# P4.1 Reference Architectures for Enterprise Big Data Use Cases

**Romeo Kienzler, Data Scientist, Advisory Architect, IBM Germany, Austria, Switzerland**
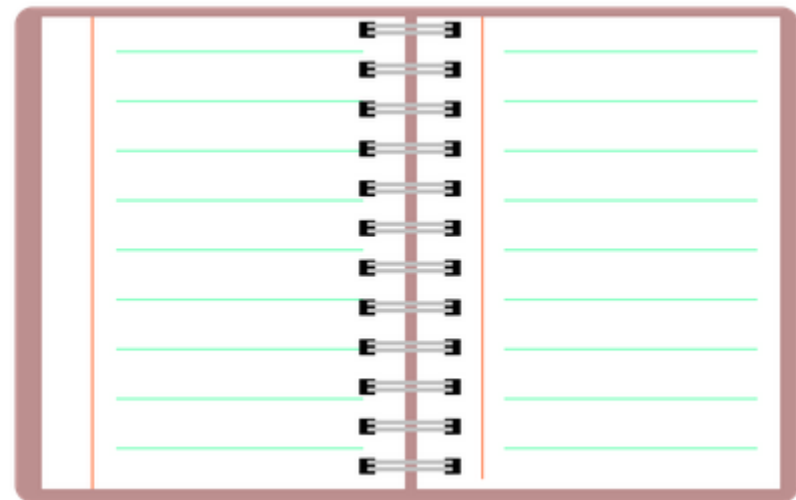
IBM Center of Excellence for Data Science, Cognitive Systems and BigData
(A joint-venture between IBM Research Zurich and IBM Innovation Center DACH)

ECSA 2014
8th European Conference on Software Architecture
Vienna, Austria, 25-29 August 2014

# Agenda

- Motivation
- Use Cases
- How Databases scale
- Evolution of Large Scale Data Processing
- Requirements and Ingredients
- Architectural Proposal

# Motivation – the World before 2000



.coms

Traditional Enterprises

MySQL, Postgres

DB2, Oracle, Teradata

# Motivation – the World after 2000


.coms


Traditional Enterprises


NoSQL


DB2, Oracle, Teradata

# Use Cases – Basic Idea

- Use ALL available data independently weather it is

# Use Cases – Basic Idea

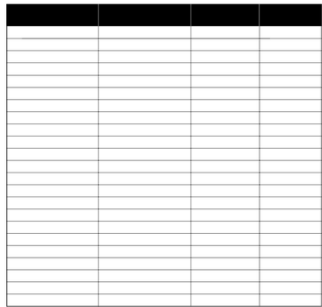- Use ALL available data independently weather it is
  - Inside

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
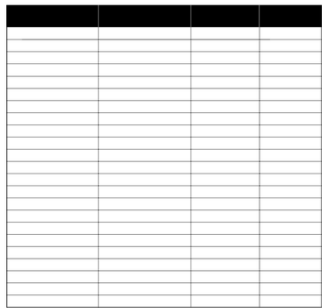    - Inside or outside your company

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
    - Inside or outside your company
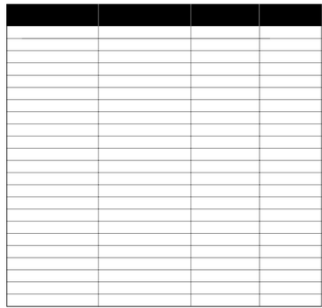    - Structured

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
  - Inside or outside your company
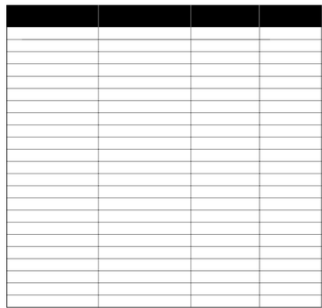  - Structured, semi-structured

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
    - Inside or outside your company
    - Structured, semi-structured, unstructured

