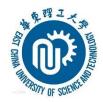
Building and Exploring an Enterprise Knowledge Graph for Investment Analysis





Tong Ruan¹, Lijuan Xue¹, Haofen Wang¹ Fanghuai Hu², Liang Zhao¹, Jun Ding²

¹ East China University of Science and Technology ² Shanghai Hi-knowledge Information Technology Corporation

Business Motivation

- In China, most securities companies provide investment bank services and investment consulting services.
- Establishment of the New Third Board, a national share transfer system for small- and medium-sized enterprises (SMEs)
- Small innovation companies can be listed on the "New Third Board" with the endorsement of securities companies.

Business Motivation

- Securities companies serve from big enterprises to small and medium-sized enterprises.
- There are about 40M companies in China. It is difficult for the securities companies themselves to gather authentic and full-fledged company information of their customers and potential customers.
- Therefore we collect company information from different sources for them and represent it in easy-to-use graphs. The "Magic Mirror" targets to help the securities companies to know and to approach their target companies better and quicker.

✓ Size: About two hundred million entities, one billion attribute value pairs, and two hundred million relations in EKG.

✓Time: It takes an hour to extract entities, three hours to extract attribute value pairs, and three hours to extract relations from various sources.

 ✓ Update: rebuilt once a month to incorporate newly added enterprise data. (Simple and require improvement)



Deployment and Business Model

✓ Sell the whole solution as services instead of software. Securities companies have customized the EKG portal and have integrated it into their own applications.

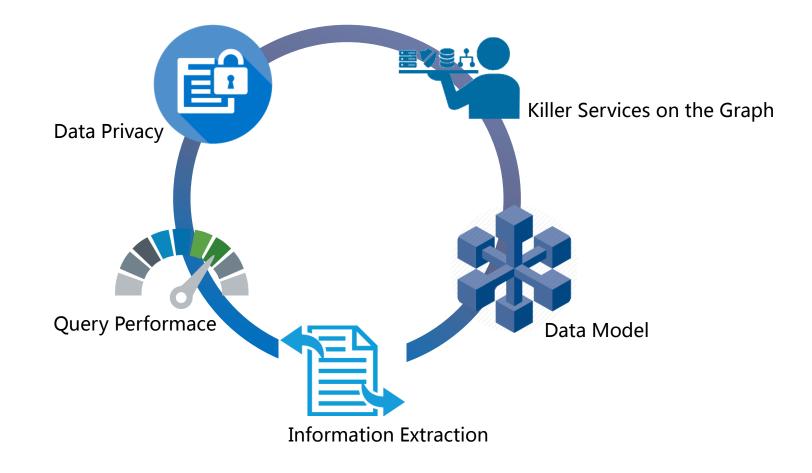
✓ Pay by times of API access Per Year

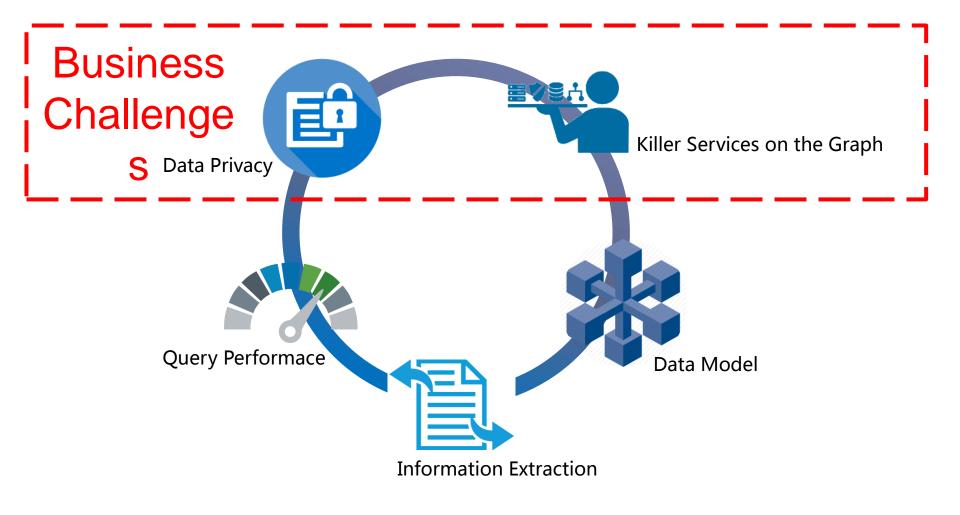
✓ General querying and graph visualization services

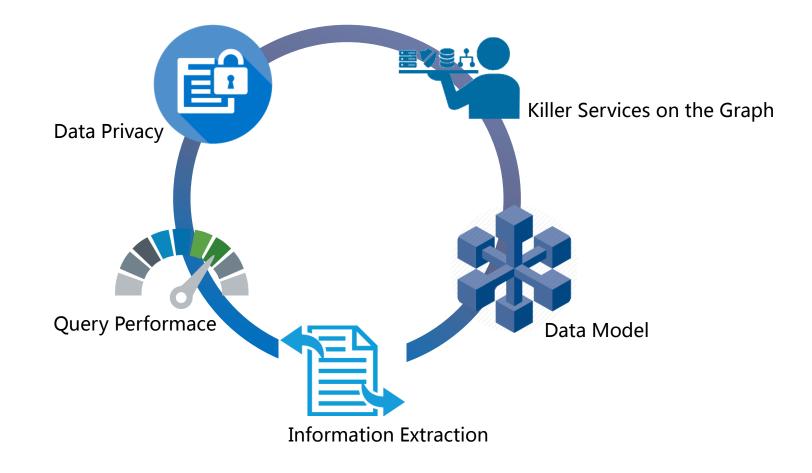
✓ In-depth analyzing services dedicated to investing requirements

