



Pentaho Data Integration 4 and MySQL

Matt Casters:

Pentaho's Chief Data Integration

Kettle Project Founder

MySQL User Conference, Tuesday April 13th, 2010

Agenda

- ***Pentaho: an introduction***
- Pentaho Data Integration
- Version 4: New features
- MySQL support in PDI
- Q&A

Pentaho Introduction

- Commercial open source alternative for business intelligence (BI)
 - Founded in 2004: Pioneer in commercial open source BI
 - Large referenceable customer base, wide range of BI/DW deployments
- Management - proven BI and open source veterans
 - Business Objects, Cognos, Hyperion, JBoss, Oracle, Red Hat, SAS, SugarCRM
- Board of Directors - deep expertise and proven success in open source
 - Bob Bearden - Executive chairman of the board (former SpringSource)
 - Larry Augustin - founder, VA Software, helped coin the phrase “open source”
 - Zack Urlocker - VP of Products, MySQL/Oracle
 - Benchmark Capital, Index Ventures, New Enterprise Associates
- Widely recognized as the leader in open source BI



Pentaho Introduction

- Complete Business Intelligence Suite
 - End-to-end coverage of all BI needs
 - Standards-based, modular, standalone or embeddable platform
- Open Source Licensing
 - Lower software acquisition costs
 - Lower Total Cost of Ownership (TCO)
- Enterprise Development Methodology
 - Transparent, detailed roadmap
 - Product roadmap and contributions managed by Pentaho
 - Core developers are Pentaho employees
 - Extensive QA
- Expert Services
 - Comprehensive Training, Consulting, Enterprise service offerings
 - Delivered by the Experts

Pentaho Introduction – Enterprise Edition

Pentaho BI Suite Enterprise Edition



Professional
Support



Software
Maintenance



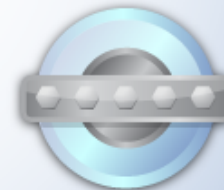
Enhanced
Functionality



Certified
Software



Product
Expertise



Software
Assurance

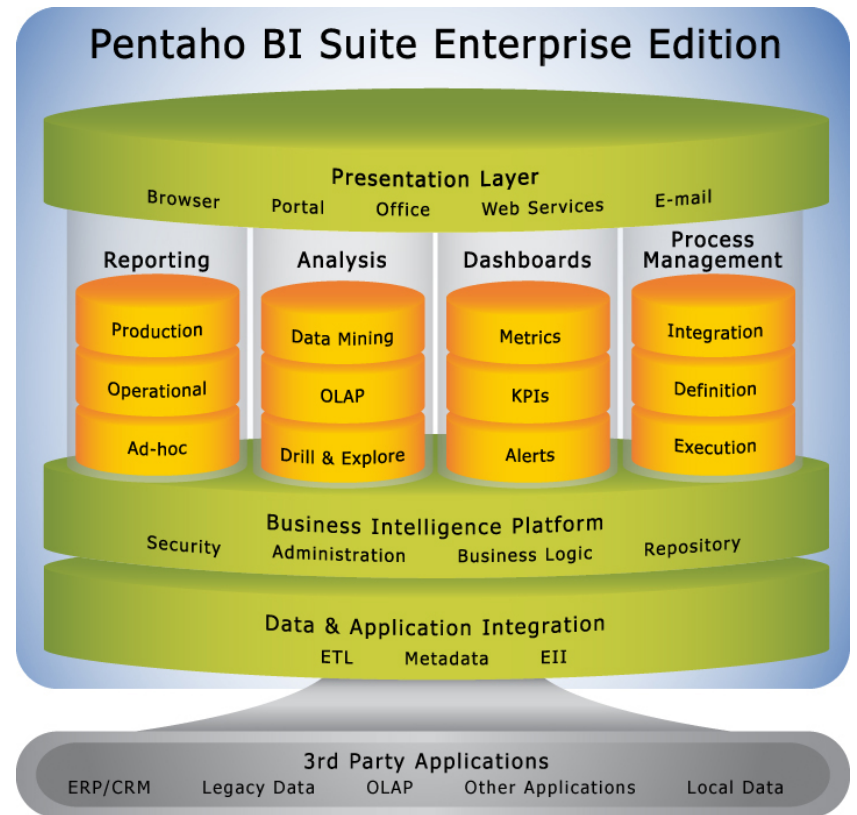
Pentaho Introduction – Deployments

- Wide range of deployments
 - Reporting
 - Data Integration / ETL
 - Dashboards
 - Full BI Suite
- Thousands of users
 - 3,000 on a single server
- Large data volumes
 - Half a terabyte of live interactive OLAP data
 - ETL loading 300K rows/second
- Sophisticated applications
 - Hundreds of dimensions
- Small deployments as well
 - 20 users, MS Access databases





Pentaho Introduction – Technology

- Componentized and modular
- Service-implemented architecture
 - Built “from the ground up” as a set of services
 - Exposed via AJAX and Web Services
- 100% Java EE server side
 - Scalable, standards-based
- Web-based, thin-client end user interfaces
- Graphical design interfaces
- Embedded process workflow engine



Pentaho Introduction – Reporting

- Access and format data from disparate sources
 - RDBMS, XML, OLAP
- Produce in popular formats
 - 
- Multiple report types
 - Operational
 - Analytical
 - Financial
 - Parameterized
- Go directly against data sources or Pentaho's centralized metadata layer



Steel Wheels
500 International Speedway, Daytona, FL 32114
(123) 456-7890 <http://www.steelwheels.com>
Run Date: 2/20/06 1:24 PM

Steel Wheels
500 International Speedway, Daytona, FL 32114
(123) 456-7890 <http://www.steelwheels.com>
Run Date: 2/20/06 1:24 PM

TO: Reims Collectables
59 rue de l'Abbaye, null
Reims, null 51100 France

Attn: Paul Henriot
Sales Rep: 1337
Terms: Net 30 days

INVOICE

Invoice #: 10121
Account Number: 353
Date: May 07, 2003

SKU	Product Description	Price/Unit	Qty Ordered	Total Price
S20_4713	2002 Yamaha YZ450 M1	\$74.85	44	\$3,293.40
S24_2360	1982 Ducati 900 Monster	\$76.88	32	\$2,460.16
S32_4485	1974 Ducati 350 MB3 Desmo	\$88.74	25	\$2,168.50
S12_2823	2002 Suzuki KX80			
S10_1678	1999 Harley Davidson Ultra			

Send Payment and Remittance Slip to:

Steel Wheels
500 International Speedway
Daytona, FL 32114

Thank you for your

REMITTANCE

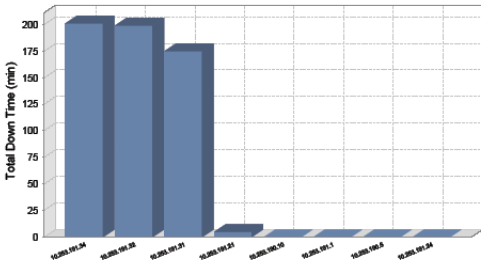
Reims Collectables
59 rue de l'Abbaye, null
Reims, null 51100 France

TopN Least Available
TopN least available managed assets

Summary Report Period: Feb 01, 2006 12:00 AM to Feb 10, 2006 5:46 PM

Number of Elements Included: 10
Availability Target: 99.0000%
Types of Outages Contributing to Downtime: Unplanned
Business Day Policy: Full Day including Weekend

