

Pivotal®

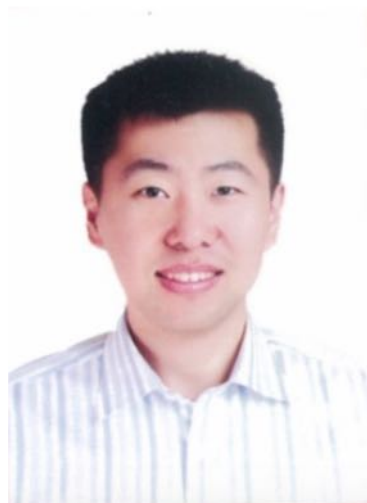
Greenplum V6及PostgreSQL DBaaS

AI/BI一体化的关键任务型HTAP MPP PostgreSQL

June 2019

韩鹏 Pivotal China

韩鹏简历



phan@pivotal.io



- 曾任 某证券交易所“新一代监察系统MPP项目” 项目经理, Pivotal资深架构师。
- 曾任 某证券交易所“新一代海量实时分析系统咨询项目”, 首席架构咨询师。
- 曾任 “中国首个银行业PaaS(招商银行Pivotal Cloud Foundry)项目” 项目经理及架构师。
- 曾在中航信(雇员4000+), 担任“大数据平台(GP+HADOOP)” 产品经理兼技术经理。
- 曾领导29人团队(任架构师及项目经理),完成中国首个从Teradata至Greenplum的企业数据仓库迁移。
- 专注于大数据及数据仓库领域**14**年, 熟悉各类分布式数据平台(MPP、NoSQL、Streaming), 兼具PaaS云平台实施经验。
- 具备**5000**余人天项目管理经验, 以及开发、运营大型信息系统经验。
- 2010年起, 持有工信部和人社部颁发的“**系统分析师**”证书。
- 吉林大学计算机科学与技术学院“计算机科学与技术”专业, 硕士。

纲要

- **Pivotal公司及产品介绍**
- **数据分析进化及数据库的演进**
- **Point-Of-Decision HTAP** (Hybrid OLTP and OLAP)
- **AI/BI一体化**
- **Greenplum on Kubernetes**
- **Postgres DBaaS**
- **QA**

Pivotal公司及产品

全球领先金融及产业公司多在Pivotal平台构建其应用和分析能力。全球开源贡献排名第三(GitHub中仅次于RedHat, IBM); 2016 Google全球唯一技术合作伙伴; 2017年, 福布斯“全球新一代爆发公司”排名第四。2018年4月, 纽交所上市公司, 总部在硅谷, 员工2500+, 近千企业用户, 年收约\$5亿, 年增 50%+, 市值约\$74亿。

Greenplum在国内外证券、银行、保险、互联网等行业拥有众多高端用户。

Greenplum社区与postgresql社区, 共同构成世界上最大的开源数据库社区。

国内Greenplum开源社区, 微信群人员1500+。



数据套件

- **Greenplum**: MPP DB
- **PostgreSQL**: OLTP
- **Gemfire**: In-Memory Data Grid



Pivotal Labs

- **Spring** (Java Framework)
- **数据科学家服务**
- **Spring Cloud**云原生应用框架咨询服务
- **敏捷**开发咨询服务



(混合)云原生应用平台

- **Cloud Foundry**: 云中立组合式 PaaS (拥有大部分财富500强客户)
- 2017 Google, Pivotal, VMware 合作推出基于K8S及BOSH的 **Pivotal Container Service(PKS)**
- **Pivotal Function Service: PFS** (Knative, Riff, etc)

15年如一日的持续开发，
铸造了**最佳版本**的大规模并行
分析PostgreSQL：
Greenplum



Greenplum在互联网中国: BAT全覆盖

阿里云 产品 解决方案 定价 ET大脑 数据智能 安全 云市场 支持与服务 合作伙伴

HybridDB for PostgreSQL

云数据库HybridDB for PostgreSQL (ApsaraDB HybridDB for PostgreSQL) 是一种在线MPP大规模并行处理数据库服务。云数据库HybridDB for PostgreSQL基于Greenplum Database开源数据库项目,并由阿里云深度扩展,支持OSS外部表、JSON数据类型、HyperLogLog预估分析等功能特性。通过符合SQL2008标准查询语法及OLAP分析聚合函数,提供灵活的混合分析能力。支持行存储和列存储混合模式,提高分析性能,同时支持数据压缩技术,降低存储成本。并提供在线扩容、性能监测等服务,用户无需进行复杂的大规模MPP集群运维管理,让DBA、开发人员及数据分析师专注于如何通过SQL提高企业的生产力,创造核心价值。

