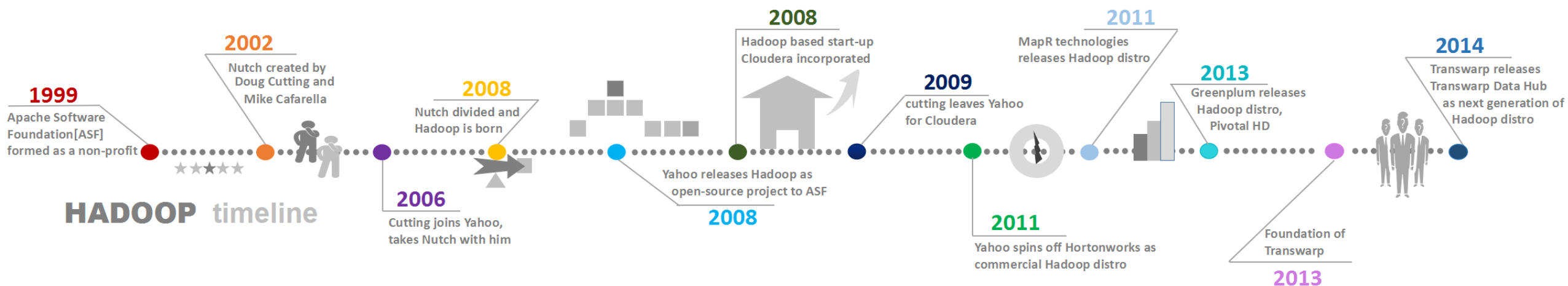


大数据技术的发展现状和最新趋势

孙元浩
星环科技
Founder & CTO
transwarp.io

Hadoop的发展历程回顾

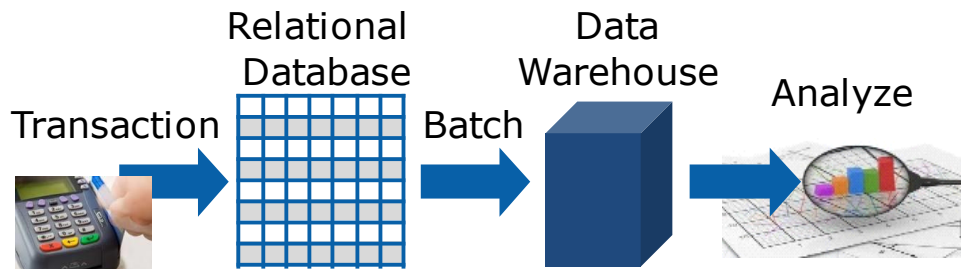


大数据技术的软件栈



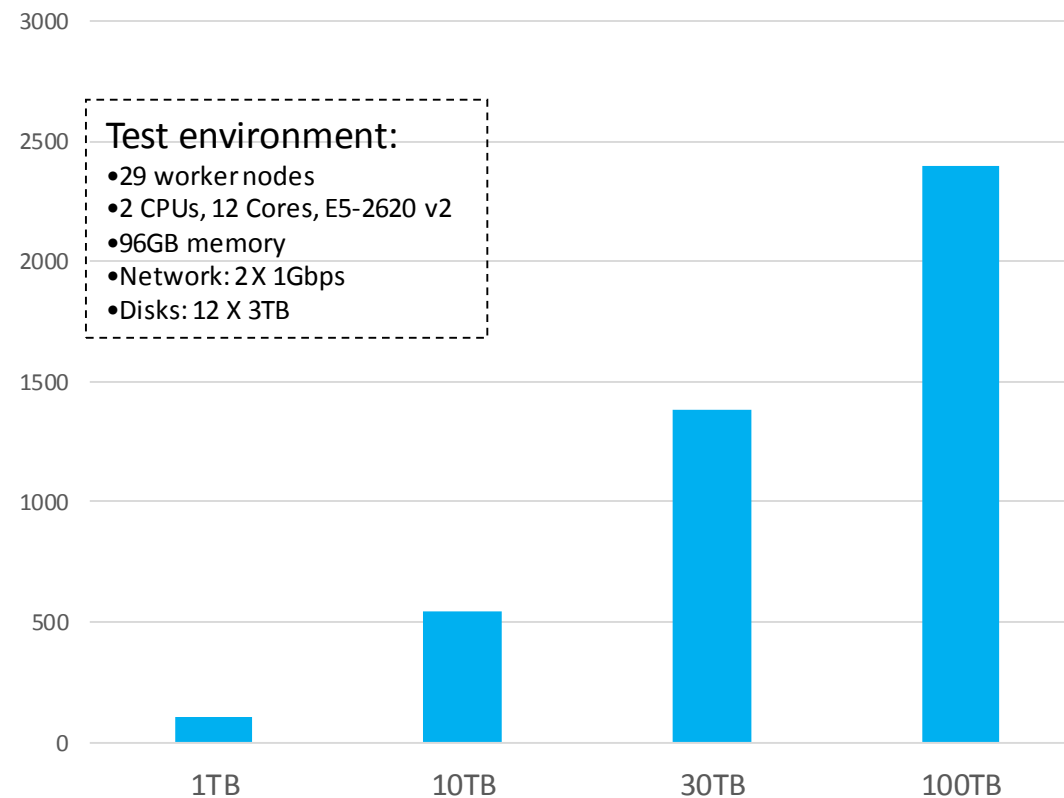
分布式计算已逐渐成为主流计算方式

Traditional Data Analysis

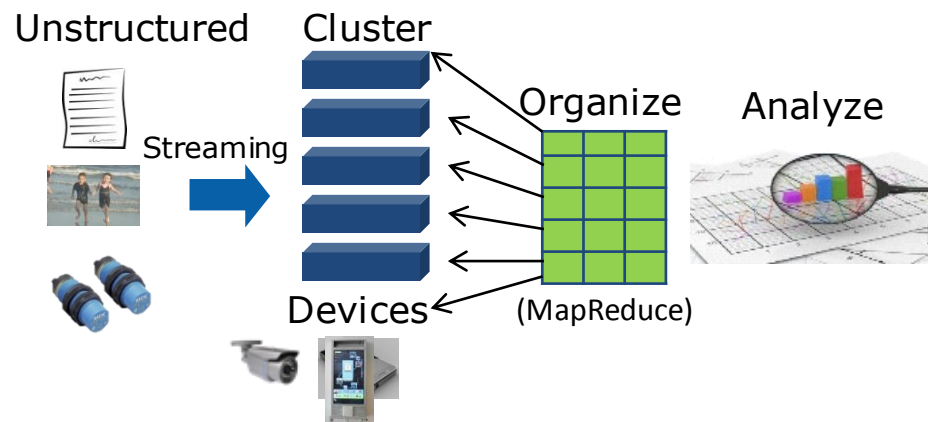


星环SQL on Hadoop已经能够高效处理100TB数据的复杂分析

Transwarp Inceptor's Performance
TPC-DS Execution Time for 99 Queries (in minutes)



Big Data Analysis



星环大数据集群已经可以在生产环境中处理20PB的数据

数据分析算法逐渐丰富，工具普及化

- R / Python语言开发 => 算法工程师，数据科学家
- 交互式挖掘 => 业务分析师，数据科学家

数据预览

预处理

特征工程

模型训练

模型上生产

- R和Midas中可以连接TDH中数据表做数据预览
- 可以对列做tag/feature的管理

Name	Type	Meaning	Statistics	Count
id	Polynomial	0	5 (1)	0 (1)
prediction	Integer	0	0	0 (0)
feature_cols_1	Real	0	0	9,200
feature_cols_2	Real	0	0	9,200



- 通过内置的分布式统计算法完成相关的预处理与数据分析
- 支持标准化，归一化，正则化，缺失值填充，数据分箱等
- 支持通过Inceptor SQL进行数据ETL处理



- 结合业务领域专家知识，以及相关算法降维，选择特征指标与维度
- 利用深度学习神经网络算法，通过升维降低特征工程维度选取难度

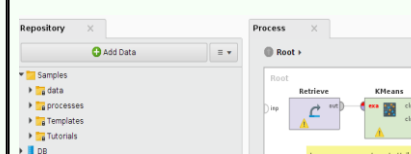
Dimension Reduction
主成分分析

Linear Regression
线性回归

Deep Learning
深度学习

.....

