

金融知识图谱建设 及其在银行业务中的应用

梁栋

- 北京松鼠山科技有限公司
- 哈尔滨工业大学智能金融联合实验室

关于松鼠山科技



松鼠山科技，是国际领先的金融科技与人工智能领域开拓者，运用大数据及人工智能相关技术服务于**以银行业、保险业与基金业为主的金融行业，致力于推动各金融垂直场景的赋能**。松鼠山科技的核心竞争力来源于前沿扎实的科研能力和传统金融的深厚经验。

其专注于金融科技领域赋能建设，其中的**业务核心团队具有15年中国金融行业信息化建设经验**。



“松鼠山”（Squirrel Hill）位于美国**卡耐基梅隆大学**校区一侧，我们的创始人以此命名公司。





专注的业务范围



经过四年多的高速发展，松鼠山科技逐渐形成了以Beeler（结构化大数据分析）为基础，Fusion 为核心，结合海量企业工商数据、专利信息、投融资信息、舆情信息、监管数据、司法数据，打通整合金融机构第一方自有数据，凭借松鼠山科技一流的工程实现能力、算法设计能力和模型建设能力，以多年金融行业的从业经验为指导，覆盖“风险控制”，“精准营销”等一系列金融关键环节的应用服务体系。

风险
控制

精准
营销



CONTENTS



CONTENTS

- 01 **知识图谱的建设背景**
- 02 知识图谱的构建过程
- 03 知识图谱的关键（知识库建设）
- 04 知识图谱的应用场景

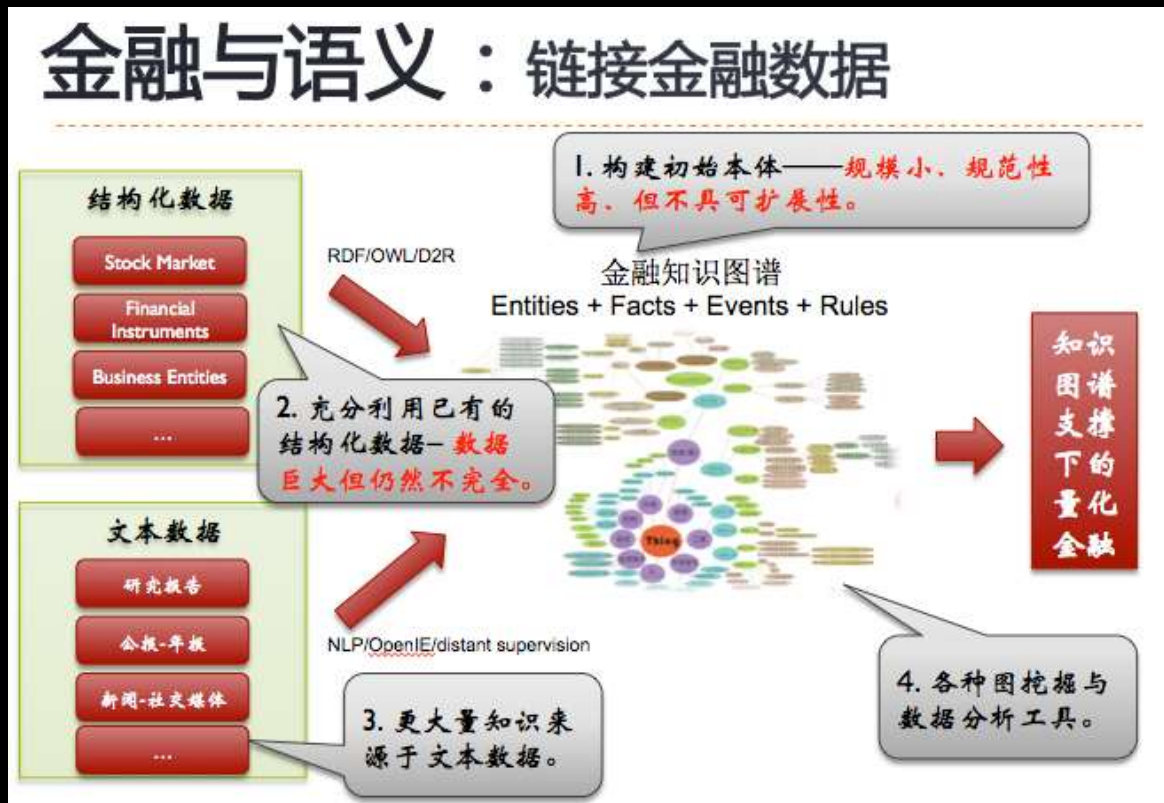
知识图谱的建设背景

– **客户认知**的问题：在外界市场纷繁的变化中，哪些因素发生了变化，对我行的客户造成了影响、价值（风险）有多大？

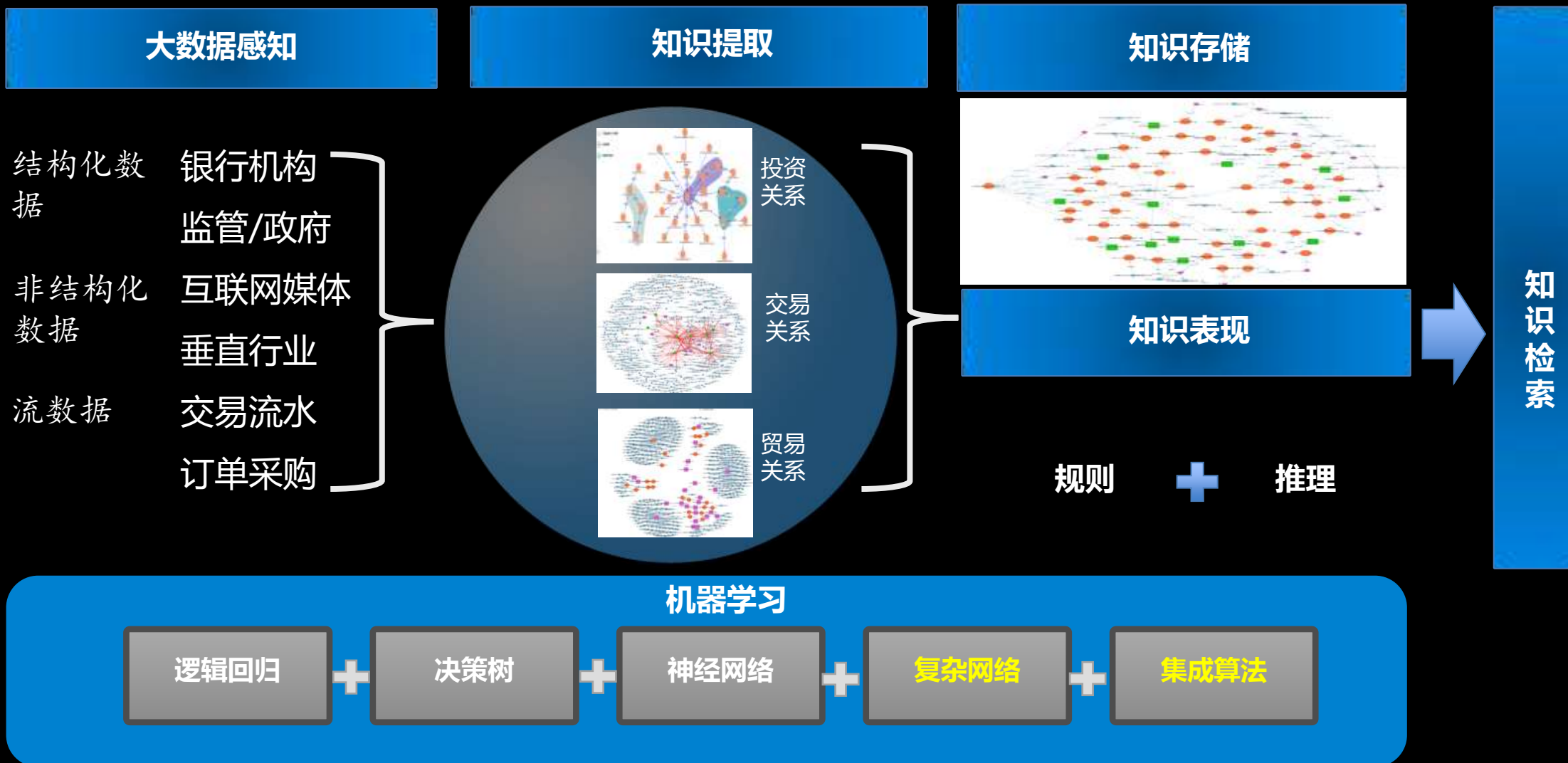
知识图谱：提供了从企业群体关系角度分析问题的能力。进一步可以拓展企业和行业、企业和客户、企业和市场价格等要素的关联传导。

