

大数据背景下电网数据质量 研究与实践

全球能源互联网研究院

2017年8月



国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

目 录

A 3D rendering of a globe showing the Americas, with three concentric light blue arcs to its right, resembling a signal or data transmission.

全球能源互联网与大数据

大数据质量提升的系统性方案

大数据清洗治理系统



人类可持续发展面临巨大的能源挑战

随着社会经济的发展，人们对于能源的需求越来越大，**能源**是现代工业社会的命脉，能源问题关乎生态环境和可持续发展。化石能源的大规模开发使用导致我们被迫面临**三大严峻挑战**，给人类生存发展带来严重威胁。

1

资源紧张

- 煤炭：只能开采**110**年
- 石油：只能开采**57**年
- 天然气：只能开采**54**年

2

环境污染

化石能源的大量开发，在生产、运输、存储、使用的各环节，对大气、水质、土壤、地貌等造成严重污染和破坏。

3

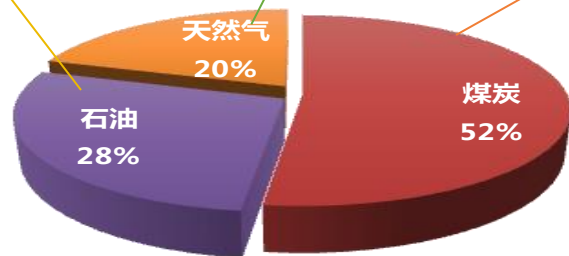
气候变暖

化石能源的碳排放是气候变暖的主因。自1850年以来，全球地表平均温升已经超过1°C。

2398亿
吨

187万亿
米³

8915亿
吨



全球化石能源剩余探明可采储量

