





Tutorial



## Using the AMBER Data Repository to Analyze, Share and Cross-exploit Dependability Data

Marco Vieira  
[mvieira@dei.uc.pt](mailto:mvieira@dei.uc.pt)  
University of Coimbra, Portugal  
*The Second International Conference on Dependability (DEPEND 2009)*  
Athens/Glyfada, Greece, June 18, 2009



The AMBER Project



- **Assessing, Measuring and Benchmarking Resilience in computer systems and components (AMBER)**
- Coordination Action supported by the European Commission in the 7th FP
- Coordinating and advancing research in resilience measurement and benchmarking in computer systems and infrastructures

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



Current challenges



- Quality of measurements
- Integration of the human and technical components of the analysis
- Dynamic and adaptive systems and networks
- Integration with the development processes

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



AMBER objectives



- State-of-the art survey
- Research agenda
- **Data repository**
- Others:
  - Dissemination events (workshops, panels, etc)
  - Benchmarking tools
  - Training material

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



This Tutorial...



## Learn how to use the AMBER Data Repository to analyze and share data from dependability evaluation experiments

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

Problems



- How to **analyze** the usually large amount of raw data produced in dependability evaluation experiments?
- How to **compare** results from different experiments or results of similar experiments across different systems?
  - Different and incompatible tools, data formats, and setup details...
- How to **share** raw experimental results among research teams?

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Current situation**

- **The situation today is not good!!!**
- Spreadsheets and other specific tools to analyze results
  - Not standard and difficult to build
- Difficult to compare data and generalize conclusions
- Researchers share final results and conclusions
  - Papers, mainly
  - Raw data is not shared

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 7

**ADR Vision and objectives**

- Vision
  - Become a worldwide repository for dependability related data
- Key objectives:
  - Provide state-of-the-art data analysis
  - Allow data comparison and cross-exploitation
  - Facilitate worldwide data sharing and dissemination
- Potential tool to increase the impact of research

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 8

**Data analysis approach**

- Repository to analyze, compare, and share results
- Use a business intelligence approach:
  - Data warehouse to store data
  - On-Line Analytical Processing (OLAP) to analyze data
  - Data mining algorithms to identify (unknown) phenomena in the data
  - Information retrieval for data in textual formats
- Adopt the same life cycle of BI data


DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 9

**Outline**

1. Business Intelligence
2. Data Warehousing & OLAP
3. Using DW to analyze dependability related data
4. The AMBER Data Repository

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 10

**1. Business Intelligence**



**What is Business Intelligence?**

- Business Intelligence (BI):
  - **Getting the right information, to the right decision makers, at the right time**
- BI is an enterprise-wide platform that supports, data gathering, reporting, analysis and decision making
- BI is meant to:
  - Fact-based decision making
  - “Single version of the truth”
- BI includes reporting and analytics

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 12

### Five classic BI questions

- What happened?
- What is happening?
- Why did it happen?
- What will happen?
- What do I want to happen?

Past  
Present  
Future

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Typical BI technologies

- ETL Tools (Extract, Transform, and Load)
- Repositories
  - Data Warehouse
- Analytical tools
  - Reporting and querying
  - OLAP
  - Data mining
- Information retrieval

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Many proprietary products

|                      |                                |                       |
|----------------------|--------------------------------|-----------------------|
| ACE*COMM             | Microsoft                      | SAS Institute         |
| Ab Initio            | Microsoft Analysis Services    | Siebel Systems        |
| Actuate              | PerformancePoint Server 2007   | Spotfire (now Tibco)  |
| ComArch              | Proclarity                     | StatSoft              |
| CyberQuery           | Oracle Corporation             | SPSS                  |
| Dimensional Insight  | Hyperion Solutions Corporation | Telerik Reporting     |
| IBM                  | Panorama Software              | Teradata              |
| Applix               | Pentaho                        | Thomson Data Analyzer |
| Cognos               | Pervasive                      |                       |
| InetSoft             | Pilot Software, Inc.           |                       |
| Informatica          | PRELYTIS                       |                       |
| Information Builders | Prospero Business Suite        |                       |
| LogiXML              | Qliktech                       |                       |
| LucidEra             | SAP Business Inf. Warehouse    |                       |
| MicroStrategy        | Business Objects               |                       |
|                      | OutlookSoft                    |                       |

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Some open source/free products

- Eclipse BIRT Project:
- Freereporting.com:
- JasperSoft:
- OpenI:
- Palo (OLAP database):
- Pentaho:
- RapidMiner
- SpagoBI:
- Weka

