# Big Data MDX with Mondrian and Apache Kylin

Sébastien Jelsch

London, 7-11-2015

- OLAP-on-Hadoop with Apache Kylin

- Features

- Apache Kylin & Mondrian

- Conclusion & Discussion

- **OLAP-on-Hadoop with Apache Kylin**

- Features

- Apache Kylin & Mondrian

- Conclusion & Discussion

**Situation**

- More and more data becoming available on Hadoop
- Limitations in existing Business Intelligence Tools
  - Limited support for Hadoop
  - Data size growing exponentially
  - High latency of interactive queries
- Challenges to adapt Hadoop for interactive analysis
  - OLAP capability on Hadoop ecosystem not ready yet

**Goals**

- Full OLAP capability and advanced functionality
- Interactive analysis in subseconds
- ANSI SQL or MDX for analysts and engineers
- Seamless integration with BI Tools
- High concurrency with thousands of end users
- Distributed and scale out architecture for large data volume

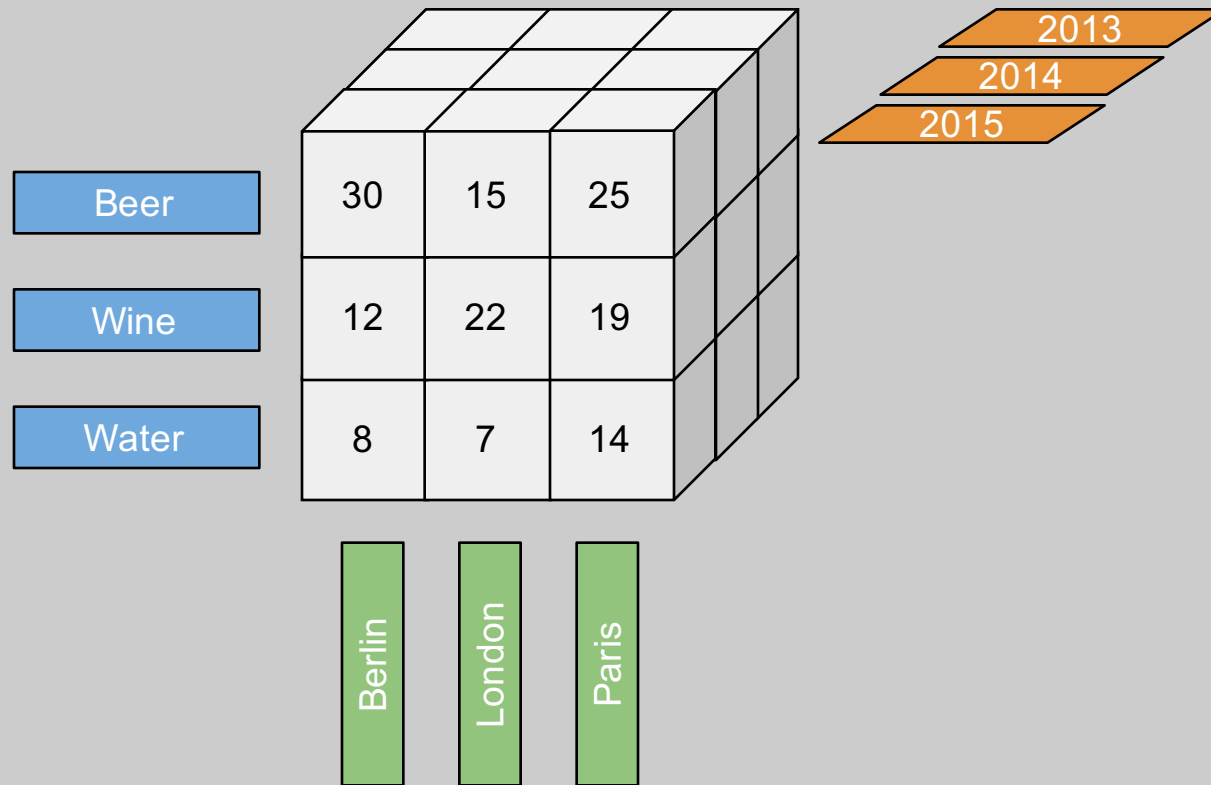**Solution:** Apache Kylin

**Extreme OLAP Engine for Big Data**

- Distributed Analytics Engine from eBay
- OLAP-on-Hadoop
- Provides SQL interface for multidimensional analysis
- Based on Hadoop ecosystem