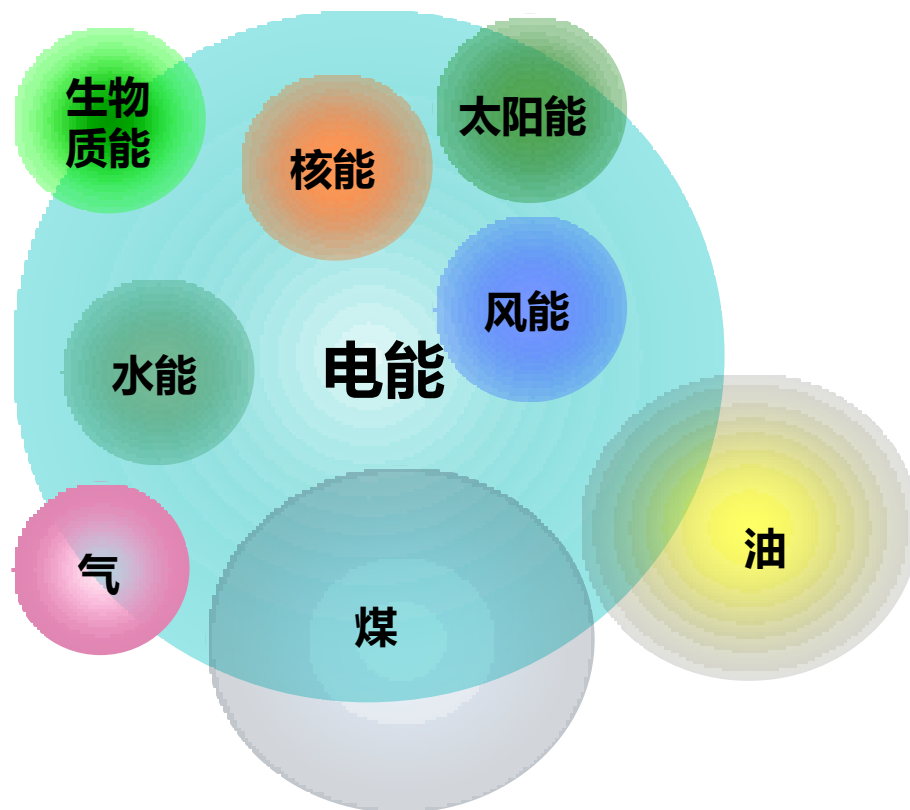


开发清洁能源体系满足人类可持续发展

清洁能源 - 电能 - 化石能源

- 能源供应侧实施**清洁替代**：
清洁能源替代化石能源
- 能源消费侧实施**电能替代**：
电能替代煤、油、气。

形成电为中心的能源格局



国际能源署（IEA）《2014能源技术展望》（ETP2014）：
推进电气化是全球能源系统的驱动力之一，可从根本上转变能源供应及终端用能

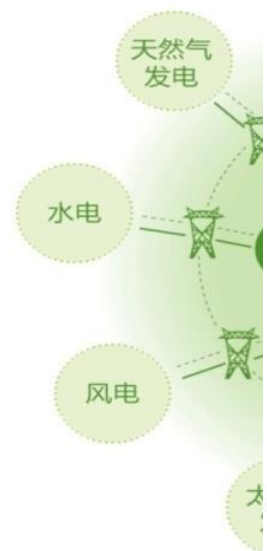


国家电网
STATE GRID

全球能源互
GLOBAL ENERGY INTER

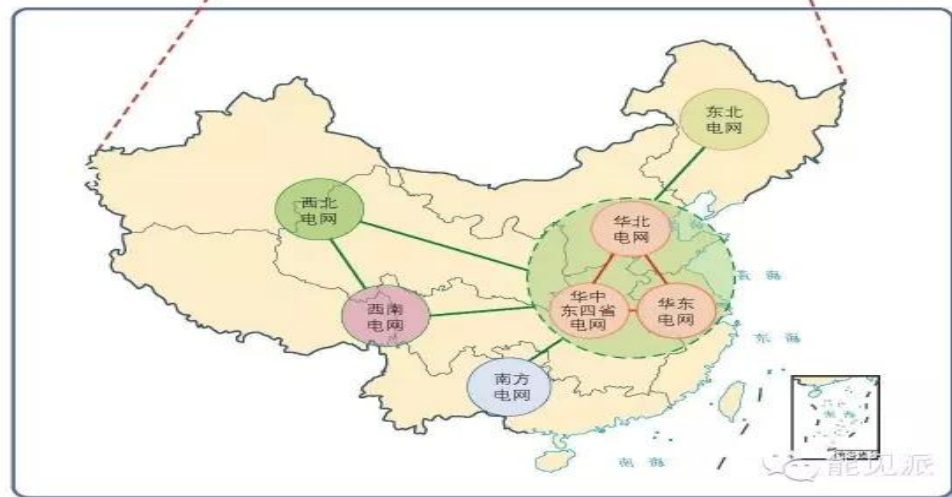
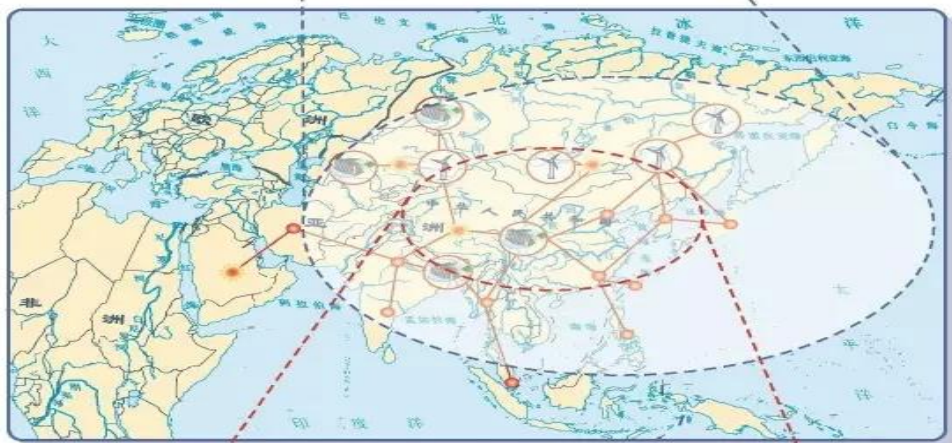
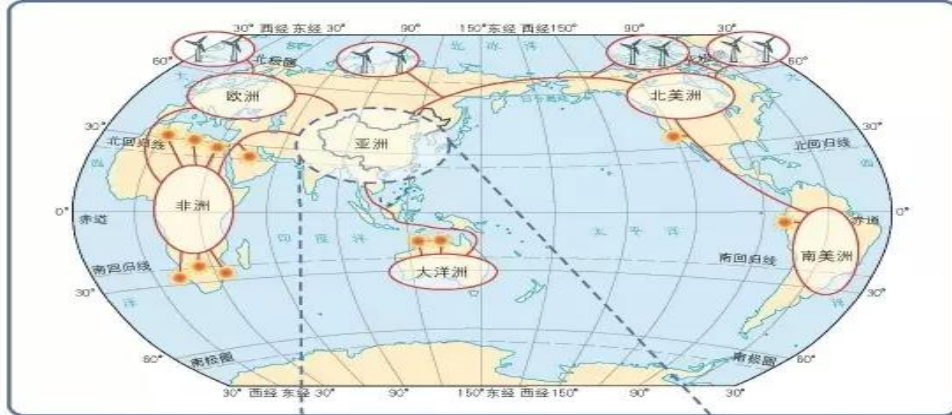
洲际互联

构建以
力需求，是1
电网必



国内互联

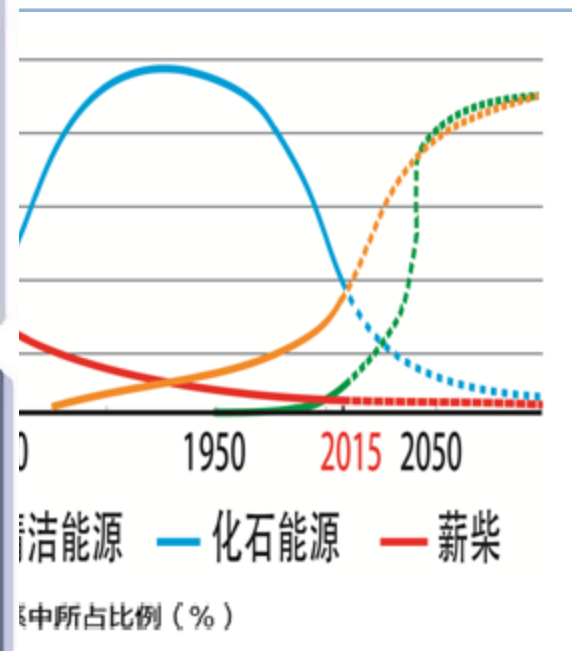
电网配置



持续发展的必由之路

和绿色方式满足全球电

网。



低碳发展



全球能源互联网实质及关键技术

全球能源互联网：

以**特高压电网**为骨干网架，以输送**清洁能源**为主导，全球互联泛在的**智能电网**。
是清洁能源在**全球范围**大规模开发、配置、利用的**基础平台**。



