

架构师峰会北京站

胡时伟

第四范式 联合创始人



促进软件开发领域知识与创新的传播



关注InfoQ官方微信
及时获取ArchSummit
大会演讲视频信息



全球软件开发大会 [北京站]

2017年4月16-18日 北京·国家会议中心

咨询热线: 010-64738142



全球架构师峰会 2016 [深圳站]

2017年7月7-8日 深圳·华侨城洲际酒店

咨询热线: 010-89880682

我们是谁&您能收获什么？

- 第四范式：使每个人都能获得AI能力（ AI for Everyone ）
- 大数据下机器学习的一些特点
- 建设一个机器学习平台产品的实践分享

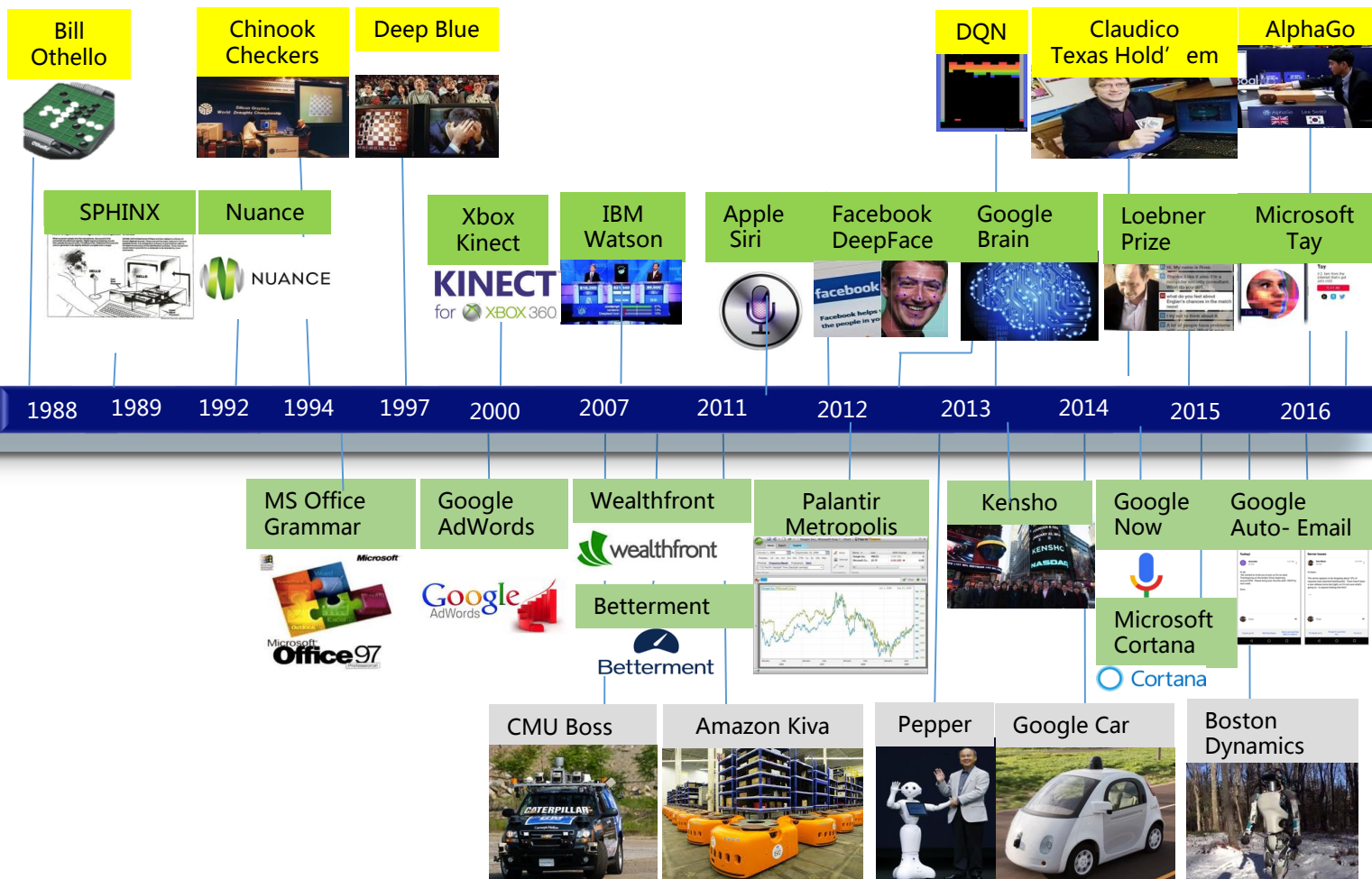
人工智能发展的主要里程碑

博弈

感知

决策

反馈



AI = 机器学习 + 大数据

存储和计算能力的发展：Intel / Nvidia / SSD / Infiniband

数据规模的变化：O2O / 物联网 / 互联网+

机器学习领域的发展：框架、人才、数据科学家

