

# 上海证券交易所博士后出站报告

## 证券市场文本挖掘技术研究与应用

报告人：白雪  
所内导师：白硕  
所外导师：朱扬勇  
联系人：卢文莹  
日期：2015. 6

# 汇报提纲

1 研究背景和意义

2 报告主要章节

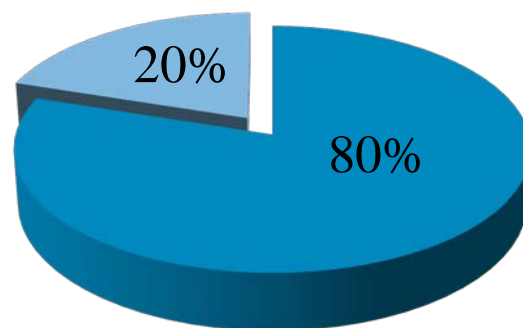
3 主要研究成果

4 启示与建议

5 工作总结

# 研究背景和意义

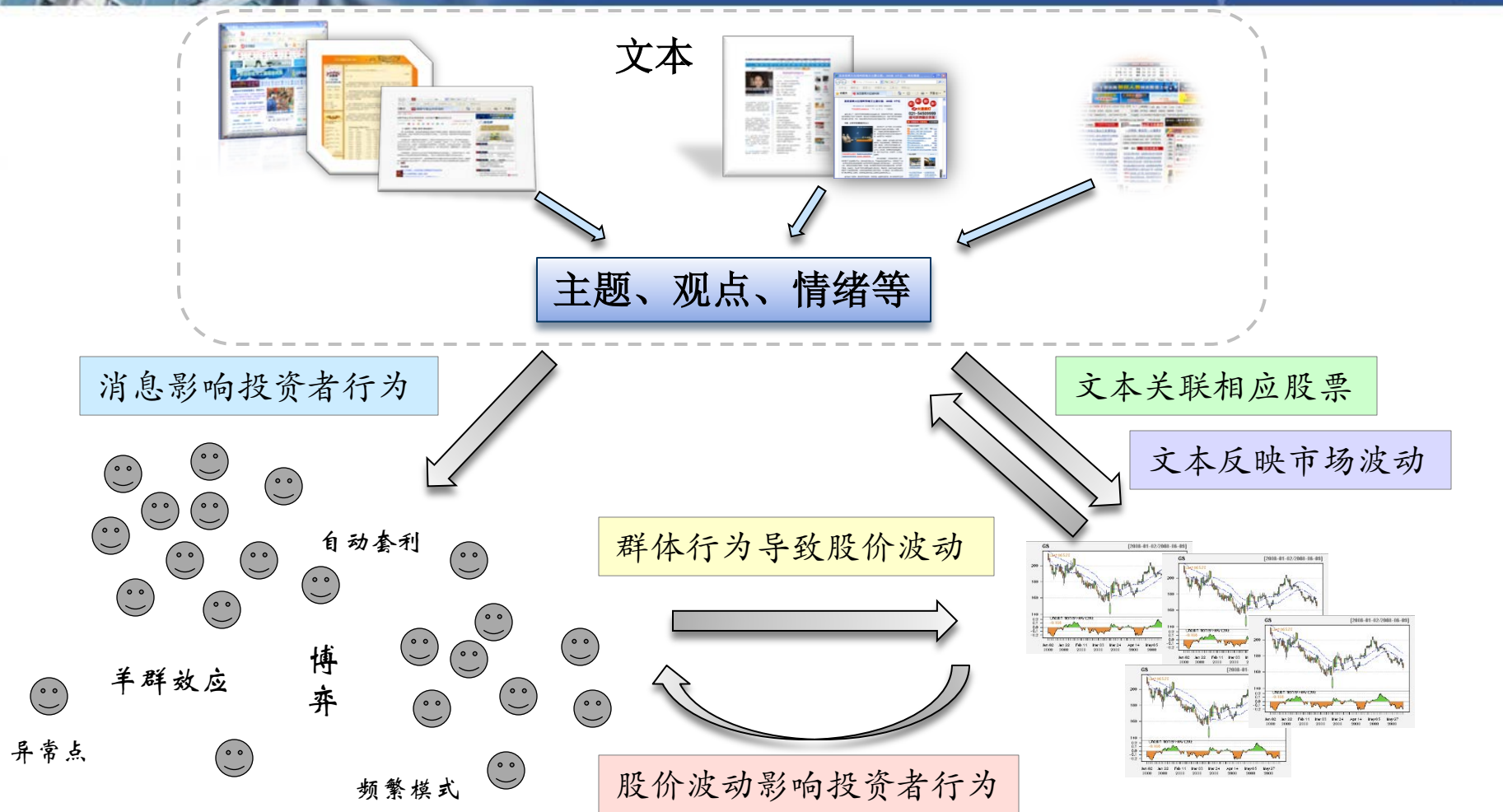
- 互联网数据大爆炸时代
  - ◆ 海量非结构化数据尚未得到有效应用
- 证券行业
  - ◆ 相关的网络数据资源越来越多
  - ◆ 对于信息的快速加工、精准反应需求也与日俱增
    - 市场情绪网络留痕——为度量、量化情绪提供了素材
    - 主题投资、事件套利
    - 舆情监控、信息追踪