- 用户通过GUI选择算法开发训练模型
- 模型编译成为DAG，由Hubble组件来调度任务
- 支持单机R算法和分布式算法训练模型

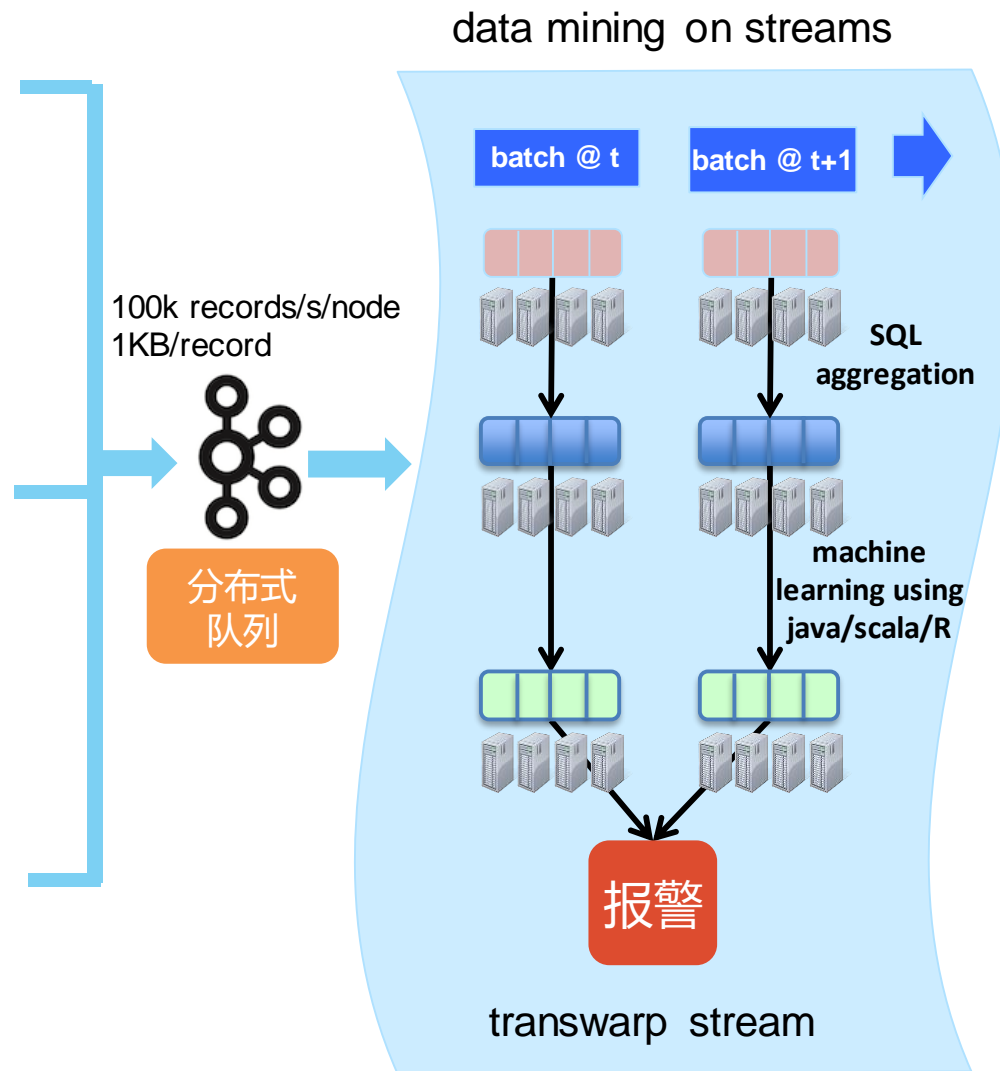


- 导出模型为PMML
- PMML模型可以转换成生产系统的代码，部署到实际业务中

```
<model key="decision_tree">
  <parameter key="featuresCol" value="features"/>
  <parameter key="labelCol" value="indexedLabel"/>
  <parameter key="predictionCol" value="prediction"/>
  <parameter key="probabilityCol" value="probability"/>
  <parameter key="rawPredictionCol" value="rawPrediction"/>
  <tree>
    <name>decision_tree</name>
  </tree>
  <root>
    <name>Internal</name>
    <count>3.0</count>
    <count>3.0</count>
  </root>
</model>
```



