# Big Data Architecture

**Guido Schmutz**

Trivadis makes IT easier.

trivadis
makes IT easier.

# Guido Schmutz

Working for Trivadis for more than 18 years

Oracle ACE Director for Fusion Middleware and SOA

Co-Author of different books

Consultant, Trainer Software Architect for Java, Oracle, SOA and Big Data / Fast Data

Member of Trivadis Architecture Board

Technology Manager @ Trivadis

More than 25 years of software development experience

Contact: guido.schmutz@trivadis.com

Blog: http://guidoschmutz.wordpress.com

Twitter: gschmutz

# Agenda

**trivadis**
makes IT easier.

# Introduction

# Big Data is still "work in progress"

Choosing the right architecture is key for any (big data) project

Big Data is still quite a young field and therefore there are no standard architectures available which have been used for years

In the past few years, a few architectures have evolved and have been discussed online

Know the use cases before choosing your architecture

To have one/a few reference architectures can help in choosing the right components

**trivadis**
makes IT easier.

# Hadoop Ecosystem – many choices ....

| Unstructured Data Sources | Analytics | SQL on Hadoop | Management / Monitoring | Workflow/Job |
|---|---|---|---|---|

| Structured Data Sources | Core | Data Storage | Serialization | Security |
|---|---|---|---|---|

# Important Properties to choose a Big Data Architecture

Latency

Keep raw and un-interpreted data "forever" ?

Volume, Velocity, Variety, Veracity

Ad-Hoc Query Capabilities needed ?

Robustness & Fault Tolerance

Scalability

…

**trivadis**
makes IT easier.

# From Volume and Variety to Velocity

Big Data has evolved …

| Past |
|---|
| Big Data = Volume & Variety |

| Present |
|---|
| Big Data = Volume & Variety & Velocity |

and the Hadoop Ecosystem as well ….

| Past |
|---|
| Batch Processing |
| Time to insight of **Hours** |

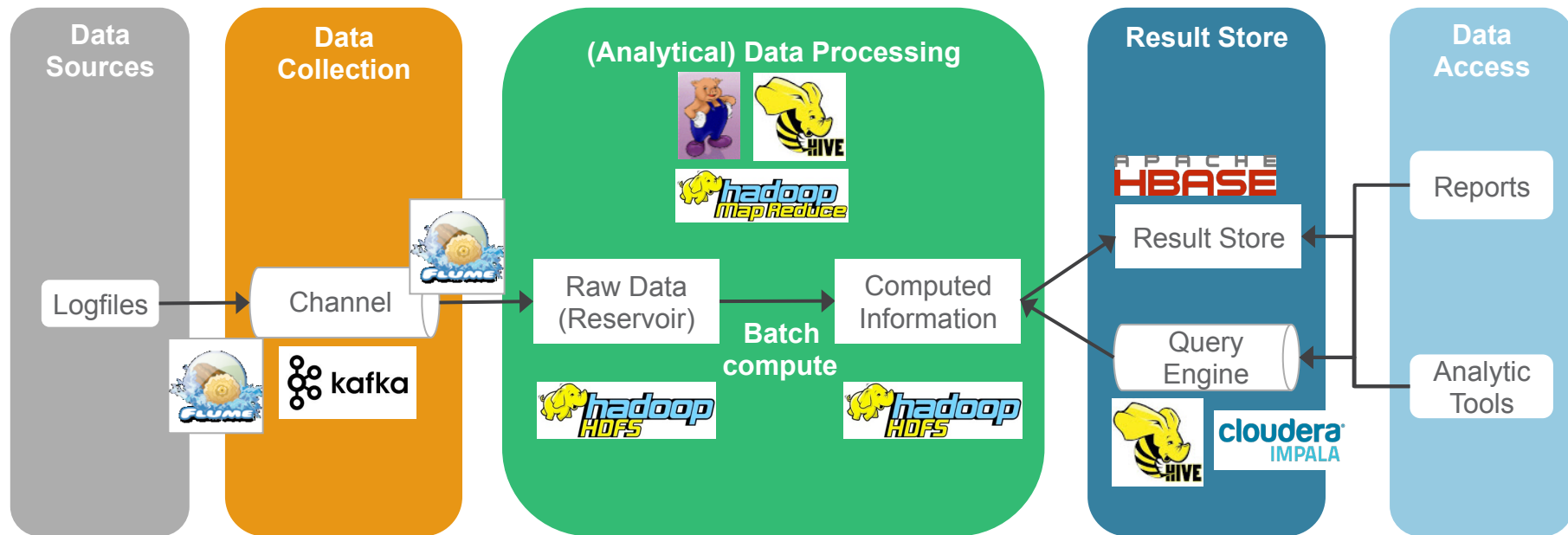| Present |
|---|
| Batch & Stream Processing |
| Time to insight in **Seconds** |

Adapted from Cloudera Blog article

**trivadis**
makes IT easier.

# Traditional Architecture for Big Data

# "Traditional Architecture" for Big Data



Data Sources: Logfiles, ERP, RDBMS, Social, Sensor, Machine, Mobile

Data Collection: Stage, Channel

(Analytical) Data Processing: Raw Data (Reservoir), Batch compute, Computed Information