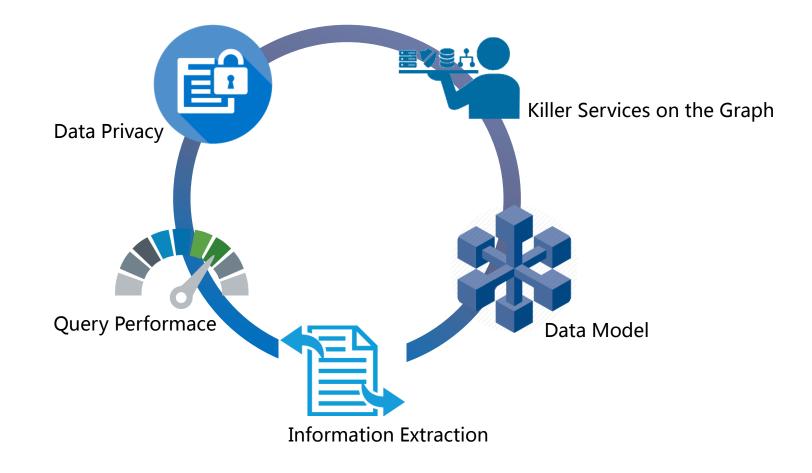
Business Model

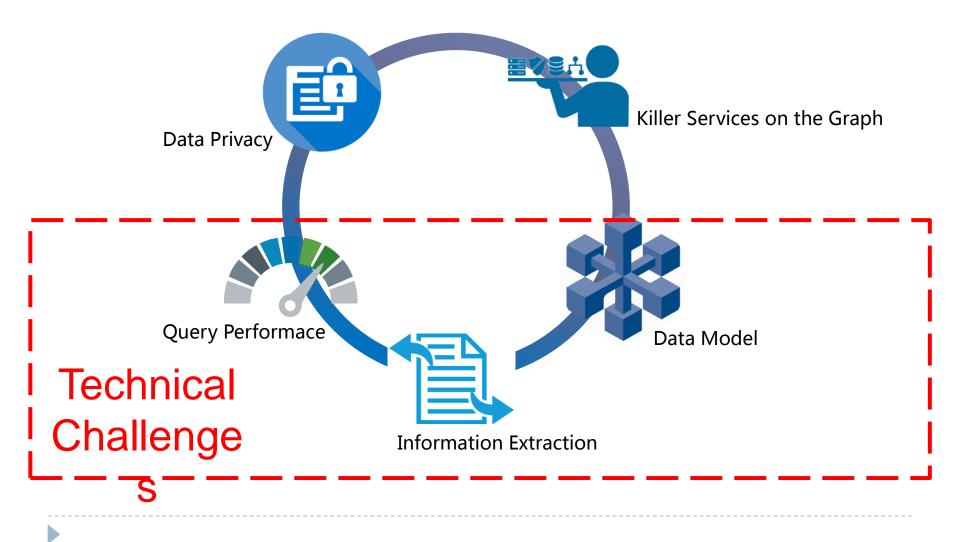


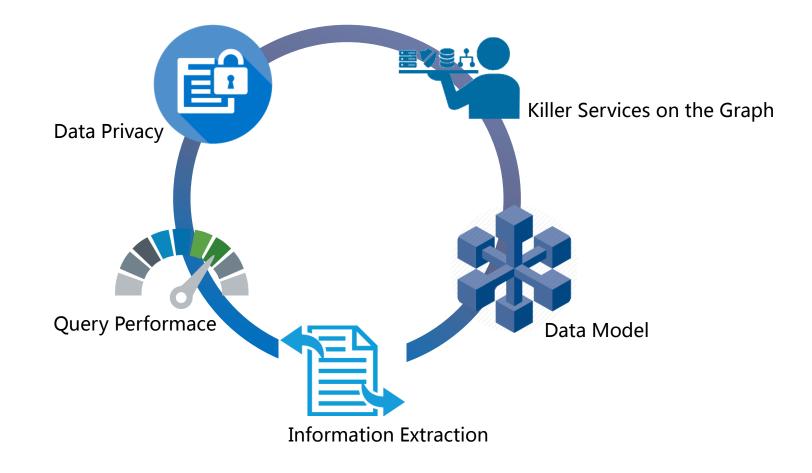


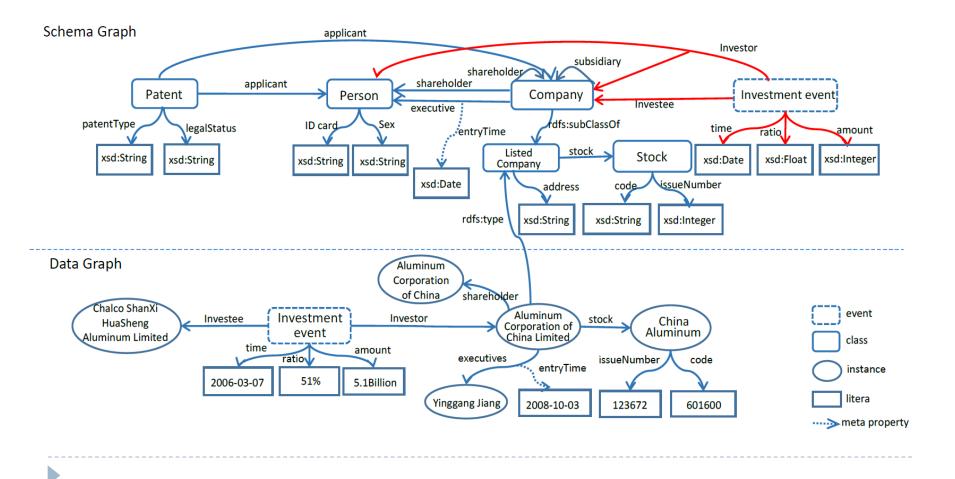


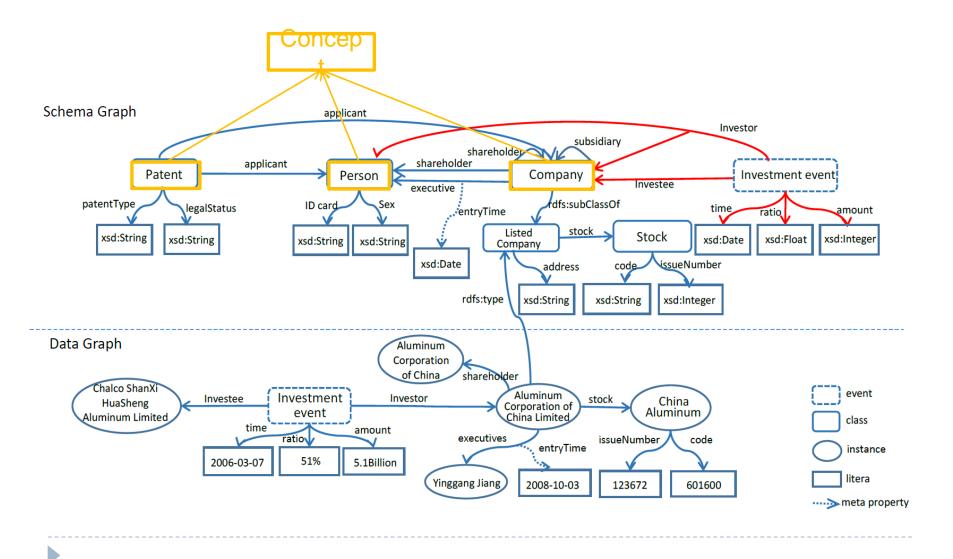


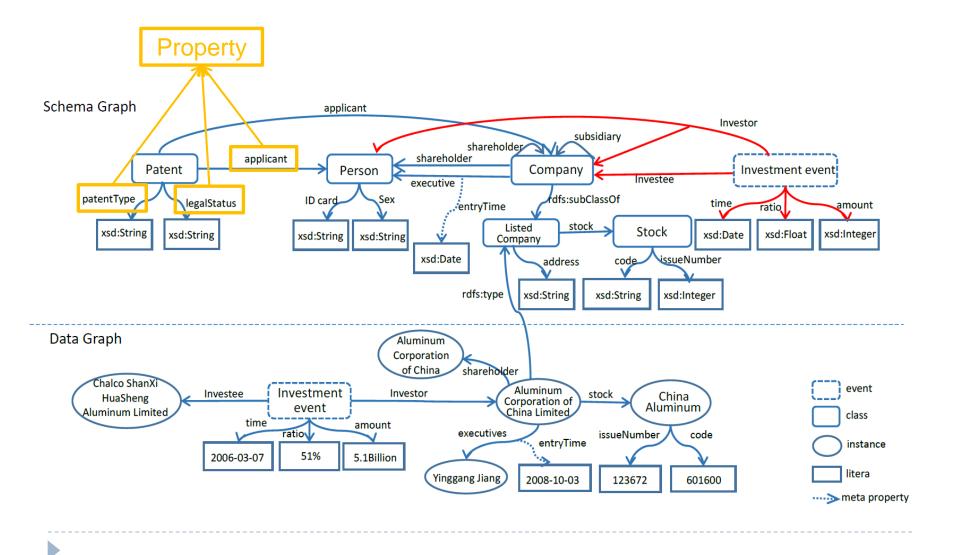


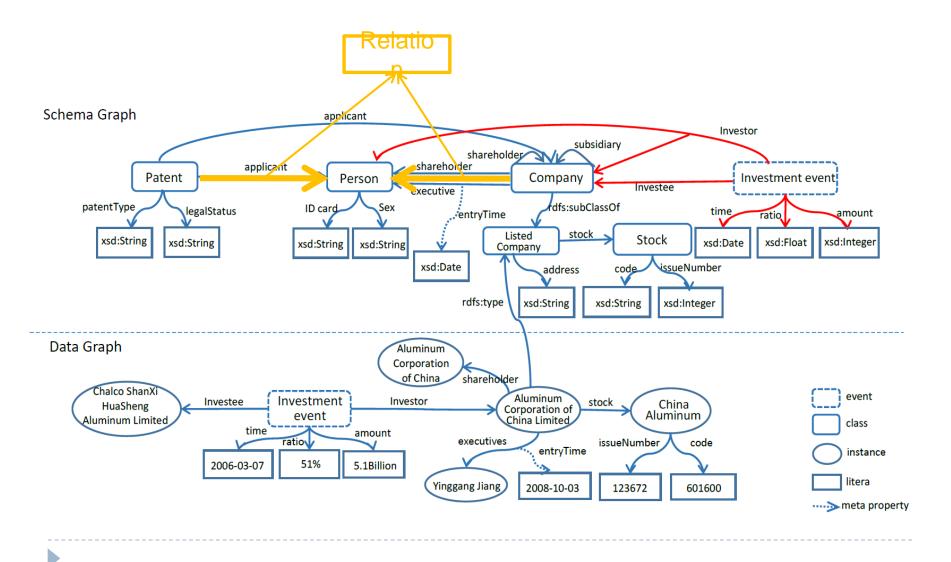


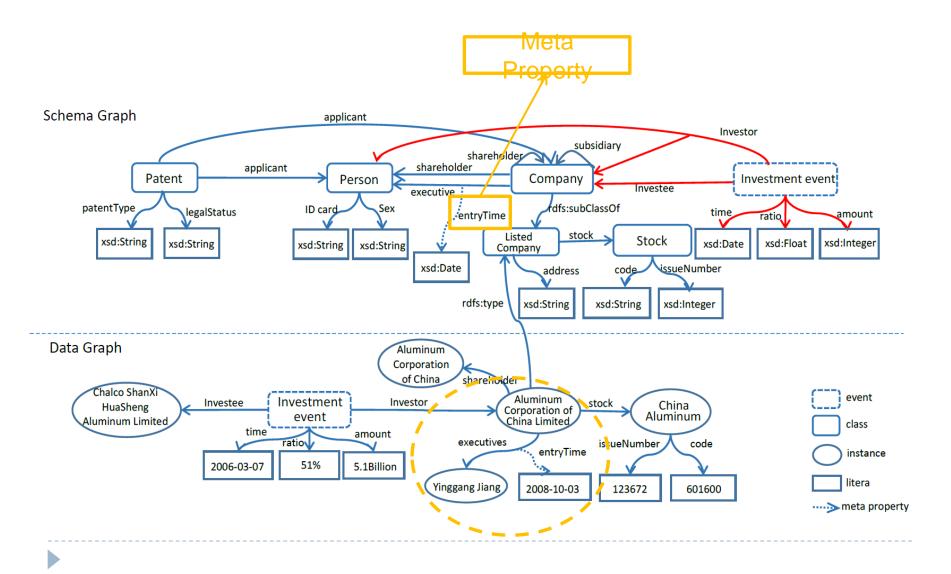


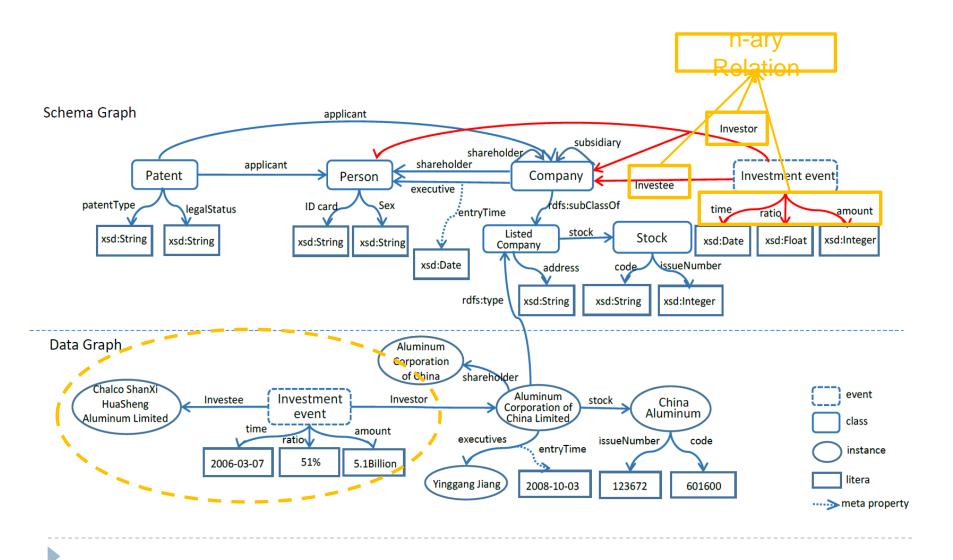


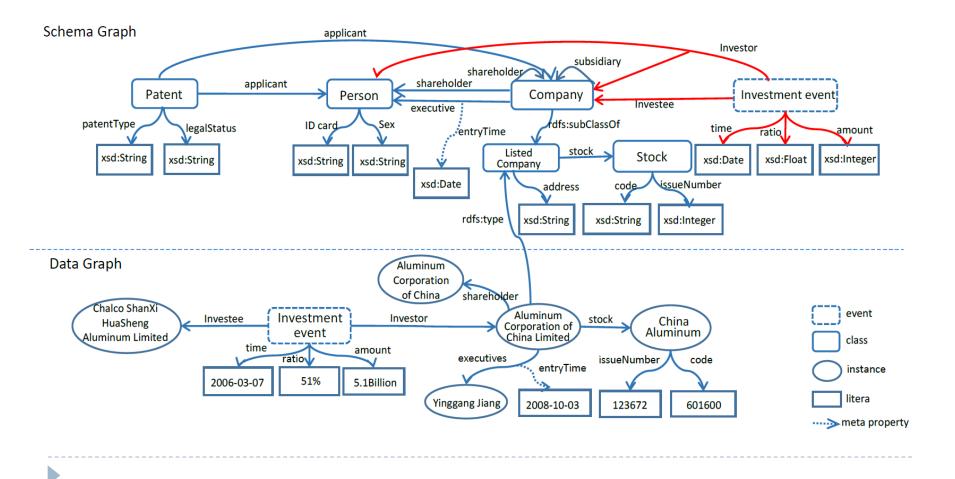












We have four major sources for constructing our EKG

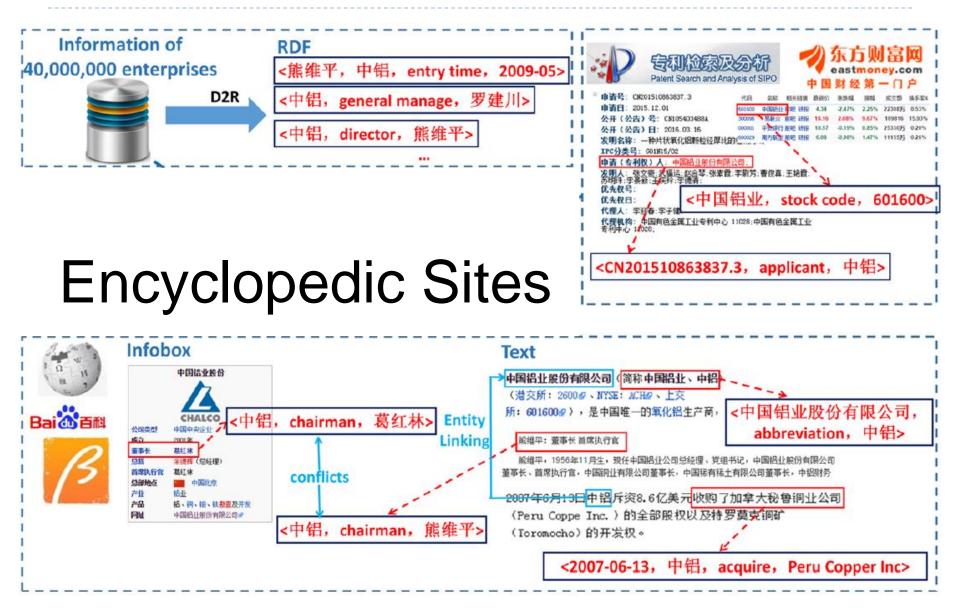