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
    - Inside or outside your company
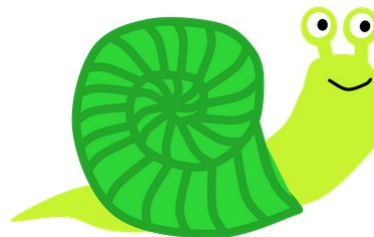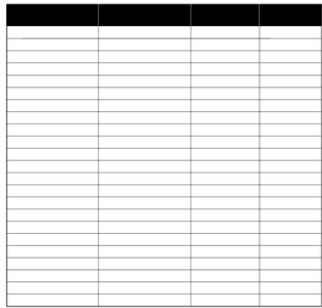    - Structured, semi-structured, unstructured or binary
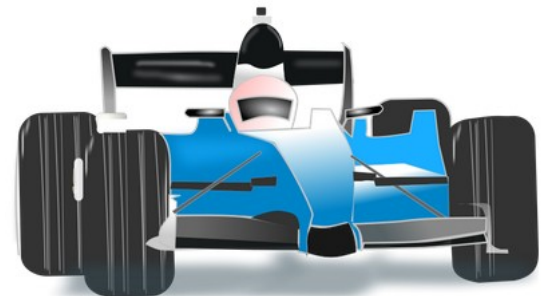
# Use Cases – Basic Idea

- Use ALL available data independently weather it is
  - Inside or outside your company
  - Structured, semi-structured, unstructured or binary
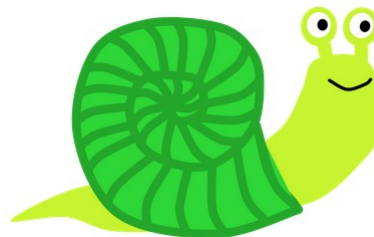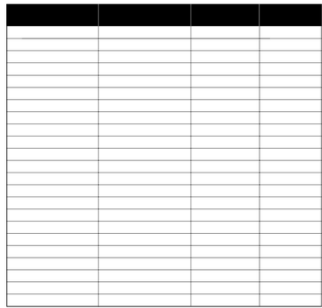  - At rest

# Use Cases – Basic Idea

- Use ALL available data independently weather it is
  - Inside or outside your company
  - Structured, semi-structured, unstructured or binary
  - At rest or in motion
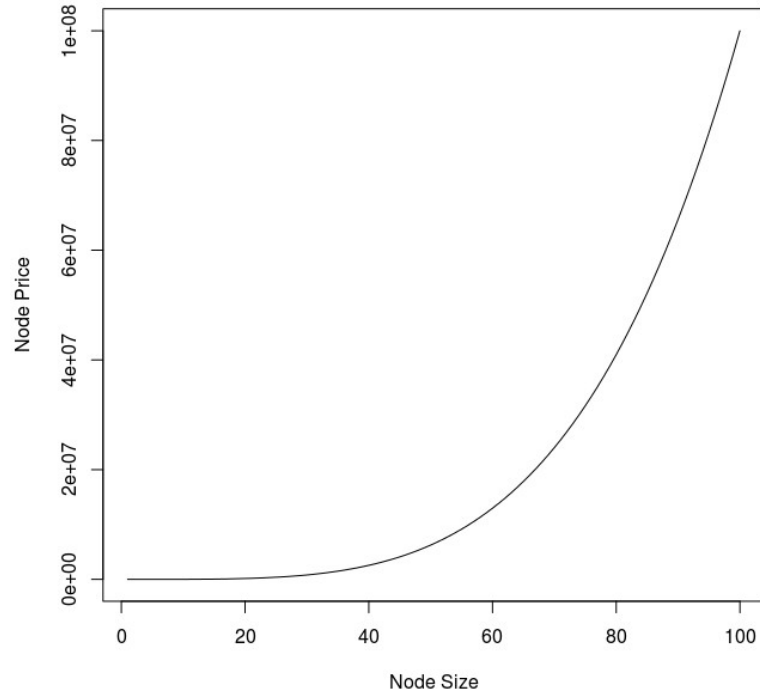
# How do Databases Scale?