Result Storage: Result Store, Query Engine

Data Access: Reports, Service, Analytic Tools, Alerting Tools

= Data in Motion    = Data at Rest

trivadis
makes IT easier.
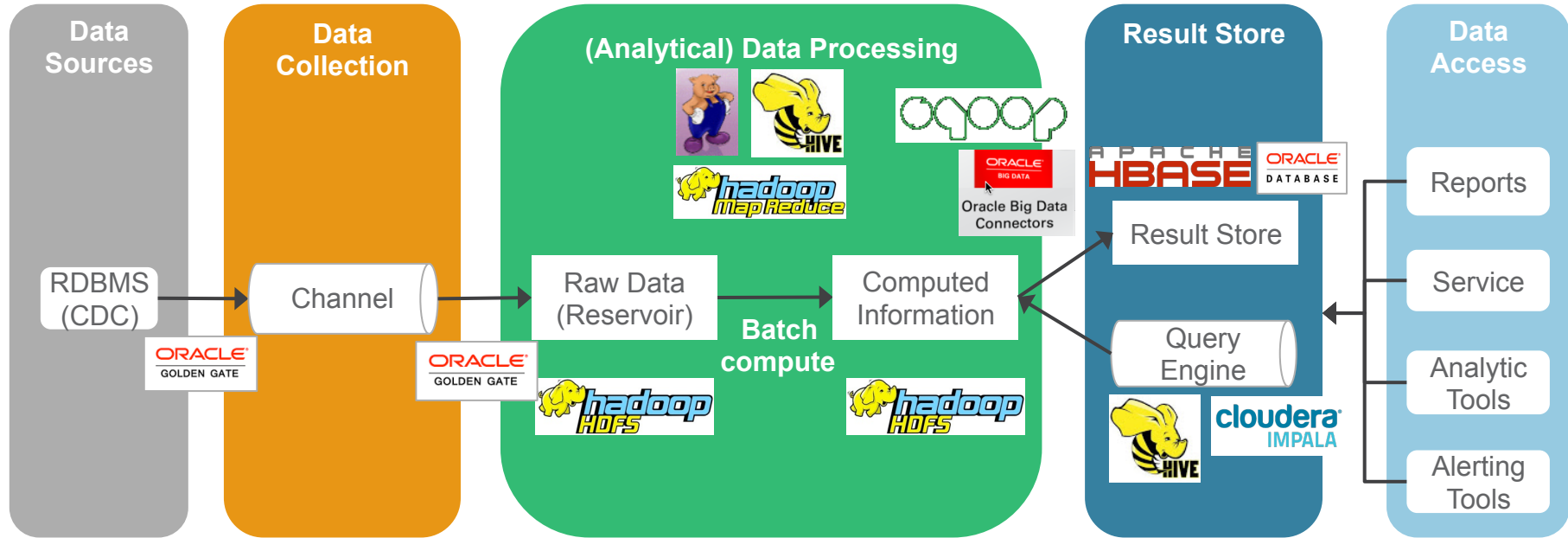
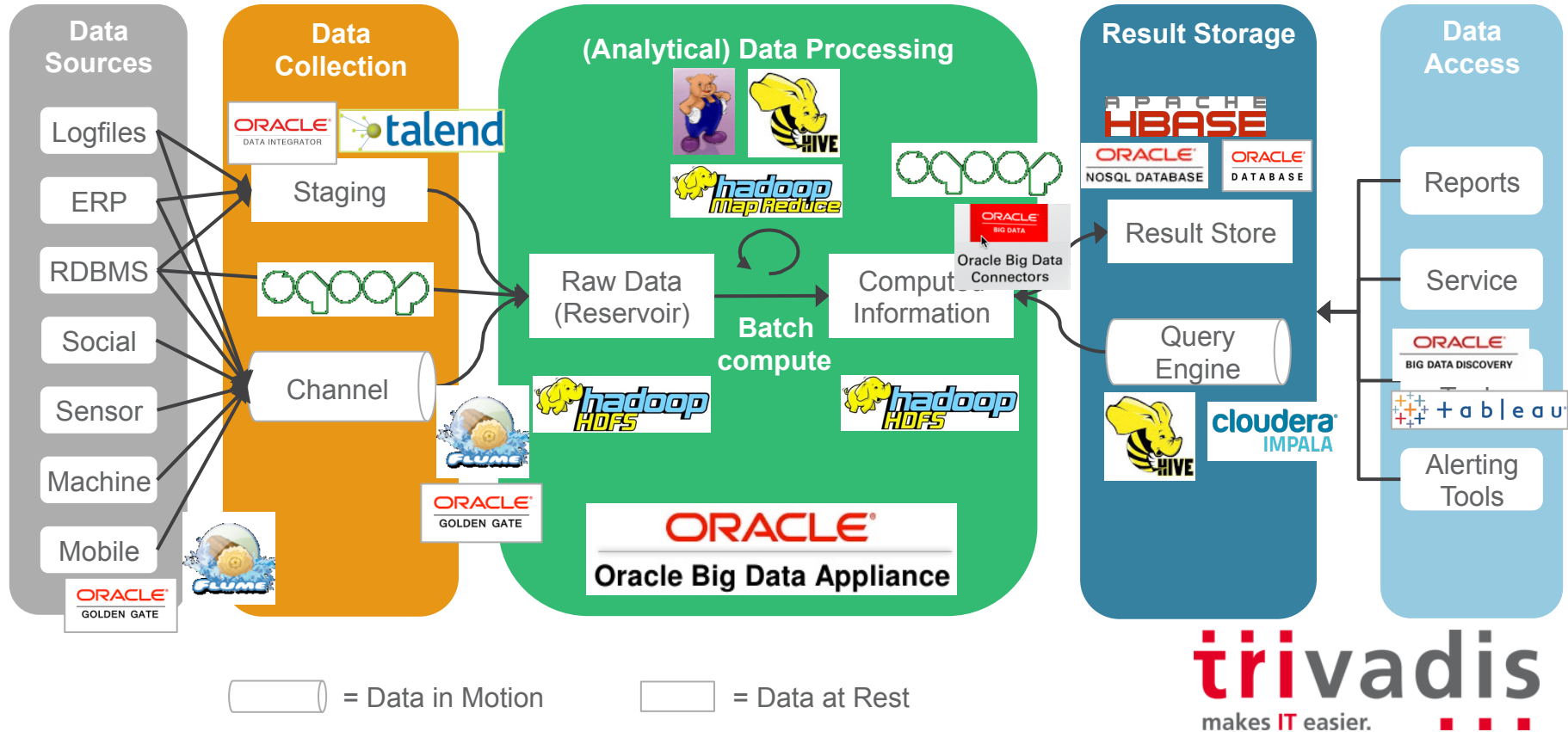# Use Case 1) – Click Stream analysis: 360 degree view of customer

# Use Case 2) – Ingest Relational Data into Hadoop and make it accessible

# Use Case 2a) – Ingest Relational Data into Hadoop and make it accessible

# Apache Spark – the new kid on the block

Apache Spark is a fast and general engine for large-scale data processing

- The hot trend in Big Data!
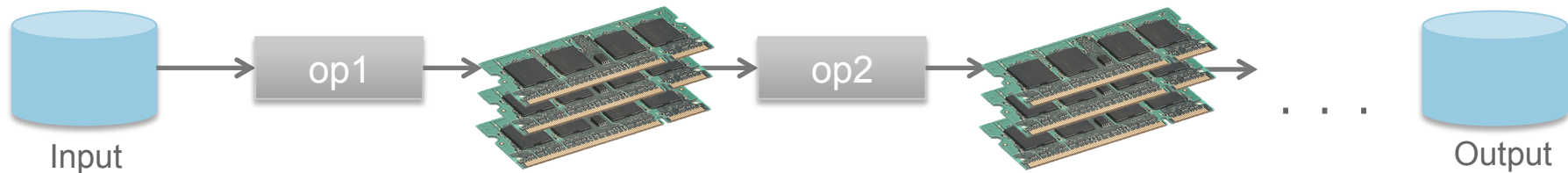- Originally developed 2009 in UC Berkley's AMPLab
- Can run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk
- One of the largest OSS communities in big data with over 200 contributors in 50+ organizations
- Open Sourced in 2010 – since 2014 part of Apache Software foundation
- Supported by many vendors

