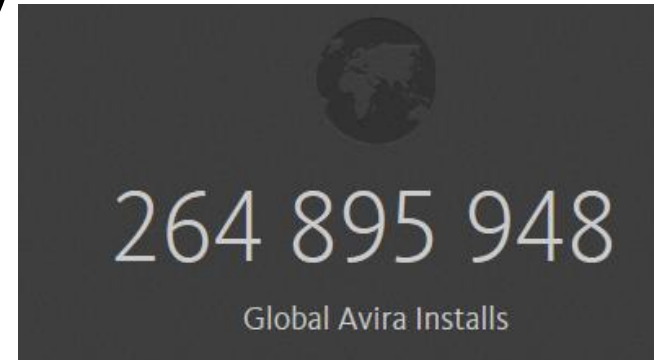# KNIME & Avira, or how I've learned to love Big Data

# Facts about Avira (AntiVir)

- 100 mio. customers

- „Extreme Reliability"

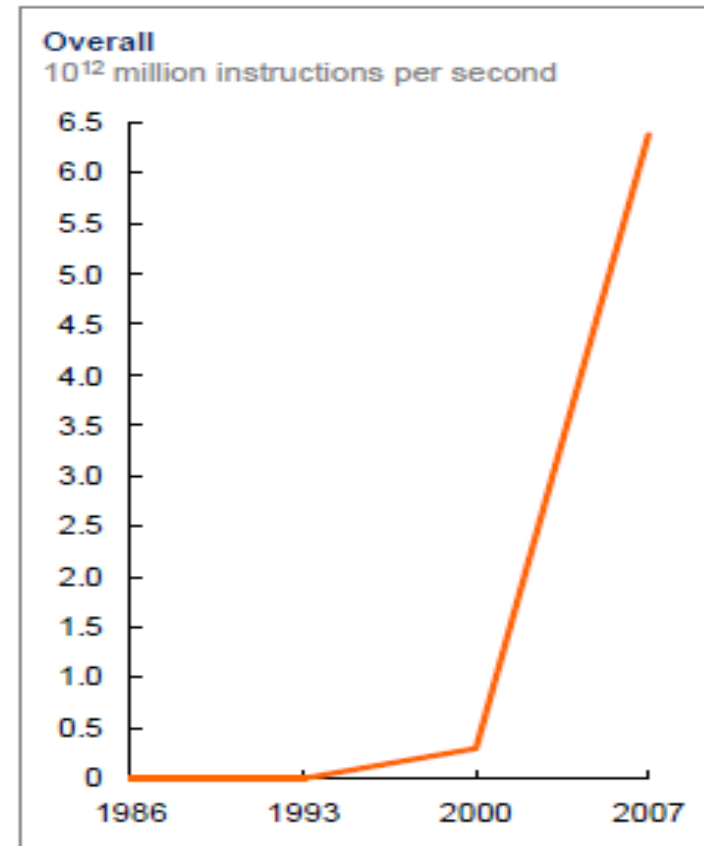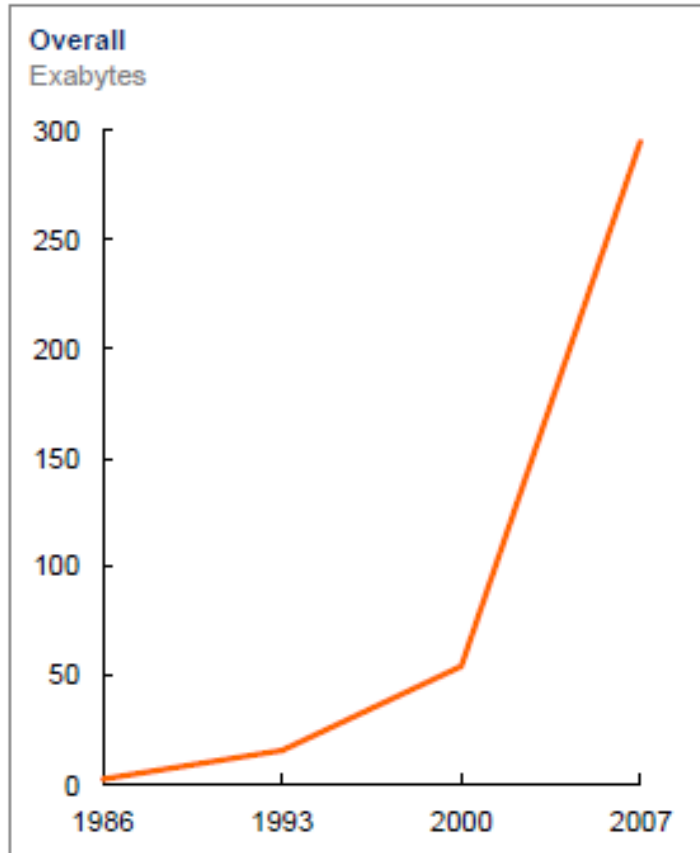- 500 employees (Tettnang, San Francisco, Kuala Lumpur, Bucharest, Amsterdam)