- Scale-out because Scale-up not possible
  - Why?
  - There is a sweet spot (global optimum) for the ideal node size
  - Determines the number of cluster nodes
  - Because node price vs node size is not linear

# How do Databases Scale, Conclusion

- To minimize overall cluster price
    - Use many node because of rather small node size
    - Use commodity hardware
        - Fault tolerance
        - Won't go into CAP theorem here → Google
    - For dynamic workloads and dynamic scale-in/out → CLOUD

# Current situation

- Current situation in Enterprises
    - BI Tools
    - SQL (42%)
    - R (33%)
    - Python (26%)
    - Excel (25%)
    - Java, Ruby, C++ (17%)
    - SPSS, SAS (9%)

- Current situation in .coms
    - Writing MapReduce Jobs
    - Using proprietary query languages
    - Code everything from scratch

# Evolution

- Programing language implementations
- Usage of high level query languages
    - Pig
    - Jaql
- SQL or SQL like
    - BigSQL
    - HQL
    - CQL
- R push back
    - BigR
    - Rhadoop
- BigData spread sheets
- BI push back
- SPSS push back

# Pig(Latin) / JAQL

```
input_lines = LOAD '/tmp/my-copy-of-all-pages-on-internet' AS (line:chararray);

-- Extract words from each line and put them into a pig bag
-- datatype, then flatten the bag to get one word on each row
words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;

-- filter out any words that are just white spaces
filtered_words = FILTER words BY word MATCHES '\\w+';

-- create a group for each word
word_groups = GROUP filtered_words BY word;

-- count the entries in each group
word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS count, group AS word

-- order the records by count
ordered_word_count = ORDER word_count BY count DESC;
