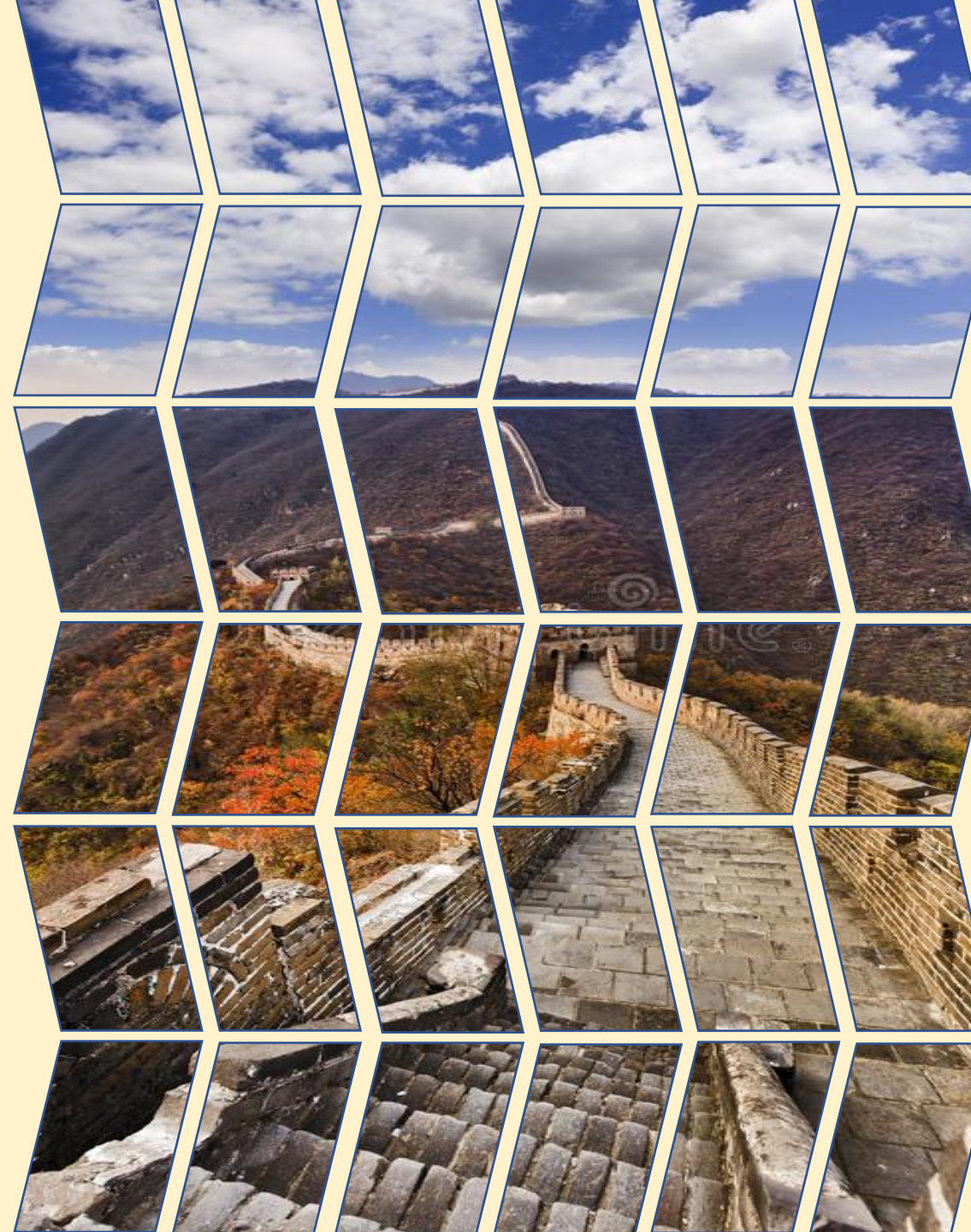# China Company Data Mining:
# Research, Scrapping, Processing, and
# Analysis for Strategic Insights

**Presented by :**
- Abdulsamet Şahin
- Subhashree Mangaraj
- Omar El Qarchaoui

Img source: Google Images

# Agenda

Img source: Google Images

# Data Sourcing

Objective of Data source search:

- Find a single or multiple data sources which can provide us with authentic company information and maximum hits.