Availability Graph



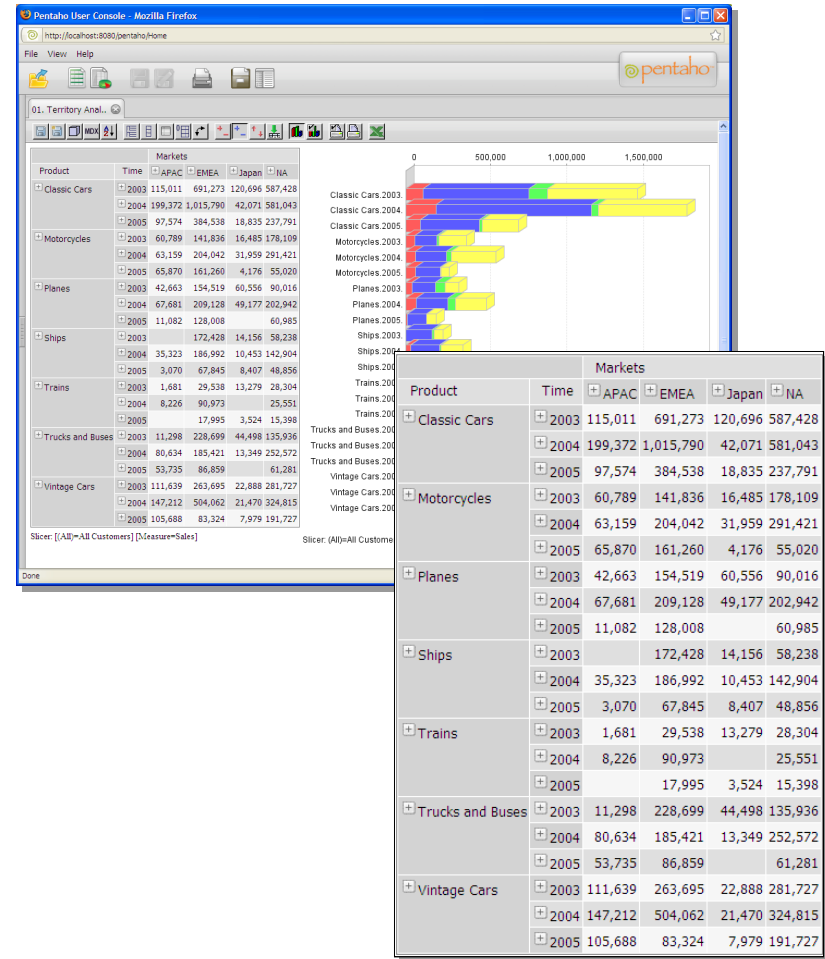
Report Details

Name	IP Address	Type	Total Downtime	No. Of Outages	Availability %
10.253.191.34	10.253.191.34	Windows Host	03h 20m 50s	11	98.568%
10.253.191.32	10.253.191.32	Windows Host	02h 18m 39s	12	98.584%
10.253.191.31	10.253.191.31	Windows Host	02h 55m 05s	15	98.732%
10.253.191.21	10.253.191.21	Windows Host	00h 05m 00s	1	99.964%
10.253.190.10	10.253.190.10	RS-3000	00h 00m 00s	0	100.000%
10.253.191.1	10.253.191.1	CaseWEEK3002	00h 00m 00s	0	100.000%
10.253.190.5	10.253.190.5	RS-3000	00h 00m 00s	0	100.000%
10.253.191.24	10.253.191.24	Windows Host	00h 00m 00s	0	100.000%

Report Generated on Feb 12, 2006 10:54:53 AM Page 1 of 1

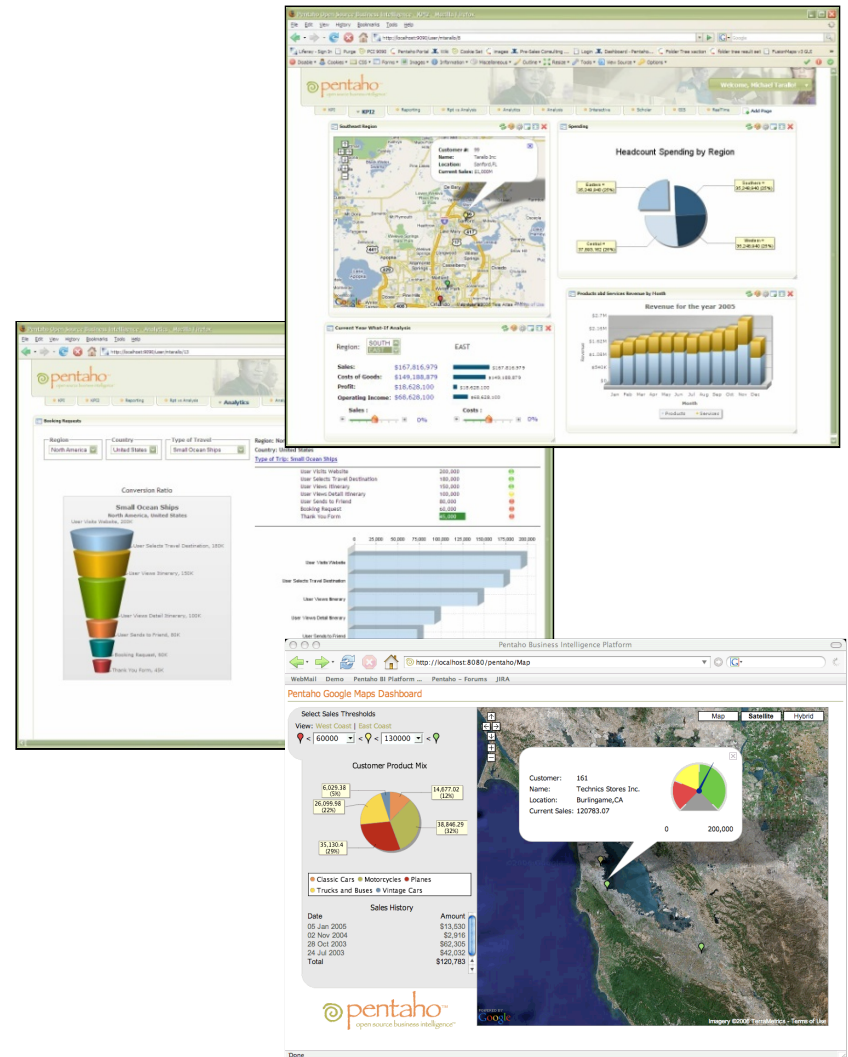
Pentaho Introduction – Analysis

- Navigate and explore
 - Ad hoc, interactive analysis
 - Drill into further detail
 - Select specific members for analysis
- View data “dimensionally”
 - i.e. Sales by region, by channel, by time period
- ROLAP architecture
 - Works with all popular open source and proprietary DBs
 - No intermediate storage
 - Aggregate table “aware” for faster analytic queries
- Design tools to build OLAP schemas and improve query performance