[立即购买](#) [产品价格](#) [进入控制台](#) [帮助文档](#)

腾讯云 产品 解决方案 定价 文档 支持 合作与生态 客户

文档平台

分析型数据库 Greenplum

文档平台 > 分析型数据库 Greenplum

使用客户端工具

最近更新时间: 2017-11-15 17:22:11

使用客户端工具

腾讯云数据库Greenplum是由Greenplum Database开源数据库项目发展而来,因此其接口协议完全兼容社区版的Greenplum Database6.3版本的PostgreSQL,可以使用Greenplum或PostgreSQL的客户端进行连接云数据库Greenplum,为更方便的使用,我们下面提供了命令行工具psql的下载和编译方式,以及其他客户端的使用。

- 免费教程
- 入门中心
- 云+社区
- 腾讯大学
- 手机管理云资源
- 联系我们

百度云 最新活动 产品 解决方案 云市场 合作与生态 帮助与支持

分析型数据库 FusionDB

百度云FusionDB

百度云FusionDB是一种在线分布式云数据库,提供了简单、快速、高性能、高可靠的大规模并行处理(MPP)数据库服务。百度云FusionDB基于开源Greenplum Database项目开发,提供了完善的监控、安全管理等功能,同时可以轻松实现在线扩展。百度云FusionDB具有良好的兼容性,可以轻松使用PostgreSQL相关生态工具,对海量数据进行分析处理,在使用方便、高效、可靠的同时,让您更专注于业务处理。

[产品特性](#) [应用场景](#) [产品架构](#) [实例规格](#)

纲要

- Pivotal公司及产品简介
- **数据分析进化及数据库的演进**
- Point-Of-Decision HTAP (Hybrid OLTP and OLAP)
- AI/BI一体化
- Greenplum on Kubernetes
- Postgres DBaaS
- Pivotal敏捷企业数据架构
- QA

新业务激发新需求

新用例 驱动新处理选择

- 应用重构
- 数据科学
- 新数据类型及数据源
- NoSQL数据处理系统的爆发

数据处理拓展 创造新竖井

- 太多具备冗余能力的系统
- 利用率、优化率低下
- 移动数据，或仅分析数据某子集
- 昨日数据英雄，今成另一数据竖井

交易型数据库联合负载 处理机会在缩减

- 交易型数据库架构有面向不同用途的优化
- 但可在单一分析型数据库，凭借高速分析型读写能力模糊此边界
- 有助于运营型和分析型负载的联合处理

分析型数据库需继续演进

- 希望使用**新数据类型和数据源**，也期望**系统更少以减少数据移动**。
- 希望**保留成熟优化器**带来的**效率和快速查询**益处。
- 与应用开发者类似，DBA们也希望得到的待遇，享受自动化和调度系统的便利：花费**更少精力控制更多DB实例**。

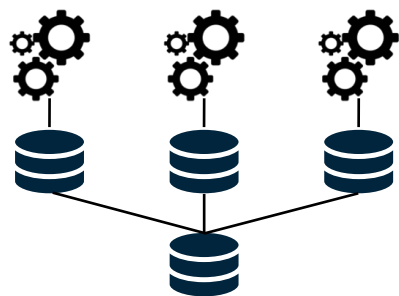
**鉴此，我们已开发Greenplum 6*：
用于企业大规模分析的并行Postgres**

