, Donald, and Adam Shook. MapReduce design patterns: building effective algorithms and analytics for Hadoop and other systems. "O'Reilly Media, Inc.", 2012.
- [4] Karau, Holden. Fast Data Processing with Spark. Packt Publishing Ltd, 2013.
- [5] Penchikala, Sri. Big data processing with apache spark. Lulu. com, 2018.
- [6] White, Tom. Hadoop: The definitive guide. "O'Reilly Media, Inc.", 2012.
- [7] Gandomi, Amir, and Murtaza Haider. "Beyond the hype: Big data concepts, methods, and analytics." International Journal of Information Management 35.2 (2015): 137-144.
- [8] Cattell, Rick. "Scalable SQL and NoSQL data stores." Acm Sigmod Record 39.4 (2011): 12-27.
- [9] Gessert, Felix, et al. "NoSQL database systems: a survey and decision guidance." Computer Science-Research and Development 32.3-4 (2017): 353-365.
- [10] George, Lars. HBase: the definitive guide: random access to your planet-size data. "O'Reilly Media, Inc.", 2011.
- [11] Sivasubramanian, Swaminathan. "Amazon dynamoDB: a seamlessly scalable non-relational database service." Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data. ACM, 2012.
- [12] Chan, L. "Presto: Interacting with petabytes of data at Facebook." (2013).
- [13] Garg, Nishant. Apache Kafka. Packt Publishing Ltd, 2013.
- [14] Karau, Holden, et al. Learning spark: lightning-fast big data analysis. "O'Reilly Media, Inc.", 2015.
- [15] Iqbal, Muhammad Hussain, and Tariq Rahim Soomro. "Big data analysis: Apache storm perspective." International journal of computer trends and technology 19.1 (2015): 9-14.
- [16] Toshniwal, Ankit, et al. "Storm@ twitter." Proceedings of the 2014 ACM SIGMOD international conference on Management of data. ACM, 2014.
- [17] Lin, Jimmy. "The lambda and the kappa." IEEE Internet Computing 21.5 (2017): 60-66.

41

41

Online courses

- <https://www.coursera.org/learn/nosql-database-systems>
- <https://who.rocq.inria.fr/Vassilis.Christophides/Big/index.htm>
- <https://www.coursera.org/learn/big-data-introduction?specialization=big-data>
- <https://www.coursera.org/learn/big-data-integration-processing?specialization=big-data>
- <https://www.coursera.org/learn/big-data-management?specialization=big-data>
- <https://www.coursera.org/learn/hadoop>
- <https://www.coursera.org/learn/scala-spark-big-data>

42

42

TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

Thank you for your attention!
Q&A

43

43

21