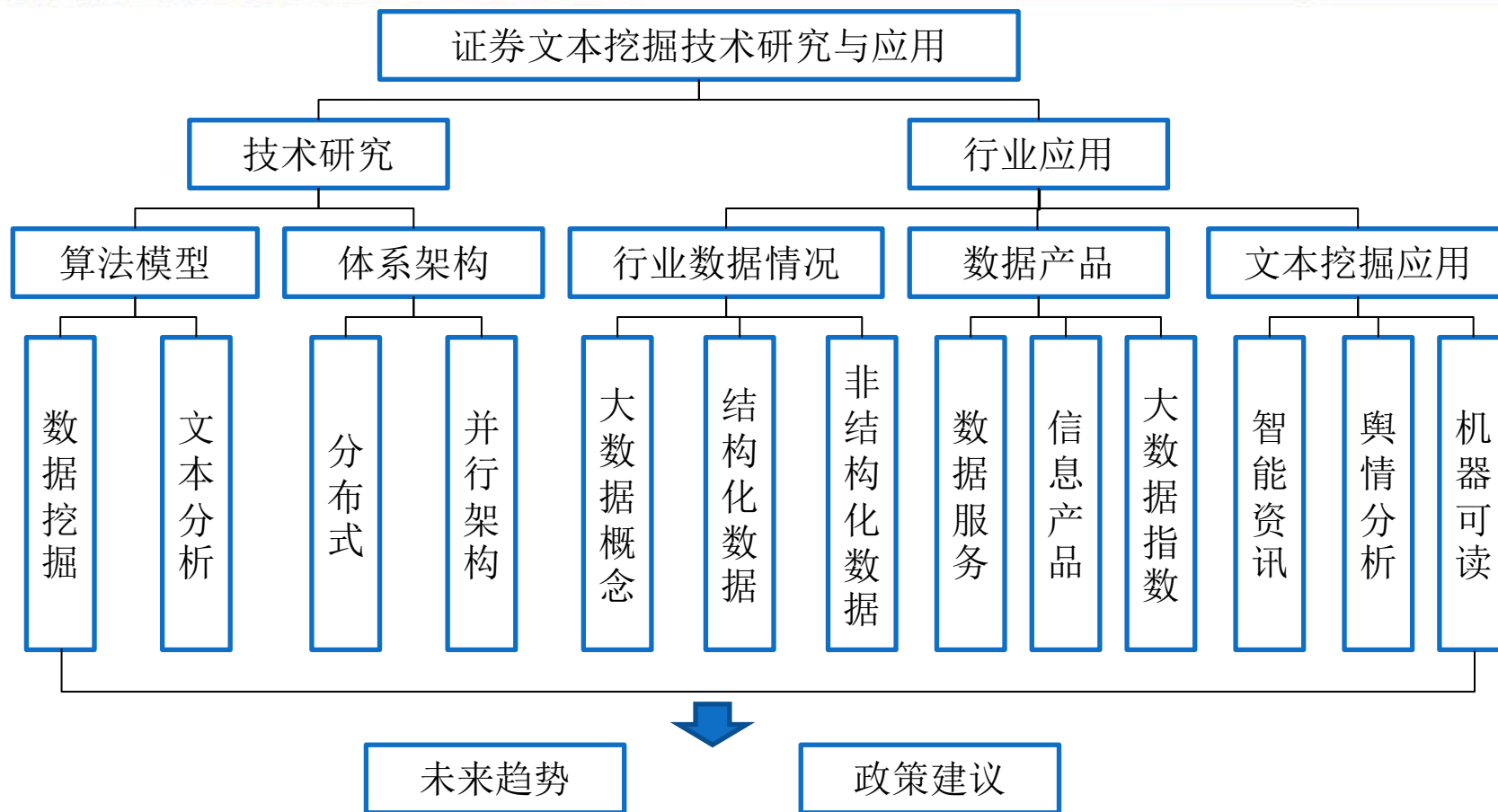
■ 非结构化数据  
■ 结构化数据

对证券行业文本挖掘产生了迫切的需求

# 文本对市场的影响



# 研究内容





# 报告章节目录

## 1 绪论

- 1.1 研究背景及意义
- 1.2 研究目的及内容
- 1.3 文本主要研究成果
- 1.4 文章组织结构

## 2 国内外相关理论及研究综述

- 2.1 情感分析
- 2.2 行为挖掘
- 2.3 关联分析

## 3 证券行业数据现状

- 3.1 大数据概念
- 3.2 证券行业数据资源
- 3.3 证券行业结构化数据
- 3.4 证券行业非结构化数据

## 4 大数据常用技术框架

- 4.1 Hadoop
- 4.2 Spark
- 4.3 Storm

## 5 数据分析基础方法与模型

- 5.1 基础概念
- 5.2 相似性度量
- 5.3 文本挖掘

## 6 证券业数据产品和服务

- 6.1 基于结构化数据的数据服务
- 6.2 基于非结构化数据的数据服务
- 6.3 证券行业大数据综合平台型服务
- 6.4 国际交易所信息产品
