



中国地质大学
CHINA UNIVERSITY OF GEOSCIENCES

“国际学生学术能力提升工程”系列讲座

丝路博士论坛

Silk Road Doctor Forum

Topic: Big Data in practice: an innovative interdisciplinary research

By: COULIBALY L. KARIM

Place: Silk Road Institute (212), Nanwangshan Campus, CUG
(南望山校区丝绸之路学院212)

Date: (2022年07月07日)

School of Geography and Information Engineering



ABOUT ME



- COULIBALY L. Karim
- 康黎明
- PhD in Surveying and mapping science and technology



- CUG, Ph.D. 2022



Coupling linear spectral unmixing and RUSLE2 to model soil erosion in the Boubo coastal watershed, Côte d'Ivoire

Lenikpoho Karim Coulibaly^a, Qingfeng Guan^{a,b}, Tchिमou Vincent Assoma^c, Xin Fan^d, Naga Coulibaly^e

^a School of Geography and Information Engineering, China University of Geosciences, Wuhan, Hubei Province, China
^b National Engineering Research Center of GIS, China University of Geosciences, Wuhan, Hubei Province, China
^c University Center for Research and Application in Remote Sensing (CURAT), University of Felix Houphouët Boigny, Abidjan, Côte d'Ivoire
^d School of Public Administration, China University of Geosciences, Wuhan, Hubei Province, China
^e Laboratory of Geoscience and Environment, UFR of Science and Environmental Management, Nanang Abrogoua University, Abidjan, Côte d'Ivoire

ARTICLE INFO

Keywords:
 Soil erosion
 RUSLE2
 Linear spectral unmixing
 climate change
 GIS
 Remote sensing

ABSTRACT

Water erosion accelerates soil degradation through land use, land cover, and climate change. Accurate modeling of soil erosion is critical for assessment of environmental variables such as nutrient loss, reduction of soil fertility, and water quality degradation. Modeling of soil erosion can provide insights to conservation planners for formulating policies to prevent land degradation. However, when used for soil erosion modeling in Geographical Information Systems (GIS), application of the Revised Universal Soil Loss Eq. (2) (RUSLE2) is realized based on the assumption that the pixels of land use data are pure and that mixed land use units within pixels can be ignored, and this opposes the accurate estimation of regional soil erosion. The methodology developed in this study includes combination of the GIS-based RUSLE2 with linear spectral unmixing (LSU) to analyze the change in vegetation cover within a pixel and to improve the spatial representation of the soil erodibility factor using climate data derived from the Boubo coastal watershed. The findings reveal that the estimated monthly erosivity density in the Boubo coastal watershed for different months varies between 0.05 and 20.86 MJ mm ha⁻¹h⁻¹ year⁻¹ in 1990 and 0.8 to 21.21 MJ mm ha⁻¹h⁻¹ year⁻¹ in 2019. The geographical soil erodibility K-factor varied between 0.008 and 0.022 t.ha.ha⁻¹MJ⁻¹.mm⁻¹. The temporal soil erodibility K_f factor was highest in May 1990 (0.194) and June 2019 (0.2). Slopes varied between 0% and 56%, with LS values exceeding 16. The deforestation rate in the Boubo coastal watershed was 65.49% from 1992 to 2019. The mean soil loss rate in June was 0.048 t/ha/month in 1990 and was 0.073 t/ha/month in 2019. Sediment yield increased from 1.09 t/ha/yr in 1990 to 2.54 t/ha/yr in 2019. Based on the RUSLE2 empirical equation, it was inferred that the estimated sediment transport capacity increased during the baseline period. Further studies should be conducted to evaluate ecosystem management based on ecosystem services and sediment deposition in this area.



CONTENTS

- 01 Introduction

- 02 Concept of Big Data

- 03 Case studies

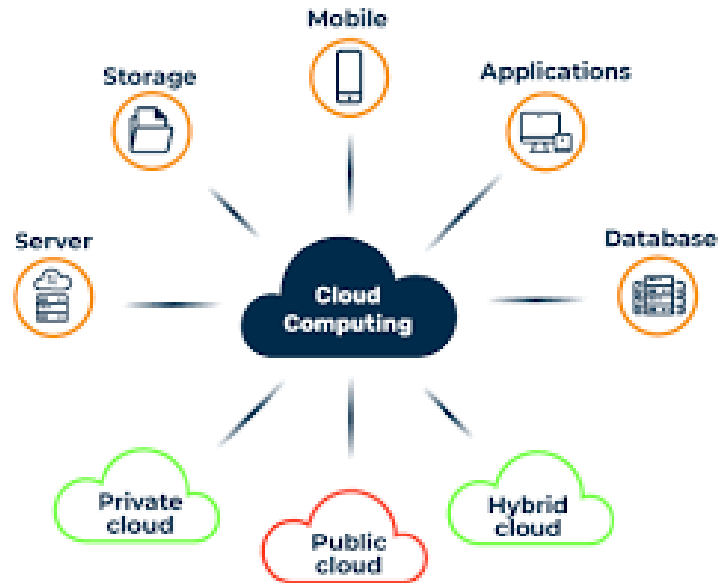
- 04 Conclusion

01

Introduction

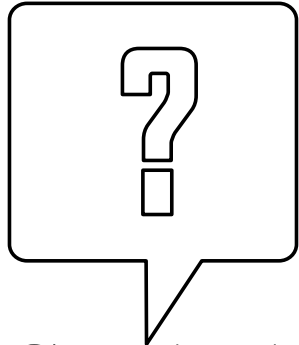


- **Background (1/3)**



Companies and government organizations are working on Big Data projects, and we thought it would be a good idea to share how Big Data is being used today to deliver real value across many industries, by both large and small businesses.