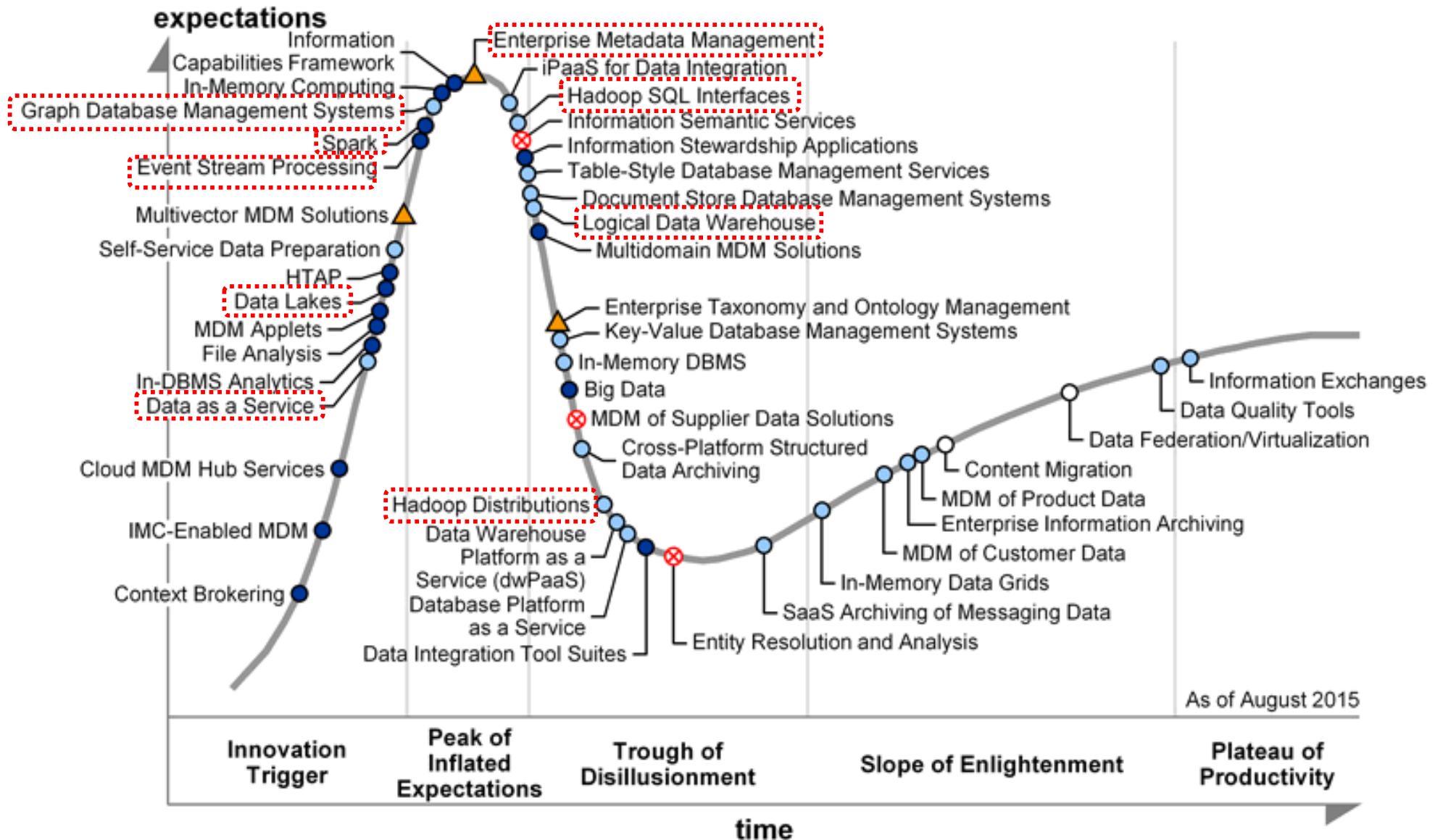
实时流处理技术推动IoT数据分析能力



1. Streaming processing and batch processing are unified in one programming model
2. SQL and its extension is the unified declarative language for device monitoring and diagnostics.
3. ANSI SQL 2003 and PL/SQL are supported on streaming events.
