

因唯安全
所以信赖
唯与同行,智御未来

京东反刷单系统

主讲人: 寿如阳

京东广告数据部架构师

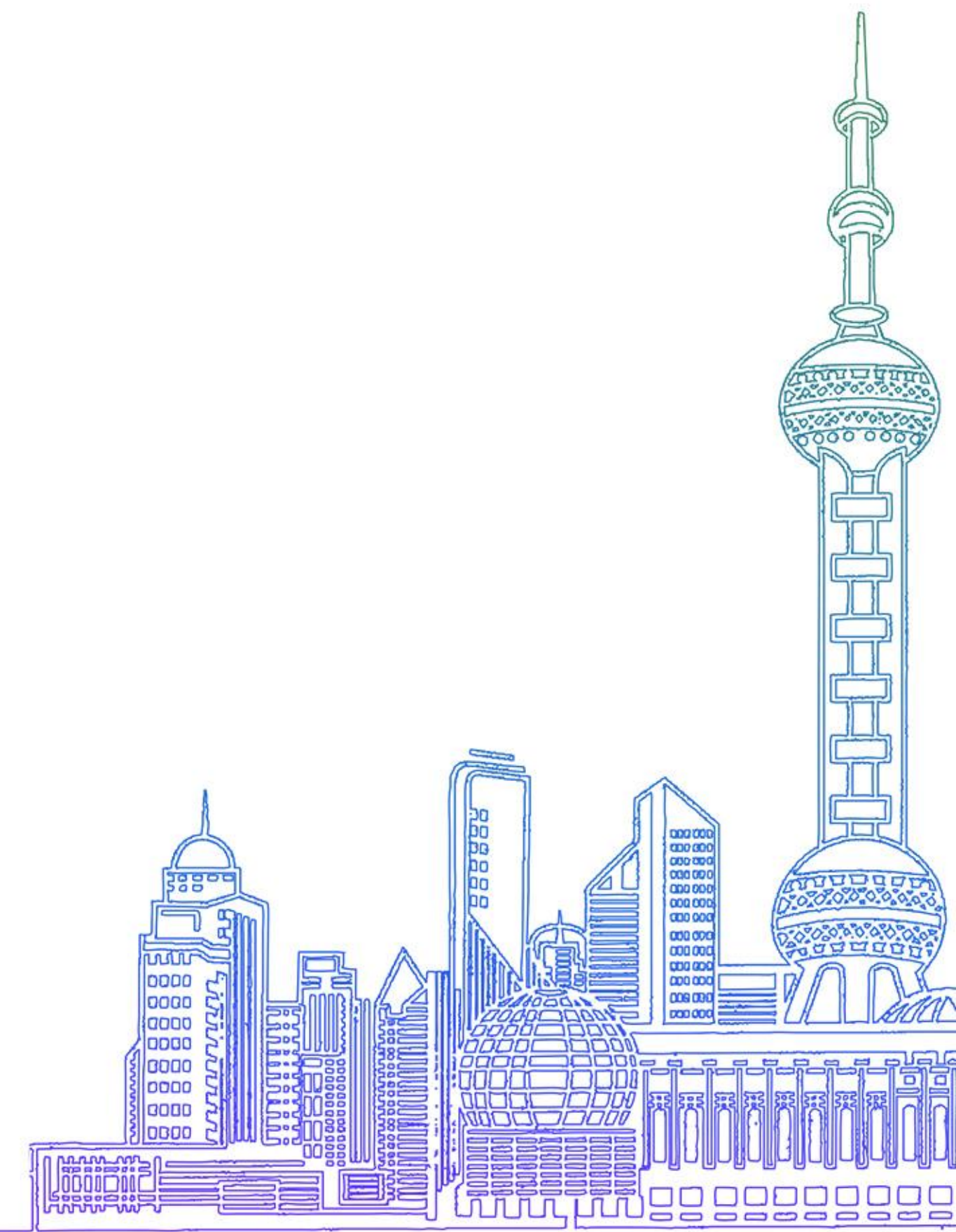
2018 唯品会第三届互联网电商安全峰会

2018 vip.com third Internet ecommerce Security Summit

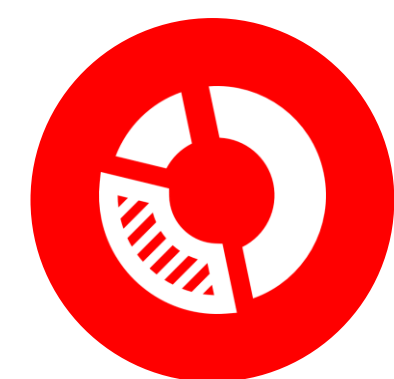
2018-5-5 上海



反刷单的需求与挑战



虚假交易的危害



商业分析

- 扰乱销售数据、转化率、毛利率等商业分析指标，歪曲业务增长真实水平，误导商业决策



客户体验

- 伪造销售数据和评价，骗取消费者对商品、品牌的认可
- 扰乱搜索排名和推荐算法的数据，使其排序失去客观性，对消费者失去参考价值



电商生态

- 在商家中引发刷单竞赛，破坏平台生态，导致商家流失
- 刷单成本廉价，妨碍正规营销业务发展，例如在线广告



刷单产业现状



规模化、市场化

- 出现各种刷单公司、刷手平台等有组织有纪律的刷单团伙，从业人员众多，产业日趋成熟
- 参与者分工明确，市场细分化，既有专营账号、软件等刷单素材，也存在提供一站式服务