具备已验证的事务处理和流数据加载能力，Greenplum可用于**交易和分析混合场景**，涵盖**传统BI到深度学习**。

同时也是首个获得**Kubernetes自动化优点**的Greenplum版本
更易安装、运维及升级

*beta available, GA expected 6/30/19

优异的Greenplum V6



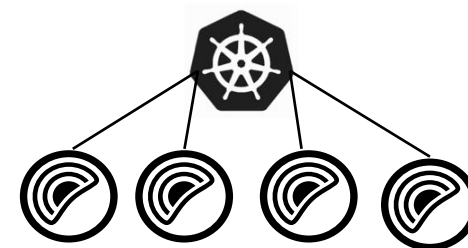
单系统联合处理 更多类型负载

点查询，快速加载，及报表长查询，数据科学探索，更大扩展及**并发**能力



更简 workflow 更强数据科学

使用Apache MADlib实现从数据科学实验到部署，**GPU和深度学习**支持即将到来，兼具自动化友好特性



Kubernetes平台上的 自动化，可重放部署

更简易地部署和管理数百实例。利用**Greenplum on PKS**，可连通PCF应用。

纲要

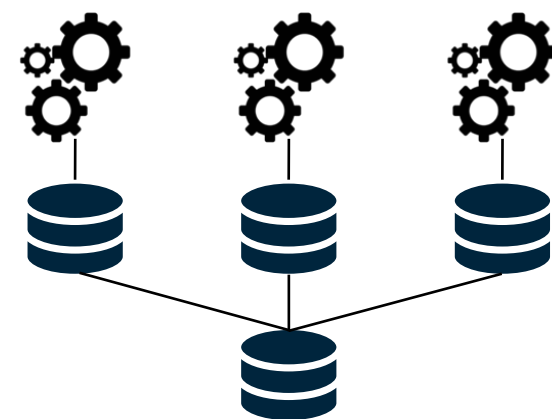
- Pivotal公司及产品简介
- 数据分析进化及数据库的演进
- **Point-Of-Decision HTAP** (Hybrid OLTP and OLAP)
- AI/BI一体化
- Greenplum on Kubernetes
- Postgres DBaaS
- Pivotal敏捷企业数据架构
- QA

若数据库可用于 从交易型到分析型的广泛类型负载， 您的应用体验将会怎样？

领先的MPP数据库，解决大规模BI分析的分布式
postgres。进一步改进对**混合型**负载的支持(特定类

型负载能力，比Greenplum V5提升高达**50**倍)

– GP V6 Beta版已就绪，预计2019年6月正式版发布。



Pivotal Greenplum Database

Point-Of-Decision HTAP的GPV6新技术支撑

更快交易型 负载能力

- 全局死锁检测器甄别等待资源进程
- 单表高并发更新和删除的行级锁
- 特定场景负载能力比GP5提升高达

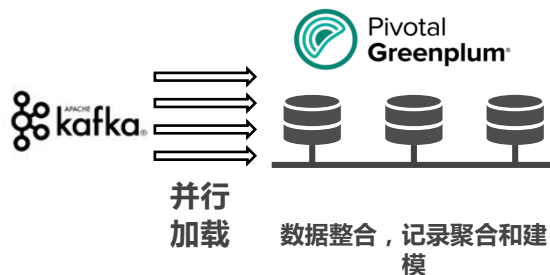
50倍

- 提供单库(ODS)运行交易型和分析型负载能力，不会因锁竞争带来性能降级
- 复制表：维度表复制至本地segments，以便更快地与事实表本地关联，避免跨Segments



Kafka Connector

- 金融级的不重不丢(Exactly-Once)
- 持续的非阻塞式数据加载
- 断点续传(Exactly-Once保障之一)
- 必要时在提交时点可自动发布SQL和UDF
- 金融业和IoT应用长期广泛
- 通过Kafka商业公司Confluent认证



其它新特性

- 合并至PostgreSQL 9.4(1万+ Code Commit)
- 相比libz和quicklz, zStd压缩/解压更快，压缩率更高
- 基于WAL日志复制的集群内部DB镜像，降低网络开销
- 数据集群规模变化时，Jump Consistent Hash算法帮助更高效重分布
- Greenplum Command Center实时执行计划的可视化，有助于一眼识别问题查询
- Parquet格式更快并行读写AWS S3数据对象
- 带GIN, Btree, Hash索引的JSONB
- 列级权限
- 磁盘配额

纲要

- Pivotal公司及产品简介
- 数据分析进化及数据库的演进
- Point-Of-Decision HTAP (Hybrid OLTP and OLAP)
- **AI/BI一体化**
- Greenplum on Kubernetes
- Postgres DBaaS
- Pivotal敏捷企业数据架构
- QA

若可简便地先并行训练模型，后以同种语言发布至生产环境，将给您的AI开发及运维带来哪些改变？

适用于Pivotal Greenplum和Postgres的，易懂的，基于SQL的统计和机器学习，开源函数库，**现已支持GPU**。
商业版包含自动模型部署特性。



