## Applied Mathematics & Information Sciences
*An International Journal*

# Real-Time Big Data Processing Framework: Challenges and Solutions

*Zhigao Zheng*[1,2]*, Ping Wang*[1,3,4,*] *and Jing Liu*[3] *and Shengli Sun* [1]

[1] School of Software and Microelectronics, Peking University, Beijing 100260, China
[2] National Engineering Research Center for E-Learning, Central China Normal University, Wuhan 430079, China
[3] School of Electronics Engineering and Computer Science, Peking University, Beijing 100791, China
[4] National Engineering Research Center of Software Engineering, Peking University, Beijing 100791, China

**Abstract:** Data type and amount in human society is growing in amazing speed which is caused by emerging new services as cloud computing, internet of things and location-based services, the era of big data has arrived. As data has been fundamental resource, how to manage and utilize big data better has attracted much attention. Especially, with the development of internet of things, how to processing large amount real-time data has become a great challenge in research and applications. Recently, cloud computing technology has attracted much attention with high-performance, but how to use cloud computing technology for large-scale real-time data processing has not been studied. This paper studied the challenges of big data firstly and concludes all these challenges into six issues. In order to improve the performance of real-time processing of large data, this paper builds a kind of real-time big data processing (RTDP) architecture based on the cloud computing technology and then proposed the four layers of the architecture, and hierarchical computing model. This paper proposed a multi-level storage model and the LMA-based application deployment method to meet the real-time and heterogeneity requirements of RTDP system. We use DSMS, CEP, batch-based MapReduce and other processing mode and FPGA, GPU, CPU, ASIC technologies differently to processing the data at the terminal of data collection. We structured the data and then upload to the cloud server and MapReduce the data combined with the powerful computing capabilities cloud architecture. This paper points out the general framework for future RTDP system and calculation methods, is currently the general method RTDP system design.

**Keywords:** Big data; cloud computing; data stream; hardware software co-design; CEP

## 1 Introduction

With the development of internet of things, various large-scale real-time data processing based on real-time sensor data are becoming the key of the construction of EPC (epcglobal network) application currently. The academia, the industry and even the government institute have already begun to pay close attention to big data issues and generated a keen interest.

In May 2011, the world-renowned consulting firm McKinley released a detailed report on big data Big data: The next frontier for innovation, competition, and productivity [1], sparking a broad discussion of Big Data. The report gave a detailed analysis of the impact of big data, key technology and application areas. In January 2012, Davos World Economic Forum released a report entitled "Big Data, Big Impact: New Possibilities for International Development"[2], raising a research boom of Big Data. The report explored how to make a better use of the data to generate good social profits in the new data generation mode, and focused on the integration and utilization of mobile data produced by individual and other data. In March, the U.S. government released "Big Data Research and Development Initiative"[3] to put the research of Big Data on the agenda, and officially launched the "Big Data development plan".

The academia started the study of big data much earlier. In 2008, "Nature" had launched a special issue of Big Data [4]; Computing Community Consortium published a report entitled "Big data computing: Creating revolutionary breakthroughs in commerce, science, and society"[5], elaborated the necessary technology to solve the big-data problem and some of the challenges faced in the context of data-driven research. In February, Science

* Corresponding author e-mail: pwang@ss.pku.edu.cn

launched "dealing with Data" special issue [6], mainly discussing the Big-Data problems in scientific study and indicating the importance of big data for scientific research. Some well-known American experts and scholars in the field of data management from the perspective of a professional study, jointly released a white paper "Challenges and Opportunities with Big Data"[7]. This white paper from an academic perspective described the generation of large data, analyzed the processing of large data and proposed a large number of challenges faced by the big data.

The hot research of Big Data doesnt mean that people have a deep understanding of big data. So far, people do not have a clear and uniform definition of big data, and remain a lot of doubt and controversy on its key technologies and its applications[8]. Whats more, an issue of concern is that the real-time processing of this massive heterogeneous stream data is a huge challenge, and there is lack of support for massive real-time data processing framework and implementation techniques. The processing of this real-time stream data is much different from that of static data. It needs to meet the extremely high data throughput and strict real-time requirements. For example, smart grid systems require a real-time monitoring for the nationwide network, automatically control the power of disaster areas in the event of storms, rain and snow disasters and other special circumstances before causing significant losses, avoid secondary disasters caused due to power issue and at the same time save energy as much as possible, and distribute energy rationally. For large area power monitoring itself involves a large data processing problem. However a regulation before the disaster needs to be made in a very period time and when the second disaster caused by power occurs, an error analysis, a fault location and a troubleshooting should be done in a very short time, otherwise it will cause huge losses to the residents as well as the country. Like the "810" blackout in United States which brought American tens of billions of losses as well as "814" in Canada which caused huge losses[9][10].

Thus, the research of real-time big data has great application prospect and research value. Because of the real-time and the large scale of data processing and other features that real-time big data requires make the study for real-time big data processing challenging, mainly in real-time, stability and large-scale etc.

Real time: Real-time processing of big data mainly focuses on electricity, energy, smart city, intelligent transportation, and intelligent medical fields. During the information processing it needs to be able to make quick decisions, and feedback relevant instructions to the sensing terminal input within a very short time delay. For instance in Fire monitoring and rescue system, its processing center needs to be able to analyze and process the data collected by sensors in the site of the incident in a very short period of time, to make integrated decision by comprehensively considering site information such as the movement of persons and the site form and meantime to

issue the corresponding instructions to the site sensing terminals, such as what extinguishing agent used for rescue, how to protect peoples safety in the site of the incident and how to help the firemen to rescue. At the same time, the information gathered by sensing terminals and instruction information must arrive information gathering or processing terminal in real time and make relevant decisions. The loss caused by that decision-making information cant be conveyed in real time is also incalculable, so real-time processing of large data is particularly important.

Stability: the areas covered by Real-time processing of big data are mostly closely related to the people such as Smart City's intelligent transportation systems and high-speed train control system, and mostly highly associated with infrastructure which also determines the real-time data processing system in a large system architecture, hardware and software equipment and other aspects must possess high stability.

Large-scale: As discussed above, real-time big data processing systems are often closely related to urban infrastructure and major national application, so its application is often a huge scale. Such as smart city intelligent transportation system, once the largest real-time data analysis and decision are made, it often aims at transportation decision-making at a provincial and municipal level or even a national level, and has a significant impact on the national life.

Recent years, many experts and scholars have made a lot of research on large data processing methods. Especially after "Nature" released "Big Data" issue in 2008, many scholars' research on big data is in full swing. Many researchers has made numerous studies on large data processing methods[11][12][13][14][15], and proposed many efficient big data processing algorithms, but also led to a lot developments in many new research areas. For example the paper [16] thought that the theory and the algorithm framework built by the whole relationship network underlying the big data is an important area after the theory of quantum mechanics; paper [17] introduced the challenges and opportunities that big data research would meet and also pointed out that the real-time processing of large data was one of the future research directions.

This paper, beginning with the nature of big data issue, briefly introduces the basic concept of large data, and analyzes characteristics of big data, challenges real-time big data processing met and the differences between real-time big data and big data and attributes to the five aspects. On this basis, combing cloud computing technology large data processing framework is outlined. We believe that Cloud Computing technology and Big Data are interdependent: Cloud Computing technology supports storage, management and data analysis for Big Data; Big Data provides an extensible platform for cloud computing. Thus this paper highlights the big data processing architecture under the cloud computing platform. It presents a data storage solution on

heterogeneous platforms in real-time big data processing system, constructs a calculation mode for big data processing and points out the general framework for real-time big data processing which provides the basis for the RTDP (Real-Time Data Processing). Section 2 in this paper gives an overview of big data; section 3 discusses the differences and challenges between big data and real-time big data; section 4 points out the current deficiencies of cloud computing technology; section 5 combine cloud computing technology with the feature of real time big data to architect the processing platform for real-time big data; section 6 gives a demo about how the RTDP system is used in smart grid system and finally summarize this article.

## 2 Big Data Overview

Big data itself is a relatively abstract concept, so far there is not a clear and uniform definition. Many scholars, organizational structure and research institutes gave out their own definition of big data [1][18][19][20]. Currently the definition for large data is difficult to reach a full consensus, the paper references Academician Li Guojies definition for big data: in general sense, Big Data refers to a data collection that cant be obtained within a tolerable time by using traditional IT technology, hardware and software tools for their perception, acquisition, management, processing and service[21]. Real-time data is a big data that is generated in real time and requires real-time processing.

According to the definition of Big Data, Big Data is characterized by volume, velocity and variety where traditional data processing methods and tools cannot be qualified. Volume means a very large amount of data, particularly in data storage and computation. By 2010 the global amount of information would rapidly up to 988 billion GB [22]. Experts predict that by 2020 annual data will increase 43 times. Velocity means the speed of data growth is increasing, meanwhile peoples requirements for data storage and processing speed are also rising. Purely in scientific research, annual volume of new data accumulated by the Large Hadron Collider is about 15PB [23]. In the field of electronic commerce, Wal-Mart's sells every day more than 267 million (267Million) products [24]. Data processing requires faster speed, and in many areas data have been requested to carry out in real-time processing such as disaster prediction and rapid disaster rehabilitation under certain conditions need quickly quantify on the extent of the disaster, the regional scope impacted and etc. Variety refers to the data that contains structured data table, semi-structured and unstructured text, video, images and other information, and the interaction between data is very frequent and widespread. It specifically includes diverse data sources, various data types, and a strong correlation between the data.

With the development of computer and network technology, as well as intelligent systems is common used in modern life, big data has become increasingly close to people's daily lives. In 2008, Big Data issue released by "Nature" pointed out the importance of big data in biology, and it was necessary to build biological big data system to solve complex biological data structure problem [25]. Paper [25] pointed out that the new big data system must be able to tolerate various structures of data and unstructured data, has flexible operability and must ensure data reusability. Furthermore, Big Data plays an important role in the defense of national network digital security, maintaining social stability and promoting sustainable economic and social development [26]. With the development of big data technology, Big Data also plays an important role in creating a smart city, and has important applications in urban planning, intelligent traffic management, monitoring public opinion, safety protection and many other fields [27].

## 3 Difference and Challenges between big data and real-time big data

Big data is characteristic by multi-source heterogeneous data, widely distributed, dynamic growth, and "data mode after the data"[28][29]. In addition to having all the characteristics with big data, real-time big data has its own characteristics. Compared with the big data, when it comes to data integration real-time big data has higher requirements in data acquisition devices, data analysis tools, data security, and other aspects. The following introduces from data integration, data analysis, data security, data management and benchmarking.

### 3.1 Data Collect

With the development of internet of things [30] and Cyber Physical System (CPS) [31], the real time of data processing requires higher and higher. Under the big data environment, numerous sensors and mobile terminals disperse in different data management system which makes data collection itself a problem. In RTDP system, its real time data collection faced makes data integration facing many challenges.