264 895 948

Global Avira Installs

13.02.2014, 06:58
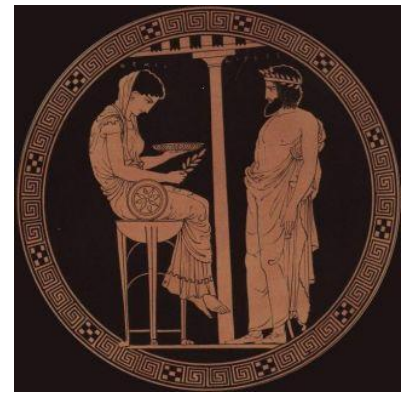
Company owner, Mr. Auerbach

# Big Data: Why did Avira decided to invest?

- Data storage has grown significantly after 2000
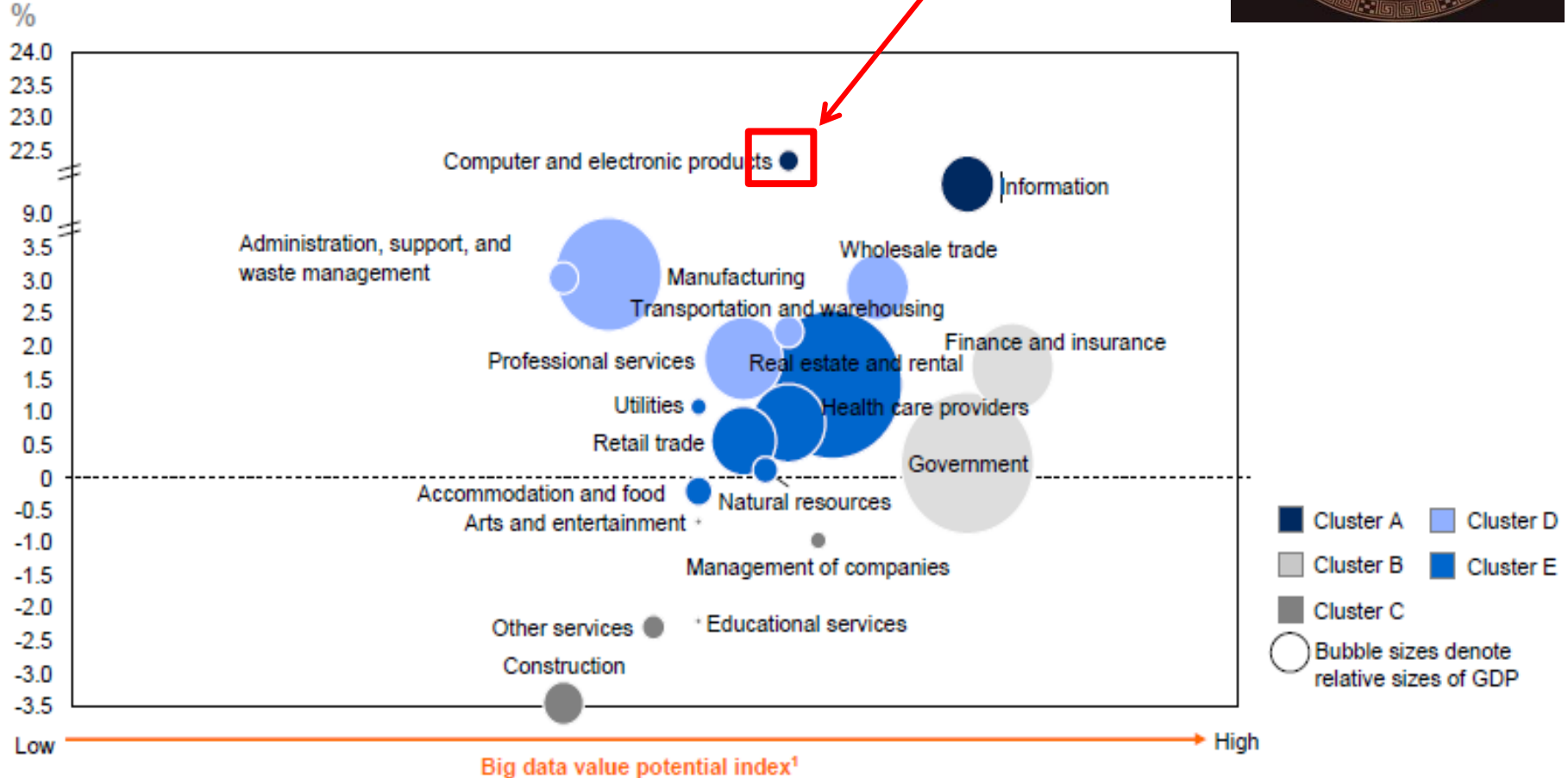
- Computation capacity has also risen sharply