大数据：完成从原始信息到企业关系、风险特征、风险事件的提取和整合。

机器学习：建立预测模型，实现风险识别、传导和预测的能力。



知识图谱的建设背景





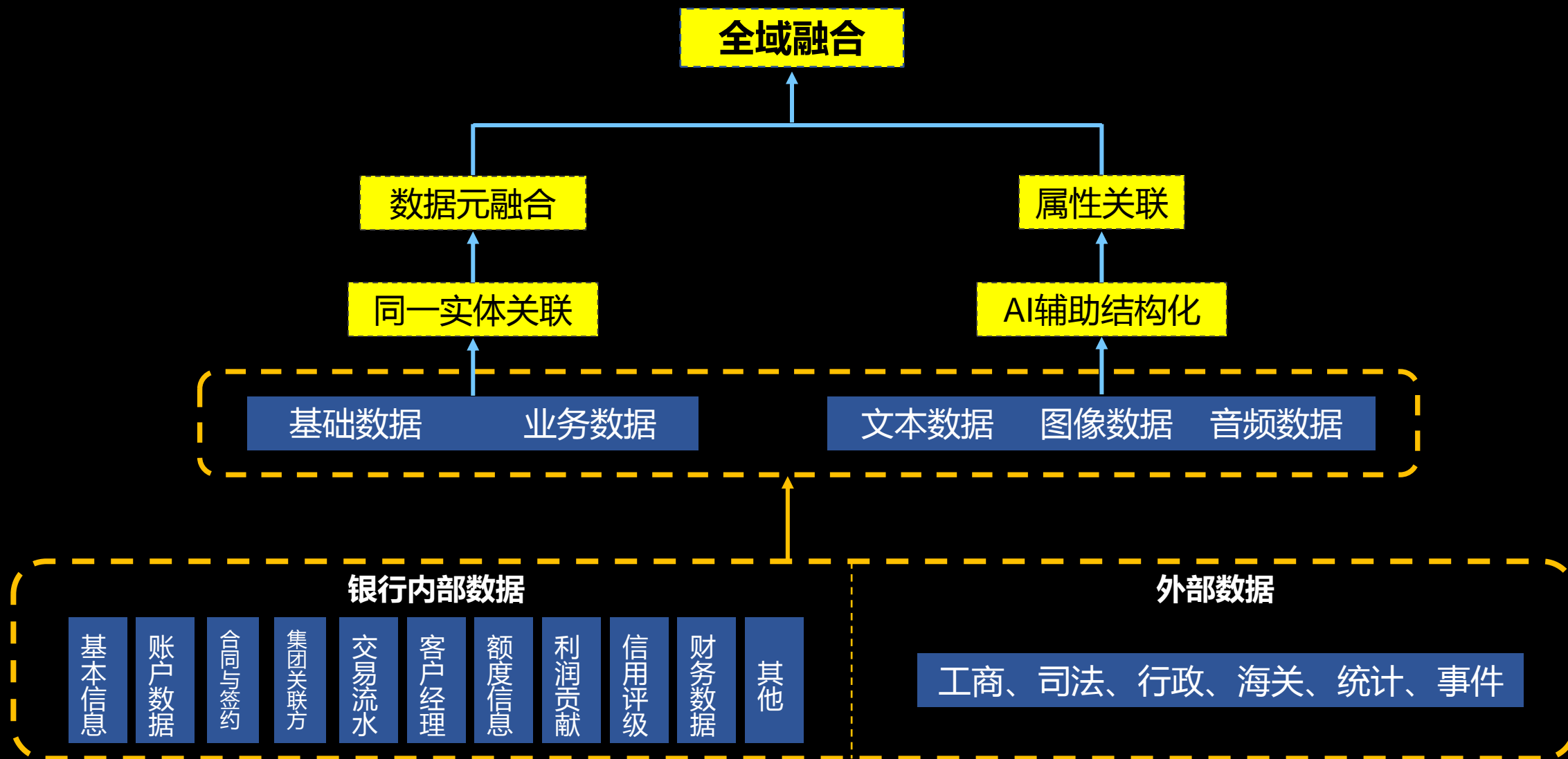
CONTENTS



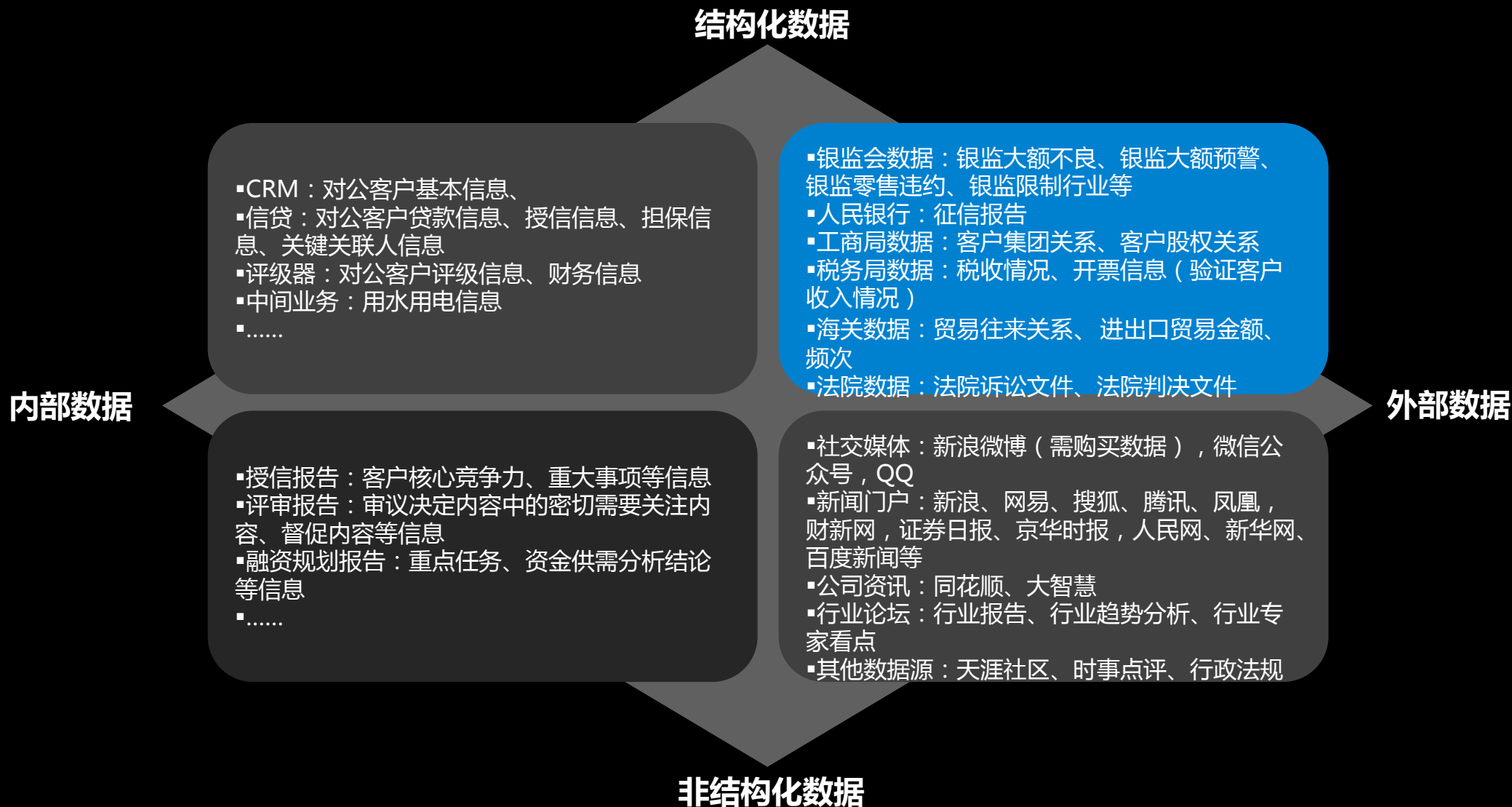
CONTENTS

- 01 知识图谱的建设背景
- 02 **知识图谱的构建过程**
- 03 知识图谱的关键（知识库建设）
- 04 知识图谱的应用场景

关键路径：数据引入（策略）



关键路径：数据引入（范围）



构建知识图谱-单一客户谱系（数据与算法）

单一客户谱系——基于大数据的精准场景服务



客户尽职调查

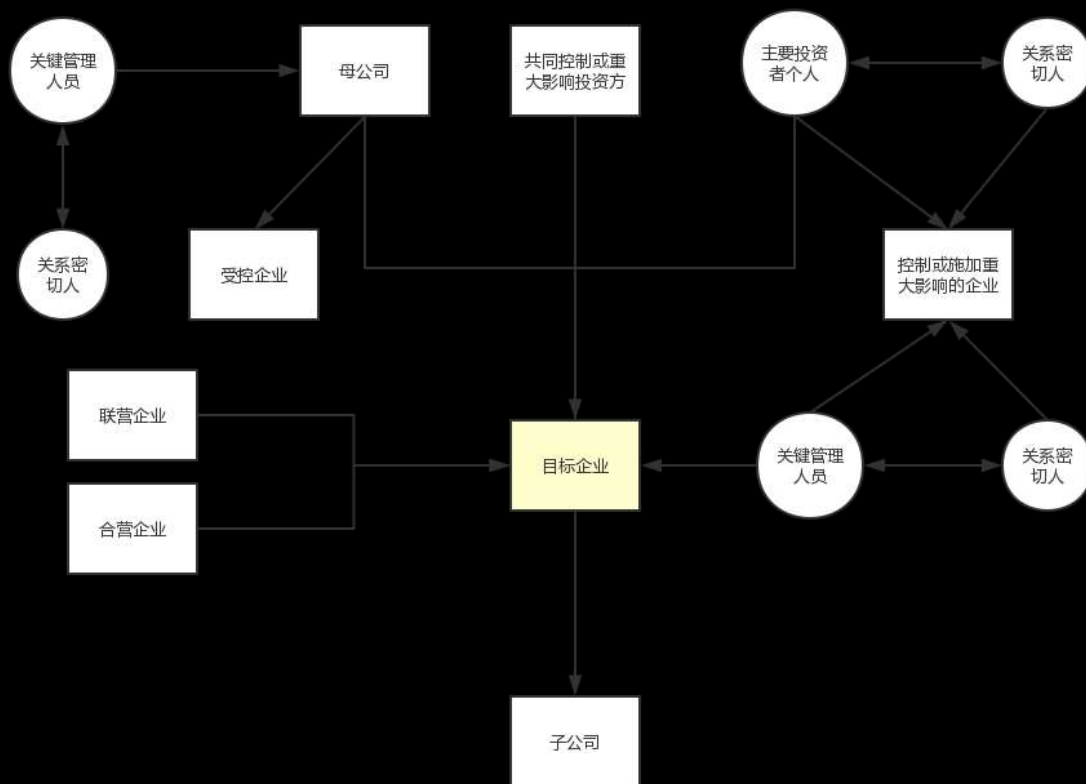
贸易背景真实性调查

财务会计审计

风险预警监控

信贷准入

根据《企业会计准则第36号——关联方披露》认定的关联方



