

中国非结构化数据中台 实践白皮书

目录

Content

01

开启数据智能时代

02

非结构化数据中台建设与挑战

03

基于非结构化数据中台的应用场景

04

展望行业趋势

开启数据智能时代

中国非结构化数据中台
实践白皮书

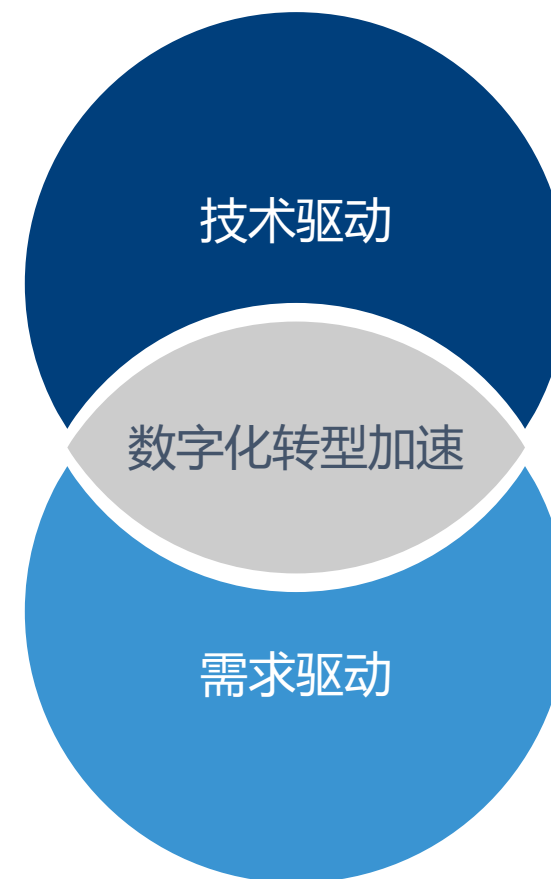
01

需求驱动

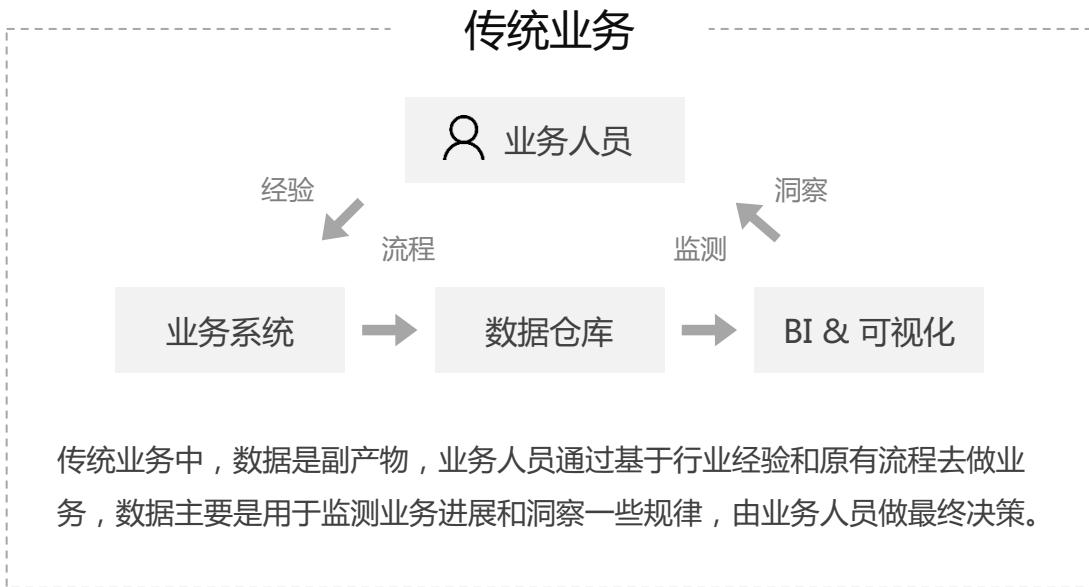
- 全球经济增速下滑，叠加疫情影响，企业面临的外部环境严峻。同时，用户需求多元化，企业战略重心由以产品为中心转变为以用户为中心。多方因素使得企业面临商业模式的重塑，企业经营由过去粗放式的流量扩张向精细化运营转变，需要借助数字化实现加速转型，实现降本增效，提升企业竞争力。

技术驱动

- 数字化转型的核心是数据。近年，互联网&移动互联网的发展产生大量数据。同时，云计算、人工智能、5G、物联网技术的发展，推动企业数据治理能力提升，使得数据驱动业务增长成为可能。

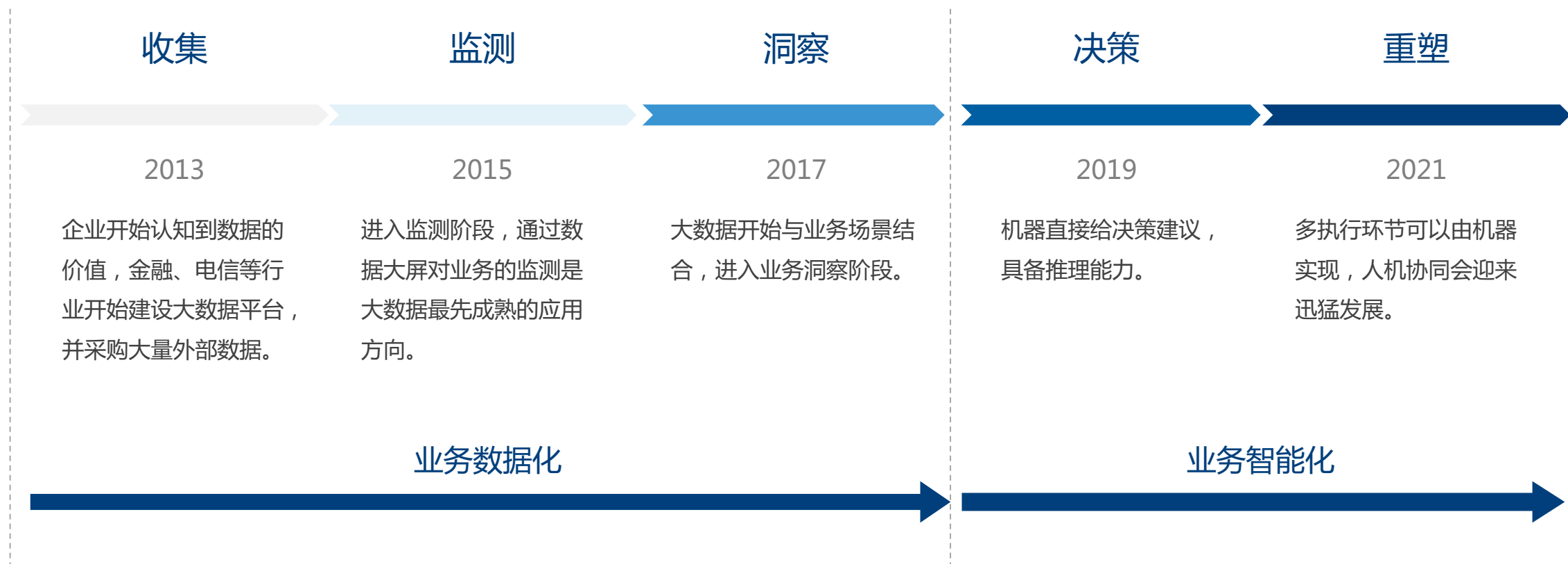


在数字化转型的过程中，数据成为新的生产资料，挖掘数据的价值，提高企业管理和决策水平已成共识，整个行业进入数据智能时代。



对比项	流程驱动	数据驱动
数据的价值	业务系统的副产物	业务系统的核心
决策方式	人工	人机配合
迭代速度	极慢	快
商业价值	低	高

从业务应用的角度，数据智能的发展经历了收集、监测、洞察、决策等四个发展阶段，数据的应用价值不断提升，逐步从业务数据化转向业务智能化。未来，数据智能将会进入人机协同的业务重塑阶段。





- 数据智能分为中台和应用场景两个核心领域。
- 中台是数据智能的核心，主要分为技术中台、数据中台和业务中台：
 - 技术中台主要由各类分析工具组成，帮助企业解决技术问题的公司，如用户行为分析、数据科学平台、BI与可视化、日志分析、NLP/知识图谱等；
 - 数据中台主要是帮助企业做数据资产化，建立数据中台的公司包含各类数据服务公司 and 数据治理公司；
 - 业务中台是基于技术和数据，结合行业应用场景，形成的模型、产品。
- 其中，数据中台是中台体系最重要的部分。



- **数据中台**汇聚企业的业务数据，包括企业经营数据、客户行为数据、设备运转数据、生态合作数据等，并赋能给各类不同的数据应用场景。
- **数据中台的价值是挖掘数据价值并服务业务场景**
- 数据中台通过自动化、智能化的数据采集与汇聚，将实时与离线数据打通关联，对数据开发深度挖掘数据价值，并开放数据服务至各业务场景中。具备汇聚整合、数据提纯加工、数据服务可视化、数据价值变现等核心能力。

数据驱动决策的前提是数据整合

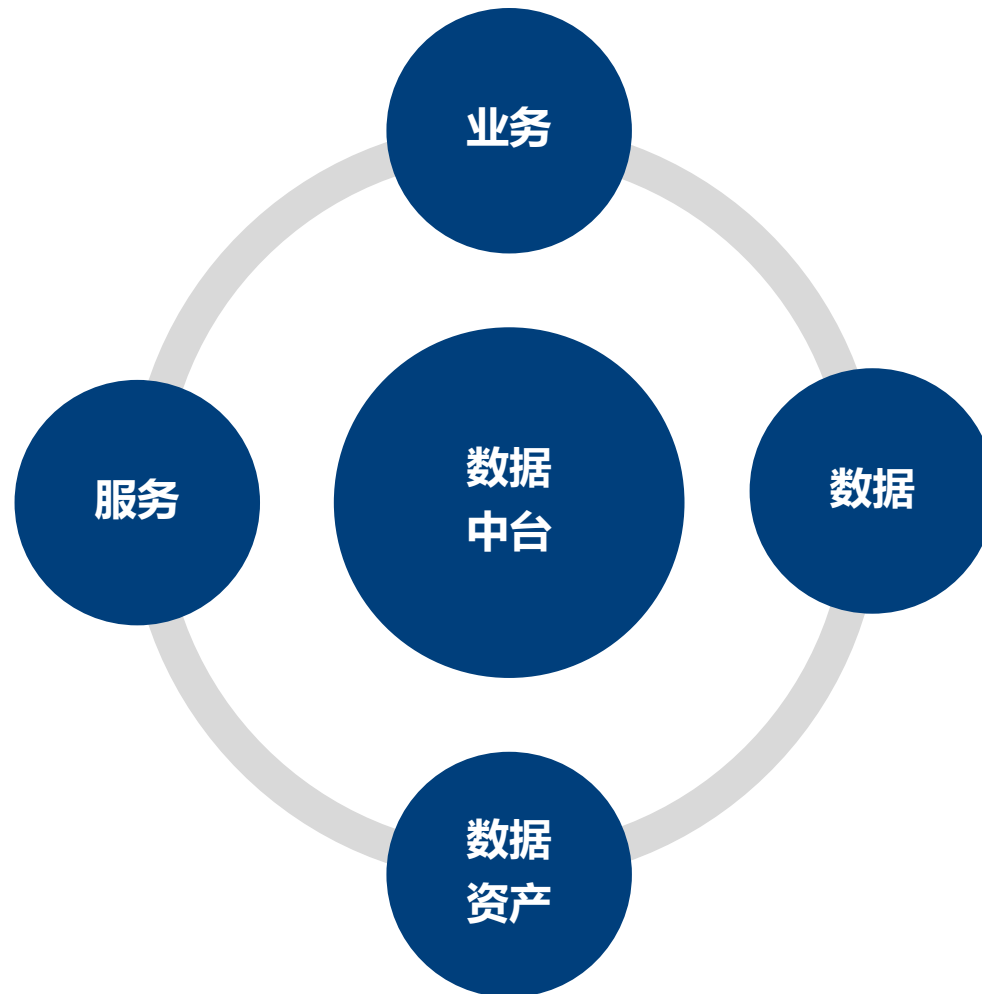
- 数据智能的标志是数据驱动决策，让机器具备推理等认知能力，大数据能够指导决策。同时完成了业务数据化进程，开始进入到业务智能化，依靠数据改变业务
- 决策需要机器具备推理能力，建立复杂关系网络，从训练模型的角度，这意味着必须要有更大规模的数据。同时，决策意味着解决的业务问题复杂性远远大于之前，因此，需要汇聚更多种类的数据。
- 信息化时代数据散落在各个系统中，数据存在脏乱差、ID不统一等问题，数据孤岛现象严重。基础设施的云化使得基础IT资源实现了统一管理和调度，数据的统一管理和调度就提上日程，成为下一个需要解决的问题。