4. Linear Algebra
5. Machine learning

Usage cases in IoT & FS:
Real-time event monitoring
Real-time dashboard & statistics
Real-time outlier detection
Real-time fraud detection

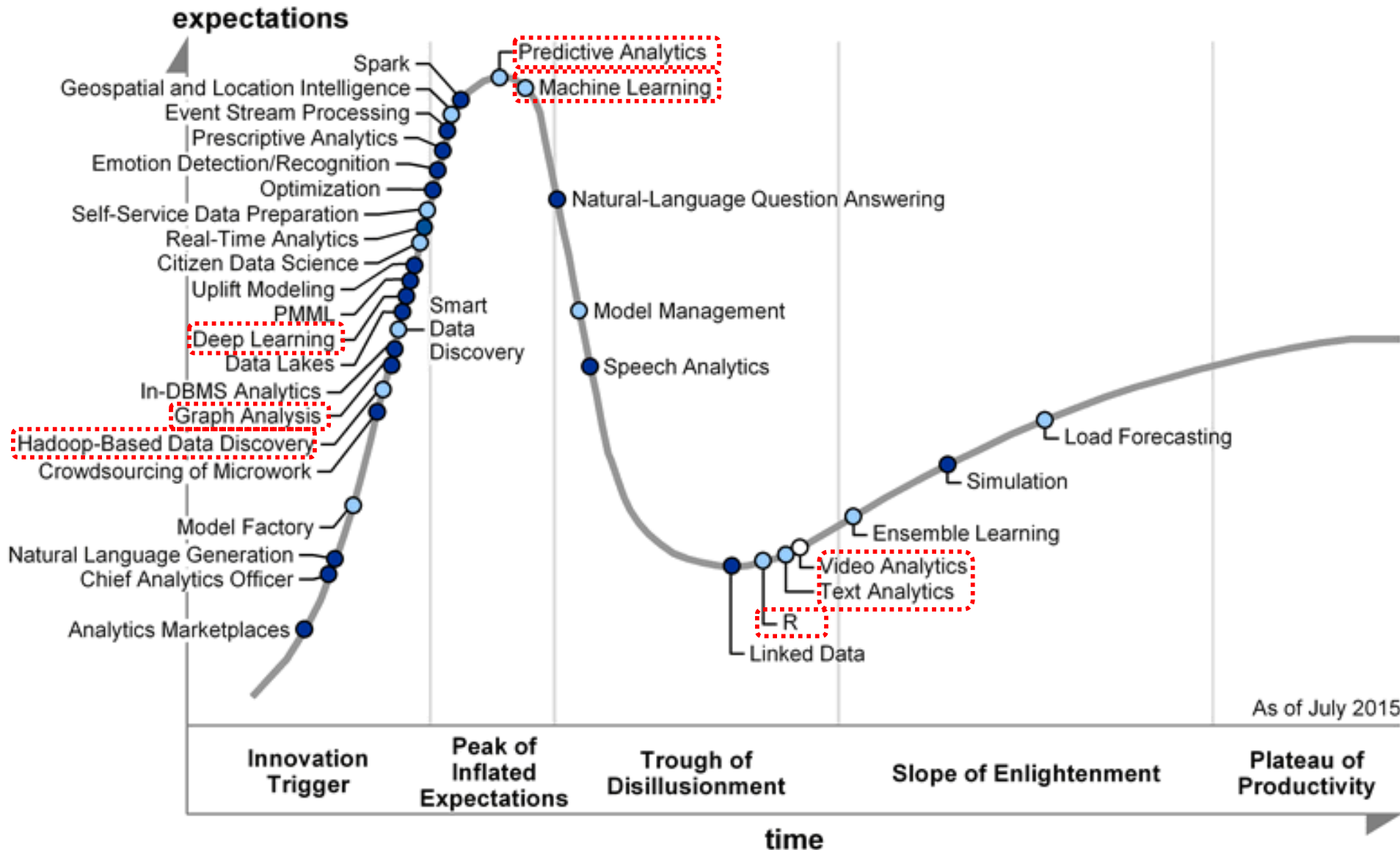
Hype Cycle for Information Infrastructure



