

Business Intelligence

Data Mining and
Data Warehousing

Dominik Ślęzak

slezak@infobright.com

www.infobright.com

Research Interests

Data Warehouses, Knowledge Discovery, Rough Sets
Machine Intelligence, Hybrid Information Technologies
Business Intelligence, Multimedia / Medical Databases

Scientific Degrees

PhD in Computer Science, Warsaw University, Poland
Faculty of Mathematics, Informatics and Mechanics
MSc/BSc in Mathematics, Warsaw University, Poland
Faculty of Mathematics, Informatics and Mechanics

Academy Record

2008 to date	York University, Canada Adjunct Professor
2008 to date	McMaster University, Canada Adjunct Professor
2006 to date	University of Regina, Canada Adjunct Professor
2004 to date	PJIT, Warsaw, Poland Adjunct Professor
2003 – 2006	University of Regina, Canada Assistant Professor
2002 – 2004	PJIT, Warsaw, Poland Assistant Professor
1995 – 2002	PJIT, Warsaw, Poland Teaching Assistant

Industry Record

2006 to date	Infobright Inc., Canada/Poland Chief Scientist (& Co-Founder)
1999 to date	QED Software Inc., Poland President



Simplicity Meets Speed:

A Breakthrough in Data Warehousing



Analytic Data Warehouse Solutions

Businesses must quickly analyze large amounts of data to respond to changes in the market and make smarter decisions. Unfortunately, decreasing performance and growing complexity of the data warehousing environment prevent businesses from seeing the full value of their information. Infobright provides solutions to help businesses become more agile without the complexity of traditional data warehousing solutions.

[Learn More](#)

BrightHouse: An Analytics Data Warehouse

BrightHouse is an analytic data warehouse designed to help businesses make better decisions, faster.

[Learn More](#)

News & Events

[RSS Feed](#)
[News Archive](#)

<p>October 5, 2007 Xerox Chooses Infobright for its Marketing Analytics Data Warehouse</p> <p>Infobright announced today that Xerox has selected Infobright to provide an Analytic Data Warehousing solution for improved customer insight.</p> <p>Learn More</p>	<p>September 6, 2007 Infobright Announces \$8M Series B Investment Lead by IDG Ventures Boston</p> <p>Infobright Inc. today announced it has secured \$8M of Series B financing.</p> <p>Learn More</p>
<p>July 30, 2007 BlueLithium Selects MySQL & Infobright</p> <p>MySQL and Infobright announce expanded partnership and joint customer.</p> <p>Learn More</p>	<p>May 14, 2007 Infobright Presents Research at JRS07</p> <p>JRS07 is the leading symposium on the foundation and application of rough sets</p>

Featured Items

Podcast

TDWI Podcast with Miriam Tuerk, President & CEO

Listen in to the B-Eye Network's interview with Miriam at the TDWI Conference in Orlando Florida.

[Listen to the Podcast](#)

Webinar

High Performance Low Cost Data Warehousing

How to scale to terabytes of data with high performance ad-hoc queries will be discussed.

[Download the Recorded Webinar](#)

Course Contents

- Business Intelligence
- Data Warehouses
 - Monday, 18:00 – 20:30
- Data Mining
 - Wednesday, 18:00 – 20:30
- Presentations
 - Saturday, 9:00 – 11:30 (++)

BI, Data Warehousing & Mining

(...) Business intelligence is a term commonly associated with data warehousing. In fact, many of the tool vendors position their products as business intelligence software rather than data warehousing software. There are other occasions where the two terms are used interchangeably. So, exactly what is business intelligence?