1. Extensive heterogeneity

In big data system, the data generated by mobile terminals, tablet computers, UPS and other terminals is often stored in cache, but in RTDP system it requires data synchronization which brings tremendous challenges to the wireless network transmission. When dealing with processing heterogeneity, big data system can use NoSQL technology and other new storage methods, such as Hadoop HDFS. But the real time requires low in this kind of storage technology, where the data is often stored once but read many times. However this kind of storage technology is far from satisfying the requirement of real-time big data system that requires data

synchronization. Due to extensive heterogeneity of big data, data conversion must be carried out during data integrations processing, however traditional data warehouse has obviously insufficient to meet the needs of time and scale that big data requires [32][33][34].

2. Data quality insurance

In the era of big data it is a phenomenon often appears that useful information is being submerged in a large number of useless information [6]. The data quality of Big Data has two problems: how to manage large-scale data and how to wash it. During the cleaning process, if the cleaning granularity is too small, it is easy to filter out the useful information; if the cleaning granularity is too coarse, it can't achieve the real cleaning effect. So between the quantity and quality it requires careful consideration and weighed which is more evident in real-time big data system. On the one hand, it requires system to synchronize data in a very short time; on the other hand, it also requires the system to make a quick response to data in real time. The performance requirements of the speed of data transmission and data analysis are increasing. Moreover the data may be filtered at a time node may become critical post processing data. Therefore, how to grasp the correlation between data and accurately determine the usefulness and effectiveness of data becomes a serious challenge.

### 3.2 Data Analytics

Data analysis is definitely not a new problem. Traditional data analysis is mainly launched for structured data source, and already has a complete and effective system. On the basis of the data warehouse, it builds a data cube for online analytical processing (OLAP). Data mining technology makes it possible to find deeper knowledge from large amounts of data. But with the arrival of the era of big data, the volume of numerous semi-structured and unstructured data rapidly grows, which brings huge impact and challenges to the traditional analysis techniques and existing processes are no longer applicable. It mainly reflects in timeliness and index design under dynamic environment.

1. Timeliness of data processing

In the era of big data, time is value. As time goes by, the value of knowledge contained in the data is also attenuation. In real-time data systems, time is required higher. For example, in a data processing of disaster analysis, real-time high-speed trains, aircraft and other high timeliness performance device, time has gone beyond economic value. Damages caused by unreasonable delay would be hard to estimate. The era of real-time big data proposes a new and higher requirement to the timelines of data processing, mainly in the selection and improvement of data processing mode. Real-time data processing modes mainly includes three modes: streaming mode, batch mode and a combination of two-a mixed processing mode. Although currently many

scholars have made a great contribution to real-time data processing mode, yet there is no common framework for real-time processing of large data. 2. Index design under dynamic environment

The data pattern in the era of big data may be changing constantly as data volume varies and existing relational database index is no longer applicable. How to design a simple, efficient and able to quickly make an adaptation has become a one of the major challenges of big data processing when data mode changes. Current solution is basically built an index by NoSQL databases to solve this problem, but they have been unable to meet the demand for real-time processing of big data.

3. Lack of prior knowledge

On the one hand, because semi-structured and unstructured data abound, it is difficult to build its internal formal relations when analyzing the data; On the other hand it is difficult for these data needed to be processed in real time to have sufficient time to establish a priori knowledge due to the coming of the data stream in the form of an endless stream.

### 3.3 Data Security

Data privacy issues associated with the advent of computers has been in existence. In the era of big data, the Internet makes it easier to produce and disseminate data, which makes data privacy problems get worse, especially in real-time processing of large data. On the one hand, it requires data transmission real-time synchronization; on the other hand, it demands strict protection for data privacy, which both raise new demands to system architecture and computing power.

1. Expose hidden data

With the appearance of the Internet, especially the appearance of social networks, people are increasingly used to leave data footprints. Through data extraction and integration technology, accumulate and associate these data footprints may cause privacy exposure. In real-time big data processing, how to ensure the speed of processing a data as well as data security is a key issue which has troubled many scholars.

2. Data disclosure conflicts with privacy protection

By hiding data to protect privacy it will lose the value of data, thus it is essential to public data. Especially by digging accumulated real-time large-scale data, we can draw a lot of useful information, which has a great value. How to ensure the balance between data privacy and data publicly is currently in research and application a difficulty and hot issue. Therefore the data privacy in the era of big data is mainly reflected in digging data under the premise of not exposing sensitive information of the user. Paper [35][35] proposed privacy preserving data mining concept, and many scholars have started to focus on research in this area. However there are is a conflict between the amount of information and the privacy of data, and thats why so far it has not yet a good solution. A

new differential privacy method proposed by Dwork may be a way to solve the protection of data privacy in big data, but this technology is still far from practical applications [36].

## 3.4 Usability issue of data management

Its challenges mainly reflect in two aspects: huge data volume, complex analysis, various result forms; numerous industries involved by big data. However analysis experts lack of knowledge of both aspects relatively. As a result, the usability of real-time big data management mainly reflects in easy to discover, easy to learn and easy to use [37]. Thus in order to achieve usability of big data management, there are three basic principles to be cared as follows:

1. Visibility

Visibility requires the use of the data and the results be showed clearly in a very intuitive way. How to achieve more methods of large data processing and tools simplification and automation will be a major challenge in the future. Ultra-large-scale data visualization itself is a problem, while real-time visualization of large-scale data will spend a lot of computing resources and GPU resources. Thus how to enhance the performance and utilization of the GPU is a very serious challenge.

2. Mapping

How to match a new big data processing technique to processing techniques and methods people have become accustomed to and achieve fast programming is a great challenge to data usability in the future. For MapReduce lacks SQL-like standard language, the researchers developed a higher level languages and systems. Typical representatives are the Hadoop Hive SQL [32] and Pig Latin [38], Google's Sawzall [39], Microsoft's SCOPE [40] and DryadLINQ [41] as well as MRQL [42], etc. But how to apply these languages and systems to real-time big data processing still remain big challenges.

3. Feedback

Feedback design allows people to keep track of their operating processes. Works about this aspect is few in Big Data field [43][44][45]. In the era of big data the internal structure of many tools is very complex. And in software debugging it is similar to Black Box debugging for the normal users and the procedure is complex as well as lack of feedback. If in the future human-computer interaction technology can be introduced in the pressure of big data, people can be more fully involved in the whole analysis process, which will effectively improve the user's feedback sense and greatly improve the ease of use.

A design meet the above three principles will be able to have a good ease of use. Visualization, human-computer interaction and data origin techniques can effectively enhance usability. Behind these technologies, massive metadata management needs our special attention [46]. So how to achieve an efficient management of the massive metadata in a large-scale

storage system will have an import impact on the usability of real-time big data.

## 3.5 Test benchmark of performance

A very important aspect for big data management is the quality assurance, especially for real-time management of large data as disaster caused by data error will be very serious and even immeasurable. The first step in quality assurance is to do performance testing. There is not yet a test benchmark for the management of big data. Main challenges faced by building big data benchmarks are as followings[47]:

1. High complexity of system

Real-time big data is highly heterogeneous in data format as well as hardware and software and it is difficult to model all big data products with a uniform model. Real-time big data system requires high timeliness which makes it hard to extract a representative user behavior in real time. Whats more, data size is very large and data is very difficult to reproduce which both make the test more difficult.

2. Rapid revolution of system

The traditional relational database system architecture is relatively stable, but the data in real-time big data processing is in a constant state of growth, and there is a certain correlation between the data, which makes the benchmark test results obtained soon not reflect the current system actual performance. In real-time big data system test results are required to be completed within a very short time delay with high accuracy, which in the hardware and software aspects is a serious challenge to the test benchmark.

3. Reconstruct or reuse existing test benchmark

Extend and reuse on the existing benchmarks will greatly reduce the workload of building a new large data test benchmark. Potential candidates standards are SWIM (Statistical Workload Injector for MapReduce) [48], MRBS [49], Hadoop own GridMix [50], TPC-DS [51], YCSB++ [52], etc. But these benchmarks are no longer applicable in real-time big data processing.

Now there are already some researches focusing on the construction of big datas test benchmark, but there is also a view which thinks its premature to discuss that currently. By tracking and analyzing the loads of seven products which are applied with MapReduce technology, Chen et al [47][53] think it is impossible to determine typical user scenarios in the era of big data. In general, building big data and real-time big data test benchmark is necessary. But the challenges it will face are a lot, and it is very difficult to build a recognized testing standards like TPC.
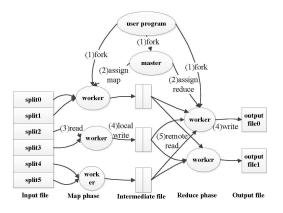
**Fig. 1:** The framework of the MapReduce Model

# 4 Shortcomings of cloud computing architecture

## 4.1 Cloud Computing Overview

Cloud computing is the product of the traditional computer technology and network technology development integration such as grid computing, distributed computing, parallel computing, utility computing, network storage, virtualization, load balancing, etc., which aims at integrating multiple relative low-cost computing entities into one perfect system with powerful computing ability via the network and with the help of the SaaS, PaaS, IaaS, MSP and other advanced business models distributing this powerful computing ability to the hands of the end user [54][55]. Cloud computing system mainly uses MapReduce model. The core design idea of MapReduce is to divide and conquer the problem and calculate on data rather than push data to calculate which effectively avoids a lot of communication costs generated during data transmission. A full pressure of MapReduce is shown in Figure 1[56]. At first, the MapReduce model partitions users original data source, and then hand them to different task areas to process. The Map task takes a series of Key/Value pairs from the input, process each with user-defined Map function to obtain middle results, and write the results into local disks. The Reduce task sorts data by the Key value after reading data from disks to put data with same key value together. Finally a user-defined Reduce function processes these ordered results and outputs the final result.

## 4.2 Shortcomings of cloud computing architecture

MapReduce model is simple, and in reality, many problems can be represented with MapReduce model. Thus MapReduce model has a high value as well as many application scenarios. But MapReduces achievement is

mainly relying on the Hadoop framework while the data processing method of Hadoop is "Store first post-processing" which is not applicable to real-time large data processing. Though currently there are some improved algorithms able to make Hadoop-based architecture almost real-time, for example, some latest technology like Cloudera Impala is trying to solve problems of processing real-time big data on Hadoop[57][58][59], the batch processing of Hadoop and its structural features make Hadoop defective in processing big data in real time[59]. Hadoops defect in real-time big data processing mainly reflects in data processing modes and application deployment. This paper will discuss these two aspects separately in the following.

### 4.2.1 Big Data Processing Mode

Big data processing mode can be divided into stream processing and batch processing [60][61]. The former is store-then-process, and the latter is straight-through-processing. In stream processing, the value of data reduces as time goes by which demanding real-time; in batch processing, data firstly is stored and then can be processed online and offline [46]. MapReduce is the most representative of the batch processing method [56].