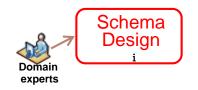


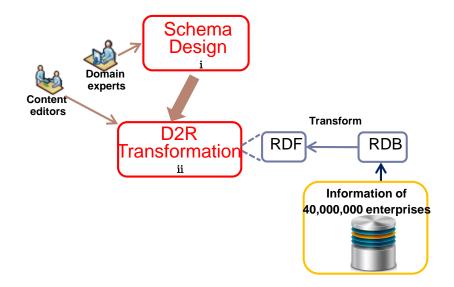
CSAIC

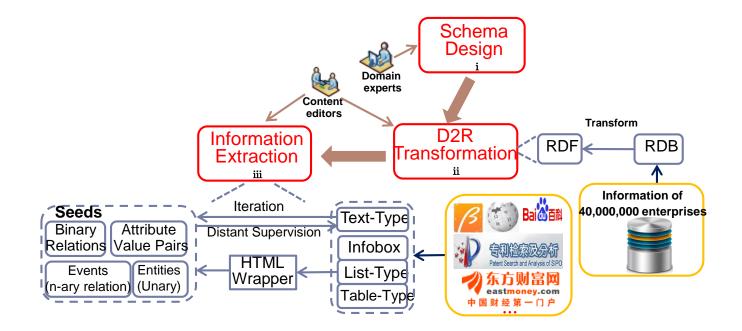


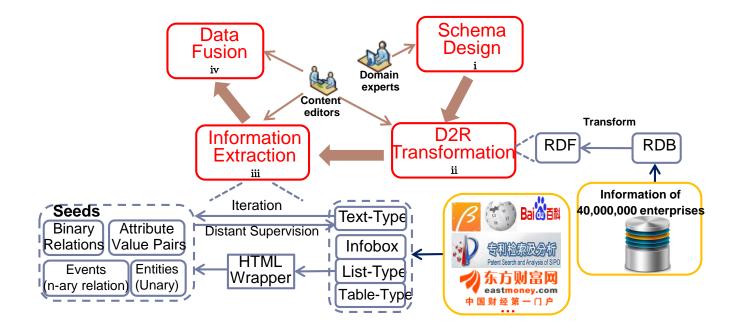
PSAN-SIPO

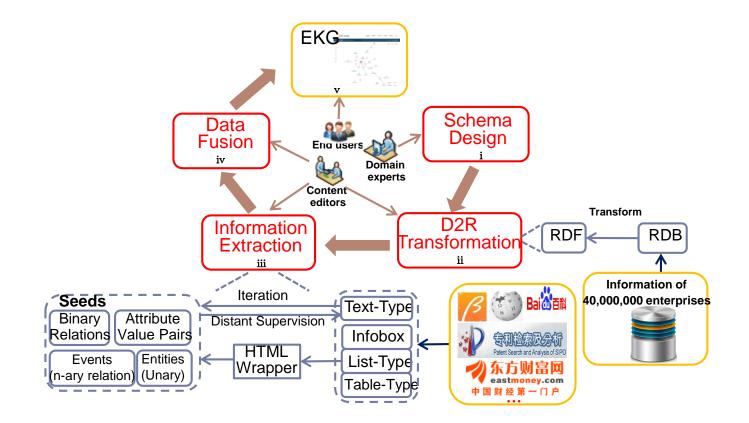


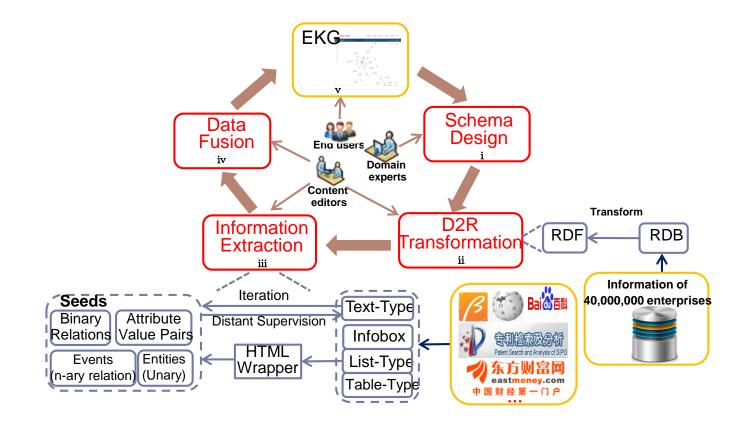








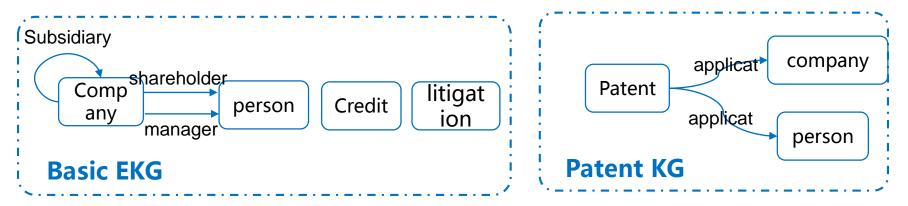




Schema Design

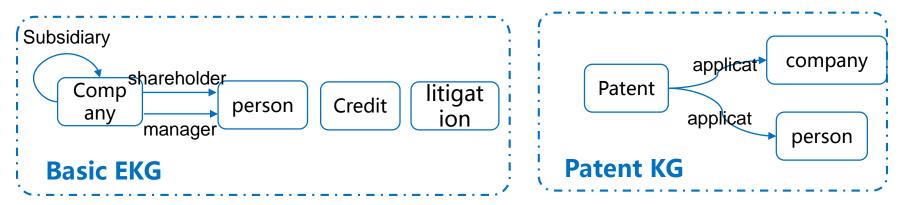
Schema Design

The 1st Iteration:

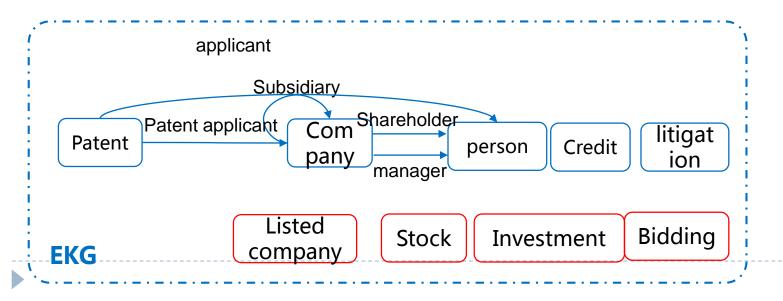


Schema Design

The 1st Iteration:



the 2nd Iteration