数据驱动业务需要数据中台

- 在企业数字化转型进程中，传统企业需要具备互联网公司快速迭代升级的能力，基于数据驱动业务发展，这需要建立一站式技术能力、统一的数据管理、快速配置开发业务的能力。
- 以阿里巴巴为代表的中台模式给传统企业提供了一条道路，各类中台会在企业内部逐步形成。因此，形成数据中台是大势所趋。

数据中台需要与业务结合， 才能真正地让数据用起来

业务产生数据，数据中台帮助企业聚合内外部数据，将原始数据转化为数据资产，快速构建高效的数据服务，使企业可以持续、充分地利用数据，以数据洞察来驱动业务决策和运营，最终提升企业决策水平和业务表现，赋能企业解决业务问题。业务产生数据，数据形成数据资产，数据资产提供数据服务，进而赋能业务，形成闭环。



在快速增长的数据中，非结构化信息占比已达80%。据Gartner估计，从2019年到2024年，非结构化数据容量预计将增加两倍。但企业现有架构通常无法应对海量非结构化数据的管理与应用。





- 非结构化数据管理之所以难，不仅因为其数量多、分散性高，还在于用户对于非结构化数据的需求是多层次的。在数据、内容、信息和知识层面分别有不同的需求。

非结构化数据管理需要革新的底层架构

- 数据整合形成数据中台，意味着大量数据治理，国内企业信息化、数据化程度不高，存在着大量文本、图像等非结构化数据。
- 非结构化数据管理需要将底层数据打通，从源头保障数据资产的复用能力，实现数字资产统一运营、全面合规、高效利用。
- 仅仅依靠数据分析技术难以解决问题，必须将计算机视觉、NLP、知识图谱等技术融入其中，借助深度学习等人工智能技术实现数据治理，进而实现知识复用与智能搜索。
- 因此，非结构化数据管理需要革新的底层数据架构，非结构化数据中台能够满足需求。

非结构化数据中台架构



- 非结构化数据中台对对象数据、元数据、索引数据进行汇集、管理，融合人工智能技术，提供先进的数据架构底座，进而通过非结构化数据赋能各行各业应用。
- 非结构化数据中台基于内容总线、内容数据湖等数据架构，能够实现智能搜索、内容窃案洞察、内容自动化等功能，应用于企业的多业务场景，包括企业内容立体安全、业务流程自动化、数据资产管理、智能知识运营等。

非结构化数据中台建设与挑战

02

中国非结构化数据中台
实践白皮书



保障：企业业务合规性

法律法规遵从
行业监管合规
组织合规内控体系建设



实现：企业业务赋能

从业务执行中抽取具有商业智能价值的信息，实现非结构化数据商业智能功能
打造企业内部业务敏捷流程



提升：企业创造力

构建知识管理体系
转型学习型组织，构建企业文化

数字资产管理

实现数字资产的全生命周期管理



数据合规体系



数据资产管理

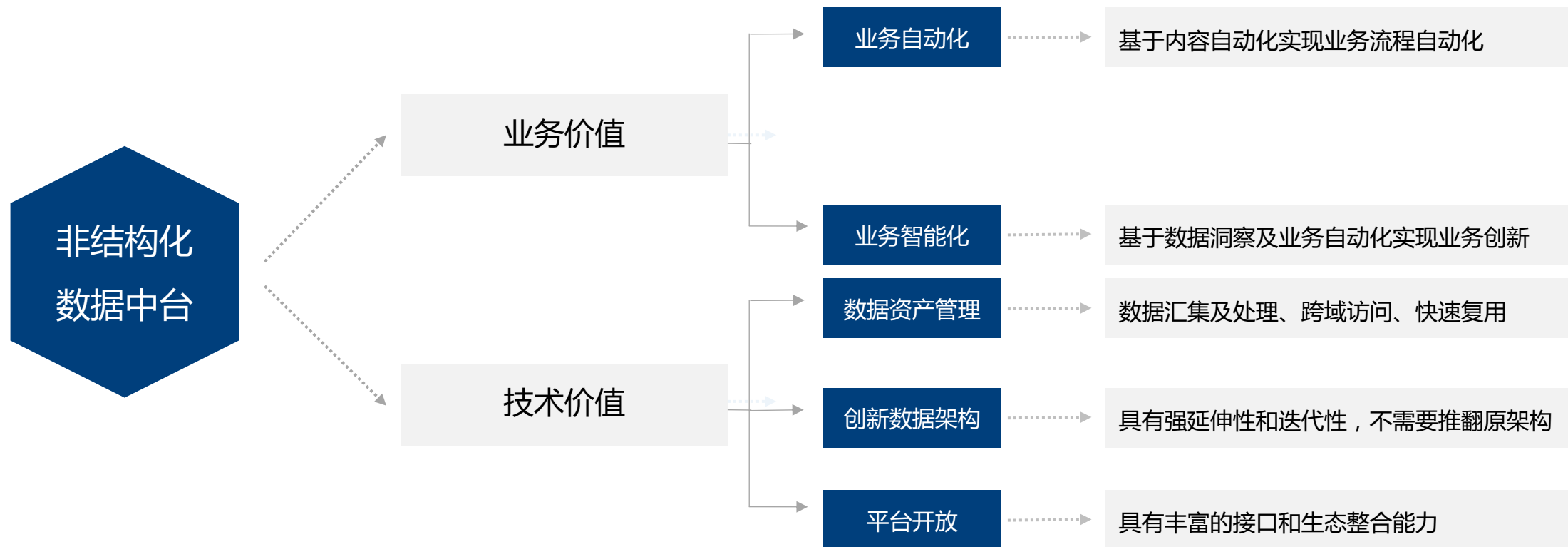


智能知识运营

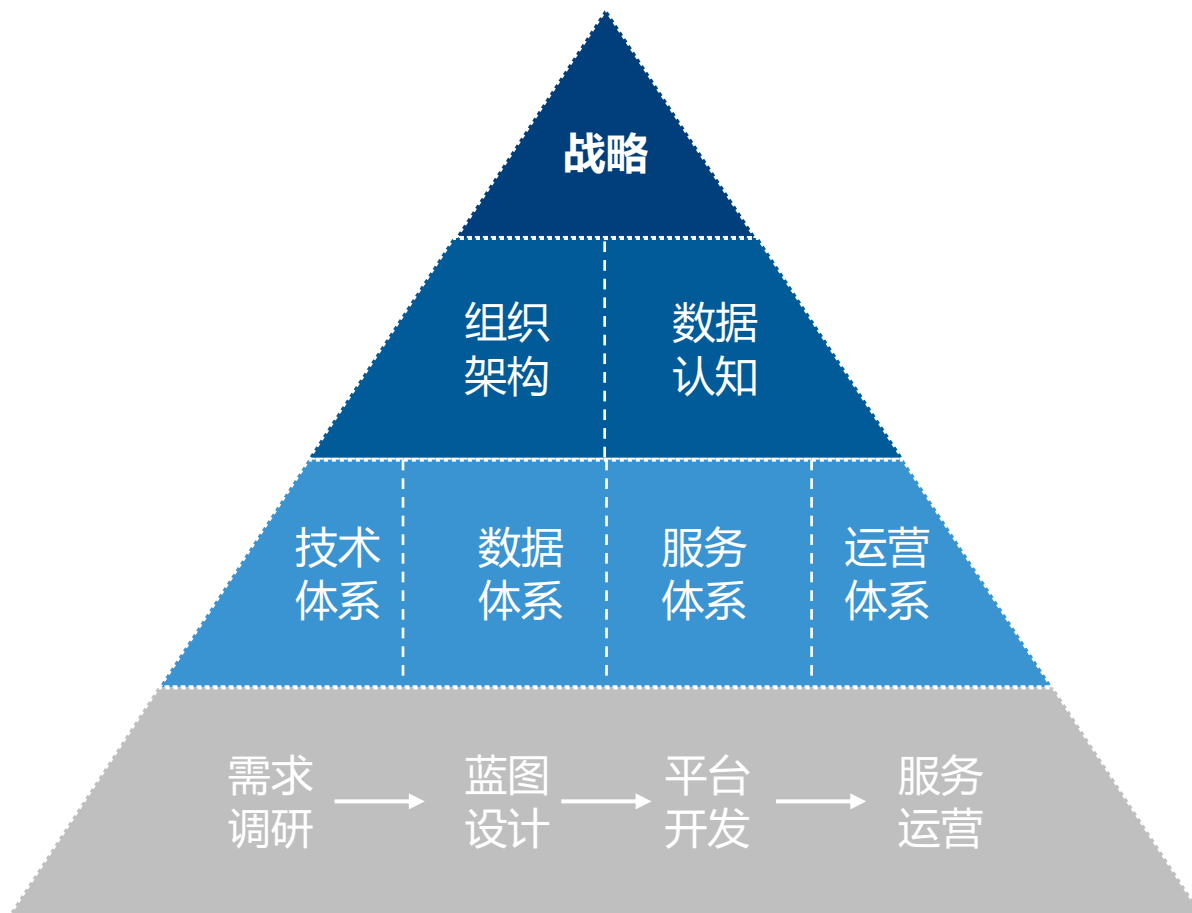


业务流程自动化

- 数据收集和使用合规体系及监管体系不完善，数据安全面临挑战。
- 非结构化数据分散于多渠道、数据种类多样、数据量大、长期保存难、使用率低。
- 业务涉及数据繁杂，整理困难，耗时耗力，缺乏自动化能力，效率低，成本高。
- 业务与内容管理脱节，企业管理缺乏知识来源与知识管理运营。



非结构化数据中台的价值主要体现在两方面：业务价值与技术价值。业务价值主要体现在赋能企业业务与商业模式创新，技术价值在于低成本实现数据治理及复用。



非结构化数据中台的建设要从战略、保障支撑、内容、步骤等方面考虑

- 战略：非结构化数据中台需要定位于企业级战略。
- 保障支撑：结构化数据中台需要企业组织架构保障和企业数据认知的支撑。
- 内容：非结构化数据中台的建设主要包括技术体系、数据体系、服务体系和运营体系等内容。
- 步骤：包括需求调研、蓝图设计、平台开发与服务运营。

企业认知挑战

战略路径

- 企业对数字化的需求明确，但是对实现数字化的战略路径不清晰
- 90%企业对非结构化数据中台不了解，大多企业没有区分结构化数据中台和非结构化数据中台



解决方案

将非结构化数据中台
定位于企业战略进行推进

组织流程

- 非结构化数据中台技术新，且处于起步阶段，落地时容易遇到阻力且效果不及预期



为非结构化数据中台提供
组织架构保障

数据认知

- 公司的业务流程与员工的数据思维不适用于非结构化数据中台



非结构化数据中台需要
企业数据认知支撑

建设过程挑战

需求调研&蓝图设计

- 业务需求不易理解，对业务场景的挖掘和理解难，和业务部门的沟通难



解决方案

建立为不同行业和
客户服务的行业纵深

平台开发

- 技术挑战多，包括算法优化以及将技术更好地深入应用于业务场景中



持续技术深化，
帮助客户实现更多的业务赋能

服务运营

- 企业对数字化的需求明确，但是对实现数字化的战略路径不清晰

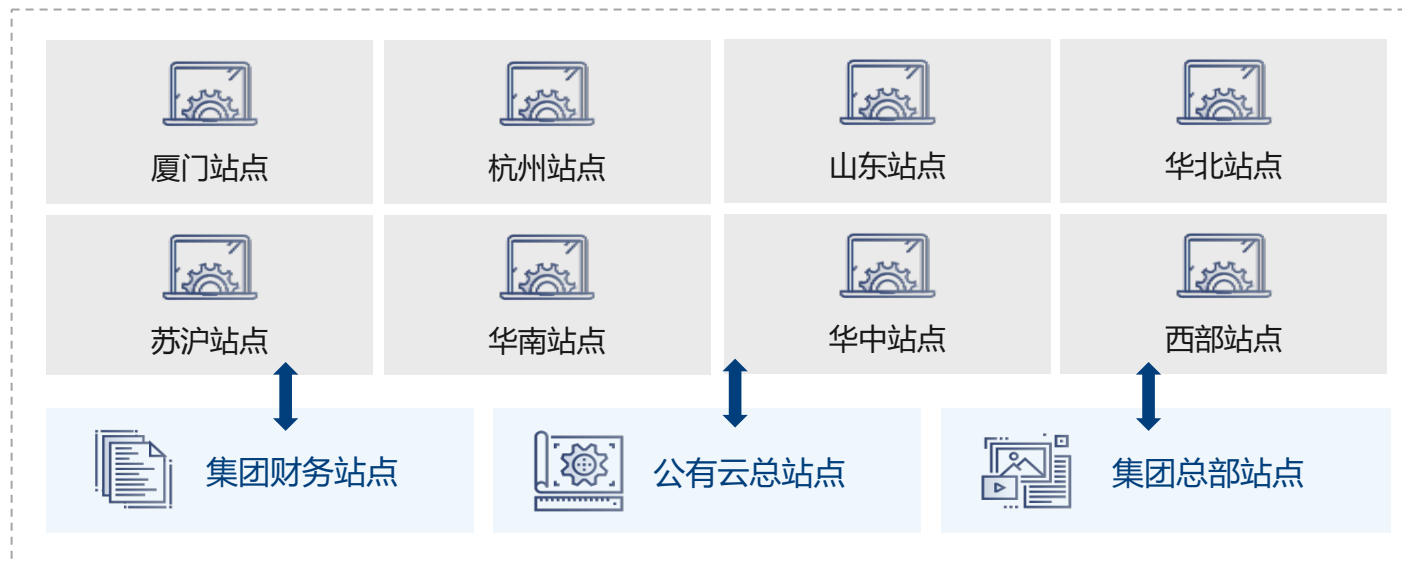


TPA交付方法论——
端对端的服务体系，构建懂行的交付体系

T (Think) -P (Plan) -A (Action) ，是从客户的数字化战略以及业务模式出发，为客户提供专业的端到端的咨询、开发和交付的服务方案，并通过大数据基础设施进行有效落地，帮助客户实现数字化战略，带来实质性的投资回报。