Plateau will be reached in:

- less than 2 years
- 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

Hype Cycle for Advanced Analytics and Data Science

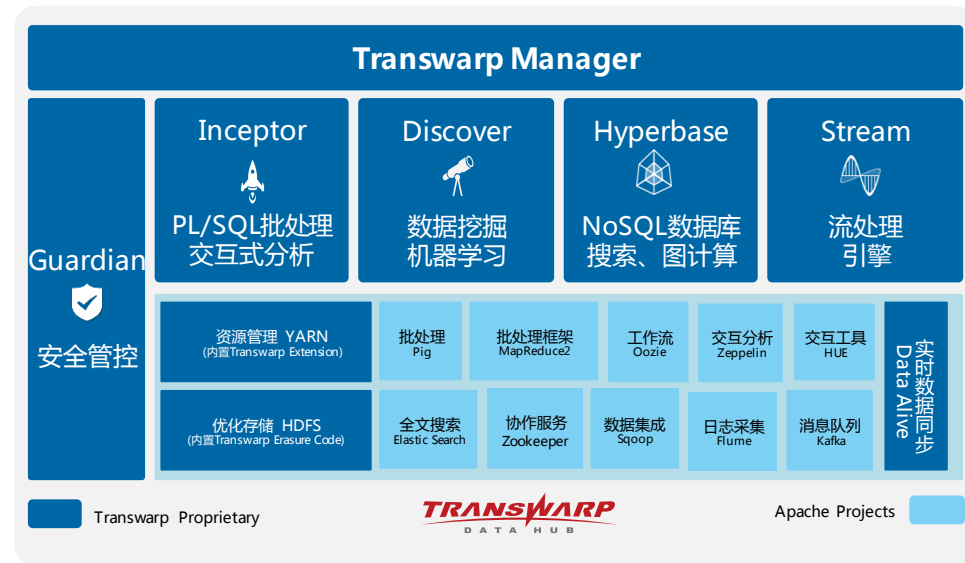
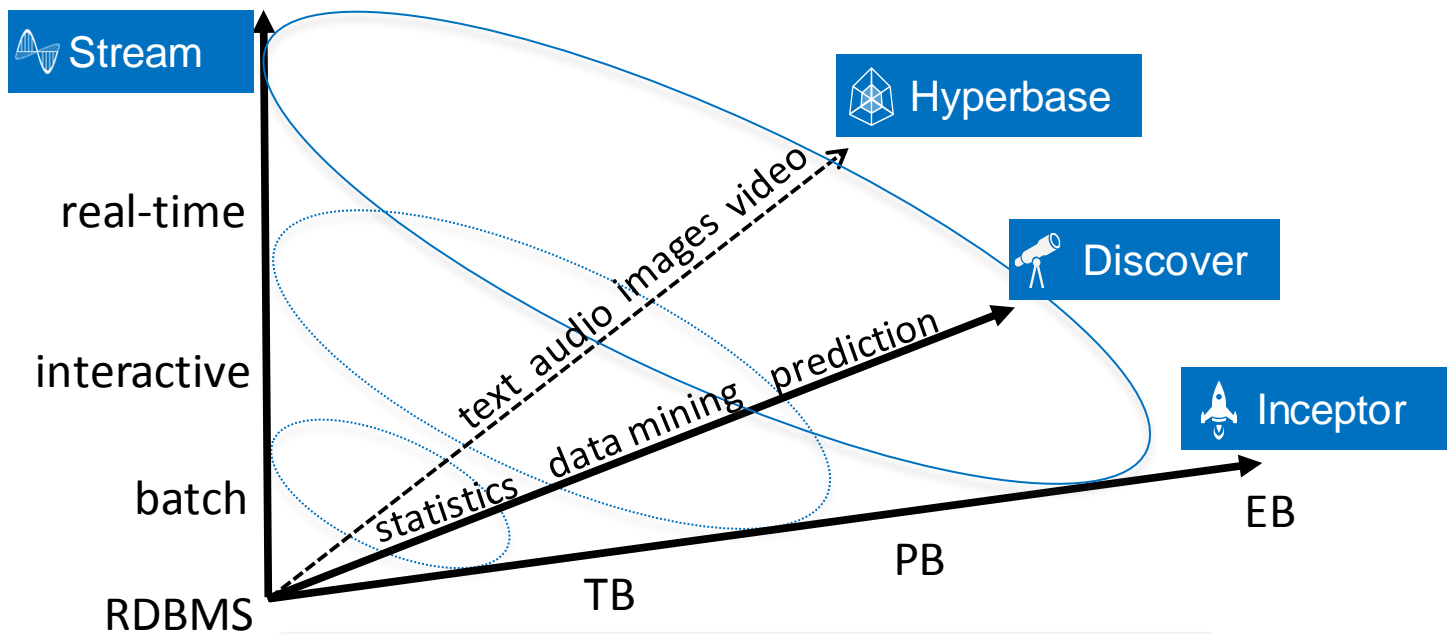