成为营销、赚钱手段

- 直接针对评价、搜索排序、推荐等流量入口作弊，ROI高
- 销量、评价作弊以外可兼顾“薅羊毛”，赚取平台优惠
- 渗透各类综合、垂直电商、以及O2O平台

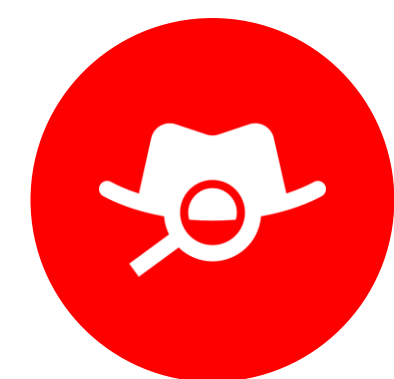


手法隐蔽、逼真、多变

- 刷单软件和工具，提高刷单效率和刷手反侦察能力
- 「人肉」刷单，模拟真实客户行，使得作弊订单的用户行为更接近正常用户
- 刷手社区互通躲避刷单检测的心得



反刷单的挑战



多维度数据引证

- 基于单个行为特点的识别方法，面对逼真的刷手行为日渐困难，需要多种维度数据上深入挖掘实体信用指标作为依据



策略的敏捷迭代

- 适应刷单手法的变化，决策识别系统能够预警并演进



快速、准确与高召回

- 对层出不穷的刷单手段，既要抓得多，又要抓得准，还要抓得快



系统需求



Fundamental



Adaptability



Reproducibility



Flexibility



Customizability



系统需求



Fundamental



Adaptability



Reproducibility



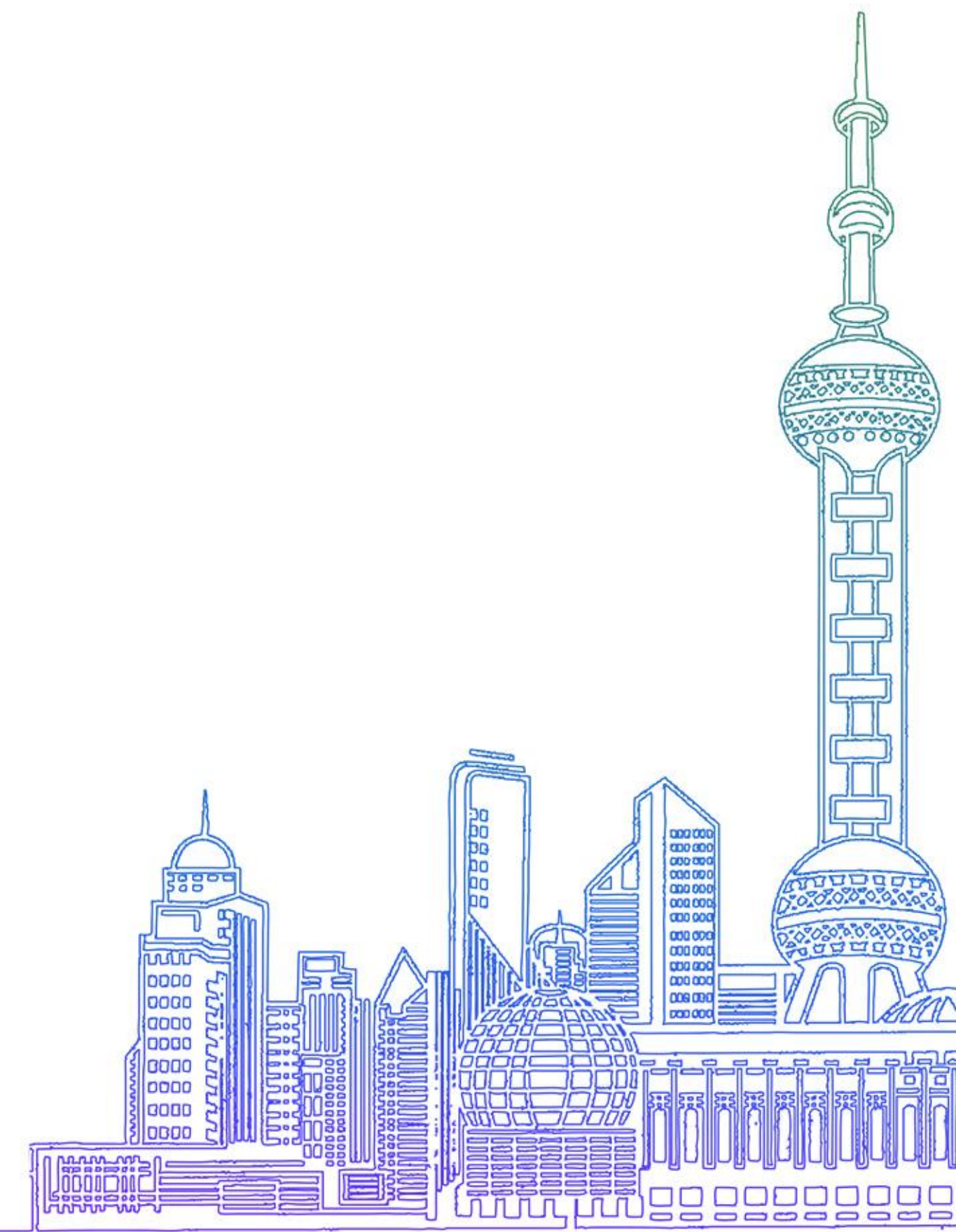
Flexibility



Customizability

《 系统基本需求 》

- 高可用性
- 可扩展性
- 低延迟



系统需求



Fundamental



Adaptability



Reproducibility



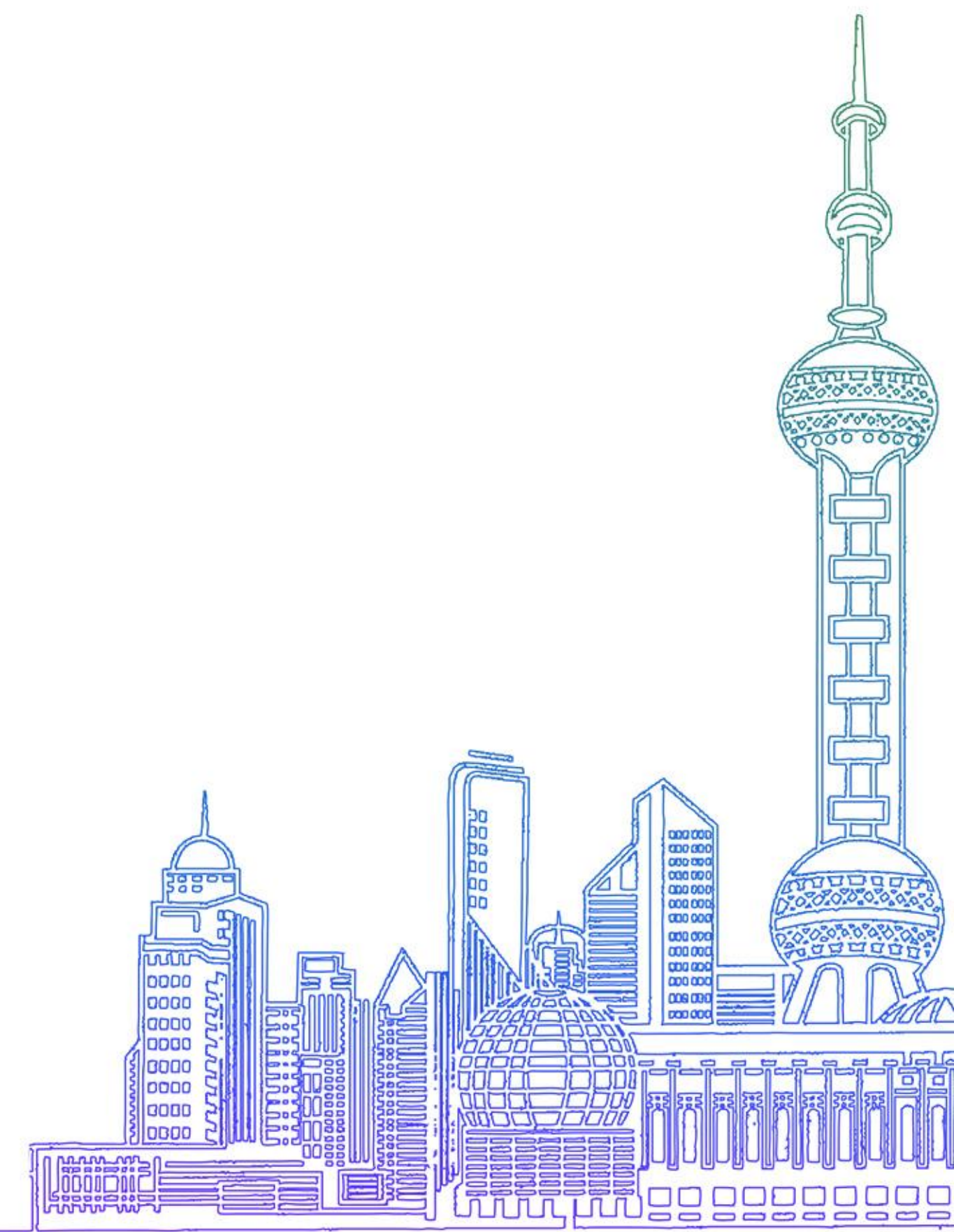
Flexibility



Customizability