Apache MADlib

美国银行业综合欺诈每年损失**\$170亿**, 构

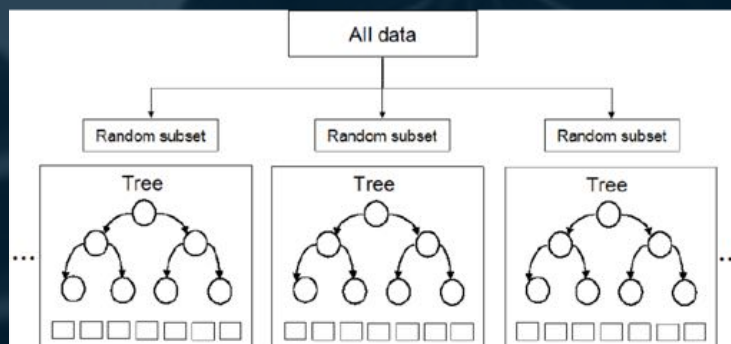
成 银行 未保全坏账 的**10-15%** Source: FICO

应对综合欺诈的典型监督学习算法

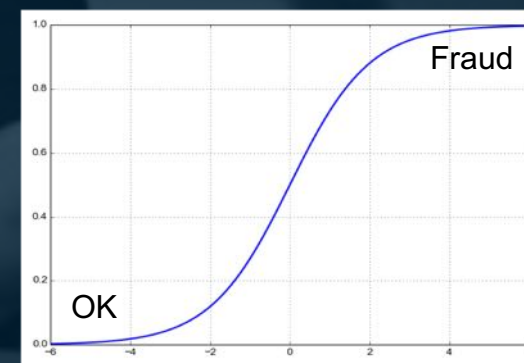
决策树
(transparency)



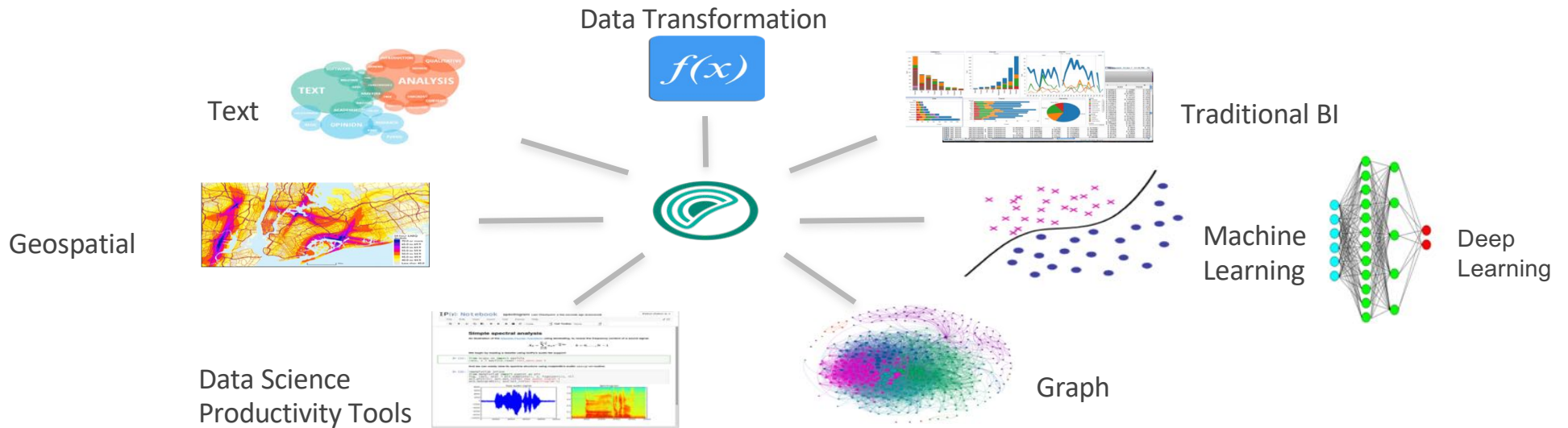
随机森林
(装配多棵树就预测结果“投票”)



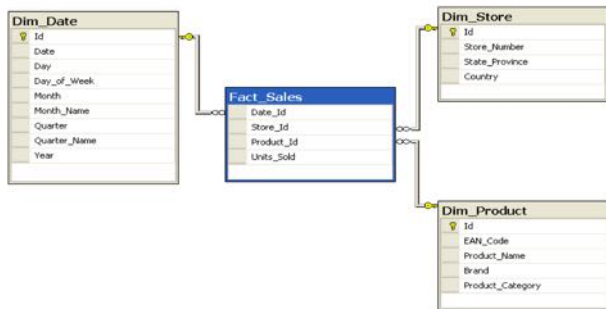
线性回归
(binary prediction)