Business intelligence usually refers to the information that is available for the enterprise to make decisions on. A data warehousing (or data mart) system is the backend, or the infrastructural, component for achieving business intelligence. Business intelligence also includes the insight gained from doing data mining analysis, as well as unstructured data (thus the need of content management systems). For our purposes here, we will discuss business intelligence in the context of using a data warehouse infrastructure. (...)



www.b-eye-network.com



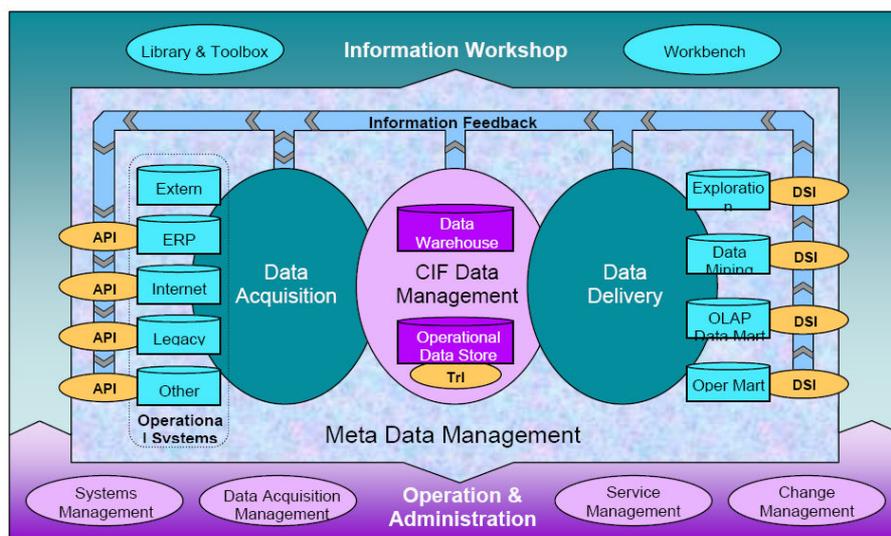
Three Categories of BI

Strategic BI provides performance metrics to management and executives. It is generally called performance management (PM). PM might be preceded by a C for corporate performance management, an E for enterprise performance management or a B for business performance management.

Tactical BI, called traditional and/or analytical, is the application of BI tools to analyze business trends, frequently comparing a specific metric (such as sales or expenses) to the same metric from a previous month or year. In most companies, there are usually a few analysts in each department who use online analytical processing (OLAP) and ad hoc query to perform this task. To date, BI tools are mostly used to analyze historical business data to discover trends or anomalies that need attention.

Operational BI delivers information to the point of business - the front lines of a business where information is used as part of an operational process. For example, when a person calls a toll-free number to speak to a customer service representative about his or her telephone bill, that representative will most likely be looking at a report about the caller's previous billing history and payment record.