适应多种 数据源

- 多种业务类型：订单、账户、支付、物流、评论……
- 不同数据形式：集群数据、流式数据、云端存储
- 业务数据的持续变化



系统需求



Fundamental



Adaptability



Reproducibility



Flexibility



Customizability

过程结果
可复现

- 需要保留数据历史快照，以便回溯算法过程
 - 业务数据和特征数据
 - 策略及系统（模型规则、参数、代码、配置）



系统需求



Fundamental



Adaptability



Reproducibility



Flexibility



Customizability

灵活的 算法系统

- 可扩展：能够随时上下线模型和规则，便于算法迭代，响应突发业务需求
- 松耦合：通用、稳定的模型算法与有针对性的场景的、易变的业务逻辑分离



系统需求



Fundamental



Adaptability



Reproducibility



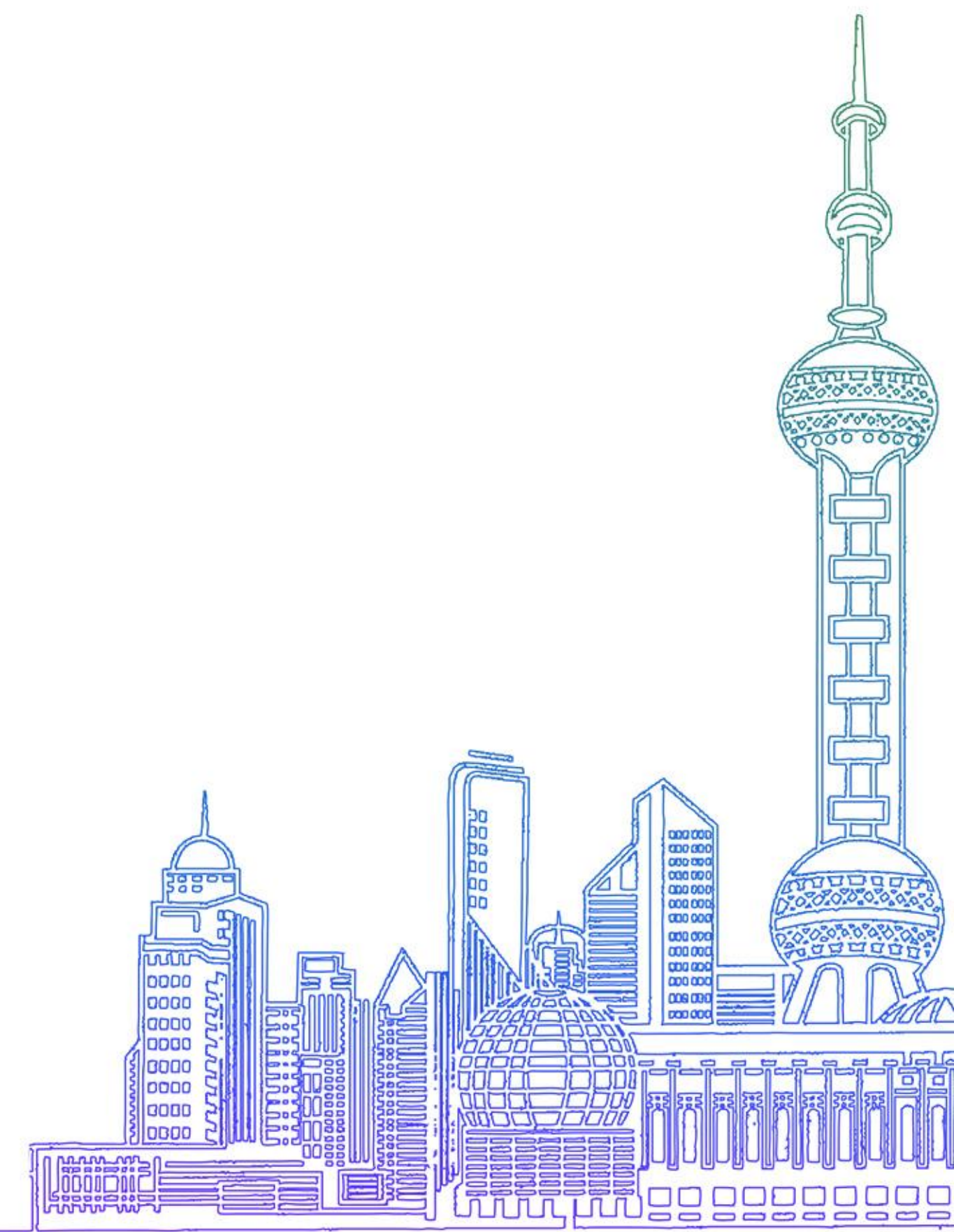
Flexibility



Customizability

定制化 下游服务

- 根据下游应用业务特点和对准召的需求定制风险指标
- 除反刷单外，用于建立用户、商家信用评分、商品质量监控等风控体系指标



数据特点

反刷单：在较长的时间跨度上，从海量持续变动的数据中挖掘刷单行为痕迹

- 生命周期长：从用户产生消费冲动到对商品发表评论，一个订单关联到的数据跨度可长达数周甚至数月
- 数据种类多：日志、买卖方属性、商品属性、交易属性、支付、物流等等数据
- 数据多变：在订单生命周期内交易数据的变动是十分常见的