- Some products from big companies can be used freely

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

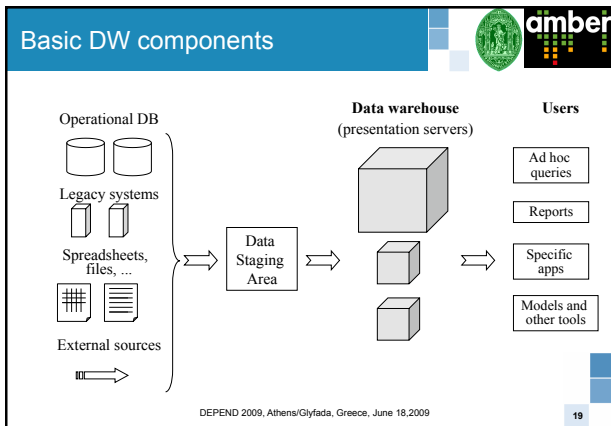
## 2. Data Warehousing & OLAP

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### What is a Data Warehouse?

- Big database that stores data for **decision support**
- Built from the operational data collected from transactional DB and other operational systems

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



- ### Data volume
- Less than 20 GBytes
    - Small dimension; runs in a PC
  - From 20 to 100 GBytes
    - Medium dimension; needs a powerful workstation
  - From 100 Gbytes to 1 TBytes
    - Large dimension; needs a powerful server, normally with parallel processing
  - More than 1 TBytes
    - Very large dimension; massive parallel processing
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009
- 20

- ### Some characteristics
- Temporal dependency
  - Non volatile
  - Target oriented
  - Data integration and consistency
  - Designed for queries
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009
- 21

- ### Temporal dependency
- The data is collected over time
    - Do not represent a specific moment
    - Represents the history
  - A temporal reference must be associated to all data in the database
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009
- 22

- ### Non volatile
- The data in the DW is never updated
  - The DW stores historic data (historic memory) collected from the operational databases
  - After being load (from the operational databases) there is only one operation:
    - Queries
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009
- 23

- ### Target oriented
- The data warehouse must only store data relevant for decision support
  - Many operational data (needed for everyday management) is not relevant for the DW
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009
- 24

### Data integration and consistency

- In an operational environment the information may be stored in different locations using different representations
- That data must be integrated and made consistent before being load in the DW

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 25

### Designed for queries

- After being load the data never changes:
  - Only queries are allowed
- DW stores a large amount of data

The data must be stored in such a way that improves performance

*Multidimensional view*  
*Partial denormalization*

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 26

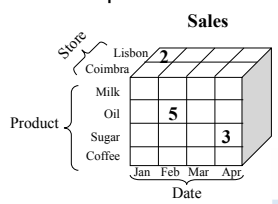
### Dimensional model

- Typical model in operational databases: E/R
- The dimensional model follows a different approach
  - Stores the same data
  - Data organization is user oriented
    - Easy to understand
    - Very good performance for queries
- Data Warehouses built over complex E/R models never succeed

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 27

### The multidimensional model

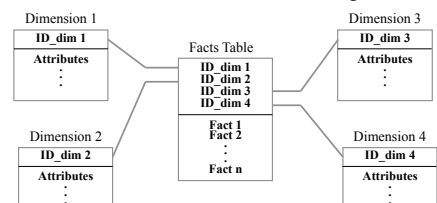
- Facts stored in a multidimensional array
- The dimensions are used to index the array
- Usually built using data from operational databases



DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 28

### Star model

- The typical dimensional model is a star structure with:
  - A central table with facts
  - Several dimensions tables describing the facts



DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 29

### Facts

- Represent the business measures
- The most useful facts are:
  - Numbers
  - Additives

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 30

### Facts table

- Comprises several numeric attributes (facts) and foreign keys to the dimensions
- Normalized table
- Relationships M:1 with the business dimensions
- Contains normally a large number of records
- Represents typically 95% of the space used by the DW

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Dimensions

- Each dimension represents a business parameter
  - Time, clients, products, etc
- Represent a entry point for the analysis of the facts
- Represent different point-of-views for the analysis of the facts

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Dimension tables

- Strongly denormalized
  - For performance
- Dimensions have hierarchies
  - Day → Month → Year → ... Contain a large set of attributes
- Typically comprise a small number of records (when compared to the facts table)

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Star schema example

The diagram shows a star schema with a central fact table 'Sale' and three dimension tables: 'Time', 'Product', and 'Store'. The 'Sale' table contains attributes: ID\_time, ID\_product, ID\_store, Units\_sold, Purchase\_cost, Sale\_value, and Num\_Clients. The 'Time' table contains: ID\_time, Day, Day\_of\_week, Week\_of\_year, Month, Trimester, and Year. The 'Product' table contains: ID\_product, Name, Type, Brand, Category, Pack, and Description. The 'Store' table contains: ID\_store, Name, Local, District, Area, and Num\_tellers.