# Motivation – Why Apache Spark?

Apache Hadoop MapReduce: Data Sharing on Disk



Apache Spark: Speed up processing by using Memory instead of Disks

# Apache Spark "Ecosystem"

## Libraries

| Spark SQL (Batch Processing) | Blink DB (Approximate Querying) | Spark Streaming (Real-Time) | MLlib, Spark R (Machine Learning) | GraphX (Graph Processing) |
|---|---|---|---|---|

## Core Runtime

**Spark Core API and Execution Model**

## Cluster Resource Managers

| Spark Standalone | MESOS | YARN |
|---|---|---|

## Data / Data Stores

| HDFS | Elastic Search | NoSQL | S3 |
|---|---|---|---|

**trivadis**
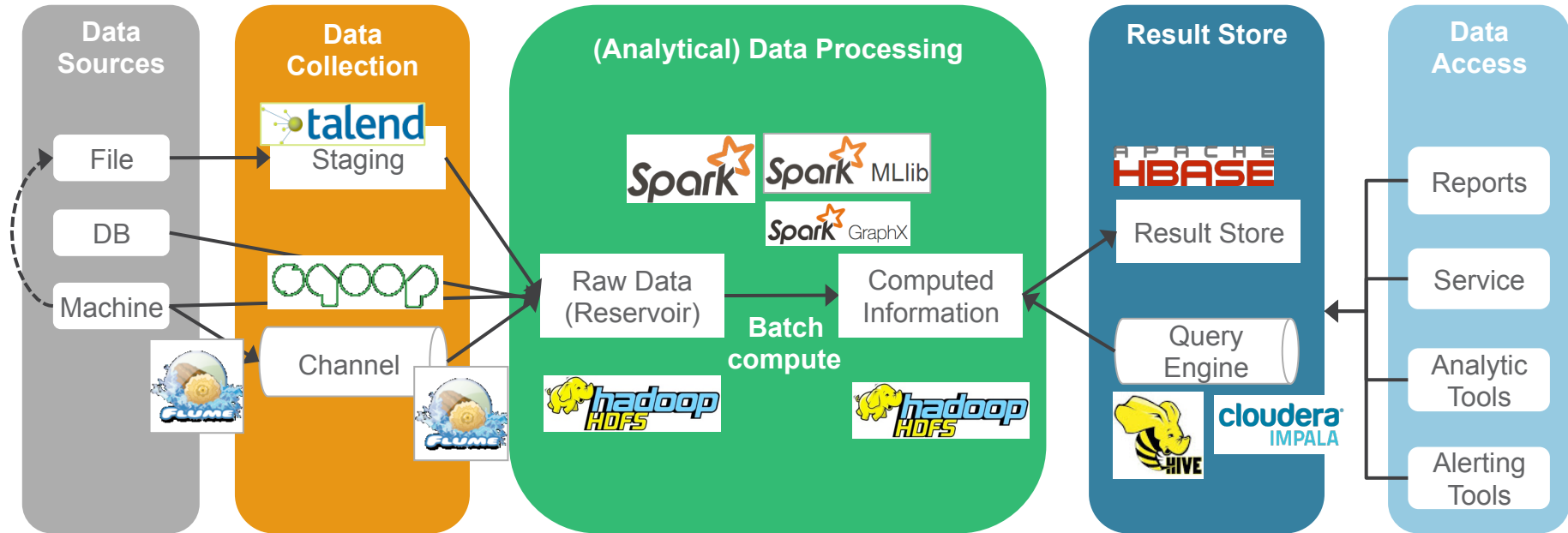makes IT easier.
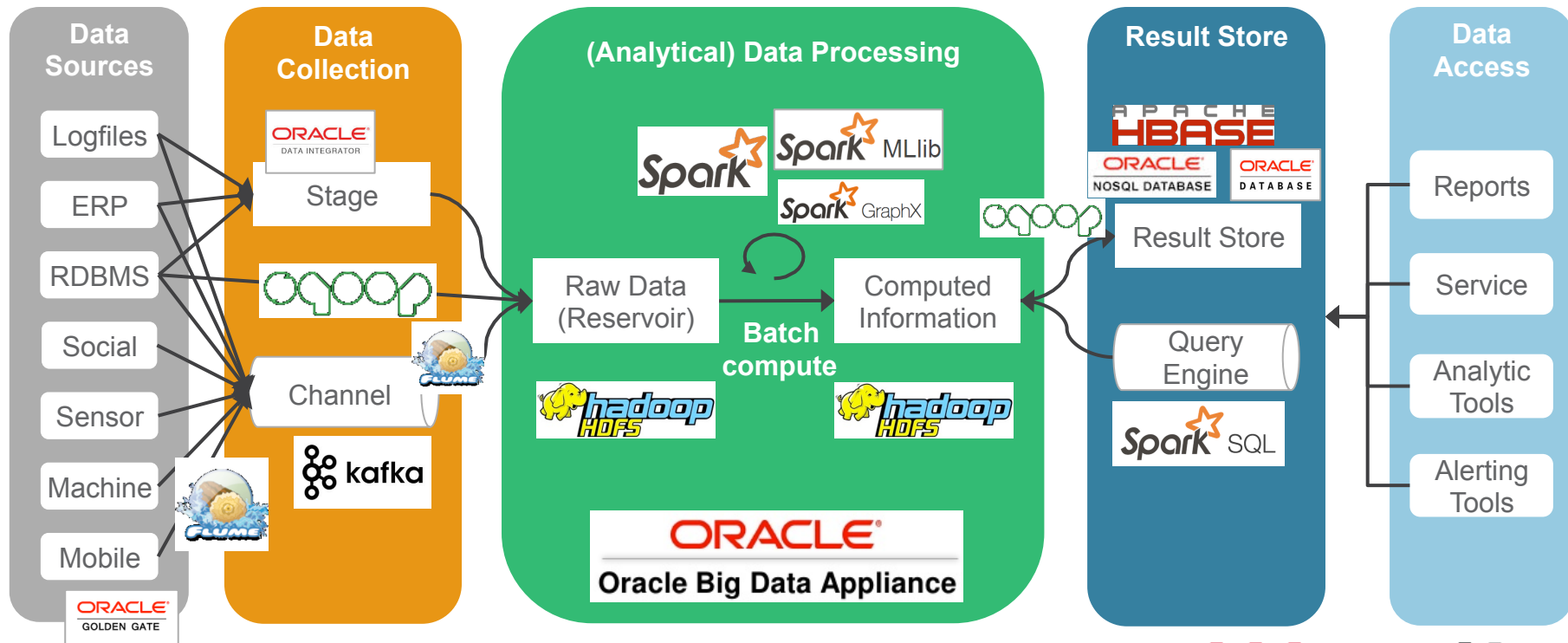
# Use Case 3) – Predictive Maintenance through Machine Learning on collected data
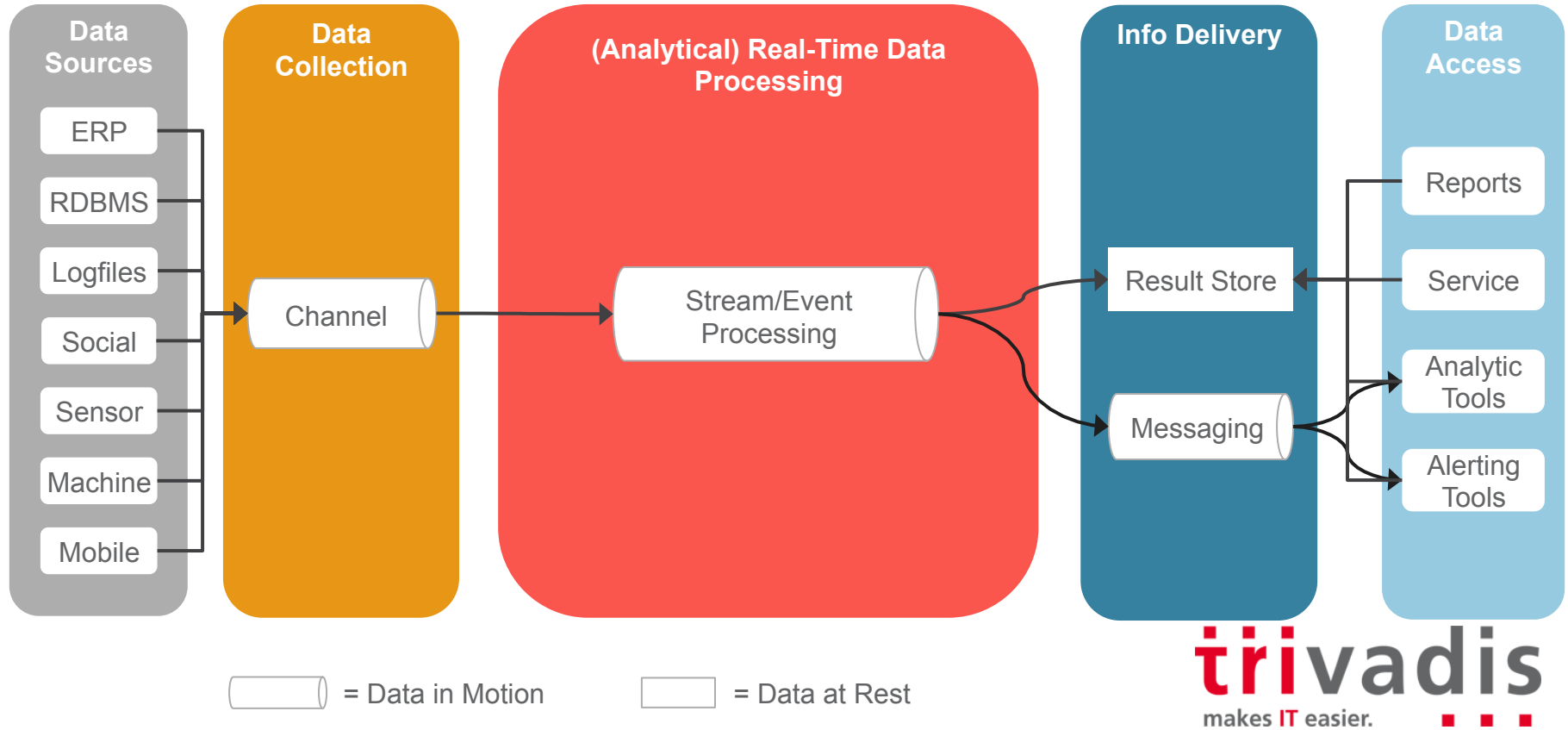
# "Spark Ecosystem" Technology Mapping