大纲

1. 机器学习产品
2. 算法与算法框架
3. 可扩展平台架构
4. 面向部署集成
5. 案例与选型

机器学习产品要解决什么问题？



业务专家：利用大数据和机器学习获得业务提升

关心：模型效果、与业务结合、可解释



数据科学家：处理数据 & 模型调研

关心：算法、灵活性、可扩展性、性能



系统管理人员：维护大量数据流 & 线上模型服务

关心：资源使用、一致性、可管理性

机器学习平台的困难？



VS



模型效果 VS 调研成本

- 大量数据导入导出&预处理
- 特征工程 & 调参



VS

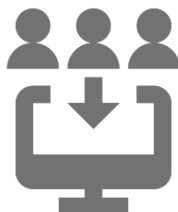


领域知识 VS 技能要求

- 问题定义和优化目标需要业务经验
- 需要懂Python / Spark / Tensorflow

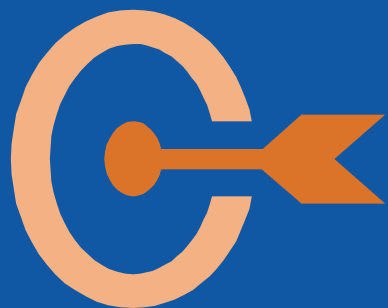


VS



投产要求 VS 运维难度

- 线上特征
- 实时预估服务



提升算法效果

聪明 VS 笨
天真无邪 VS 博览群书