构建知识图谱-图计算模型

-基于图连通性的关联模型

- **单源最短路径模型**

在复杂关系图中，找到两个不同节点最近的关联路径

- **基于权重的最强关联搜索**

在有权重的关系图中，找到两个不同节点多条路径中，权重和最大的一条路径（最重要）

- **多种社团发现模型**

在复杂关系图中，找到关联关系联系紧密的一个节点群，主要用于 欺诈团伙发现，利益共同体发现，企业集团/系发现等场景

- **桥发现模型**

在复杂关系图中，找到多个关联团中的关键联系节点

- **传导计算模型**

在复杂关系图中，某一个节点发生状态改变，根据关联关系，寻找状态改变影响的范围和关联节点影响的程度。

- **自主设计的带参数的局部结构检索模型**

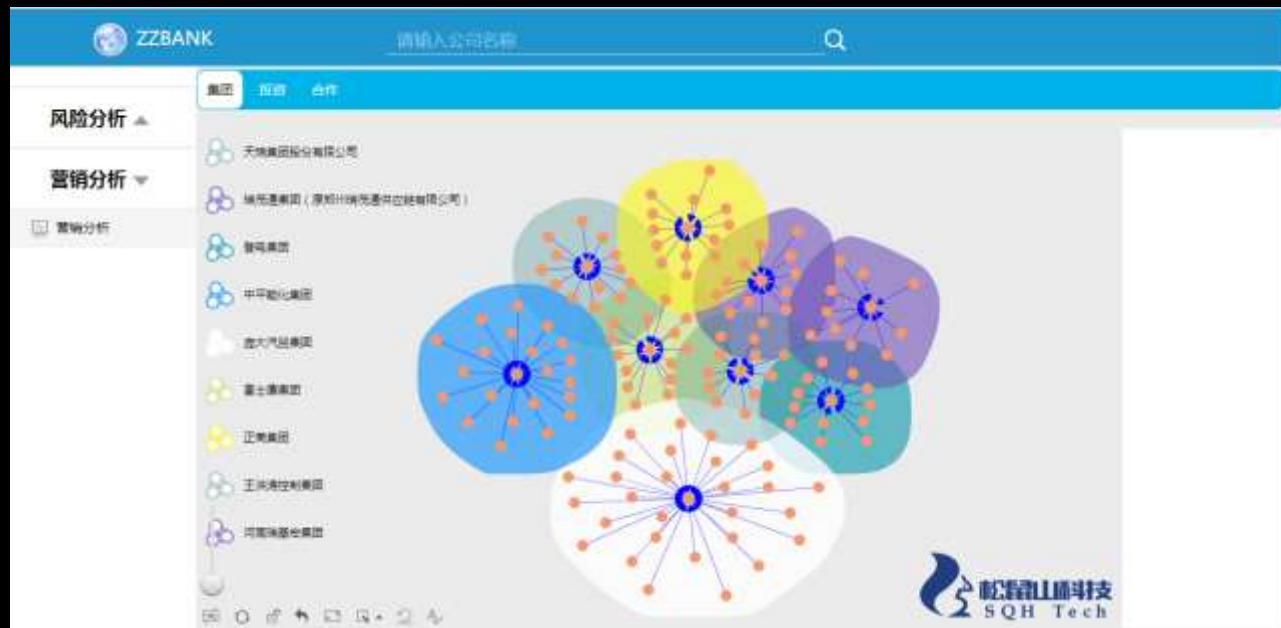
比较难于理解，举例说明：查找销售额大于10（数值用户定义）亿的企业的客户，并且必须是高科技企业。