Plateau will be reached in:

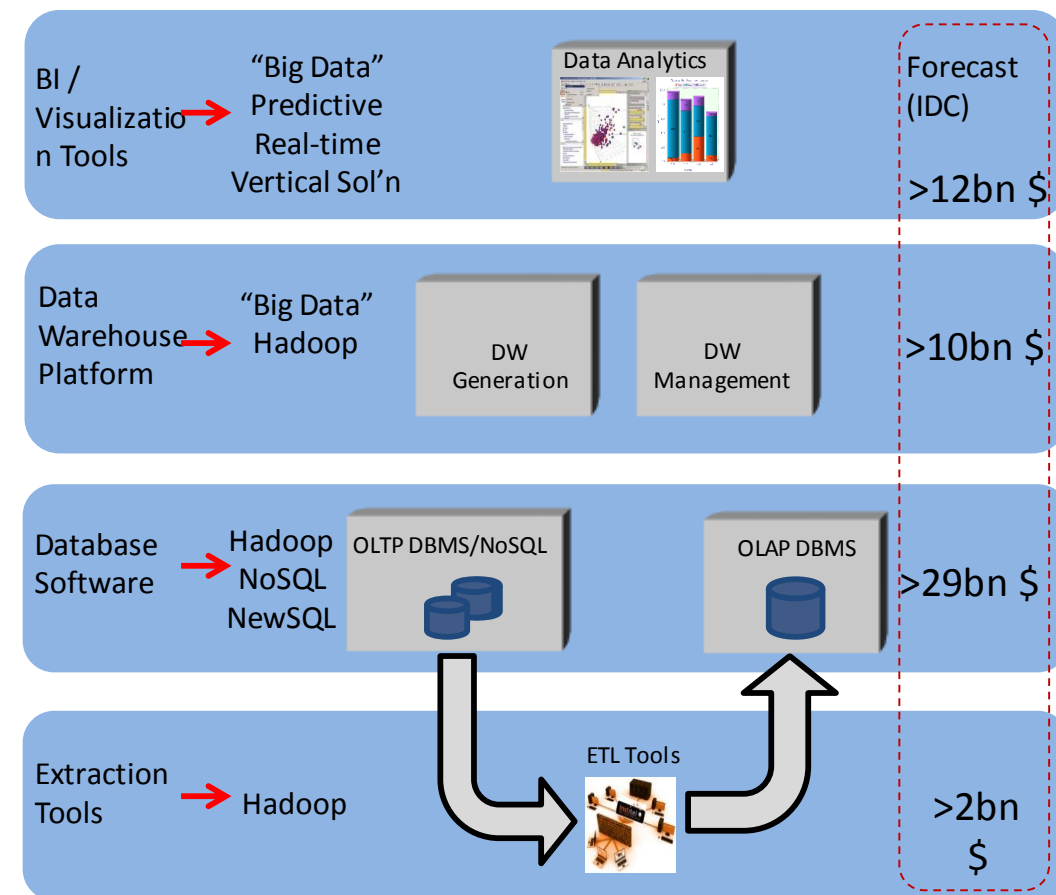
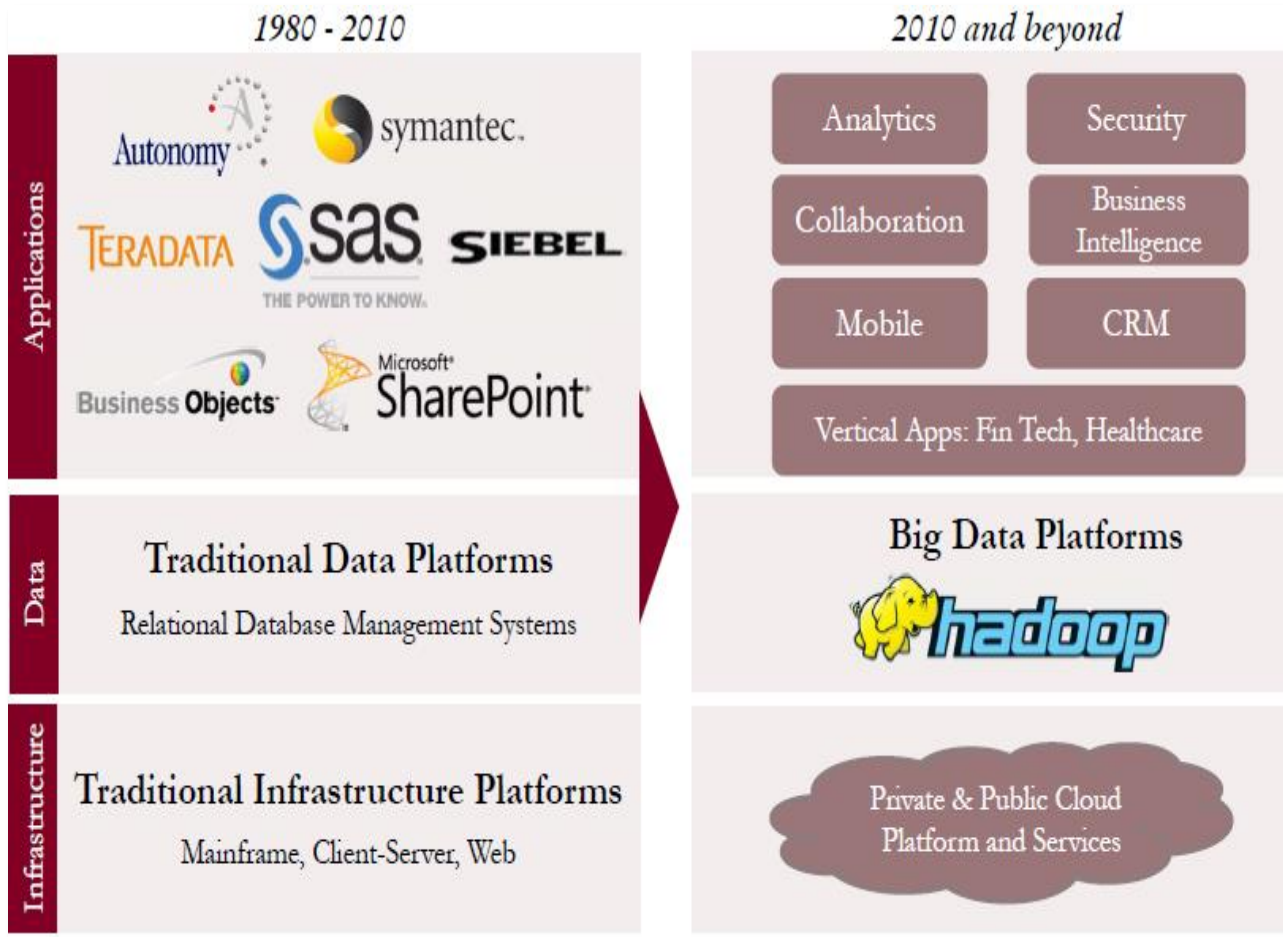
- less than 2 years
- 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

后Hadoop/Spark时代的技术发展趋势

新技术在四个维度上加速创新	
数据量 Volume	驱动力来自结构化数据的不断增加，企业需要性价比更高的技术方案
	SQL on Hadoop快速成熟，Big Data Database将替代传统relational database。传统的MPP分析型数据库将消失。
类型 Variety	驱动力来自对多种数据类型数据（文本、图片、音频、视频）的存储和分析需求
	深度机器学习技术（如TensorFlow）快速发展并得到应用，基于GPU / FPGA的加速技术逐渐普及。
速度 Velocity	从离线处理进步到实时数据处理，特别是IoT的广泛部署，推动了对实时计算的需求
	Flink, Apex, SqlStream, Internana, ParStream, Transwarp Stream等新技术，融合了批处理和流处理，提供强大易用的低延时实时计算能力，将逐渐取代现有流处理技术。
价值 Value	从历史统计发展到预测性分析。大数据的真正意义在于从数据中发现价值。
	数据挖掘、机器学习、图计算等产品和工具将日益普及，使用门槛将极大降低，普通业务人员很快能够自助进行分析建模。