- **Background (2/3)**



BIG DATA

BIG DATA is an interdisciplinary topic. Several researchers struggle to find an innovative topic. BIG DATA can be applied in a variety of topics, and researchers can explore this technology to publish innovative papers.





1.1 Background of the topic



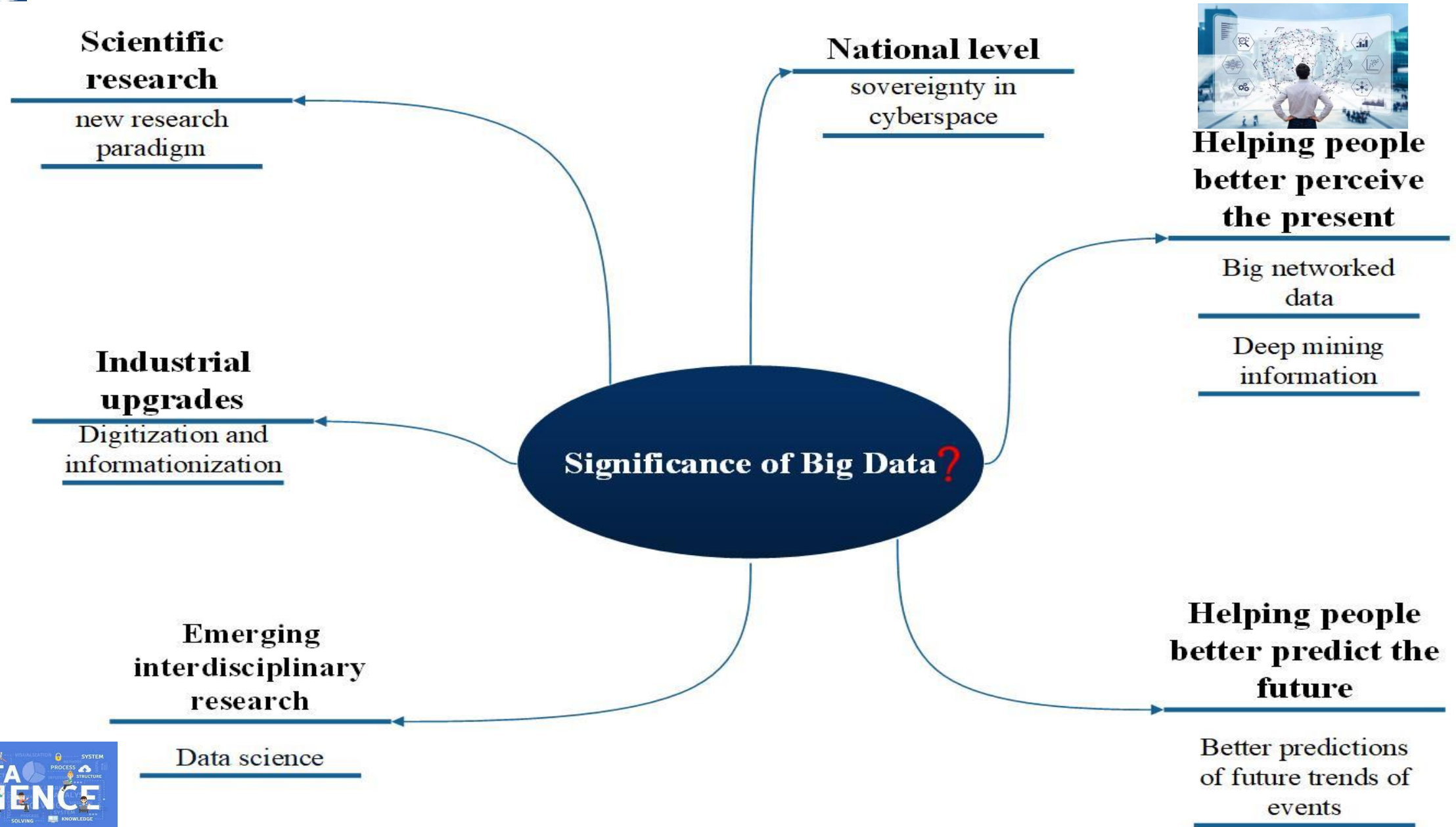
- **Background (3/3)**

Big data has become an essential tool for today's top organizations in their ongoing efforts to better understand their customers. The virtual tomes of big data coursing through any given company's system offer vast and unfathomable opportunities for a significant increase in profits, improved customer satisfaction, and edging out competition.

This presentation aim to provide insights to International students of China University of Geosciences (CUG) for understanding Big Data concept and the Practices of Big Data in different fields.

At the end of this presentation CUG students researchers should be able to:

- Define a Big data;
- Understand Big data in practice;
- Connect Big data analytics to their field and develop standard quality to get a job in china.





1.4 Why this topic?



Most Silk Road presentation have focused on how to write a paper while also find an innovative topic could be challenging.

Besides this topic concern all CUG'ers and can be applied to all field. Many students and researchers would like to work in china, we explain how this innovative topic could give them quick access to publications and we also provides some suggestions for the integration to Chinese job market.

Where do you want to work?

What do you want to do in your research?

How can you use the Big Data in your field?



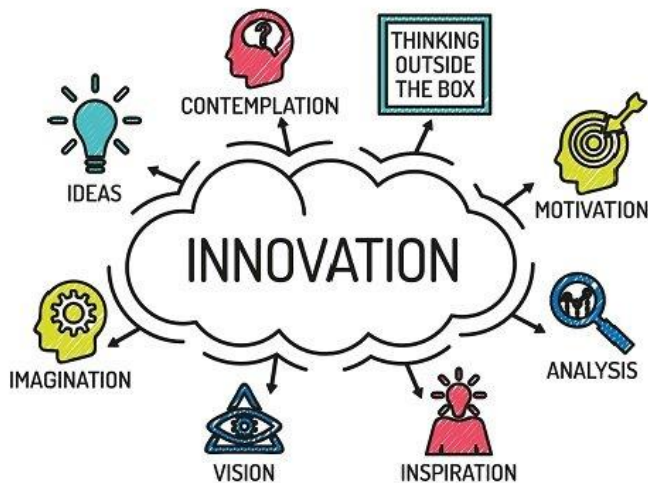
How your research will help the scientific community and/or industries?