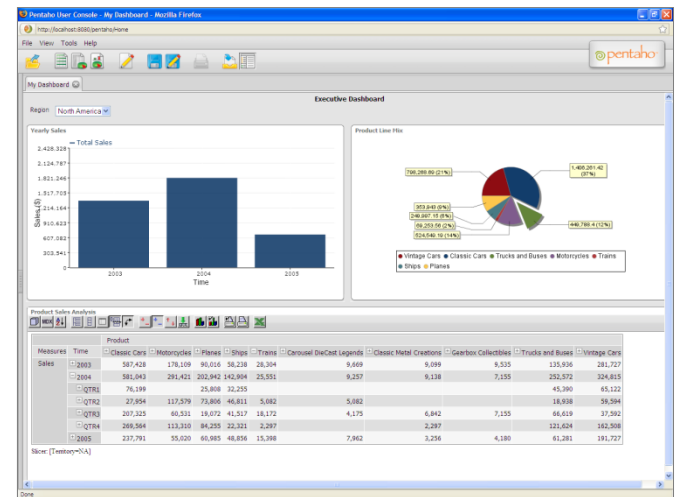
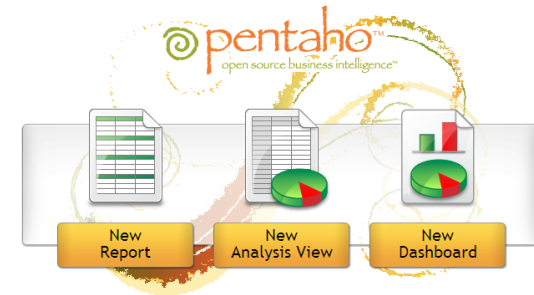
Pentaho Introduction – Dashboards

- Gain visibility into your organization's key performance indicators (KPIs)
 - Monitor top-level performance and drill into supporting detail
 - Illuminate metrics for quick insight into business activities
 - Track exceptions and receive alerts
- Leverage the full Pentaho BI Suite
 - Comprehensive auditing of user activity, performance and data access
 - Context-sensitive drilling to reports and analysis views
 - Integrated security, scheduling, alerting, portal integration
- Integrate with 3rd-party and custom applications



Pentaho Introduction – Dashboard Designer

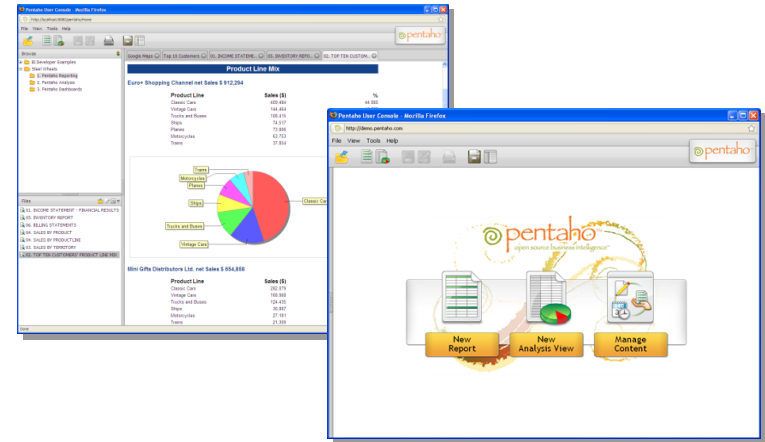
- Web-based end user dashboard creation
 - From Pentaho User Console
 - “Zero training”
- Template and theme-based creation
- Incorporate reports, analysis views, Adobe Flash-based charts and other Pentaho content
- Create new charts and interactive data grids from scratch
 - Pentaho metadata – no SQL required
- Filter controls



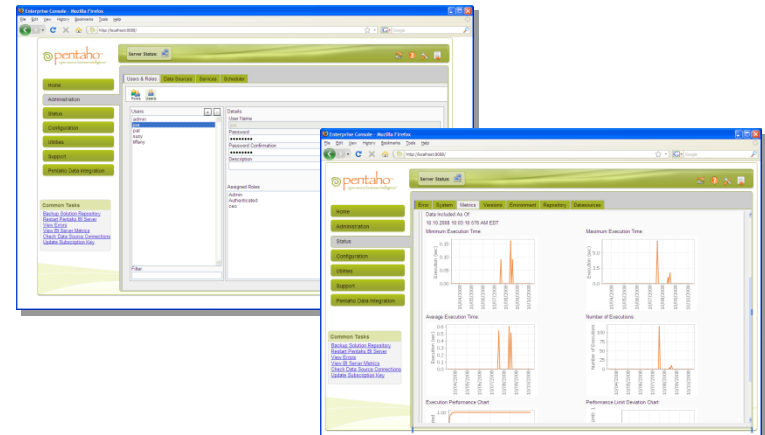
Pentaho Introduction – BI Platform

- Provides critical services for end users
 - Easy access to business information
 - Intuitive scheduling
 - Delivery over the web or via email
 - Alerting and notification
- Provides critical services for administrators
 - Centralized thin-client administration
 - Data source and security management
 - Auditing and Performance monitoring
 - Enterprise security integration
 - Definition and execution of business rules
 - Integration points with 3rd party applications

Pentaho User Console

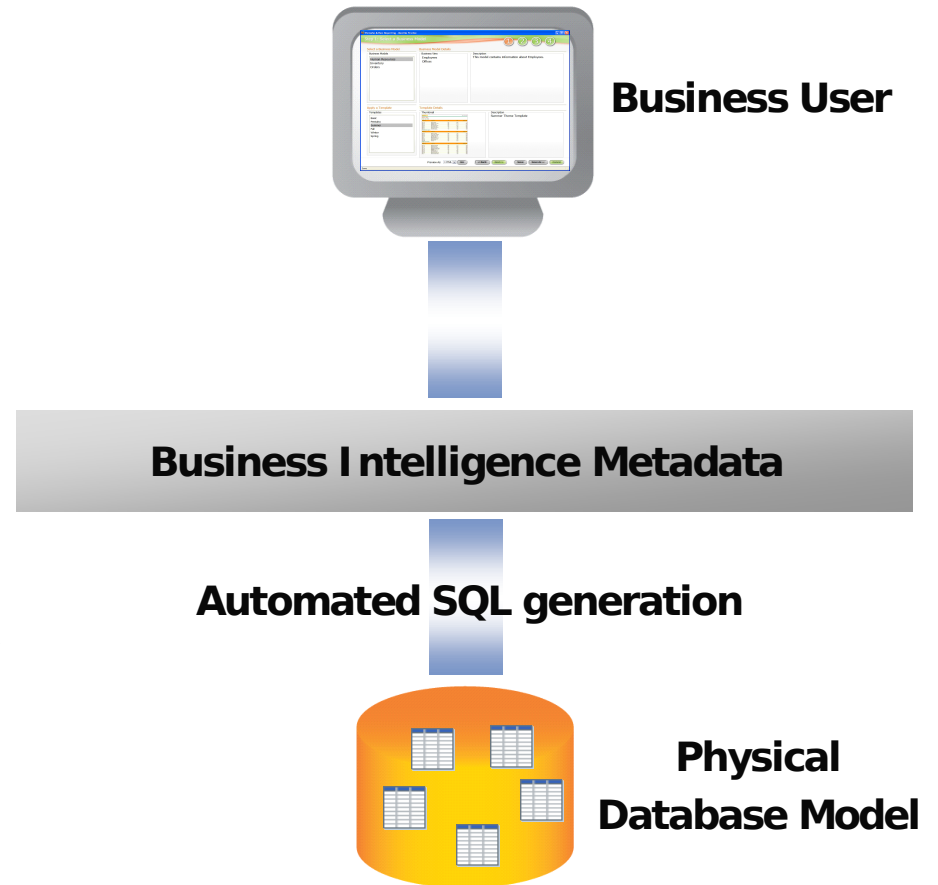


Pentaho Enterprise Console



Pentaho Introduction – Metadata

- Provides an abstraction layer between source systems and business user concepts
- Graphical design environment for defining metadata model
- Data presented to business users in business terms
- Allows business users to create their own ad hoc reports based on centralized business rules, without any technical skills or knowledge of SQL
- Changes to physical database do not impact reports or analytic views