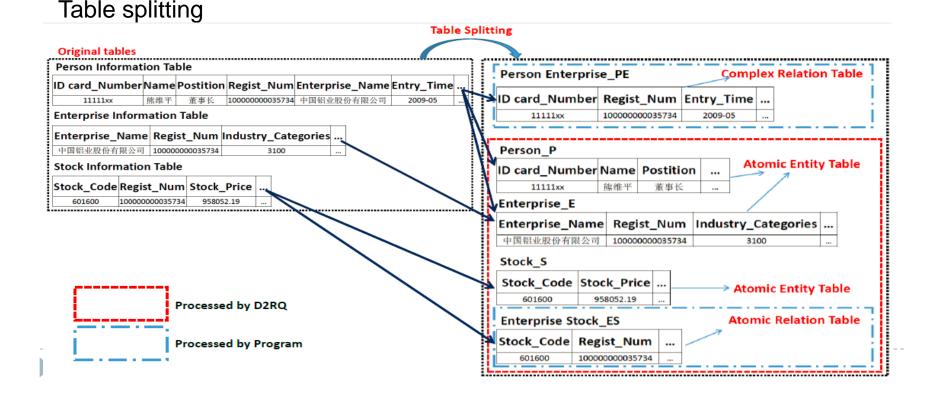
D2R Transformation

Difficulties:

- Table is lack of standardization
- Meta-property mapping
- Existing D2R tools can not solve complex mapping relationship

Solution:

- Table splitting
- Basic D2R transformation by D2RQ
- Post processing



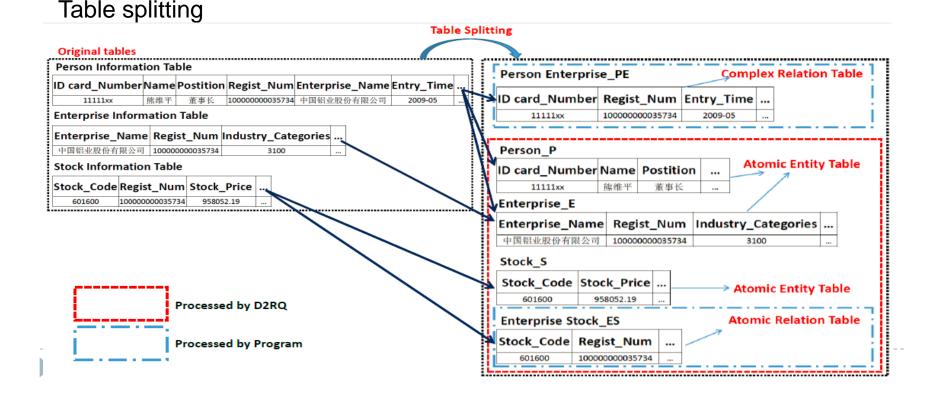
D2R Transformation

Difficulties:

- Table is lack of standardization
- Meta-property mapping
- Existing D2R tools can not solve complex mapping relationship

Solution:

- Table splitting
- Basic D2R transformation by D2RQ