一代宗师 VS 走火入魔

确保模型效果 – 充分使用尽可能多的数据

经验风险：模型对于训练数据分类结果的误差

置信风险：模型对于未知数据分类结果的误差

样本不足的情况下，VC维越高，越容易过拟合

样本充足的情况下，VC维越高，模型效果越好

-> 如何获得足够的样本数据：使用更多的表和字段，3维特征

-> 如何获得足够的计算能力：分布式机器学习

VC维 = 机器学习的**智商**

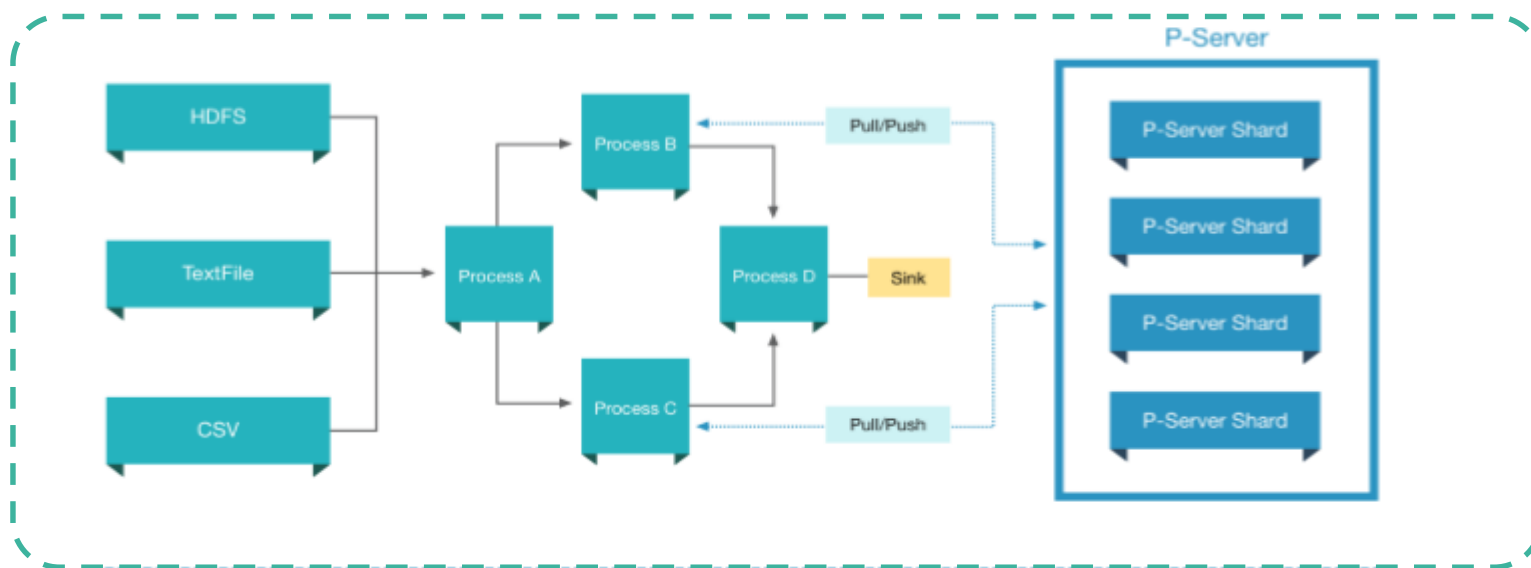
大规模机器学习框架GDBT

C++ 14 / 兼具运行效率和开发效率

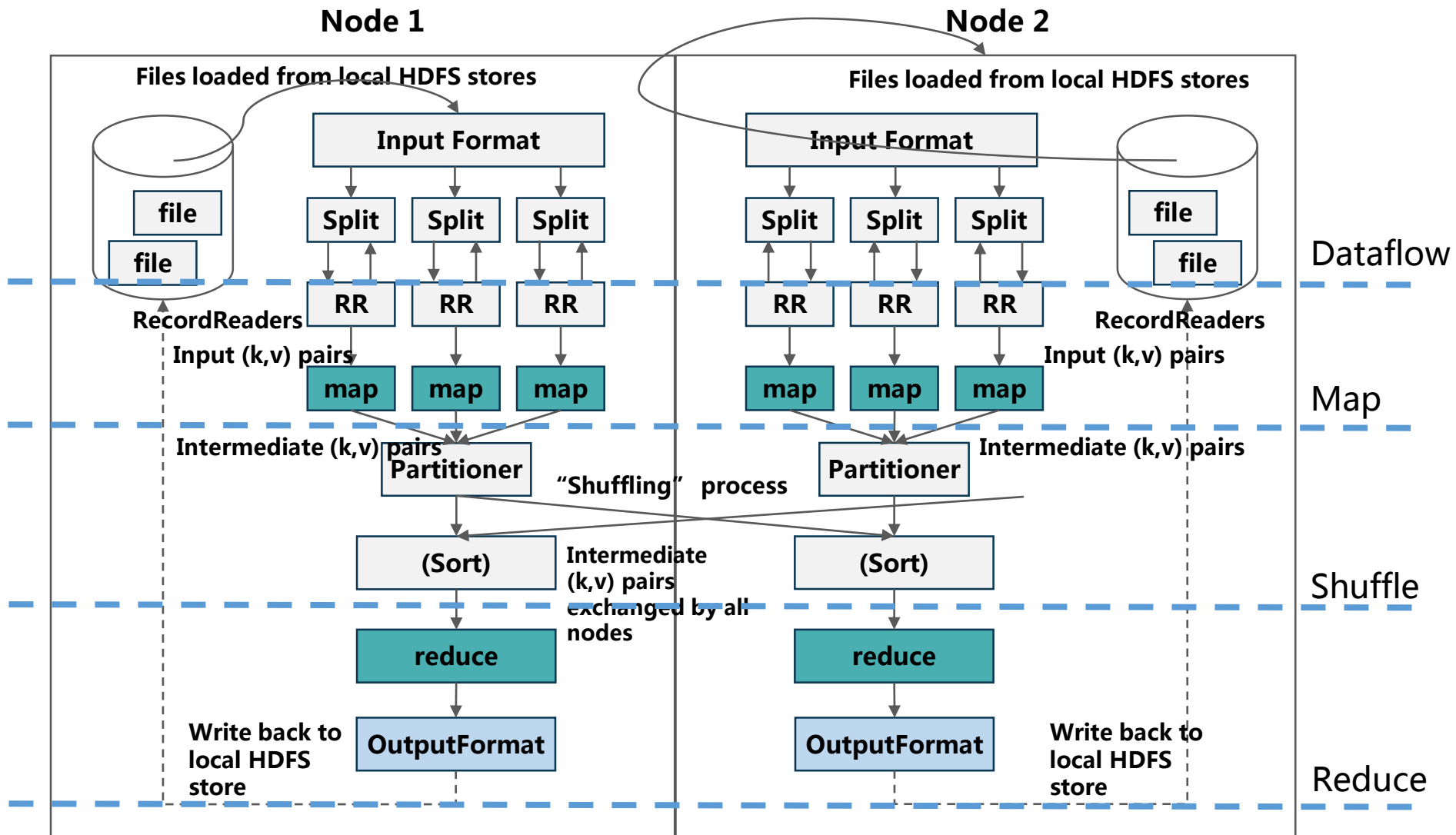
机器学习过程抽象，隐藏分布式细节

数据流与学习过程的紧密结合

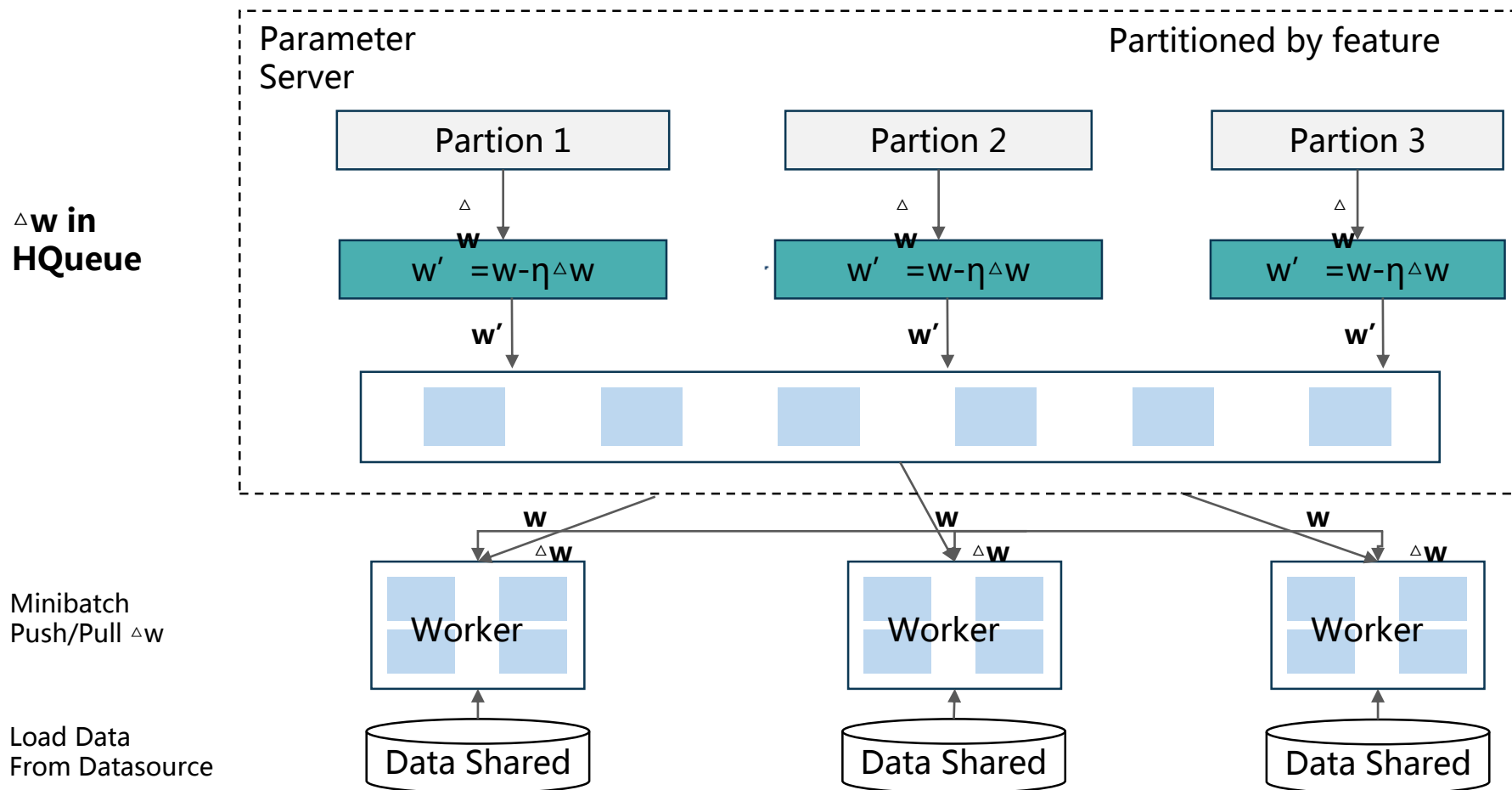
面向实际客户问题的算法包



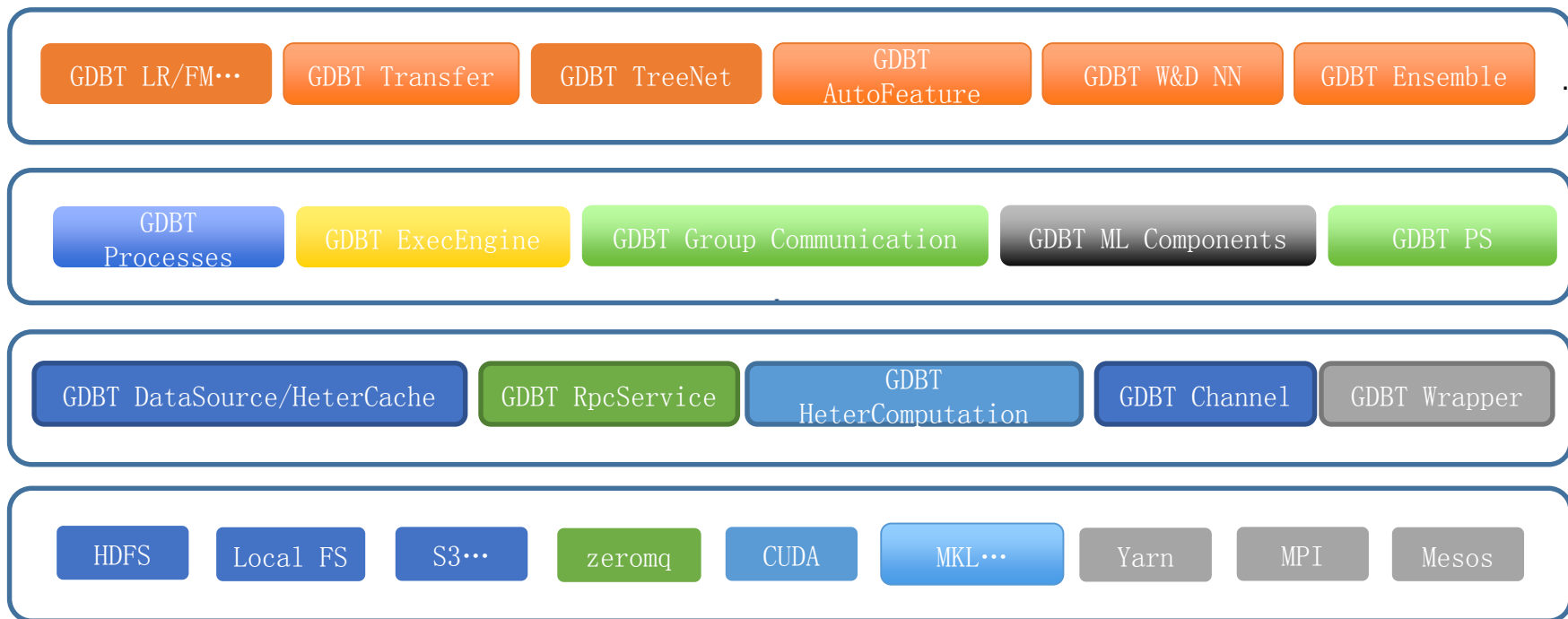
MR/Spark ML 计算模型



GDBT 计算模型