- 6.5 南方-新浪大数据指数

## 7 证券行业文本挖掘应用

- 7.1 投资综合型社区
- 7.2 文本信息资讯服务
- 7.3 专业文本挖掘服务
- 7.4 机器可读新闻技术专题调研
- 7.5 证券行业文本挖掘相关研究报告
- 7.6 行业应用现状与问题探讨

## 8 应用案例——证券知识图谱关键技术研究

- 8.1 知识图谱概述
- 8.2 证券领域知识图谱应用前景
- 8.3 系统架构与主要模块
- 8.4 关键技术

## 9 大数据在监管中的应用与建议

- 9.1 国际应用实践
- 9.2 证监会与深交所应用现状
- 9.3 对上交所的启示

# 数据、算法和架构

## 证券行业数据现状

### 结构化数据

行情数据

成交数据

指标数据

### 非结构化数据

公司数据

网络数据

研报数据

## 数据挖掘算法

文本挖掘

情感分析

关联分析

聚类分析

...

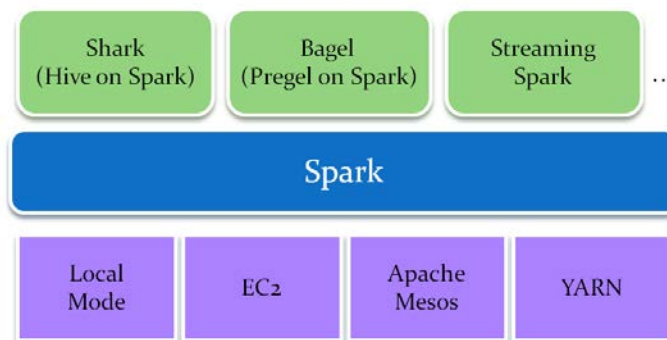
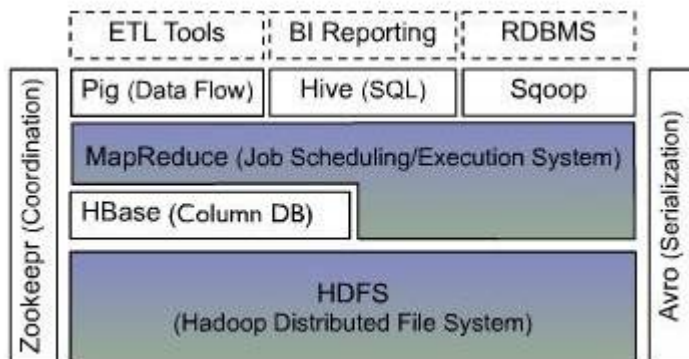
## 大数据基础架构