Data Source Selection:

- www.cec.org.cn

- https://companylist.org/China/Keywords/Registered\_Companies/

- https://www.listofcompaniesin.com/china/registered/

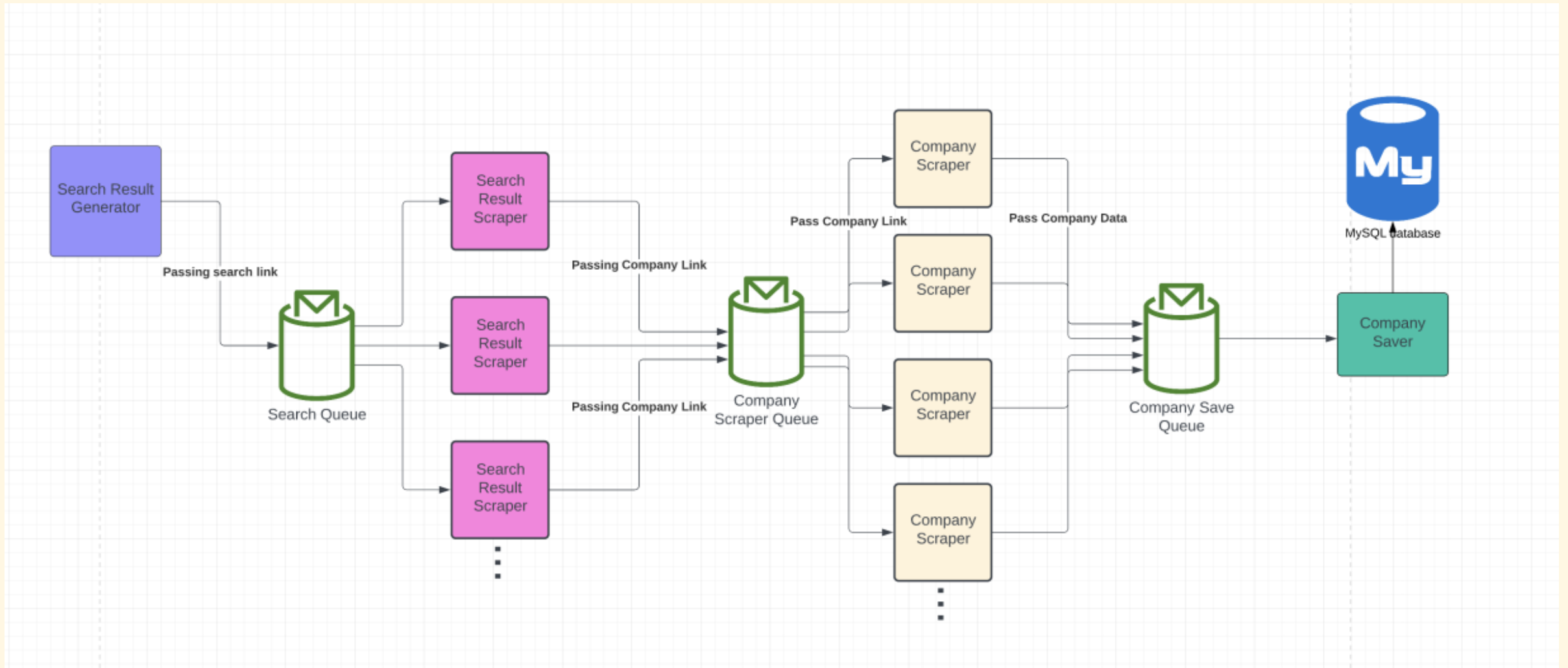- https://www.gongsi.com.cn/

# Gongsi.com.cn

## Features:

- Based on public data and available to the public.
- Contains a database of over 200 million enterprises.
- No rate limit in place.