Are you using data in your field? Are the data meeting the requirements to conduct a Big data research ?



Are you the right researcher to develop a Big Data for an Industry or a government?

What is the next project after your study?



02

BIG DATA CONCEPT



There is no universally accepted definition of "Big Data." (Chebbi et al., 2015).

Big data is a combination of structured, unstructured data, and semi-structured that may be mined for information and used in machine learning projects, predictive modeling, and other advanced analytics applications.



Types of Big Data

Structured



UnStructured



Semi-Structured



Big data is a new idea, and it has got numerous definitions from researchers, organizations, and individuals.

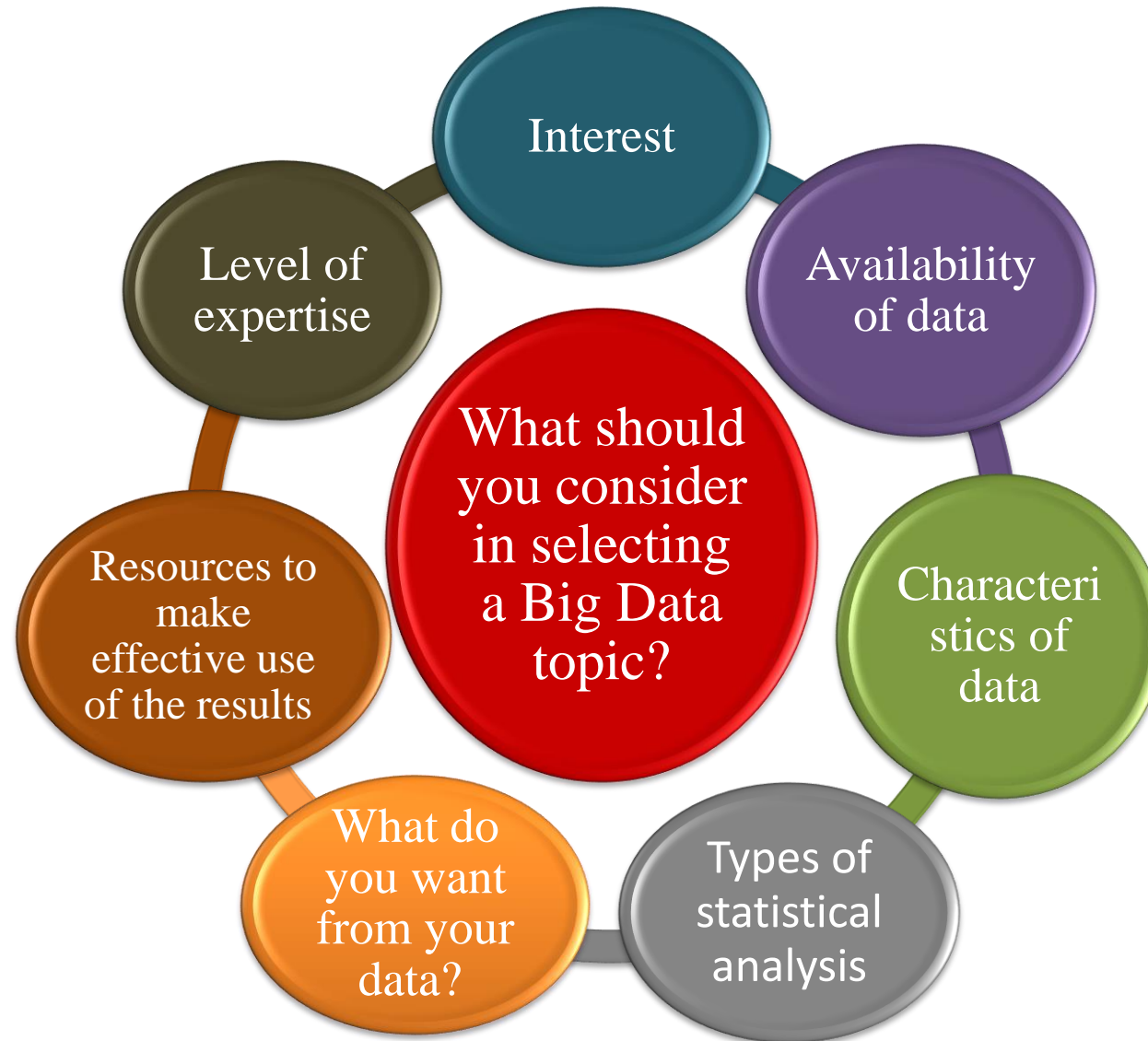
Some authors reviewed the development of the V's of Big Data (Arockia Panimalar.S et al., 2017)

Table 1: Characteristics Of Big Data

S No	Big Data Characteristics	Elucidation	Description
1	volume	Size of data	Quantity of collected and stored data. Data size is in TB
2	velocity	Speed of data	The transfer rate of data between source and destination
3	value	Importance of data	It represents the business value to be derived from big data
4	variety	Type of data	Different type of data like pictures, videos and audio arrives at the receiving end
5	Veracity	Veracity of data	Accurate analysis of captured data is virtually worthless if it's not accurate

There is no ideal format and structure of Big data for every use case.