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Low level queries

The diagram is identical to the one in slide 34. Below it is the following SQL query:

```
select avg (sale_value x units_sold)
from sale, time, product
where JOIN_TABLES
group by brand, month
```

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### User interfaces

- Explore data in Data Warehouses
  - Typical OLAP tools
    - Access the relational engine using SQL
    - Data presentation using tables, graphics, reports, etc
    - Targeted for ad-hoc queries
  - Other tools
    - Data mining
    - Modeling

The screenshot shows a software interface with a table of data on the left and a bar chart on the right. The table has columns for 'Description', 'Frequency', and 'Recovery Time'. The bar chart has 'Recovery Time' on the x-axis and 'Frequency' on the y-axis.

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Queries - Slice and Dice

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

37

### Drill-Down & Roll-Up

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

38

### Time: Drill-Down & Roll-Up

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

39

### Steps for the design of the star model

1. Identify the business process/activity
2. Identify the facts
3. Identify the dimensions
4. Define the data granularity
  - Day, Week, Month, ...
  - Product, Category, ...
  - Store, City, ...

*Do not forget that the model depends on the data available (operational databases, files, etc)*

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

40

### Example – Retail sales

- Set of stores belonging to the same enterprise
- Goal: Analysis of sales
- Each store has several departments (food, hygiene and cleaning, etc)
- Sells thousands of products
- Products are identified using a unique number

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

41

### Retail sales – Business data

- Where to collect the data?
  - POS - point of sales
  - Operational database
- What to measure?
  - Sales
- Goals?
  - Maximize the profit
  - Maximum sales price possible
  - Lower costs
  - More clients

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

42

### Retail sales – Facts

- Examples of relevant decision support facts:
  - Number of units sold
  - Acquisition costs
  - Sale value
  - Number of clients that bought the product
- **Question:** is it possible to obtain base data (from the operational system) for these facts?

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Retail sales – Dimensions

- Main dimensions:
  - Product x Store x Time
- Are there other relevant dimensions?
  - Supplier? – Promotions? – Client?
  - Employee responsible for the store on that day?
- It is normally possible to add extra dimensions
- All the dimensions have a 1:M relationship with the facts

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Retail sales

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Granularity

- Example: record the daily sales for all products
  - Analyze in detail (price, quantity, etc) the products sold every day, in each store, ...
- Retail sales granularity:
  - Products x Store x Promotion x Day
- The granularity defines the detail of the DW and has a strong impact in the size
- The granularity must be adjusted to the analysis requirements

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Retail sales – Details

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Retail sales – Details

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



### Retail sales – Details

**Product**

- ID\_product
- description
- full\_description
- SKU\_number
- package\_size
- brand
- subcategory
- category
- department
- package\_type
- diet\_type
- weight
- weight\_unit\_of
- units\_per\_retail
- units\_per\_shippi
- cases\_per\_pallet
- shelf\_width\_cm
- shelf\_height\_cm
- shelf\_depth\_cm
- .....

**Store**

- ID\_store
- name
- store\_number
- store\_street\_address
- city
- store\_county
- store\_state
- store\_zip
- sales\_district
- sales\_region
- store\_manager
- store\_phone
- store\_FAX
- floor\_plan\_type
- photo\_processing\_type
- finance\_services\_type
- first\_opened\_date
- last\_remodel\_date
- store\_sqft
- grocery\_sqft
- frozen\_sqft
- meat\_sqft
- .....

- Must characterize the stores as seen by the business management
- Must contain the attributes that are relevant for posterior queries
  - Includes geographical attributes (localization)
  - Includes time attributes (opening date, ...).

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 49

### Retail sales – Details

**Product**

- package\_type
- diet\_type
- weight
- weight\_unit\_of
- measure
- units\_per\_retail
- case
- units\_per\_shipping
- case
- cases\_per\_pallet
- shelf\_width\_cm
- shelf\_height\_cm
- shelf\_depth\_cm
- .....

**ID time**

- date
- day\_of\_week
- day\_number\_in\_month
- day\_number\_overall
- week\_number\_in\_year
- week\_number\_overall
- Month
- quarter
- fiscal\_period
- year
- holiday\_flag
- .....

**Promotions**

- ID\_promotion
- number
- name
- type\_price\_red
- type\_advertisement
- type\_poster
- Type\_coupons
- promotion\_cost
- start\_date
- end\_date
- .....

**Sales**

- sales\_district
- sales\_region
- store\_manager
- store\_phone
- store\_FAX
- floor\_plan\_type
- photo\_processing\_type
- finance\_services\_type
- first\_opened\_date
- last\_remodel\_date
- store\_sqft
- grocery\_sqft
- frozen\_sqft
- meat\_sqft
- .....