# Traditional Architecture for Big Data

- Batch Processing

- Not for low latency use cases

- Spark can speed up, but if positioned as alternative to Hadoop Map/Reduce, it's still Batch Processing

- Spark Ecosystems offers a lot of additional advanced analytic capabilities (machine learning, graph processing, …)
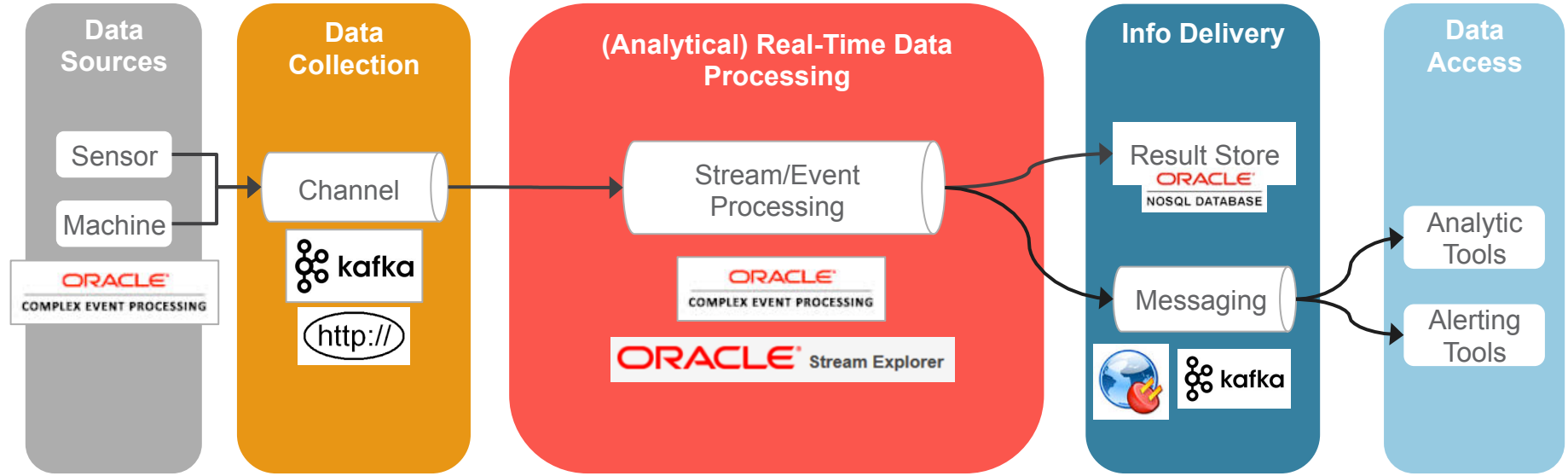
**trivadis**
makes IT easier.