Corporate Information Factory



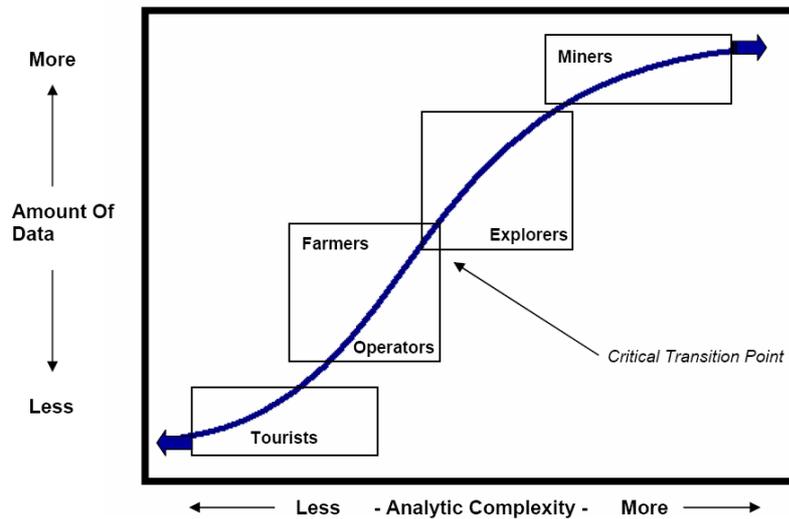
Five Communities that use the CIF

- **Farmers** usually come from the management or business planning groups. Typical farmers may be the financial analysts responsible for reporting on revenues and costs. They may be sales and product analysts determining how well a product is selling in some part of the world. They may be the people who track campaigns or promotions from week to week or they may be the analysts who monitor the budget vs. actual reports.
- **Tourists** generally come from the executives or from very technical resources. Executives are among the most critical users not just because of their need to get information directly from the CIF – they can get information through other means. They are critical because their satisfaction – or dissatisfaction – can significantly impact the future viability of the CIF.

Five Communities that use the CIF

- **Operators** are the most common set of users of the CIF. Sometimes operators make special requests for information, but most of the time they need current detailed information on a scheduled basis, so they rely heavily on standardized queries. As information storage has increased via computers, more and more operators have adopted the view that all the information they need is in the system somewhere. If only they could get someone to give them access to it – and now! As a result, many of today's operators must continually struggle to consolidate and evaluate current information from disparate sources. The good news is that an alternative exists for effectively satisfying operators' needs for current, detailed, enterprise-wide, consolidated tactical information – the Operational Data Store.

The BI User Curve



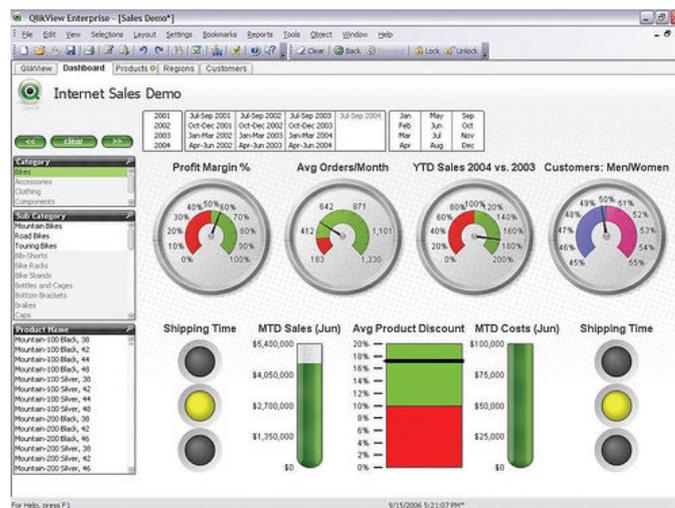
Five Communities that use the CIF

- **Explorers** are the original corporate "out of the box" thinkers. They are individuals who do not look at the life and commerce of the corporation in the standard ways. In some cases these insights are very valuable...
- **Miners** are equipping themselves with tools to find rare and valuable nuggets of information from mountains of data found in data warehouses and data mining data marts. Once settled with appropriate tools, they are patient in their search, determined to find the few insights that can make their organizations rich. They analyze the data to make sure they get meaningful correlations. Miners are thorough. If the data looks odd for some reason, they check it out, since data outliers might indeed be just what they're looking for (finding an instance of fraud, for example).