SOURCE: Hilbert and López, "The world's technological capacity to store, communicate, and compute information," *Science*, 2011

# Reason 1: Because the McKinsey prophets said it may bring us more money ☺

Historical productivity growth in the United States, 2000–08
%



SOURCE: US Bureau of Labor Statistics; McKinsey Global Institute analysis

# Reason 2: Because we are sitting on the top of several dozens of terabytes of anonymized customer data.
## The only challenge was to link this data and converge towards a single customer profile

Documented Datasources, unique ID & data volume

Protection Cloud uID = rndsnr; Volume = 10gb/day

LogBox (Product Download & VDF Update): uID = rndsnr & license no; Volume = 30 gb/day

Global Mailing System: uID = license no.; Volume = 7,5 mio Emails/month

eshop.avira.com: uID = license no.; Volume = 10 gb/day

Licensing system: uID = license no.; Volume = registration data 100 mio users

Cleverbridge Shopping Cart: uID = license no.; Volume = 10gb/day

Website Site Catalyst: Uid not yet defined; Volume = 65 mio page loads & clicks/month

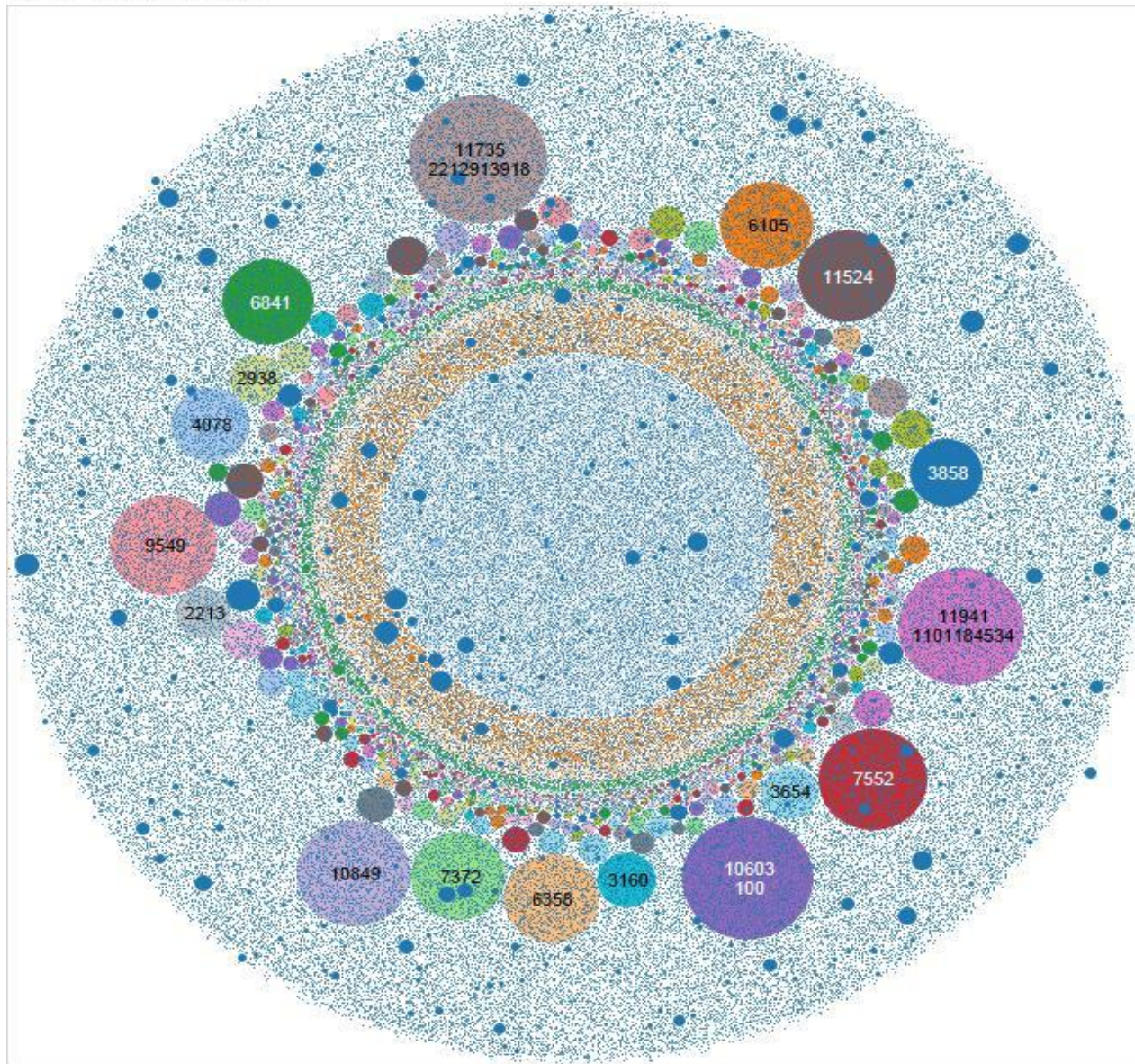Notifier: Uid not yet defined; Volume = 20 mio impressions/day