Hadoop

Yarn

Spark

Storm



# 行业文本挖掘现状调研

- 证券行业文本挖掘应用可归纳为三大类

- ◆ 投资综合性社区

StockTwits

雪球

- 各种重要消息发布、频繁互动、思想碰撞的交流平台

- ◆ 文本信息资讯服务



THOMSON REUTERS

Bloomberg

朝阳永续  
SUNTIME

- 金融资讯端文本分类、机器可读等服务

- ◆ 专业文本挖掘服务



SmogFarm

- 互联网专业情感分析服务

- 大数据指数

- ◆ 南方-新浪大数据



- ◆ 百度百发





# 机器可读新闻技术调研

## ● 机器可读新闻



**Bloomberg**

- ◆ 基于信息编码技术、文本挖掘、自然语言处理、情感分析、统计学等算法技术；
- ◆ 以电脑可读的语言编写，**低延迟**瞬时传递关键信息；
- ◆ 服务于高频交易、程序化交易、实时指数发布、套利交易、风险管理等。

```

16:39:37.923944 IP 172.27.71.6.www > 10.10.10.75.1103: F 52:253(201
ack 56 win 32768
0x0000: 4500 0001 0076 0000 3506 111b ac1b 4706 E...lv.....G
0x0010: 0a0a 0a0b 0050 004f 27aa b5fa 036d 6ebc ...K.P.O.....mn.
0x0020: 5018 8001 12d1 0000 4554 5450 2f51 2e31 F...0...HTTP/1.1
0x0030: 2010 0000 4b04 4b04 0a43 6163 6865 2d43 200...OK...
0x0040: 6f6e 7472 6266 3a20 6e6f 2a63 6163 6865 control:n-co-cache
0x0050: 0d0a 4d73 6749 643a 2030 0a04 636f 6874 ...MsgID: 0..Cont
0x0060: 6f6e 742d 4c65 8e67 7468 3a20 3537 0d0a ent-length: 87..
0x0070: 4370 7474 686e 742d 5479 7065 3a20 6170 Content-Type: ap
0x0080: 706c 6963 6174 674 696f 6d2f 6e63 7465 742d plication/cont
0x0090: 7374 7265 616d 0d0a 0a0a 3c45 3a20 4d52 stream....<E>=M
0x00a0: 4242 3e3e 2d45 3e3e 0000 0000 10a0 a9d1 O/>=<E>=M
0x00b0: e930 0000 0000 0000 0000 0000 0000 0000 ..E.....
0x00c0: 0000 0000 0000 0000 0000 0000 0000 0000 ..E.....
0x00d0: 0000 4050 0000 0000 0000 e843 8000 7e43 0000 ..C...C..C
0x00e0: 0040 0000 0000 0000 0000 0000 0074 4300 0000 ..C...C..C
0x00f0: 4d 0000 0000 0000 0000 0000 0000 0000 ..C...C..C
Reserved string
Null terminating byte
Binary payload type (0x03)
Economics Binary Header - number of indicators (4 bits), here 4
Economics Binary Header - protocol version (4 bits), here 1
Indicator flags
Actual value
Economics Binary Header - previous revised value
Set mode

```

[illegible]