In EPC applications based on real-time sensor data, real-time data stream is one of the most important ways to generate large data. Dramatic increase in the number of data processing objects and the amount of data for each data object results in an expanding historical data. Under this condition, conflicts between the requirement in real-time processing data stream of large-scale historical data and the defect of computing and storage ability become the new challenges in cloud computing and the EPC area. In paper [62][62], this problem is defined as a scalability issue of data stream processing. In order to support the storage and computing of a large-scale data, a relatively classic architecture currently used is multi-core cluster computing architectures with multiple CPU and a four storage structures–cache, memory, storage, and distributed storage as Figure 2. In this architecture, multi-core CPUs on nodes constitute local computing resources; compared to the distributed storage, memory and external memory on nodes constitute the local storage with a high speed.

However, existing MapReduce[63] methods like Hadoop, Phoenix [64] belong to a batch processing to the persistent data, which need to initialize runtime environment, repeatedly load and process large-scale data, execute synchronously Map and Reduce and transfer large amounts of data between nodes in each processing. When process data stream arrived continuously in batch processing method, if process a small-scale data batch every time, the system cost will be too large and timeliness will be limited; if system waits for arriving batch to reach a certain size, it will increase
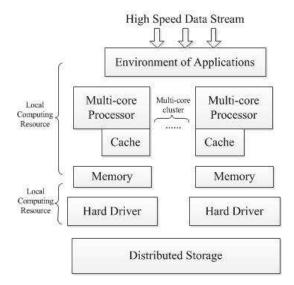
**Fig. 2:** Supporting Environment

processing delays which also cant meet the needs of system in real time. Therefore, for real-time large-scale data processing needs in high-speed data stream, how to design its processing model needs to be reconsidered.

#### 4.2.2 Application deployment

In no shared cluster architecture, using MapReduce [63] programming model to solve the conflict between large-scale data processing needs and the insufficient of computing and storage capacity if the core technology of cloud computing. MapReduce provides a high-level programming model, through a simple programming interface supports parallel processing large-scale data that can be divided and blocks task scheduling, data storage and transmission and other details to the programmers, whose Programming granularity is higher. User-written MapReduce programs compiled in the Master node through the virtual machines will be deployed to Node.

However, real-time big data capture terminal are often some small sensor node with small memory which cant be deployed Java Virtual machine with large capacity. And in the real-time sensing EPC network, data capture sensor nodes has often platform heterogeneity and varying computing capability. For example, the DSP(Digital Signal Processor) used for a digital station to compression encode and decode digital audio and video, the STB chip codec, RF chip for wireless applications, luxury car electronic chip, ASIC (Application Specific Integrated Circuit) used to deal with complex event in medical field and multi-core CPU and GPU for scientific computing and etc. Thus, in real-time big data processing system Hadoop's MapReduce model cant be used to a unified deployment and how to design a computing model fit for real-time big data processing platform to deal with

platform heterogeneity and varying computing capability faces new challenges.

### 4.3 Other architectures

Except the Hadoop-based real-time big data processing architecture, researchers already design an architecture to deal with a streaming data based on the way to process the stream-oriented data, noticing that batch processing in Hadoop cant meet the feature of real-time big streaming data. For example, Twitters Storm processing mode, Apaches Spark and LinkedIns In-stream.

#### 4.3.1 Spark

Spark, as an advanced version of Hadoop, is a cluster distributed computing system that aims to make super-big data collection analytics fast. As the third generation product of Hadoop, Spark stores the middle results with internal storage instead of HDFS, improving Hadoops performance to some extent with a higher cost. Resilient Distributed Dataset, RDD, is an abstract use of distributed memory as well as the most fundamental abstract of Spark, achieving operating the local collection to operate the abstract of a distributed data set. Spark provides multiple types operations of data set which is called Transformations. Meantime it provides multiple actions, brought convenience to the upper development. Its more flexible than Hadoops programming model. However because of RDD, Spark is no more applied to applications of updating the status in asynchronous fine-grained and application models of incremental changes, which is an important application direction in real-time big data processing system.

#### 4.3.2 Storms

Storm cluster has some similarity with Hadoop. The difference is that its Job in MapReduce running in Hadoop cluster and Topology in Strom. Topology is the highest-level abstract in Storm. Every work process executes a sub-set of a Topology, which consists of multiple Workers running in several machines. But naturally the two frameworks are different. Job in MapReduce is a short-time task and dies with the tasks ending but Topology is a process waiting for a task and it will run all the time as system running unless is killed explicitly. In Storm cluster, it also has Master node and Worker node. Master node on which the background control program Nimbus runs charges the distribution of codes in cluster, the task allocation and the monitoring of Worker nodes status, and a Supervisor node runs on each work node. Supervisor monitors the task allocated by Nimbus of Work node which it belongs to and start/shut a
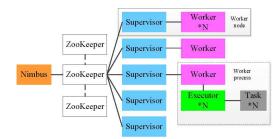
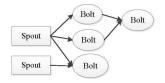**Fig. 3:** Physical Architecture of Storm



**Fig. 4:** Physical Architecture of Storm

Worker if needed. Here is the physical architecture of Storm as shown in Figure 3.

Nimbus is a Master node, mainly charging the submission and allocation of tasks and the monitor of cluster. Nimbus writes the task will allocate to Supervisor to Zookeeper, and coordinates it through Zookeeper. At the same time, in Zookeeper it stores public data like the heartbeat information and configuration information. However Supervisor charges accepting the tasks Nimbus allocates to itself, and manages its own Workers.

When processing streaming data, Storm mainly uses Stream, a key abstract, which is a Tuple sequence without boundary. Strom provides some primitive to transfer a Stream to a new Stream in a distributed and reliable way, among which the most basic primitives are Spout and Bolt, namely the input and output of a flow. Spout is a proactive role. Usually it will read data from an external data source (queues, databases, etc.), and then packaged into Tuple form, and then sent to the Stream. Bolt is a passive role, which can perform filtering, function operation, Join and operating a database or any other operations. It can process input Stream and generate a new output Stream. The network composed with Spout and Bolt will be packaged into a Topology, the relationship shown in Figure 4.

Through the agreement Preamble, to some extent Storm supports multi-language processing, but in essence, its the target language calls the interface JVM provides through Preamble. Currently the Storm team only provides the interface of Ruby, Python and Fancy version to the Preamble protocol, and doesnt achieve a real multi-language support. Except this, Storm still has some shortcomings, like data can only be transferred between Spout and Bolt in customized Topology and it is likely to

cause a waste of resource that cant be transmitted across the Topology.

In the current environment of relatively developed IOT, the highly heterogeneity of sensor data collection terminal means the diverse content of Spout, so in this case it will cause severe resource waste if processing real-time large-scale data with Storm. And the integration of the platform itself is an important constraints of the development of the real-time large data processing system.

## 5 Real-time big data processing framework

In addition to powerful computing ability, real-time big data processing system must have strong timeliness which means it must quickly respond to the request from system terminals in a very short time delay. So at first, real-time big data processing system must have powerful computing ability for big data. A traditional method to process big data is to rely on the powerful computing capabilities of the cloud computing platform to achieve, while for the timeliness it must rely on the ability of the rapid data exchange between system's internal and nodes. Paper [71] gave a conceptual description on real-time big data processing system, and divided system into data collection and storage, model building, model validation and deployment, real-time processing and model updating five stages from the point of application scenario.

This paper, according to the demands on the computing ability of real-time big data processing system and the timeliness, divides the RTDP (Real-Time Data Processing) framework into four layers–Data, Analytics, Integration and Decision from a functional level. Shown in Figure 5.

1. Data

This layer mainly charges for data collection and storage, but also including data cleaning and some simple data analysis, preparing data for Analytics. At the terminal of data collection, it needs to manage all terminals. For example, the FPGA commonly used in Data Stream Management System, DSMS; the ASIC used in Complex Event Processing, CEP; and CPU and GPU (Graphic Processing Unit) in batch processing system represented by MapReduce. Data storage module is responsible for the management of large-scale storage systems. Thanks to the heterogeneity of real-time data sources and the large data processing platform, RTDP systems can handle data from various data sources, including Hadoop for unstructured storage, the data warehouse system for structured storage and analysis, SQL databases, and some other data source system.

2. Analytics

This layer is the core of RTDP system and the critical layer to determine the performance of RTDP system. This layer is mainly responsible for data structure modeling, data cleansing and other data analysis processing, preparing data for the algorithm integration layer.
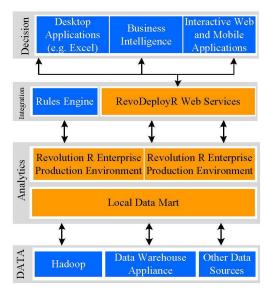
3. Integration

**Fig. 5:** Architecture of Real-Time Data Processing



**Fig. 6:** Adaptive Task Allocation of Data Stream Management System

extraction points of interest, select the characteristic function to determine the data formats and extract useful information from data marts, and several steps in which the data feature extraction for unstructured text data, etc. of data is very important, therefore, makes the feature extraction for data collection and storage is an important part of the process.

### 5.1.1 Data Collection and Data Base

RTDP systems heterogeneous platforms and performance makes RTDP system's data source contains a variety of ways, according to the data processing mode can be roughly divided into the CEP, DSMS, DBMS, based on a variety of ways such as MapReduce batch for each treatment have their different data acquisition techniques, such as remote medical field for surgical treatment of complex event processing scenarios for data acquisition ASIC, decoding audio and video coding in an FPGA, etc. Thus, during the data collection and management there are certain rules that must be collected on the side of the device identification, and can be based on different device programming overhead deployment and management nodes.

Based on the above analysis, the big data analysis problems need to be resolved first is the data acquisition side data preprocessing and data flow control for high-speed data stream management paper [65] proposed an adaptive massive real-time data flow management system adaptively according to the data flow and data distribution node preprocessing task, shown in Figure 6.

To enhance the data stream processing capabilities of DSMS, can be pre-distributed caching and reuse method of the intermediate results to avoid each data stream arrives historical repetition processing overhead and makes the data stream localization, reducing data between nodes transmission overhead for localized data stream processing, you can use event-driven stage of processing architecture[66] (Staged Event Driven Architecture, SEDA), using the thread pool technology to reduce cost of initializing each treatment, and by dividing the stage and at asynchronous transfer data between stages, eliminating data synchronization between stages based on

This layer plays a connecting role in RTDP system. In this layer it combines many common data processing algorithm packages. Depending on the scene it calls the appropriate algorithm for data analysis and data display, provides technical support for Analysis layer and at the same time provides a decision support and theoretical basis for Decision layer. Meanwhile the layer also needs to identify the device in data collection layer according to the rules been set and deploys applications.

4. Decision

This layer makes decisions with the results of data analysis which is the highest layer of data processing system as well as the ultimate goal of data analysis process. RTDP is a procedure involving numerous tools and systems interact with each other iteratively. At every level, the definition of "Big data" and "Real time" is not immutable. They have their own unique meaning at every level due to the functional association at each level. The four layers will be general process of RTDP in the future as well as the basic framework of the RTDP in this paper. Here we are going to discuss each layer in detail from the functionality, processing methods, related tools and deployment aspects of the system.