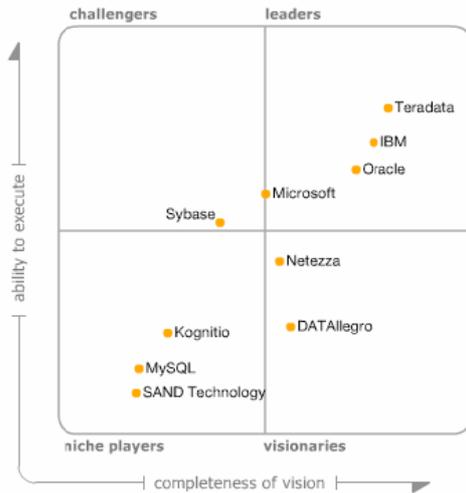
Presentation Topics

For Saturday, January 12

BI Tools and Rules



Gartner Magic Quadrants



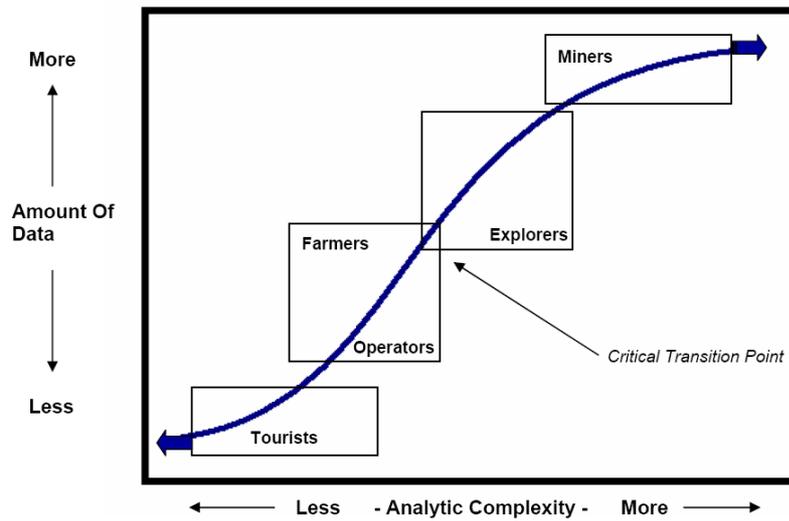
- Gartner's Magic Quadrant for Data Warehouse Database Management Systems
- MySQL is an Open Source Based Organization

Open Source in BI / DM / DW

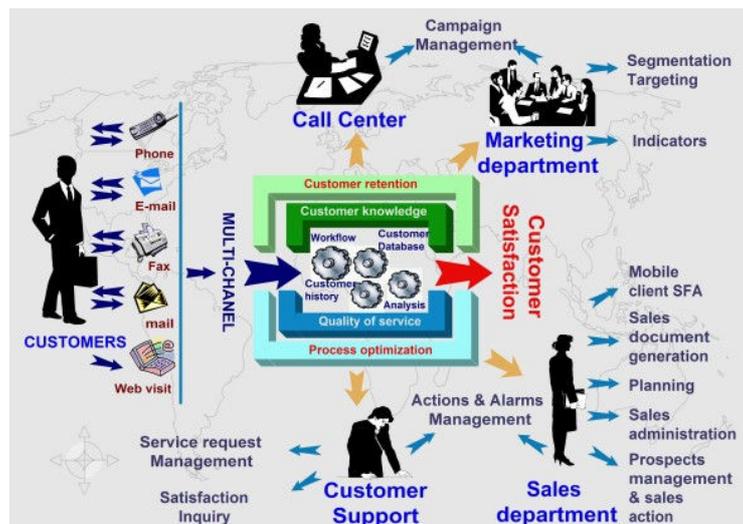
Home	Products	Services	Community	About	Get Support
<p>About Us</p> <p>About Pentaho Customer Successes Meet the Team News, Blogs, Links Events & Webcasts Tribal Store WebEx Event Center Careers Contact Us</p> <p>News</p> <p>News releases Press/Analyst Coverage Blogs Community Announcements</p>		<p>Pentaho Acquires Weka Project</p> <p>World's Most Popular Open Source BI Suite Adds Best-in-Class Open Source Data Mining and Predictive Analytics</p> <p>Orlando FL – Sept. 19, 2006 - Pentaho Corp., creator of the world's most popular open source business intelligence (BI) suite, today announced that it has acquired the Weka open source data mining project developed by the University of Waikato in New Zealand. Weka's advanced analytics equip users with the tools to identify complex patterns, associations and correlations that are critical to evaluating business performance and supporting strategic planning initiatives, particularly in consumer-driven industries such as retail, financial services, communications and consumer packaged goods.</p> <p>Effective immediately, the Weka project is downloadable from www.pentaho.com/download/ as well as from www.sourceforge.net under the GNU General Public License. In addition, Pentaho has launched the first-ever live forum for interaction among Weka project community members at http://community.pentaho.com and has updated its own Data Mining technical forum on www.pentaho.com to support Weka users.</p> <p>The Pentaho and Weka teams will soon announce a road map covering their plans for enhancements including extension of Weka's analytic techniques to additional vertical industries, and integration of Weka with Pentaho's platform, data integration, and information delivery capabilities.</p> <p>"Through the acquisition of the Weka project, Pentaho continues to show that there is no capability or application in the business intelligence market that can't or won't be addressed</p>			