- Post processing

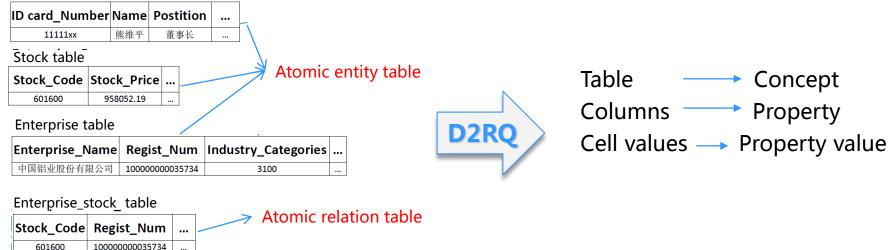


D2R Transformation

Basic D2R transformation by D2RQ

We write a customized mapping file in D2RQ to map fields related to atomic entity tables and atomic relation tables into RDF format.

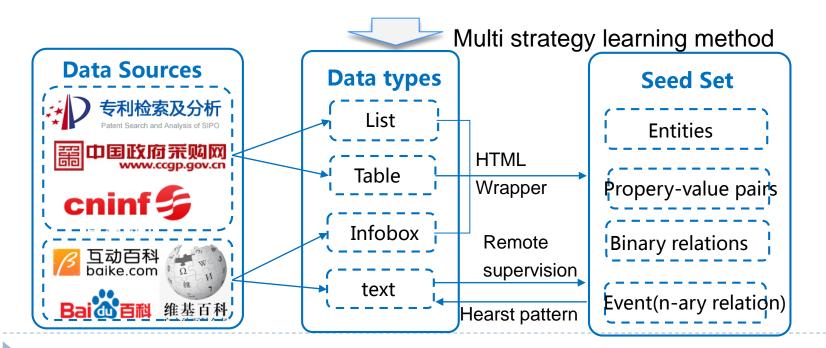
Person table



Post processing

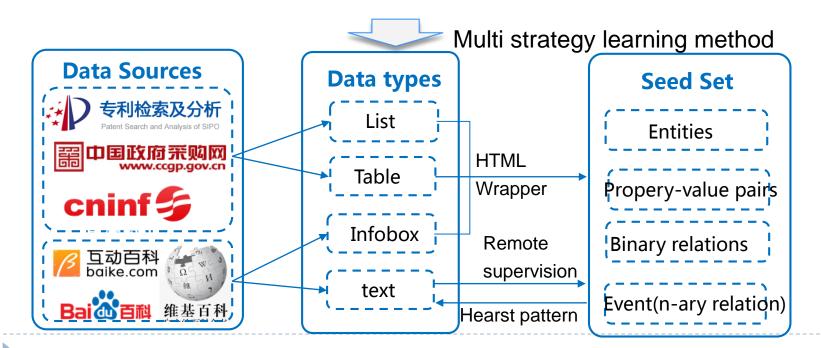
- For the complex relationship tables, and complex entity tables:
 - Meta-property mapping
 - Conditional taxonomy mapping
 - Conditional class mapping