BI/AI 一体化分析平台



结构化



半结构化



非结构化(<=1GB/件)





日用精品算法: 机器学习, 图分析, 统计等

监督学习

Neural Networks
Support Vector Machines (SVM)
Conditional Random Field (CRF)
Regression Models

- Clustered Variance
- Cox-Proportional Hazards Regression
- Elastic Net Regularization
- Generalized Linear Models
- Linear Regression
- Logistic Regression
- Marginal Effects
- Multinomial Regression
- Naïve Bayes
- Ordinal Regression
- Robust Variance

Tree Methods

- Decision Tree
- Random Forest

非监督学习

Association Rules (Apriori)
Clustering (k-Means)
Principal Component Analysis (PCA)
Topic Modelling (Latent Dirichlet Allocation)

最近邻居

- k-Nearest Neighbors

图分析

All Pairs Shortest Path (APSP)
Breadth-First Search
Hyperlink-Induced Topic Search (HITS)
Average Path Length
Closeness Centrality
Graph Diameter
In-Out Degree
PageRank and Personalized PageRank
Single Source Shortest Path (SSSP)
Weakly Connected Components

工具函数

Columns to Vector
Conjugate Gradient
Linear Solvers

- Dense Linear Systems
- Sparse Linear Systems

Mini-Batching
PMML Export
Term Frequency for Text
Vector to Columns

抽样

Balanced/ Random/ Stratified Sampling

时序分析

- ARIMA

数据类型和转换

Array and Matrix Operations
Matrix Factorization

- Low Rank
- Singular Value Decomposition (SVD)

Norms and Distance Functions
Sparse Vectors
Encoding Categorical Variables
Path Functions
Pivot
Sessionize
Stemming