- **模式匹配模型**

在图中查找某种固定的拓扑结构，场景：循环担保发现等。

- **紧密中心度模型**

发现局部网络中，发现最中心的节点（最容易到达其他所有节点）



构建知识图谱-企业集团谱系（数据方法体系3）

对于集团客户的核定标准银监会《银行集团客户授信业务风险管理办法》给出了可参考借鉴的定义。

我们认为集团客户的形式复杂多样，在界定其具体范围时，应自始至终遵循实质重于形式和“穿透制”的原则，从控制与施加重大影响的视角出发，重点分析隐匿在复杂关系链路中的施控与被控关系。

《中国银行业监督管理委员会关于修改（商业银行集团客户授信业务风险管理指引）的决定》——2010年第4号

第一章 总则

第三条 本指引所称集团客户是指具有以下特征的商业银行的企事业法人授信对象：

- （一）在股权上或者经营决策上直接或间接控制其他企事业法人或被其他企事业法人控制的；
- （二）共同被第三方企事业法人所控制的；
- （三）主要投资者个人、关键管理人员或与其近亲属（包括三代以内直系亲属关系和二代以内旁系亲属关系）共同直接控制或间接控制的；
- （四）存在其他关联关系，可能不按公允价格原则转移资产和利润，商业银行认为应当视同集团客户进行授信管理的；