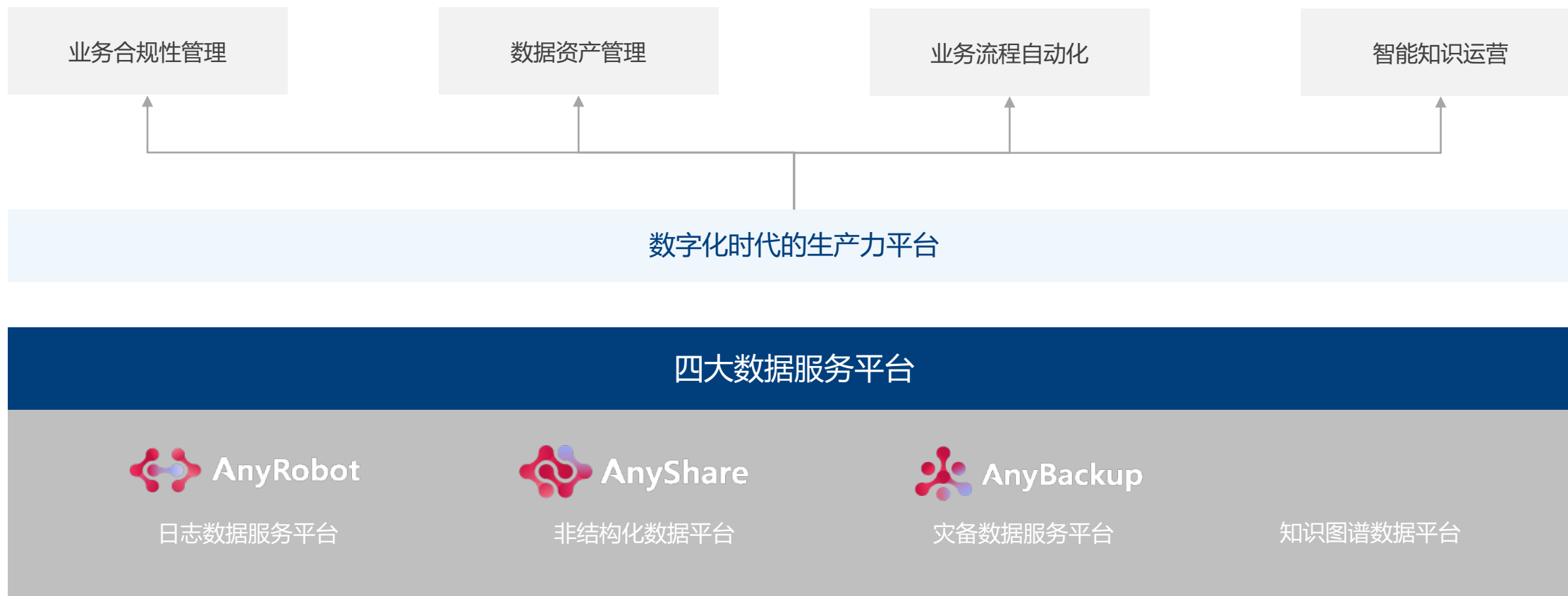
- 文档体系与规范梳理，数字资产有序管理
- 总部与分支机构的有序文档共享协作
- 集成知识应用，企业知识管理的底座

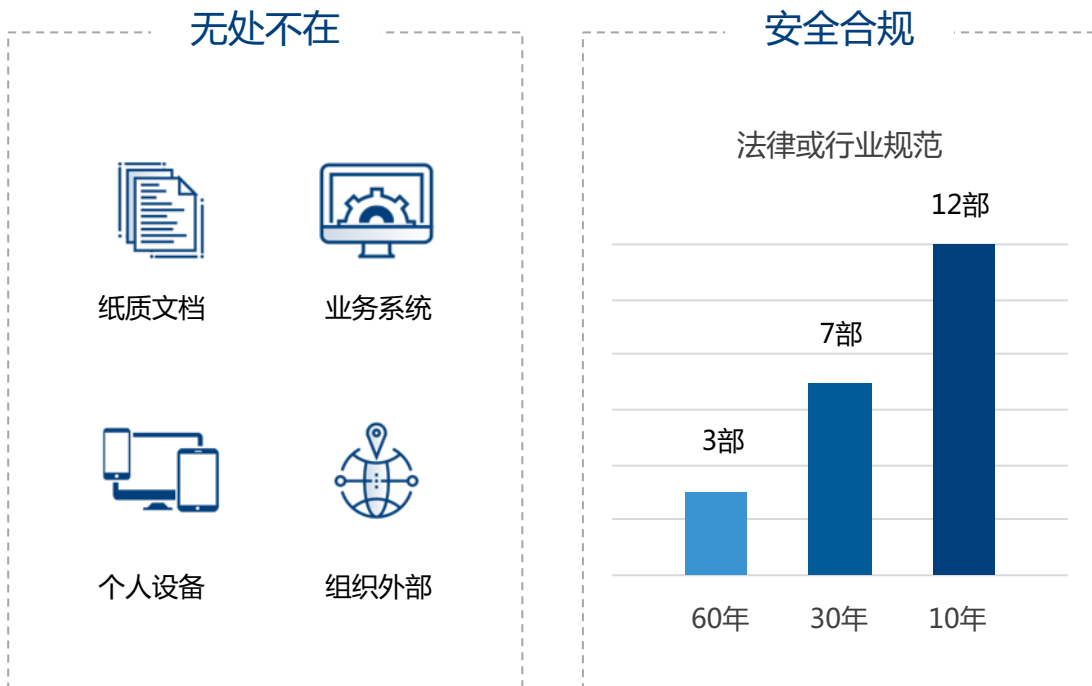


基于非结构化数据中台的应用场景

03

根据客户需求，基于非结构化数据中台的应用场景有：数字资产管理、业务流程自动化、智能知识运营、业务合规性管理等。





- 非结构化数据面临组织合规内控体系、行业监管合规与法律合规问题。
- 但非结构化数据存储分散，数据安全把控难，内容合规审核难，数据与内容安全面临挑战。

挑战1：非结构化数据备份性能挑战

非结构化数据具有数量多，且以海量小文档组成，传统的数据备份方式在备份海量非结构化数据时，会遇到非常明显的性能挑战，使得备份无法有效实施。

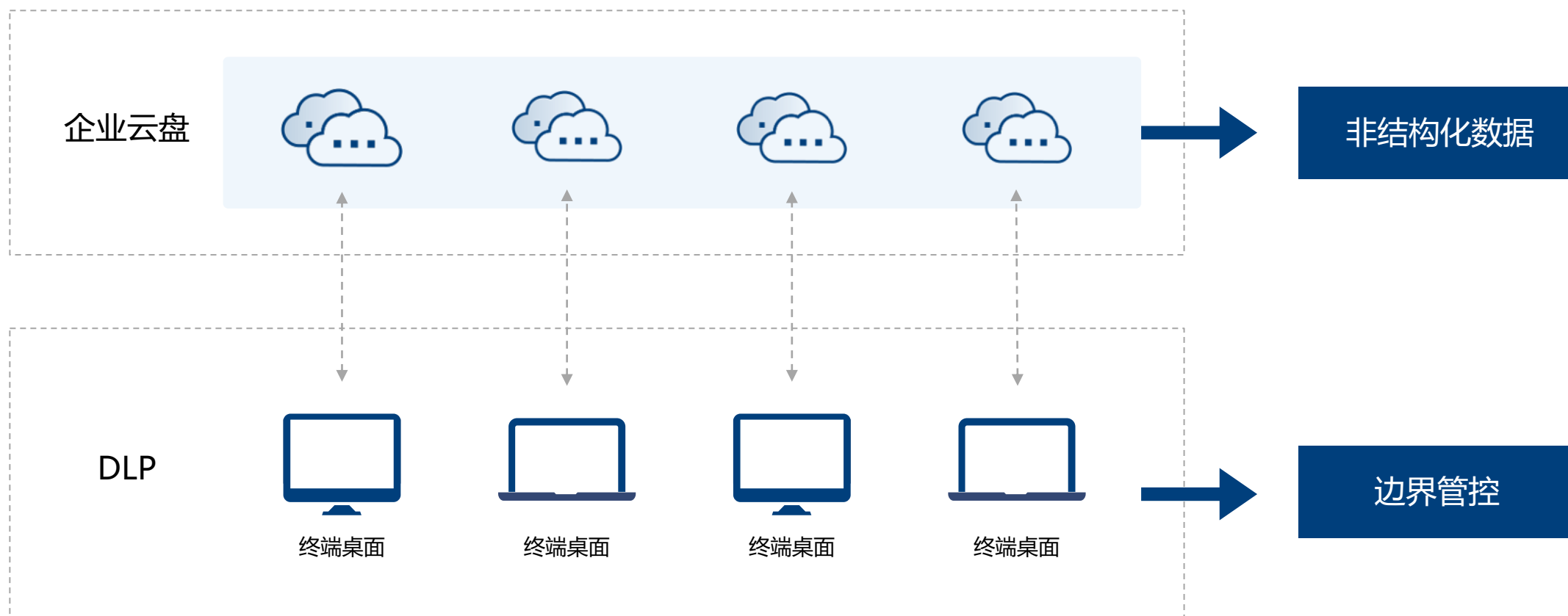
挑战2：内容全方位管控难

企业安全合规是多维度的，目前企业在内容上传下载过程中没有统一的方式进行访问边界控制以及非法或敏感内容识别措施，无法从多个层次上进行全面管控。

挑战3：数据追溯难

非结构化数据存储分散，企业在数据上传过程中没有进行安全管控，造成数据追溯难。

传统的安全方案是企业云盘+DLP，仅解决了文档层面的安全问题，但没有解决在非结构化数据的多层次方面的问题。



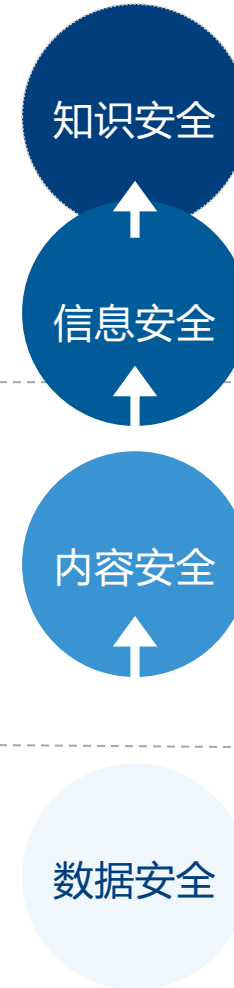
立体安全体系才能够满足业务合规性

- 企业需要搭建包括数据安全、内容安全、信息安全、知识安全等在内的立体安全体系。
- 非结构化数据中台将数据汇集，同时本身有强大的生态整合能力，可以整合第三方等应用或程序帮助客户做好数据安全合规。