关键技术

- 电源技术：提高清洁能源开发效率和经济性
 - 风力发电、太阳能发电、海洋能发电及分布式电源
- 电网技术：提高输电距离和容量，能源资源全球配置
 - 特高压、海底电缆、超导输电、直流电网、微电网
- 储能技术：保障清洁能源大规模发展和电网安全运行
- 信息通信技术：实现电网智能化、互动化和运行控制

特征：网架坚强、**广泛互联**、高度智能、开放互动



国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

智能电网广泛采用先进信息技术：

- 自动预判、识别各类电网故障和风险;
- 适应各类集中式、分布式清洁能源大规模接入和大范围配置要求；
- 满足用户多样化、智能化用电需求，构建公共服务平台，促进智能家居、智能社区、智能交通、智慧城市发展。
- 2016年2月29日，国家发改委、能源局以及工信部联合下发《关于推进“互联网+”智慧能源发展的指导意见》（下简称《指导意见》），以促进能源和信息深度融合，推动能源互联网新技术、新模式和新业态发展，推动能源领域供给侧结构性改革和能源革命。

先进信息技术 - 让电网更加智能



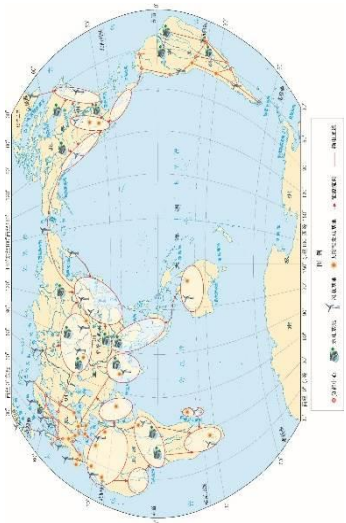


国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

信息技术创新应用 - 大数据



全球能源互联网

大数据是全球能源互联网的必然产物

大数据产品是全球能源互联网的重要服务形态

大数据是支撑全球互联网运行控制的必要技术

大数据

知识发现与价值挖掘

跨专业、跨领域整合分析

海量数据处理能力





国家电网
STATE GRID

全球能源互联网研究院

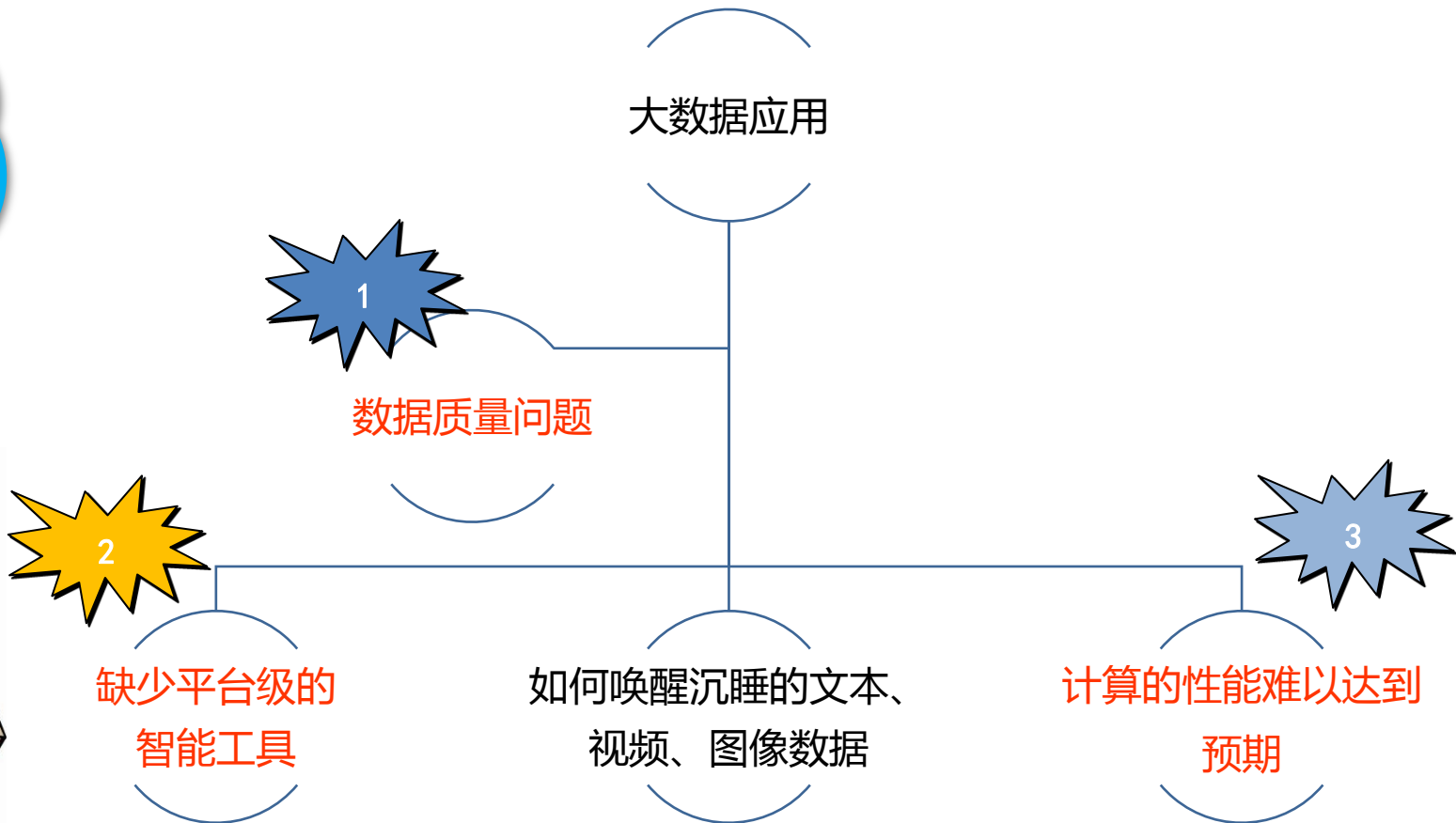
GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