搜索



购买



支付



运输



关注



结算



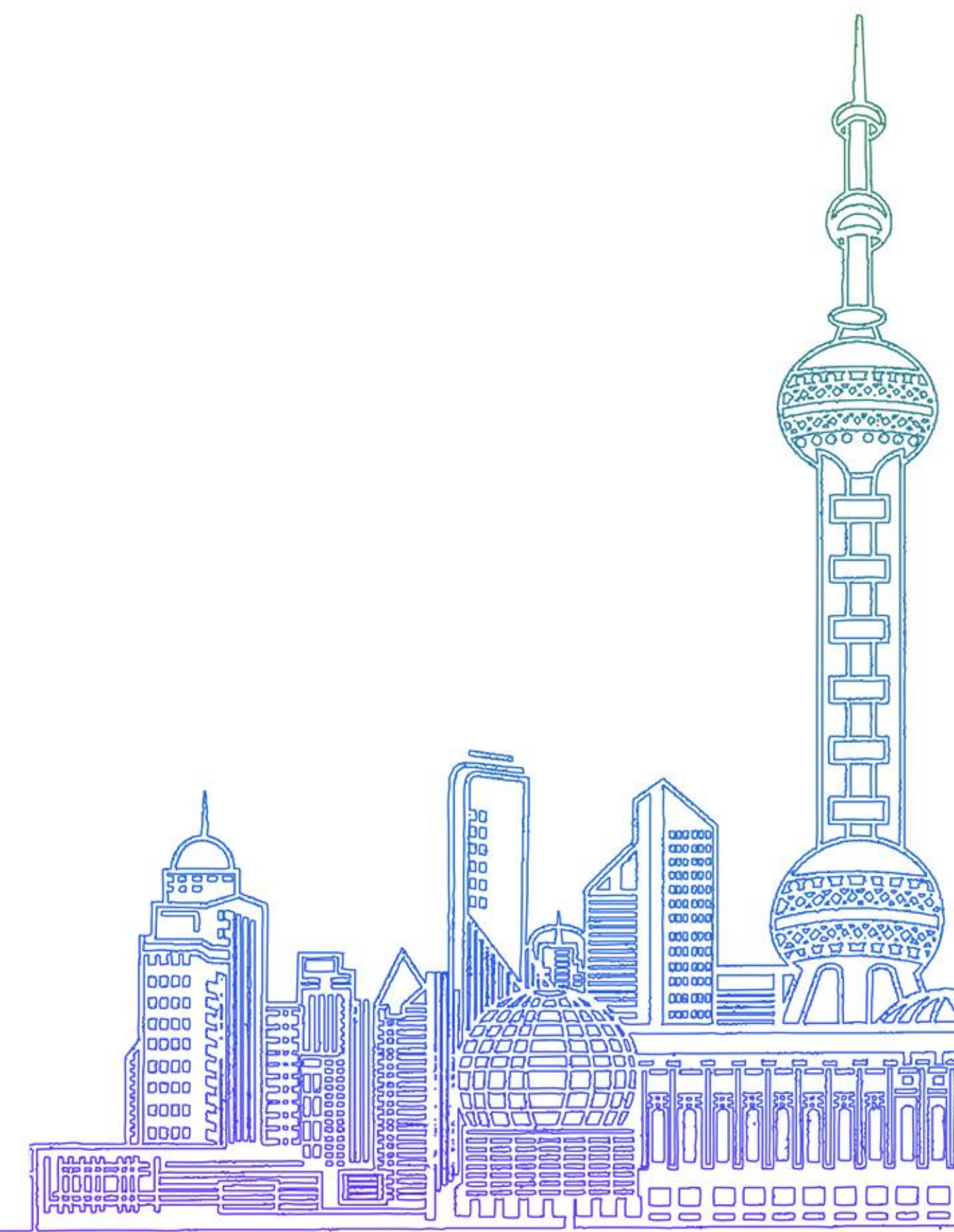
发货



评价



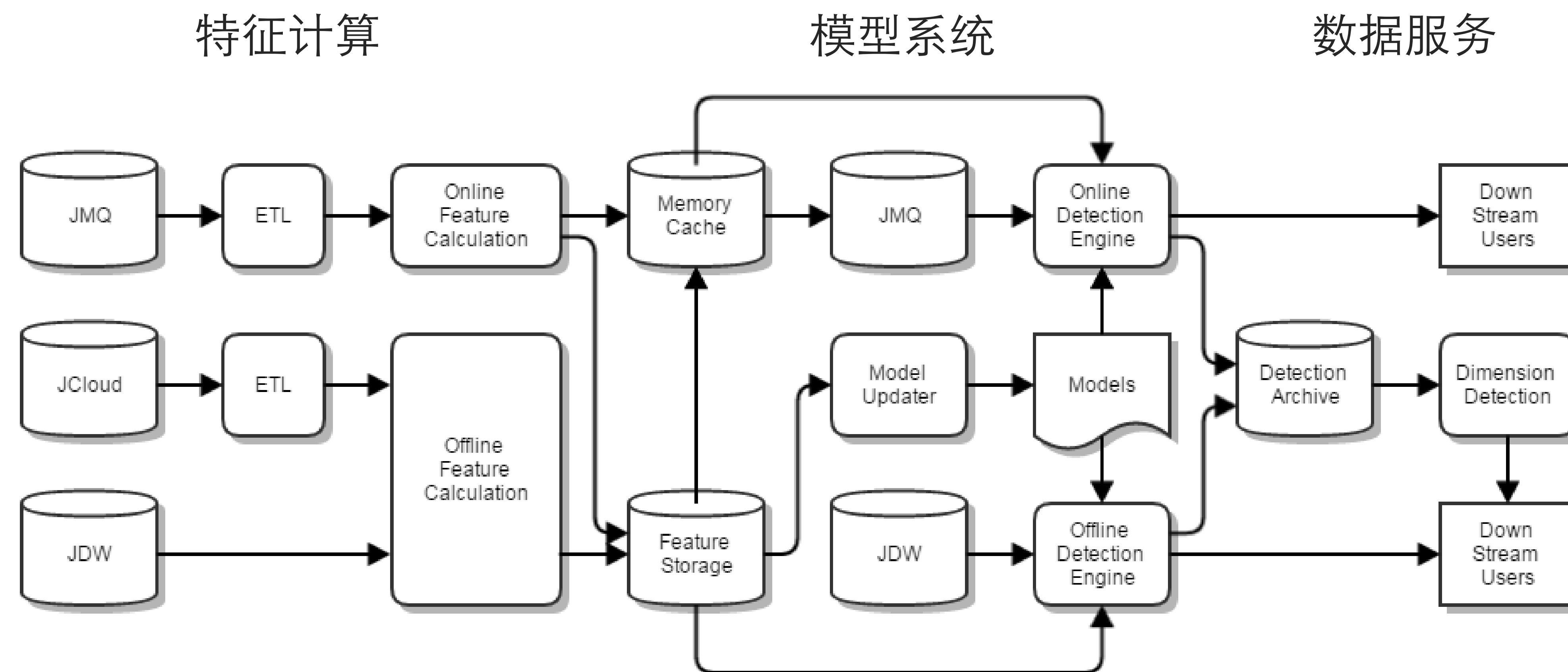
系统架构设计实践



京东反刷单系统架构

Hadoop Stack + Spark

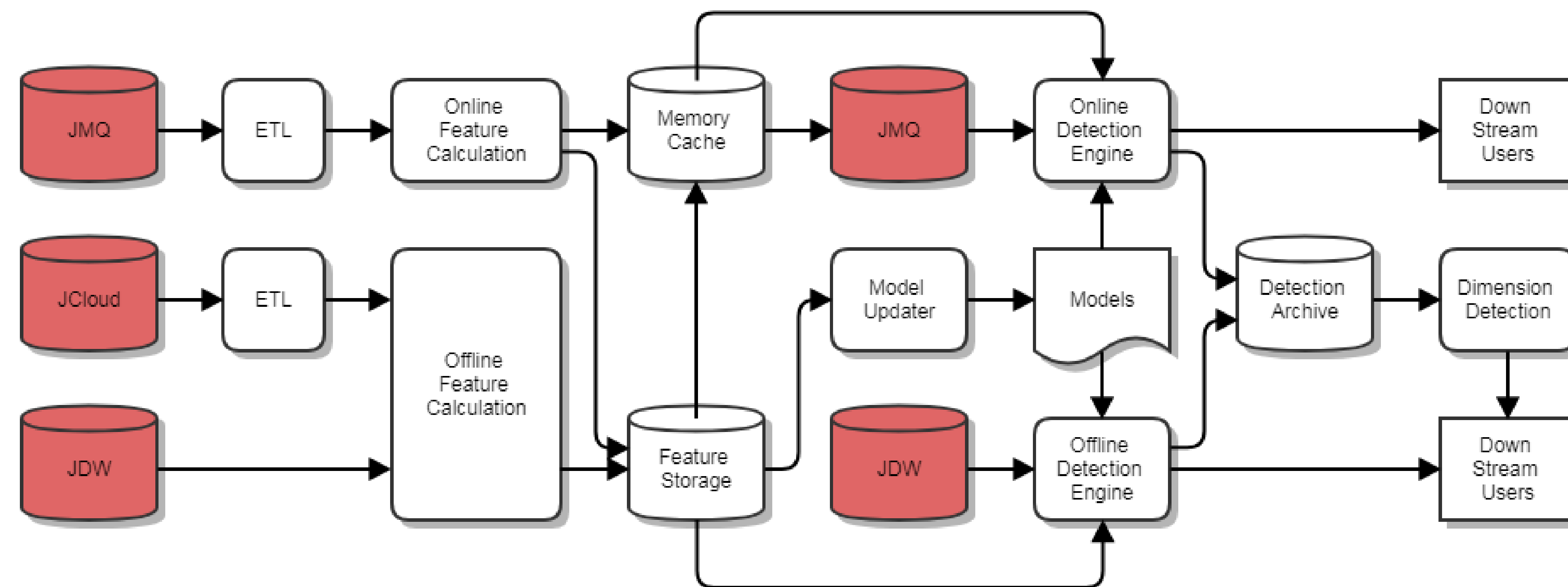
- 根据加工数据源和数据计算任务的特点选择适合的数据处理技术
- 精简选择，用尽可能少的框架的解决复杂的需求
- 紧跟大数据处理技术发展的最近成果



京东反刷单系统架构

数据预处理

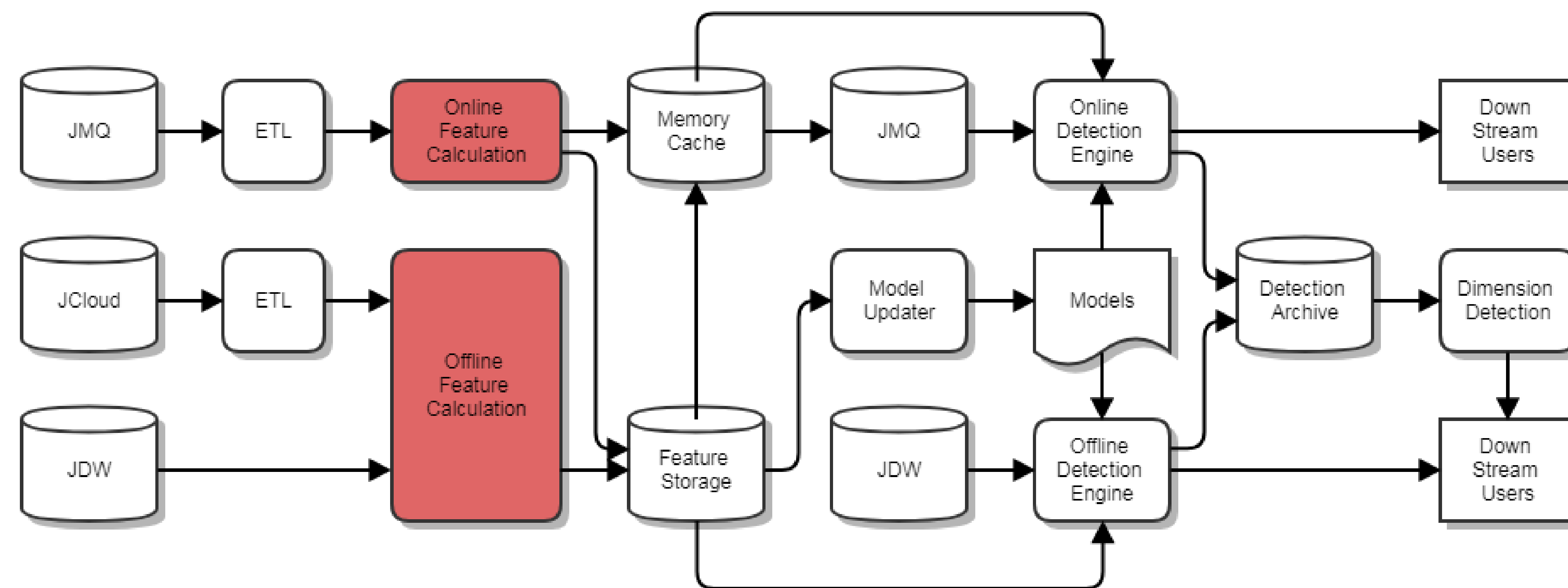
- 离线数据
 - 数据源：JDW（HDFS + Hive）和京东云
 - ETL：Hive + Pig
- 实时数据
 - 数据源：JMQ（Kafka）
 - ETL：Camus + Kafka
- 作业调度：Oozie