- Characterizes the existing promotions
- In this example there is only one dimension related to promotions
- Represents a very important dimension
  - Managers want to know the impact of promotions in the sales in order to target new promotions to specific products, stores and time

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 50

### More than one star

Store

**Sales**

- ID\_time
- ID\_product
- ID\_store
- units\_sold
- purchase\_cost
- sale\_value
- num\_clients

Time

**Stock**

- ID\_time
- ID\_product
- ID\_warehouse
- quant\_available
- quant\_out
- purchase\_cost
- last\_sell\_price

Warehouse

Product

- Two or more stars can be connected using one or more dimensions
- Shared dimensions must be conform
  - Contain consistent data when considering each star
- Drill across:** query that crosses more than one start

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 51

### Several stars

**Orders**

- Dimension: Time
- Dimension: Component
- Dimension: Supplier
- Dimension: Contract

**Sales**

- Dimension: Time
- Dimension: Component
- Dimension: Client
- Dimension: Contract

**Stocks**

- Dimension: Time
- Dimension: Component
- Dimension: Warehouse

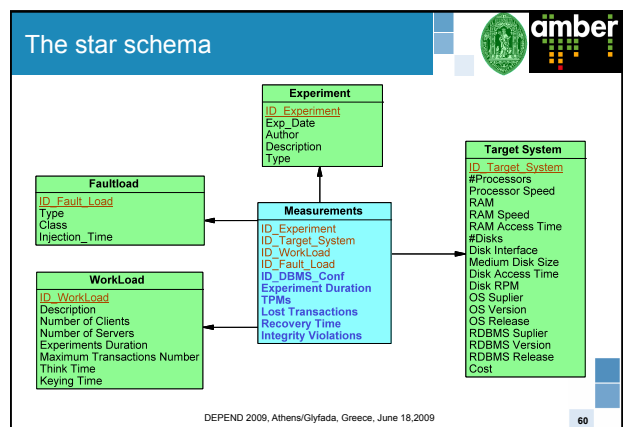
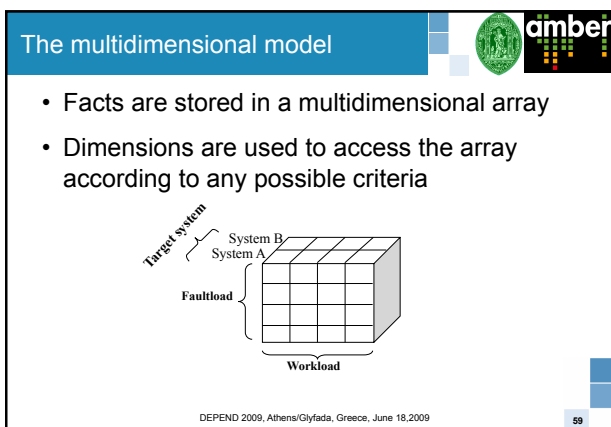
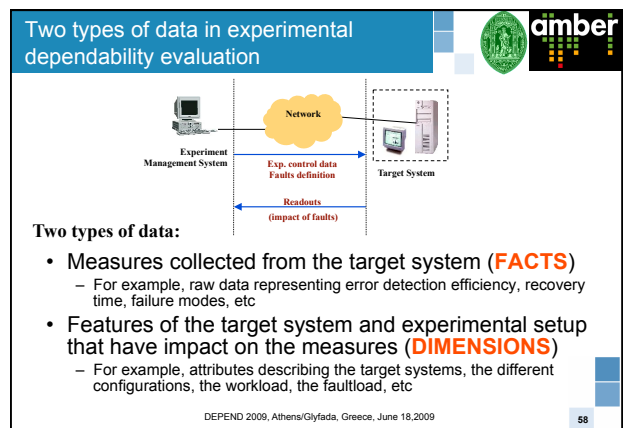
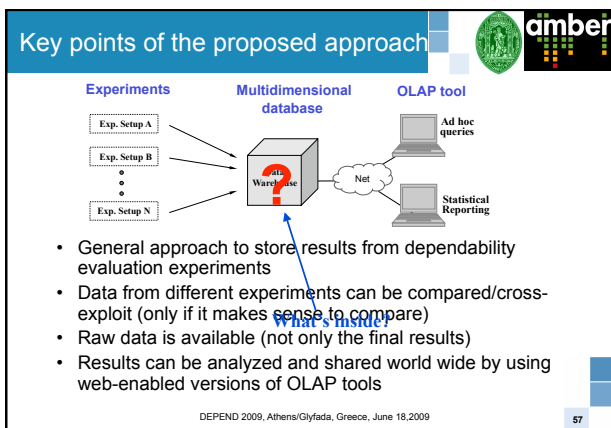
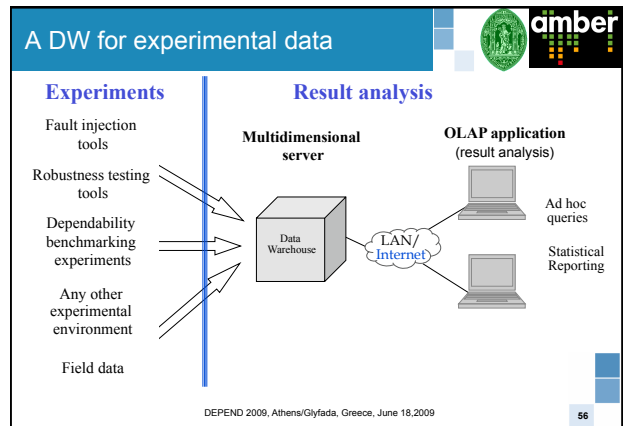
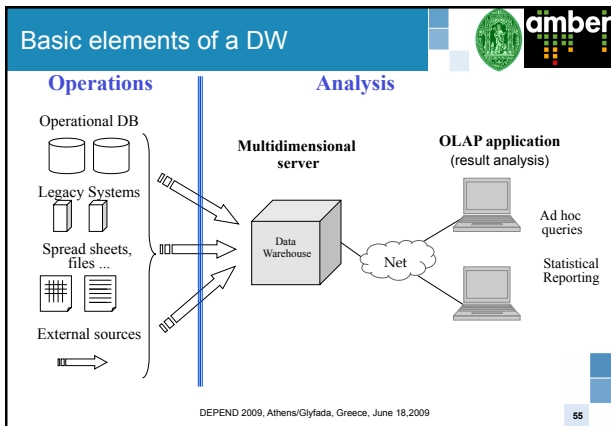
DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 52