## Data Coverage:

- Basic company information

- Detailed company information

- Main staff

- Branches

- Shareholders

- Change history

- Investments

- Annual Reports

innoscripta

# Workflow of Scraper

# Technical Specifications

Features:

- Created using Python3

- Utilizes RabbitMQ to hold tasks and facilitate communication between workers.

- Thanks to RabbitMQ, the scraper can be easily scaled.

- Utilizes MySQL to store scraped data.

- Designed to run continuously.

- Multiple workers can scrape 1 million companies in a single day.

# RabbitMQ

**⊔RabbitMQ**™   RabbitMQ 3.9.7   Erlang 24.1.2

**Overview**   **Connections**   **Channels**   **Exchanges**   **Queues**   **Admin**

## Queues

▼ **All queues (3)**

Pagination

Page `1 ▾` of 1  - Filter: [_____]   ☐ Regex **?**

| Overview | | | | Messages | | | Message rates | | | +/- |
|---|---|---|---|---|---|---|---|---|---|---|
| **Name** | **Type** | Features | **State** | **Ready** | **Unacked** | **Total** | **incoming** | **deliver / get** | **ack** | |
| **company_data** | classic | D Args | ■ running | 0 | 0 | 0 | 6.8/s | 11/s | 11/s | |
| **company_link** | classic | D Lim Ovfl Args | ■ running | 86,326 | 21 | 86,347 | 0.00/s | 11/s | 11/s | |
| **search_page** | classic | D Args | ▢ idle | 491,887 | 0 | 491,887 | 0.00/s | 0.00/s | 0.00/s | |

▶ **Add a new queue**

**HTTP API**   **Server Docs**   **Tutorials**   **Community Support**   **Community Slack**   **Commercial Support**   **Plugins**   **GitHub**   **Changelog**

# Search Result Generator

Features:

- Starting point of the scraper.

- https://www.gongsi.com.cn/search/bs1in52eyg5eyl6

- Since there is a pagination limit, filters are used to reduce the number of pages.

- It uses combinations of province, category, and establishment year filters to generate search result.

- It will regenerate when then queue is empty.

# Search Result Scraper

Features:

- The scraper collects company links by scraping the search result pages one by one, which are then scraped by the "Company Scraper" worker.

- To prevent overloading the RabbitMQ, this scraper will stop working when there are 100,000 companies unscraped.

# Company Scraper

Features:

- The scraper will retrieve tasks from the previous worker (Search Page Scraper).

- It will scrape all the required data and pass it to the "Company Saver" worker.

- The scraper can be easily scaled by running as many workers as necessary.

- It is possible to run the scraper from other servers.

# Company Saver

Features:

- This worker retrieves company data from the "Company Scraper" worker.

- It performs various cleaning operations on the data.

- The data is then stored in the database.

- Updates the company if already exists.