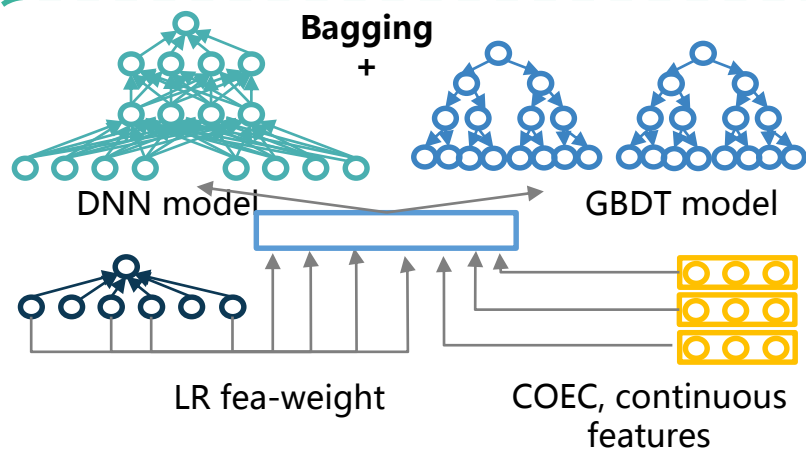


GDBT Not Only Parameter Server

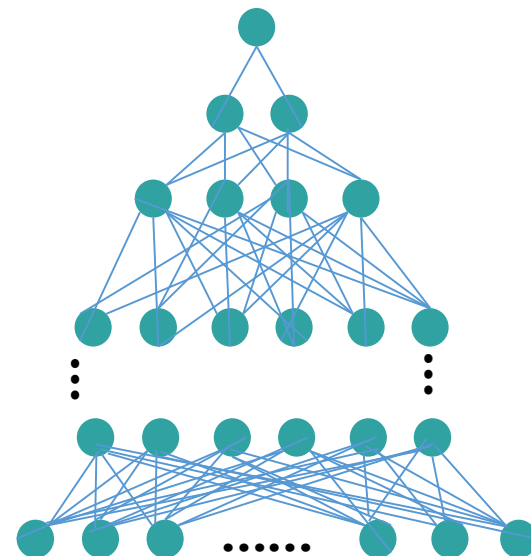


存储 | 计算 | 通讯 | 灾备 | 开放接口 | 场景优化

- 机器学习问题并非0和1问题
- 需要尽可能利用离散&连续特征
- 面向客户场景（模型稳定性）



Deep Sparse Network
(第四范式新一代深度学习模型，2015)



开发新算法只需要一百到几百行代码（LR、FM）

无须关心分布式细节，就可获得分布式算法

支持LossFunction/算法数据流的定制





降低成本与门槛

培养一个合格的AI人才需要**6-10**年的时间

--杨强 AAIL Fellow , 第四范式首席科学家

培养一个合格的AI人才可增加经济收益**500-1000万**美元

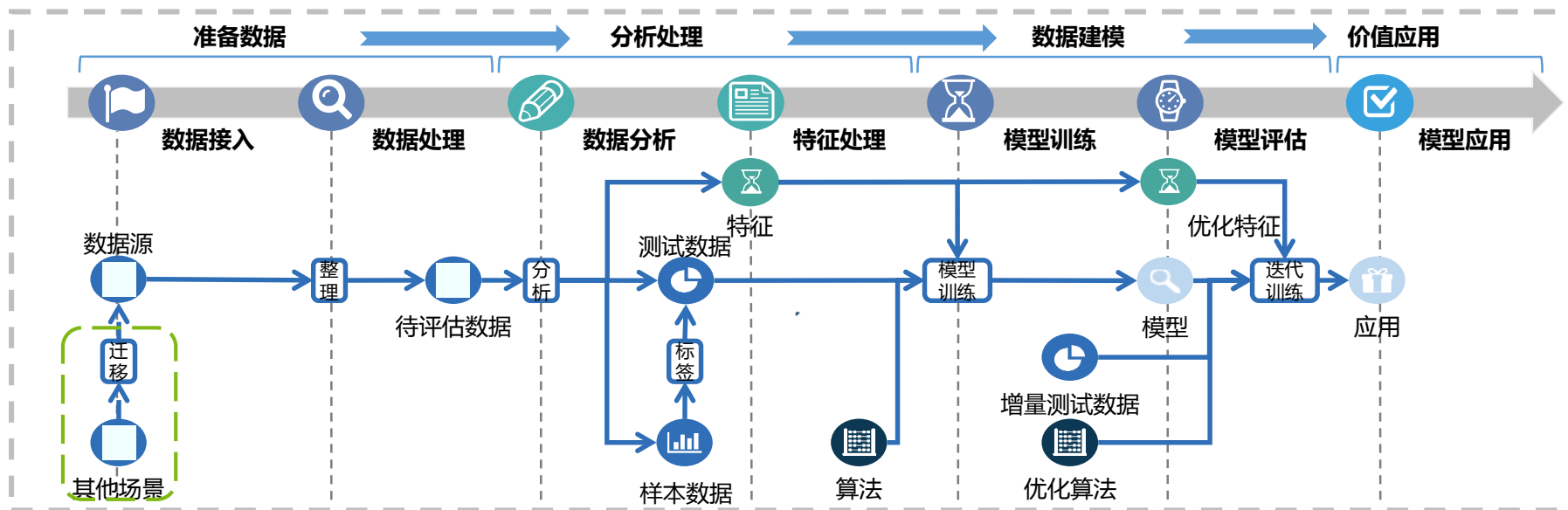
--Andrew Moore 卡耐基梅隆大学计算机学院院长在
白宫AI发展听证会上的讲话