统计

Descriptive Statistics

- Cardinality Estimators
- Correlation and Covariance
- Summary

Inferential Statistics

- Hypothesis Tests

Probability Functions

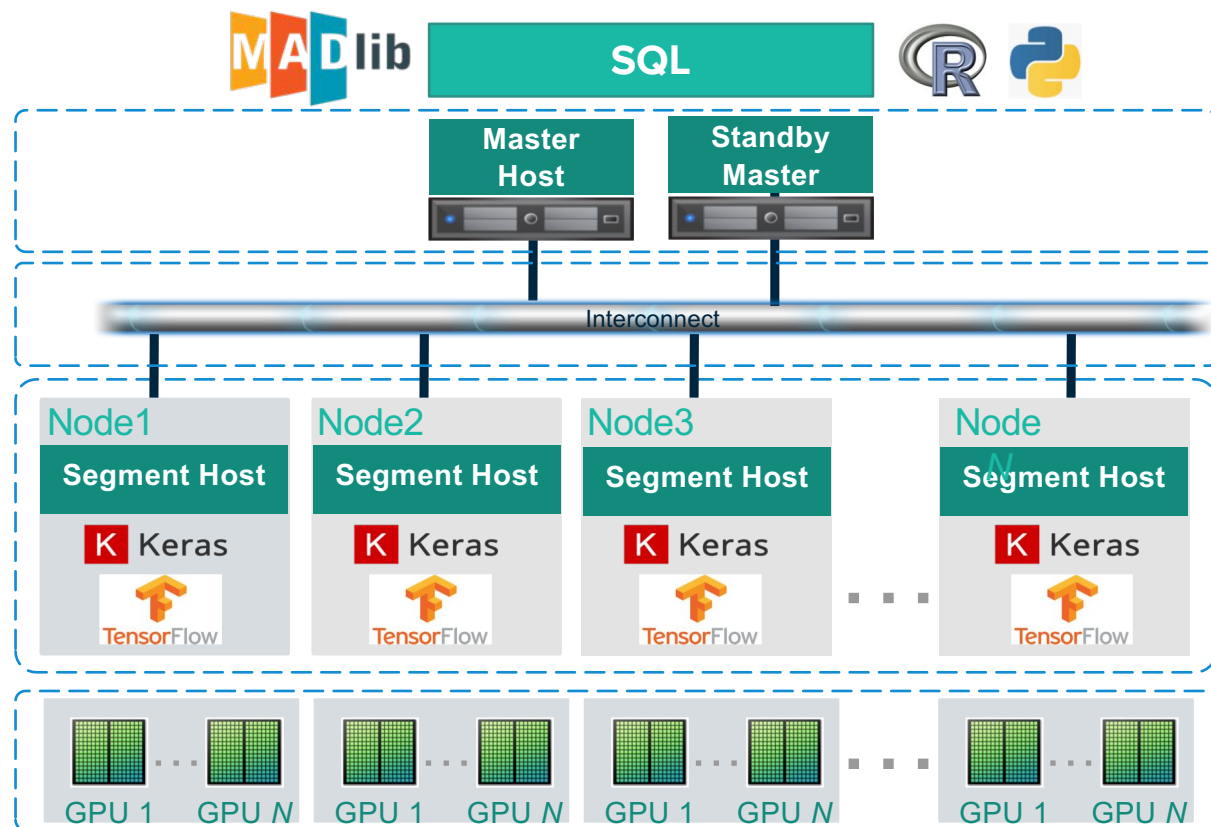
模型选择

Cross Validation
Prediction Metrics
Train-Test Split

GPV6: MPP融合用于机器学习的GPU算力

MADlib新能力

- 用于图像分类的卷积神经网络(CNN)
- 支持以Tensorflow为后台的Keras
- GPU算力加速模型训练和推断
- 模型管理 Model management
 - 简化多版本和多模型管理
 - 比较不同模型性能，择优选用



过程语言表述
(R,Python)转译为
MADlib中的
SQL或UDFs

以Tensorflow为
后台的Keras运行
在每个segment

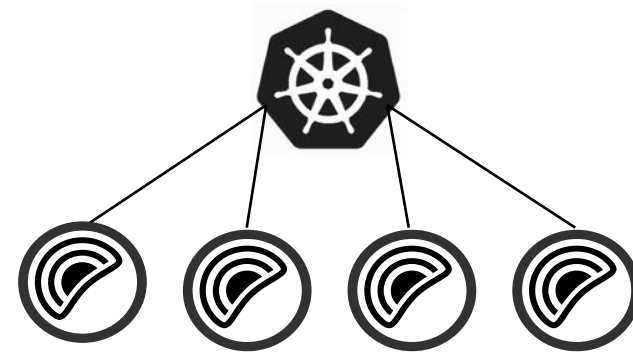
分布计算至每个
segment的多个
GPU

纲要

- Pivotal公司及产品简介
- 数据分析进化及数据库的演进
- Point-Of-Decision HTAP (Hybrid OLTP and OLAP)
- AI/BI一体化
- **Greenplum on Kubernetes**
- Postgres DBaaS
- Pivotal敏捷企业数据架构
- QA

若可简便调度和管理数据库实例，
则DBA的工作体验将会怎样？

利用先进云平台**自动化特性**管理大量实例，提供更大弹性，及减少重复性工作



Pivotal Greenplum for Kubernetes

Greenplum
及补充工具



Docker 镜像



CI 管线



- Pivotal认证过的
- 配置
 - 依赖项
 - 安全
 - 网络
 - 集成
 - ...



菜单

Greenplum组件:

- 机器学习
- 空间地理
- 图分析
- 文本分析
- Kafka Connector
- Greenplum Workshop

过程语言:

- Python
- R
- Java

数据科学Notebooks:

- Jupyter
- Zeppelin

外部服务:

- LDAP: ...
- 备份 / 恢复: ...

部署拓扑:

- 规模 & Failover选择 :...
- 存储选择:...



Greenplum
Operator



Greenplum
实例

用户可指定存储类:

- Kubernetes 提供灵活性
- 存储类在增长:
 - 偏好性能的本地化
 - 偏好灵活的远程
 - 其它特性的方案, 如动态增长
- 用户可按需择优使用



Azure Disk



AWSElasticBlockStore



(formerly ScaleIO)



Pivotal

Greenplum for Kubernetes 优益



速(Speed)

数分钟部署
可精确重放
自服务



省(Savings)

自动任务
可运行于任意K8S
任意基础设施



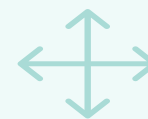
安(Security)

预加固
预置网络
安全Docker镜像



稳(Stability)

自动恢复
更快升级, 打补丁
构建CI / CD管线



弹(Scalability)