Difficulties:



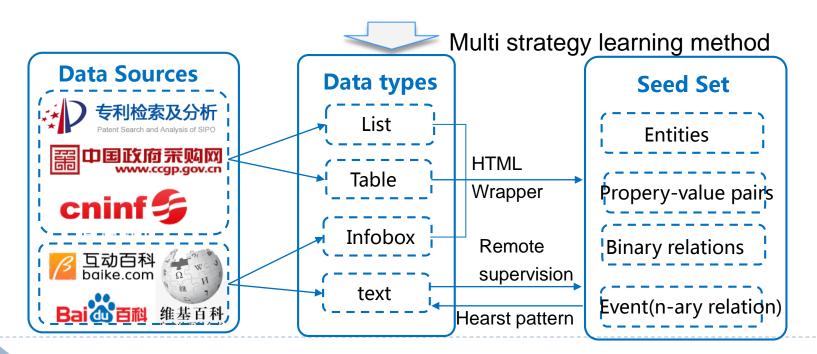
- Difficulties:
 - Multiple data source: 調中国政府 宗购网 か ち利检索及分析 www.ccgp.gov.cn か ち利检索及分析 E 潮资讯

维基百科



- Difficulties:

 - Various of target data types:



互动百科 baike.com

维基百科

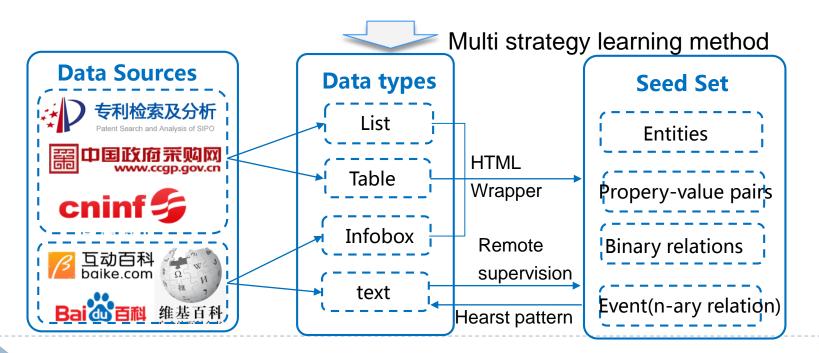
Baida音科

巨潮资讯

- Difficulties:

 - Various of target data types:

Different types of entities(Map, List, Range.....)



互动百科 baike.com

维基白科

Bai伽西和

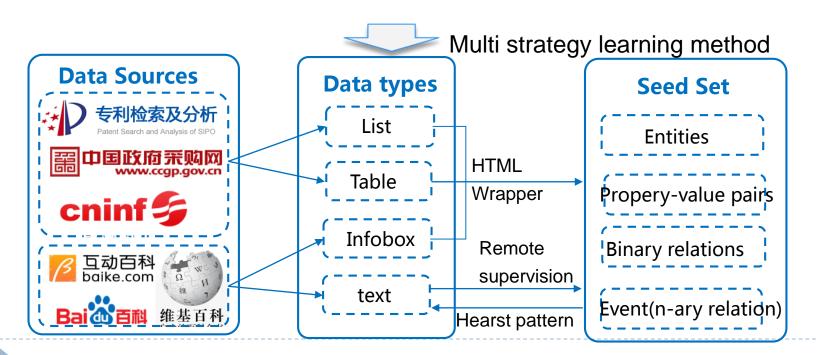
巨潮资讯

- Difficulties:

 - Various of target data types:

Different types of entities(Map, List, Range.....)

binary relations, attribute value pairs,



互动百科 baike.com

Bai 🔥 百利

巨潮资证

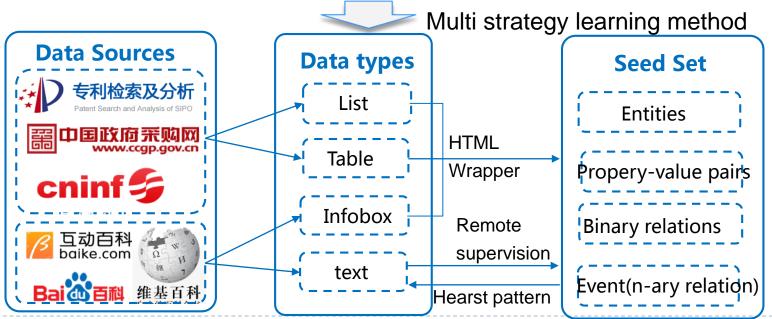
- Difficulties:

 - Various of target data types:

Different types of entities(Map, List, Range.....)