Mixed Workload Challenge



CRM from BI Perspective



Technology-Marketing Misalignment for CRM

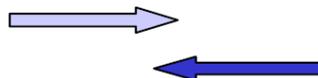
- (...) CRM software was seen as a 'silver bullet' solution that would change the way companies were organized and run to focus on the customer.
- CRM software systems, when deployed, instead resulted in massive data overload that companies were not prepared to handle. In the worst cases, no attention was paid to helping business users take advantage of the systems, and the software was simply put on the shelf.
- The vaunted 'closed loop' promise really became a closed loop breakdown in software-only CRM implementations. (...)

Technology-Marketing Misalignment

Technology



Logical
Organized
Standardized



Marketing



Constant Change
Tactical
Strategic

Technology-Marketing Misalignment

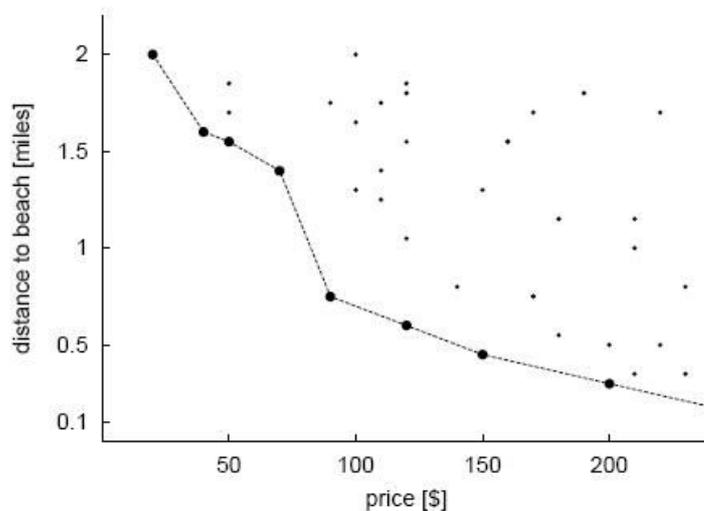
- The questions the IT professional tends to ask:
 - What storage, management, and access systems are most efficient and effective?
 - What tools, programming, and protocols need to be in place to keep the system running (nearly) all the time?
 - What analytic or viewing tools allow the best 'cut' at the data in my system without disrupting everything or slowing the system to a halt?
- The questions marketers want to answer:
 - What is the right balance of marketing spending?
 - What mix of tactics and communication channels deliver the highest return?
 - Why does my customer choose my competitor's product over mine?

Five Communities – Marketing

- Farmers: a sizeable group focused mostly on monitoring and keeping track of the dynamic activities of marketing. This includes: tracking (campaigns, advertising, brand / customer loyalty,) and receiving reports from explorers.
- Tourists: are rare in marketing, as most analysis filters down to the marketing farmers, who then condense reports to the chief executive of marketing.
- Operators: these are the brand and marketing managers, who attempt to track marketing allocations, pay back on campaigns, and success rates of advertising. Their analysis is very current, and if they have a fault it is in not using technology successfully to track, trend, and evaluate short, medium, and long-term results. This aspect of marketing is one of the least developed from a technology and BI perspective.