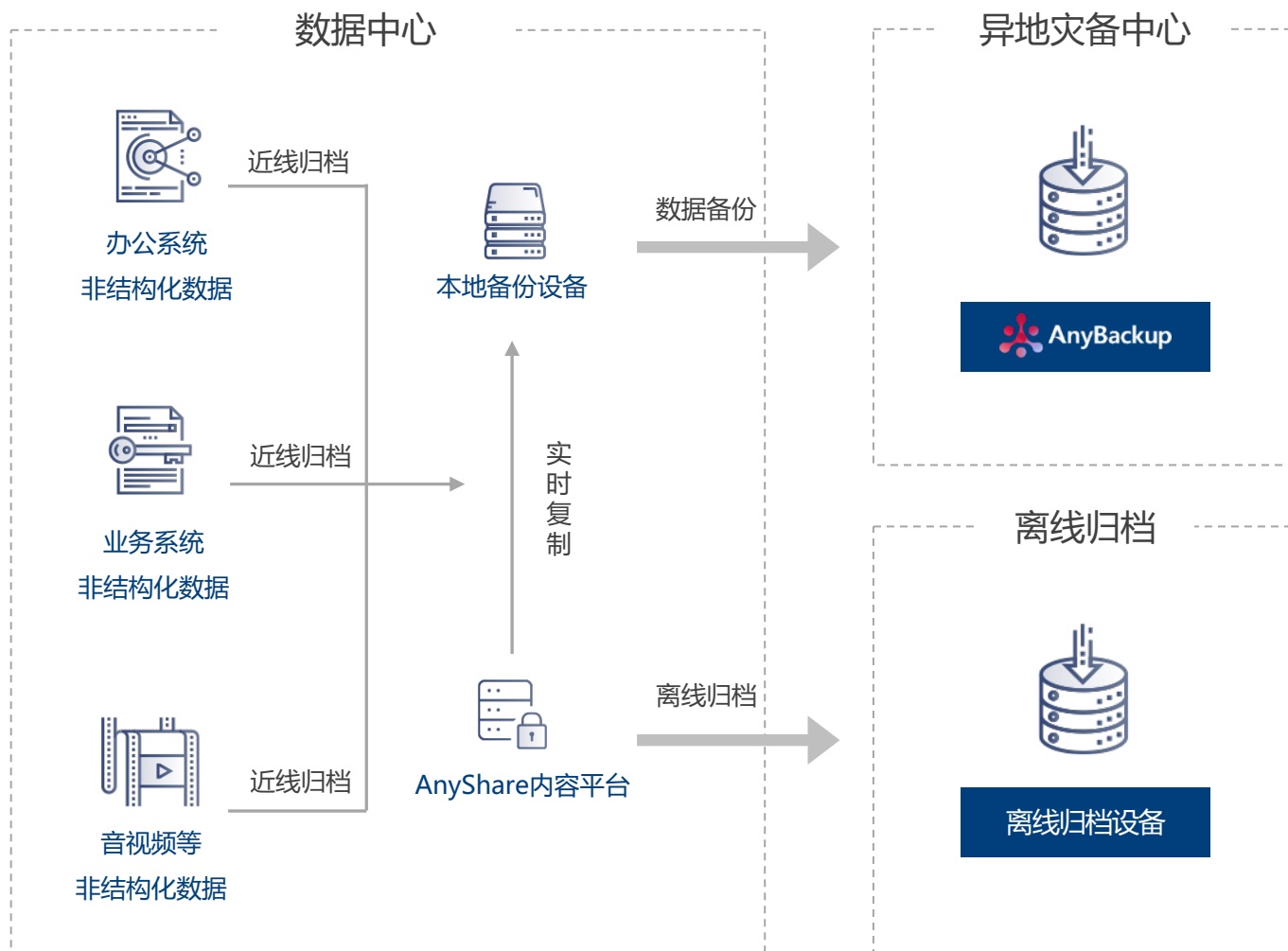
- 《个人信息保护法》,2020
- 《数据安全法》,2020
- 《上市公司信息披露管理办法》,证监会
- 《通用数据保护条例》,欧盟,2018 (GDPR)

- 《电子公文归档管理暂行办法》, 2003
- 《商业银行资本管理办法》, 2013
- 《金融企业业务档案管理规定》, 2015
- 《药品生产质量管理规范》, 2010 (GMP)