Expert knowledge help to design and implement Big Data



1. STRONG TECHNICAL EXPERTISE

2. ABILITY TO ANALYZE SPECIFIC BUSINESS REQUIREMENTS

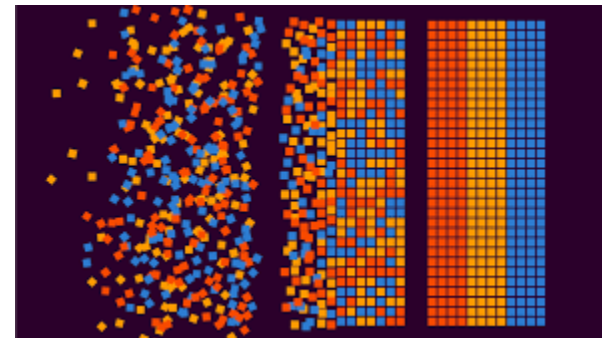
3. BIG DATA CLOUD SOLUTIONS

4. MACHINE LEARNING AND DATA MINING

5. PROBLEM-SOLVING SKILLS

6. INNOVATIVENESS

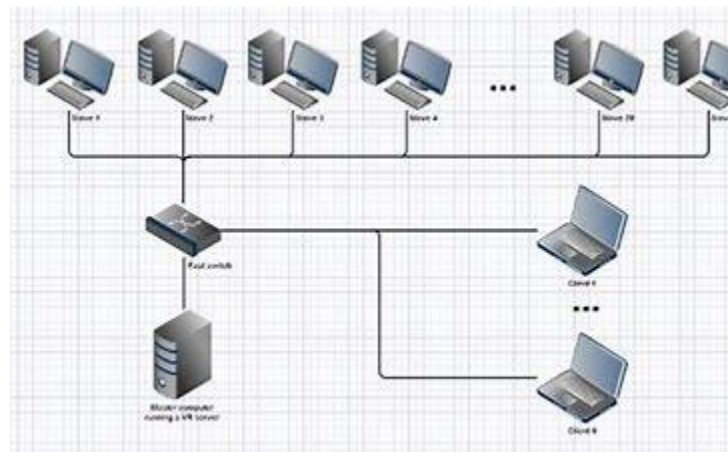
7. STATISTICAL AND QUANTITATIVE ANALYSIS



Big Data is useless unless we can turn it into insights. To accomplish this, we must collect and analyze data.

The quantity of data that could previously be saved in databases was constrained; the more data there was, the slower the system got. New methods that **enable us to store** and **analyze data** across several databases, in dispersed locations, connected by networks, can now overcome this.

Distributed computing means huge amounts of data can be stored (in little bits across lots of databases) and analysed by sharing the analysis between different servers (each performing a small part of the analysis).



Big Data can improve products performance and deliver strategic value in cities, telecoms, sports, gambling, fashion, manufacturing, research, motor racing, video gaming, zoo, and everything in between.

Big Data, more than any other current trend, will have an impact on everyone and everything we do.





Challenges from the design of processing systems to analysis methods, as well as a number of open problems in scientific research.

Some of these challenges are caused by big data's characteristics, others by its current analysis models and methods, and still others by the limitations of current data processing systems.

Big data can be applied in your field

Your Big Data expertise start today



03

BIG DATA IN PRACTICE



High-performance Spatial Computational Intelligence Lab @ CUG



3.1 Geospatial Applications

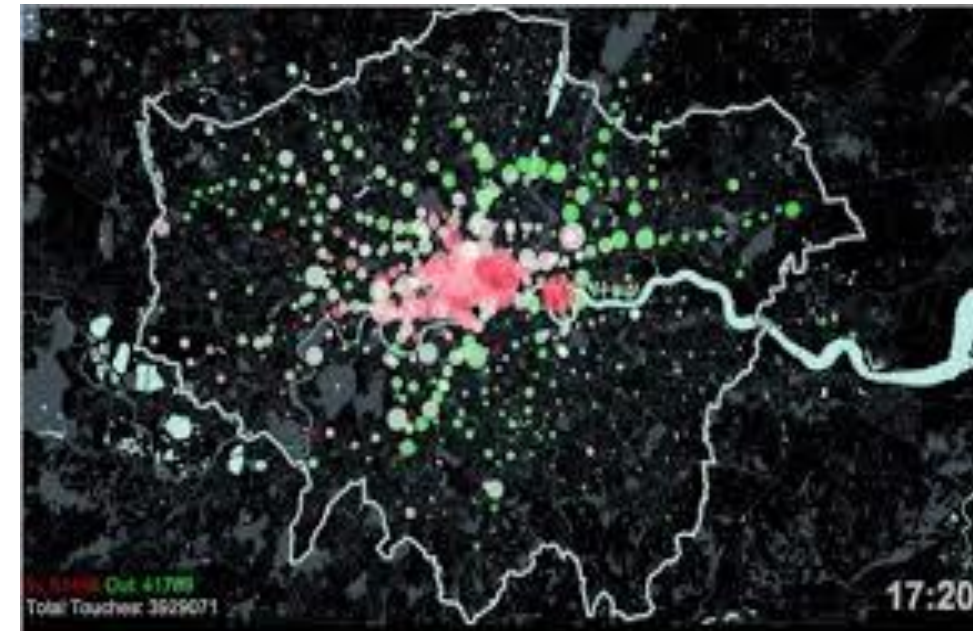


3.1.1 What problem is big data helping to solve?

Big Data Computing for Geospatial Applications aims to capture the latest efforts on **utilizing, adapting, and developing new computing approaches, spatial methods, and data management strategies** to tackle geospatial big data challenges for supporting applications in different domains, such as climate change, disaster management, human dynamics, public health, and environment and engineering.