# Streaming Analytics Architecture for Big Data
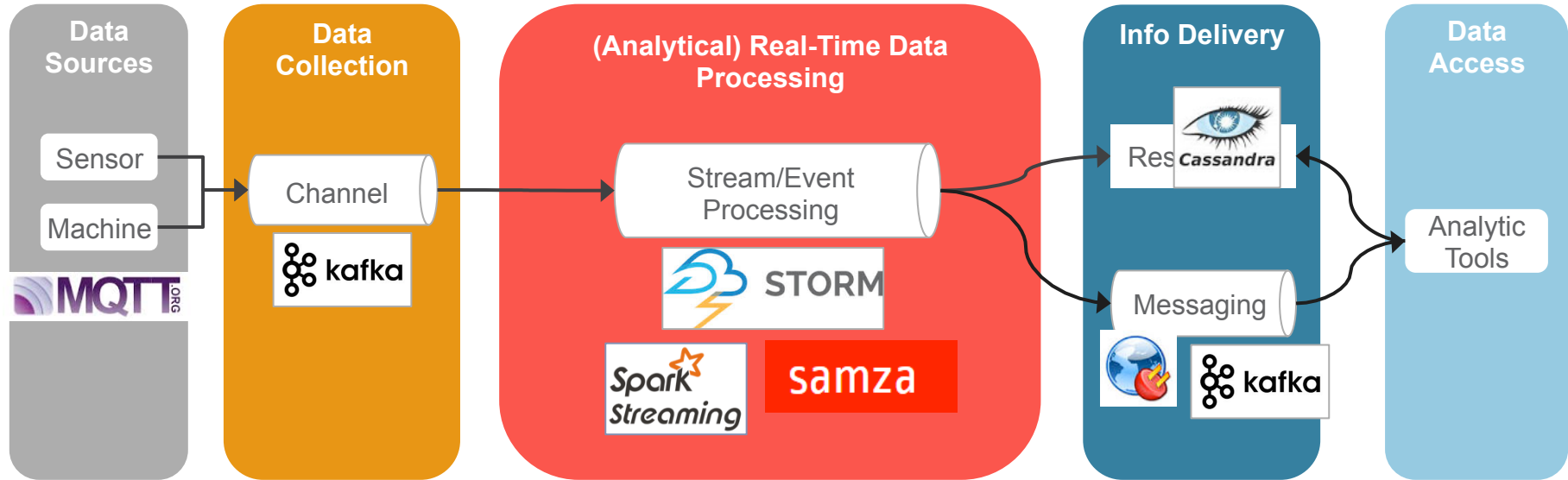
# Streaming Analytics Architecture for Big Data

aka. (Complex) Event Processing)

# Use Case 4) Alerting in Internet of Things (IoT)

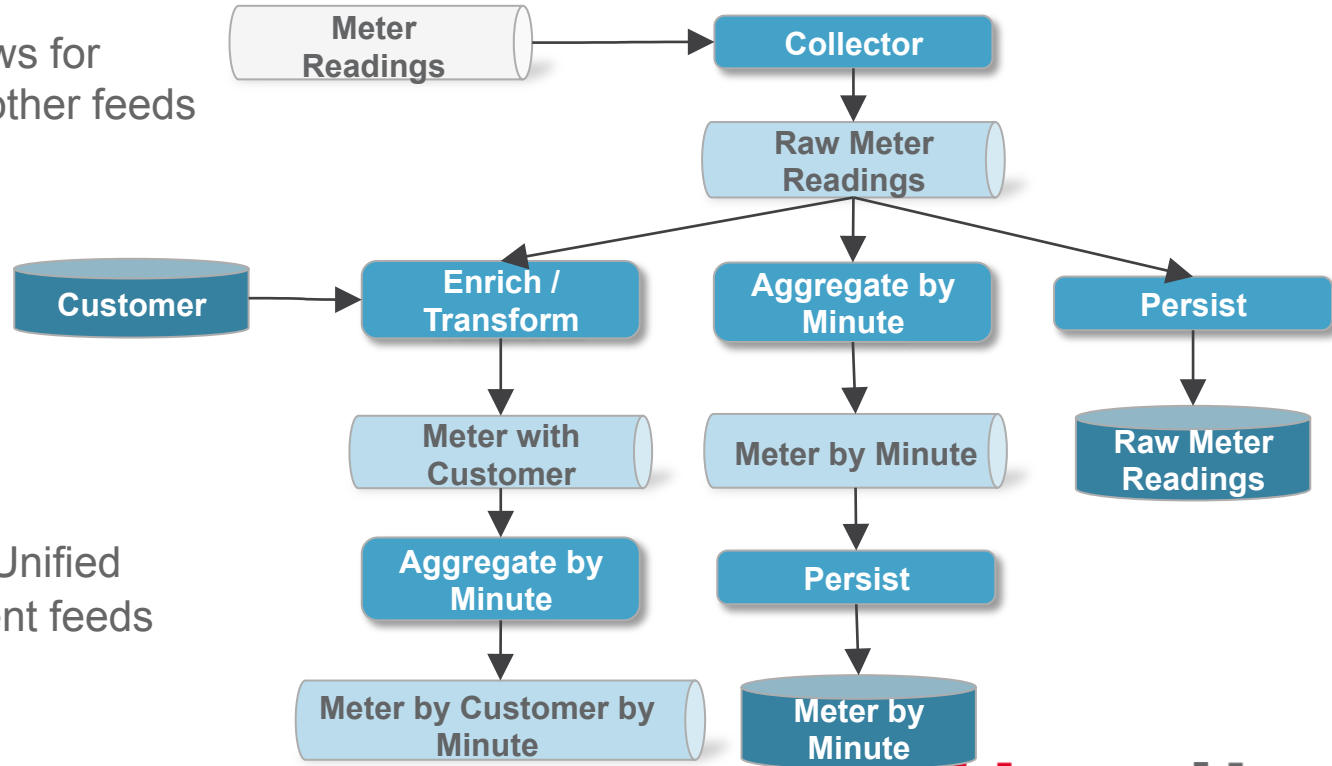# ■ Use Case 5) Real-Time Analytics on Sensor Events

# Unified Log (Event) Processing

Stream processing allows for computing feeds off of other feeds

Derived feeds are no different than original feeds they are computed off

Single deployment of "Unified Log" but logically different feeds

```
Meter Readings → Collector
                    ↓
              Raw Meter Readings
              ↙        ↓        ↘
Customer → Enrich /   Aggregate by    Persist
          Transform    Minute           ↓
              ↓          ↓         Raw Meter Readings
        Meter with   Meter by Minute
        Customer         ↓
              ↓        Persist
        Aggregate by     ↓
         Minute      Meter by Minute
              ↓
    Meter by Customer by Minute
```