- 《网络安全法》, 2017
- 《信息安全等级保护管理办法》, 2007
- 《涉及国家秘密的信息系统分级保护管理办法》,2005



- 个人隐私数据
- 企业敏感信息
- 非法内容管控
- 内容边界安全
- 内容访问审计
- 文档的生命周期管理
- 数据生命周期管理
- 数据备份恢复



✓ 数据可进行在线或离线备份：

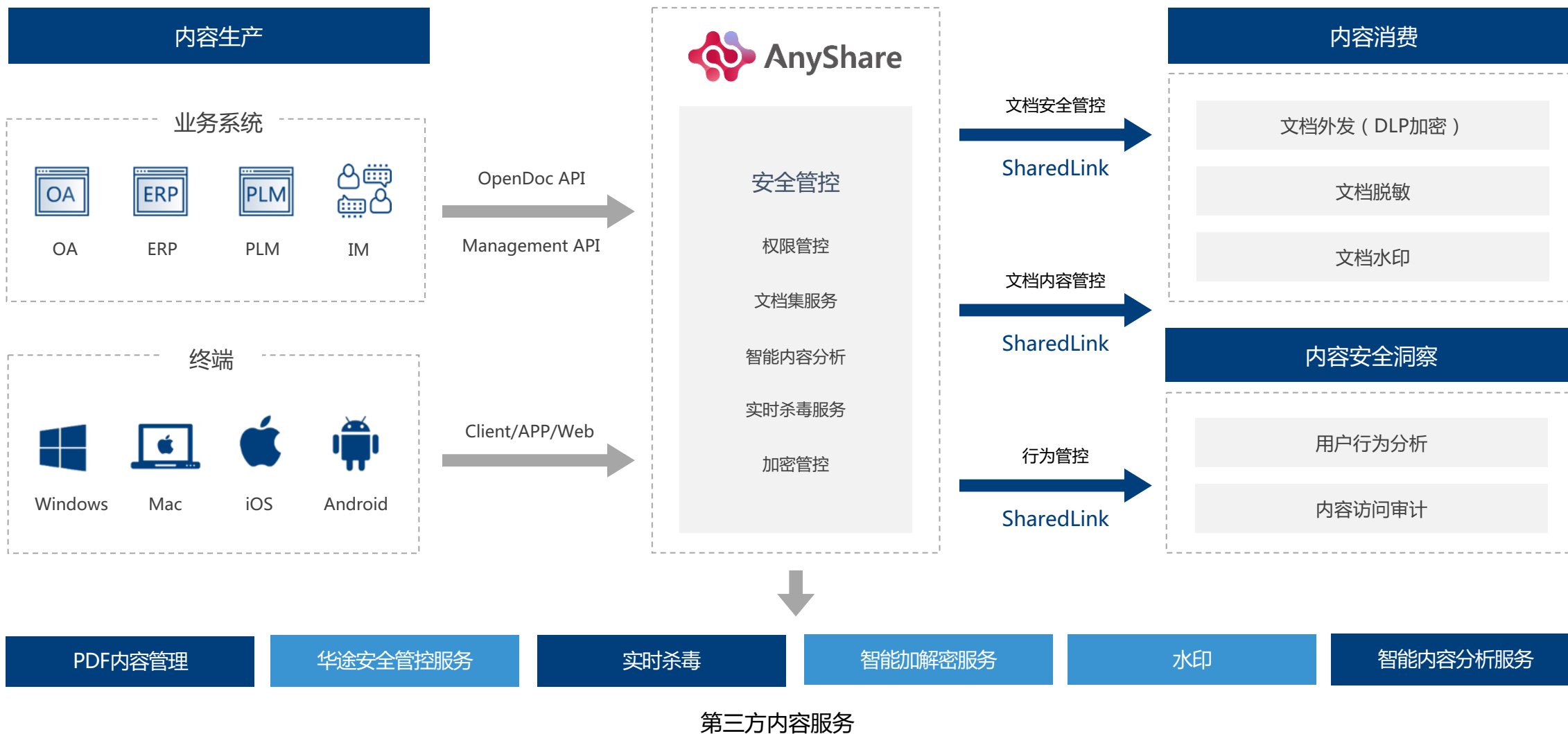
温数据近线归档到内容平台，随时可查询、可回调。

冷数据离线备份到AnyBackup 一体机。

✓ 电子文件归档一式3套。一套位于现有存储，一套位于内容平台中提供利用，另一套位于异地灾备中心保管。

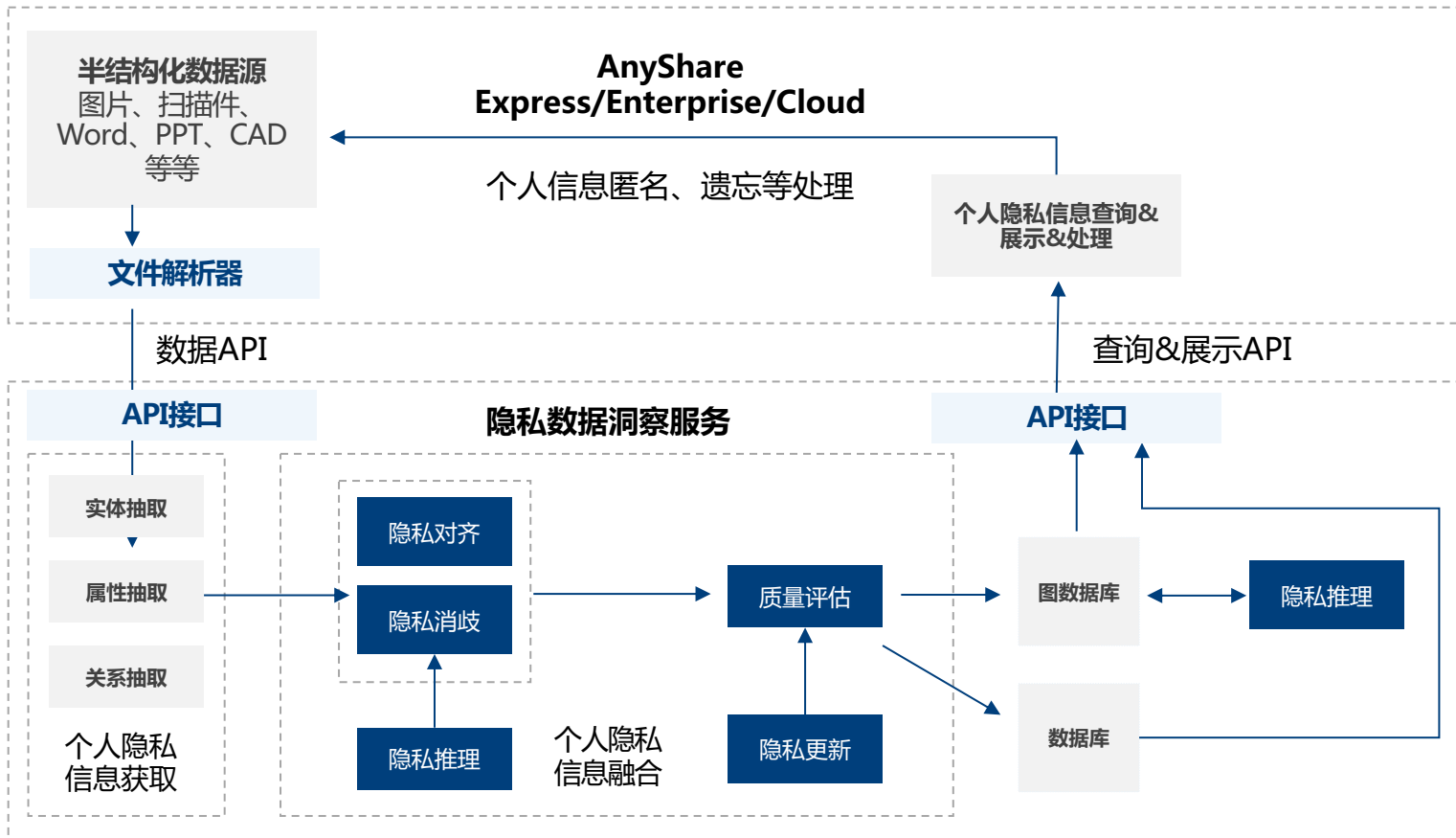
✓ 小文件备份，比业界方案高出 **50倍以上**的备份性能。

非结构化数据中的知识安全与合规



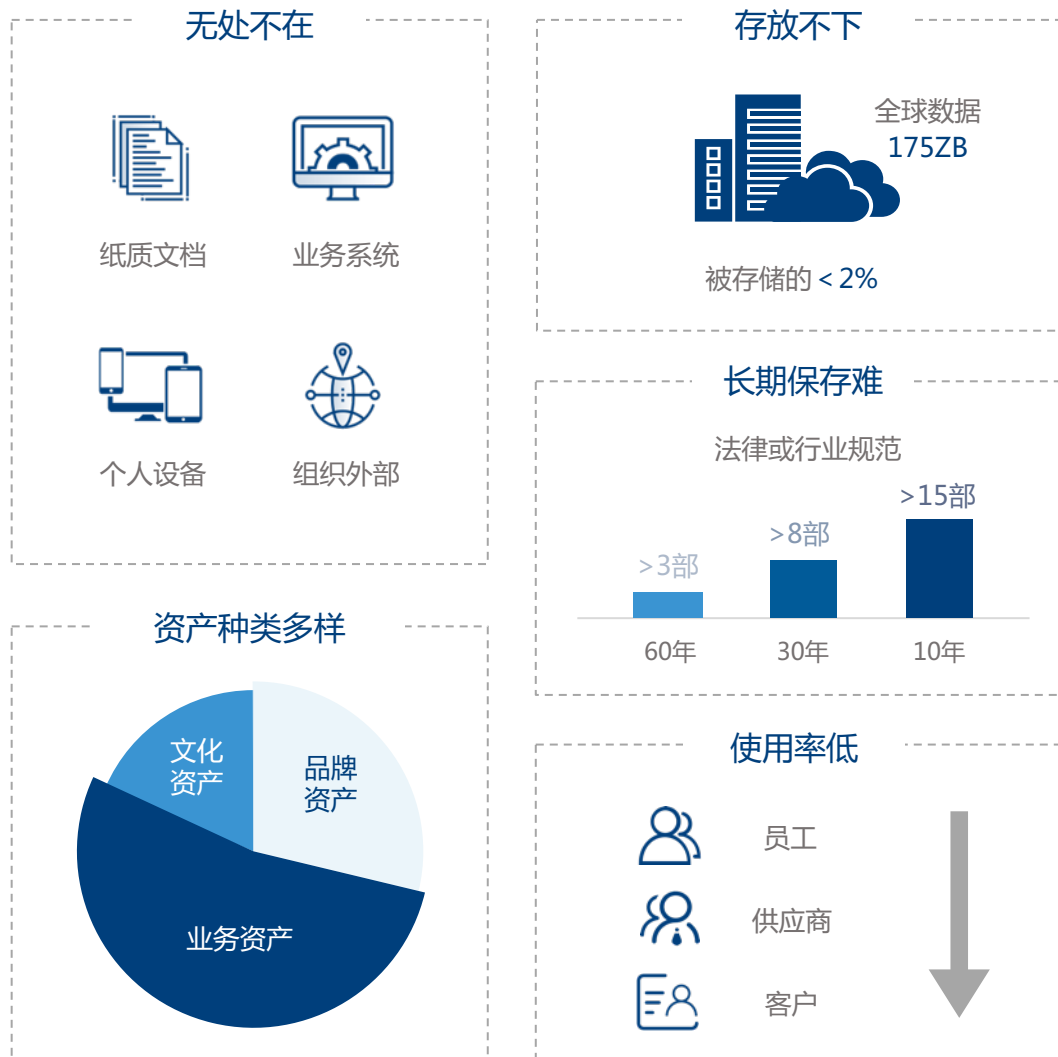


知识图谱 (Knowledge Graph) ，在图书情报界称为知识域可视化或知识领域映射地图，是显示知识发展进程与结构关系的一系列各种不同的图形，用可视化技术描述知识资源及其载体，挖掘、分析、构建、绘制和显示知识及它们之间的相互联系。

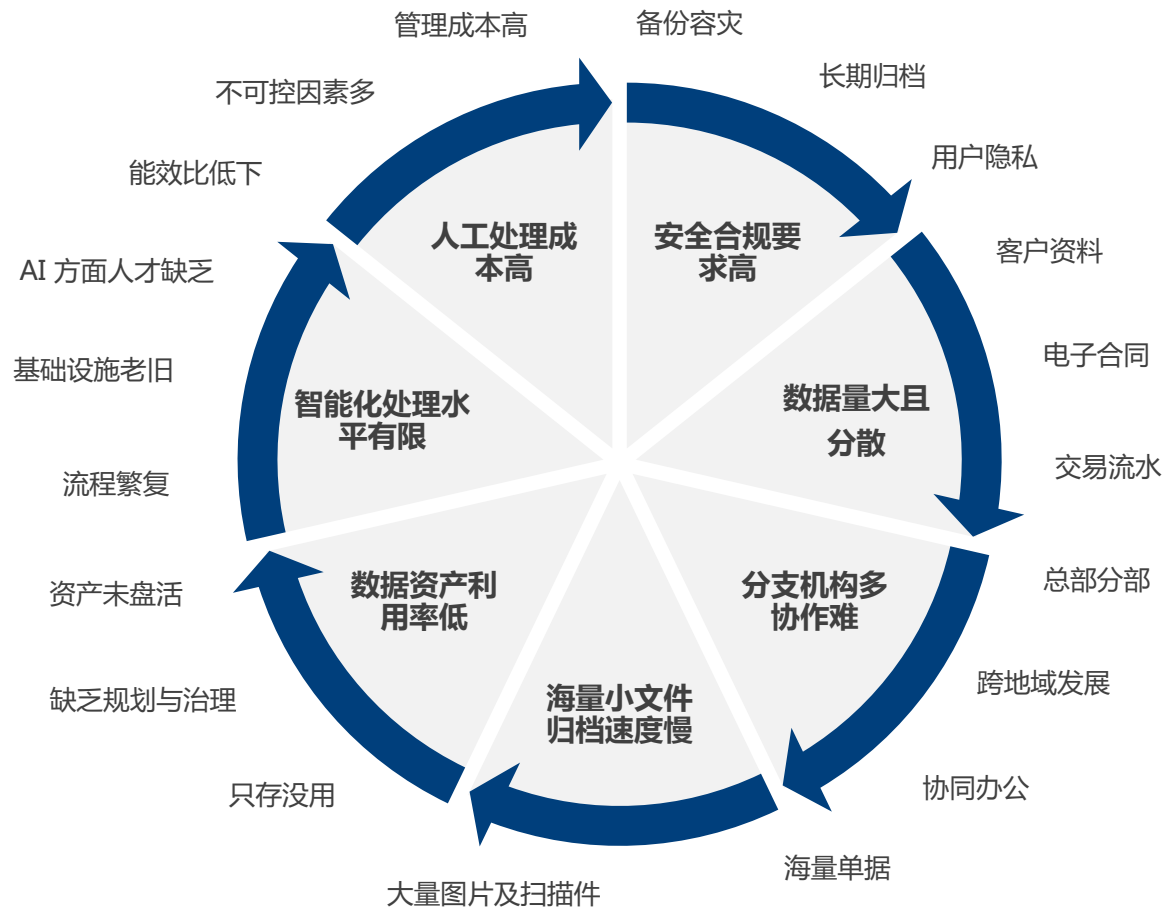


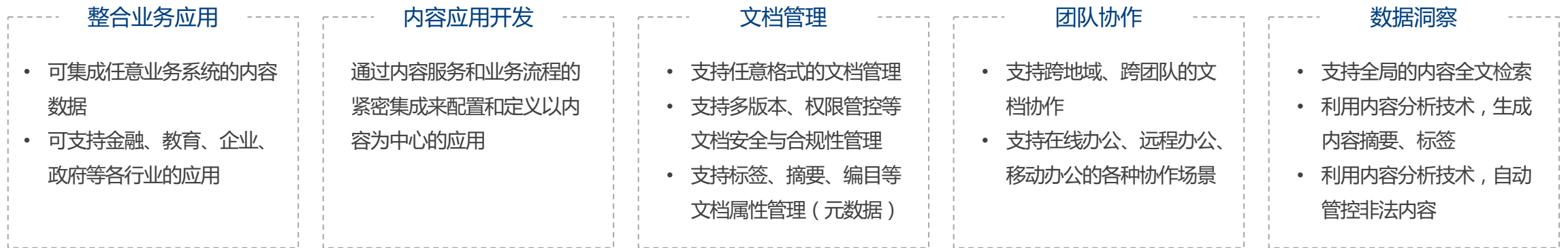
个人隐私数据洞察项目中采用知识图谱技术的目的：

- ①一段文字中出现人的部分信息，需要推理出隐含属性。例如“姚沁蕾是姚明和叶莉的女儿”，不光有人物之间关系，还有性别等隐含信息。
- ②主要是用来围绕人将所有相关信息进行组织、呈现其各种实体及之间关系。并用于后续GDPR等法律场景下信息的查找和提取。

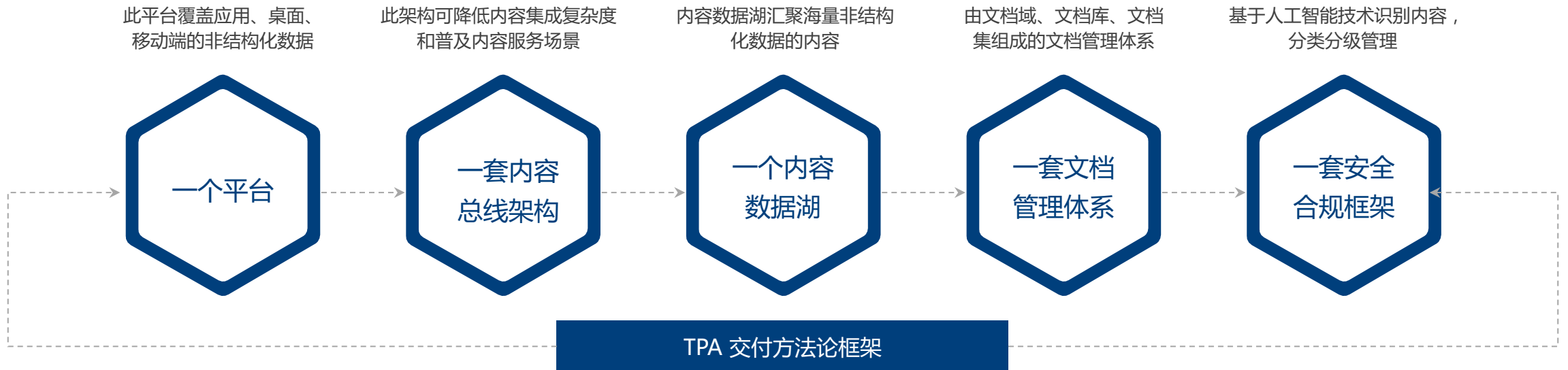


银行业在非结构化数字资产管理方面的挑战





智能内容云



传统文档管理体系

挑战1：分散存储，无法统一管理

- 传统文档管理体系相对杂乱，建设文件仓库但仍然随意存放
- 知识数据脱离业务流程，陈旧过时
- 存储资源重复建设，可扩展性差

挑战2：内容搜索效率低

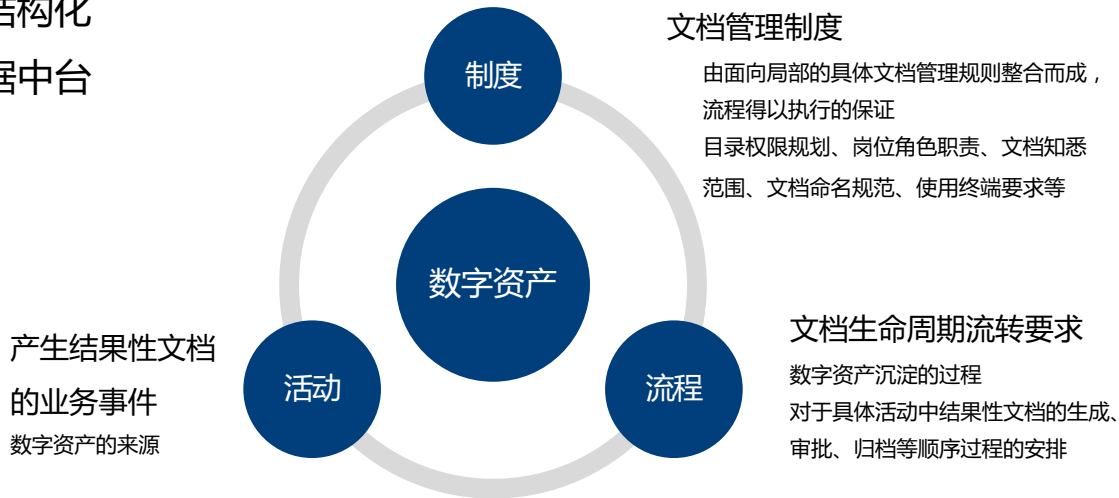
- 传统业务系统难以提供内容搜索能力
- 不支持图片、特定格式的内容搜索
- 搜索的文档版本过时
- 不支持元数据、标签等高级搜索能力