## 5.1 Data Layer

Since the data collected by sensors is rough and messy, and original data often contain too much useless data, modeling and data analysis for the tremendous difficulties, so the data collection process must be preliminary data analysis and filtering. first need to extract the data features, integrated data sources,
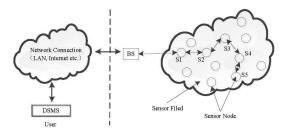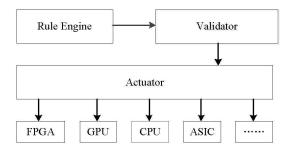
**Fig. 7:** Structure of Logical Management Adapter

this idea, the paper [62] proposed a large-scale high-speed data stream real-time data processing methods RTMR (Real-Time MapReduce). However, this method still exist some of the challenges: how to make effective use of the CPU's processing power; how to support local storage for intermediate results of high concurrent access.

To solve the above problem, this paper in the paper [65] proposed adaptive massive real-time data stream management system based on the combination of Logical Management Adapter (LMA) based on the underlying sensor data processing mode for dynamic management and control, including the adapter logical rules engine, validator and actuator of three parts, the structure shown in Figure 7.

Which rule builder using Embedded Software Component corresponding to different data processing to generate the corresponding node management rules , rules and validators for validation and matching nodes , and finally the actuator to the underlying FPGA, ASIC, GPU and other data collection node to manage and deploy the upper task execution , for example, in the wisdom of urban transportation system which , in essence, is to monitor the traffic in RTDP real-time video processing environment in which real-time monitoring system, because the data is in the form of a stream of data for processing , so you can use FPGA or ASIC video information collection , the bottom of the FPGA or ASIC large gathering real-time data stream composed of data , the corresponding data can be used for flow management DSMS .Complex Event Processing of how efficiently a plurality of basic event has a more complicated complex composite of semantics , including consideration of constraints between events , and even in some applications to continue to generate more complex event detection high-level composite events. several active database related work discussed for the basic model of the composite event detection[67], including : a model based on finite automata , Petri net -based model, based on matching tree model and based on directed graph models these models are also the basic model CEP problems , compared with other methods CEP technology has a tense , relevance , semantic richness , heterogeneity and mass and other characteristics[68] in the wisdom of the medical system, remote surgery system is essentially in RTDP

environment management and control of the CEP , because the remote surgery for surgical precision is very high , the system control engineering is fairly complex, suitable for surgery ASIC data collection , and for such elaborate and complex events are CEP can be used for data analysis and management.

An important issue in the data management process is in the organization of the data, a logical data structure will have a significant impact on the subsequent calculation of large data processing lies in a difficult, complex and heterogeneous hardware and software platform, design a reasonable data structure is resolved heterogeneous platform data storage and processing of a basic work is the most important step in how to organize in order to facilitate real-time processing of large data, there is no clear framework, the existing database systems and MapReduce systems are first collecting data and for processing the data, not suitable for real-time processing of the data, taking into account the real-time data are time-related, so we consider the data warehouses and data stream processing means [69]. Real-time data of the data has been in a state of flux until the data is stored and [70], so that the data stream is expected to become a real big way data is organized.

### 5.1.2 Data Storage and Cleaning

In RTDP, data comes from a wide range of sources, unstructured and structured data mixed, so Hadoop and other unstructured storage system in RTDP system has a natural advantage, but Hadoop itself does not achieve full real-time requirements, which determines our in real-time using Hadoop big data storage process Hadoop first need to solve real-time problems in the framework of the proposed RTDP use of multi-level storage architecture to solve the problem, its architecture is shown in Figure 8.

In RTDP multi-level storage system data through a lot of the local server first preliminary processing, and then uploaded to the cloud server for in-depth analysis and processing. Such architectural approach to solve the data filtering is how to determine the relevance of the issue of data is an important means, Since the real time processing of large data nodes need to collect data in the shortest possible time for rapid processing, but also need to filter out unwanted data, but the data collection process, we can confirm the current data be collected for post data key input, for data-dependent judgment is an extremely complex task.

In RTDP architecture, the local server preliminary data processing, the data collection terminal for rapid response in a very short period of time for the short delay the processing of requests for rapid response, and will not be able to determine the data-dependent data and present without processed data uploaded to the cloud server, the use of cloud computing power for subsequent analysis and processing work due to a limited number of local node, it is generally a PC on the local server capable of
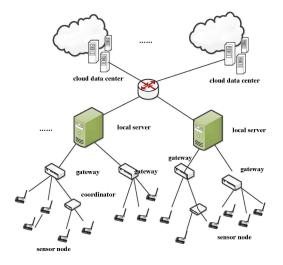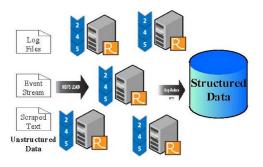
**Fig. 8:** Multi-Level Data Storage Model



**Fig. 9:** Structured Data

working. Computing resources and the reorganization of the local server can have the cloud computing resources of rationing server, make full use of local and cloud computing capabilities in RTDP architecture because data collection in real time, so multi-level storage system performance bottleneck is the network transmission speed. Use what network to ensure mass real-time data transmission is a major challenge. Article about the problem of transmission network will be discussed in detail in 5.3.

Data collection terminal data collected through the multi-level storage system, preliminary analysis of the processing of final upload to the cloud server, stored in a Hadoop cluster, or other data storage tool, so that the text and other unstructured data storage is relatively easy, but RTDP system during data processing and data presentation often for structured data processing more convenient, while in the RTDP systems and traditional data warehouse system is also required when performing data exchange using structured data, for which Smith proposed a model using the R language, will unstructured data into structured data approach to solve this problem [71].
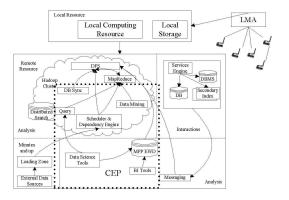


**Fig. 10:** Schematic of Storage and Analysis

## 5.2 Data Analysis Layer

This layer mainly carry through big data analysis and process, including data characteristic extraction, sample extraction, change of variables, model evaluation, model optimization, and model correction, as with common data warehouse doing big data processing. Its primary target is to found a robustness and easily comprehensive prediction model. Besides, the feature of RTDP systems instantaneity and big data processing decide that an impeccable RTDP system must be quickly, flexible and with good computing power as well as data reappear. Establishment of a robustness RTDP model is based on comprehensive understanding of needs and on the basis of a comprehensive analysis of the data is based on repeated comparisons of various models, verification basis.

To guarantee the flexibility and instantaneity of a RTDP system, in this thesis tasks in the RTDP frame are controlled prior according to time requirements, tasks with lowest time delay requirements have the highest priority, and the priority can be adjusted during real time process. Thus the system is divided into three modules: data storage system, analytical calculation system and ordering system. The data storage system is mainly using multi-level storage systems various storage mode, the analytical calculation system includes many RTDP algorithm packages, the ordering system task sorting section. The system structure is as shown in Figure 10.

### 5.2.1 Date Processing Algorithm

Along with architectural patterns, algorithm framework also plays an important role in RTDP systems computation results. In recent years, many research have done in big data processing algorithm, but the research about real-time big data has not been taken into account [72]. Cheng Yusheng[73] analyzed the reason why equivalent matrix rule extraction algorithm is so inefficient when dealing with big data set, then put forward a serial carry chain rule based extraction block
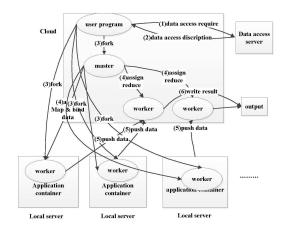
**Fig. 11:** New Big Data Processing Model

matrix algorithm. In the algorithm, condition attribute and decision attribute equivalent matrices merged into one matrix, thus lowered the scale of the equivalent matrices greatly. Besides, it transformed big data set into serial carry chain computing processes in many subsystems in the computing process, reflected the divide and rule ideology in artificial intelligence field with good practicability and high efficiency. In order to solve the knowledge acquisition problem in big data and incremental data, Shan et al. [74] and Wang Yaying[75] proposed an incremental knowledge acquisition inductive learning algorithm, which to some extent adapted big datas knowledge acquisition problem. But Mapreduce-Algorithm [76] proposed by google pushed big data process into application time.

MapReduce has strong data processing capabilities, if MapReduce technology can be applied to real-time big data processing, undoubtedly it will brought the dawn to real-time big data processing. However, in MapReduce cluster, data is stored in the form of files on each node, while in real-time big data system, data comes from different heterogeneous terminals, and is real-time transmitted. On the other hand, data matching problem in real-time data caused by data source terminals heterogeneity is very serious, and in MapReduce key-value pairs are relative stable. In order to solve these problems, paper [77][77] proposed an improved MapReduce model fitted for real-time big data, as is shown in Figure 11. Paper [78] analyzed the challenges of big data processing in mobile phone clients, and built a big data processing scheme on mobile data collecting system LDCC (Lausanne Data Collection Campaign).

According to real-time requirements of data processing in RTDP model, analyze of data should be divided into local processing and cloud processing. Among them the local server dispose data collected in data acquisition port, it mainly does basic operation such as data cleansing and data structured analysis. While the cloud server conduct big data analysis and process, offer

technical support to decision-making. Data Collection Terminal firstly preprocess data and then submit it to server and deploy the corresponding processing program to local server. The local server then does the Map operation according to data format, then data in heterogeneous nodes will be mapped in different servers, thus avoid data matching problem. The server nodes do the Reduce operation to the data that have been mapped in local servers, then return results to the output terminal. In order to guarantee veracity in data processing, the system supports rolling back action when abnormal thing happens.

Because data processing on local servers are classified according to different types, real-time data processing can be supported in RTDP model. Local computed results update to cloud servers, and make up computed results generated by MapReduce. Thus not only take advantage of the great computing power of cloud computing, but also guaranteed the instantaneity of data processing. For example, in medical wisdom, sphygmomanometers and other instruments can be used to make primary diagnosis, then data can be uploaded to medical service system, on one hand primary diagnosis results can provide a reference to doctors for further diagnose, on the other hand, the data amount uploaded can be reduced after processed by the terminals, so narrower bandwidth shall be needed and therefor provide convenient to data transmission.

### 5.2.2 Calculation Implementation Method

In previous chapters, a two layers calculation mode has been proposed, first, the local server choose local node management and calculation procedures on the local node management, and simple data cleansing and structured modeling according to LMA. Unstructured data collected by data collector will be transformed into structured data and then uploaded to cloud memory systems and mapped to different management servers. Superstardom makes use of the computing power of cloud terminal to carry through real-time computation and analyze.

In recent years, numerous scholars have done a lot of research to the upper layer processing algorithm [77], and have made great achievements, but the problem how to make real-time processing aimed at the bottom layer and how to combine calculation modes has been unsolved [57]. This thesis takes advantage of LMA to do data cleansing as well as managing bottom layer data collection sensors, in the upper layer, an improved MapReduce will be adopted to do large-scale analysis and process. Such a system architecture to some extent, improve the system's flexibility and heterogeneity, and has a good treatment efficiency. For example, in using an FPGA for video encoding and decoding processing can be used during the data collection terminal DSMS for analysis and filtering, in this mode corresponds to an FPGA acquisition DSMS in a Node, DSMS FPGA

collected a large number of data flows flow cytometric analysis performed at the same time to upload data to the cloud server, cloud server will each FPGA configuration tasks performed Reduce Node operation, the calculation results compiled through the LMA choose the right way back to the data collector.