### Questions

?

DEPEND 2009, Athens/Glyfada, Greece, June 18,2009 53

### 3. Using DW to analyze dependability data



**Basic elements of the proposed approach**

The experimental setups are used as they are. You can use your favorite dependability evaluation tool and do the experiments in the usual way. It's necessary...

- To know the format of the raw results
- To have access to the results

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Basic elements of the proposed approach**

**Loading applications**

- General purpose loading applications
- Some transformations in the data are normally necessary for consistency

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Basic elements of the proposed approach**

**Data warehouse**

- Raw data is available in a standard star schema (facts + dimensions)
- Results from different experiments are compatible and can be compared/analyzed together, then they are stored in the same star schema (or in schema that share at least one dimension)
- If results are from different unrelated experiments then they are stored in a separated schema

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Basic elements of the proposed approach**

**Analysis**

- Commercial OLAP tools are used to analyze the raw data and compute the measures. These tools are designed to be used by managers: *very easy to use :-)*
- Just need an internet browser to analyze the data

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Steps needed to put our approach into practice**

- Definition of the adequate star schema to store the data. Create the tables in the data warehouse
- Use general-purpose loading application to define the loading plans for each table in the star schema
- Run the loading plans to load the star tables with the raw data collected from the experiments
- Every time a new experiment is done corresponding loading plans are run again to add the new data to the data warehouse
- Analyze the data: calculate measures, find unexpected results, analyze trends, etc

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Example: Recovery and Performance Evaluation in DBMS**

- Tuning of a large DBMS is very complex
- Administrators tend to focus on performance tuning and disregard the recovery features
- Administrators seldom have feedback on how good a given configuration is
- A technique to characterize the performance and the recoverability in DBMS is needed

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### The Approach

- Extending existing performance benchmarks to evaluate recoverability features in DBMS
- Include a **faultload** and new **measures**
- **Faultload** based on operator faults
- **Measures** related to recovery:
  - Recovery time
  - Data integrity violations
  - Lost transactions

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Operator faults injection and recovery