大数据技术应用困局

大数据技术应用
远比想象的
困难！



toopen.com.cn 插画师





国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

目 录



全球能源互联网与大数据

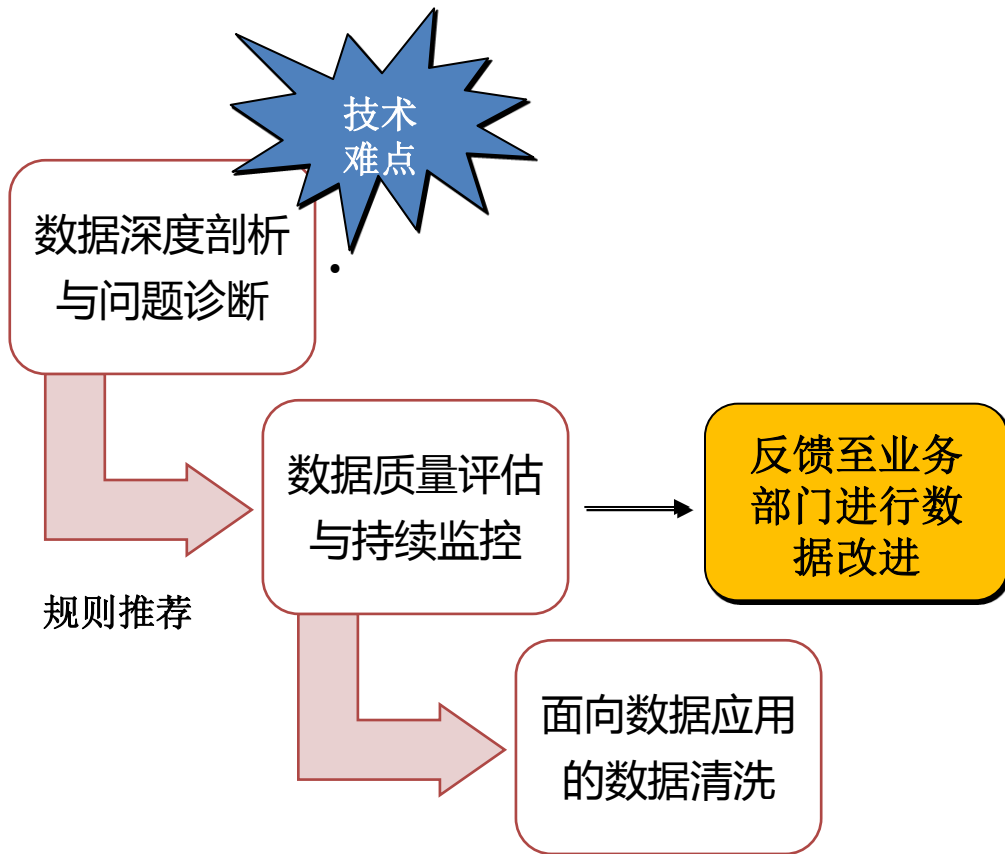
大数据质量提升的系统性方案

大数据清洗治理系统



大数据质量提升的系统性方案

数据质量 (Data Quality) 是数据分析结论有效性和准确性的基础，也是最重要的前提和保障。“数据质量”尚无统一定义，但学术界普遍认可以下评价维度：



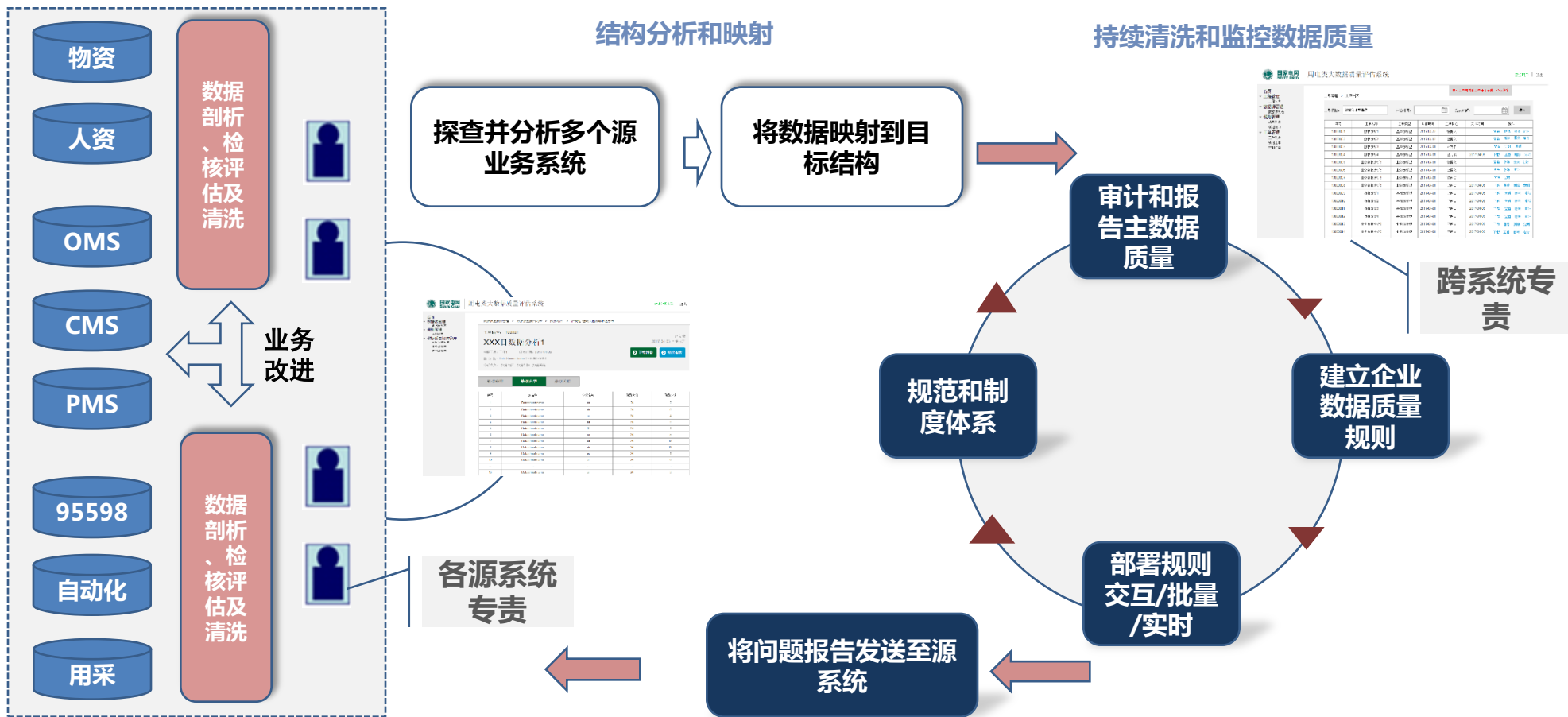
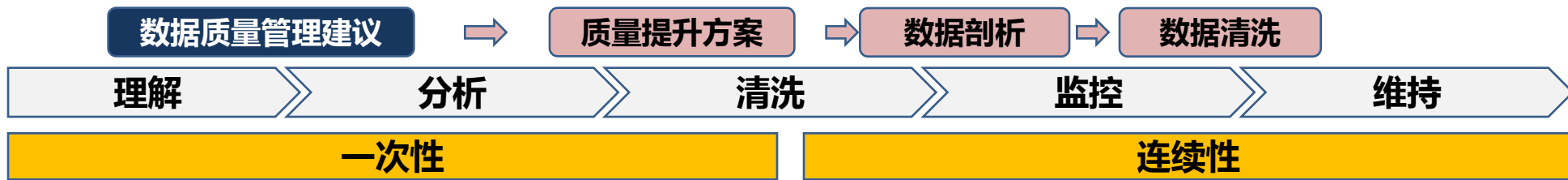