- If the company already exists, it is updated with the new data.
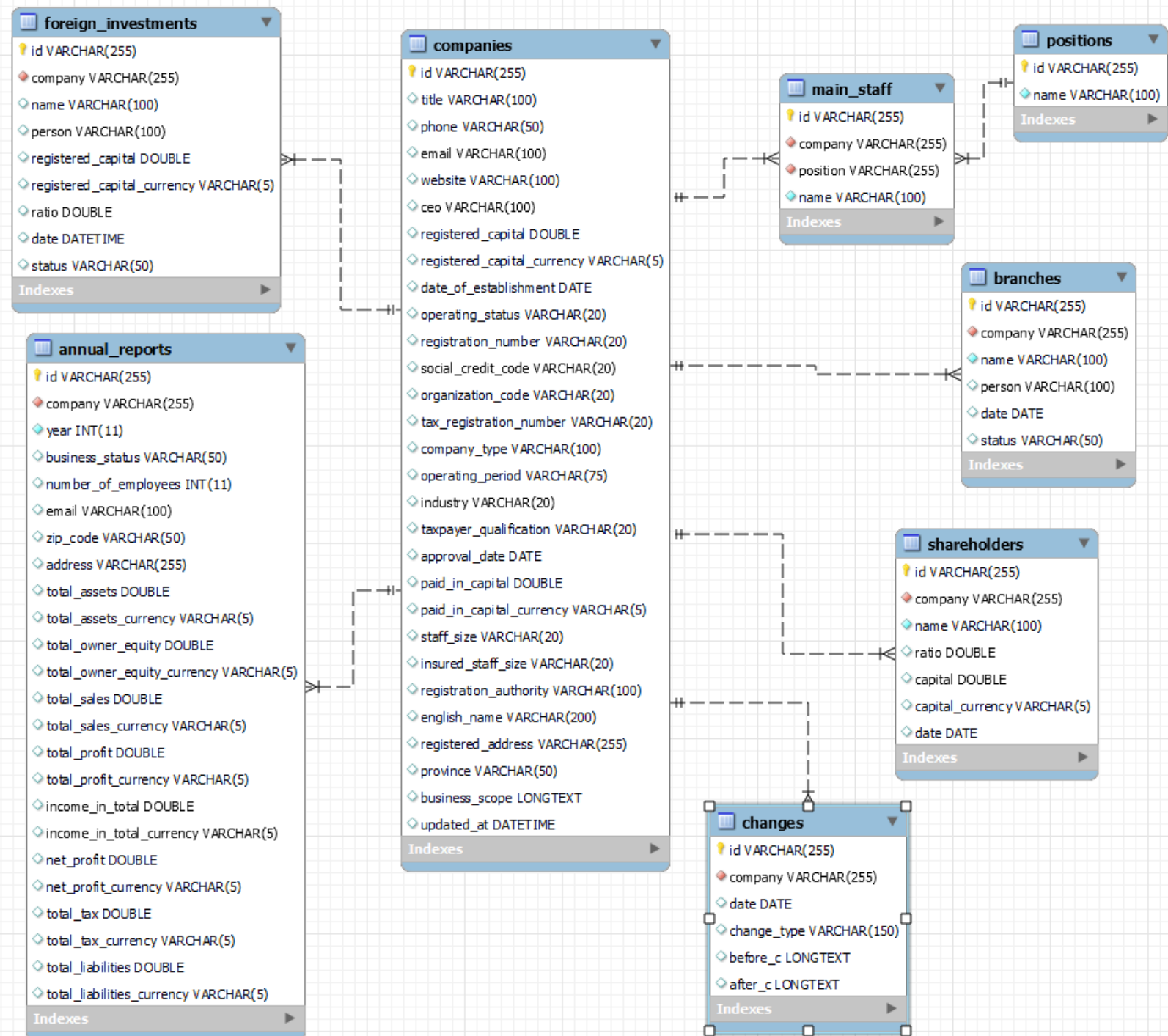
INNOSCRIPTA

# Data Pre-processing

Features:

- Replacing strings indicating a lack of data with "None,"

- Removing non-numeric characters from numeric values,

- Converting currency values to a unit currency metric,

- Adding currency indicator fields to address currency representation inconsistencies

- Cleaning up main staff data by separating multiple positions into a new table and storing main staff positions as a many-to-many relationship.

- The data is described as fairly clean after undergoing these pre-processing steps.