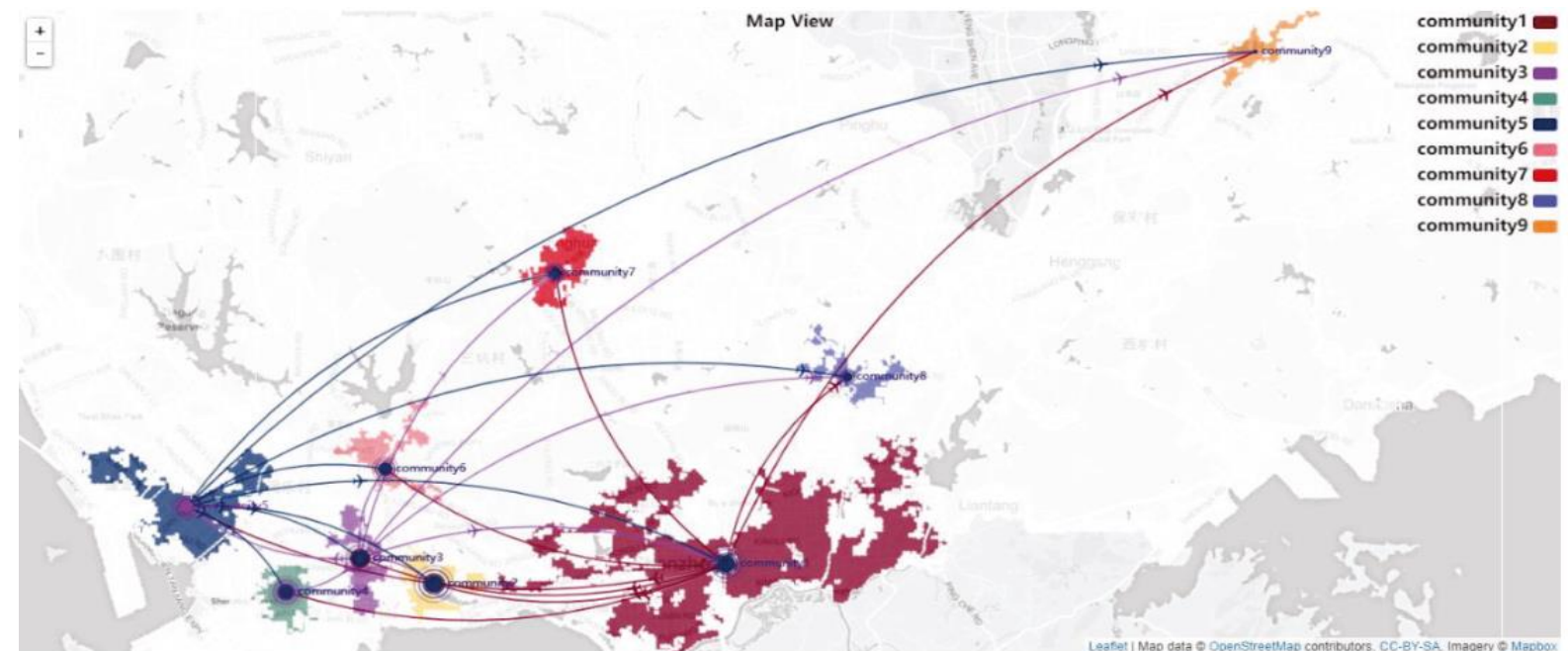
3.1.2 How is big data used in practice?

- 1- Geospatial data preprocessing
- 2- Overlay analysis
- 3- Land-use change prediction
- 4- Global scale terrain analysis
- 5- Human mobility (pattern discovery)
- 6- Disaster management (earthquake mitigation)
- 7- Geospatial problem solving
- 8- Geospatial big data management and searching (climate data).



3.1.3 What were the results?

For example, Zhang et al 2019 used advanced machine learning methods to identify time-varying mobility patterns based on smart card data and other urban data. The proposed approach delivers a comprehensive solution to pre-process, analyze, and visualize complex public transit travel patterns.



3.1.4 what data was used?

In human mobility research we can use **smart card automated** to collect large volume of travel data at the individual level. We also can use **bus trajectory data**, **public transit network**, and **road network data**.



3.1.5 What are the technical details? (Human mobility)

With more than **6 million records collected** for each day, the size of the dataset for the week amounts to **6.5 GB**. Each day, the bus trajectory dataset has approximately **63–73 million GPS records**.