# 文本挖掘应用调研启示

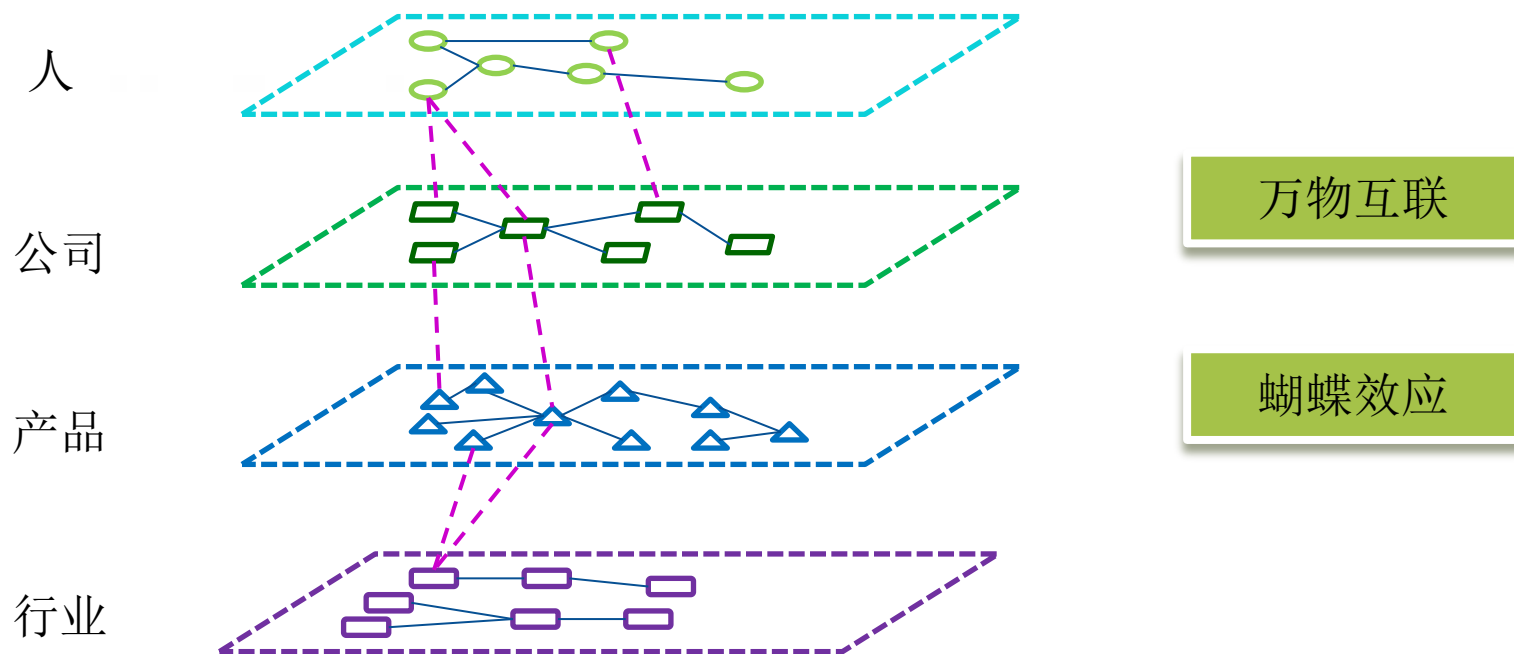
- 证券行业文本挖掘有着广泛的需求及市场前景
- 仍存在很多问题：
  - ◆ 证券行业深度依赖信息精准度 VS 文本挖掘？
  - ◆ 网络舆情的准确度、可信度
  - ◆ 数据孤岛、作用链传导
  - ◆ 大数据交互印证

# 文本挖掘应用调研启示

- 建立证券行业知识图谱是证券文本语义理解和语义搜索的关键基础技术，为证券领域文本分析、舆情监控、知识发现、模式挖掘等提供了坚实支撑
  - ◆ 行业专业中文词典
  - ◆ 业务知识分析及行业特性储备
  - ◆ 证券行业实体对象及属性的定义与提取
  - ◆ 行业内实体关联关系网络拓扑图的建立
  - ◆ 行业语义分析等等

发现潜在关联、行业推理、传导链、量化影响程度、提高分析精准度.....

# 基础工作之一——证券行业知识图谱



E.g., 突发洪水，导致一类原料短缺，这与所持股票有何关系？又将预示着哪些投资套利机会？



# 证券行业知识图谱

- 终极目标：建立一个全谱系的上市公司关联图谱



证券行业  
知识图谱

监管

追溯异动产生根源

.....