STORE ordered_word_count INTO '/tmp/number-of-words-on-internet';
```
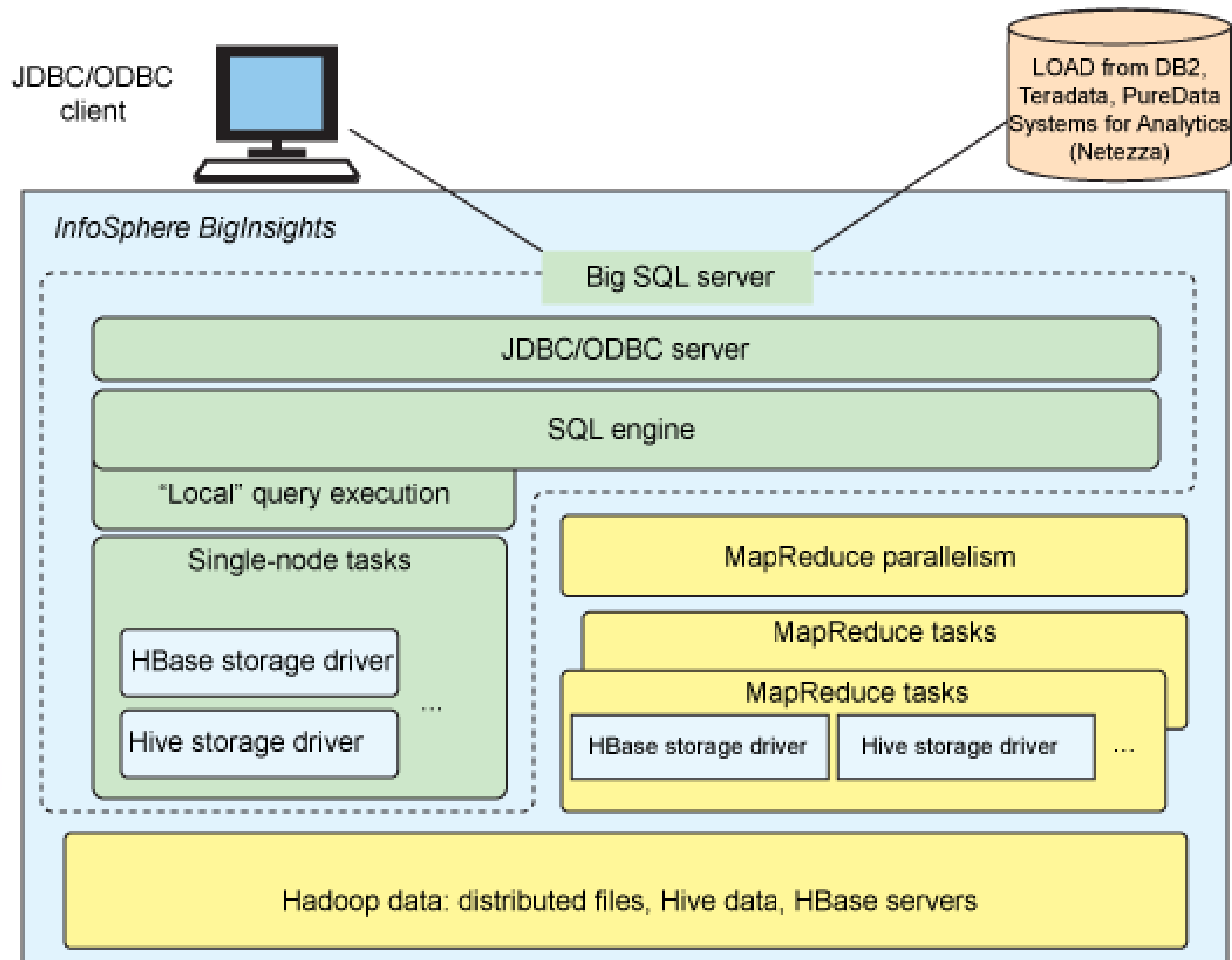
```
'home/biadmin/Documents/tweets-short.txt ',
read({type: 'hdfs', location: file,
                    inoptions: {format: 'org.apache.hadoop.mapred.TextInputFormat',
                    converter: 'com.ibm.jaql.io.hadoop.converter.FromJsonTextConverter'}});
tweets.statuses -> transform { $.created_at,
            tweet_id: $.id_str,
            $.geo,
            user_followers_count: $.user.followers_count,
            $.lang,
            text: $.text };
> write(del("/user/root/tweets.del", schema = schema {created_at
                    tweet_id,
                    geo,
                    user_followers_count,
                    lang
```
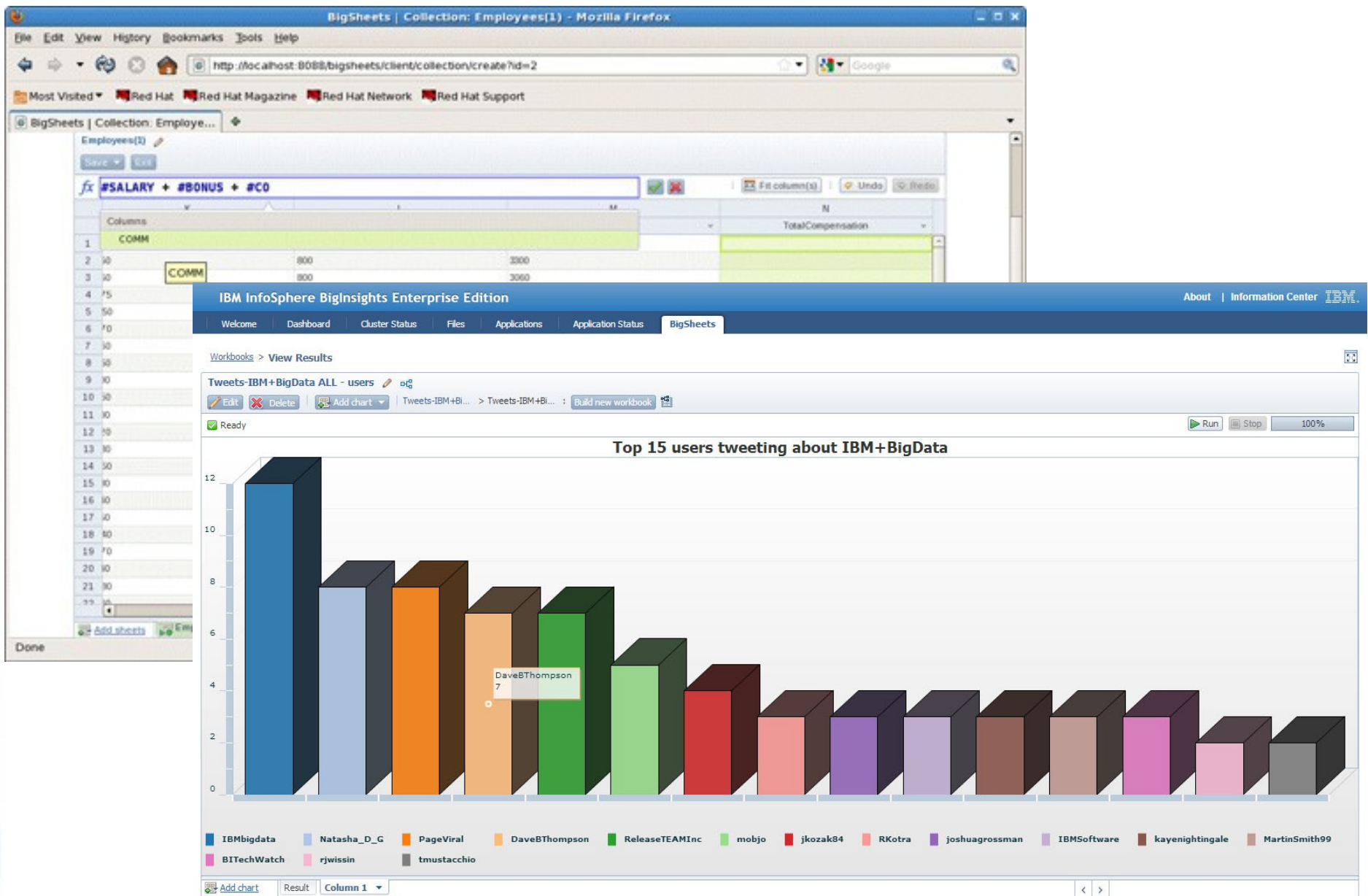
JAQL in HADOOP
BRIEF INTRODUCTION