# Database MySQL

**foreign_investments**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔷 name VARCHAR(100)
- 🔷 person VARCHAR(100)
- 🔷 registered_capital DOUBLE
- 🔷 registered_capital_currency VARCHAR(5)
- 🔷 ratio DOUBLE
- 🔷 date DATETIME
- 🔷 status VARCHAR(50)

Indexes

**companies**
- 🔑 id VARCHAR(255)
- 🔷 title VARCHAR(100)
- 🔷 phone VARCHAR(50)
- 🔷 email VARCHAR(100)
- 🔷 website VARCHAR(100)
- 🔷 ceo VARCHAR(100)
- 🔷 registered_capital DOUBLE
- 🔷 registered_capital_currency VARCHAR(5)
- 🔷 date_of_establishment DATE
- 🔷 operating_status VARCHAR(20)
- 🔷 registration_number VARCHAR(20)
- 🔷 social_credit_code VARCHAR(20)
- 🔷 organization_code VARCHAR(20)
- 🔷 tax_registration_number VARCHAR(20)
- 🔷 company_type VARCHAR(100)
- 🔷 operating_period VARCHAR(75)
- 🔷 industry VARCHAR(20)
- 🔷 taxpayer_qualification VARCHAR(20)
- 🔷 approval_date DATE
- 🔷 paid_in_capital DOUBLE
- 🔷 paid_in_capital_currency VARCHAR(5)
- 🔷 staff_size VARCHAR(20)
- 🔷 insured_staff_size VARCHAR(20)
- 🔷 registration_authority VARCHAR(100)
- 🔷 english_name VARCHAR(200)
- 🔷 registered_address VARCHAR(255)
- 🔷 province VARCHAR(50)
- 🔷 business_scope LONGTEXT
- 🔷 updated_at DATETIME

Indexes

**main_staff**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔴 position VARCHAR(255)
- 🔷 name VARCHAR(100)

Indexes

**positions**
- 🔑 id VARCHAR(255)
- 🔷 name VARCHAR(100)

Indexes

**branches**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔷 name VARCHAR(100)
- 🔷 person VARCHAR(100)
- 🔷 date DATE
- 🔷 status VARCHAR(50)

Indexes

**annual_reports**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔑 year INT(11)
- 🔷 business_status VARCHAR(50)
- 🔷 number_of_employees INT(11)
- 🔷 email VARCHAR(100)
- 🔷 zip_code VARCHAR(50)
- 🔷 address VARCHAR(255)
- 🔷 total_assets DOUBLE
- 🔷 total_assets_currency VARCHAR(5)
- 🔷 total_owner_equity DOUBLE
- 🔷 total_owner_equity_currency VARCHAR(5)
- 🔷 total_sales DOUBLE
- 🔷 total_sales_currency VARCHAR(5)
- 🔷 total_profit DOUBLE
- 🔷 total_profit_currency VARCHAR(5)
- 🔷 income_in_total DOUBLE
- 🔷 income_in_total_currency VARCHAR(5)
- 🔷 net_profit DOUBLE
- 🔷 net_profit_currency VARCHAR(5)
- 🔷 total_tax DOUBLE
- 🔷 total_tax_currency VARCHAR(5)
- 🔷 total_liabilities DOUBLE
- 🔷 total_liabilities_currency VARCHAR(5)

Indexes

**shareholders**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔷 name VARCHAR(100)
- 🔷 ratio DOUBLE
- 🔷 capital DOUBLE
- 🔷 capital_currency VARCHAR(5)
- 🔷 date DATE

Indexes

**changes**
- 🔑 id VARCHAR(255)
- 🔴 company VARCHAR(255)
- 🔷 date DATE
- 🔷 change_type VARCHAR(150)
- 🔷 before_c LONGTEXT
- 🔷 after_c LONGTEXT

Indexes

**innoscripta**

# Tables

| Tables | Description |
|---|---|
| companies | General info table for companies. 'id' column is the primary key, columns for company details such as contact information, financial data, and registration details, and a datetime column to track updates |
| annual_report | contains financial information related to different companies for different years |
| branches | contains information related to different branches of a company |
| changes | information related to changes that occurred in different companies. |
| Foreign_investmen | contains information about the foreign investments made by a company. |
| Main_Staff | information about the staff of a company |
| positions | information about the positions of C-Level employees in a company. |
| shareholders | contains information about the shareholders of a company. |