The platforms isomerism and the nodes mutual ability can be greatly improved through LMAs logic management to the gather port. For example, FPGA are used in streaming data processing systems, compared with ASIC, FPGA is slower and cant afford complex designs. But FPGA has its own advantages, it has lower power dissipation, can be quickly finished product, can be modified to correct errors in the program and has cheaper cost. Otherwise, in most cases, real-time big data processing scenes need the chip to be fast programmed, and data collection terminals sensor nodes dont need high processing ability, therefor, FPGA should be enough for the earlier stage, when a scene is relatively mature one application can easily be transferred from FPGA chip to the ASIC chip, for example, the smart power grids real-time big data processing model. Besides, FPGA facilitate changes in the program can quickly modify some subtle errors, to facilitate real-time data processing system in a large process control during some dynamic adjustment and simple error handling. At the same time, FPGAs low power dissipation makes it suitable for real-time big data processing, for sensors in the system are deployed in various environments, low power dissipation has a great influence on the systems robustness and safety. For example, in the smart grid system, data acquisition nodes tend to be widely distributed across the country, its energy consumption is enormous, and the environment is also facing the complex, and in some extreme circumstances, data collection terminals is important for system security effects, such as iron and steel smelting, chemical manufacturing sites, high energy consumption equipment may cause some disasters and have an important influence on system security and environmental security.

### 5.2.3 Complex Event Processing

After determining the mode of calculation methods we also have to realize the model task allocation, as well as prioritization and other operations to make detailed and accurate control, so it relates to a number of complex event processing issues that is, how efficiently a plurality of basic events complex compound has a more complex semantic events, including consideration of constraints between events, and even in some applications to continue to detect complex event to generate a higher level of complex event[79]. Some active database related work discussed for the basic model of the composite event detection [67], including: a model based on finite automata, Petri net-based model, based on matching tree model and a model based on directed graph of these CEP

model is the basic model of the problem, compared with other methods CEP technology has a tense, relevance, semantic richness, heterogeneity and mass and other characteristics [68], here we will have a brief introduction to some the basic models of CEP issues.

(1) Automaton model. Due to the simple expression of complex events with regular expressions have a similar form, and spatial interactions and events have a causal relationship with the local power grid model, with a strong ability of temporal evolution [80], so you can use the automation model to implement event expression. ODE [81] was the first to propose using automatic machine model composite event detection system. Composite event in any one basic event arrives, the automatic machine will transition from a state to the next state, when the automaton enters an acceptable state, then the composite event has occurred because of the simple automata model is not reversible, some of the basic events in the match had not re-visit after the event, so if event and time to consider the link between the values, the need to introduce additional data structure to hold the time information , then extended to form the automaton model. Additionally, the automatic machine during the transition can be added in the transition predicate more complex numerical limits on the time or conditions to design some special automaton model for the application of certain situations.

(2) Petri net model. Petri gateway Note interval endpoints because the calculation and reasoning, so Key Petri net and testing complex event which represents basic event input position, the output location represents the composite event, the event represents a composite intermediate calculation process by entering the Token calculated Warp guard function that computes Warp is raised up and mark the position of the node, marking the last node in the sequence occurs when the composite event detection mechanism is through its incremental marked Petri net description of the position. Active database system SAMOS[82] and monitoring systems HiFi [83] both used this model.

(3) Matching tree model. Based on tree matching technology is mainly by matching tree[84] of the structure to achieve complex event filtering, basic event as matching tree leaf node levels as a composite event matching intermediate node tree root corresponding nodes are filtered out of the composite of the matching tree root means to achieve a complex event detection. READY [85] and Yeast [86] systems use this technique.

(4) Directed graph model. Directed graph model is similar to matching tree model, a directed graph model uses a directed acyclic graph (DAG) represents the composite event. Nodes are used to describe events, edges represent the synthesis of the event rules. Penetration zero for the event input node, the node is zero the output of said composite of intermediate nodes for each level of composite events event tag node can also be simply described composition rules, the node event occurs, node
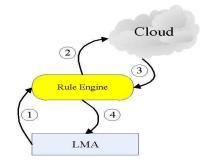
**Fig. 12:** Application Deployment Process

rule is triggered. Sentinel system [87] and the EVE system [88] used this model.

According to the above analysis, based on automata and Petri Nets composite event detection is only matched by event order of arrival, and tree-based or graph filtering do not consider the basic sequence of events or the timing due to a number of events may occur path, the first event did not occur in the case it is possible to filter a second time, resulting in unnecessary overhead, so in practice have some limitations [89]. And these two methods are is the basic event filtering, and composite event detection as different steps to deal with, without taking into conducting composite event detection but also the need for basic event filtering conditions. Paper [90][90] in CompAS systematic basis, considering the basic events and for the integration of complex event processing by the detection complex event while selectively filtering basic events, reducing the response delay of the composite event detection, so this method could become a complex event processing system for RTDP in the future.

## 5.3 Integration Layer

### 5.3.1 Application Arrangement

The integration layer is a connecting link between the preceding and the following layers in RTDP system, the task of this layer is completed under the deploy of the rule engine, on the one hand algorithms for the data analysis layer provides the necessary libraries and algorithm package, requires the integration phase of the algorithm based on data analysis layer needs to provide appropriate resources, advance scheduling and allocation of resources, including the scheduling algorithm on the other hand also need and data access layer interacts LMA underlying data acquisition devices, while access to the cloud through the rules engine application's compiled code, then copy the codes of application program to LMA, and complete the arrangement of application programs. The processing procedure is as shown in Figure 12.

First, LMA sends bottom layers equipment information to the rule engine, rule engine for local data acquisition device information device can deploy applications based on the way to the cloud application requested resources, relying on cloud computing power real-time application is compiled and compiled Java, .exe and other documents sent to the rules engine, by the rules engine to LMA, LMA will get the final compiled application files deployed to the appropriate node, then the arrangement of application programs shall be completed.

This thesis relies on middle management institution LMA to arrange application programs of the nodes, avoids the disadvantage in Hadoop architecture to install high-capacity visual machines, therefor it will be suitable to manage sensing nodes in heterogeneous platforms. At the same time, relying on powerful cloud computing capabilities for program compilation, avoiding the drawbacks of native compilation is not compatible, is conducive to cross-platform data application deployment. Actually, in real life RTDP systems, data acquisition terminal sensing equipment already with heterogeneous, and the nodes computing ability is relatively low.

### 5.3.2 Network Convergence

In real-time big data processing system, real-time data processing mechanism is the integration layer and decision-making of the decision, which is the decision-making system administrators and other decision-makers. Some big data analysis system in real-time decision-making phase and data collection phase systems using the same hardware, but it is different data systems in which a data processing method to the data from the data mart layer varies. this leads to the bottom of the data collection and storage of the speed and processing speed upper mismatch problem to solve this problem from both aspects to be considered, on the one hand is the system architecture and processing algorithms, the other is the data transmission method, as already discussed RTDP architecture model and CEP, DSMS other modes respectively FPGA, GPU, ASIC and other technology solutions RTDP system data processing problems, in this chapter we discussed the solutions of data transmission problems.

For RTDP systems instantaneity requirements, making a RTDP output transmission system is an important part of data acquisition and data processing is a combination of important link is raised to the processing based on data of equal importance. RTDP system because of the presence of iso-node, in order to ensure data we need to revolutionize the traditional real-time data transmission mode[91] how to use a variety of techniques to achieve these heterogeneous nodes in the network real-time, safety and reliability, there are still many scientific problems to be solved[91][92]. Article from heterogeneous network integration and networking for

RTDP new network technologies and high network reliability and QoS three aspects outlined.

### 5.3.3 Heterogeneous network convergence and networking mode

In RTDP system, the network architecture and networking systems that affect system availability and efficiency of the key factors. Developed rapidly in recent years, mobile AdHoc networks, mesh networks, sensor networks and other new network technology, combined with traditional wired networks, cellular networks, in order to build a large real-time data transmission network to provide the foundation for building real-time transmission network due when not completely abandon the existing legacy network infrastructure, next-generation networks must be a mixture of a variety of network technologies in complex networks. In the future each one has a physical component network module should be able to at any time, any place convenient access to the network.

M2M (Machine-to-Machine) network is the current Internet extension and development of a new trend [93][94][95]. In order to achieve anytime, anywhere connectivity, M2M community focusing on how the device through a variety of wireless access technologies (PAN/LAN/WAN) access network and interconnected [96]. RTDP system node characteristics of autonomy and automation requirements RTDP part of the network should have flexible access and ad hoc networking abilities, the ability to easily access and exit the network anytime, anywhere As mobile devices M2M network supports massive multi-level multi -scale networking, and provides anytime, anywhere and flexible network access and other features, including massive equipment access, high reliability, and enhance access priority, low energy consumption, micro- burst transmission for low mobile or fixed device optimization, monitoring and security, large address space, group control, time control of transmission, the transmission time delay tolerance, one-way data transmission, low latency transmission, sporadic transmission [97][98] based on M2M networks can build a unified, flexible and high-capacity network public platform. Therefore, M2M network may be developed to support the real-time transmission backbone network technology [99].

### 5.3.4 New network technique for RTDP

Now widely used network technology is not designed specifically for RTDP These network technologies are based on the "best effort" thinking, in order to optimize the target point to point connection, the timing of which there are a lot of variability and stochastic behavior, therefore, the real-time high system requirements are often forced to use a dedicated network technologies such as CAN, FlexRay, LIN, MAP and other bus

technology[100][101][102], but these networks are limited geographical area network , while in large RTDP systems, many transmission of signals and control commands are global transport require high reliability , so it is necessary to clarify these current network technology in what may be , and how to use the network in large RTDP addition also need to study suitable RTDP new network architecture and networking.

The latest developments in wireless sensor and actuator networks (wireless sensor/actuator network, WSANs) refers to a group of sensors and actuators interconnected via a wireless medium , and can perform distributed sensing and action network[103]. WSANs can observation of the physical world, data processing, data-based decision-making and perform the appropriate action, is considered to be the next one of the key technologies to build RTDP[104], in constructing RTDP networks play an important role[105].

### 5.3.5 Networks high reliability and QoS

RTDP system through the complex network of large data processing depth, has an important role for decision analysis, but also indirectly for the physical world with a deep and broad impact, and many analysts and other decision-making and implementation of closed-loop control commands are transmitted over the network, so its safety, reliability and quality of service becomes extremely important[100].

Different types of RTDP applications have different QoS requirements, such as telemedicine surgical system, due to the delayed signal generated mm -level errors can cause a fatal accident, when in fact demanding aspects of QoS in the smart grid need for rapid diagnosis of the fault zone and fault recovery, the economic loss and the processing time is directly related to the performance of different applications have different requirements on the network, which are provided on the network QoS expectations are different, QoS parameter set defines methods may also different, so starting from the application requirements, careful study of the QoS parameters to determine the appropriate CPS set various parameters to determine their priorities and values to guide RTDP development of network technology.