前款所指企事业法人包括除商业银行外的其他金融机构。

商业银行应当根据上述四个特征结合本行授信业务风险管理的实际需要确定单一集团客户的范围。

指引
解读

控股母公司

控股子公司

共同被第三方控制

关键管理人员或
与其近亲属

施加重大影响

遵循
实质
重于
形式
的原则



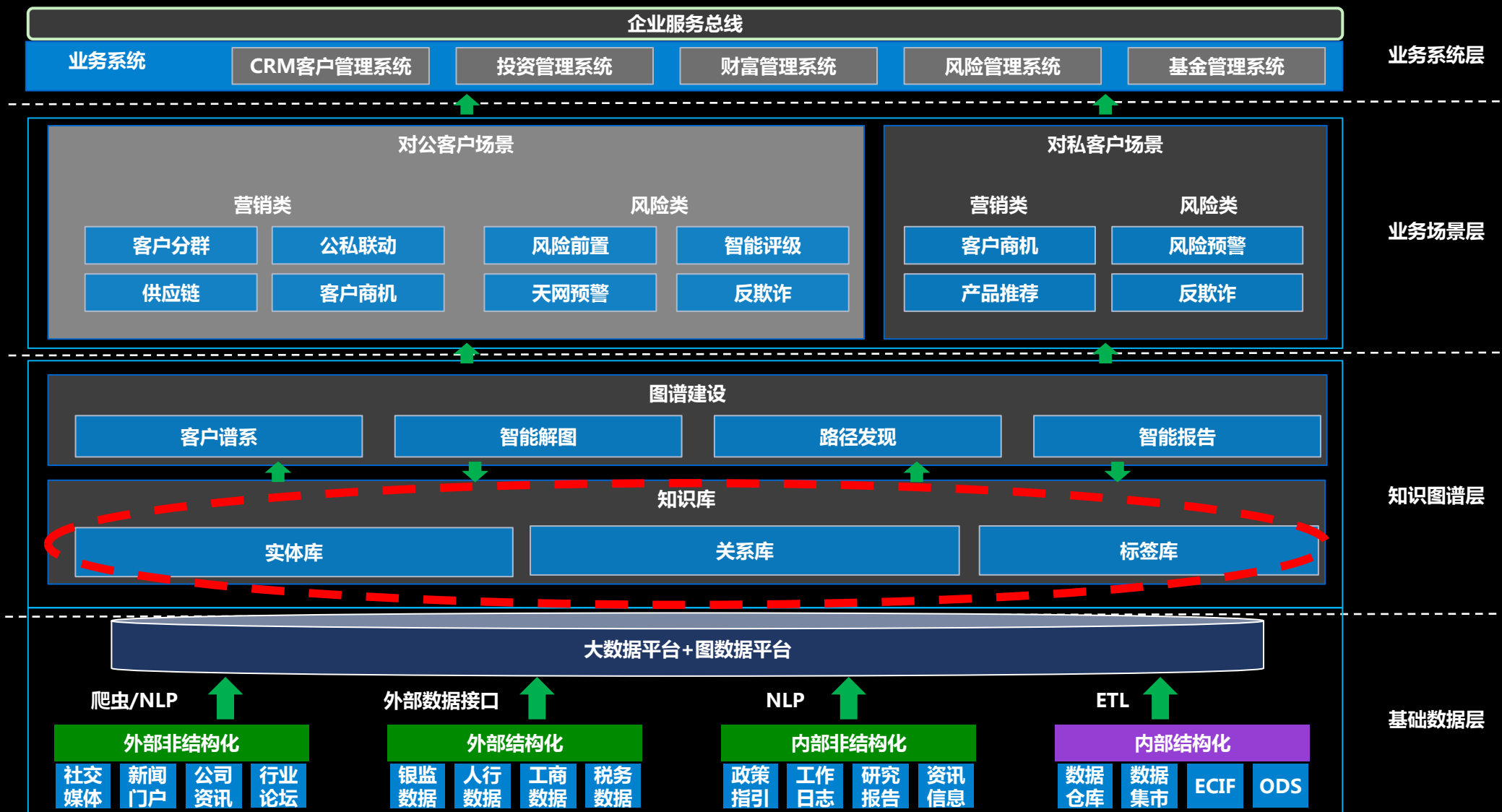
CONTENTS



CONTENTS

- 01 知识图谱的建设背景
- 02 知识图谱的构建过程
- 03 **知识图谱的关键（知识库建设）**
- 04 知识图谱的应用场景

知识图谱的关键路径（概览）



知识库建设-标签库体系（1）

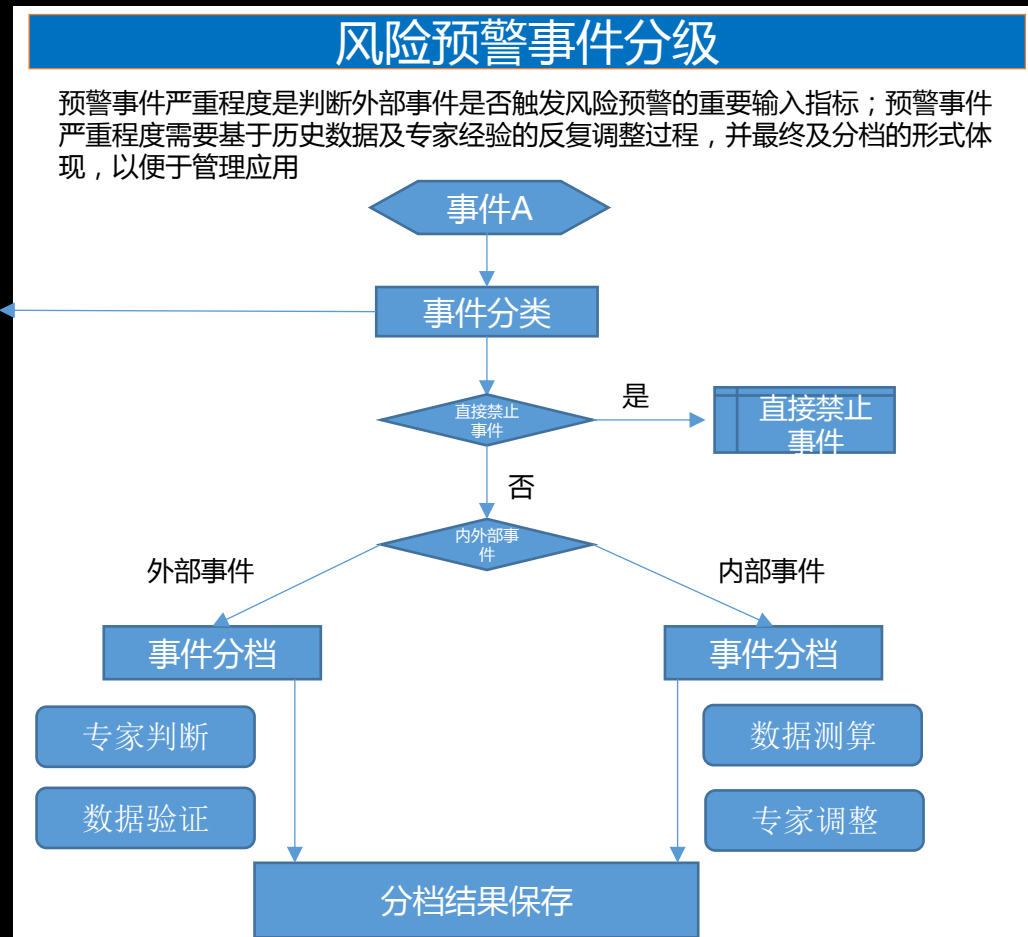
企业标签库体系（200+）

基础信息

生命周期	公司轮廓	产品持有	交易行为	价值贡献	经营特征	往来关系	风险特征
生存周期	企业性质	产品偏好	周期频率	直接贡献	贸易数据	资金往来特征	风险评级
年限分层 发展阶段	国有企业	产品持有偏好 产品风险偏好 产品流动偏好 产品持有组合 持有产品数量 持有产品额度 首次购买产品	客户购买频率分 层	存款贡献度 贷款账户贡献度 债券承销贡献度 票据承销贡献度 远期结售汇贡献度 资产托管贡献度	进口总额分层 出口总额分层 主要贸易对象国 进出口主要商品	往来户标识 是否我行客户 资金方向 结算工具 结算金额分层 前N大上下游企业	客户评级 债项评级
价值周期	企业规模		交易场景		经营活动		违约历史