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Experimental setup

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Steps towards data analyzes

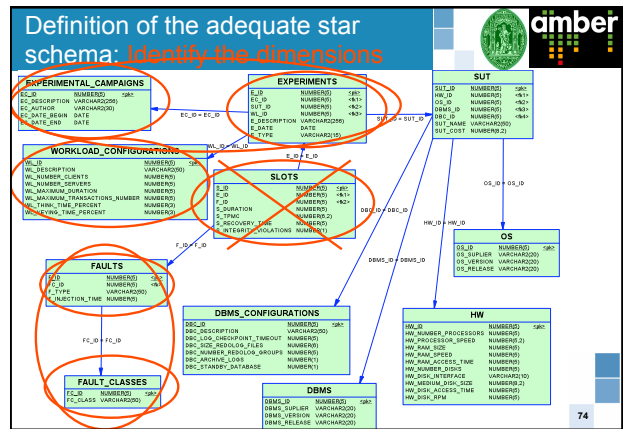
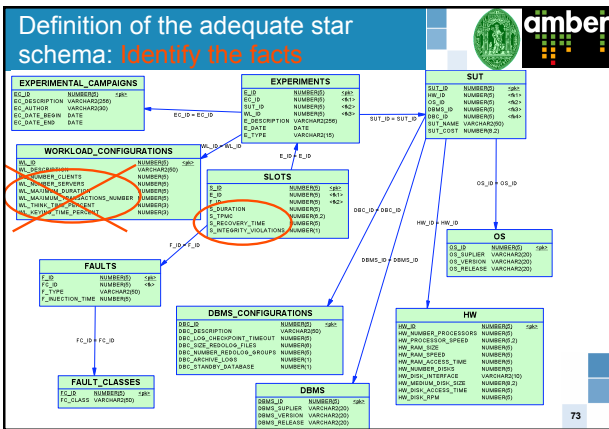
1. Definition of the adequate star schema
  - a. Identify the process/activity
  - b. Identify the facts
  - c. Identify the dimensions
  - d. Define the data granularity
2. Load the data
3. Analyze the data

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

### Definition of the adequate star schema: Identify the process/activity

- Experiments to characterize the performance and the recoverability in DBMS
- Includes a faultload and new measures
- Faultload based on operator faults
- Measures related to recovery

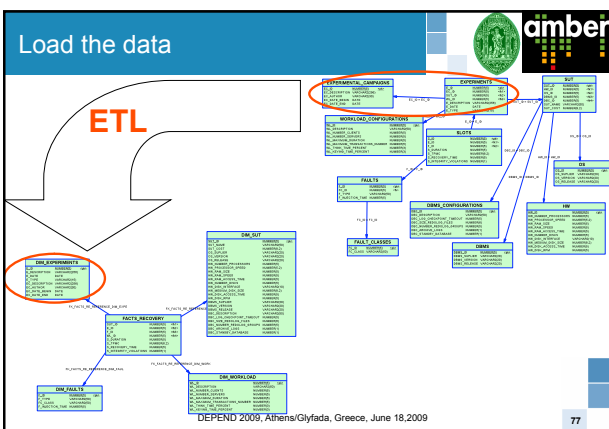
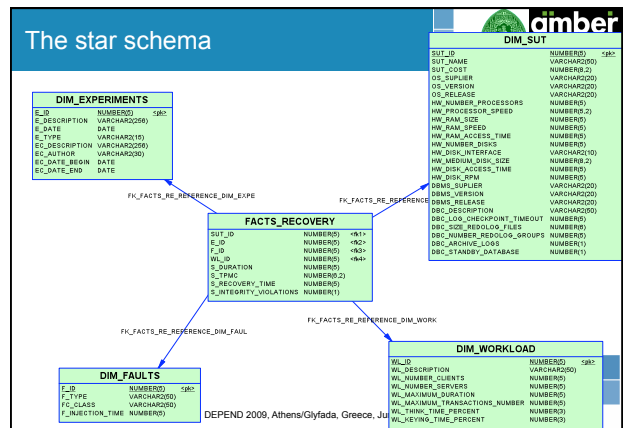
DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



### Definition of the adequate star schema: Define the data granularity

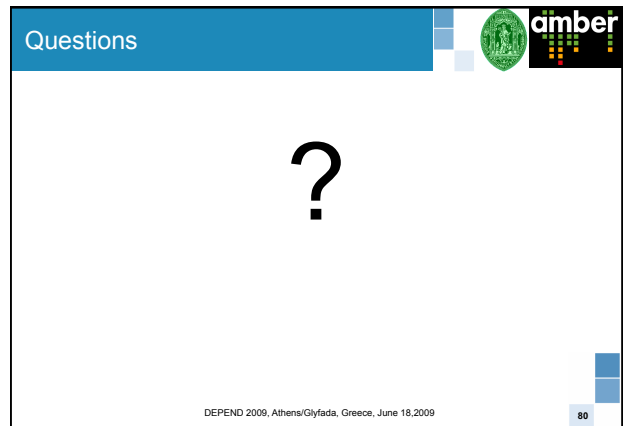
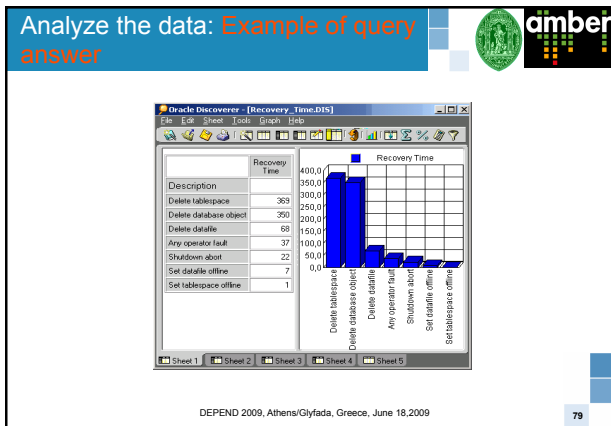
- Performance and recovery results
  - Per experiment
  - Per SUT
  - Per workload
  - Per fault type