- Explorers: in marketing, due to the nature of its filtering and creative processes, exploration and discovery is a major part of BI. And the community is extended from internal practitioners to outside consultants, agencies, and specialists. Many of the explorers' roles have evolved over the years into specific, non-technology based activities, such as: focus group testing of new product introductions, or surveys of customer attitudes.
- Miners: as marketing databases have grown, and the introduction of e-commerce has taken hold, miners perform a particularly specialized role in marketing. Their link in the process is segmenting, classifying, and using predictive modeling to assist marketer's efforts. They are also, however, the most out of touch with marketer's daily concerns. Two primary functions they carry are to answer or address marketing explorer's questions and to provide support to direct marketer's targeting.

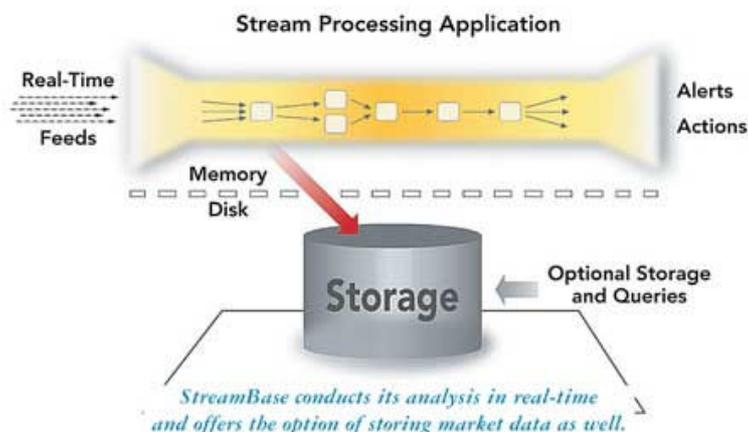
New Queries for Analytics / BI



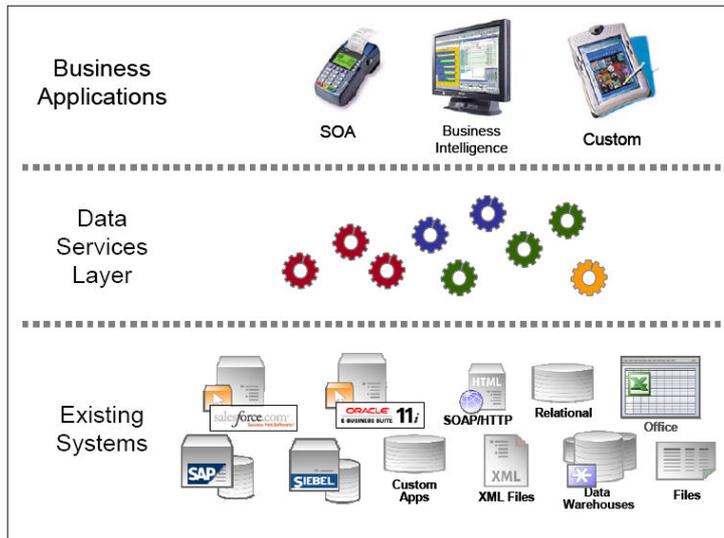
New Queries for Analytics / BI

(...) Forecasting future events based on historic data is useful in many domains like system management, adaptive query processing, environmental monitoring, and financial planning. We describe the Fa system where users and applications can pose declarative forecasting queries – both one-time queries and continuous queries – and get forecasts in real-time along with accuracy estimates. Fa supports efficient algorithms to generate execution plans automatically for forecasting queries from a novel plan space comprising operators for transforming data, learning statistical models from data, and doing inference using the learned models. (...)