Now widely used in a variety of network QoS technology, and not enough to guarantee RTDP real-time, high reliability requirements. RTDP network QoS problems of difficulty comes from the RTDP has the inherent characteristics: RTDP network is a complex, heterogeneous converged networks; RTDP in for massive data processing and RTDP widespread large number of dynamic systems, uncertainties present RTDP network QoS issues yet to be carefully and systematically studied for the realization of RTDP integration of heterogeneous network QoS, also need to address the following problem[106]:

(1) How to get to meet the application's QoS request QoS routing calculations necessary information. (2) How to build to meet the QoS request path. (3) How to maintain the path set up short, requires a unified framework to meet the complex network environment, various types of RTDP application QoS requirements.

## 5.4 Decision Making Layer

In fact, decision making layer includes two parts of concepts, first, the test and update the model, the second is to provide managers for decision making. RTDP system during the process of data processing, with the flow of data, the data at different times with a certain variability, and between data also has a certain relevance. therefore change with time and depth data processing, data analysis layer data model created may not meet the current needs, so we need to keep the data processing while the update data and update the data model to adapt to changes in the data on the other hand, decision support layer is the highest level of RTDP system, the purpose is to carry out data processing related decisions, so the layer must visualize the generated output results in order to provide decision-makers to manage related decision-making activities. Next a function overview will be made about model validation and decision support.

### 5.4.1 Model Validation

In order to guarantee that the model is designed to be entirely correct, with a certain fault tolerance, and is stable, we need to model validation tests repeated. Model validation the basic objective is to ensure that the model is really effective to run the model validation phase requires re-extract new data model in an established and validated after the operation centralized data for comparison. If the model is running correctly, it can be deployed to the production environment. Otherwise, the model needs to be repeated error checking until the model is correct and can be stable operation. Accordingly, this part of the model actually contains updates and bug fixes two aspects.

Model is updated in the data flow process according to the data processing requirements of the model update and adjust, and bug fixes are in the model mismatch occurs under the existing data processing an automated solution is a model bug fixes extremely complex task, requiring accurate positioning of the error, when necessary, also need systems that can perform a certain degree of auto repair. and this process is difficult to use mathematical models to express, it has also become a hot research and difficult[79] in the system there are many rules already developed a number of commercial systems for large data run business rules, such as IBM's ILOG real-time analysis using the R language deployment, so how fast the future RTDP system model validation and
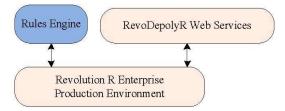


**Fig. 13:** Application Deployment Process

deployment in to some extent, you can refer to existing business systems models and methods, based on the existing expert system rules engine to quickly build automated modeling and analysis, real-time adjustment of the data, analysis, and model deployment.

### 5.4.2 Decision Support

Decision support is the ultimate goal of data analysis, decision support in part requires the use of a large number of visualization tools for data analysis results show different dimensions. Data presentation forms including business intelligence systems, desktop office systems and mobile terminal systems, etc. Use the tool includes data warehousing systems and graphical processing tools.

## 6 Application Examples

And application domain related RTDP research has also got a lot of attention, and have made some preliminary research. IBM in 2010 formally proposed "smart city" vision[107], and since then a large real-time data and research to get smart city widespread attention as urban infrastructure studies intelligent transportation, smart grid, and urban services related to medical wisdom has also been a great development. Various applications have their own unique characteristics, many key issues still unresolved, many of the existing areas most relevant studies also remain in the laboratory phase. An example of RTDP smart grid smart grid system in the application overview will be made below.

Smart grid through a large number of real-time sensor sensing grid operation state changes, can provide faster dynamic evaluation of real-time fault diagnosis and energy tracking[108] In covering provinces, municipalities and even nationwide dynamic regulation of energy to achieve distribution energy production with reasonable dispatch, for improving energy efficiency, improving urban infrastructure plays an important role.

Energy system for the country and the importance of social life determines the smart grid should be highly flexible, autonomous ability, intelligence, scalability, efficiency, predictability, computing security, and many other features. Numerous academics, researchers also
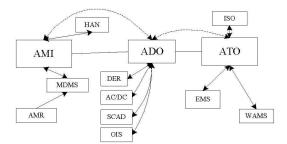
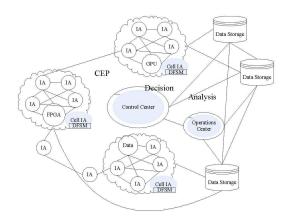**Fig. 14:** Technical Composition and Functionality of Smart Grid



**Fig. 15:** Architecture of Distributed Smart Grid

made a lot of research, the paper [109] proposed signal processing can be affected by interference problem into a model checking problem, and design RT_PROMELA grid disturbance model, using RT_SPIN tool has been verified. Paper [110] for the smart Grid reliability issues studied by the AC transmission equipment for error injection approach the equipment failure mode analysis. Failure of the equipment and on its influence and power grid reliability were assessed. literature indicates that the current construction of smart grid integration of applications and technology solutions, including smart meters, communication networks, metering database management (MDMS), customer premises network (HAN), customer service, remote turn on or off, as shown in Figure 14, smart grid technologies has reached a certain degree of availability and flexibility, but in the power of real-time deployment, management, fault detection and recovery, there are still deficiencies[111].

Current research areas of smart grid applications shortcomings mainly in pharos measurement technology and wide area measurement techniques[108]; dimensional, dynamic visualization dispatching automation technology; access to renewable energy and grid technology[112]; advanced the infrastructure and automated meter reading system[113]; demand response and demand side management technologies; advanced distribution operation technology; distributed generation technologies[114], micro-grid technology and power storage technology and other aspects.

The numerous technical difficulties are all within the scope of the whole network and data related to monitoring and real-time calculation, so this paper that the solution to the current field of smart grid of the key issues, we must rely RTDP processing frameworks for large-scale network-wide real-time data combining the proposed RTDP framework to build a new type of smart grid architecture, shown in Figure 15.

Fast simulation and modeling is the core software of ADO, including risk assessment, self-healing and other advanced control and optimization software system for the smart grid to provide support and predictive mathematical ability, in order to achieve improved grid stability, safety, reliability and operational efficiency. Distribution fast simulation and modeling need to support

network reconfiguration, voltage and reactive power control; fault location, isolation and restoration of electricity; when the system topology changes following the security re-tuning four self-healing capabilities. Above function interconnectedness, resulting DFSM become very complicated, for example, either a grid reconstruction requires a new relay with voltage regulation or the new festival program also includes functions to restore power. DFSM via distributed intelligent network agents to achieve organizational boundaries across geographical boundaries and intelligent control system in order to achieve self-healing capabilities of these intelligent network agents, able to collect and exchange information and systems (such as the following such electrical protection operation) local control decisions, while according to the system requirements to coordinate these programs.

# 7 Conclusion and Future Work

Real-time data processing for large current technology is undoubtedly a huge challenge, there is lack of support for massive real-time processing of large data frame and platform real-time processing of large data processing compared to conventional static data with high data throughput and real-time requirements. cloud computing technology in order to solve massive data processing and developed a series of techniques, however, cloud computing is very suitable for mass static, long-term without the written data has a good effect, but it is difficult to achieve real-time processing.

In this paper, on the basis of cloud computing technology to build a kind of real-time processing of large data frame, the model proposed RTDP four architecture, and hierarchical computing model. RTDP system in order to meet real-time requirements, and to consider different system platforms RTDP structural characteristics, the paper also presents a large data storage for real-time multi-level storage model and the LMA-based application

deployment methods of data collection terminal based on the different ways of data processing were used DSMS, CEP, batch-based MapReduce other processing mode, depending on the environment in which the sensor data acquisition and the desired type of difference data collected were used FPGA, GPU, CPU, ASIC technology to achieve data collection and data cleansing and, through a structured process the data structure modeling, uploaded to the cloud server for storage, while the washed structured data on the local server for Reduce, combined with powerful computing capabilities cloud architecture for large-scale real-time computing with MapReduce.

This thesis indicates generally that the basic framework for future RTDP system and basic processing mode, but there are still many issues that need further study. The main point are as follows:

1. How to determine the appropriate mode of calculation in a RTDP system, how to determine the data processing mode and approach is a key factor in determining system performance , so the calculation mode and how to determine the appropriate method of calculating the design of the future core of the work;

2. Calculation models and how to achieve unity between computing technology is currently used mainly batch calculation mode and streaming processing , data computing model in determining how to design the corresponding calculation after the manner and with what kind of hardware implementation is the next big real-time data processing priority;

3. How to ensure the network transmission speed and QoS (Quality of Services); now widely used in a variety of network QoS technology, RTDP not sufficient to ensure a real-time, high reliability requirements. RTDP network QoS issues RTDP with difficulty from the inherent characteristics, so to guarantee QoS of the real-time RTDP sex have a significant impact;

4. How to ensure the system's physical time synchronization. RTDP system involves many interactions between systems and tools, software used for real-time marker approach does not meet the future RTDP high real-time requirements, the interactive how to ensure data during physical time synchronization is the future research directions;

5. How to ensure the correctness of the data processing. Error detection mechanism and automatically repair the computer has long been a difficult area of research, how to handle the data detection and error diagnostic and system repair is a huge project.

RTDP is a complex project involving many disciplines and techniques to be thorough in all aspects of research, pointed out that the article provides an overview of future research directions, and this is our future research subject.

## Acknowledgements

## Conflicts of Interest

The authors declare that there is no conflict of interests regarding the publication of this paper.

## References

[1] Brown B, Chui M, Manyika J. Are you ready for the era of 'big data[J]. McKinsey Quarterly, 2011, 4: 24-35.

[2] Tene O, Polonetsky J. Big data for all: Privacy and user control in the age of analytics[J]. Nw. J. Tech. & Intell. Prop., 2012, 11: xxvii.

[3] Lohr S. The age of big data[J]. New York Times, 2012, 11.

[4] Lynch C. Big data: How do your data grow? [J]. Nature, 2008, 455(**7209**):2829

[5] Bryant, RE, Katz, RH, Lazowska, ED, Big-Data Computing: Creating revolutionary breakthroughs in commerce, science, and society [R], Ver. 8, (2008), Computing Research Association, Computing Community Consortium.

[6] Jonathan T. Overpeck, Gerald A. Meehl, Sandrine Bony, and David R. Easterling. Climate Data Challenges in the 21st Century [J]. Science, 2011, 331(**6018**):700-702

[7] Labrinidis, Alexandros and Jagadish, H. V. Challenges and opportunities with big data [J]. Proc. VLDB Endow. 2012, 5(**12**): 2032-2033.

[8] WANG Shan, WANG Hui-Ju, QIN Xiong-Pai, ZHOU Xuan. Architecting Big Data: Challenges, Studies and Forecasts [J].Chinese Journal of Computer. 2011, 34(**10**): 1741-1752.