挑战3：业务系统重建数据难以整合

- 非结构化数据的管理模块需要同步重建
- 非结构化数据迁移成本高
- 难以再对接，进一步加剧数据孤岛

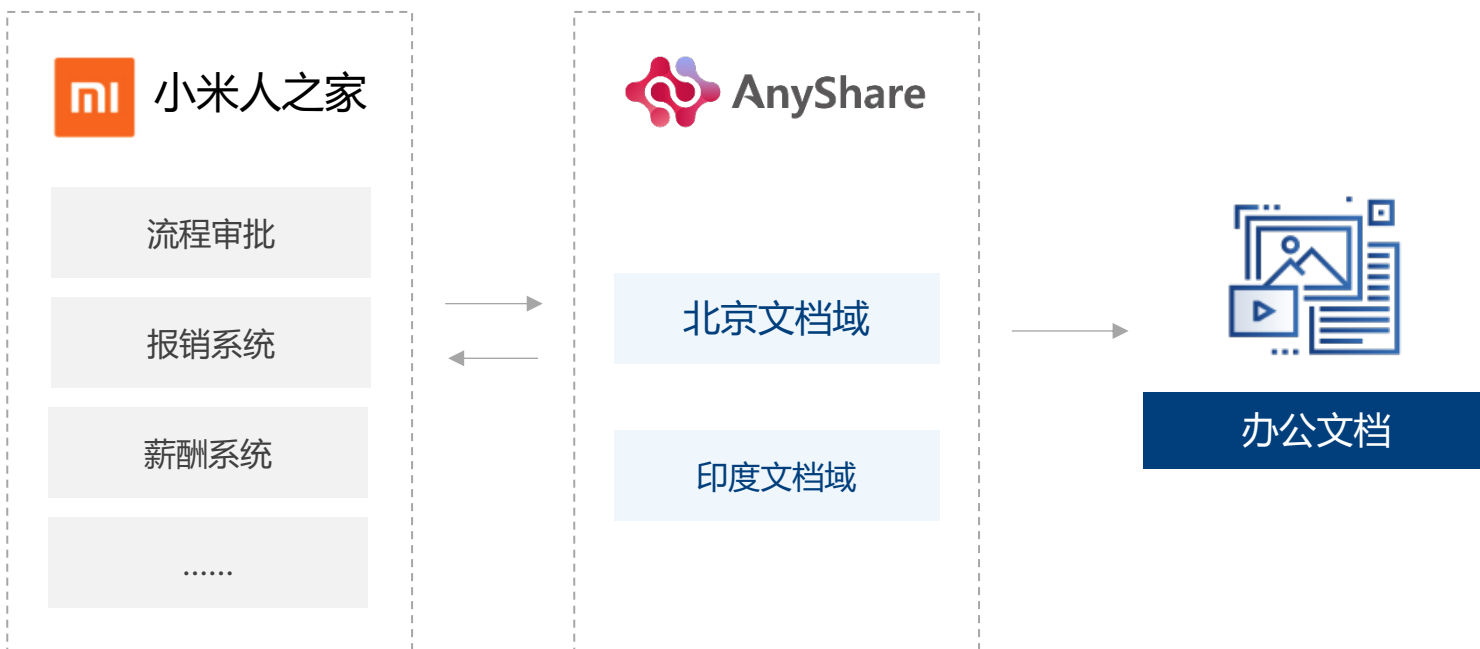
管理和运营数字资产

非结构化数据中台



围绕结果型文档的数字资产管理

- 全面梳理数字资产定义、活动、流程，完成管理制度优化
- 基于非结构化的数据中台将数据打通，从源头保障数字资产复用能力。
- 基于非结构化的数据中台能够定义数字资产类型，缩短平均查找时间，提升运营效率，提高关键活动中的文档协作。
- 重构底层框架，架构可延展，更新成本低。



- 全球部署，服务全球小米人及3000家海内外供应商
- 全球协作，安全、可靠、易管理的“小米企业网盘”
- “小米人之家”背后的统一内容管理平台

例：保险业/银行业

客户多年前已经上了IBM CM 系统，用于影像数据的存储从14年到现在已经存储了10TB 左右的数据，主要通过企业NAS存储。近两年数据增长较快，每年数据增长在5TB左右，企业数据各式各样，数据增长来源于档案、信贷、无纸化等业务，目前已经对接20个业务系统。

痛点

- 数据不规范。合作方均按自己的方式提供的数据，仍然存在不规范的数据，数据归集标准各式各样，数据分类不一致、业务术语不一致、数据量级分类不一致，整理困难；需要大量人力处理。进行客户标准与公司标准比对，才可能实现数据的整合和应用。
- 无法识别图片。
- 信息图片数据无法自动搜索。企业有6000万图片的存量，自动搜索难。

需求

- 降低出错率、提升生产效率、降低人力成本；操作可监控、短期内产生效益。
- 需要简单方便的OCR识别平台，可以针对不同类别的数据，配置相对应的模板识别。
- 针对多种影像类型数据，实现海量图片的分类及查找。

非结构化数据中台VS.传统表格数据处理

ifenxi |

传统表格
数据处理



收集合作方数据
数据组织方式不一致
数据分类不一致
业务术语不一致
数据量级分类不一致



按特定标准转换处理
全靠人工
全靠经验
全靠理解
比对完再进行核对



整合为标准数据输出
建立标准对应
按照标准汇集所有数据
为后端 BI 提供汇集数据接口
提供数据协作进行数据再处理

表格数据自动化流程

非结构化
数据中台

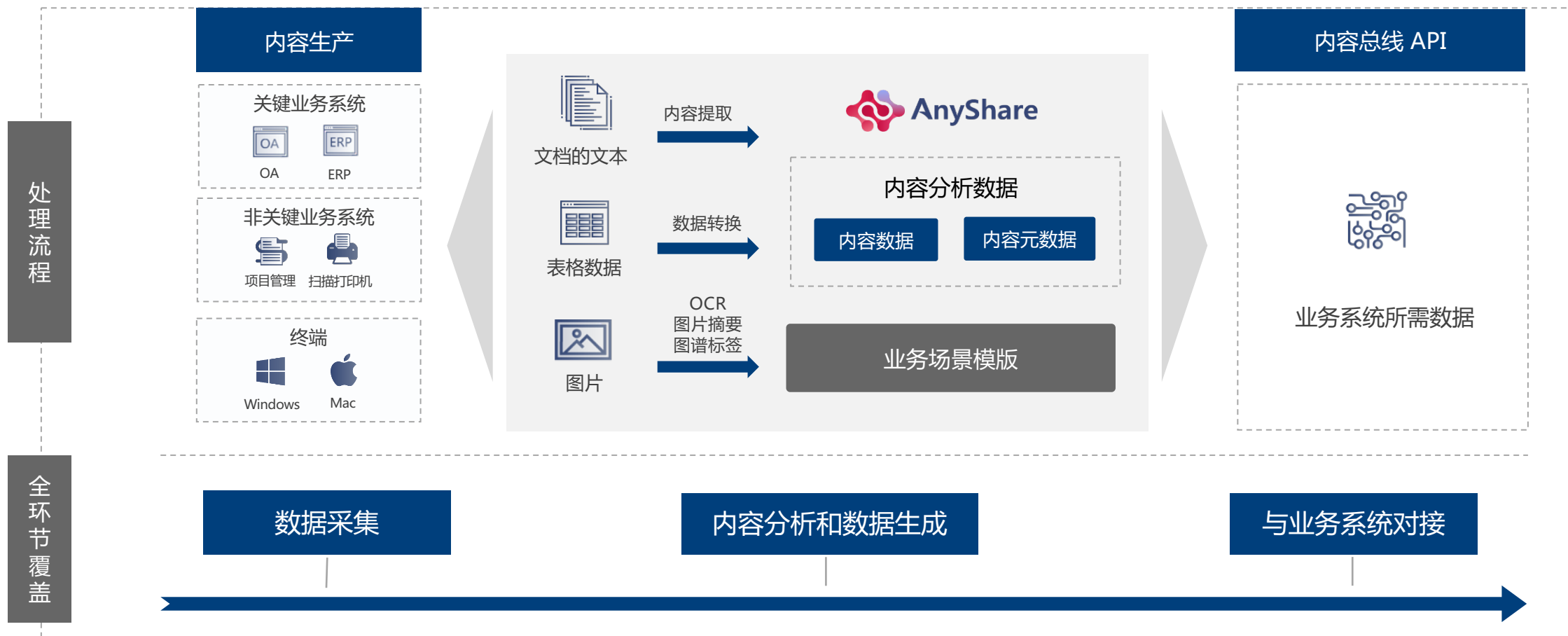
数据转换

智能比对

数据校验

智能表格数据平台

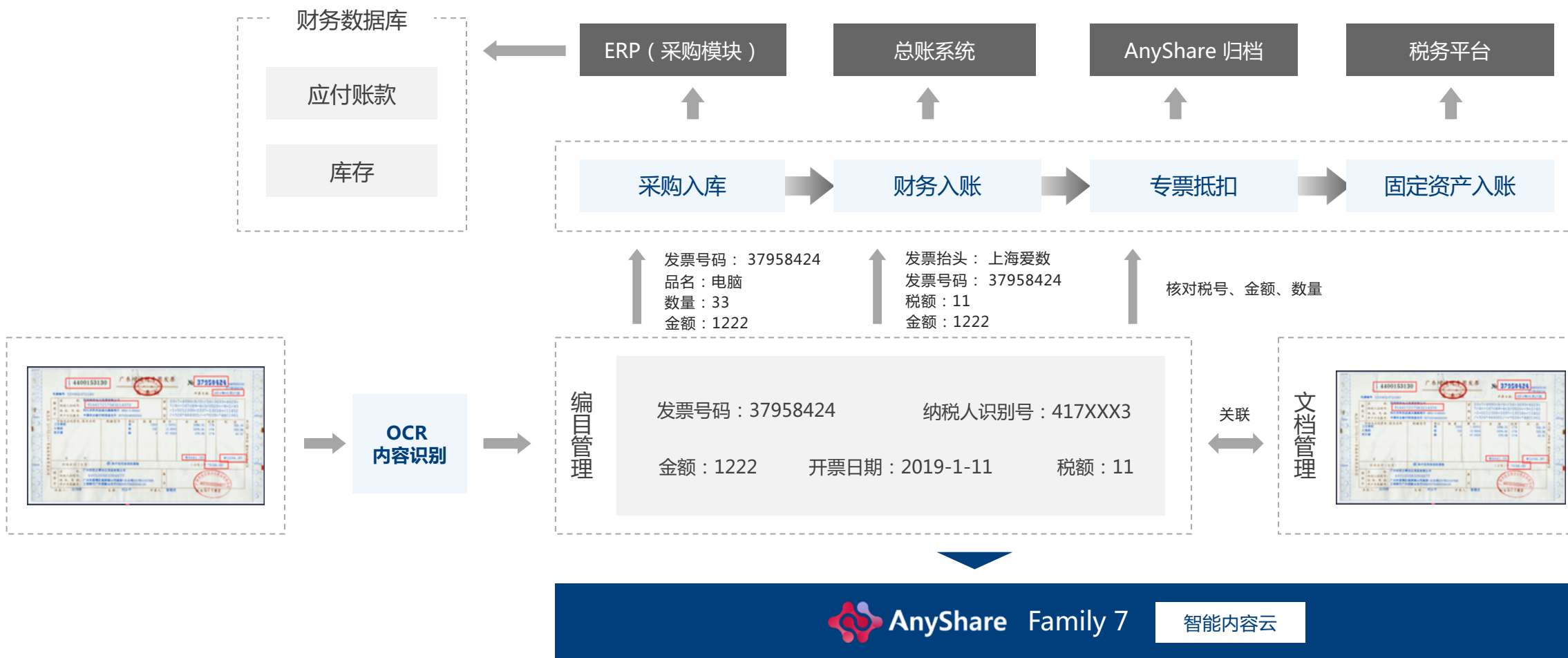
- 非结构化数据中台和业务系统对接，数据可复用，实现持久自动化；可定制标准表单进行规范数据的收集。
- 集成人工智能技术，可将内容转化为文字或生成标签，利用智能比对和数据校验模块实现流畅的业务流转，提升业务效率。
- 标准数据全部整合到智能表格数据平台上，通过API提供标准数据格式；利用数据多人协作进行数据再处理。
- 同时数据只存在于非结构化数据中台中，不浪费存储空间。



搭建非结构化数据中台实现内容自动化，核心技术为源数据提取、SmartSheets与OCR。

案例：采购场景实现票据内容自动入账

OCR内容识别+ 内容自动化



知识运营贯穿企业的经营全局：建立知识复用体系，提升企业的知识传承与复用，实现降本增效。



传统的知识管理

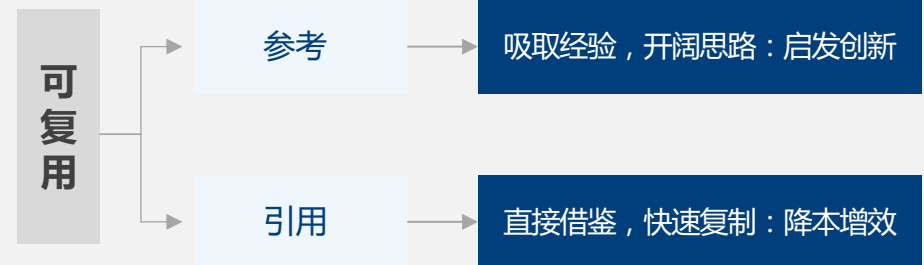
1. 重心放在社区、激励等应用层面，单纯依靠运营驱动，没有资料，不活跃，没有人气；
2. 不重视知识库文档管理/内容管理的建设，导致缺乏来源，知识库难以维护！