客户分层 资产分层 负债分层 交易分层 核心客户 潜力升优质 普通升潜力 优质向下迁徙	大型企业		交易行为偏好 转账行为偏好 交易时间偏好				逾期次数 违约次数
	治理结构		业务客群				渠道偏好
	双层结构 <td rowspan="2">主账户客群 存款客群 融资客群 中间业务客群 小微客群</td> <td>渠道使用频率</td> <td>银监大额预警 银监大额不良</td>	主账户客群 存款客群 融资客群 中间业务客群 小微客群	渠道使用频率	银监大额预警 银监大额不良			
集团标识	转账汇款		潜在贡献	分布特征	外部事件		
交往周期 <td>所属集团<td>频繁大额转出</td><td>资产收入价值分层 抵押物价值 担保合同金额 贷款授信额度</td><td>往来户行业分布 往来户地域分布 与我行核心企业往来 特征 与集团内关联企业往 来特征 上下游往来户变化特 征</td><td>高管潜逃</td></td>	所属集团 <td>频繁大额转出</td> <td>资产收入价值分层 抵押物价值 担保合同金额 贷款授信额度</td> <td>往来户行业分布 往来户地域分布 与我行核心企业往来 特征 与集团内关联企业往 来特征 上下游往来户变化特 征</td> <td>高管潜逃</td>	频繁大额转出				资产收入价值分层 抵押物价值 担保合同金额 贷款授信额度	往来户行业分布 往来户地域分布 与我行核心企业往来 特征 与集团内关联企业往 来特征 上下游往来户变化特 征
	关联企业	资金流量		定期全额转出 <th>经营活动事件</th> <th>关联关系</th>	经营活动事件	关联关系	