Pentaho Introduction – Data Mining

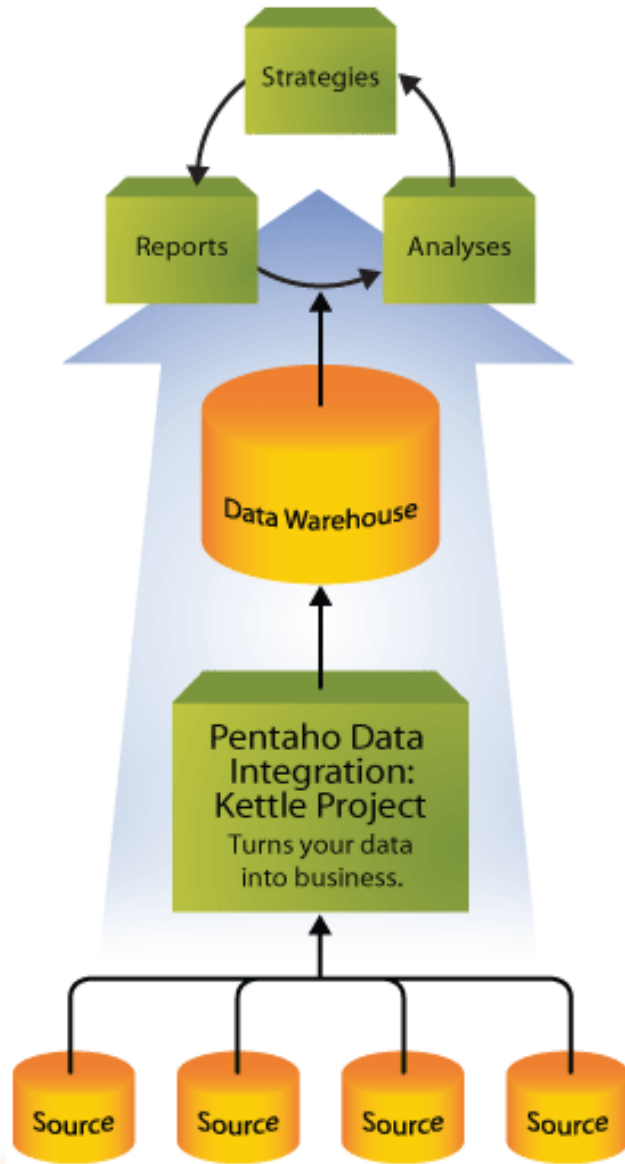
- Take BI to the next level with predictive analytics
- Gain insight into hidden patterns and relationships
- Discover indicators of future performance
- Exploit correlations to improve organizational performance
- Embed recommendations in reports, dashboards, or custom applications



Agenda

- Pentaho: an introduction
- ***Pentaho Data Integration***
- Version 4: New features
- MySQL support in PDI
- Q&A

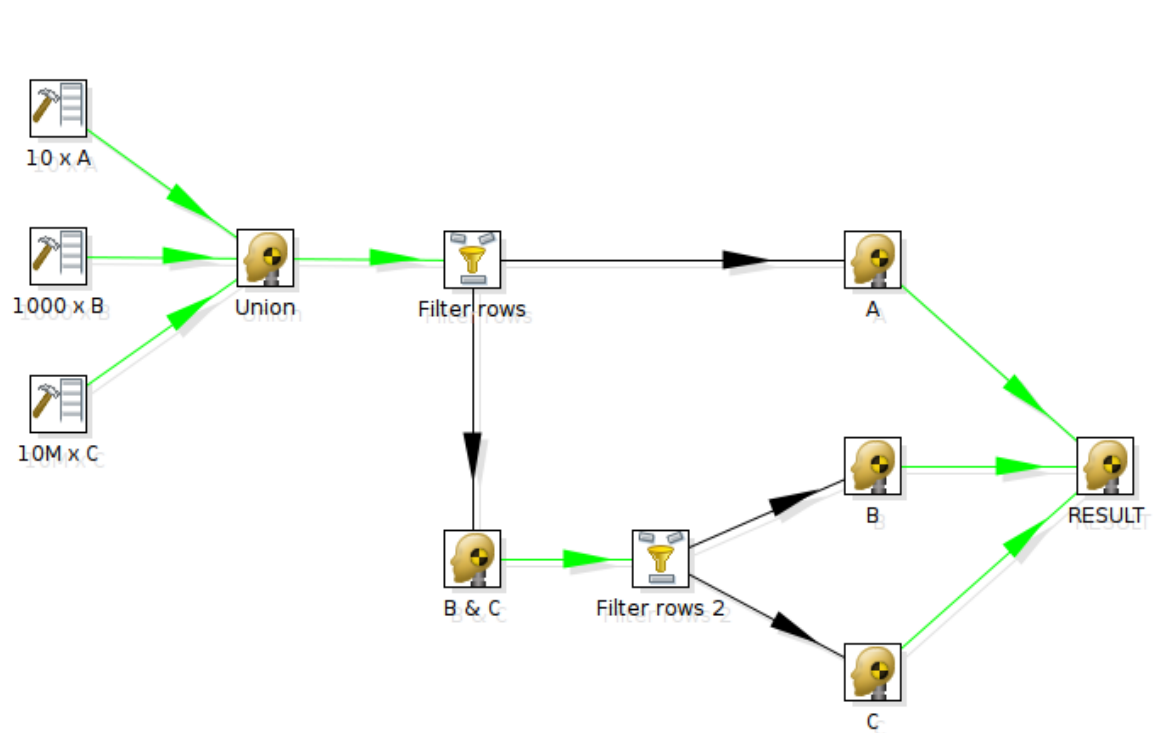
Pentaho Data Integration for BI



Business Intelligence!
That's what we do.