IPM: Uid not yet defined = 12,5 mio impressions/day

# Reason 3: Because we want to perform customer profilling & next best offer marketing to increase our margin

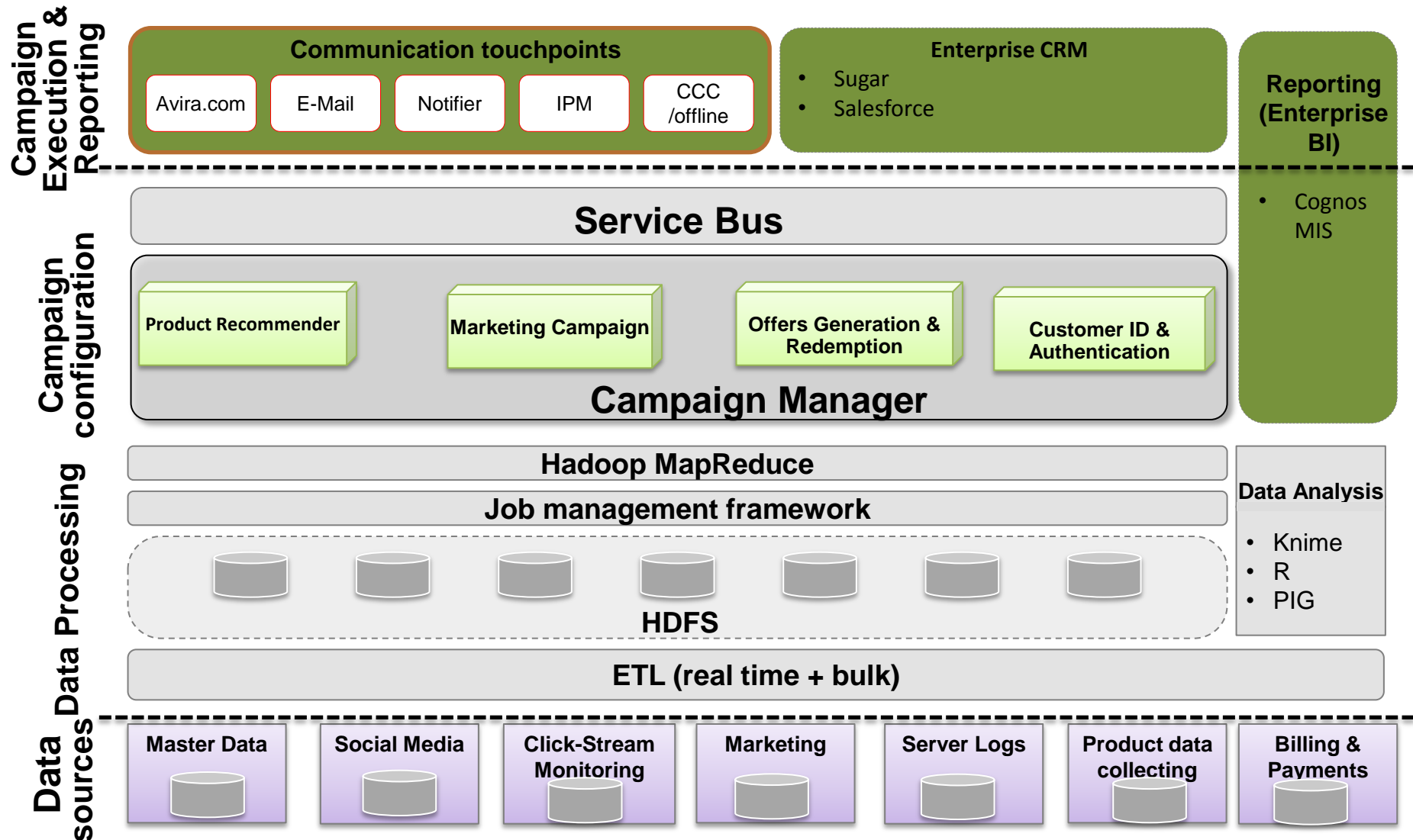| Business objectives | Technology requirements |
|---|---|
| <ul><li>Deliver the optimal price, features & messaging to each customer in order to maximize EBIT from sales of new and renewal licenses</li><li>Protect long-term margins by making each offer timely and unrepeatable (e.g., unique to a specific customer, product, event)</li><li>Learn the underlying mechanics of features and price-elasticity on the level of customer cells</li><li>Enable business to optimize campaign portfolio (i.e., über-algorithm traffics campaigns)</li><li>A/B testing in all customer touch-points</li></ul> | <ul><li>Ability to link a specific offer in a specific touch-point to specific customer.</li><li>Machine-learning over all design dimensions to continually improve performance of the application</li><li>Causal reporting to maximize the learning effect in the organization from algorithmic approaches to marketing and automation</li><li>Create customer, product and behavioral tables from Avira's raw data within the dev. environment</li><li>Create machine-learning algorithms optimizing the offer (price, features) per user-session</li><li>Implement the services on development platform and place in listening mode to train</li><li>Setup of the Hadoop framework (HDFS & MapReduce) & Couchbase, KNIME & Impala</li></ul> |

