

Home / Information Systems

etl-architecture.pdf

 School

Universidad TecMilenio*

*We are not endorsed by this school

 Course

CS MISC

 Pages

76

 Upload Date

Jan 29, 2024

Uploaded by pep.perez on coursehero.com

 Helpful  Unhelpful

ETL ARCHITECTURE

↑ ↓ Page 1 of 76



**"I wish I didn't
wait so long."**

APPLY NOW

ETL Architecture Options - Overview

Search for

12/08/2004

ETL Architecture Options - Guidelines

be tool dependant, client environment dependant or it can be based on the best practices available. Explain different plug-ins available with the ETL tool. Explain the command language used by the ETL tool.

- Provide different option of architecture by which the tool can be configured and implemented effectively in the organization, the option may be based on setup of repository system. Provide pros and cons for each of the options. Limit the number of option to max 4 and select the one appropriately which will suite the client's requirement and environment.
- Make sure that you are handling the geographically distributed architecture is right. Make sure that the architecture suggested is flexible enough to expand as needed.
- Provide different option for Disaster recovery solutions and Highly Availability solutions

12/08/2004

ETL Architecture Options - Guidelines

• **Architecture Option 1**

Explain the architecture with the diagram and provide pros and cons. Explain how this option can be implemented in client environment. Probably this can be a distributed architecture (explained in ETL Framework session)

• **Architecture Option 2**

Explain the architecture with the diagram and provide pros and cons. Explain how this option can be implemented in client environment. Probably this can be a centralized architecture (explained in ETL Framework session)

• **Architecture Option 3**

- **Recommendation for ETL Infrastructure**

Explain the recommended architecture option for ETL infrastructure and justification. Explain how client can be maximum benefited from the recommended architecture option in their current environment. Also provide small roadmap client regarding how they should grow with this option.

ETL Design Considerations

- Modularity
- Consistency
- Flexibility
- Speed
- Heterogeneity
- Meta Data Management

12/08/2004

ETL Design Considerations

- **Modularity**

ETL systems should contain modular elements which encourages reuse them easy to modify when implementing changes

Flexibility

ETL systems may be appropriate to accomplish some transformations in the target system and some on the source data system; others may require the development of custom applications

ETL Design Considerations

- **Speed**

ETL systems should be as fast as possible.

- **Heterogeneity**

ETL systems should be able to work with a wide variety of data in different formats.

ETL Architectures - Types

12/08/2004

ETL Architectures - Types

▼ **Based on the TWO OF THE PROCESSES.**

- Traditional Architecture
- Conformed Architecture

12/08/2004

ETL Architectures - Types

Homogenous Architecture

A homogenous architecture for an ETL system is one that involves only source system and a single target system

Features

Single data source

Data is extracted from a single source system, such as an OLTP

type.



Light data transformation

No data transformations are required since the incoming data is format usable in the data warehouse.

12/08/2004

- **Light structural transformation**

Because the data comes from a single source, the amount of structural changes as table alteration is also very light

- **Simple research requirements**

The research efforts to locate data are generally simple: if the data is in the source System, it can be used. If it is not, it cannot

12/08/2004

Heterogeneous Architecture

A heterogeneous architecture for an ETL system is one that extracts data from multiple sources



Features



Multiple data sources



More complex development

The development effort required to extract the data is increased because there are multiple source data formats for each record type.



Significant data transformation

Data transformations are required as the incoming data is often in an incompatible format usable in the data warehouse.

12/08/2004

Heterogeneous Architecture



Substantial research requirements to identify and match data

Heterogeneous Architecture

ETL Architectures - Types

Traditional ETL Architecture

Create individual ETL processes for each source system

A traditional ETL architecture would create one ETL process to perform logic necessary to transform the source data into its target destination.

ETL Architectures - Types

Traditional ETL Architecture

Advantages

- There are fewer ETL processes to create and maintain

Disadvantage

- Here each individual ETL process redundantly perform many of the same tasks.
- During development, this results in an increase in the amount and complexity of ETL code to create and test.
- Upon Implementation, modifications require more effort as changes must be made in multiple places

12/08/2004

ETL Architectures - Types

Conformed ETL Architecture

12/08/2004

Conformed table

Advantages

- **Modularization of ETL processes**

The creation of smaller, less complex ETL processes makes troubleshooting problems and creating future enhancements easier.

- **Reusability of post-conform processes:**

- Enforcing referential integrity (looking up foreign key assignments) performing inserts and/or updates to the final target is made simpler.

12/08/2004

- Extensibility

Confirmed architecture allows us to add new source systems.

- Disadvantages

- There are more objects and processes that must be created and

12/08/2004

Based on Repository

- 3 types of ETL Architectures based on repository
 - Single Domain Distributed Repository Architecture
 - Single Domain Single Repository Architecture
 - Distributed Domain Architecture

12/08/2004

Single Domain Distributed Repository Architecture

Single Domain Single Repository Architecture

Distributed Domain Architecture

12/08/2004

12/08/2004

The ETL process flow Types

- TRANSFORM THEN LOAD
- LOAD THEN TRANSFORM
- TRANSFORM WHILE LOADING

12/08/2004

The ETL process flow Types

TRANSFORM THEN LOAD

The data is manipulated outside the database to cleanse and sort it, with loaded into the database

12/08/2004

Transform and load are

- If the data is transformed outside the database, the external tools used to do this may not scale as effectively as the database does and will become a bottleneck.
- Depending on the architecture, the external mechanism has to control the flow of the ETL process and provide recovery and restart ability.

12/08/2004

The ETL process flow Types
LOAD THEN TRANSFORM

The main risks and disadvantages of Load then Transform are

- Here, extra disk storage for the staging tables is required.
- The transformation process is interrupted by storing not only intermediate but also the original raw data from the source systems in the Database.