	客户发展状态 客户关系时间 存款持有周期 贷款持有周期 理财持有周期	控股比例 <td rowspan="2">收款金额分层 下月到期资产分层 下月到期存款分层</td> <th>交易趋势</th> <td rowspan="2">重大战略发布重大投 资项目 经营业绩预警 上市公司违规</td> <td rowspan="2">供应链圈子 担保圈子</td>	收款金额分层 下月到期资产分层 下月到期存款分层	交易趋势	重大战略发布重大投 资项目 经营业绩预警 上市公司违规	供应链圈子 担保圈子	
上市企业		金额变化趋势					
	是否上市企业						

知识库建设-事件库体系

客户风险事件预警中事件分类与分级是重心，不同类别事件对客户影响不同，同一类型事件，不同的严重程度对客户的影响不同。





CONTENTS



CONTENTS

- 01 知识图谱的建设背景
- 02 知识图谱的构建过程
- 03 知识图谱的关键（知识库建设）
- 04 **知识图谱的应用场景**

输入

根据单一关联生成企业图谱
根据各种关联关系生成企业图谱

关联关系

股权关系
交易关系
担保关系
行业关系
.....

传导

在关联关系图谱中，部分信用主体发生变化，该类变化沿着图谱网络以一定效能进行传导

传导前

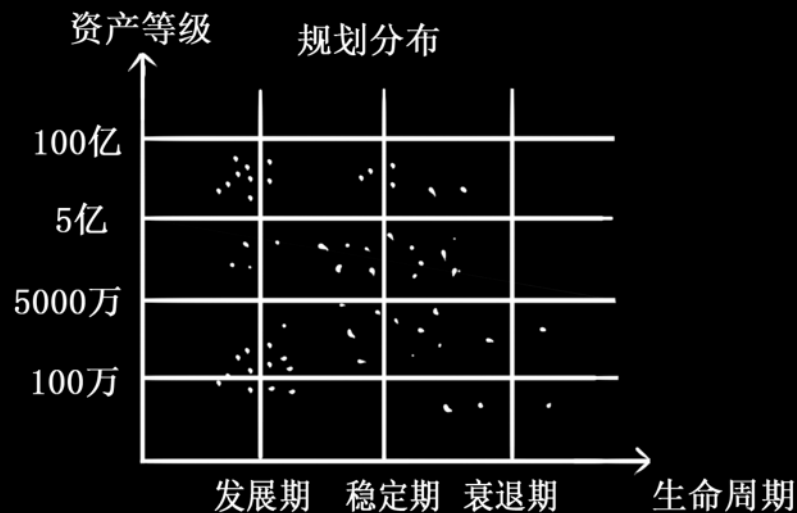
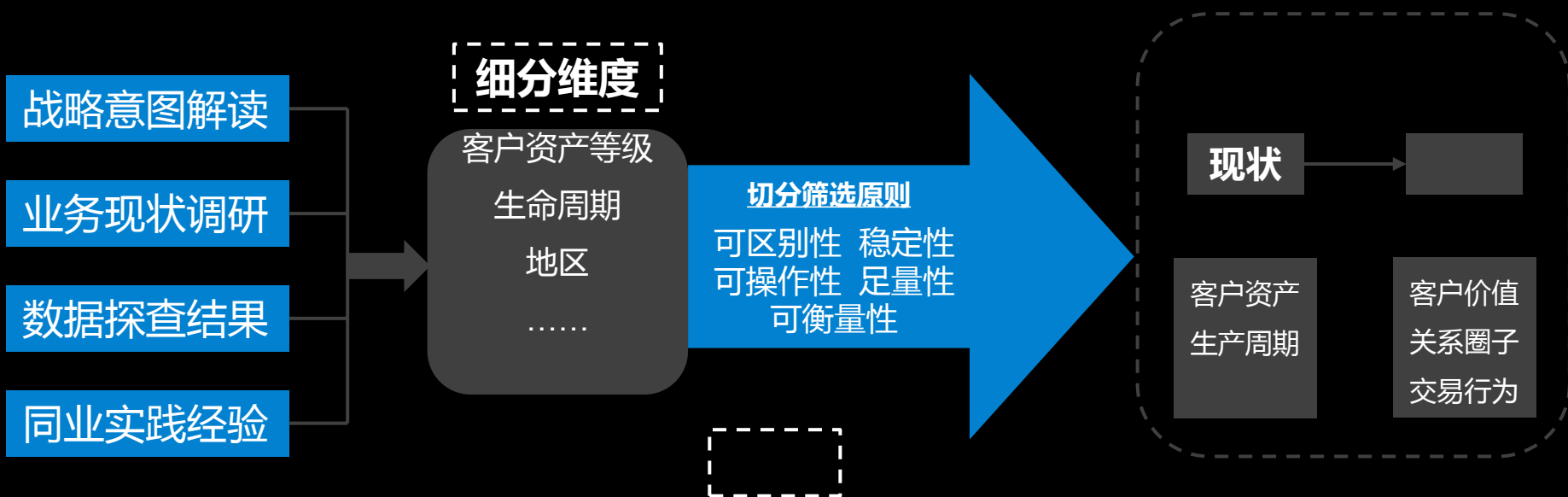
A
B*3
C
D
E
F*4