国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

大数据质量提升的系统性方案





国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

大数据质量提升的系统性方案

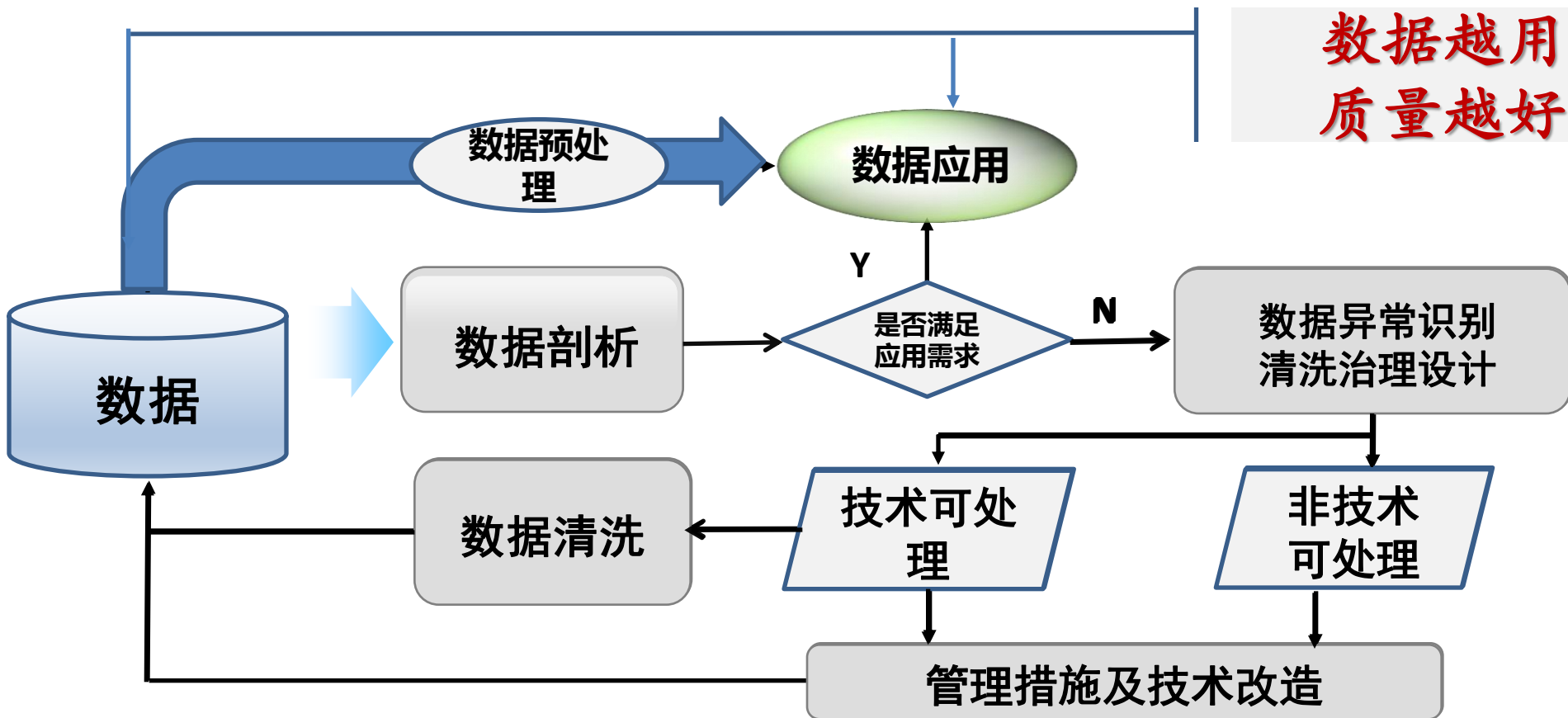
数据质量管理建议

质量提升方案

数据剖析

数据清洗

数据越用
质量越好