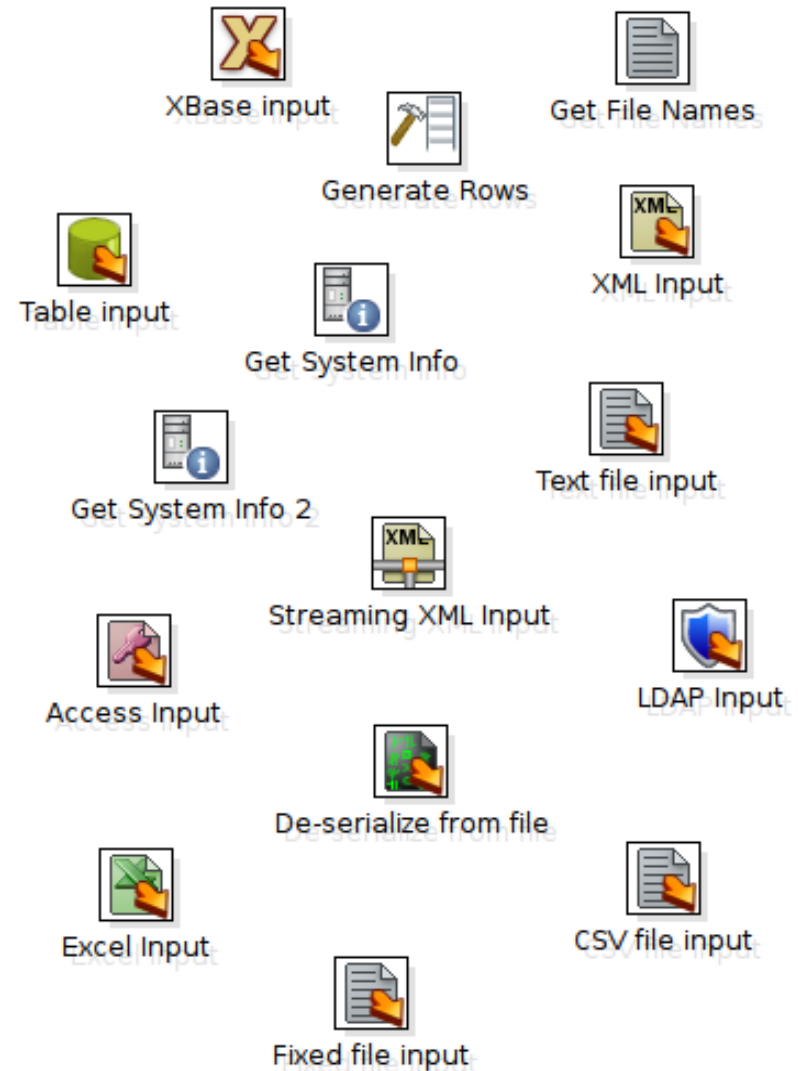
Pentaho Data Integration - Kettle

Kettle
Extraction
Transportation
Transformation
Loading
Environment



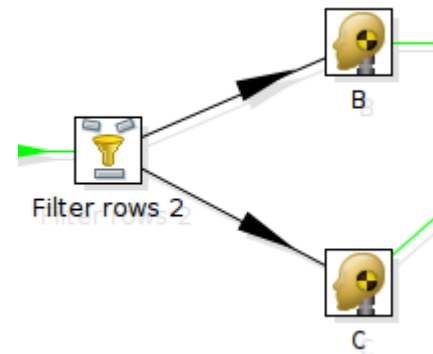
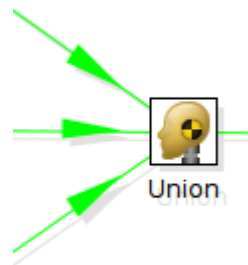
Pentaho Data Integration – Extraction

- Extract data from :
 - 35+ database types
 - MySQL, PostgreSQL, SQLite, ...
 - Oracle, SQL Server, etc
 - Text files
 - XML files
 - XLS files
 - Xbase files (dBase, Foxpro, etc)
 - File systems information
 - Generated data
 - MS Access files
 - LDAP
 - Geo-data
 - ...



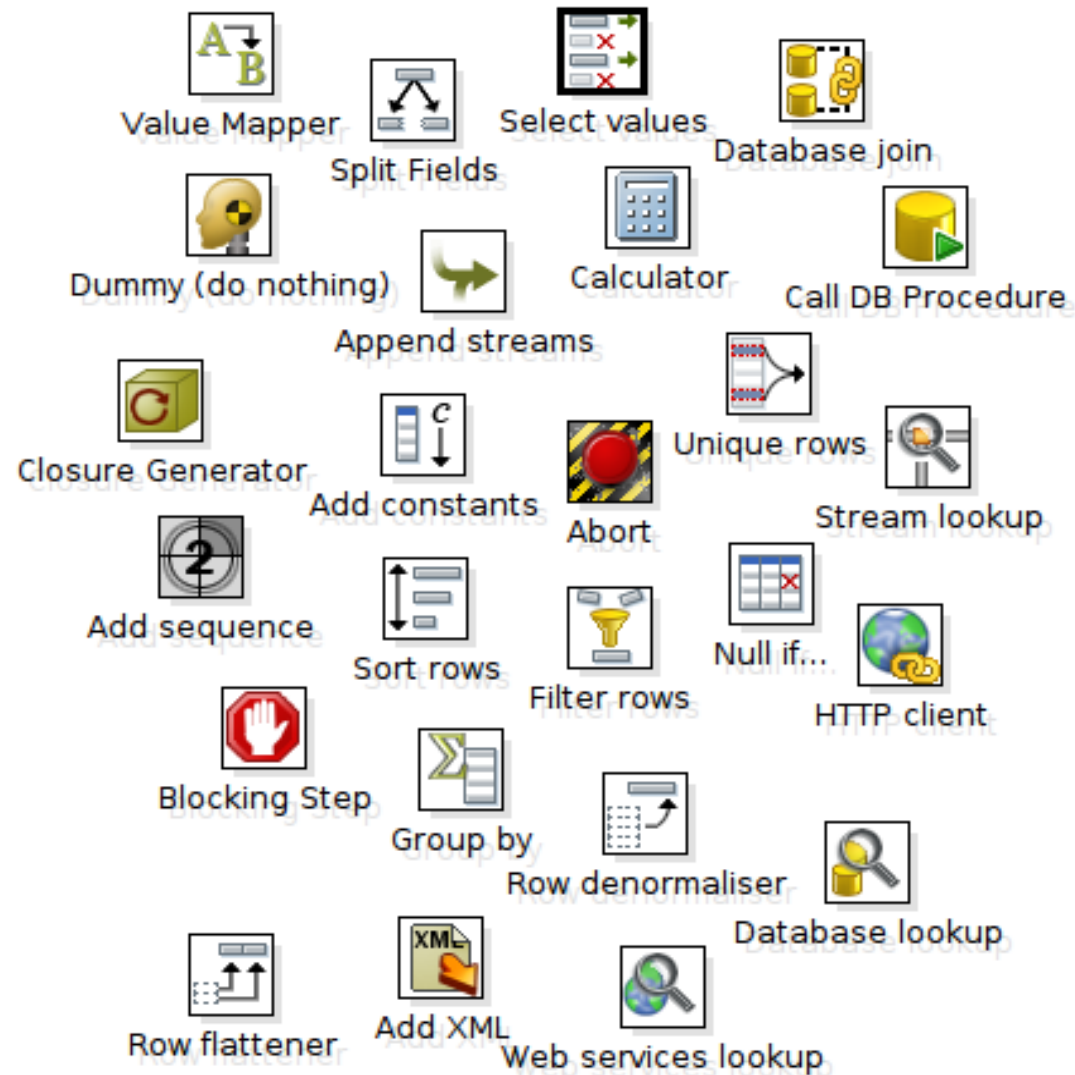
Pentaho Data Integration – Transportation

- Transportation of data
 - Engine based data transfer (no code generator)
 - Very flexible pathways:
 - splitting
 - partitioning
 - merging
 - joining
 - duplicating
 - clustering (MPP)