输出

根据传导结果，以资产规模，违约逾期，营业收入等维度反映并输出企业传导后风险情况

传导后：信用主体风险情况

A*1
B*2
C*1
D*2
E*5
F*6





我们的合作伙伴



Solution Partner

Neo4j亚太区合作伙伴
中国金融行业总代理商

- 结合世界上最领先最流行的图形数据库产品，打造银行、PE/VC、保险、财富管理、证券等细分金融行业的专业知识图谱，建立起各实体关联图，并将其直接关联、间接关联的各种实体、概念相联系，帮助各参与者洞察市场脉搏。
- 目前为止，中国金融行业企业版用户（中国平安集团）的唯一服务商。



关系数据库
AgensGraph(PostgreSQL)

原生单机存储

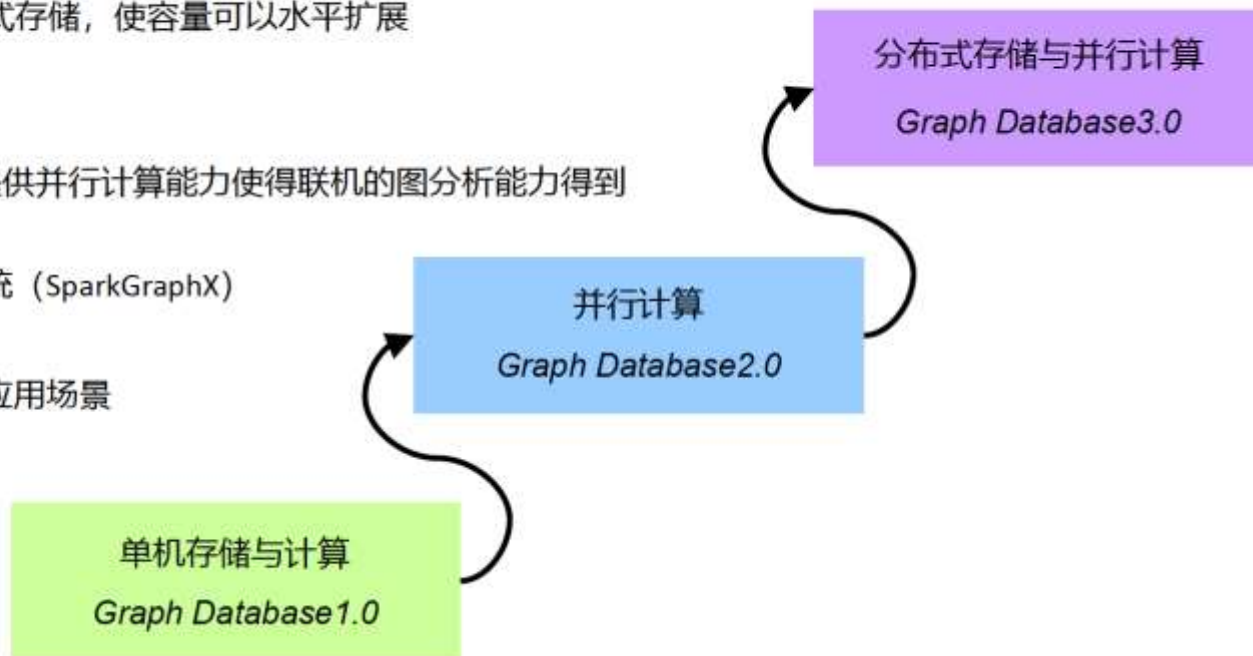
Neo4j

分布式存储
DSE(Cassandra)
JanusGraph/Titan

原生分布式存储
ArangoDB(Raft + RocksDB)
Dgraph(Raft + Badger)
TigerGraph



- 分布式存储
 - Neo4j 存储受限于单机容量
 - 基于HBase/Cassandra或原生分布式存储, 使容量可以水平扩展
- 计算并行化
 - JanusGraph 单点计算
 - ArangoDB / Dgraph / TigerGraph ,提供并行计算能力使得联机的图分析能力得到大幅度提升。
 - 并行计算能力不再局限于离线系统 (SparkGraphX)
- 数据一致性 (Raft)
 - 提供了更好的数据一致性, 拓宽应用场景
- 底层存储引擎优化
 - HBase -> RocksDB -> Badger(SSD)
 - 大幅提升数据更新及查询能力





感谢您的聆听