国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

大数据质量提升的系统性方案

数据质量管理建议

质量提升方案

数据剖析

数据清洗

主数据模型设计优化

业务逻辑及架构顶层设计优化

业务数据库结构优化

促进

数据质量问题



数据应用及共享

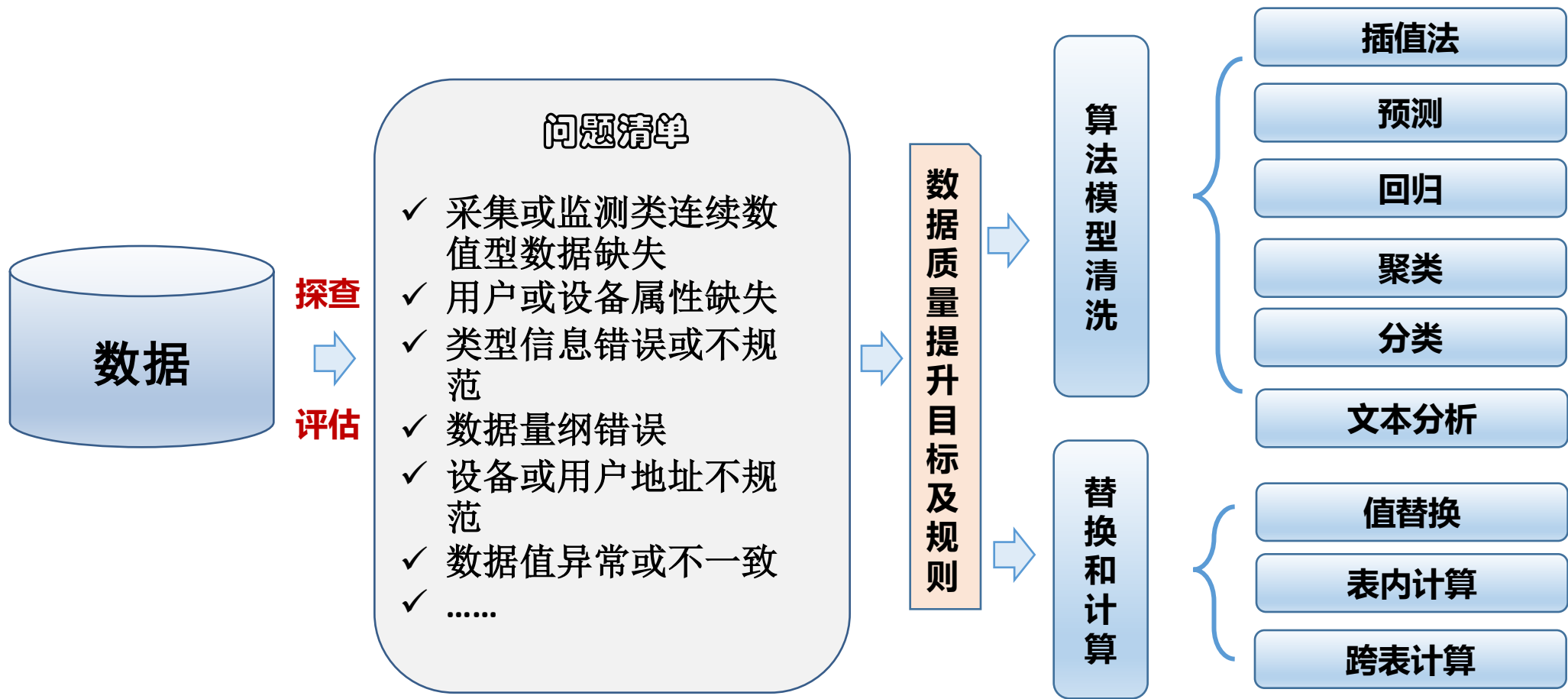
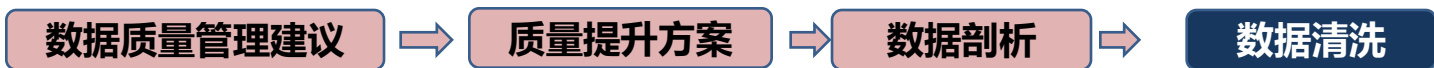