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009



### Analyze the data: Example of query construction

Screenshot of Microsoft Excel 'Edit Worksheet' dialog box. The 'Available' list contains 'Fault Load', 'Description', 'Type', 'Injection Time', 'Experiment Duration', 'Tpm', 'Recovery Time', 'SUM'. The 'Selected' list contains 'Fault Load', 'Description', 'Measurements', 'Recovery Time', 'AVG'.



## 4. The AMBER Data Repository

**AMBER Repository vision and objectives**

- Vision
  - Become a worldwide repository for dependability related data
- Key objectives:
  - Provide state-of-the-art data analysis
  - Allow data comparison and cross-exploitation
  - Facilitate worldwide data sharing and dissemination
- Potential tool to increase the impact of research

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Potential use**

- Research team level
  - Perform the analysis of data in an efficient way
  - Efficient dissemination of the results of the team
- Project level
  - Sharing and cross-exploitation of results from different project teams
- World wide
  - Common repository to store and share data
  - Many teams are performing dependability evaluation but there are no results available at the web

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Data analysis approach**


- Repository to analyze, compare, and share results
- Use a business intelligence approach:
  - Data warehouse to store data
  - On-Line Analytical Processing (OLAP) to analyze data
  - Data mining algorithms to identify (unknown) phenomena in the data
  - Information retrieval to access data in textual formats
- Adopt the same life cycle of BI data
- Use technology already available for DW, DM & IR

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009

**Steps**

1. User registration
2. Multidimensional analysis
3. Definition of the loading plans
7. Load the data
8. Definition of data ownership policies
9. Analysis of the data

- Analyze DBench-OLTP results using OLAP



DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 85

**User registration**

- ADR users must undergo a registration procedure
- Provide identification information that is verified by the ADR support team
  - To filter malicious users
- Contact information is used to get in touch with the potential repository user
- To access the repository users must authenticate

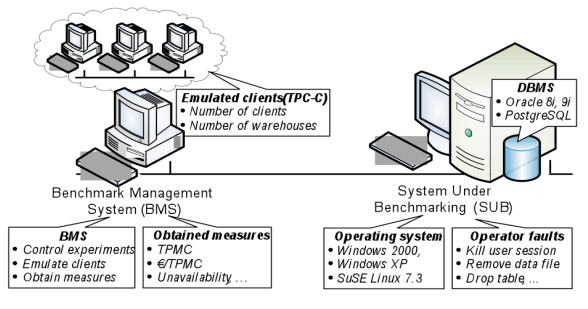
DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 86

**Multidimensional analysis**

- Design an adequate multidimensional data model
- User has the required expertise to design the data model ☺
  - Send to the ADR support team the SQL scripts needed to create the database tables
- The ADR team helps the user defining the model
  - The user only needs to explain us the experimental setup and the format of the data collected

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 87

**The DBench-OLTP benchmark**



DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 88

**Format of the raw data**

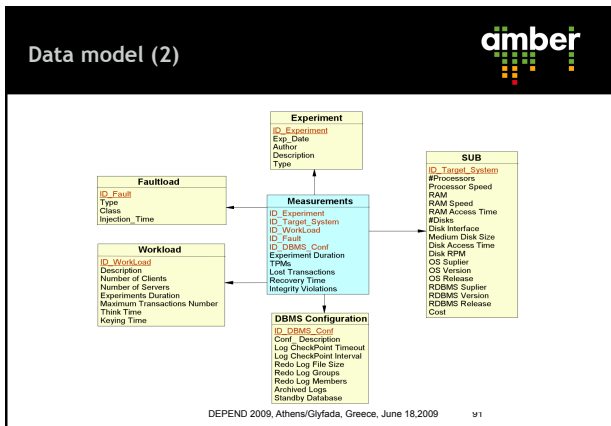
- Raw data collected by DBench-OLTP is composed of tens of CSV files (one from each run)
- Each row contains data from an injection slot
  - Identification, duration, number of transactions executed, data integrity errors discovered, type of fault injected, moment of fault injection, workload used, etc)
- A text file describes the experiment and the characteristics of the SUB

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 89

**Data model (1)**

- Key steps:
  - Identification of the facts that characterize the problem under analysis
  - Identification of the dimensions that may influence the facts
  - Definition of the granularity of the data stored in the star schema