Open Source on: 1. October 2014
Accepted into incubation: 25. November 2014
Current version: 1.1 (25. October 2015)

## OLAP Cube

# Agenda

- OLAP-on-Hadoop with Apache Kylin

- **Features**

- Apache Kylin & Mondrian

- Conclusion & Discussion

# Apache Kylin: Cube Designer

# Apache Kylin: Cube Designer

# Apache Kylin: Cube Designer

# Apache Kylin: Cube Designer

# Apache Kylin: Cube Designer

# Apache Kylin: Monitoring

- OLAP-on-Hadoop with Apache Kylin

- Features

- **Apache Kylin & Mondrian**

- Conclusion & Discussion

**SQL returns 2-dimensional result set**

For more dimensions SQL was not designed

**Wish**:

- Multidimensional  result set
- Consider  hierarchies  and levels in the data

⟹   Query Language: <u>MDX</u>

**Mondrian**

- OLAP Engine
- Transforms MDX queries into SQL
- Multidimensional representation of data
- Integrated into Saiku / Pentahos Business Analytics Platform
- Expandable through SQL dialects
  e.g. MySQL, Postgres, Hive, Impala, ...

inovex

OLAP Client

↕ MDX

Mondrian Schema → XML → Mondrian | Kylin Dialect

Measures

Dimensions

Hierarchies

Levels

Attributes

↕ JDBC

**Apache Kylin**
HBase, Cuboids ...

**Work done:**

- Kylin dialect created
- Optimized Kylins JDBC driver
- Bugs fixed to get Mondrian working with Kylin

**TBD:**

- Integrate Kylin dialect into Mondrians official code*
- Make every MDX query executable

**Successful tests**:**

- Current Saiku and Mondrian 4.4
- Current Saiku and Mondrian 3.x (not tested very well)

*  Pull Request: https://github.com/pentaho/mondrian/pull/480

** Github Project: https://github.com/mustangore/kylin-mondrian-interaction

```java
public class KylinMondrianOlap4J {
  public static void main(String[] args) throws ClassNotFoundException, SQLException {
    Class.forName("mondrian.olap4j.MondrianOlap4jDriver");

    Connection connection = DriverManager.getConnection(
        "jdbc:mondrian:"
      + "Jdbc=jdbc:kylin://{YOUR_URL}:7070/{YOUR_PROJECT_NAME}};"
      + "JdbcDrivers=org.apache.kylin.jdbc.Driver;"
      + "JdbcUser={YOUR_USER};"              // Default: admin
      + "JdbcPassword={YOUR_PASSWORD};"      // Default: KYLIN
      + "Catalog=file:/absolute/path/to/your/mondrianSchema.xml;");

    // We are dealing with an OLAP connection. we must unwrap it.
    final OlapConnection olapConnection = connection.unwrap(OlapConnection.class);

    // Prepare a statement.
    final OlapStatement olapStatement = olapConnection.createStatement();

    // We use the utility formatter.
    RectangularCellSetFormatter formatter = new RectangularCellSetFormatter(false);

    // Your MDX Statement
    String mdxStatement = "{YOUR_MDX_QUERY}";

    CellSet cellSet = olapStatement.executeOlapQuery(mdxStatement);

    // Print out.
    PrintWriter writer = new PrintWriter(System.out);
    formatter.format(cellSet, writer);
    writer.flush();
  }
}
```