**trivadis**
makes IT easier.

# Streaming Analytics Technology Mapping

# Streaming Analytics Architecture for Big Data

The solution for low latency use cases

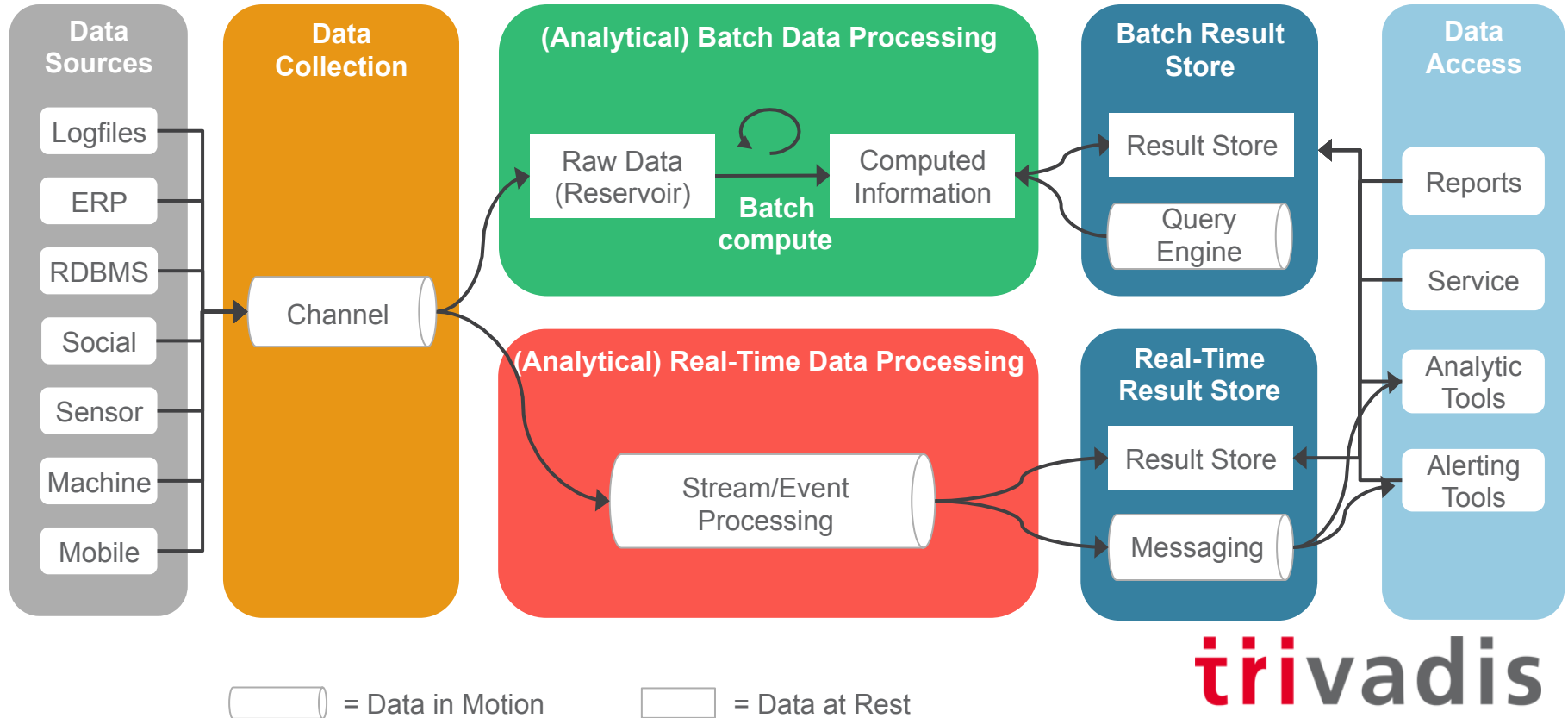Process each event separately => low latency

Process events in micro-batches => increases latency but offers better reliability

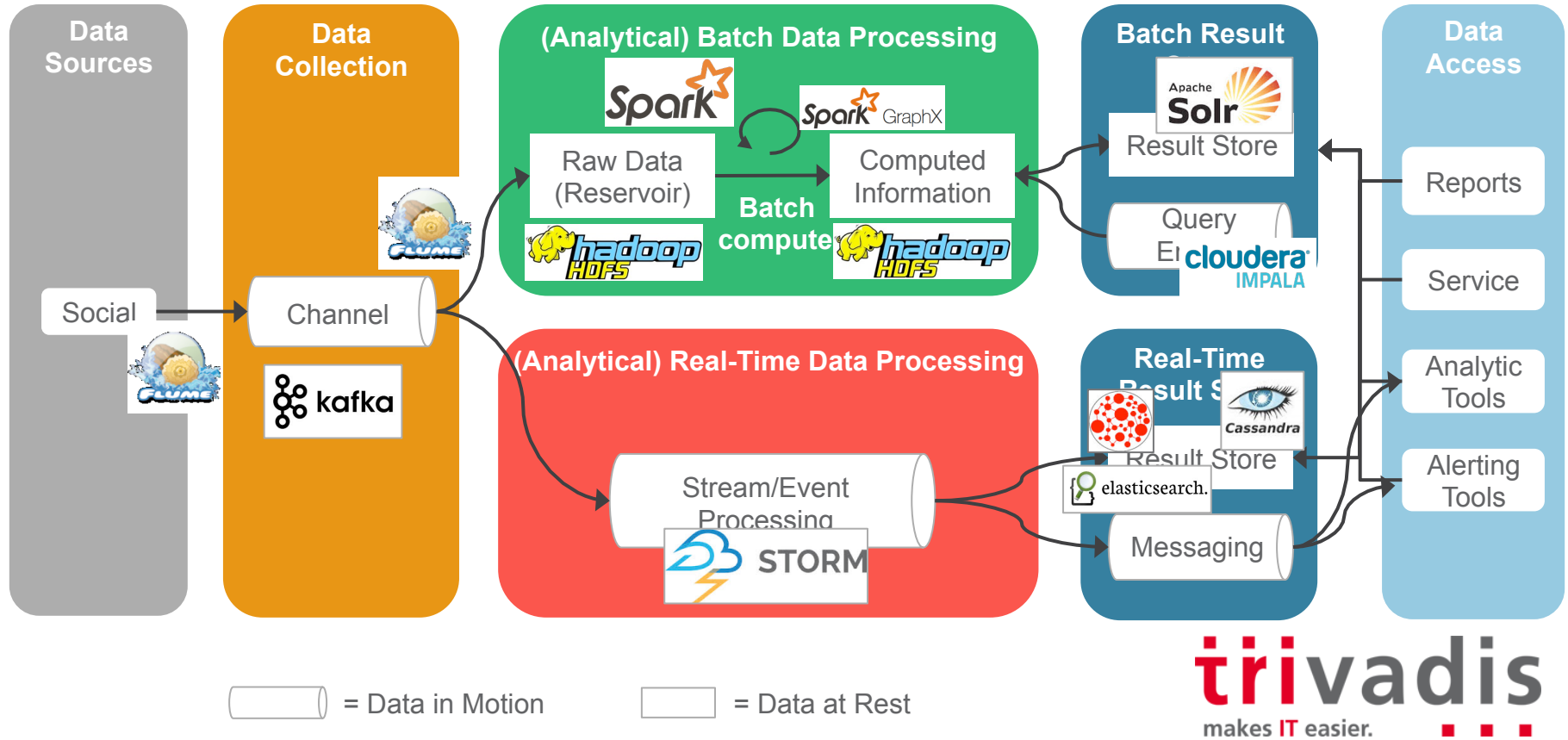Previously known as "Complex Event Processing"