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 90



- ### Definition of the loading plans
- Data extraction
    - SQL scripts to extract data from the CSV files to a temporary database schema (data staging area)
  - Data transformation
    - SQL scripts transform the data into an adequate format
  - Data load
    - SQL scripts to load the transformed data into the data warehouse
  - Loading plans documented and stored in the ADR
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 92

- ### Load the data
- Executing the loading plans created before
  - If new data becomes available we just need to rerun the plans
    - e.g., if the benchmark is executed in other systems
  - The documentation of the DBench-OLTP includes papers and technical reports
    - This is considered as part of the DBench-OLTP data
    - It is loaded to the repository and made available to the potential readers of the data
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 93


- ### Data ownership policy
- Data ownership policies of ADR are divided in two main groups
    - Private data
    - Proprietary data
    - Collaborative data
  - For the DBench-OLTP data we have decided to use a collaborative approach
    - Allows other potential users of the benchmark to compare their results with the ones available in the ADR
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 94

- ### Analysis of the data
- On-line Analytical Processing (OLAP) tools
    - Support the analysis in a very flexible way
    - Provide high query performance and easy, intuitive data navigation
  - Oracle Business Intelligence Discoverer Plus (ODP)
    - Commercial tool included in Oracle Business Intelligence package
    - Widely used by industry Used freely for research purposes under an Oracle Academy Agreement
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 95

- ### OLAP Wizard
- Selection of query type (crosstab or table) and characteristics (title, graph, text area, etc)
  - Selection of measures and dimensional attributes
  - Setting the query layout
  - Selection of the fields to be used to sort the results
  - Creation of parameters used to filter data
- DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 96



### Some results



|                        | Oracle | PosgreSQL |           | Integrity Errors |
|------------------------|--------|-----------|-----------|------------------|
| Tpmc with Faults       | 1263   | 644       | Oracle    | 0                |
| Étpmc                  | 20     | 7         | Microsoft | 0                |
| Server Unavailability  | 212    | 460       | RedHat    | 0                |
| Clients Unavailability | 13341  | 23240     | PosgreSQL | 0                |
|                        |        |           | RedHat    | 0                |


|           | Server Unavailability |        | Clients Unavailability |        |
|-----------|-----------------------|--------|------------------------|--------|
|           | Microsoft             | RedHat | Microsoft              | RedHat |
| Oracle    | 194                   | 264    | 11473                  | 19077  |
| PosgreSQL |                       | 460    |                        | 23240  |


|        | Tpmc With Faults |        |
|--------|------------------|--------|
|        | Microsoft        | RedHat |
| Oracle | 1271             | 1240   |

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 97

### Quick demo...



- Murphy's law... ☺



DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 98

<http://www.amber-project.eu>




## Do you have data?

## Share Them!

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 99


### Questions



# ?

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 100


### Generic bibliography



- Ralph Kimbal, Margy Ross, "The Data Warehouse Toolkit: The Complete Guide to Dimensional Modeling" (Second Edition), Ed. J. Wiley & Sons, Inc, 2002.
- Ralph Kimbal, "The Data Warehouse Lifecycle Toolkit", Ed. J. Wiley & Sons, Inc, 2001.


DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 101

### ADR bibliography



- Madeira, H., Costa, J., Vieira, M., "The OLAP and Data Warehousing Approaches for Analysis and Sharing of Results from Dependability Evaluation Experiments", International Conference on Dependable Systems and Networks, DSN-DCC 2003, San Francisco, CA, USA, June 2003
- Pintér, G., Madeira, H., Vieira, M., Pataricza, A., Majzik, I., "A Data Mining Approach to Identify Key Factors in Dependability Experiments", Fifth European Dependable Computing Conference (EDCC-5), Budapest, Hungary, April 2005

DEPEND 2009, Athens/Glyfada, Greece, June 18, 2009 102

ADR bibliography 

- Pintér, G., Madeira, H., Vieira, M., Majzik, I., Pataricza, A. , "Integration of OLAP and Data Mining for Analysis of Results from Dependability Evaluation Experiments", International Journal of Knowledge Management Studies (IJKMS), Volume 2 – Issue 4 – 2008, Inderscience Publishers, July 2008
- Vieira, M., Mendes, N., Durães, J., Madeira, H. , "The AMBER Data Repository", DSN 2008 Workshop on Resilience Assessment and Dependability Benchmarking (DSN-RADB08), Anchorage, Alaska, June 2008
- Vieira, M., Mendes, N., Durães, J. , "A Case Study on Using the AMBER Data Repository for Experimental Data Analysis", SRDS 2008 Workshop on Sharing Field Data and Experiment Measurements on Resilience of Distributed Computing Systems, Naples, Italy, October 2008

DEPEND 2009, Athens/Cityfada, Greece, June 18, 2009

103