# Reason 4: Because it's fun (isn't this beautiful?)



Prem_non_germany

Counts of rndsnr and license. Color shows details about counts of rndsnr. Size shows sum of total updates. The marks are labeled by counts of rndsnr and license.

# Our business/architectural vision



**Campaign Execution & Reporting**

Communication touchpoints
- Avira.com
- E-Mail
- Notifier
- IPM
- CCC /offline

Enterprise CRM
- Sugar
- Salesforce

Reporting (Enterprise BI)
- Cognos MIS

**Campaign configuration**

Service Bus

Campaign Manager
- Product Recommender
- Marketing Campaign
- Offers Generation & Redemption
- Customer ID & Authentication

**Data Processing**

Hadoop MapReduce

Job management framework

HDFS

ETL (real time + bulk)

Data Analysis
- Knime
- R
- PIG

**Data sources**

- Master Data
- Social Media
- Click-Stream Monitoring
- Marketing
- Server Logs
- Product data collecting
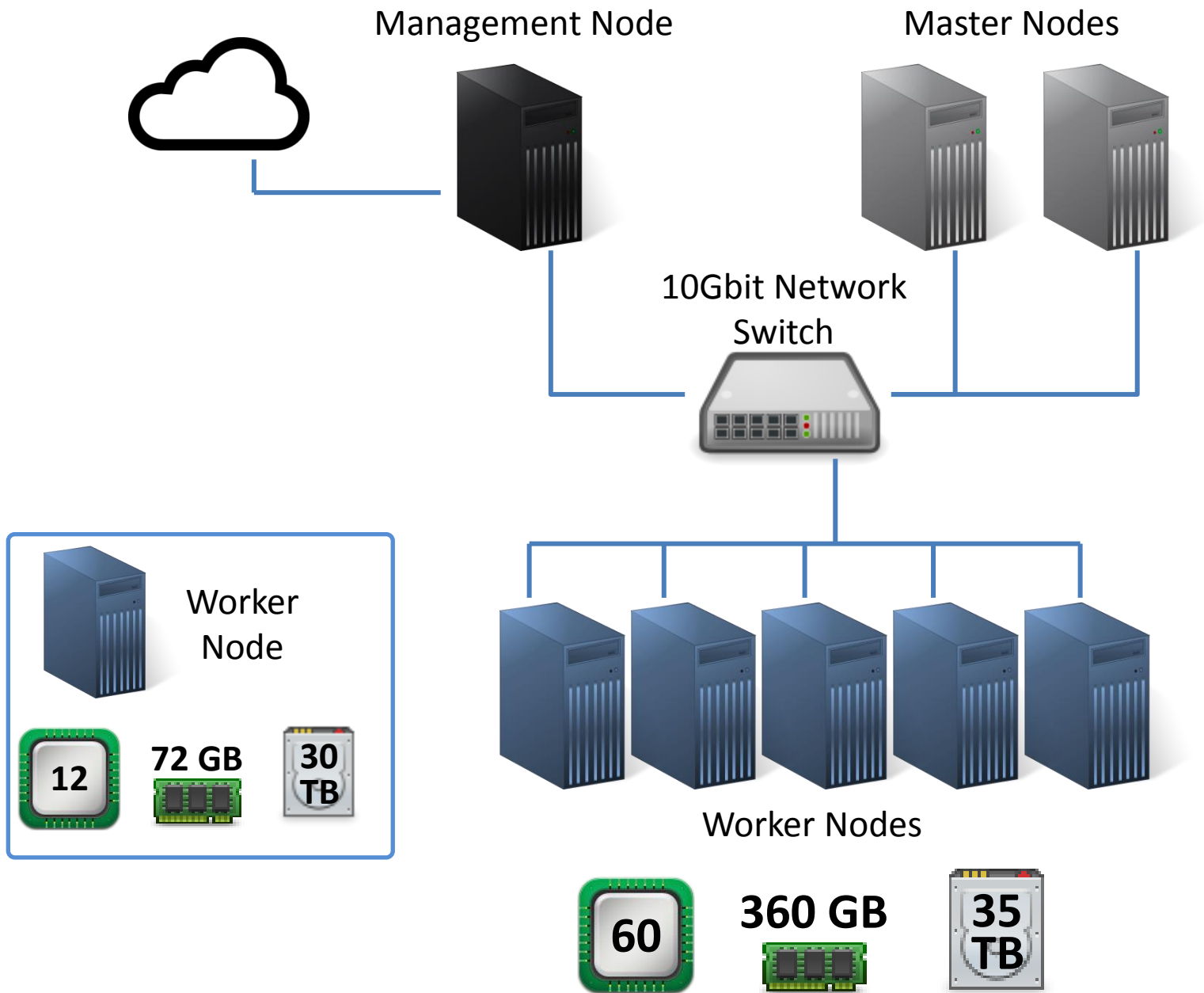- Billing & Payments

# Our Daily Data

- Website logs: 5.000.000 lines

- Installation logs: 2.200.000 lines

- InProduct Messaging: 43.000.000 lines

- Download/Updater logs: 2.000.000.000 lines