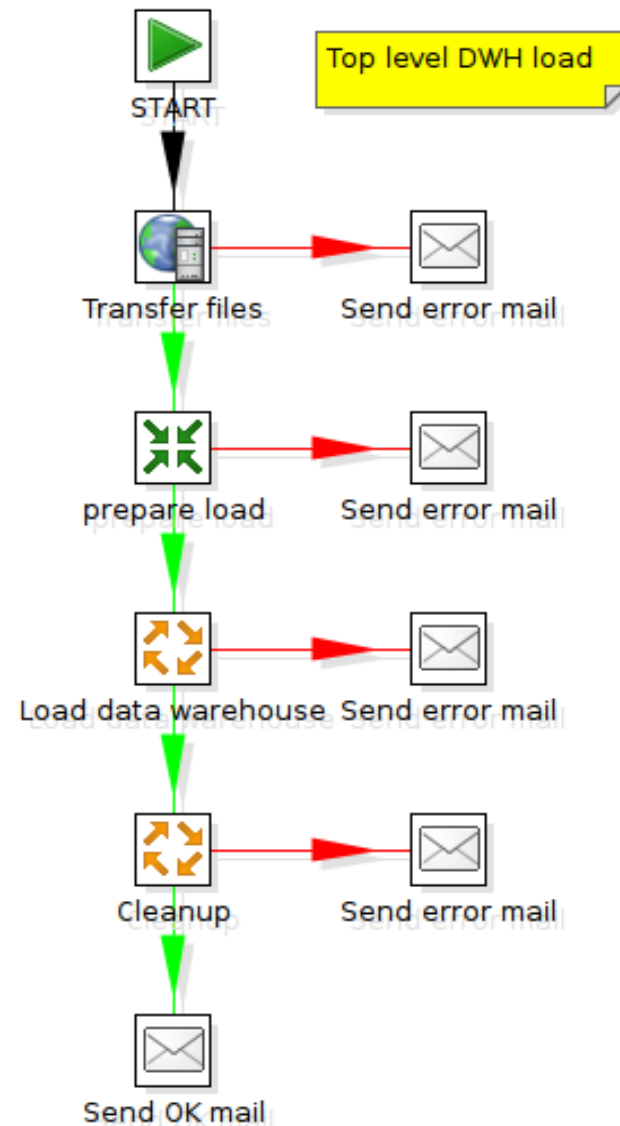
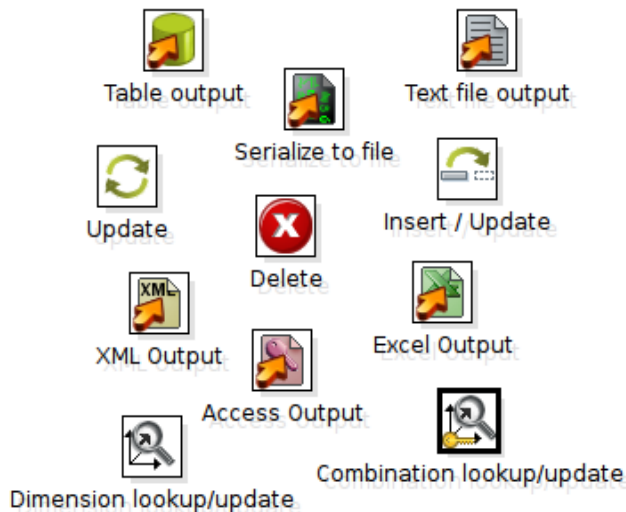
Pentaho Data Integration - Transformation

- Flexibly transform data
 - Looking up data
 - databases
 - files
 - memory...
 - Calculating
 - Scripting
 - JavaScript, SQL, RegExp
 - Splitting
 - Mapping
 - Selecting
 - Filtering
 - Pivotting ...



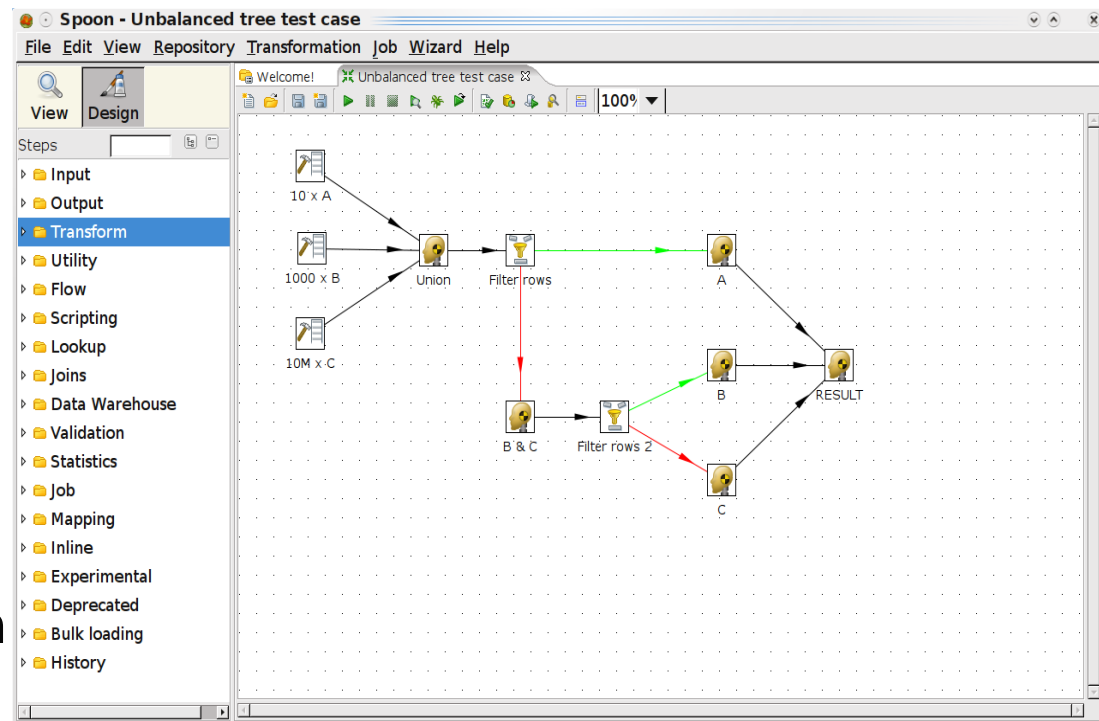
Pentaho Data Integration - Loading

- Load data into a target format
 - Database loads
 - Data warehouse population
 - Partitioned loading
 - Bulk loading
 - Parallel loading
 - Clustering



Pentaho Data Integration – Environment

- Full GUI called “Spoon” to edit every option in Kettle
 - Drag & Drop
 - Debugger
 - Rich GUI
- Command line tools
 - execute jobs
 - execute transformations
- Web server
 - clustering
 - remote execution
- Programming API for Java
- Plugin eco-system
- ...



Pentaho Data Integration – Community

- Paying Pentaho customers
- Large and small corporations
 - All possible sectors
- Lone rangers & Hobbyists
- All regions on Earth
- Meet on our Forum : +30,000 posts in 3 years
- Use our JIRA case tracking systems
- Download more than 10,000 copies of Kettle per month

ohloh popular!

<http://www.ohloh.net/projects/3624?p=Kettle>



<http://www.softpedia.com/progClean/Kettle-Clean-80094.html>

Pentaho Data Integration – use-cases

- Load data from text files and store it into a database **[demo]**
- Export data from database to text-file or more other databases
- Data migration between database applications
- Exploration of data in existing databases (tables, views, ...)
- Information improvement using lookups
- Data cleaning
- Application integration
- Data warehouse population
- Application integration
- Report data generation
- ...