京东反刷单系统架构

特征计算

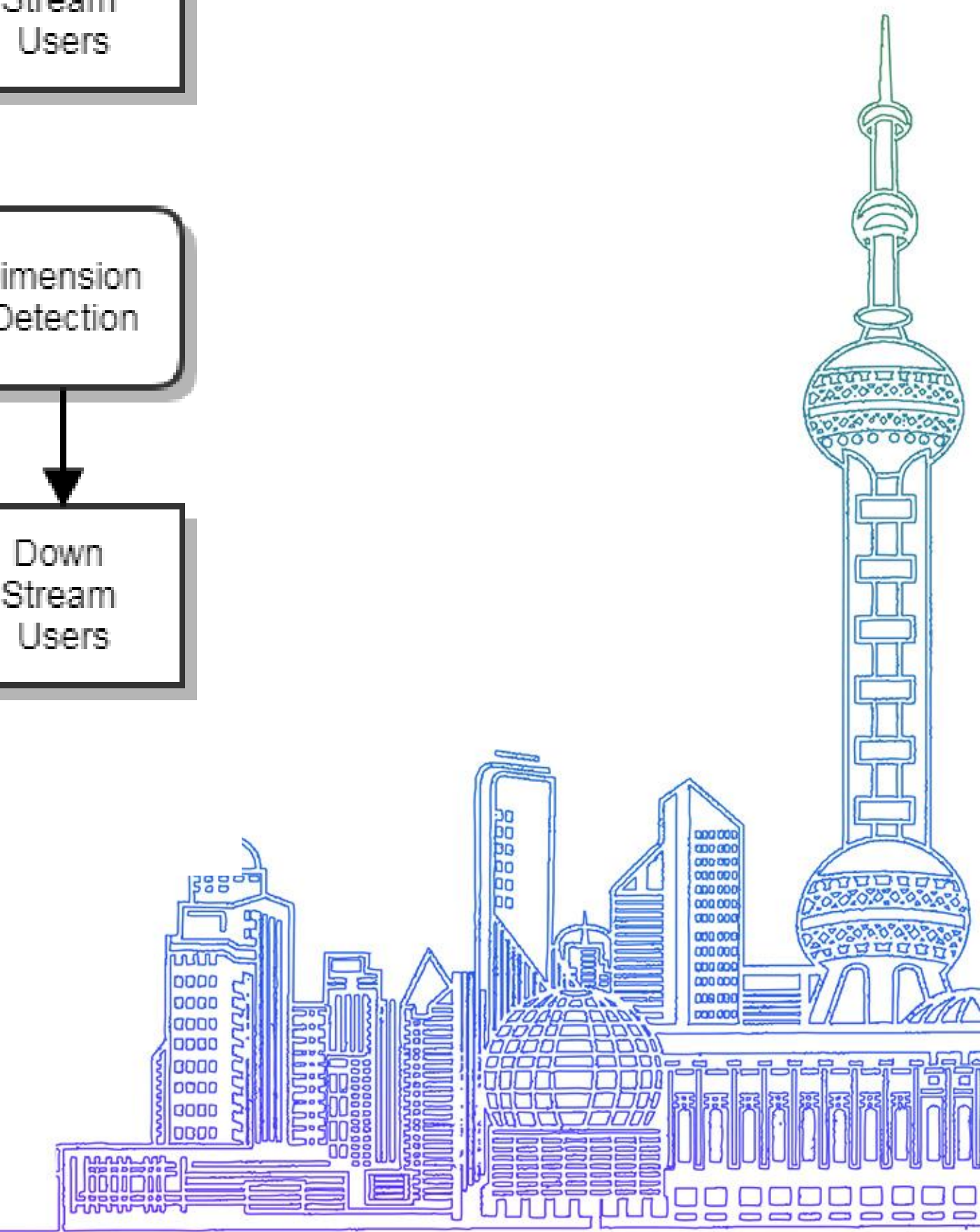
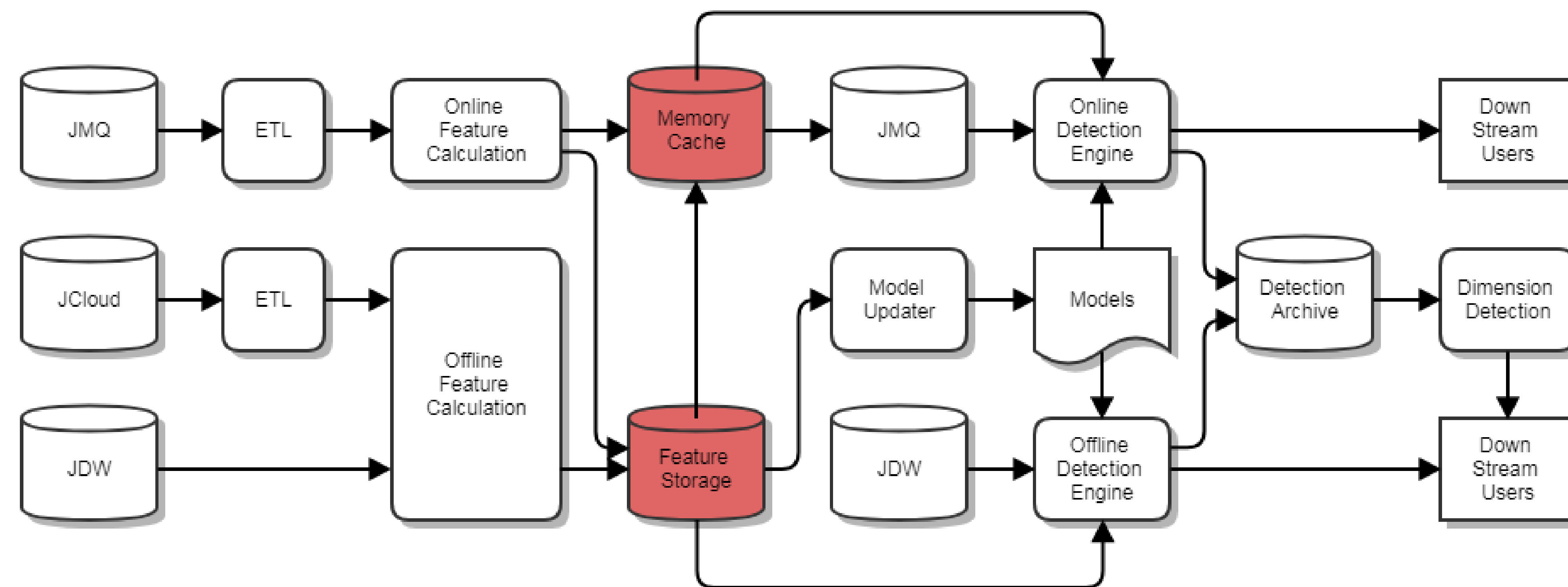
- 离线特征
 - 初级特征：特征工厂
 - Nebula（特征计算语言）
 - 高阶特征
 - 图模型：Spark GraphX
 - 无监督聚类、分类：Spark ML
 - 序列分析等：自实现
- 在线特征
 - 时间窗口：Spark Streaming