**60GB of daily compressed data**

or around 2TB of monthly data

# Finding the Right Tool

- Our data is:
  - Unstructured; Messy
  - Coming from all kinds of sources:

    log files, log tables, relational databases

- We want to:
  - Gather and store historical data
  - Process huge amounts of it
  - Support both batch and real-time operations

Management Node

Master Nodes

10Gbit Network Switch

Worker Node

12 | 72 GB | 30 TB

Worker Nodes

60 | 360 GB | 35 TB

# cloudera®
## Ask Bigger Questions

## CDH

| BATCH PROCESSING (MapReduce, Hive, Pig) | ANALYTIC SQL (Impala) | SEARCH ENGINE (Cloudera Search) | MACHINE LEARNING (Spark, MapReduce, Mahout) | STREAM PROCESSING (Spark) | 3RD PARTY APPS (Partners) |

### WORKLOAD MANAGEMENT (YARN)

## STORAGE FOR ANY TYPE OF DATA
### UNIFIED, ELASTIC, RESILIENT, SECURE (Sentry)

| Filesystem (HDFS) | Online NoSQL (HBase) |

### DATA INTEGRATION (Sqoop, Flume, NFS)

# Using The Right Tool

# MapReduce

- Full control over how the data is processed
- Works on structured and unstructured data
- Good for very complex business logic