Keep the data moving / Data in Motion instead of Data at Rest => raw events are (often) not stored

**trivadis**
makes IT easier.

# Lambda Architecture for Big Data

# "Lambda Architecture" for Big Data

# Use Case 6) Social Media and Social Network Analysis



**Data Sources**

**Data Collection**

**(Analytical) Batch Data Processing**

Raw Data (Reservoir)

Batch compute

Computed Information

**Batch Result**

Result Store

Query Engine

**Data Access**

Reports

Service

Analytic Tools

Alerting Tools

Social

Channel

kafka

**(Analytical) Real-Time Data Processing**

Stream/Event Processing

STORM

**Real-Time Result S**

Result Store

elasticsearch.

Messaging

Cassandra

= Data in Motion    = Data at Rest

trivadis
makes IT easier.

# Lambda Architecture for Big Data

Combines (Big) Data at Rest with (Fast) Data in Motion

Closes the gap from high-latency batch processing

Keeps the raw information forever

Makes it possible to rerun analytics operations on whole data set if necessary
=> because the old run had an error or
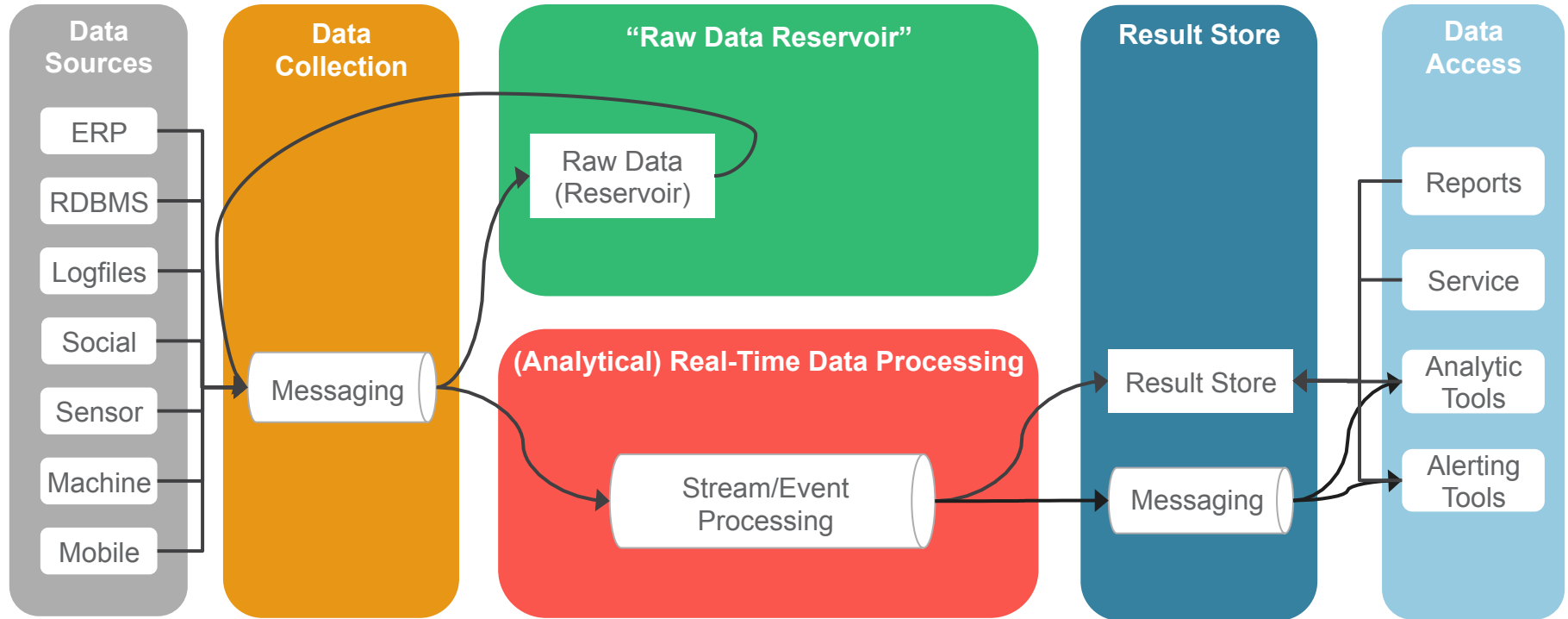=> because we have found a better algorithm we want to apply

Have to implement functionality twice
- Once for batch
- Once for real-time streaming

# „Kappa" Architecture for Big Data

trivadis
makes IT easier.