京东反刷单系统架构

特征管理

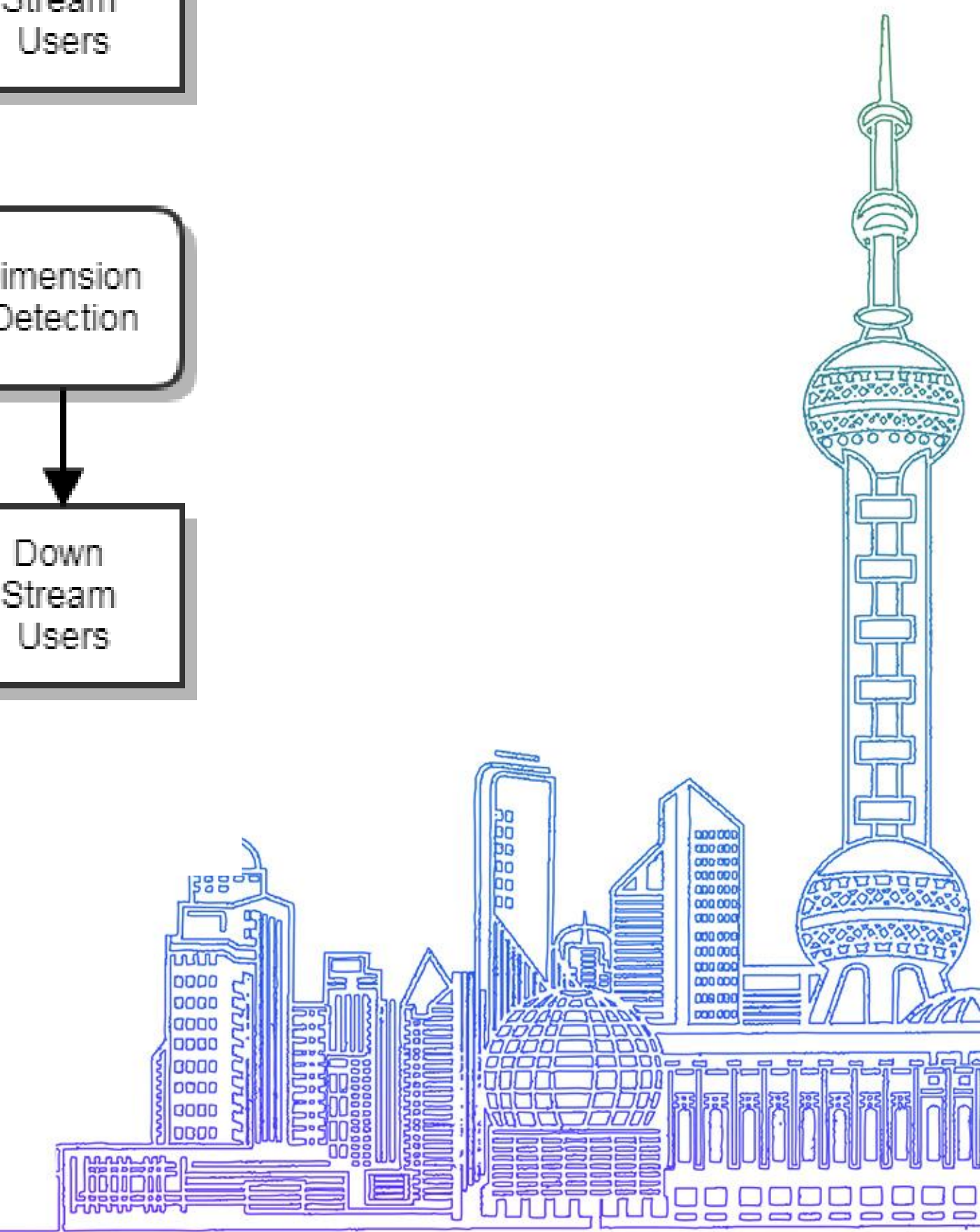
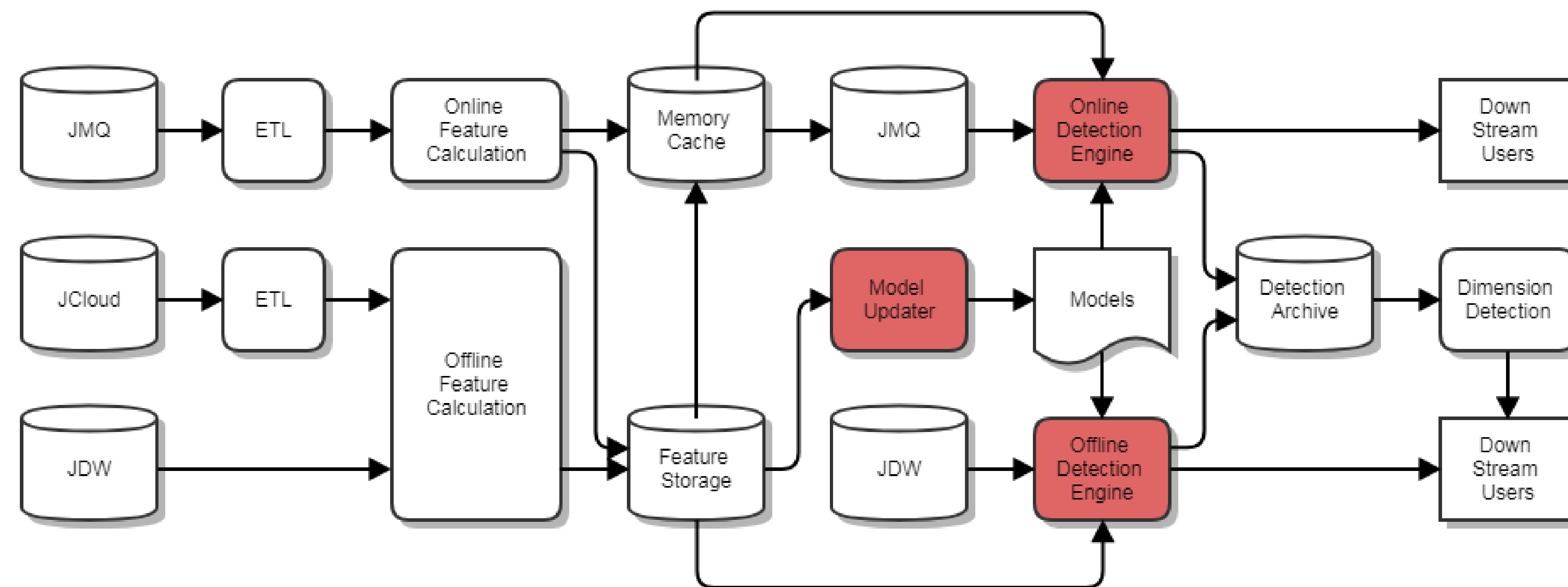
- 离线特征存储：特征仓库（HDFS + Hive）
 - 服务离线识别
 - 服务模型训练和更新
 - 提供特征共享
- 在线特征索引：Redis



京东反刷单系统架构

模型与决策引擎系统

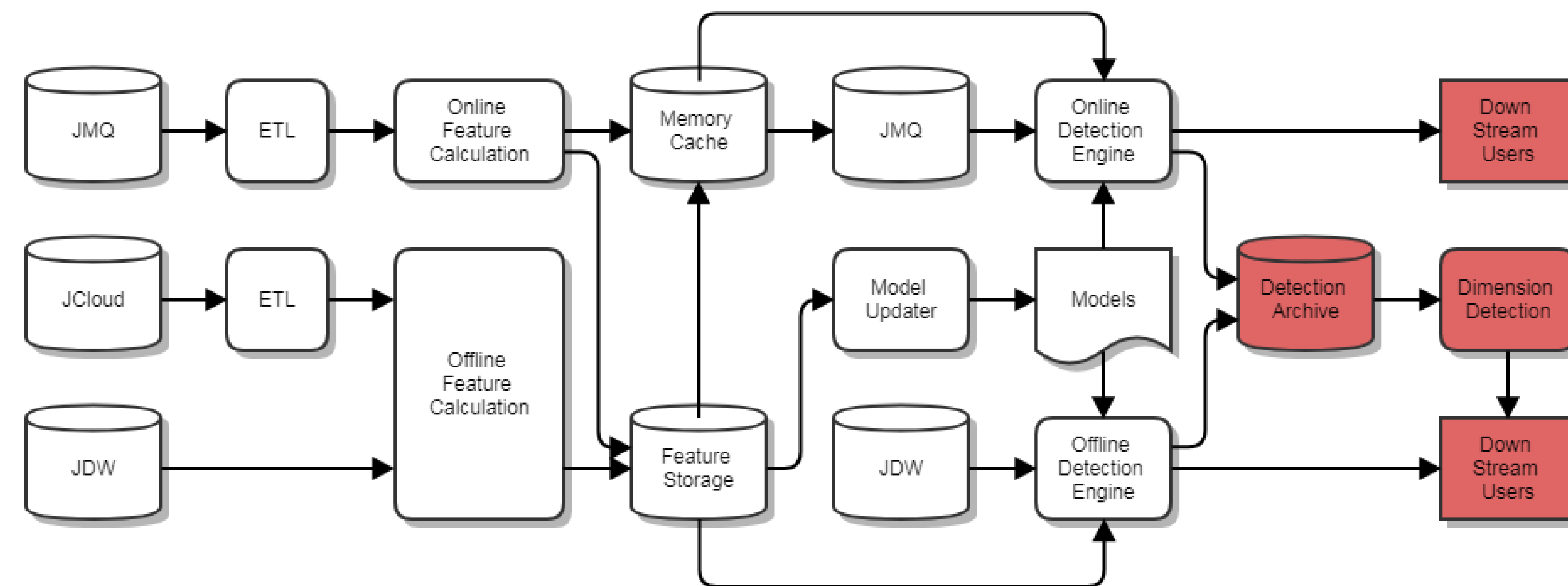
- 算法模型引擎：Rigel (Spark)
 - 机器学习算法：Spark ML
 - 深度学习方法：TensorFlow on Spark
 - Rigel (模型配置语言)
- 业务规则引擎：Saiph (Spark)
 - 规则：JBoss Drools



京东反刷单系统架构

结果归档与下游服务

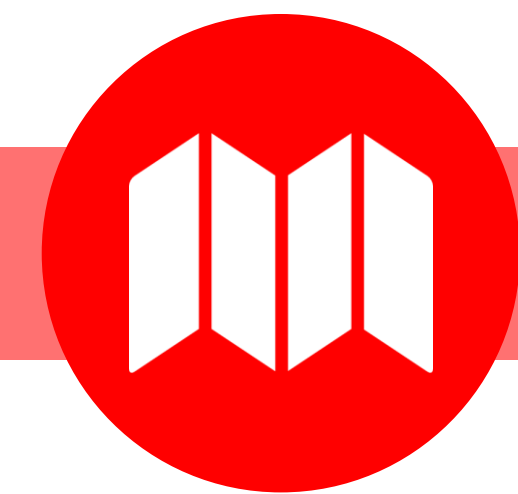
- 历史归档: HDFS
- 结果推送
 - 实时请求: JSF (RPC)
 - 消息推送: JMQ (Kafka)
 - 离线接口: JDW (HDFS + Hive)
- 产品化服务
 - OLAP: Presto