- You have to write Java code
- Restricted to the MapReduce programming model
- Some things are difficult to code (JOINs, custom sorting)

# Hive

- You write SQL-like queries

- Great for ad-hoc queries, data exploration

- Very fast development
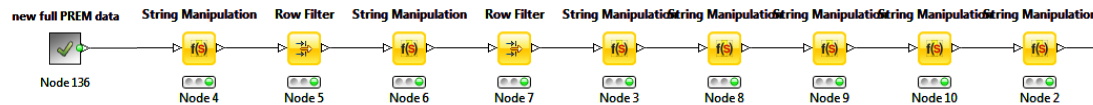

- Works only on structured data

-  Gets ugly if the business logic is complex

# Impala

- Like Hive but A LOT faster
- Runs directly in memory
- Delivers almost real-time results

- Limited to in-memory processing
- Unreliable

# KNIME and Hadoop at Avira

# Knime for Understanding the Data



MapReduce            Hive / Impala            KNIME

# Some Big Data Practicalities

## Use Crisp-DM !



**Big Data Methods Never Needed**

Project Understanding
- What exactly is the problem, the expected benefit?
- What should a solution look like?
- What is known about the domain

Data Understanding
- What data do we have available?
- Is the data relevant to the problem?
- Is it valid? Does it reflect our expectations?

**Decision:   Would Big Data methods add value?**

Does data suit problem? — no → Cancel Project
partially

**Big Data Methods Possibly useful**

Data Preparation
- Which data should we concentrate on?
- How is the data best transformed for modeling?
- How may we increase the data quality?

Modeling
- What kind of model architecture suits the problem best?
- What is the best technique/method to get the model?
- How well does the model perform technically?

Technical Quality Improvable? — Likely
Unlikely

Evaluation
- How good is the model in terms of project requirements?
- What have we learned from the project?

Revise Objective

**Decision:  Would Big Data methods add value?**

Business Objective Achieved? — no → Close Project
partially
Success

**Big Data Methods Possibly useful**

Deployment
- How is the model best deployed?
- How do we know that the model is still valid?

Guide to Intelligent Data Analysis, Berthold, 2010

# Why did we decided to go with KNIME?

- The dark side of the moon: a typical symptom for home-grown business applications.



KNIME helped us to tie a „knot" for the multiple uncorrelated data points and create customer 360 tables

# Why did we decided to go with KNIME?

It helped us move from code based data mining towards workflow based analytics; Analytics for everyone, easy to explain to all management/company levels

# Some results: using KNIME & Tableau we've managed to perform forensycs and license outlier analysis

**Prem_non_germany**



Our notorious „hacker friend" working for the overall good of torrent visitors

This is me using Avira in Sep 2013

Counts of rndsnr and license.  Color shows details about counts of rndsnr.  Size shows sum of total updates .  The marks are labeled by counts of rndsnr and license.

# Running k-means in KNIME to identify relevant clusters for Germany by looking at their antivirus software update behaviour

Graph of clusters by time of day % usage:

# Running k-means in KNIME to identify relevant clusters for Germany by looking at their antivirus software update behaviour
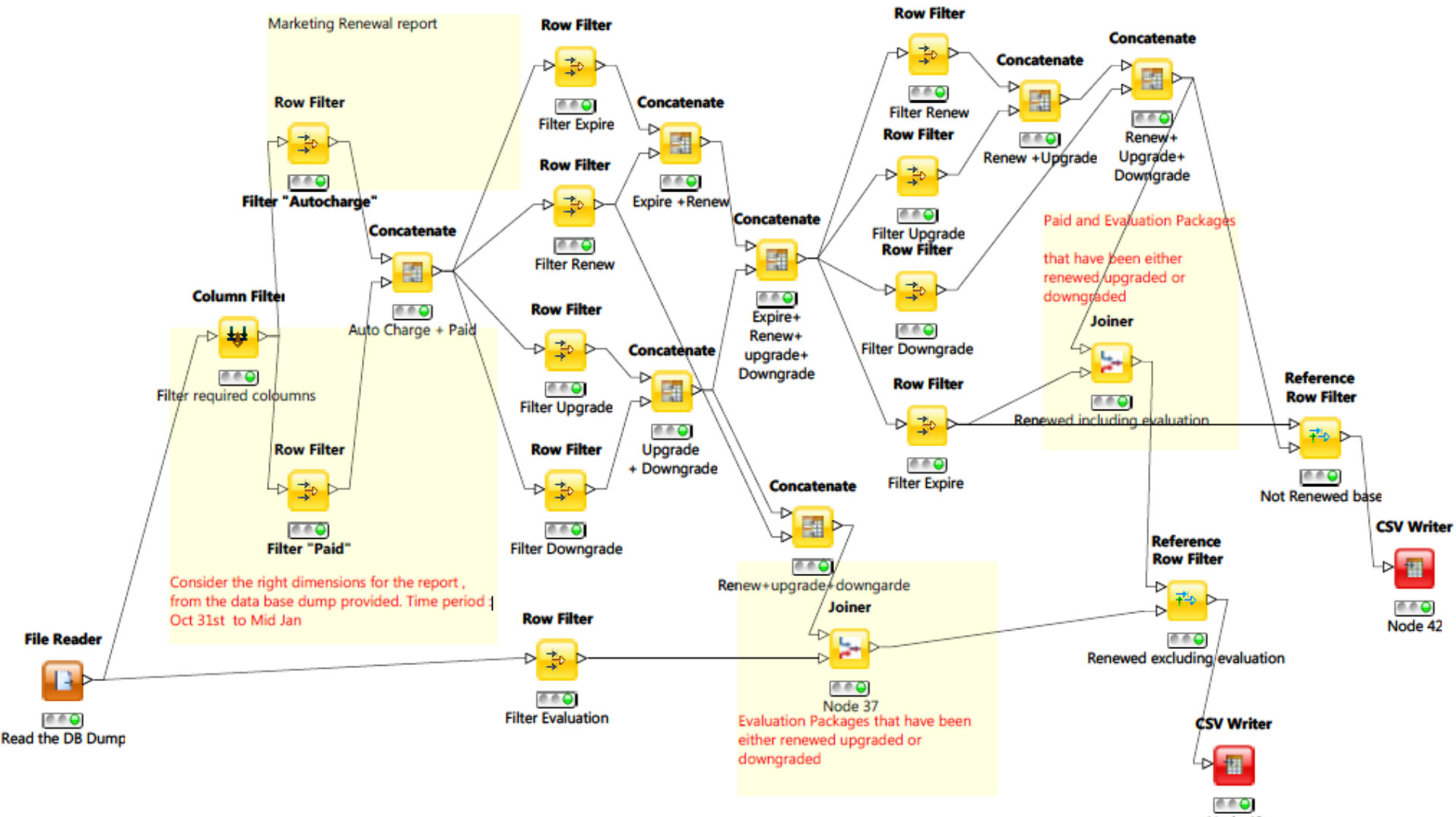
Graph of clusters by daily % usage

# Using KNIME to identify the real License Renewal pattern of our customers

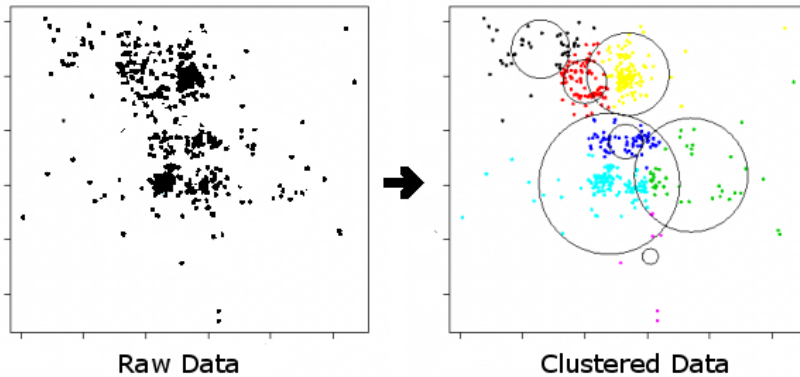The timing of renewing an Avira license in %

"Reason 6" = 12%

# Using KNIME to do standardized reporting of our license renewal metrics

# Next steps

- Identify unknown groups of customers by allowing the machine to find patterns in data for creating special association rules/product recommendations & next best offer; test & train in KNIME, real-time model execution in Couchbase;



Raw Data → Clustered Data

NBO via Email

NBO via In App Messaging

# Thank you.



radu.pastia@avira.com/ phil.winters@knime.com/ florin.veringa@avira.com