痛点

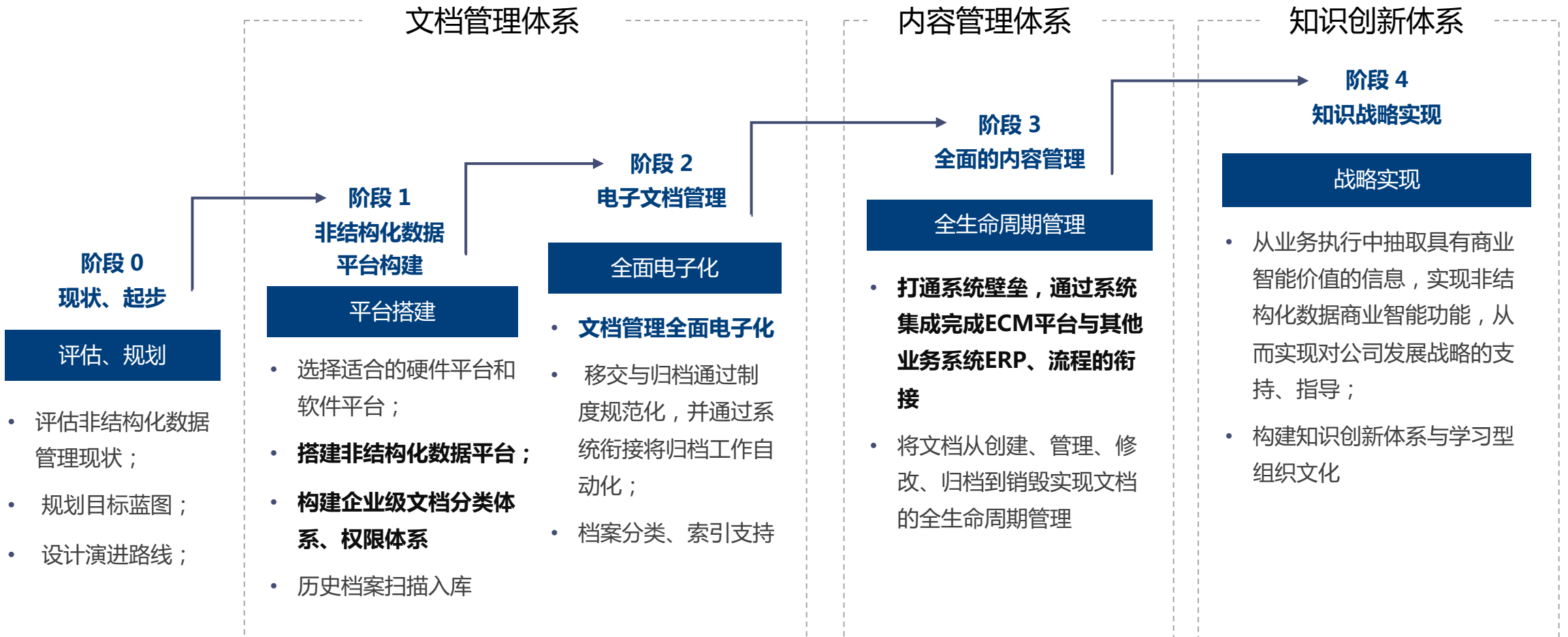
- 面对海量数据增长，传统知识管理的底层架构很难解决知识库的文档/内容海量存储和扩展的挑战。
- 知识不完整，只涉及办公/培训文档，不涉及业务流程中的各环节；
- 数据割裂，且没有探索性的知识分析的能力，从而无法做到知识的更新或持续反馈，进而无法实现真正的知识沉淀；
- 业务更新与知识更新不同步，知识管理滞后，员工知识维护和更新的动力不强。

基于非结构化数据中台的内容管理

- 在底层将数据重构，易拓展。
- 利用人工智能&知识图谱实现内容洞察以及知识服务，通过探索式分析对知识进行融合推荐，并不断训练反馈，让知识逐步得到认可并广泛推广，将知识管理和运营建立起来。
- 能够实现精准的知识搜索。