# "Kappa Architecture" for Big Data

# Kappa Architecture for Big Data

The solution for low latency use cases

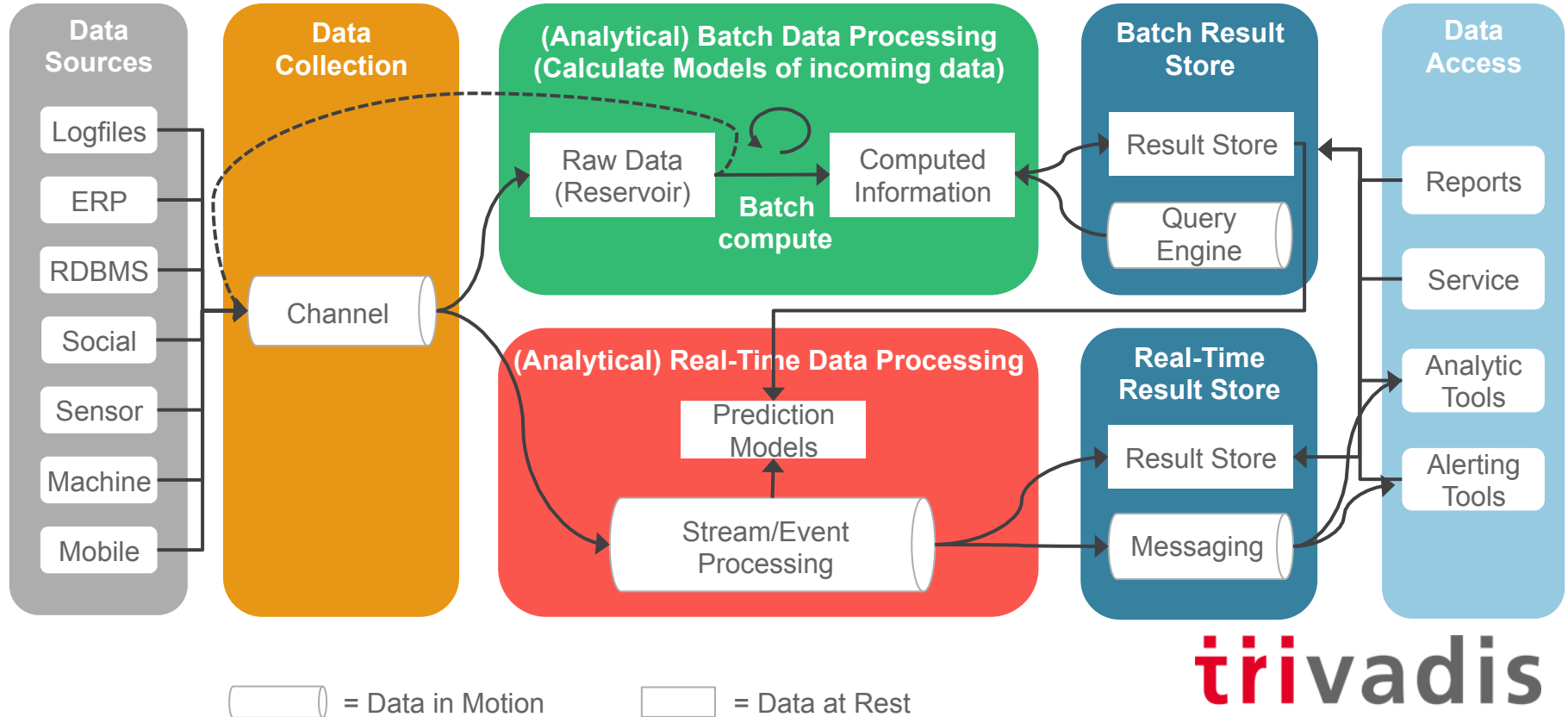Process each event separately => low latency

Process events in micro-batches => increases latency but offers better reliability
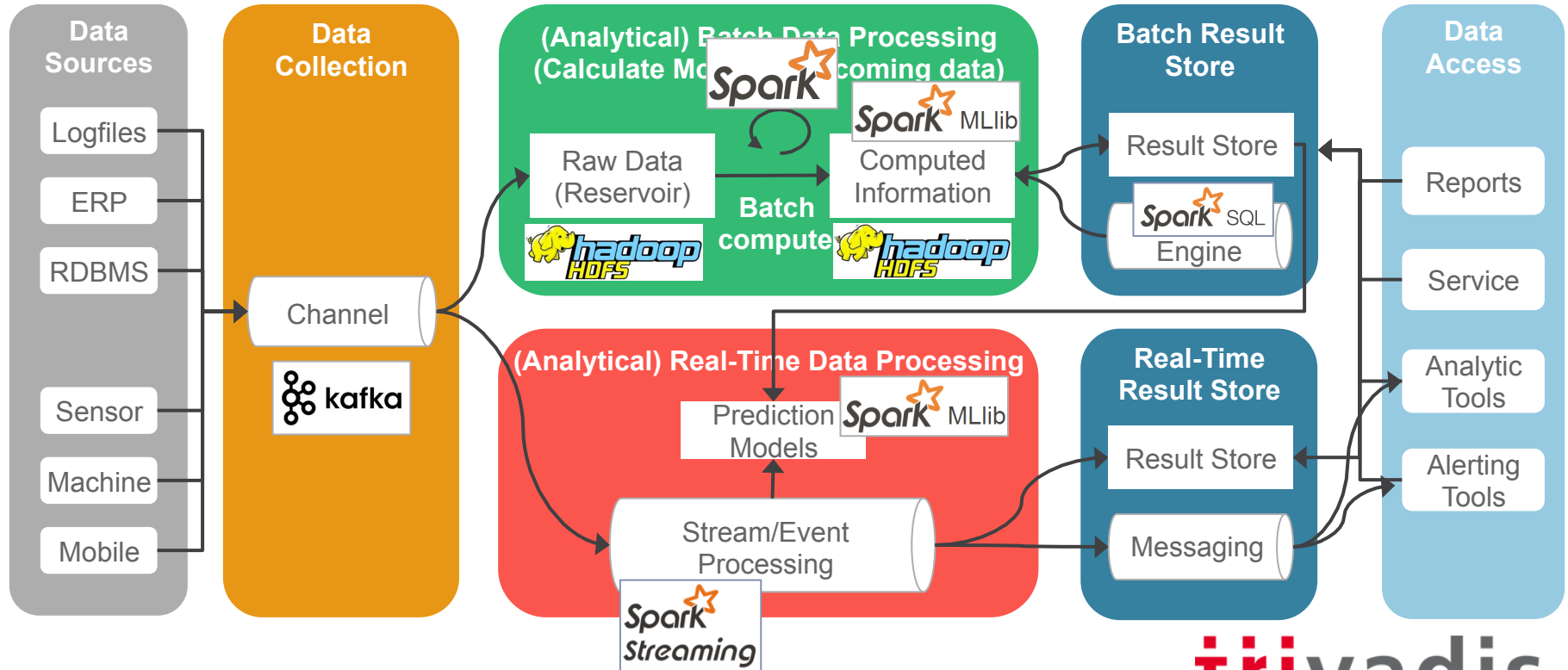
Previously known as "Complex Event Processing"

Keep the data moving / Data in Motion instead of Data at Rest

trivadis
makes IT easier.

# „Unified" Architecture for Big Data

# "Unified Architecture" for Big Data

# Use Case 7) Fraud Detection

# Summary

# Summary

Know your use cases and then choose your architecture and the relevant components/ products/frameworks

You don't have to use all the components of the Hadoop Ecosystem to be successful

Big Data is still quite a young field and therefore there are no standard architectures available which have been used for years

Lambda, Kappa Architecture are best practices architectures which you have to adapt to your environment

**trivadis**
makes IT easier.

Trivadis
makes IT
easier.

**Name Referent**
**Titel Referent**

**Tel. +00 00 000 00 00**
**vorname.name@trivadis.com**

Trivadis
makes IT
easier.

**trivadis**
makes IT easier.