- 业务专家
- 数据科学家
- **机器学习系统开发人才**

- BI
- SQL
- **Oracle DB/DB2**

- 算法创新：较少的需要参数手动调整
- 自动特征工程：利用DSN，同时获得千人千面 + 可推理的效果
- 高级特征工程算子：序列事件特征、社交关系特征
- 特征和模型可解释性：辅助建模人员更有效率工作
- Transfer Learning (IN PROGRESS)：如何打破全局意义上的数据分割

Prophet – 对模型的全生命周期管理



易使用
交互式的图形化界面能快速的完成业务问题转化和建模过程的定义

数据科学家/业务专家

高效率
提供多种系统化实验, 并提供自动的优化和调参功能

团队协作
为不同的团队角色的提供针对性的功能和与之对应的协作方式, 同时提供不同角色的培训服务

高效能
自主知识产权的专利算法和计算框架提供高效的计算能力和精准的应用效果

开发者/系统工程师

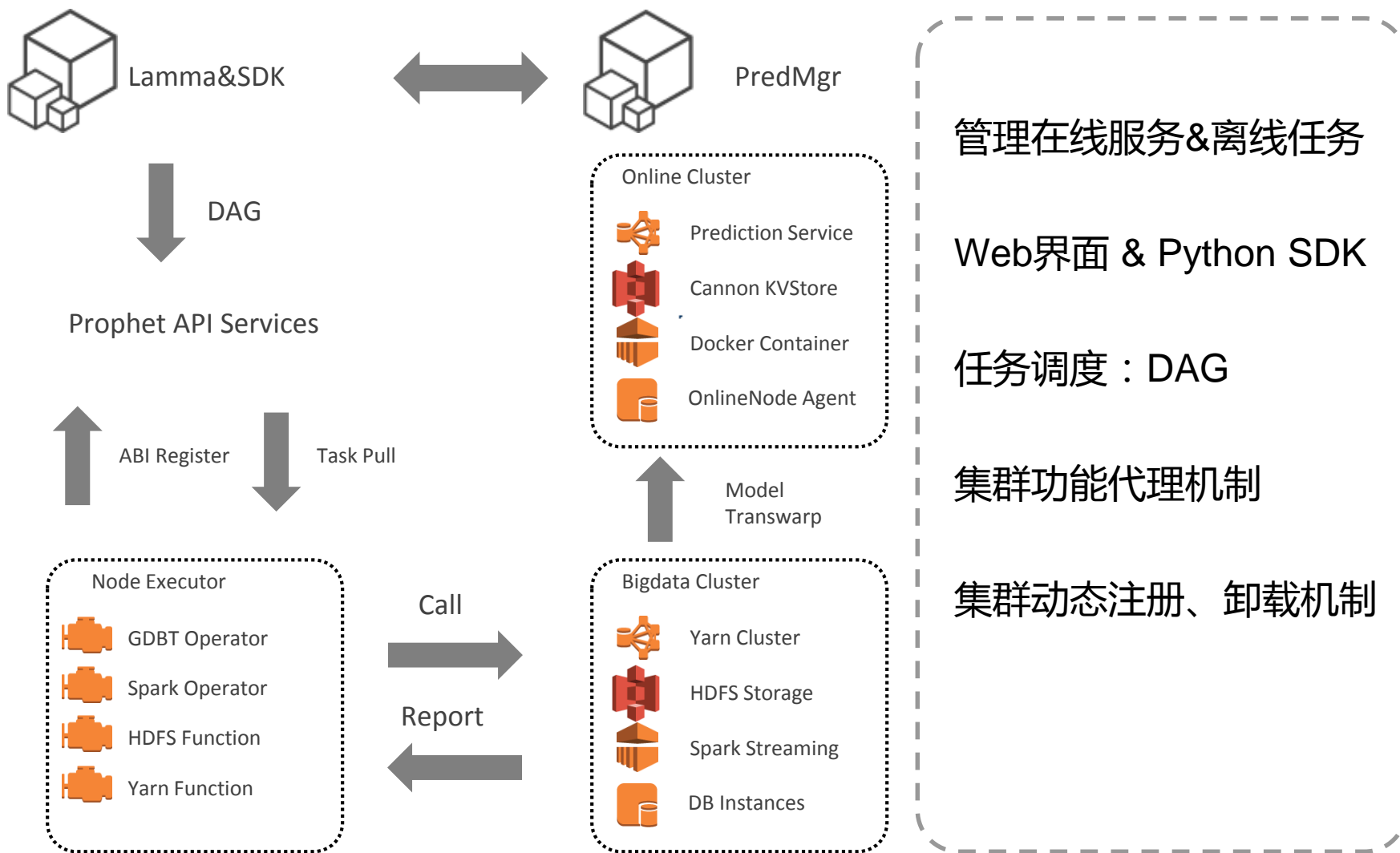
快速定制
作为通用开发平台, 开发者可快速依托平台的组件库和架构完成专属的人工智能业务系统的定制和对接

多功能多语言
支持Python、R、SQL等多功能语言和用户习惯的使用方式

高可扩展
提供多语言的SDK, 帮助开发者在此基础上完成二次开发和扩展使用

高处理能力
大规模分布式的底层架构, 满足高业务复杂度和数据量的存储和处理需求

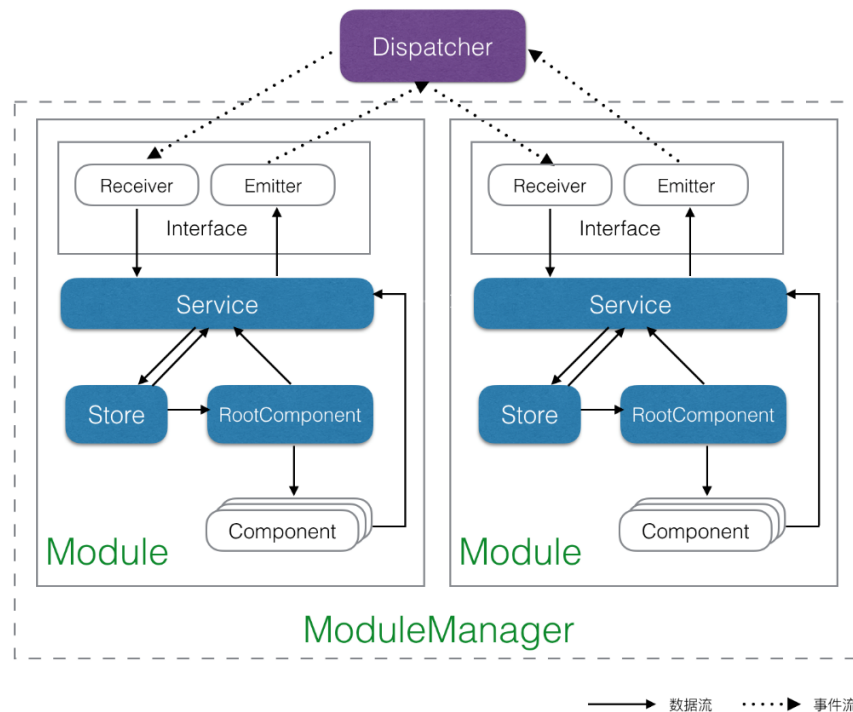
Prophet – 系统架构



4Paradigm
第四范式



- 使用场景以PC为主
- ReactJS (核心框架)
- Lamma-Flux (数据流框架)
- Lamma-Parts (组件框架)