弹性容量
分别扩展计算和存储

纲要

- Pivotal公司及产品简介
- 数据分析进化及数据库的演进
- Point-Of-Decision HTAP (Hybrid OLTP and OLAP)
- AI/BI一体化
- Greenplum on Kubernetes
- **Postgres DBaaS**
- Pivotal敏捷企业数据架构
- QA

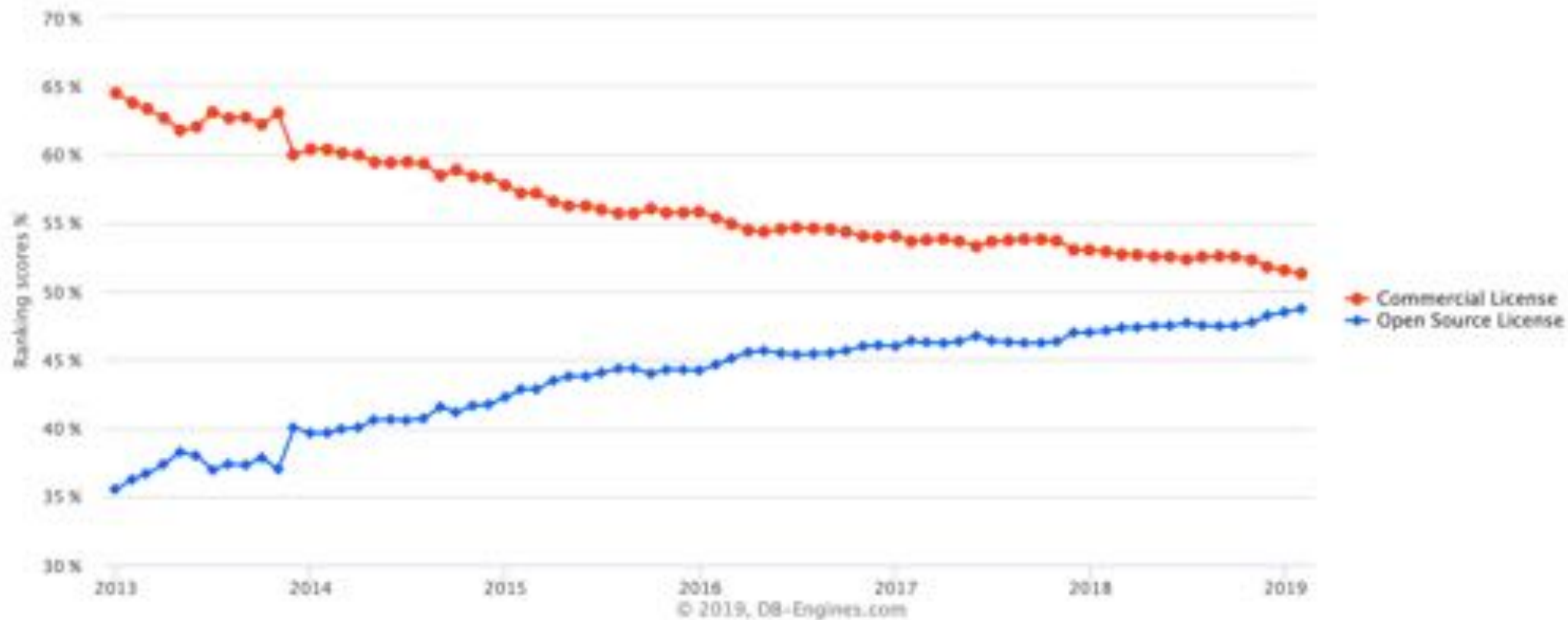