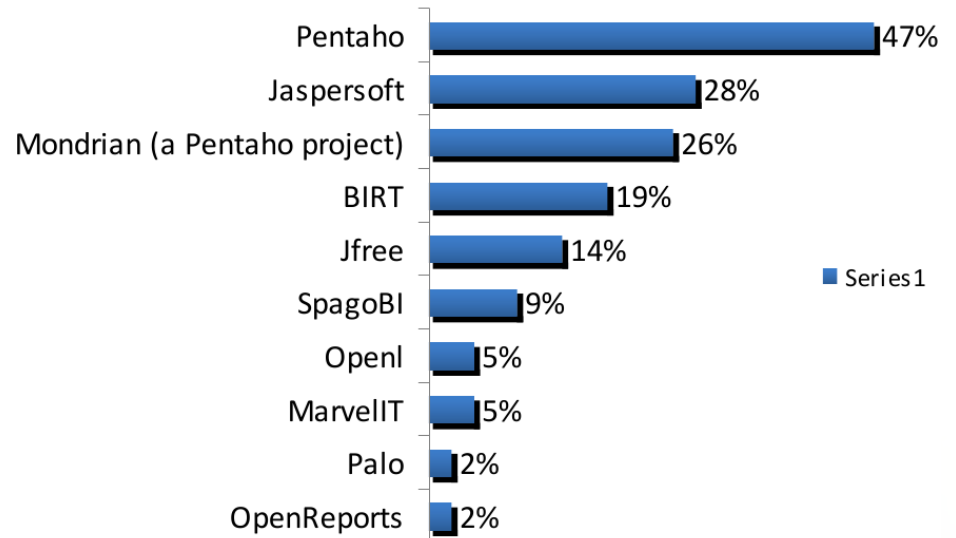
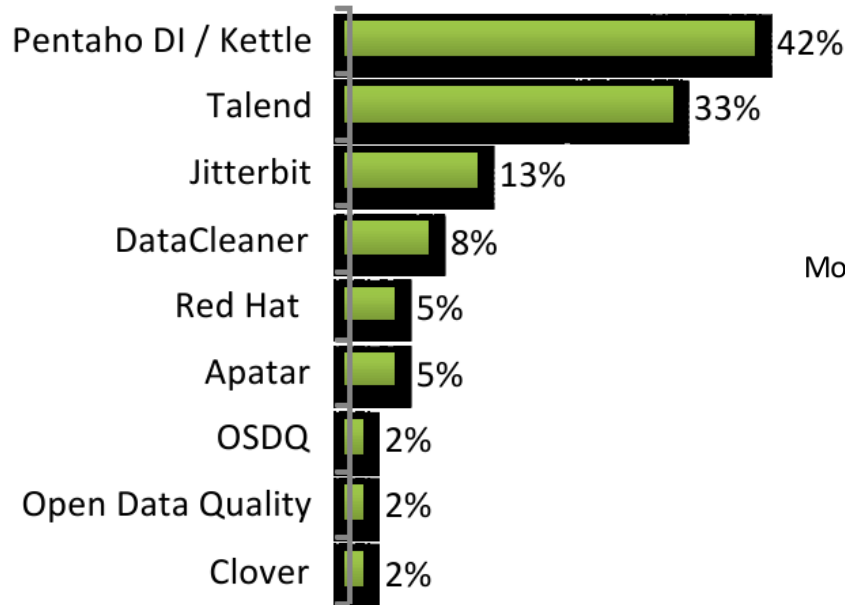
Pentaho Data Integration – Adoption

- Wide range of production deployments
 - Small and medium-sized companies
 - Large enterprises
- Rapid product evolution
 - Driven by Pentaho investment
 - Includes significant community contributions
 - “Contribution-friendly” architecture
 - Natural fit for additional data sources, targets and transformations



Pentaho Data Integration – Adoption

- Most deployed open source data integration solution. Independent study by Mark Madsen of Third Nature and the BeyeNETWORK
- Download free study at pentaho.com



Pentaho Data Integration – Links

- Homepage: <http://kettle.pentaho.org>
- Forum: <http://forums.pentaho.org/forumdisplay.php?f=69>
- Case tracker: <http://jira.pentaho.org/browse/PDI>
- Continuous Integration Server: <http://ci.pentaho.com/job/Kettle>
- Wiki : <http://wiki.pentaho.org/display/EAI>
- IRC Channel: ##pentaho (on Freenode)
- Mailing list: <http://groups.google.com/group/kettle-developers>
- My blog: <http://www.ibridge.be>
- My coordinates: mcasters at pentaho dot org

Agenda

- Pentaho: an introduction
- Pentaho Data Integration
- **Version 4 : New features**
- MySQL support in PDI
- Q&A

Version 4: New features - Visualisation

Demo

- New welcome screen
- Mouse-over slide-outs for icons
- Hop creation
- Improved error handling configuration
- New perspectives support for Agile BI visualisations, modelling, scheduling, etc.

Version 4: New features - Running jobs

- Drill down into running job entries
- Visual indicators of running and completed job entries
- Success and failure mini-icons
- Mouse over completion mini-icons shows details of execution results
- Log capturing of completed job entries

Version 4: New features - Running transformations

- Drill down into running transformation job entries and mappings
- Row input/output sniff testing: see what rows are passing (demo)
- Remote input/output sniff testing on a Carte server

Version 4: New features - Better logging

- Reduced memory consumption
- Incremental log updates
- Global log buffer size limit for long running jobs/transformations
- Interval logging
- Auto clean-up of old log records
- Log record time-outs & execution lineage
- Log record colour coding in Spoon (blue and red for error lines)
- Step and job entry level Logging
- Execution lineage logging
- Renaming individual columns
- Global configuration options for all log tables

Version 4: New features - Plugins

- Unified plug-in architecture
- Easier deployment and packaging
- Step, job entry, partitioner, database type, spoon perspective, life-cycle, ... : all pluggable
- --> MySQL 5.1 plugin

Version 4: New features - Repositories

- Allowing for 3rd party repositories like the Pentaho Unified Enterprise Repository
- Removed dependencies to relational database repository (still supported though)
- Added support for repositories capable of team-development (file locking)
- Added support for repositories capable of fine-grained security repositories
- Added support for repositories capable of storing and retrieving revision history

Version 4: New features – New steps

- SAP Input
- Data Grid
- OLAP Input (Mondrian, Palo, SSAS, SAP B/W)
- Palo Cell Input/Output, Dimension Input/Output
- Salesforce Delete, Insert, Update, Upsert
- Add fields changing sequence (group sequence)
- User Defined Java Class: create your own plugin in Java on the fly in a step
- Send information using Syslog: Send a message to a Syslog server.
- Java Filter
- Memory Group By
- Farrage streaming bulk loader
- Teradata Fastload Bulk loader
- Experimental steps like Get table names, Email messages input, ...

Agenda

- Pentaho: an introduction
- Pentaho Data Integration
- Version 4 : New features
- ***MySQL support in PDI***
- Q&A