3.1 Geospatial Applications



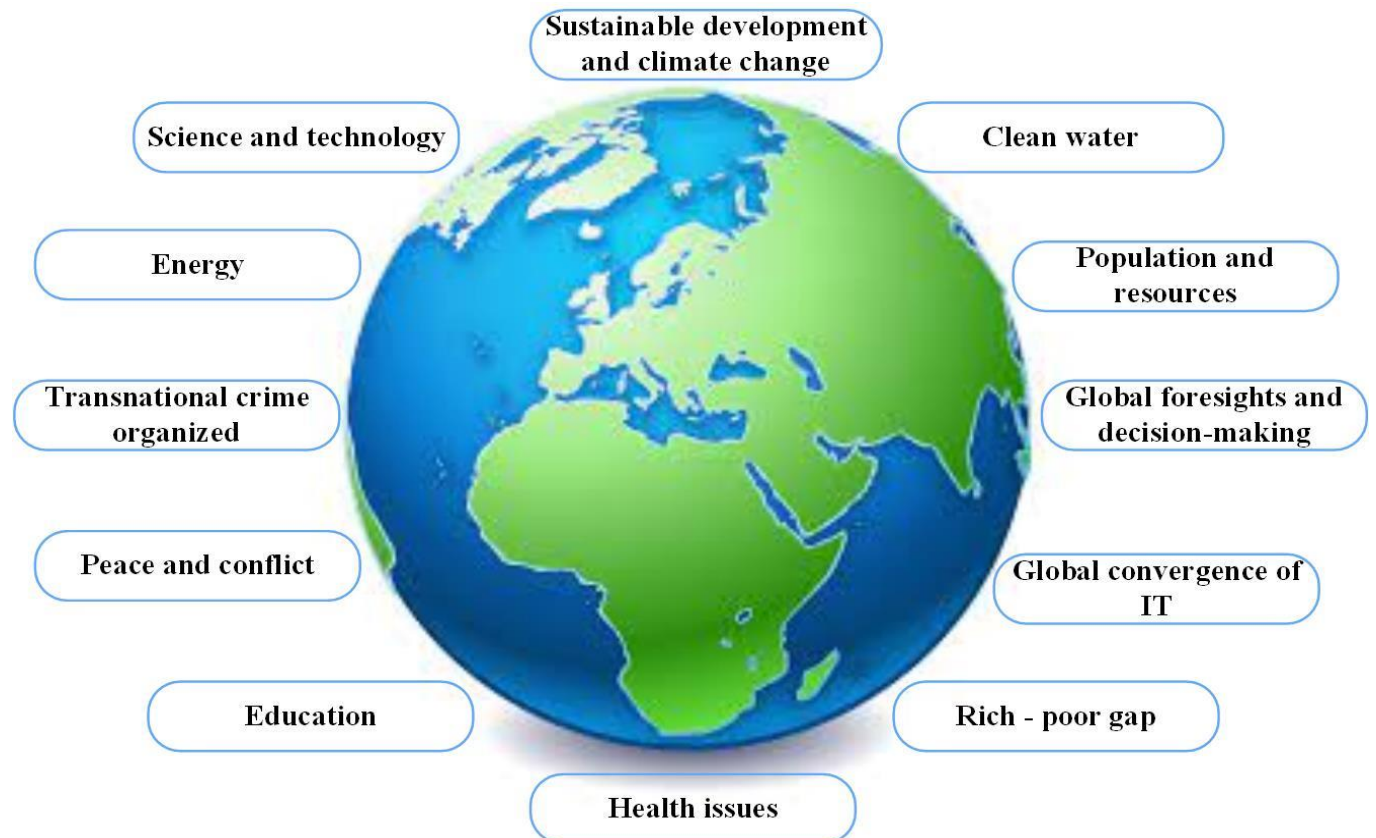
3.1.6 Any challenges that had to be overcome?

- Build on this integration of knowledge and skills across the disciplines of GIScience and computational science, which has been termed cyber literacy for GIScience;
- More efforts should be devoted to identifying geospatial applications of great impact and benefiting from the integration of geospatial methods and parallelization in the big data era;
- Further research directions should focus on improving and optimizing the performance of big data frameworks from different aspects.

3.1.7 what are the key learning points and takeaways?

Integrated studies may imply managing and visualization of multiple complex systems. BIG

DATA could help to handle big geospatial data produce every minutes in the world.



3.2.1 what problem is big data helping to solve?

The search for hydrocarbons necessitates massive amounts of manpower, equipment, and energy. With the average deep-water oil well costing \$100 million or more to drill, **it's critical that drilling takes place in the best locations.**



3.2.2 How is big data used in practice?

Previously a survey might have involved **a few thousand readings** being taken, today it will typically involve **over a million**. This data is then uploaded to analytics systems and compared with data from other drilling sites around the world.

Big Data is used at Shell to monitor the performance and condition of their equipment.

Across its logistics, distribution and retail functions, Big Data is amalgamated from many external sources, including local economic factors and meteorological data, which go into complex algorithms designed to determine the price we end up paying at the pumps.



3.2 SHELL



3.2.3 What data was used?

Shell collects data that allows them to calculate the probable size of oil and gas resources by monitoring seismic waves below the surface of the earth. The exact nature of these measurements and the analytics is a closely guarded **commercial secret**.

3.2.4 what are the technical details?

Shell use fiber-optic cables and sensor technology developed by Hewlett-Packard to carry out their surveys of potential drilling sites. The data is stored and analyzed using Hadoop infrastructure running on Amazon Web Service servers. Data volumes are also an industry secret.

It's estimated that so far Shell has generated around **46 petabytes** through their data-driven oilfield programme. Shell is also known to have worked with IBM and movie special effects experts at DreamWorks to produce the visualization tools that give analysts 3D and 4D representations allowing them to explore forecasted reserves.

3.2.5 Any challenges that had to be overcome?

The huge increase in the amount of data being generated at oilfields means that increasingly advanced analytics must be developed in order to more efficiently determine valuable signals amongst the background “noise” of the data. **Large-scale system upgrades** were needed as existing analytics platforms were not capable of carrying out the predictive analytics necessary to accurately make forecasts from the Big Data being generated.



3.2 SHELL



3.2.6 What Are The Key Learning Points And Takeaways?

Until science and society evolve to the point where we have reliable alternatives, the world is dependent on fossil fuel. With the difficulty of finding new reserves rising along with the cost of extraction, **Big Data holds the key to driving efficiency and reducing costs of extraction and distribution.**

Walmart is the world's largest retailer, with over **two million employees** and **20,000 stores** in **28 countries**. With operations on this scale it's no surprise that they have long seen the value in data analytics. In 2015, the company announced that it was building the world's largest private data cloud, capable of processing **2.5 petabytes of data per hour**.

3.3.1 What problem is Big Data helping to solve?



Supermarkets compete on price, customer service, convenience

Getting the right items to the right people at the right time poses significant logistical challenges.

3.3.2 how is big data used in practice?

@WalmartLabs

Walmart developed *The Data Cafe hub*

In 2011, with a growing needs of customers, Walmart established @**WalmartLabs** and their fast Big Data Team to research and deploy new data-led initiatives across the business.

Sales across different stores in different geographical location can also be monitored in real-time.

Another initiative is **Walmart's Social Genome Project**, for example, this project monitors public social media conversations and attempts to predict which products people will buy based on their conversations.