[9] Lu Weixing, Shou Yinbiao, Shi Lianjun. WSCC DISTURBANCE ON AUGUST 10, 1996 IN THE UNITED STATE[J]. Power System Technology.1996, 20(**9**):40-42.

[10] P. Jeffrey Palermo. The August 14, 2003 blackout and its importance to China [J]. EAST CHINA ELECERIC POWER. 2004, 32(**1**): 2-6.

[11] Li Cuiping, Wang Minfeng. Excerpts from the Translation of Challenges and Opportunities with Big Data [J]. e-Science Technology & Application, 2013, 4(**1**):12-18.

[12] Dobbie W, Fryer Jr R G. Getting beneath the veil of effective schools: Evidence from New York City[R]. National Bureau of Economic Research, 2011.

[13] H. V. Jagadish, Johannes Gehrke, Alexandros Labrinidis, Yannis Papakonstantinou, Jignesh M. Patel, Raghu Ramakrishnan, Cyrus Shahabi. Big data and its technical challenges [J]. Communications of the ACM, 2014, 57(**7**):86-94.

[14] Flood M, Jagadish H V, Kyle A, et al. Using Data for Systemic Financial Risk Management[C]. Proceedings of The 5th biennial Conference on Innovative Data Systems Research (CIDR 2011). 2011: 144-147.

[15] Genovese Y, Prentice S. Pattern-based strategy: getting value from big data[J]. Gartner Special Report G, 2011, 214032: 2011.

[16] Albert-Lszl Barabsi. The network takeover. Nature Physics, 2012, 8(**1**): 14-16.

[17] Labrinidis A, Jagadish H V. Challenges and opportunities with big data[J]. Proceedings of the VLDB Endowment, 2012, 5(**12**): 2032-2033.

[18] Lohr S. How big data became so big[J]. New York Times, 2012, 11.

[19] Gattiker A, Gebara F H, Hofstee H P, et al. Big Data text-oriented benchmark creation for Hadoop[J]. IBM Journal of Research and Development, 2013, 57(**3/4**): 10: 1-10: 6.

[20] Chen M, Mao S, Liu Y. Big data: A survey[J]. Mobile Networks and Applications, 2014, 19(**2**): 171-209.

[21] Li Guojie, Cheng Xueqi. Research Status and Scientific Thinking of Big Data [J]. Bulletin of the Chinese Academy of Sciences. 2012, 27(**6**): 647-657.

[22] Yadagiri S, Thalluri P V S. Information technology on surge: information literacy on demand[J]. DESIDOC Journal of Library & Information Technology, 2011, 32(**1**):64-69.

[23] Cohen J, Dolan B, Dunlap M, Hellerstein JM, Welton C. MAD skills: New analysis practices for big data [J]. PVLDB, 2009,2(**2**):14811492.

[24] Randal E. Bryant & Joan Digney. Data-Intensive Supercomputing: The case for DISC [R].2007.10: 1-14.

[25] John Boyle. Biology must develop its own big-data systems. Nature. 2008, 499(**7**): 7.

[26] Wang Yuan-Zhuo, Jin Xiao-Long, Chen Xue-Qi. Network Big Data: Present and Future [J].Chinese Journal of Computer. 2013, 36(**6**):1125-1138.

[27] Zou Guowei, Cheng Jianbo. TheApplication of BigDataTechnologyto SmartCity [J]. POWER SYSTEM TECHNOLOGY, 2013, 4: 25-28.

[28] QIN Xiong-Pai, WANG Hui-Ju, DU Xiao-Yong, WANG Shan. Big Data Analysis-Competition and Symbiosis of RDBMS and MapReduce [J]. Journal of Software. 2012, 23(**1**): 32-45.

[29] Tan Xiongpai, Wang Huiju, Li Furong, et al. New Landscape of Data Management Technologies [J]. Journal of Software. 2013, 24(**2**): 175-197.

[30] CHEN Hai-Ming, CUI Li, XIE Kai-Bin. A Comparative Study on Architectures and Implementation Methodologies of Internet of Things [J]. Chinese Journal of Computers. 2013, 36(**1**): 168-188.

[31] Lee E A, Seshia S A. Introduction to embedded systems: A cyber-physical systems approach[M]. Lee & Seshia, 2011.

[32] Thusoo A, Sarma J S, Jain n, et al. Hive-A petabyte scale data warehouse using Hadoop [C]. Proc. of ICDE 2010. Piscataway, NJ: IEEE, 2010: 996-1005

[33] Abouzied A, Bajda-Pawlikowski K, Huang Jiewen, et al. HadoopDB in action: Building real world applications [C]. Proc. of SIGMOD 2010, New York: ACM, 2010: 1111-1114.

[34] Chen Songting. Cheetah: A high performance, custom data warehouse on top of MapReduce [J]. PVLDB, 2010, 3(**2**): 1459-1468.

[35] Agrawal R, Srikant R. Privacy preserving data mining [C]. Proc. of SIGMOD 2000. New York: ACM, 2000: 439-450.

[36] Dwork C. Differential privacy [C]. Proc. of ICALP 2006. Berlin: Springer, 2006: 1-12.

[37] Norman D A. The Design of Everyday Things [M].New York: Basic Books. 2002.

[38] Olston C, Reed B, Srivastava U, et al. Pig Latin: A not-so-foreign language for data processing [C]. Proc of SIGMOD 2008, New York:ACM, 2008:1099-1110.

[39] Pike R, Dorward S, Griesemter R, et al. Interpreting the data: Parallel analysis with Sawzall [J]. Scientific Programming, 2005, 13(**4**): 277-298.

[40] Chaiken R, Jenkins B, Larson P-A, et al. SCOPE: Easy and efficient parallel processing of massive data sets [J]. PVLDB, 2008, 1(**2**): 1265-1276.

[41] Isard M, Yu Y. Distributed data-parallel computing using a high-level programming language [C]. Proc. of SIGMOD 2009. New York: ACM, 2009: 987-994.

[42] Fegaras L, Li C, Gupta U, et al. XML query optimization in MapReduce [C]. Proc. of WebDB 2011. New York: ACM, 2011.

[43] Morton K, Balazinska M, Grossman D. Para Timer: A progress indicator for MapReduce DAGs [C]. Proc. of SIGMOD 2010. New York: ACM, 2010: 507-518.

[44] Morton K, Friesen A, Balazinka A, et al. KAMD: Estimating the progress of MapReduce pipelines [C]. Proc. of ICDE 2010. Piscataway, NJ: IEEE, 2010: 681-684.

[45] Huang Dachuan, Shi Xuanhua, Ibrahim Shadi, et al. MR-scope: A real-time tracing tool for MapReduce [C]. Proc. of HPDC 2010. New York: ACM, 2010: 849-855.

[46] Meng Xiaofeng, Ci Xiang. Big Data Management: Concepts, Techniques and Challenges [J]. Journal of Computer Research and Development. 2013, 50(**1**): 146-169.

[47] Chen Y. We dont know enough to make a big data benchmark suite-an academia-industry view[C]. Proc. of WBDB, 2012.

[48] Chen Yanpei, Ganapathi A, Griffith R, et al. The case for evaluating MapReduce performance using workload suites [C]. Proc. of MASCOTS 2011. Piscataway, NJ: IEEE, 2011: 390-399.

[49] Sangroya A, Serrano D, Bouchenak S. Mrbs: A comprehensive mapreduce benchmark suite[R]. LIG, Grenoble, France, Research Report RR-LIG-024, 2012.

[50] Tan J, Kavulya S, Gandhi R, et al. Light-weight black-box failure detection for distributed systems[C]. Proceedings of the 2012 workshop on Management of big data systems. ACM, 2012: 13-18.

[51] Zhao J M, Wang W S, Liu X, et al. Big data benchmark-big DS[M]. Advancing Big Data Benchmarks. Springer International Publishing, 2014: 49-57.

[52] Patil S, Polte M, Ren K, et al. YCSB++: Benchmarking and performance debugging advance features in scalable table stores [C]. Proc. of SoCC 2011. New York: ACM, 2011.

[53] Chen Y, Alspaugh S, Katz R. Interactive query processing in big data systems: A cross-industry study of MapReduce workloads [J]. PVLDB, 2012, 5(**12**): 1802-1813.

[54] Liu Peng. Cloud Computing [M]. Beijing: Electronic Industry Press, 2011.

[55] Fingar P. Dot Cloud: The 21st Century Business Platform Built on Cloud Computing. Beijing: Electronic Industry Press.2010.

[56] Dean J, Ghemawat S. MapReduce: Simplified data processing on large clusters [C]. Proc. of OSDI 2004. Berkeley, CA: USEUIX Association, 2004: 137-150.

[57] Patrick Mellen Wendell. Scalable Scheduling for Sub-Second Parallel Jobs [D]. University of California at Berkeley.2013.3.

[58] Xiongpai Qin, Biao Qin, Xiaoyong Du, Shan Wang. Reflection on the Popularity of MapReduce and Observation

of Its Position in a Unified Big Data Platform [J]. Web-Age Information Management Lecture Notes in Computer Science. 2013, 7901: 339-347.

[59] JeongJin Cheon, Tae-Young Choe. Distributed Processing of Snort Alert Log using Hadoop [J]. International Journal of Engineering and Technology.2013,5(**3**):2685-2690.

[60] Chang R M, Kauffman R J, Kwon Y O. Understanding the paradigm shift to computational social science in the presence of big data[J]. Decision Support Systems, 2014, 63: 67-80.

[61] Chen J, Chen Y, Du X, et al. Big data challenge: a data management perspective[J]. Frontiers of Computer Science, 2013, 7(**2**): 157-164.

[62] Yuan ZHAO, Zhuo-Feng, FANG Jun, MA Qiang. Real-Time Processing for High Speed Data Stream over Large Scale Data[J]. Chinese Journal of Computers. 2012.35(**3**). 477-490.

[63] Dean J, Ghemawat S. MapReduce: Simplified data processing on large clusters. ACM Communication, 2008, 51(**1**):107-113.

[64] Ranger C, Raghuraman R, Penmestsa A, Bradski G, Kozyrakis C. Evaluating MapReduce for multi-core and multiprocessor systems [C]. Proc. of the 13th International Conference on High-Performance Computer Architecture (HPCA 2007). Phoenix, USA, 2007: 13-24.

[65] Shirin Mohammadi, Ali A. Safaei, Fatemeh Abdi, Mostafa S. Haghjoo. Adaptive data stream management system using learning automata [J]. Advanced Computing: An International Journal. 2011, 2(**5**): 1-14.

[66] Welsh M, Culler D, Eric Brewer E. SEDA: architecture for well-conditioned, scalable Internet services [C]. Proc. of the 18th ACM Symposium on Operating System Principles (SOSP 2001). Lake Louise, Banff, Canada, 2001:230-243.

[67] ZANG Chuanzhen, FAN Yushun. COMPLEX EVENT PROCESSING OF REAL TIME ENTERPRISES BASED ON SMART ITEMS [J]. CHINESE JOURNAL OF MECHANICAL ENGINEERING.2007,42(**2**):22-32.