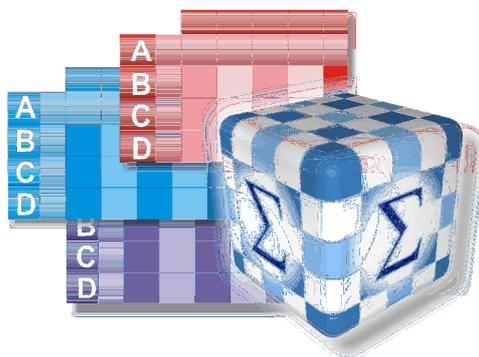
Data Stream Processing



Between DW and BI



OLAP – Extensions and Mining



(...) OLAP cannot help marketers distinguish or understand the 'who, what, and why' of their customers beyond simple dimensional representations of historical data (...)

OLAP Rules (E. F. Codd, 1993)

- **Multidimensional conceptual view.** OLAP operates with CUBEs of data that represent multidimensional construct of data. Event though the name implies three dimensional data, the number of possible dimensions is practically unlimited.
- **Transparency.** OLAP systems should be part of an open system that supports heterogeneous data sources.
- **Accessibility.** The OLAP should present the user with a single logical schema of the data.
- **Consistent reporting performance.** Performance should not degrade as the number of dimensions in the model increases.
- **Client/server architecture.** Should be based on open, modular systems.
- **Generic dimensionality.** Not limited to 3-D and not biased toward any particular dimension. A function applied to one dimension should also be able to be applied to another.
- **Dynamic sparse-matrix handling.** Related both to the idea of nulls in relational databases and to the notion of compressing large files, a sparse matrix is one in which not every cell contains data. OLAP systems should accommodate varying storage and data-handling options.
- **Multiuser support.** OLAP systems should support more than one user at the time.
- **Unrestricted cross-dimensional operations.** Similar to rule of generic dimensionality; all dimensions are created equal, and operations across data dimensions should not restrict relationships between cells.
- **Intuitive data manipulation.** Ideally, users shouldn't have to use menus or perform complex multiple-step operations when an intuitive drag-and-drop action will do.
- **Flexible reporting.** Save a tree. Users should be able to print just what they need, and any changes to the underlying financial model should be automatically reflected in reports.
- **Unlimited dimensional and aggregation levels.** The OLAP cube can be built with unlimited dimensions, and aggregation of the contained data also does not have practical limits.

Interaction and Visualization

The screenshot shows a window titled "Postgres95 Online Aggregation Interface". It contains a table with the following data:

	Speed	major	AVG	Confidence	Interval
		1	2.23507	95	0.11801
		2	2.52418	95	0.0847007
		3	2.67458	95	0.146706
		4	2.87143	95	0.0847007
		5	3.12746	95	0.18659
		9	2.89143	95	0.146706

Below the table, there is a "Skip Factor" set to 250, a "Cancel All" button, and a progress bar showing "28% done".

MPP Architectures in DW

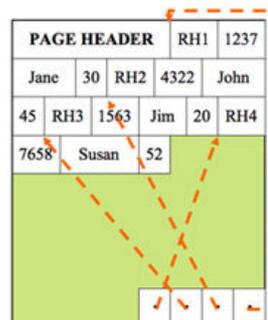


Column Oriented Databases

Standard record layout for row stores

Employee Table

RID	SSN	Name	Age
1	1237	Jane	30
2	4322	John	45
3	1563	Jim	20
4	7658	Susan	52
5	2534	Leon	43
6	8791	Dan	37



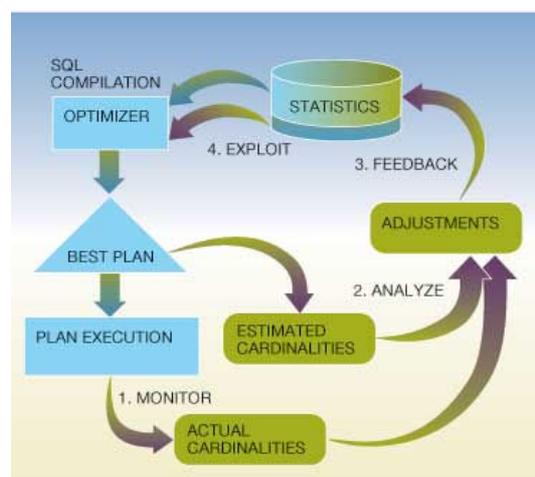
- Records are stored sequentially
- Offsets to start of each record at end of page