# Results and Findings : 903685 total number of companies

| Number of companies scraped with | Count |
|---|---|
| contact information | 450833 |
| website information available | 42461 |
| CEO information available | 902021 |
| addresses available | 901570 |
| information on the number of employees | 147376 |
| industry information available | 905870 |
| atleast one contact info and ceo information | 451457 |
| information on C-Level personnel | 523753 |
| information regarding company branches | 46049 |
| change history available | 290241 |
| shareholder information | 404433 |
| information on their foreign investments | 34063 |
| information on the balance sheet and revenue and profit | 373353 |

# Analysis

Geographical Analysis:

- There are 30 provinces

- The top 7 province where maximum of the companies reside are given below

# Number of Employees Analysis:

- There are 7 categories in which the companies are divided as per staff size

- Majority of the companies have staff size less than 5

- Only 74 companies having 10,000 employees or more

- Fewer mid sized companies as number of companies with staff sizes in the range of 100-999 employees is relatively low



Chart Title

# Industry Analysis:

- There are 94 categories of industries

- We present here the top 20 industries and number of companies in that category

- The largest number of companies falls within the stevedoring and transport agency industry

## Operating Status Analysis:

- Removing Null, the data has 12 statuses

- More than 50 percent of companies in the database are on 'Survival' status
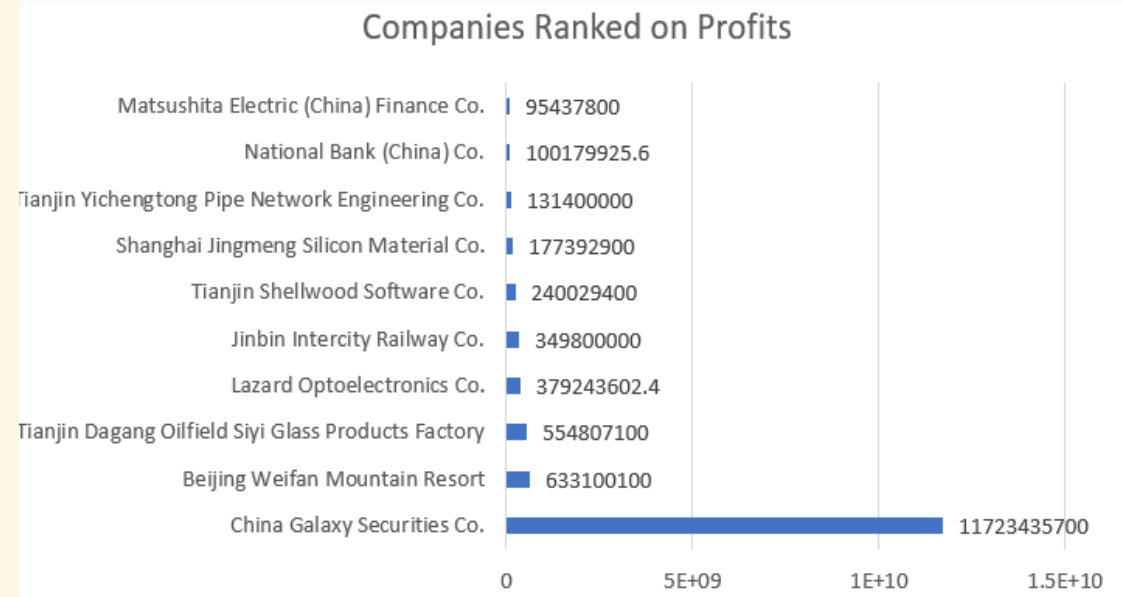
- Followed by companies in 'cancelled' status