大数据质量提升的系统性方案





大数据质量提升的系统性方案





国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

目 录

A 3D rendering of a globe showing the Americas, with three concentric light blue arcs to its right, resembling a signal or energy wave.

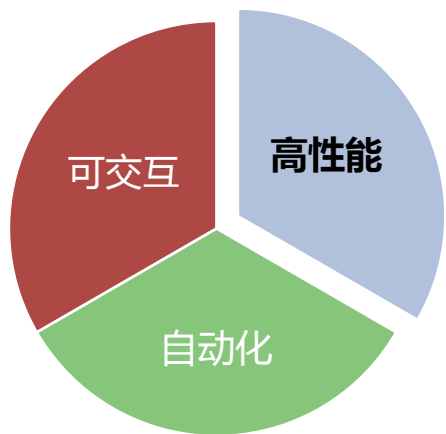
全球能源互联网与大数据

大数据质量提升的系统性方案

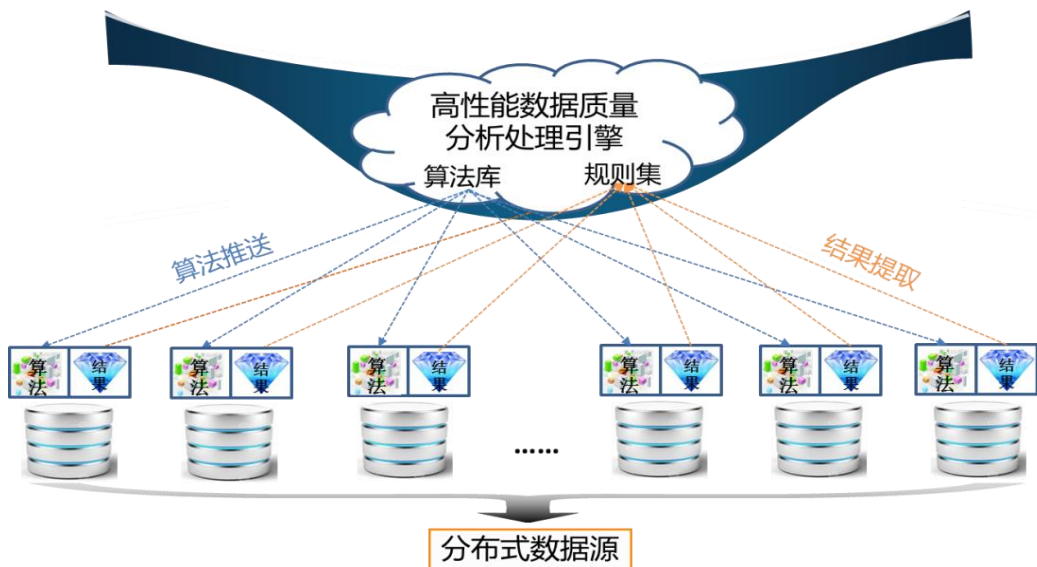
大数据清洗治理系统



大数据清洗治理系统-功能特性



- 多源异构数据源的集成
- 灵活配置规则，内嵌基于业务逻辑的清洗模型并支持扩展
- 基于大数据的批量数据清洗，可持续调度执行
- 灵活部署，高效执行



数据源管理

新建数据源

选择数据源

数据库名: servername

用户名: username

密码: password

主机: only specify if not default

数据库: database

测试连接

任务调度管理

任务调度列表

任务ID	任务名称	常用依赖ID	设置依赖任务ID	定时时间	操作
4	清洗任务测试	rele_id	<div>rele_id</div> <div>4</div> <div>3</div> <div>2</div> <div>1</div>	03:00	<div>未开始</div> <div>定时执行</div> <div>手动执行</div>