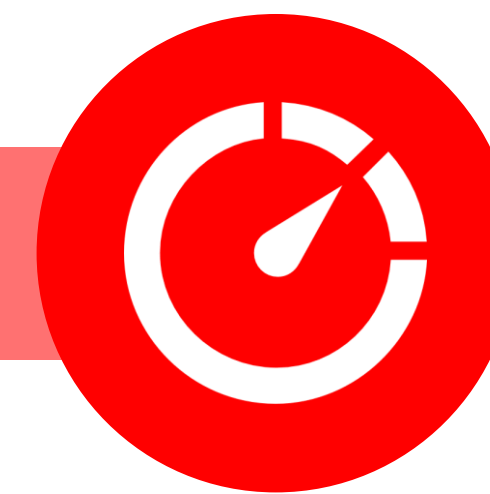
架构实践



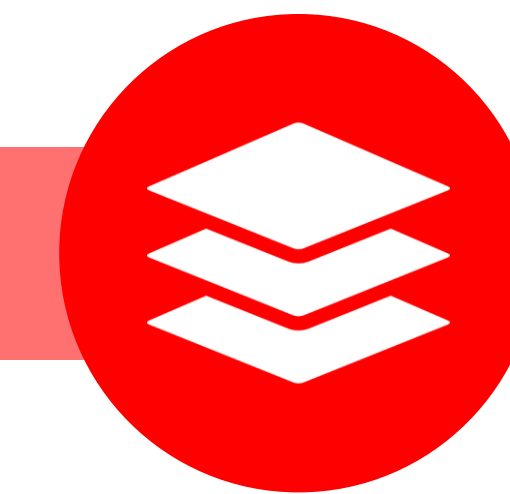
Fundamental



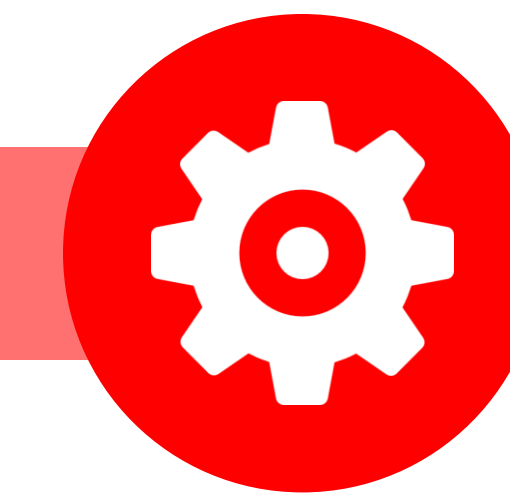
Adaptability



Reproducibility



Flexibility



Customizability



高可用

- 系统模块化，各模块分布式、配置化、主从灾备、监控



低延迟

- 定义数据优先级、低优先级数据降级处理
- 数据变化期内多次识别同一订单，尽早发现刷单



架构实践



Fundamental



Adaptability



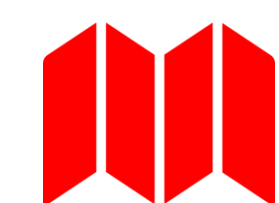
Reproducibility



Flexibility



Customizability



多种业务
类型

- 特征工厂，抽象业务数据依赖、特征计算过程、计算作业调度，使特征计算配置化，减少开发量



不同数据
形式

- ETL将所有数据统一到集群数据（HDFS）和基于消息队列（Kafka）的流数据



初级特征

对象的时间跨度以及筛选条件，如限于过往半年的订单记录，限于移动端日志

按空间筛选后，聚合的字段，如账户名、商品标识符



按维度聚合后，群组上的统计方式，如计数、均值、方差、信息增益

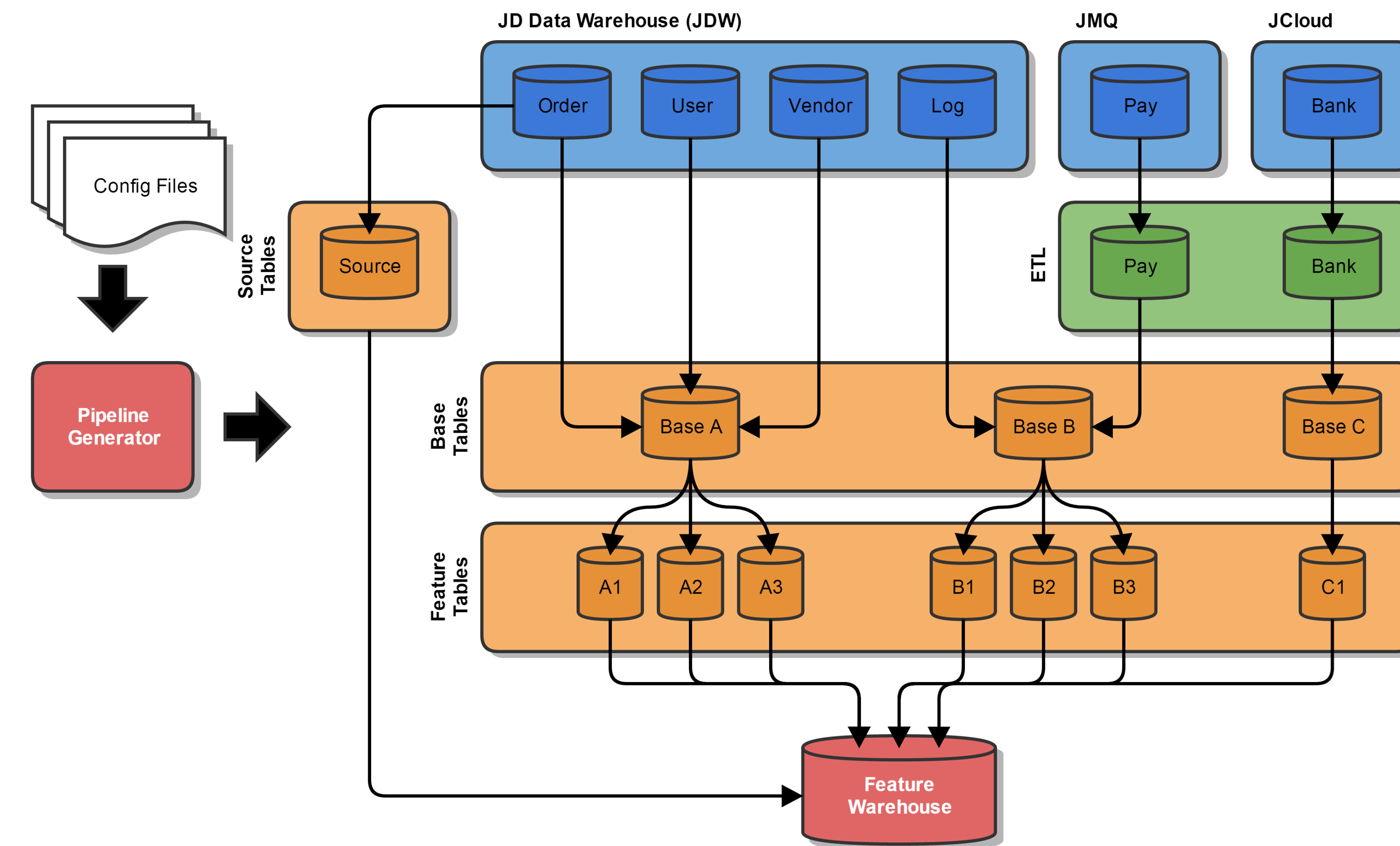
最后将测度按照维度关联到订单记录上时，目标的范围，如仅适用于当天的订单记录，仅适用于自营业务的订单记录



特征工厂：通用特征计算框架

初级特征计算框架Nebula

- 提供上述要素的配置语言表达形式，包括数据源定义，数据源加工，聚合字段定义，测度计算，关联识别订单操作等等
- 由配置语言自动构建特征计算脚本（Pig）和作业调度任务（Oozie），并根据作业历史执行性能优化调度