- 界面组件模板化开发

```
{
  "taskType": "DataSplitAtom",
  "enableGroup": false,
  "nodeTemplates": [
    {
      "name": "DataSplitAtom", "label": "数据拆分", "tag": [ "DataSplit" ],
      "inputs": { "type": "data", "slots": [ { "type": "data" } ] },
      "outputs":
        { "type": "data", "slots": [ { "type": "data" }, { "type": "data" } ] },
      "config": { "basic":
        {
          "method": { "content": 0,
            "widget": {
              "name": "DropDown",
              "isVisible": true,
              "order": 1,
              "candidates": [
                { "label": "按比例拆分数据", "value": 0, "isDefault": true },
                { "label": "按规则拆分数据", "value": 1 },
                { "label": "先排序后拆分数据", "value": 2 } ],
              "label": "拆分方式",
              "isParent": true,
            }
          }
        }
      }
    }
  ]
}
```



数据拆分

拆分方式

按比例拆分数据

按规则拆分数据

先排序后拆分数据

随机拆分

随机种子

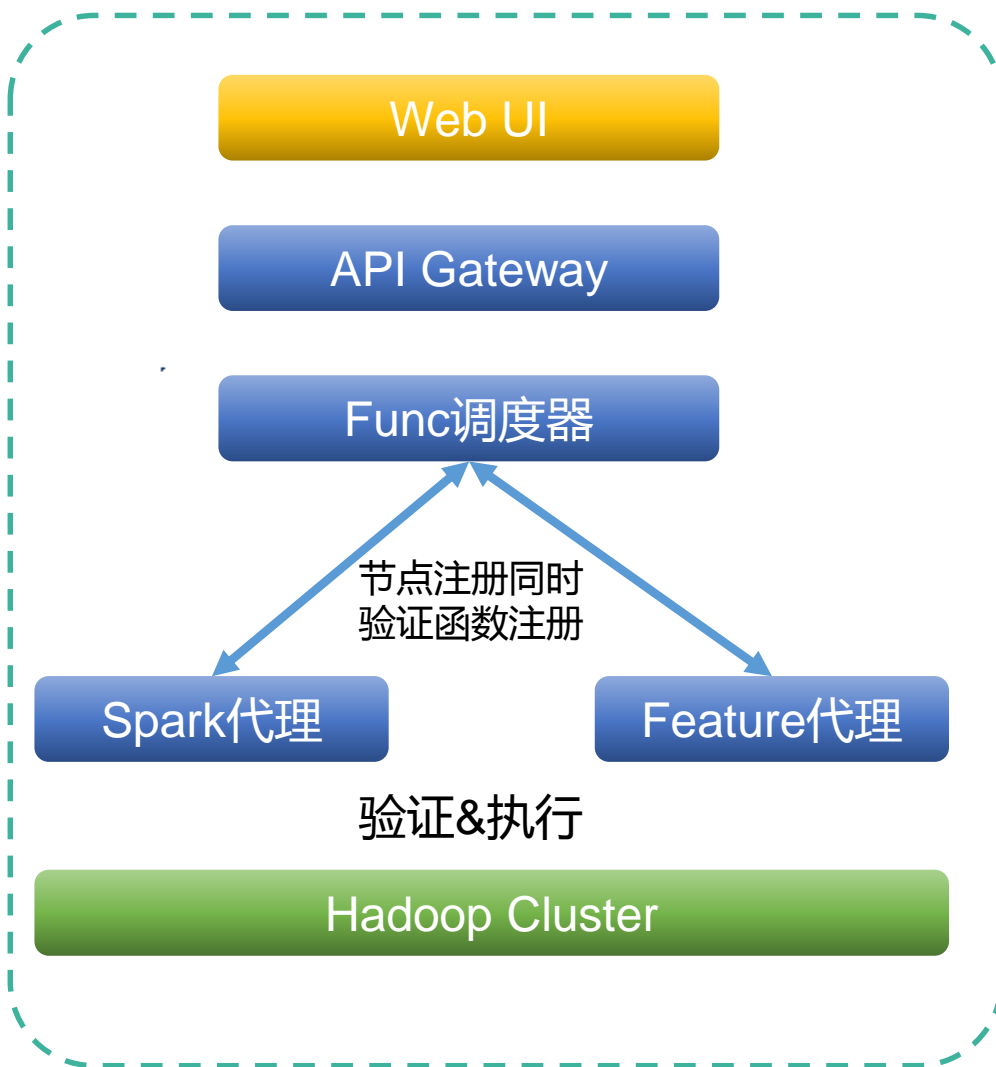
0

分层拆分

col_1

备注

- 服务器端语法推断和验证



SDK – 更快的调研或生产

- Web的优点：直观、可视化
- Web的缺点：操作复杂，不利于重复任务（例如For循环）

```
1 import prophet
2
3 client = prophet.from_env()
4 client.login('username', 'password')
5
6 ws = client.WorkSpace()
7
8 # Prepare training data
9 source_table = client.Table('sales_20161128')
10
11 # Split source data for training and prediction
12 t_train, t_predict = client.split(source_table, ratio=0.6)
13
14 # FeatureExtract
15 i_train = client.fe(t_train, script='fe_script_file')
16 i_predict = client.fe(t_predict, script='fe_script_file')
17
18 # Train
19 model = client.lr_train(i_train)
20
21 # Predict
22 t_predict = client.predict(model, label=True, columns=['col_1', 'col_2'])
23
24 # Model eval
25 report = client.eval(t_predict)
26
27 # Run jobs
28 ws.run(report)
29
30 print(report.auc)
31
```