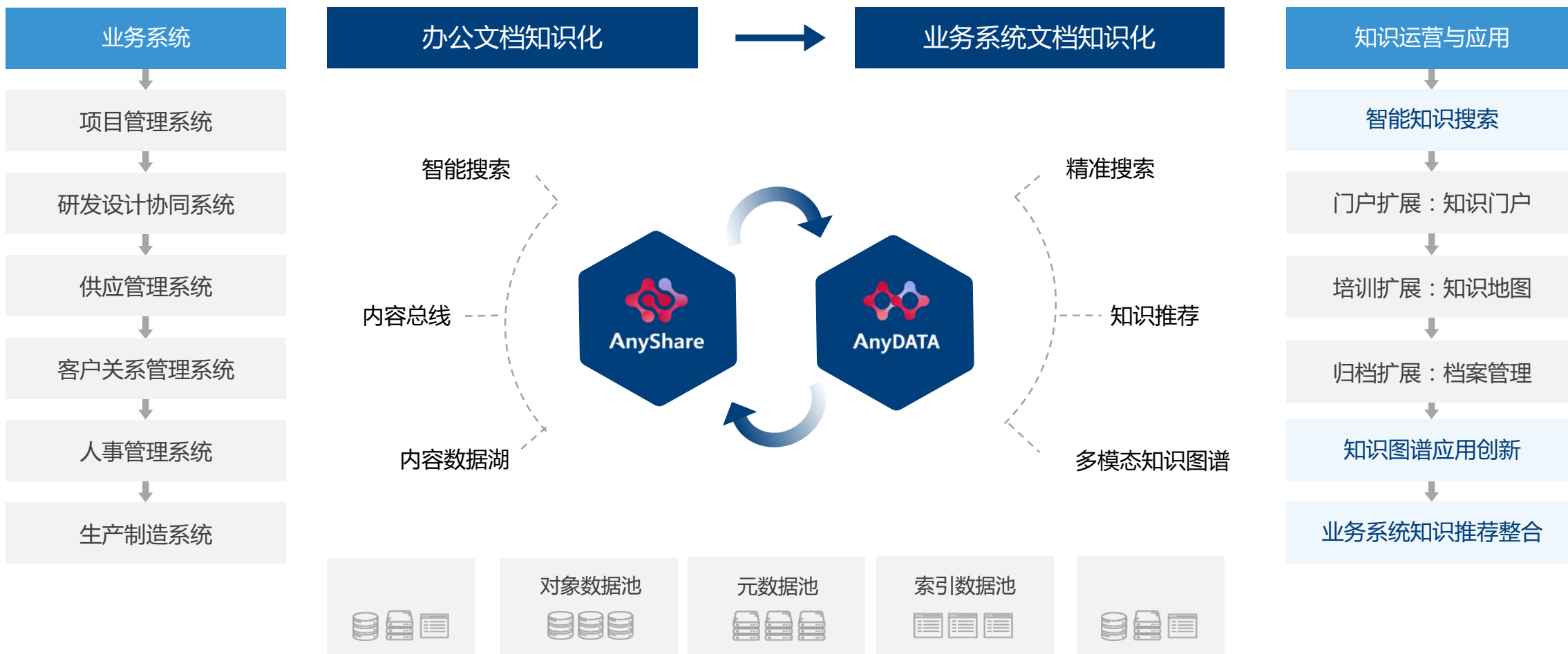


数字化企业的知识战略蓝图

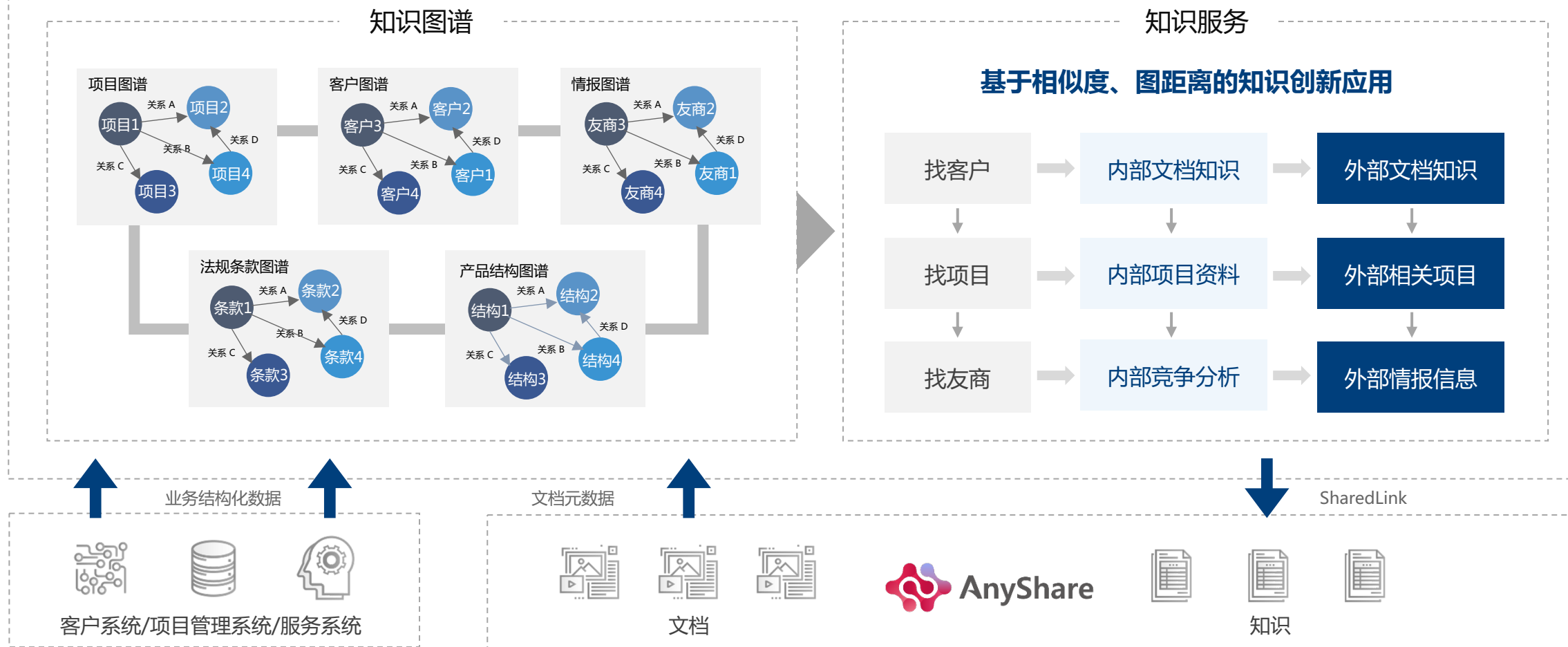


知识运营的解决方案

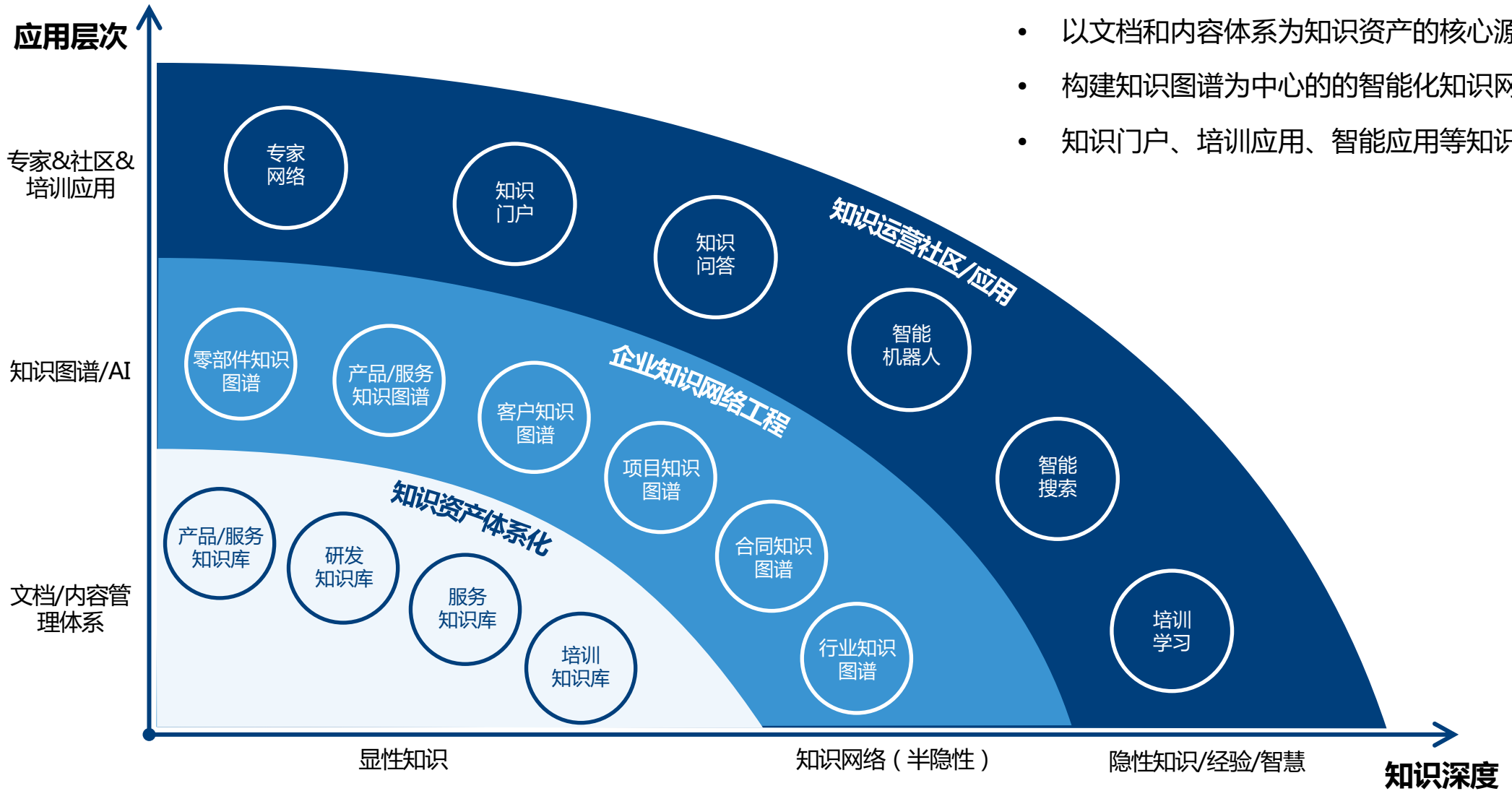
基于多模态知识图谱的智能知识运营解决方案



AnyDATA 多模态知识图谱平台



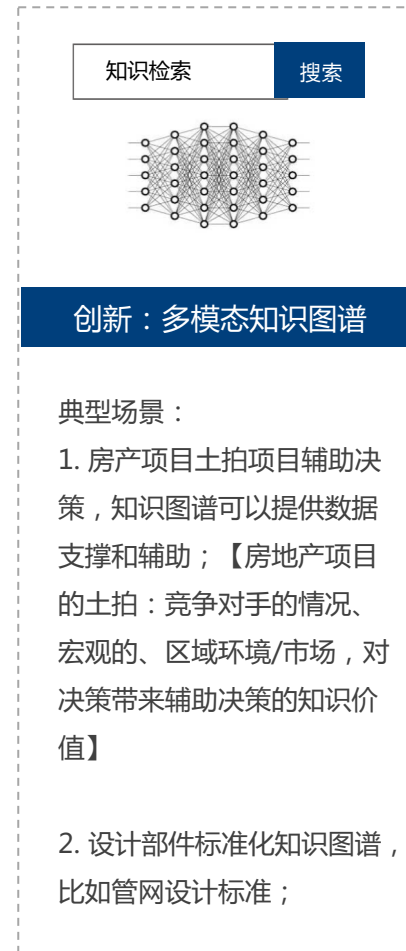
目标：通过内容管理体系和知识图谱，创新智能知识运营体系



- 以文档和内容体系为知识资产的核心源头；
- 构建知识图谱为中心的智能化知识网络；
- 知识门户、培训应用、智能应用等知识运营社区；

案例：设计企业的项目知识运营

建立贯通项目和技术运营的知识运营底座



- 从项目管理、设计协同、图纸出版到最后项目归档的内容知识流转复用；
- 设计协同过程基于项目管理，项目完结自然形成项目设计知识资料，可供新项目参考复用；
- 档案管理、知识门户、培训学习可轻量级扩展，也可以与成熟系统对接，实现一个开放、统一、可复用的知识运营体系；
- 扩展知识图谱，服务于特定场景的精准知识搜索、探索式分析与推荐场景，包括土地拍卖、设计部件标准化知识库等。

展望行业趋势

中国非结构化数据中台
实践白皮书

04

新行业需求



传统行业数字化转型加速伴随着非结构化数据中台的应用行业边界扩大

新业务场景需求



非结构化数据来源于企业各个业务场景，随着企业对数据应用的需求增加，更多的业务场景需求被激发。

人机协同需求



机器无法完全取代人类，但数据智能辅助并赋能员工、实现人工智能是行业新的诉求。

中台能力输出

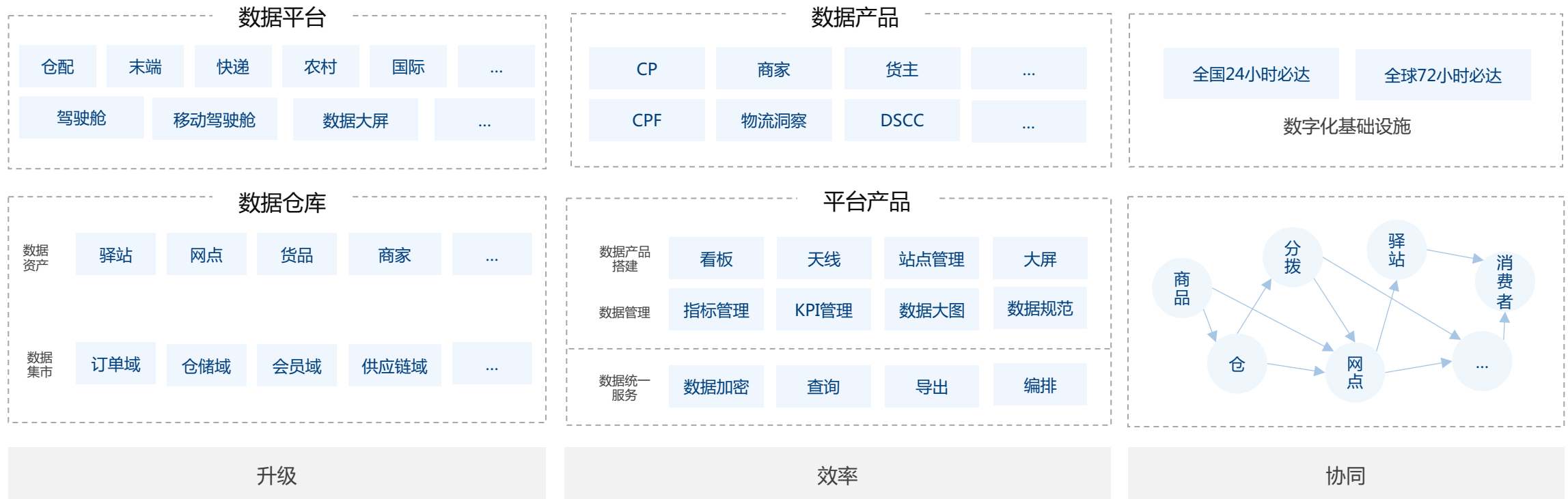


随着企业非结构化数据中台的成熟，企业也将中台能力输出，赋能产业上下游企业。

案例：菜鸟数据中台能力输出行业协同

菜鸟数据中台通过“平台产品”加上“数仓”支撑菜鸟各个业务线的数据化运营工作，具有良好的数据化运营，但仍然需要与生态伙伴合作提升物流网络效率。

菜鸟将其数据、工具、方法论输出给生态合作伙伴，共同提升整个物流网络的效率，最终实现其使命：全国 24 小时必达，全球 72 小时必达。



非结构化数据中台并非企业数据的中转站，是能够实现智能推荐、领域图谱、AI 决策为一体的平台。因此，新技术与非结构化数据中台的融合是未来赋能企业业务的必然趋势。



机器学习

- 针对图像、视频等媒体信息的人工智能、机器学习技术



知识图谱

- 在不同行业领域，用可视化技术描述知识资源及其载体，挖掘、分析、构建、绘制和显示知识及它们之间的相互联系。

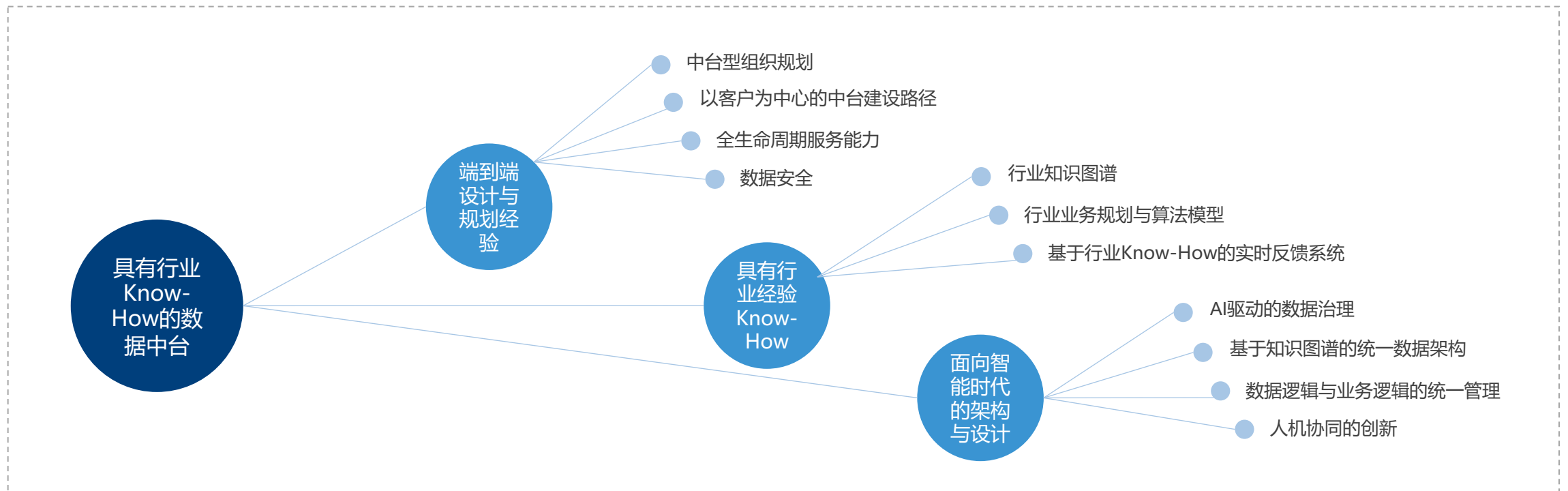


5G与边缘计算

- 5G将成为线下数据新基础设施，5G网络将打破跨场景互联最关键的技术瓶颈
- 边缘计算加速企业侧数据处理的速度。

知识图谱技术，是新一代数据中台最核心的技术，因为通过知识图谱可以实现探索式的分析。任何两个数据节点之间都可以建立关系，并进行分析、关联、探索，就像人的大脑一样。

知识图谱之于中台，融合了从智慧到知识的过程，包括规则、算法、推理等程序性支持，有了这些支撑，才能丰富数据中台相关的数据服务能力。同时，知识图谱需要中台来为其提供完美适配的环境。



案例：蚂蚁AI能力

行业&场景



金融大脑



- 蚂蚁AI能力的背后是大量人工智能技术的融合，包括知识图谱、NLP、强化学习、视觉分析等，其数据中台结合新技术共同组成金融大脑。
- 基于金融大脑支撑大量银行、保险、证券基金等多行业的场景应用智能化，包括智能风控、智能信贷、智能保险等。



爱分析是一家中国领先的产业数字化研究与咨询机构，成立于中国数字化兴起之时，致力于成为决策者最值得信任的数字化智囊。

凭借对新兴技术和应用的系统研究，对行业和场景的深刻洞见，爱分析为产业数字化大潮中的企业用户、厂商和投资机构，提供专业、客观、可靠的第三方研究与咨询服务，助力决策者洞察数字化趋势，拥抱数字化机会，引领中国产业数字化升级。

- 数字化成熟度评估

基于研究、数据和案例调研积累，对比行业数字化基准水平，评估企业当前数字化成熟度，诊断数字化转型面临的困难与挑战，辅助制定业务与市场策略，实现业绩增长。

- 行业最佳实践

针对企业的特定业务场景，深入研究行业同类别公司及最佳实践案例，辅助业务与决策，优化业绩增长策略。

- 研讨与交流

参与我们的线上/线下研讨会，与专家学者、业内同行、数字化厂商，共同探讨行业数字化进程、技术应用趋势与最佳实践案例。

- 厂商遴选建议

基于您的需求，凭借对新兴技术和应用的系统研究、对供应商全面而充分的调研，秉承专业、客观、中立的原则，提供精准的供应商遴选建议。

地 址：北京市朝阳区兆维华灯大厦A1区1门2层2017室

联系邮箱：marketing@ifenxi.com



上海爱数信息技术股份有限公司成立于 2006 年，是领先的大数据基础设施提供商，致力于为政府、公共事业及企业的数字化转型赋能，帮助各行各业的客户在数字化浪潮中充分释放数据价值，实现即时、随时、实时的数据服务。经过多年的沉淀与积累，目前有 1300+ 名员工，1000+ 家合作伙伴，总部位于上海，在全国各地设有分支机构及办事处。在政府、企业、金融、教育、医疗等行业内已获得 20,000+ 家客户的认可。未来，爱数将继续保持以客户需求为核心，坚持技术创新，以领先的产品与方案、卓越的服务，加速全球企业数字化创新，释放无尽的数据潜力。

地 址：上海市联航路 1188 号浦江智谷 8 号楼 2 层 A 座

邮 编：201112

咨询热线：021-5422 2601

服务热线：400-880-1569

传 真：021-54222601-8800

客服邮箱：support@aishu.cn



Aug. 2020



ifenxi

专注产业数字化