发掘打击内幕交易

市场

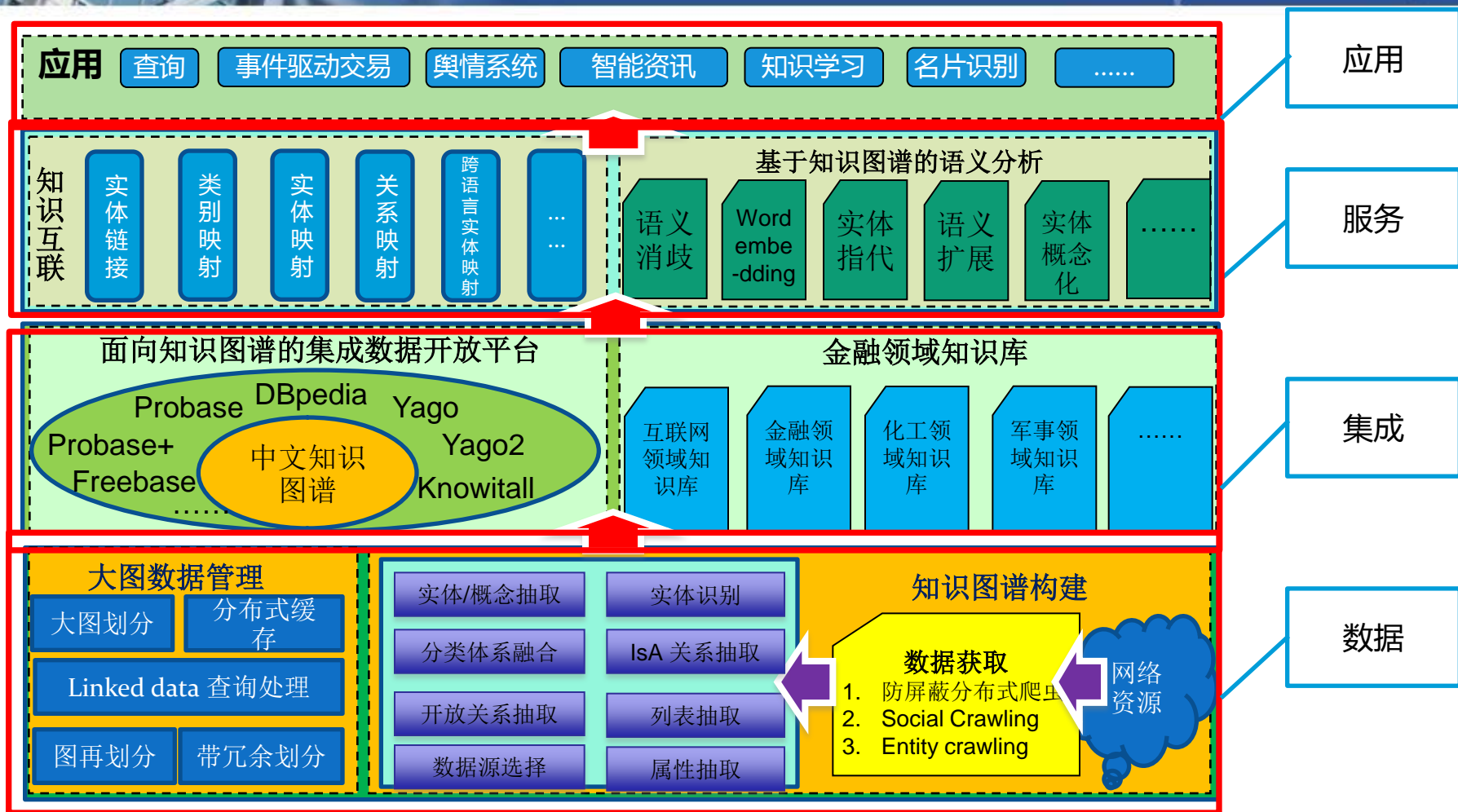
关联事件套利

.....

个性化智能化资讯

- ◆ 相关联合研究项目：证券知识图谱关键技术研究

# 知识图谱系统架构设计



# 总结与建议

- 给监管部门提供了全新的风险评估与管理手段
- 提供了高效、便捷、准确的信息情报分析方式
- 有助于还原相关金融市场事件的因果链条
- 大数据的体量、维度均为分析的准确性提供了保障

舆情监控

打击内幕  
交易

搭建大数据服  
务体系

投资者资  
格认定

加强对拟上市  
公司的审查

建立大数据指  
标指数



上海證券交易所  
SHANGHAI STOCK EXCHANGE

# 工作总结

- 发表文章：

- ◆ 《证券行业文本挖掘技术应用现状与探讨》 交易技术前沿
- ◆ 《证券行业数据及大数据服务探讨与展望》 交易技术前沿
- ◆ 《机器可读新闻技术与服务》 交易技术前沿
- ◆ 《金融时间序列数据挖掘实例分析》 交易技术前沿
- ◆ 《时间序列显著模式挖掘及其在证券市场的应用展望》 交易技术前沿等。

- 翻译：

- ◆ 《MESA：异地可复制、近实时、可扩展数据仓库》
- ◆ 《知识图谱的探索性搜索》等

- 书籍：

- ◆ 《证券市场国际化技术实践》



# 工作总结

- 参与课题：
  - ◆ “基于知识图谱的沪港两地情绪指数研究”，上证联合研究课题，课题协调人。
- 所内工作参与：
  - ◆ 参与证监会《行业信息安全规划》的制定；
  - ◆ 参与所内舆情系统的试用、问题反馈，并参与证监会互联网信息采集及舆情监测系统需求统筹工作会议；
  - ◆ 负责《交易技术前沿》的编辑、征稿、审稿、出版等工作；
  - ◆ 参与沪港通检查小组，对券商沪港通技术方面的准备情况进行了解、调研与学习；
  - ◆ 负责每周《金融科技动态》的资料采集，编辑和发布工作。



请各位老师批评指正

谢谢！