SDK

Web

共用

Prophet API Service

Prophet Backend

Cluster Computing

Distributed Storage

工程团队和算法团队的粘合剂：缩短新技术产品化流程

模型调研过程更有效率，无人值守

提供前后端打通功能：训练过程可视化、进度和错误

可上线





面向部署集成

大客户IT三件事
安全、稳定、规范

企业产品运维三件事
标准、灵活、自动化

资源抢占

Troubleshooting麻烦

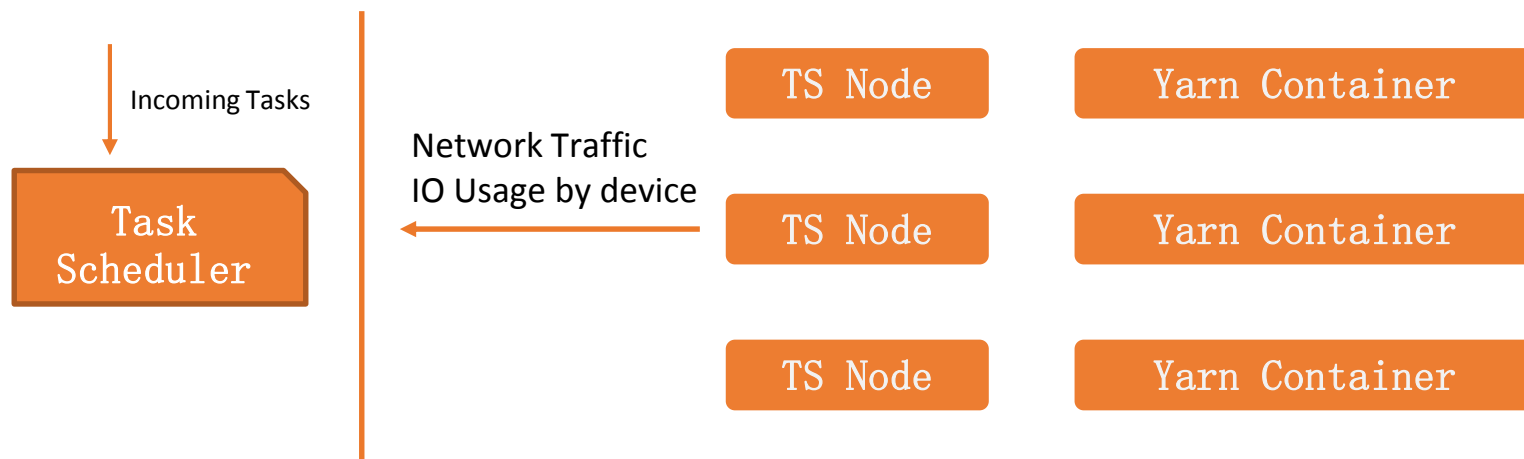
大数据集群兼容性

上线困难



Problem :

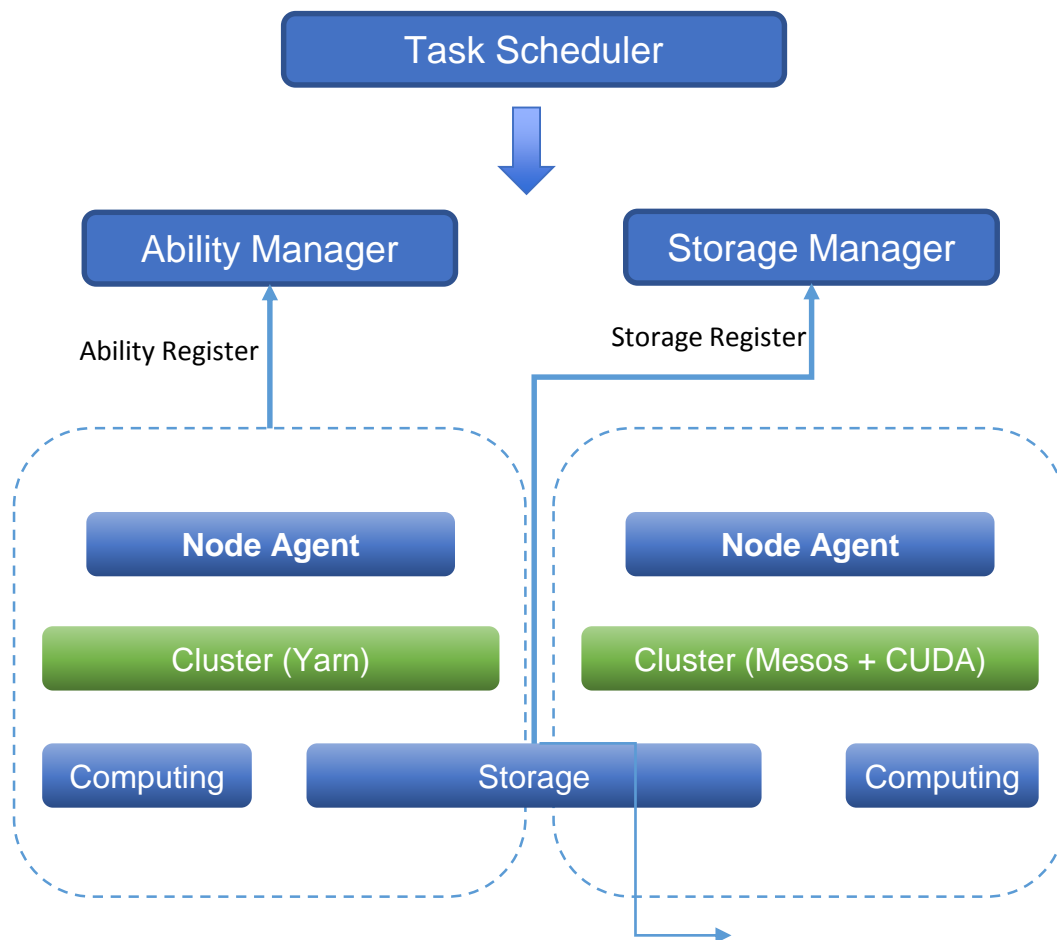
- 机器学习任务的灾备设计与ETL不同
- 局部独占是通常较优的调度策略
- 除了Yarn默认的vCPU/内存以外，网络带宽、IO也是重要考量因素



Dango – Yarn on Yarn

- 全功能调度

- 计算和存储分离可能
- 根据Ability调度任务
- 多集群灾备



Problem :