# Agenda

inovex

- OLAP-on-Hadoop with Apache Kylin

- Features

- Apache Kylin & Mondrian

- **Conclusion & Discussion**

- Extremely fast and scalable OLAP Engine
- OLAP-on-Hadoop
- Depends on Apache Hadoop infrastructure
- MOLAP Cube
- Incremental refresh of cubes
- Integration into existing BI Tools
- MDX queries with Mondrian possible (ongoing work)
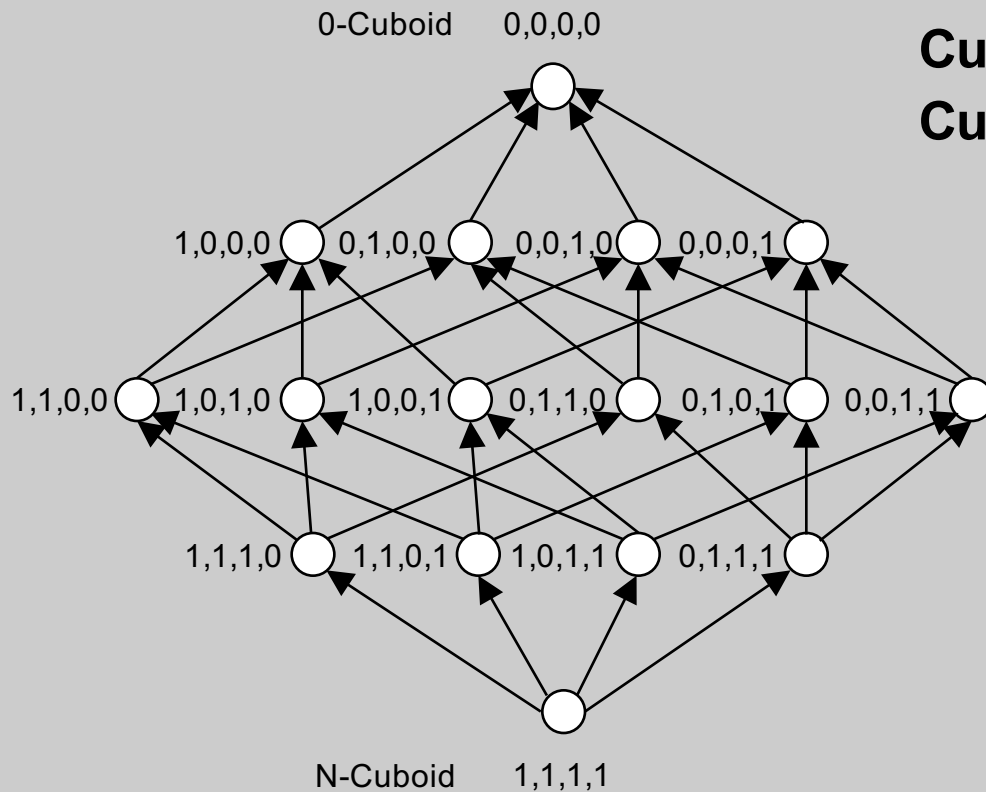
# Thank you for your attention

**Contact**

Sébastien Jelsch
Big Data Scientist


inovex GmbH
Office Karlsruhe
Ludwig-Erhard-Allee 6
76131 Karlsruhe


Tel: +49 176 - 45786280
E-Mail: sjelsch@inovex.de
Twitter: @inovexgmbh | @Mustangore