大数据清洗治理系统-核心功能

解决的问题

- 1、替换标准内容
- 2、主表与码表关联获取码值内容
- 3、统一字段格式及内容
- 4、空值字段赋值处理
- 5、特殊字符替换
- 6、多列间逻辑运算或拼接
- 7、删除重复行记录

基于规则的清洗治理

规则

表达式

解决的问题

- 1、提取有规律的信息（如提取数字）
- 2、逻辑运算或简单统计信息
- 3、数据精度标准化
- 4、日期标准化、结构化（如分解年月日及季度等）
- 5、字符串大小写，去空格、特殊字符
- 6、正则替换（如替换数字开头的内容）
- 7、特殊格式数据解析（如URL、json格式数据的解析）
- 8、补足位数，简单加密字段内容

解决的问题

- 1、地址标准化
- 2、地址补全、地址去重
- 3、错误行政区划所属关系纠正、行政区划错字纠正
- 4、门牌楼号正则化
- 5、基于地址丰富数据维度（增加市、区、镇）

基于模型的清洗治理

文本型算法模型

数值型算法模型

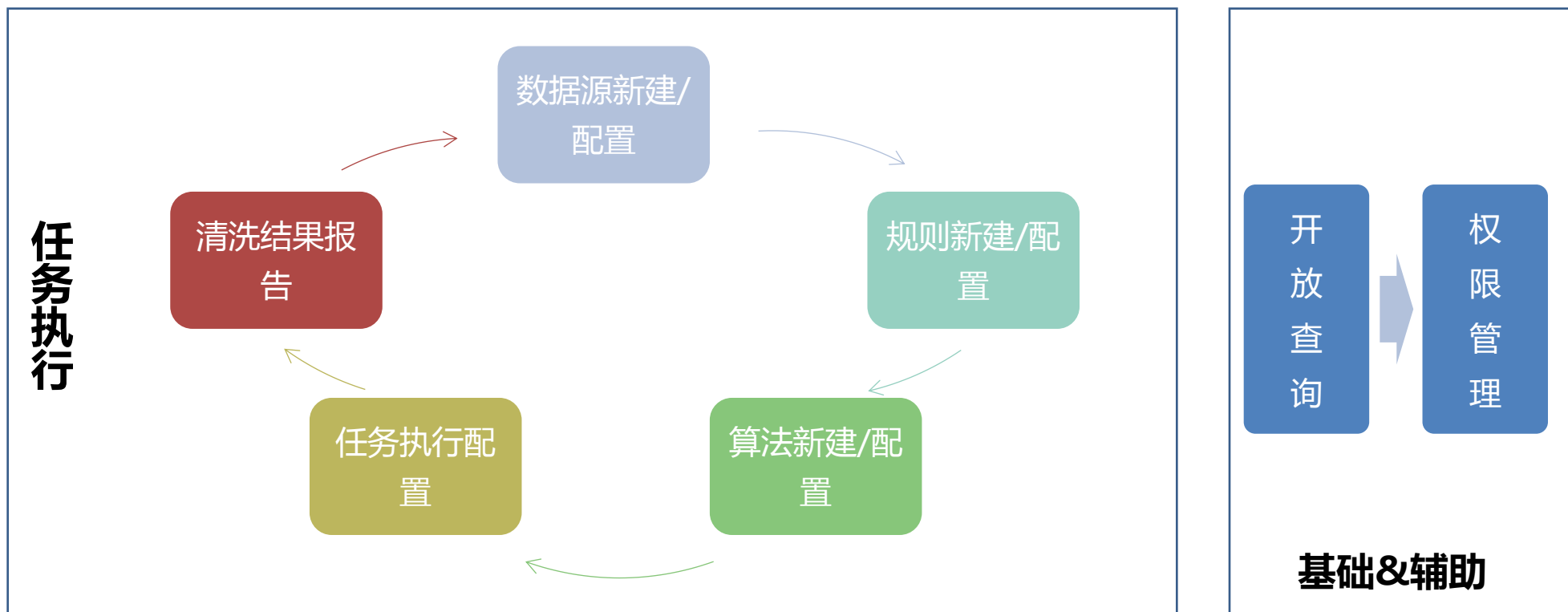
解决的问题

- 1、最优匹配搜索CBMS：适用于数据矩阵的缺失值填补,对连续缺失有优势，能够清洗平均连续丢失10个数值型数据
- 2、改进的交替式最小二乘优化IALS：适用于时间序列类数据的缺失值的预测填补
- 3、时间序列ARIMA模型与卡尔曼滤波KARM：适用于时间序列数据的缺失值填补,对分散性缺失比较有优势，能够填补分散性缺失50%的数值型数据
- 4、集合模型：将多个独立缺失值还原算法的结果相结合，返回更准确的估算结果



大数据清洗治理系统-应用流程

通过简单的配置操作即可完成专业而持续的数据清洗治理，使经过简单训练的人员轻松完成繁杂的数据处理。





国家电网
STATE GRID

全球能源互联网研究院

GLOBAL ENERGY INTERCONNECTION RESEARCH INSTITUTE

谢谢！