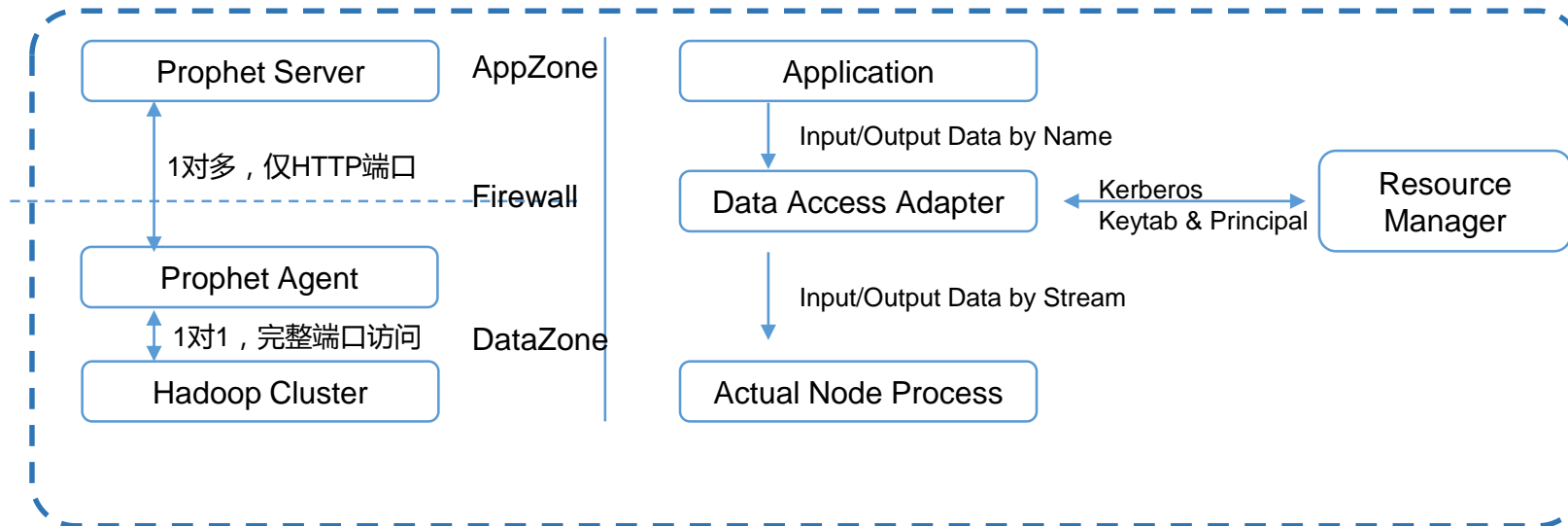
- Hadoop默认UI不友好（域名、端口、操作方式）
- 分布式任务的TroubleShooting需要经验
- 小错误导致的时间浪费（大型人物半途终止）

Solution :

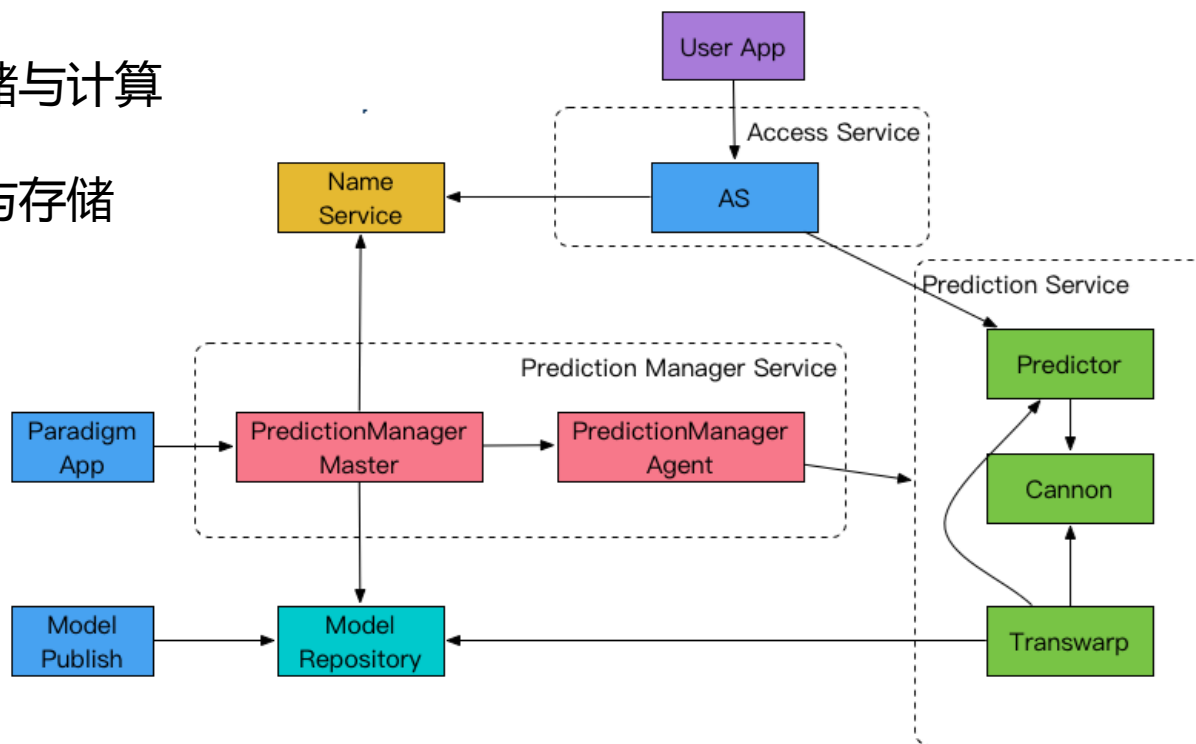
- LogStreaming / WebViewer / LogDownloader
- 对日志的关键条目进行分析并展示到UI
- 执行计划预先推断

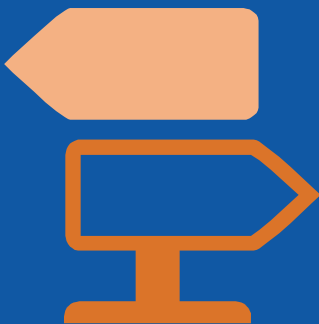
Problem :

- 企业通常已有商业版本的Hadoop集群，开启安全机制（Kerberos）
- 安全机制导致的数据服务器到应用服务器有限端口开放
- 多集群管理



- 线下DAG图到线上DAG图的自动转换
- 自动容器化部署与资源调度
- 分布式在线模型存储与计算
- 时间窗口特征计算与存储





经验&选型参考

哪些业务最适合开始机器学习实践？

传统金融（有历史数据、有业务干预点）：

- 推荐类：千人千面营销方案、产品组合推荐
- 定价类：因人而异的服务组合和定价策略
- 风险类：新户风险评分，贷后风险评分，欺诈识别

新兴互联网企业：

- 业务闭环可打通（购买行为、评价、反馈）
- 有一定的数据规模，更重要的是数据累积速度
- 基础设施建设（日志、物料库、效果分析系统）

- 学习目的 OR 生产目的？
- 是否具备足够的样本规模？
- 是否需要平台化管理？
- 是否需要线上实时应用？



高成本
高收益

低成本
低收益

高成本
低收益

THANKS

Email: kav#sjtuer.net



[北京站]