# Foreign  Investment Analysis:

- Top 10 companies with highest foreign investments



Invested capital

# Balance Sheet and Revenue Analysis:

## Companies Ranked on Profits

| Company | Value |
|---|---|
| Matsushita Electric (China) Finance Co. | 95437800 |
| National Bank (China) Co. | 100179925.6 |
| Tianjin Yichengtong Pipe Network Engineering Co. | 131400000 |
| Shanghai Jingmeng Silicon Material Co. | 177392900 |
| Tianjin Shellwood Software Co. | 240029400 |
| Jinbin Intercity Railway Co. | 349800000 |
| Lazard Optoelectronics Co. | 379243602.4 |
| Tianjin Dagang Oilfield Siyi Glass Products Factory | 554807100 |
| Beijing Weifan Mountain Resort | 633100100 |
| China Galaxy Securities Co. | 11723435700 |

## Companies Ranked on Assets

| Company | Value |
|---|---|
| Beijing Shougang Hotel Development Co. | 2734717500 |
| Xinli Energy Development Co. | 2832452200 |
| Beijing Construction Engineering Civil Engineering Co. | 4110632156 |
| Tianjin Jinran Thermal Power Co. | 7139610000 |
| Lazard Optoelectronics Co. | 8724234824 |
| Jinbin Intercity Railway Co. | 10324700000 |
| Matsushita Electric (China) Finance Co. | 12894957800 |
| National Bank (China) Co. | 20267078697 |
| Tianjin Dagang Oilfield Siyi Glass Products Factory | 27963604200 |
| China Galaxy Securities Co. | 4.73359E+11 |

## Companies rankled on Income

| Company | Value |
|---|---|
| China News Development Co. | 998883000 |
| Shanghai Jingmeng Silicon Materials Co. | 1109659300 |
| Huadian (Beijing) Cogeneration Co. | 1478857500 |
| Shanghai Ming Hua Engineering & Construction Co. | 1494259640 |
| Lazard Optoelectronics Co. | 1678112785 |
| ConocoPhillips China Ltd. | 2740350000 |
| Beijing Construction Engineering Civil Engineering Co. | 3787891581 |
| China Construction Technology Group Co. | 3832444598 |
| China People's Property and Casualty Insurance Company... | 5529981300 |
| China Galaxy Securities Co. | 18877386900 |

## Companies Ranked on Total Tax paid

| Company | Value |
|---|---|
| China Union Qian Yuan Real Estate Fund Management Co. | 40399693 |
| China Construction Bank Corporation Shanghai Songjiang Sub-branch | 68280000 |
| Beijing Huadu Brewery & Food Co. | 88045200 |
| Huadian (Beijing) Cogeneration Co. | 112693700 |
| China Construction Technology Group Co. | 114752665.3 |
| Shanghai Ming Hua Engineering & Construction Co. | 127410467.9 |
| Guoxin Securities Company Limited Beijing Branch | 161031900 |
| ConocoPhillips China Co. | 216960000 |
| Tianjin Dagang Oilfield Siyi Glass Products Factory | 327651500 |
| Beijing Xin Yan Tong Hui Cleaning Service Co. | 1529678000 |

# Real Time Power BI Dashboard:

- Objective is to analyse and track the complete details of each company in a single report.

- Connected to database for real time update of data

- Information displayed based on selected company id

# Future Work

## Pre-trained language model for tranlation

- The "facebook/nllb-200-distilled-600M" model is a highly accurate and efficient pre-trained language model that can be customized for specific languages or domains to improve its translation quality.

# Future Work

Stream processing(Kafka/Redpanda) and real-time databases (Rockset/Pinot)

- Low latency data access

- High processing,

- High scalability, and

- Flexibility for changing requirements

- Ideal for applications that require real-time data processing

innoscripta

# Conclusion

Data Sourcing

Data Scrapping

Database storage

Real-time scrapping

Analysis

BI Dashboard

Recommendations

# Thank you