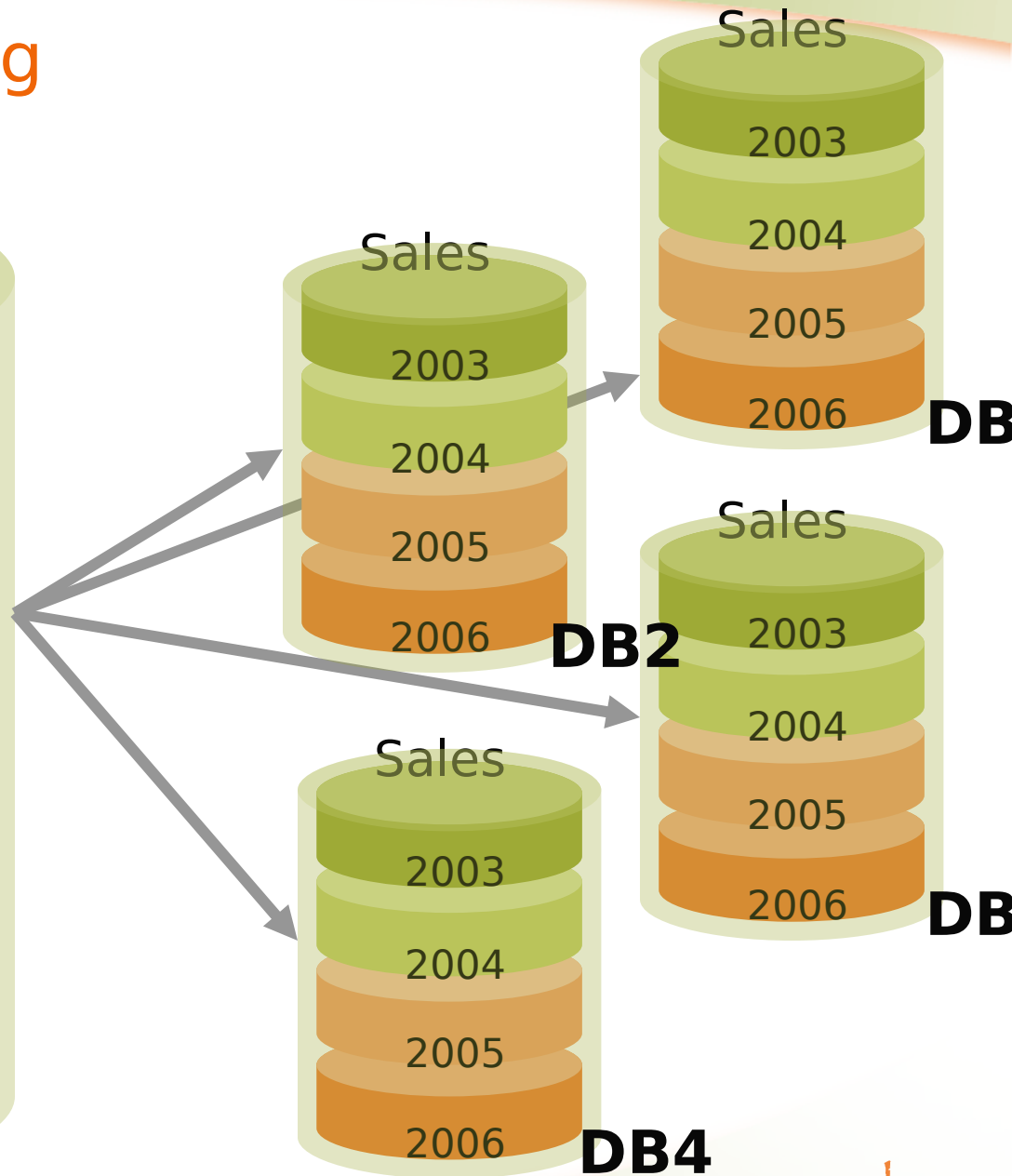
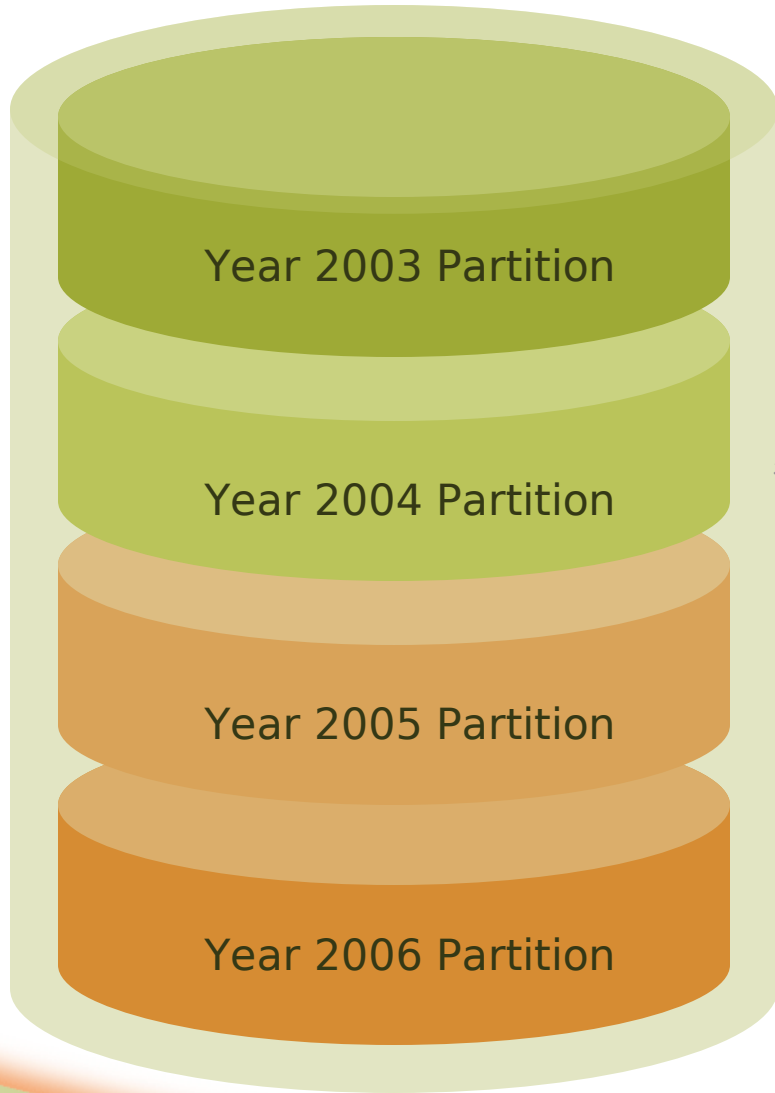
MySQL Support in PDI

- JDBC/ODBC Driver Integration
- Reading: MySQL Result Streaming (cursor emulation) support
- Writing: MySQL dialects for data types
- Job entry: Bulk Loader of text files for MySQL
- Job entry: Bulk writer to a text file for MySQL

Database Partitioning (Sharding) *Demo*

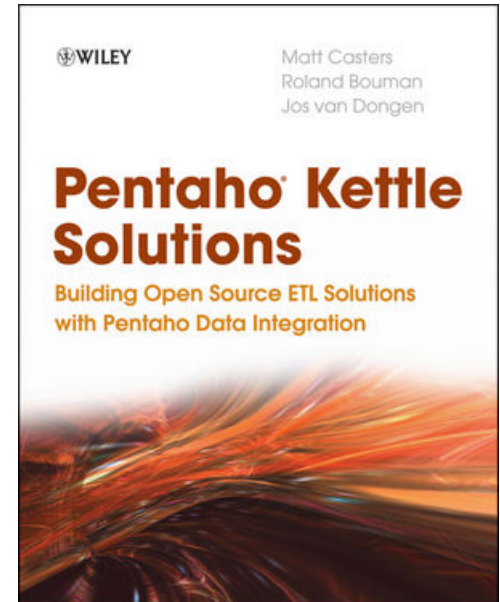
Database partitioning

Sales table



Questions and Closing

- Other Pentaho related User Conference information:
- Collapsing BI from Months to Minutes (Agile BI)
 - Jared Cornelius
 - Ballroom H
 - 11:55am Tuesday April 13th
- MySQL Binary Log Analysis With Pentaho BI
 - Robert Booth
 - Ballroom B
 - 5:15pm Wednesday April 14th
- The Pentaho Booth 516 in the Exhibition Hall



ETA: September 2010