Hadoop及其生态系统将重构数据处理市场

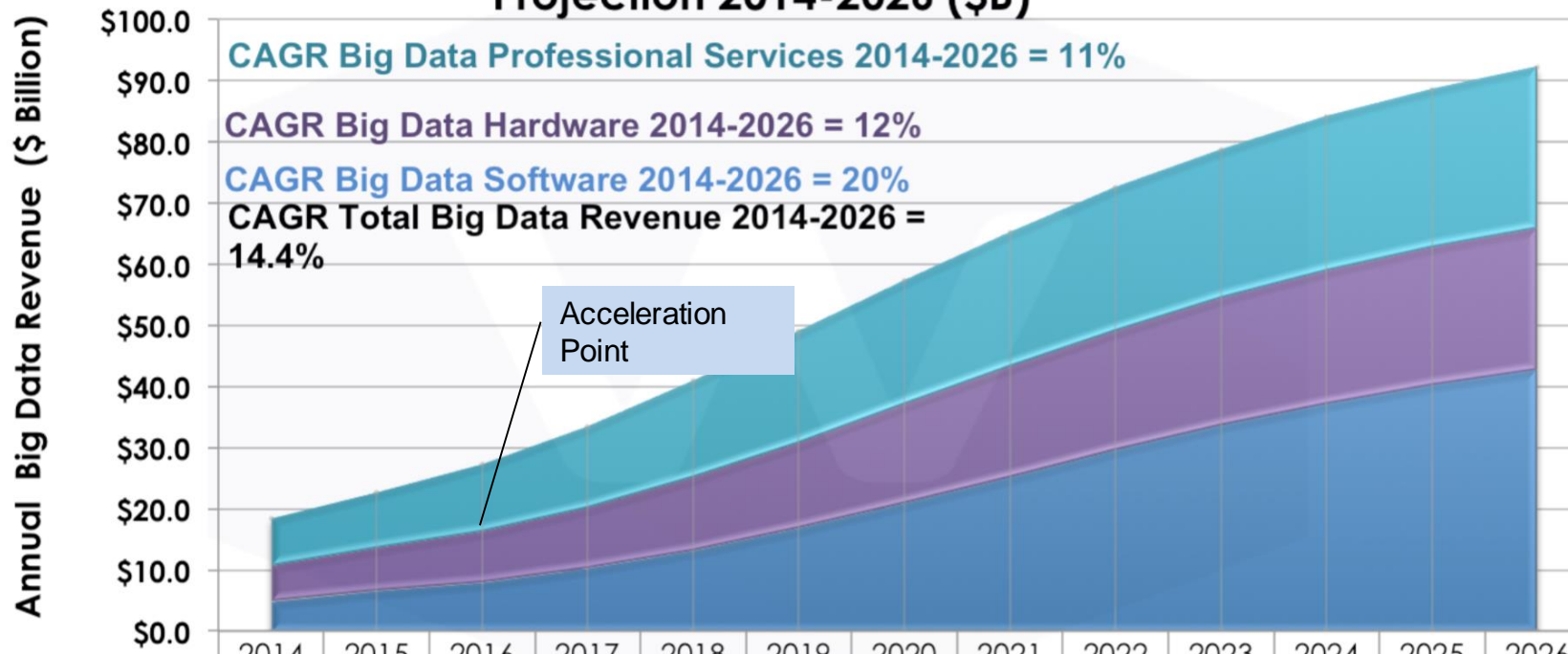


大数据产业链日益繁荣



2016年是Hadoop技术大规模应用的战略转折点

**Wikibon Big Data Software, Hardware & Professional Services
Projection 2014-2026 (\$B)**



未来五年大数据市场将以每年30%的速度增长

	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026
Professional Services	\$7.6	\$9.1	\$11.1	\$13.4	\$15.8	\$18.2	\$20.3	\$22.0	\$23.3	\$24.3	\$25.1	\$25.8	\$26.3
Big Data Hardware	\$5.7	\$6.9	\$8.2	\$9.9	\$11.8	\$14.0	\$16.1	\$18.0	\$19.6	\$20.8	\$21.8	\$22.5	\$23.1
Big Data Software	\$4.9	\$6.6	\$8.0	\$10.2	\$13.2	\$16.8	\$20.9	\$25.2	\$29.5	\$33.5	\$37.1	\$40.2	\$42.7
Total Big Data	\$18.3	\$22.6	\$27.3	\$33.5	\$40.8	\$49.0	\$57.3	\$65.2	\$72.4	\$78.7	\$84.0	\$88.5	\$92.2

Source: © Wikibon Big Data Project, 2016

The image features a stylized red logo for 'TRANSWARP' centered over a background of a planet's horizon and a bright sun. The logo consists of the word 'TRANSWARP' in a bold, italicized, sans-serif font. A red swoosh underline runs beneath the text, with a sharp, upward-pointing triangular spike at the 'W'. The background shows a curved horizon of a planet with a blue and green surface, and a bright yellow sun on the left side, creating a lens flare effect.

TRANSWARP