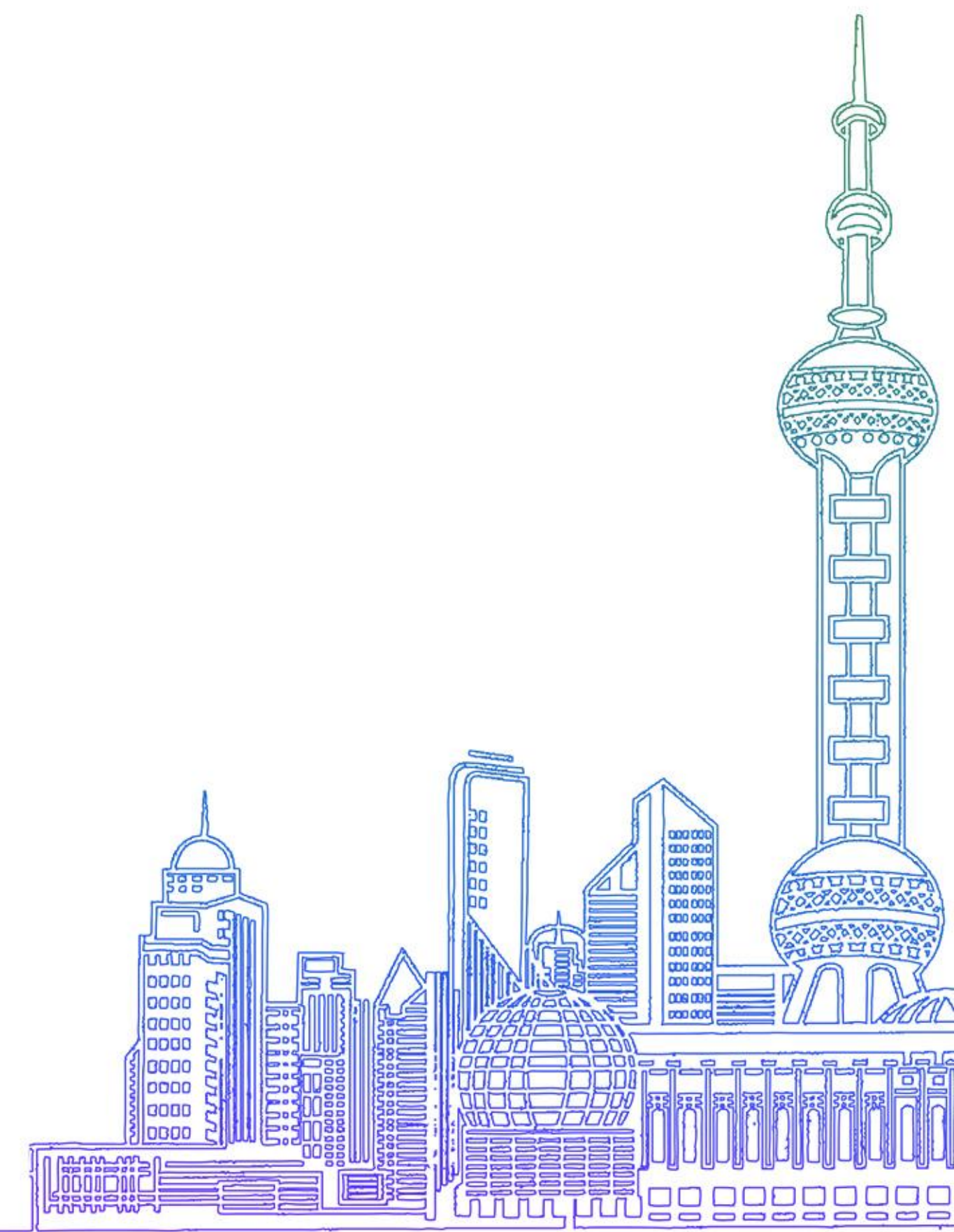


架构实践



保留现场历史

- 结果涵盖所有特征数据、模型识别结果、模型版本和配置等，结合结果归档可由系统中任意模块开始重现、回放数据计算和算法检测过程
- 可在重现过程中加入新的数据、特征、模型等对比更改前后结果



架构实践



Fundamental



Adaptability



Reproducibility



Flexibility



Customizability



可扩展

- 模型配置化涵盖Spark ML、Spark GraphX、TensorFlow、Drools等主流机器学习、深度学习、规则引擎框架，为算法实现提供更多选择

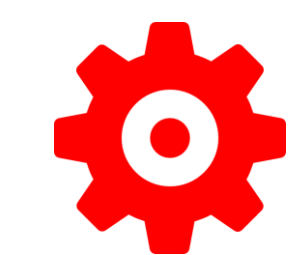


解耦

- 系统在模块层级分离模型算法（机器学习、深度学习）和业务规则，防止模型结果和规则之间的交叉影响以及由此带来模型训练的困扰

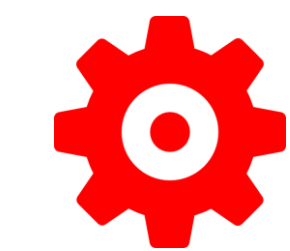


架构实践



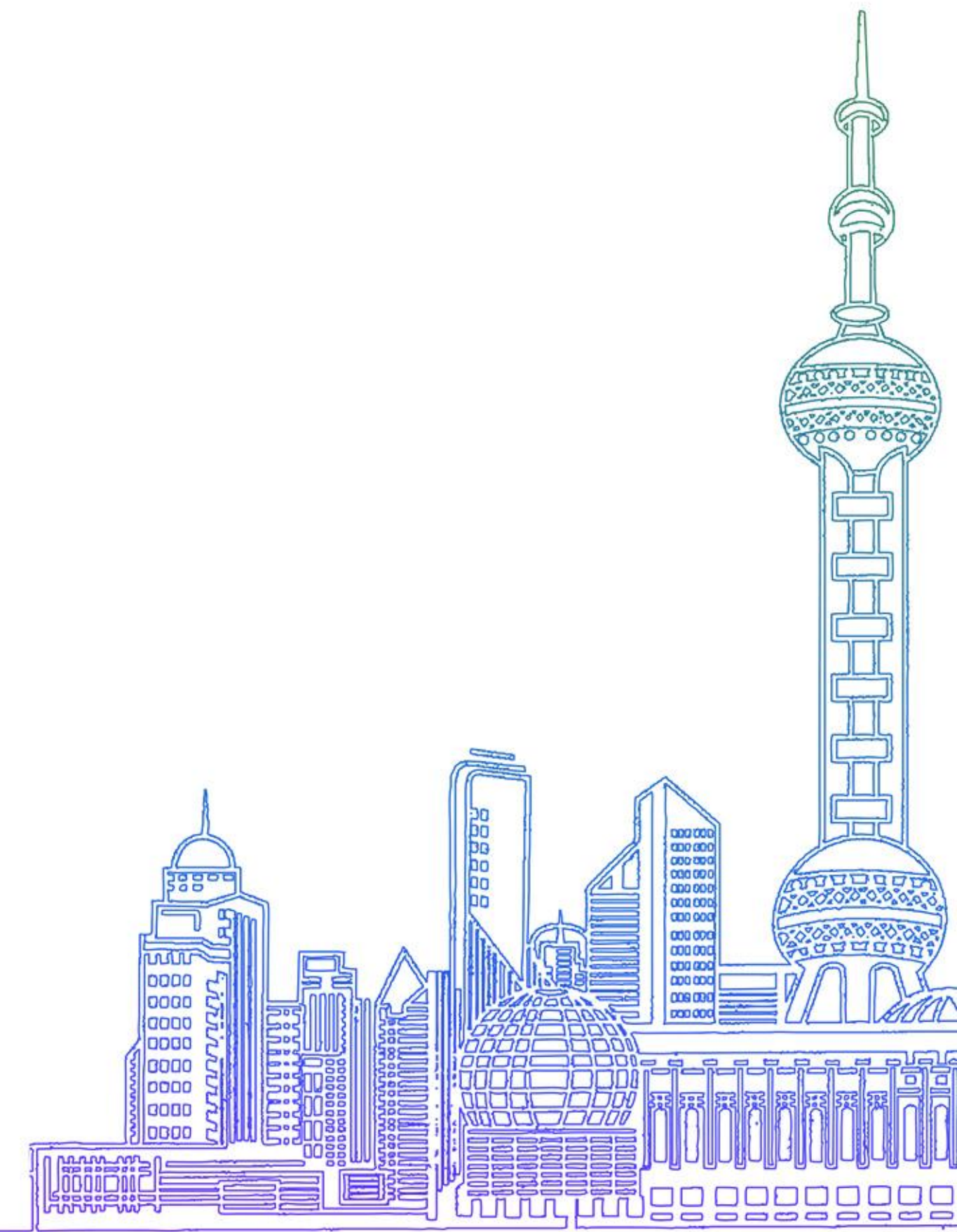
定制化

- 反刷单数据产品化平台，下游业务利用反刷单结果定制自己的风控指标和规则



数据保密性

- 基于Presto实现的OLAP服务在配置端增加自研聚合维度探查模块，防止明细数据暴露给下游使用方



谢谢观看！

2018 唯品会第三届互联网电商安全峰会

2018 vip.com third Internet ecommerce Security Summit

2018-5-5 上海