# BigSQL

JDBC/ODBC client

LOAD from DB2, Teradata, PureData Systems for Analytics (Netezza)

## InfoSphere BigInsights

Big SQL server

JDBC/ODBC server

SQL engine

"Local" query execution

Single-node tasks

HBase storage driver

Hive storage driver

...

MapReduce parallelism

MapReduce tasks

MapReduce tasks

HBase storage driver

Hive storage driver

...

Hadoop data: distributed files, Hive data, HBase servers

# Big Spreadsheets

# SPSS

# SPSS



**IBM SPSS Modeler**

Modeler Client → Stream File → Modeler Server

**Analytic Catalyst**

Analytic Catalyst Tablet Client

Analytic Catalyst Browser Client

Big Data Request

SQL / UDF

**IBM SPSS Analytic Server**

Analytics

Hadoop Job

Relational Database

**IBM InfoSphere BigInsights & Other Hadoop Distributions**

# Business Intelligence

# Requirements Summary

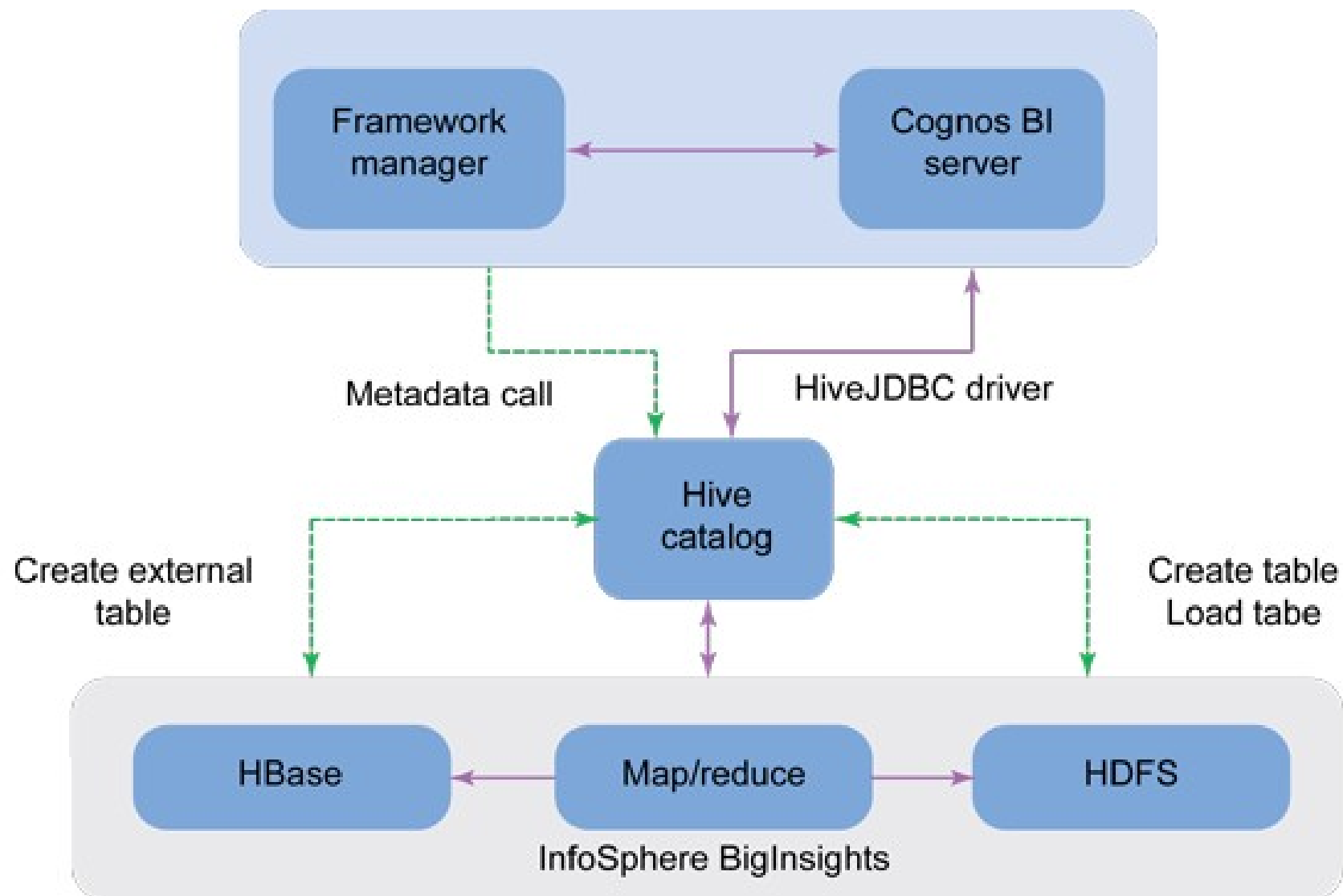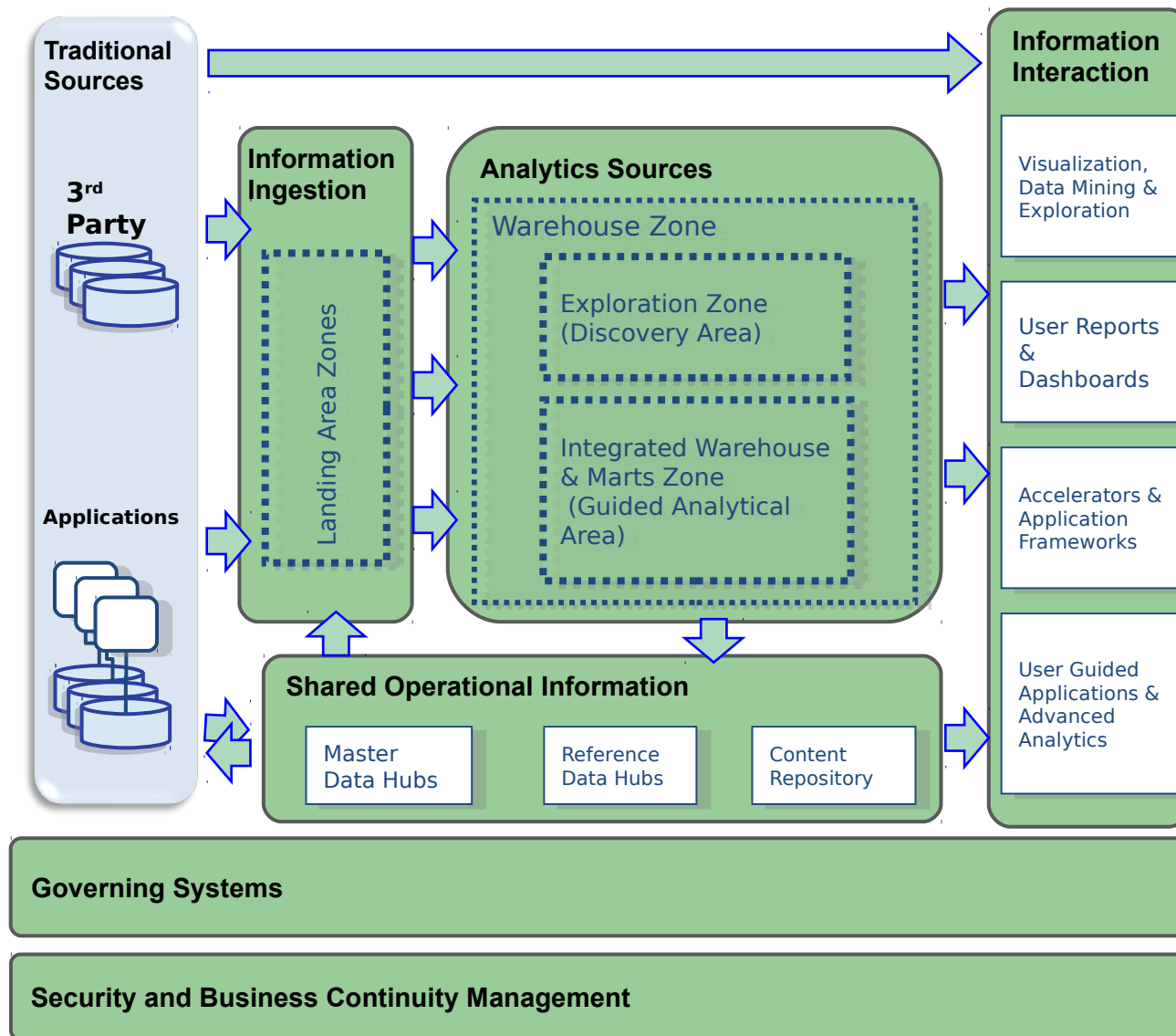- Fault tolerance
- Dynamic and elastic scale-in and out
- Processing data of all types
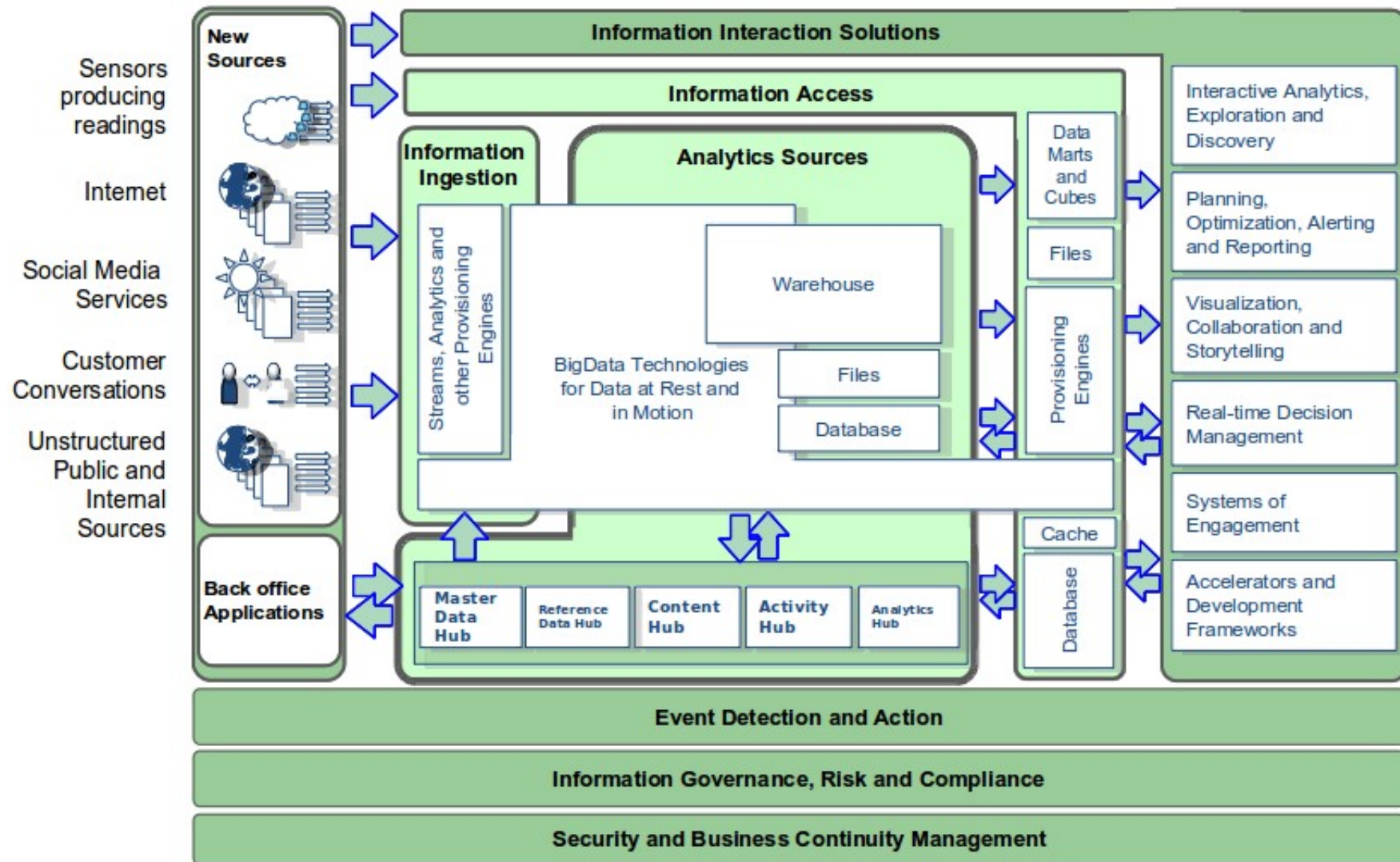- Use familiar ways of working with data

# Ingredients

- NoSQL DB
- Cloud
- Push-back from Business Applications

# IBM Reference Architecture (current)



**Traditional Sources**

**3rd Party**

**Applications**

**Information Ingestion**

Landing Area Zones

**Analytics Sources**

Warehouse Zone

Exploration Zone (Discovery Area)

Integrated Warehouse & Marts Zone (Guided Analytical Area)

**Shared Operational Information**

Master Data Hubs

Reference Data Hubs

Content Repository

**Information Interaction**

Visualization, Data Mining & Exploration

User Reports & Dashboards

Accelerators & Application Frameworks

User Guided Applications & Advanced Analytics

**Governing Systems**

**Security and Business Continuity Management**

# IBM Reference Architecture (transition)

# 15 minutes Discussion



Twitter: @romeokienzler