binary relations, attribute value pairs,

event(n-ary relation), synonym extraction



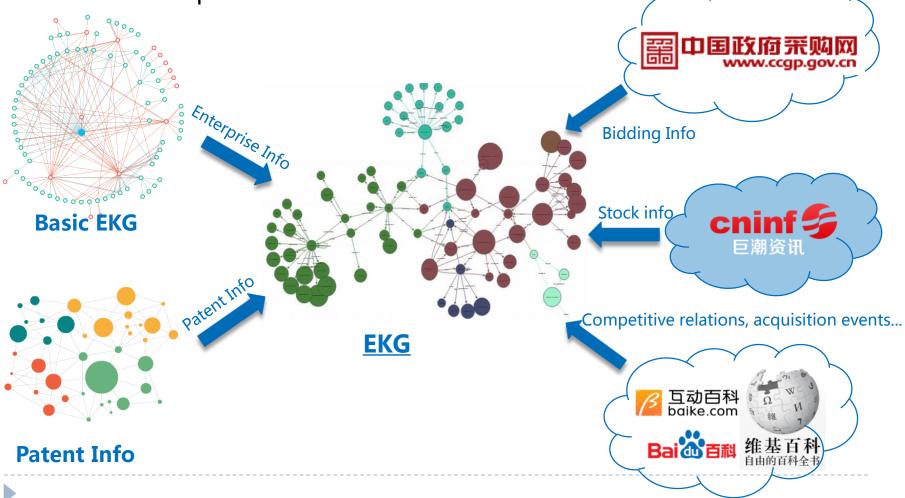
互动百科 baike.com

Bai 🔥 百利

巨潮资证

Data Fusion with Instance Matching

- Entities such as companies, people are aligned
- Data conflict problem



5 Storage Design and Query Optimization

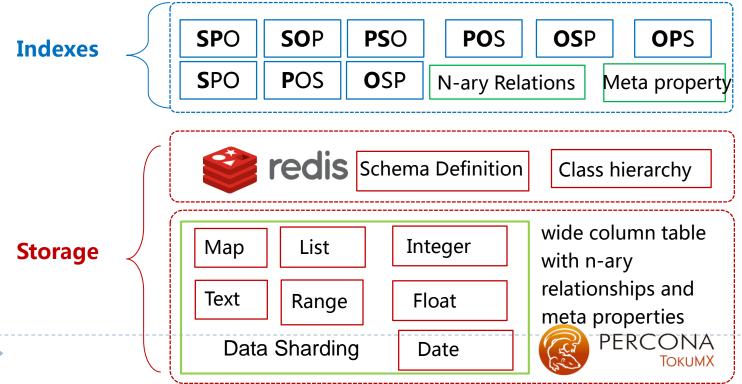
Storage Desgin

- TokuMX+ Redis
- Data Types : List, Map.....
- Store n-ary relations in the same row of a table

(Wide column table).

Query optimization

- Nine Indexes and indexes on meta properties and
- Cache schema Data in Redis
- Data Sharding for different data type of the property value
 - Support query on n-ary relation efficiently



5 Storage Design and Query Optimization

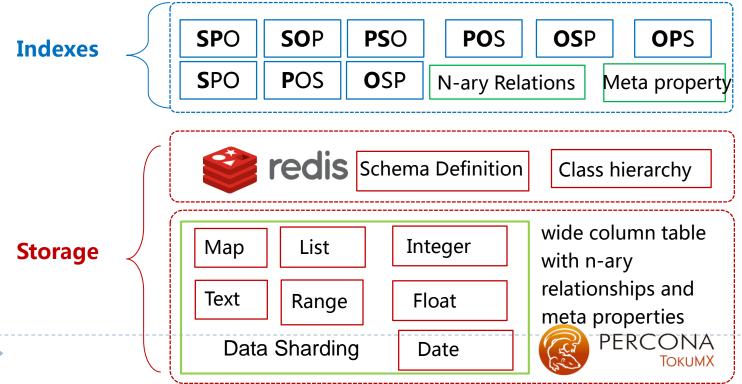
Storage Desgin

- TokuMX+ Redis
- Data Types : List, Map.....
- Store n-ary relations in the same row of a table

(Wide column table).

Query optimization

- Nine Indexes and indexes on meta properties and
- Cache schema Data in Redis
- Data Sharding for different data type of the property value
 - Support query on n-ary relation efficiently



5 Storage Design and Query Optimization

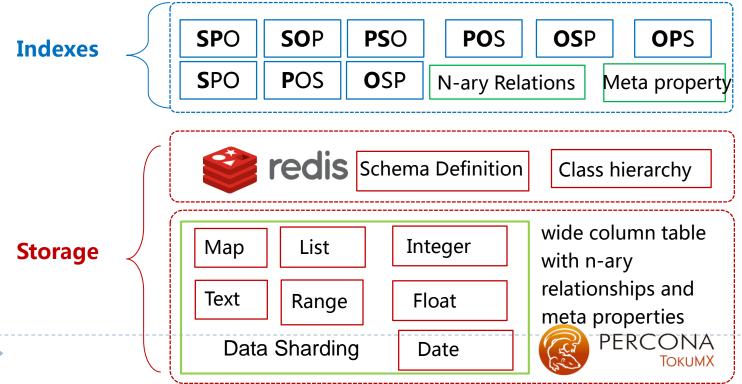
Storage Desgin

- TokuMX+ Redis
- Data Types : List, Map.....
- Store n-ary relations in the same row of a table

(Wide column table).

Query optimization

- Nine Indexes and indexes on meta properties and
- Cache schema Data in Redis
- Data Sharding for different data type of the property value
 - Support query on n-ary relation efficiently



Usage Scenarios

- An Investor hears about a company on big data area called Long Credit Beijing Corp. Ltd.
- He wants to know the detailed information of the enterprise for the investment decision in the future.
- ✓ He uses our "Magic Mirror" to have a look.



Usage Scenarios (1. General Overview)

- Firstly, the investors have a glance at basic information of Long Credit
- Then they can have a general overview of the financial status and innovation strength of Long Credit with the Key Performance Indicators (KPI) Module.



Usage Scenarios (1. General Overview)

- Firstly, the investors have a glance at basic information of Long Credit
- Then they can have a general overview of the financial status and innovation strength of Long Credit with the Key Performance Indicators (KPI) Module.



Usage Scenarios (1. General Overview)

- Firstly, the investors have a glance at basic information of Long Credit
- Then they can have a general overview of the financial status and innovation strength of Long Credit with the Key Performance Indicators (KPI) Module.



