Decision Support Benchmarks

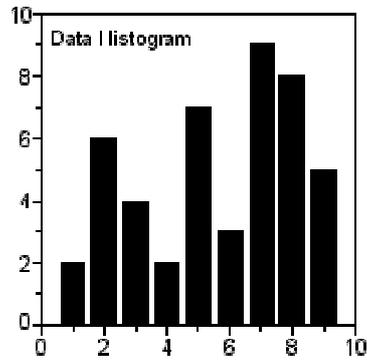
TPC-H → TPC-DS

TPC Transaction Processing
Performance Council

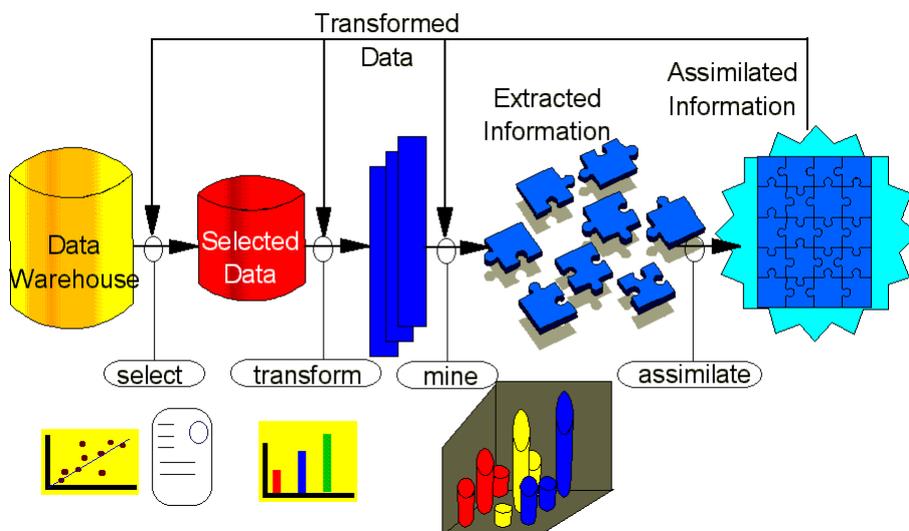
Query Optimization



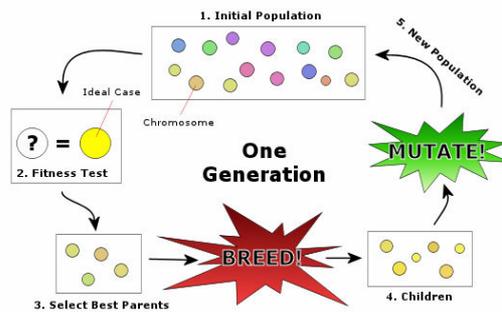
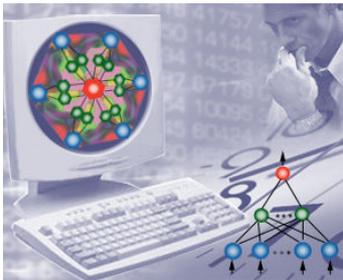
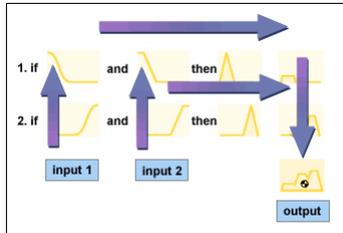
Histograms in DW



Data Mining for Databases



Artificial Intelligence for BI / DW



DZIĘKUJĘ !!!

slezak@infobright.com