3.3.3 What were the results?

According to Walmart, the Data Cafe system has reduced the time it takes from detecting a problem in the numbers to proposing a solution from an average of two to three weeks to around 20 minutes.

3.3.4 What data was used?

The Data Cafe has a database with **200 billion rows of transactional data** that is regularly updated, and that **only covers the last few weeks** of activity.

It also takes data from 200 other sources, such as meteorological data, economic data, telecoms data, social media data, gas prices, and a record of events occurring near Walmart stores.

3.3.5 What are the technical details?

The real-time transactional database at Walmart is **40 petabytes in size**. Despite the large number of transactional data, it only includes data from the most recent weeks, as this is where the true value in real-time analysis can be discovered. On **Hadoop**, data from the chain's stores, internet divisions, and corporate entities is centralized (a distributed data storage and data management system).

3.3.6 Any challenges that had to be overcome?

Walmart's rapid growth necessitated a large intake of new employees, and finding the right people with the right skills proved difficult.

3.3.7 What are the main takeaways and learning points?

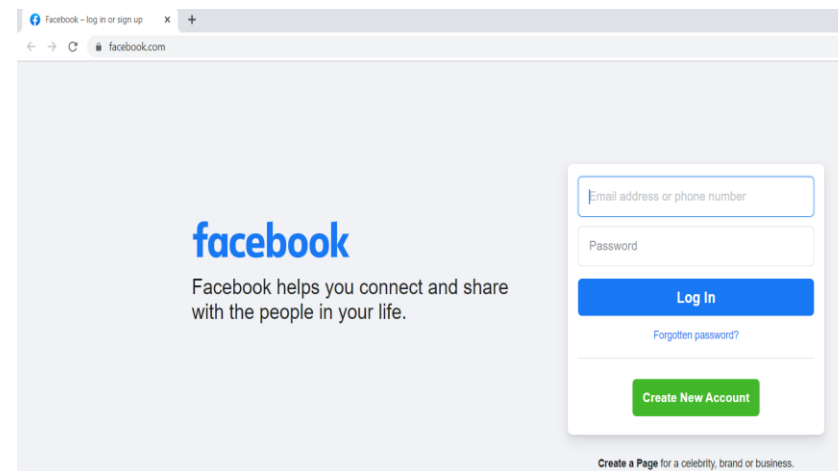
Supermarkets are large, fast-changing businesses that are complex organisms with many individual subsystems. This makes them an ideal business for implementing Big Data analytics.

How Facebook use big data to understand customers?

Facebook is still the world's largest social network. Everyone uses it to stay in touch with friends, share special occasions, and plan social gatherings. Every day, millions of people use it to read the news and interact with others brand names and make purchasing decisions.



The money used to pay their 10,000-plus employees and keep their services online comes from businesses that pay to access the data Facebook collects on us as we use their services.



3.4.1 What problem is big data helping to solve?

To survive, businesses must sell **goods** and **services**.

Big companies can use TV, Newspaper, radio and display advertising



A small business just starting should think carefully about the most efficient way of spending its limited marketing budget.

These companies can't afford to cover all the bases, so tools that can help them work out who their customers are..., and where to find them, can be hugely beneficial.



3.4.2 How is big data used in practice?

Facebook often have full access to demographic data.

Facebook is used to match them with companies which offer products and services that, statistically, they are likely to be interested in.

How many friends
we have?



where we live,
work, play?

what are the particular movies, books
and musicians we like?

What we do
in our spare
time?

3.4.3 What Were The Results?

Facebook's strategy of selling advertising space by leveraging their massive wealth of consumer data resulted in them having a 24 percent share of the US online display ads market in 2014, and **generating \$5.3 billion in revenue from ad sales**. By 2017, it was about 27 percent share worth more than **\$10 billion**.

3.4.4 what data was used?

Facebook, together with its users, **generates its own data**. Users upload 2.5 million pieces of content every minute.

This content is analyzed for information about us that advertisers can use to segment us. Additionally, they interact with other people's content as well as data stored in Facebook's own databases, which include business listings and databases of films, music, books and TV shows. Whenever we "Like" and share this content, it learns a little bit more about us.

3.4.5 What Are The Technical Details?

Facebook is the most visited Web page in the world after Google's search engine. It is said to account for around 10% of all online traffic. Of course, a Web service of this size requires a huge amount of infrastructure.

Its custom-designed servers, built using **Intel** and **AMD chips**, and power-saving technology to help cut down the huge cost of keeping so many machines running 24/7.

The designs for the server systems have been made available as open-source documentation.



3.4.6 Any challenges that had to be overcome?



Big company in The world rigidly abided by the terms of their privacy and data-sharing policies, the most watertight policies in the world are powerless in the face of data loss or theft, such as **hacking attacks**.

Gaining the trust of users is essential. Aside from data thefts and such illegal activity, users can become annoyed simply by being subjected to adverts they aren't interested in, too frequently.

3.4.7 What are the key learning points and takeaways?

Facebook has revolutionized the way we communicate with each other online by allowing us to build our own network and choose who we share information about our lives with.

This information is **extremely valuable to advertisers** because it allows them to precisely target their products and services to people who are statistically more likely to want or need them.



3.5 GOVERNMENT



Using big data to run a country

An example of the US government is explain in this section:

Barack Obama was called “The Big Data president” by The Washington Post, after committing to a \$200 million investment in data analytics and a pledge to make as much government-gathered data as possible available to the public.

For better or worse, Obama's presidency has coincided with the massive explosion in data collection, storage, and analysis known as Big Data.



3.5 GOVERNMENT



3.5.1 What problem is big data helping to solve?

Administering US big economic power and its population of 300 million-plus people clearly takes a colossal amount of effort and resources.