[68] GU Yu, YU Ge, ZHANG Tiancheng. RFID complex event processing techniques. Journal of Computer Science and Frontiers,2007,1(**3**):255- 267.

[69] McAfee A, Brynjolfsson E, Davenport T H, et al. Big data[J]. The management revolution. Harvard Bus Rev, 2012, 90(**10**): 61-67.

[70] Kitchin R. The real-time city? Big data and smart urbanism[J]. GeoJournal, 2014, 79(**1**): 1-14.

[71] Milke Barlow. Real-Time Big Data Analytics: Emerging Architecture [M]. O'Reilly. 2013, 2. First Edition.

[72] Kshetri N. Big data s impact on privacy, security and consumer welfare[J]. Telecommunications Policy, 2014, 38(**11**): 1134-1145.

[73] CHENG Yu-sheng, ZHANG You-sheng, HU Xue-gang. Matrix blocks computation with serial carry chain for rule extraction in massive data sets [J]. JOURNAL OF UNIVERSITY OF SCIENCE AND TECHNOLOGY OF CHINA. 2009,39(**2**):196-203.

[74] Shan N, Ziarko W. An incremental learning algorithm for constructing decision rules[C]. Proceedings of the RSKDp93. London: Springer-Verlag,1993: 326-2334.

[75] WANG Ya-ying, SHAO Hui-he. A kind of two types of incremental inductive learning decision system [J]. Information and Control, 2002 , 29(**6**): 521-525.

[76] Jeffrey Dean, Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters [C]. OSDI'04:

[77] Jing Liu et al. An open, flexible and multilevel data storing and processing platform for very large scale sensor network [C]. Proc. of ICACT 2012. Seoul Korea: IEEE, Feb. 2012: 926-930

[78] Juha K. Laurila, Daniel Gatica-Perez, Imad Aad, etal. The Mobile Data Challenge: Big Data for Mobile Computing Research [C]. Proceedings of the Workshop on the Nokia Mobile Data Challenge, in Conjunction with the 10th International Conference on Pervasive Computing. 2012.1-8.

[79] DING Jian, BAI Xiao-min, ZHAO Wei, FANG Zhu, LI Zai-hua, ZHONG Wu-zhi. Fault Information Analysis and Diagnosis Method of Power System Based on Complex Event Processing Technology. 2007,27(**28**):40-45.

[80] Ke Chang-Qing, Ouyang Xiao-Ying. The Advances in Modeling Urban Spatial Change Based on Cellular Automata [J]. JOURNAL OF NANJING UNIVERSITY (NATURAL SCIENCES).2006, 42(**1**):103-110.

[81] Gehani N, Jagadish H. Ode as an active database: constraints and triggers[C]. Proc of VLDB, 1991:327- 336.

[82] Gatziu S, Dittrich K R. SAMOS: an active object-oriented database system [J]. IEEE Bulletin on Data Engineering, 1992, 15(**4**):23-26.

[83] Rizvi S, Jeffrey S, Krishnamurthy S, et al. Events on the edge[C]. Proc of SIGMOD, 2005:885- 887

[84] WANG Jin-Ling, JIN Bei-Hong, LI Jing, SHAO Dan-Hua. Data Model and Matching Algorithm in an Ontology-Based Publish/Subscribe System [J]. Journal of Software. 2005, 16(**9**):1625-1635.

[85] Gruber R E, Krishnamurthy B, Panagos E. The architecture of the READY event notification service[C]. Proc of the 19th IEEE Int Conf on Distributed Computing Systems Middleware Workshop, 1999.

[86] Krishnamurthy B, Rosenblum D S. Yeast: a general purpose event-action system [J]. IEEE Transactions on Software Engineering,1995,21(**10**):845-857

[87] Chakravarthy S, Mishra D. Snoop: an expressive event specification language for active databases, UF-CISTR-93-007[R]. University of Florida, Gainesville,1993

[88] Geppert A, Tombros D. Event-based distributed workflow execution with EVE, ifi-96.05[R]. University Zurich, 1996

[89] LUO Hai-bin, FAN Yu-shun, WU Cheng. Overview of Workflow Technology [J]. Journal of Software. 2000, 11(**7**):899-907.

[90] Hinze A. Efficient filtering of composite events [C]. Proc of the British National Database Conf, 2003:207- 225.

[91] Hu Yafei, Li Fangmin, Liu Xinhua. CPS: Network System Framework and Key Technologies [J]. JOURNAL OF COMPUTER RESEARCH AND DEVELOPMENT. 2010,47(**z2**): 304-311.

[92] Koubaa A, Andersson B. A vision of cyber-physical Internet 2009 [c]//Proceeding of the 8th International Conference on Real-Time NetworkPiseataway, NJIEEE, 200925-30.

[93] Saxon L A. Ubiquitous wireless ECG recording: a powerful tool physicians should embrace[J]. Journal of cardiovascular electrophysiology, 2013, 24(**4**): 480-483.

[94] S. Whitehead. Ado p ting Wireless Machine-to-Machine Technology [J]. IEEE Computing and Control Eng. J., 2004, 15(**5**):40-46.

[95] G. Lawton. Machine-to-machine technology gears up for growth [J]. IEEE Transactions on Computer, 2004, 37(9):12-15.

[96] Geng Wu, Shilpa Talwar, Kerstin Johnsson, Nageen Himayat, and Kevin D. Johnson. M2M: From Mobile to Embedded Internet [J]. IEEE Communications Magazine, 2011,4: 36-43.

[97] Taleb T, Kunz A. Machine type communications in 3GPP networks: potential, challenges, and solutions[J]. Communications Magazine, IEEE, 2012, 50(**3**): 178-184.

[98] Chang K, Soong A, Tseng M, et al. Global wireless machine-to-machine standardization[J]. IEEE Internet Computing, 2011(**2**): 64-69.

[99] Chen K C, Lien S Y. Machine-to-machine communications: Technologies and challenges[J]. Ad Hoc Networks, 2014, 18: 3-23.

[100] Davis R I, Kollmann S, Pollex V, et al. Schedulability analysis for Controller Area Network (CAN) with FIFO queues priority queues and gateways[J]. Real-Time Systems, 2013, 49(**1**): 73-116.

[101] Wang C C, Chen C L, Hou Z Y, et al. A 60 V Tolerance Transceiver With ESD Protection for FlexRay-Based Communication Systems[J]. Circuits and Systems I: Regular Papers, IEEE Transactions on, 2015, 62(**3**): 752-760.

[102] Rogers F A, Lin S S, Hegan D C, et al. Targeted gene modification of hematopoietic progenitor cells in mice following systemic administration of a PNA-peptide conjugate[J]. Molecular Therapy, 2012, 20(**1**): 109-118.

[103] T.W. Carley, M.A. Ba, R. Barua, D.B. Stewart, Contention-free periodic message scheduler medium access control in wireless sensor/actuator networks, in: Proc. of Real-Time Systems Symposium, Cancun, Mexico, December 2003.

[104] Park K J, Zheng R, Liu X. Cyber-physical systems: Milestones and research challenges[J]. Computer Communications, 2012, 36(**1**): 1-7.

[105] Wu Fangjing, Kao Yufen, Tsong, Yuchee. From wireless sensor networks towards cyber-physical systems[J]. Journal of Pervasive and Mobile Computing, 2011, 7(**4**), 397-413

[106] Tommaso Melodia, Dario Pompili, Vehbi C. Gungor, and Ian F. Akyildiz. Communication and Coordination in Wireless Sensor and Actor Networks. IEEE TRANSACTIONS ON MOBILE COMPUTING, 2007, 6(**10**):1116-1129

[107] Eckman, B.; Feblowitz, M.D.; Mayer, A.; Riabov, A.V., Toward an integrative software infrastructure for water management in the Smarter Planet [J]. IBM Journal of Research and Development, 2010, 54(**4**):1-20.

[108] CHEN Shu-yong, SONG Shu-fang, LI Lan-xin, SHEN Jie. Survey on Smart Grid Technology [J]. Power System Technology.2009,33(**8**):1-7.

[109] Sun Y, McMillin B M, Liu X, et al. Verifying noninterference in cyber-physical systems the advanced electrical power grid [C]. Proc of QSIC. Piscataway,NJ: IEEE,2007:363-369.

[110] Mulligan G. The 6LoWPAN architecture[C]. Proceeding of ACM EmNets. New York, ACM, 2007, 78-92.

[111] HU Xue-Hao. Smart Grid-A Development Trend of Future Power Grid [J].Power System Technology. 2009,33(**14**):1-5

[112] Mark McGranaghan, et al. Renewable systems Interconnection Study: Advanced grid planning and operation[R].2008.

[113] YU Yi-xin, LUAN Wen-peng. Smart Grid [J]. POWER SYSTEM AND CLEAN ENERGY. 2009,25(**1**):7-11

[114] YuYi-Xin. Technical Composition of Smart Grid and its Implementation Sequence [J]. SOUTHERN POWER SYSTEM TECHNOLOGY. 2009, 3(**2**):1-5

**Zhigao    Zheng** (M12-M12)    was    born in    Muzidian    Village, Macheng City, in 1988. He received the M.eng degree in Information Science and Engineering from Peking University, Beijing, China, 2014. From July 2014 to now, he was a Research Assistant with the National Engineering Research Center for E-learning and Collaborative & Innovative Center for Educational Technology at Central China Normal University. He is the author of more than five articles. His research interests include distributed data stream analysis, cloud computing and Big Data. He became a Member (M) of CCF in 2012, and Member (M) of ACM in 2012.

**Ping    Wang** (M90) is a professor with National Engineering Research Center for Software Engineering and School of Software and Microelectronics at Peking University, China. He is the Director of intelligent Computing and Sensing Laboratory (iCSL), Peking University and the Co-Director of PKU-PSU Joint Intelligent Computing and Sensing Laboratory (PKU side). He received his doctorate degree in Computer Science from the University of Massachusetts in 1996. His research interests include intelligent computing and its applications, Internet of Things, and related information and network security issues. Dr. Wang has authored or co-authored over 30 papers in journals or proceeding such as IEEE Transactions on Dependable and Secure Computing, Ad Hoc Networks, Computer Networks, IEEE International Conference on Web Service (ICWS), IEEE Wireless Communications and Networking Conference (WCNC), and Communications of the ACM. He was co-recipient of one Best Paper and two Outstanding Paper Awards from different professional societies.

**Jing Liu** is a lecturer at School of Software and Microelectronics, Peking University, His research interests include Internet of Things and Cyber-Physical Systems (CPS).



**Shengli Sun** was born in Hunan Province, China, in 1979. He received his Ph.D. degree in computer science from Fudan University, Shanghai, China, in 2008. He was with IBM China Development Laboratory, focusing on automatic software testing as a visiting student from June, 2003 to June, 2004. He is currently an Associate Professor in the School of Software and Microelectronics, Peking University. His research interests include data mining and knowledge discovery, query optimization and data warehouse. Since 2006, Prof. Sun has published over 20 papers in various journals and conference proceedings. He is supported by The Natural Science Foundation of Jiangsu Province under Grant No. BK2010139.