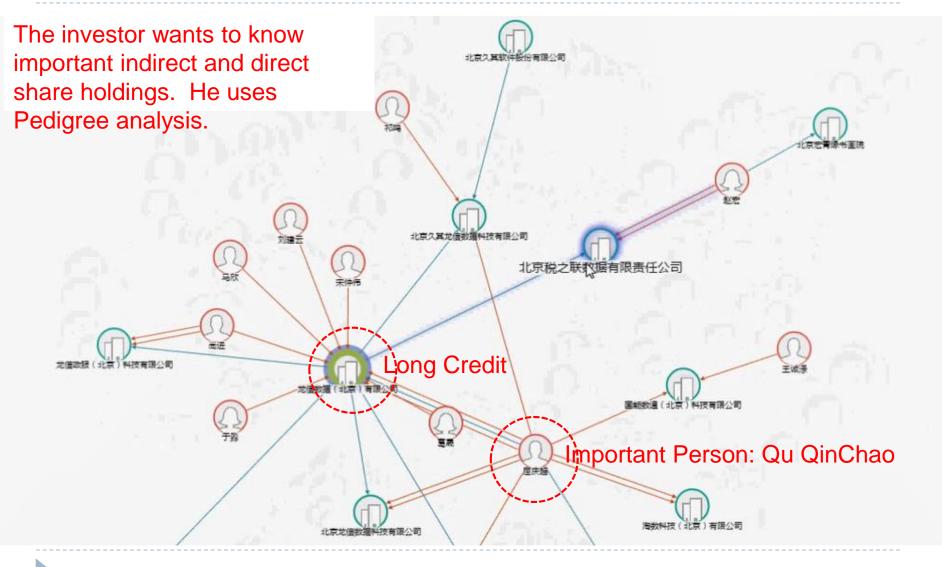




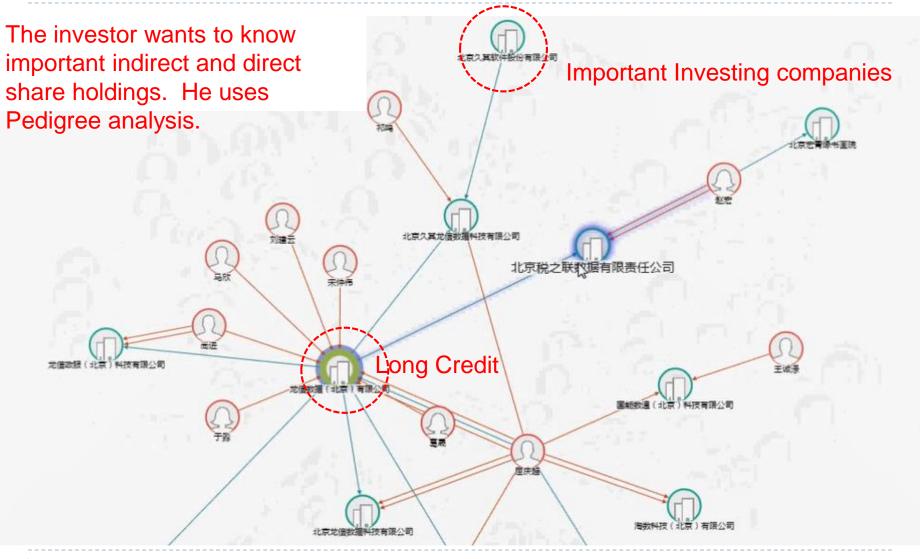




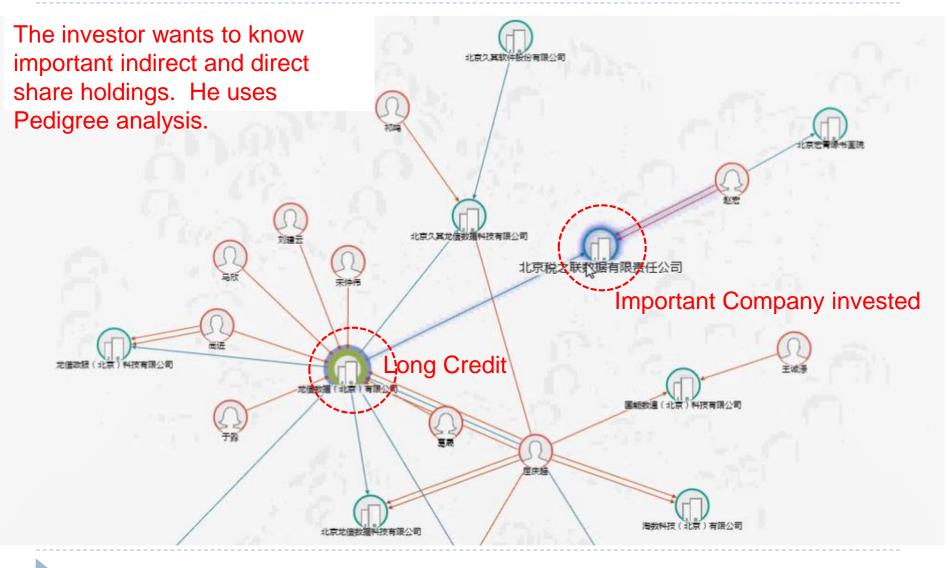
Usage Scenarios (3. Pedigree Analysis)



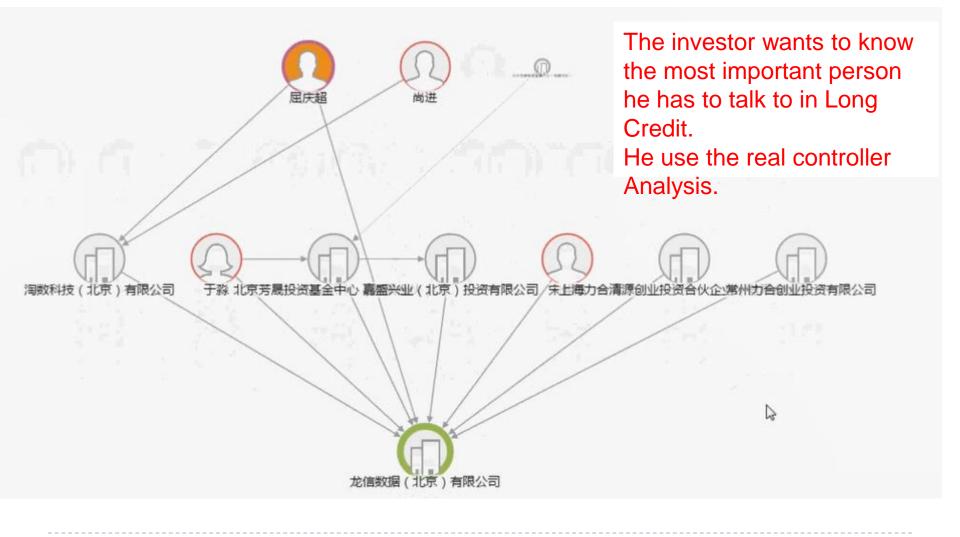
Usage Scenarios (3. Pedigree Analysis)



Usage Scenarios (3. Pedigree Analysis)



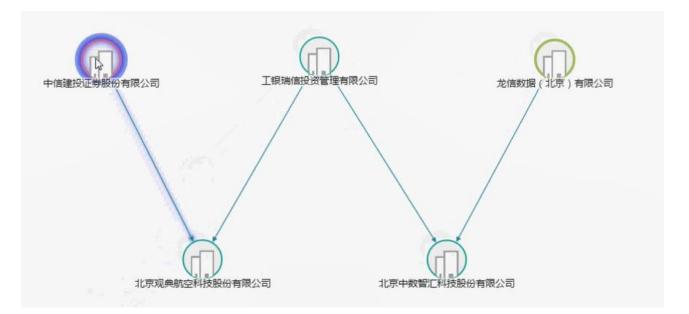
Usage Scenarios (4. Real Controller)



Usage Scenarios (5. Path Discovery)

Investors would like to know how to target the company. Therefore they look for the link between their own company and the target company..

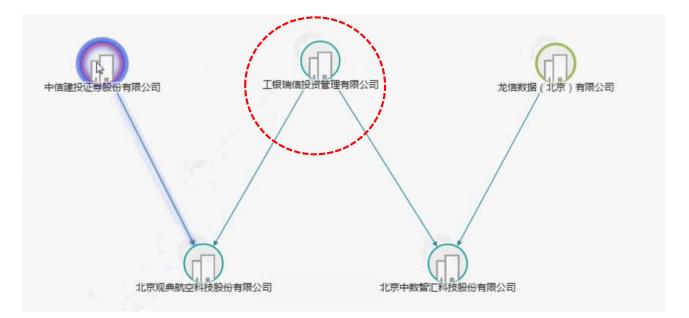
Path discovery can help users find the shortest path among concerned companies.



Usage Scenarios (5. Path Discovery)

Investors would like to know how to target the company. Therefore they look for the link between their own company and the target company..

Path discovery can help users find the shortest path among concerned companies.



Future Work

- In the future, we plan to add more data sources to the KG, such as tax and invoice information per month.
- We will also try to monitor the change of shareholders as well as share ratios
- We could develop interesting applications such as "Control intention recognition to warn the current controller of the company.

Thanks!