inovex

**Cube**: All combinations
**Cuboid**: One single combination

**Number cuboids growing exponentially**

**Problem**: Number of Cuboids grows exponentially

**Example**:

  Cube with 30 dimensions

  Number of Cuboids: $2^{30} > 1$ billion


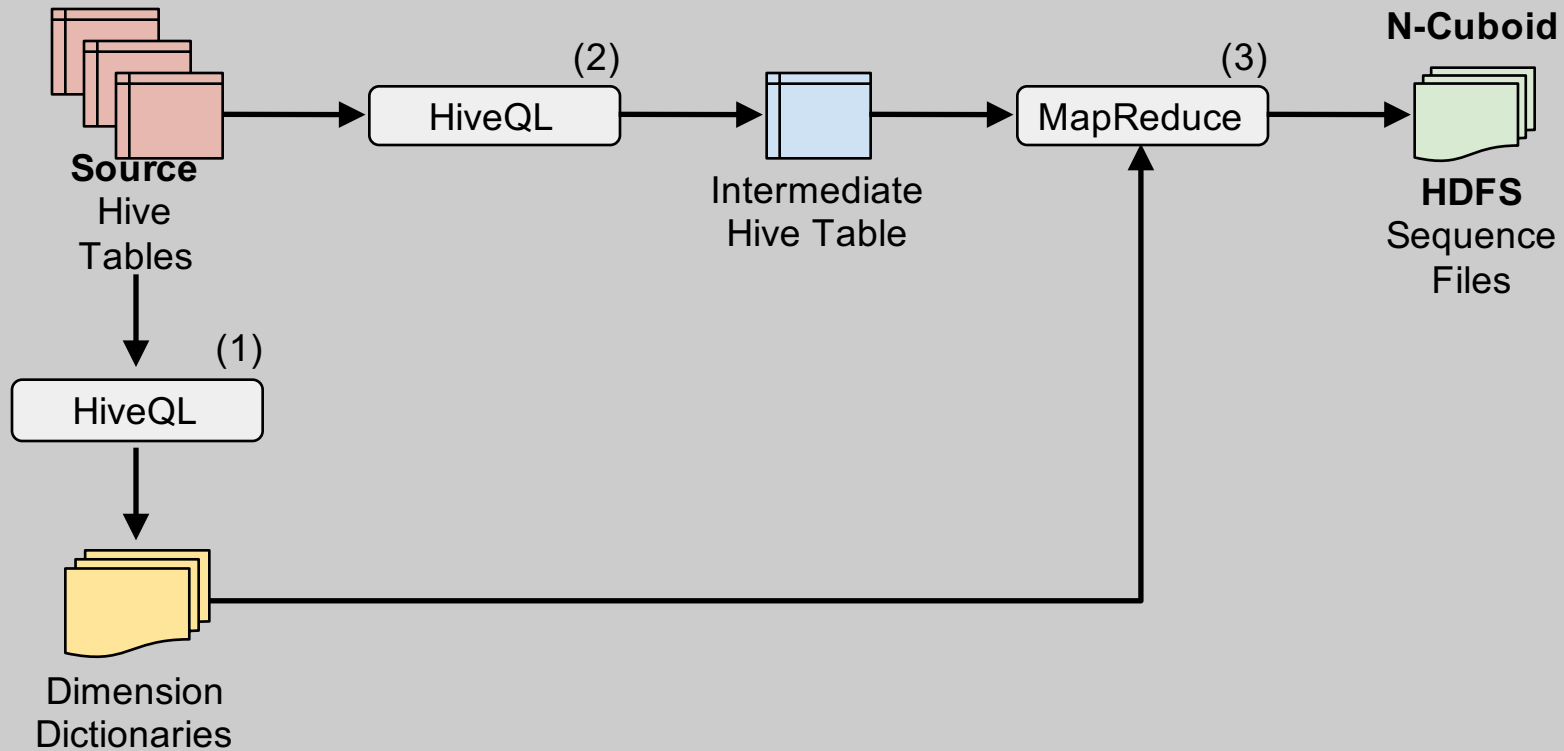**Solution**: <u>Partial Cube</u>

Classificate the OLAP Cube in Aggregation Groups

**Example**:

  30 dimensions splitted into 3 groups of 10 dimensions

  Number of Cuboids: $2^{10} + 2^{10} + 2^{10} = 3072 \ll 1$ billion

# Apache Kylin: Cube Build Process