Pivotal®

Pivotal Postgres DBaaS

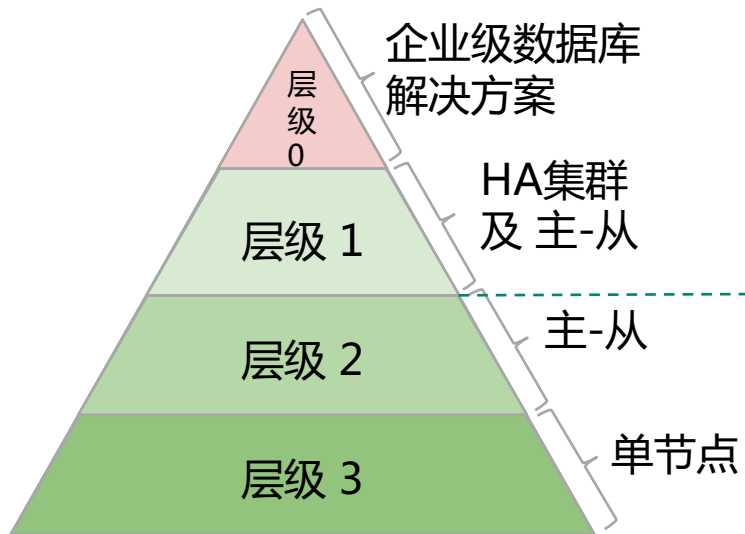


数据库市场的持续开源化变迁

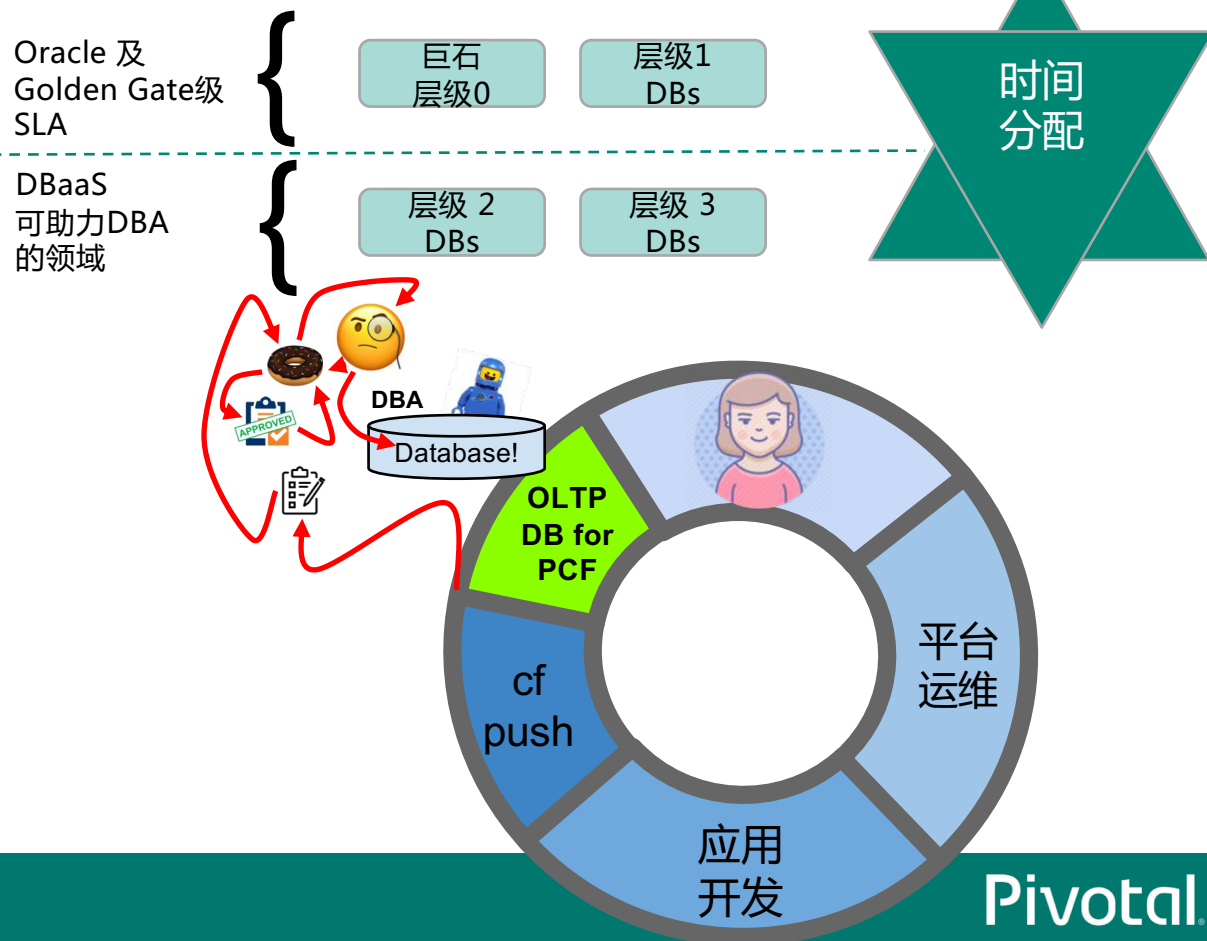
Popularity trend



应用层级



理想运维人员时间分配





Pivotal®

Greenplum/PostgreSQL DBaaS
伴你升腾新高度