The Government have responsibility for **national security, economic security, healthcare, law enforcement, disaster recovery, food production, education ...**

To solve several problem, the US Government appointed the country's first ever **chief data scientist: D. J. Patil**, and before taking it up he had been employed at the Department of Defense, where he analyzed social media attempting to detect terrorism threats. He has also held positions at LinkedIn, Skype, PayPal and eBay.

3.5.2 How is big data used in practice? (1/2)

The US Government **Collect a large number of data-driven strategies** among their numerous departments and agencies, including networks of automated license plate recognition (ALPR) scanners and monitors tracking the flow of car, train and air passengers **to discern where infrastructure investment is needed.**

In education, Because more and more learning at schools and colleges is taking place online, those in charge of setting education policy can gain a better understanding of how the population learns and assess the level of education and skills in a specific geographical area, allowing for more efficient resource planning and deployment.

3.5.2 How is big data used in practice? (2/2)

In healthcare, social media analysis is used by the Centers for Disease Control (CDC) to track the spread of epidemics (for example COVID 19) and other public health threats. Also the National Institutes of Health launched the Big Data to Knowledge (BD2K) project in 2012 to encourage healthcare innovation through data analytics.

In security, The CIA were also partly responsible, through investments, for the rise of predictive security specialists Palantir, which use predictive data algorithms **to fight international and domestic terrorism and financial fraud.**

In Agriculture, the Department of Agriculture carry out research and scientific analysis of farming and food production, based on Big Data gathered in fields and farmyards.

3.5.3 What Were The Results?

In 2014, it was reported that a predictive fraud prevention system used by US administrators of the Medicare and Medicaid services had **prevented \$820 million in fraudulent payments** being made since its introduction three years earlier.

After a thorough review of the methodologies, either currently adopted or planned, US concluded: “Big Data tools offer astonishing and powerful opportunities to unlock previously inaccessible insights from new and existing datasets.”

“Big Data can fuel developments and discoveries in health care and education, in agriculture and energy use, and in how businesses organize their supply chains and monitor their equipment. Big Data holds the potential to streamline the provision of public services, increase the efficient use of taxpayer dollars at every level of government and substantially strengthen national security.”

3.5.4 What data was used?

The US Government monitor, collect and analyze a vast volume and variety of data, both through their own agencies, such as the Food and Drug Administration, county and law enforcement, and with a wide range of third-party partners.

Climate and meteorological data

Food production data from agriculture

Statistics on crime and security threats

Population movement data from camera networks and demographic research

3.5.5 what are the technical details?

The open-source software WordPress and CKAN are used to build and maintain the interface that makes the data available to the public. US developed a public data portal

Data.gov – the online portal where, by US Government decree (the 2009 Open Government Directive), all agencies must make their data available, has grown from 49 datasets at launch to close to 190,000 datasets today. The biggest individual contributors are NASA (31,000 datasets uploaded), the Department of the Interior (31,000) and the Department of Commerce (63,000).

3.5.6 Any challenges that had to be overcome?

Without doubt, the single biggest challenge facing the US Government in their mission to collect and analyze data has been **public trust** and **security**.

There have been numerous requests for more openness into government data collecting, which, when done without the consent of the people it is being collected from, can result infrequently seen as simply "spying" by those who are its targets.

This was undoubtedly the stimulus for Obama's Open Data Initiative as well as ongoing efforts to increase public understanding of the work carried out by Patil and the Office of Science and Technology Policy.

3.5.7 What are the key learning points and takeaways?

Big Data has enormous potential for driving efficiencies that could improve people's lives all over the world, so **it's critical that governments learn how to handle Big Data** in a way that doesn't cause discomfort or suspicion among their citizens.

US administration, have come to the conclusion that the potential benefit of Big Data outweighs the potential negative impact.

This is evident by the continued increase in investment in data collection and analytics, as well as the concerted efforts being made by politicians to playdown our fears by pointing to increased transparency and accountability.

04

CONCLUSION





Conclusion



Every field need a data scientist. Professional **big data developers** are mostly valued when they have a strong technical background and great problem solving skills.

Developing and managing data in your field can offer to you a great carrier.



Recommendations

If you want to get a job and be successful in China we suggest you to :

- Learn Chinese language and Chinese culture (if your participate to school activities you may learn more about how to work with Chinese. If you make Chinese friends your may learn more the culture and the language.....)
- Read books (not only articles) to improve your knowledge and to find innovative solutions in your field;
- Be on time and serious when your are doing your work;
- Evaluate yourself (what will be your adding value in the company? The lab? Are you the right person to do the job).



REFERENCES



- Arockia Panimalar.S, Varnekha Shree.S and Kathrine.A, V., 2017. The 17 V's Of Big Data. International Research Journal of Engineering and Technology, 4(9).
- Chebbi, I., Wadii Boulila and Farah, I.R., 2015. Big Data: Concepts, Challenges and Applications. Lecture Notes in Computer Science. 10.1007/978-3-319-24306-1.
- Li, Z., Tang, W., Huang, Q., Shook, E. and Guan, Q., 2020. Big Data Computing for Geospatial Applications. International Journal of Geo-Information.
- Marr, B., 2016. Big Data in practice: how 45 successful companies used big data analytics to deliver extraordinary results.
- Zhang, T., Wang, J., Cui, C., Li, Y., He, W., Lu, Y. and Qiao, Q., 2019. Integrating Geovisual Analytics with Machine Learning for Human Mobility Pattern Discovery. ISPRS International Journal of Geo-Information, 8(10). 10.3390/ijgi8100434.



Silk Road Institute (212), Nanwangshan Campus, CUG
20. 06. 2022 Wuhan., Hubei, China

Lenikpoho Karim, Coulibaly

Professor Qingfeng Guan

High-Performance Spatial Computational Intelligence Lab

School of Geography and Information Engineering

China University of Geosciences, Wuhan, China

THANK YOU!



High-performance Spatial Computational Intelligence Lab @ CUG