The ETL process flow Types

TRANSFORM WHILE LOADING

The raw data is selected directly from a stream of data from the product or flat files, transformed by applying one or more table functions to it, and written to the database

12/08/2004

ETL Extraction Methodology

12/08/2004

Important considerations for extraction

- The extraction method to choose is highly dependent on the source system from the business needs in the target data warehouse environment.
- The estimated amount of the data to be extracted and the stage in the ETL also impact the decision of how to extract, from a logical and a physical

12/08/2004

Methods to Extract data

- Logical Extraction Method

 - There are two kinds of logical extraction:

 - Full Extraction
 - Incremental Extraction

- Physical Extraction Method

 - There are two kinds of physical extraction

 - Online Extraction
 - Offline Extraction

LOGICAL EXTRACTION METHODS

- Full Extraction

- The data is extracted completely from the source system

- Incremental Extraction

- Only the data that has changed since a well-defined event backlog will be extracted.

- Two ways to accomplish

- Change data capture technique

12/08/2004

Entire tables extraction

- Change data capture technique

- Extract only the most recently changed data

- Tables from the source systems are extracted to the data warehouse staging area, and these tables are compared with a previous extract from the source system to identify the changes.

Physical Extraction Methods

- Online Extraction

- The data is extracted directly from the source system itself.

- Offline Extraction

- The data is not extracted directly from the source system but is explicitly outside the original source system

Extract data in two ways

- **Extraction Using Data Files**

- Most of the database systems provide mechanisms for exporting or unloading data from the internal database format into flat files

- **Extraction Via Distributed Operation**

- Using distributed-query technology, one database (oracle) can query tables located in various different source systems, such as another database (oracle)

12/08/2004

Transportation

- Transportation is the operation of moving data from one system to another.
- The most common requirements for transportation are in moving data from:
 - A source system to a staging database or a data warehouse database
 - A staging database to a data warehouse
 - A data warehouse to a data mart

12/08/2004

Three basic choices for transporting data in warehouses

- Transportation Using Flat Files

- Transportation through Distributed Operations
- Transportation Using Transportable Table Spaces

12/08/2004

● Transportation Using Flat Files

- The most common method for transporting data is by the transmission of flat files, using mechanisms such as FTP or other remote file system protocols.

Advantages

- Source systems and data warehouses use different operating systems and database systems, using flat files is the simplest way to transport data between heterogeneous systems with minimal transformation.
- When transporting data between homogeneous systems, flat files are often the most efficient and most easy-to-manage mechanism for data transfer.

12/08/2004

Transportation through Distributed Operations

- Distributed queries, either with or without gateways, can be an effective method for extracting data.

Transportations Using Transportable Table Spaces

- Using transportable table spaces, Oracle data files (containing table data and almost every other Oracle database object) can be directly transported from one database to another.



Disadvantage



Source and target systems must be running Oracle8i (or higher). They must be running the same operating system, must use the same character set, and must have the same memory architecture.

Loading Methodology

12/08/2004

Loading Mechanisms

- SQL*LOADER
- External Tables
- OCI and Direct Path API's
- Export /Import

12/08/2004

Staging Area STAGING AREA

12/08/2004

Definition

- A place where raw data is brought in, cleaned, combined, archived, and one or more data marts.
- It is also used to get data ready for loading into a presentation server.

12/08/2004

Area.

12/08/2004

Data Staging Area Roles

- Integrates data from many application source systems so there is one central System wide enterprise view of the data.

Characteristics of Staging Area

- Facilitates moving data from different sources on different schedules.
- Provides a place to check data cleanliness and correctness.
- Is the enterprise-wide integration of data.
- Data is stored at the lowest level of detail available.

Need For Staging Area

- Data used in the data warehouse is extracted from the data sources, cleaned and transformed into the data warehouse schema.
- The data is checked for consistency and referential integrity.
- Promotes effective data warehouse management.
- Data transformation in the data source systems can interfere with OLTP performance.
- Allows DW professionals to assess the data quality problems before the data is loaded to the warehouse.

12/08/2004

12/08/2004

Creating Staging Area

- Create tables and other database objects to support
 - the data extraction
 - cleansing
 - transformation operations required to prepare the data for load data warehouse.

12/08/2004

Creating Staging Area

- It should include tables to contain the
 - incoming data
 - tables to aid in implementing surrogate keys
 - tables to hold transformed data.
- Design will depend on
 - the diversity of data sources
 - the degree of transformation necessary to organize the data for loading
 - the consistency of the incoming data.

12/08/2004

Data Staging Techniques

- Surrogate key creation and maintenance
- Processing Slowly Changing Dimensions
- Combining from Separate Sources
- Data Cleaning
- Processing Names and Addresses
- Validating One-to-One and One-to-Many Relationships
- Fact Processing
- Aggregate Processing

12/08/2004

Data Staging Techniques

- Surrogate key creation and maintenance
 - create a surrogate key .
 - every DW key should be a surrogate key

Type 1: Overwrite the Value

"Rewriting History" - no history .

Type 2: Add A Dimension Row

"

Data Staging Techniques

- Combining from Separate Sources

- Dimensions are derived from several sources. the merge opera on the some criteria.

- Data Cleaning

- Data cleaning may involve checking the spelling of an attribut checking the membership of an attribute in a list

- Names and addresses have been cleaned and put into standard

Data Staging Techniques

- Fact Processing

- The incoming fact records will have production keys, not data keys. The current correct correspondence between a production data warehouse key must be looked up at load time

- Aggregate Processing

- Each load of new fact records requires that aggregates be calculated and augmented. It is very important to keep the aggregates synchronized with base data at every instant

12/08/2004

Data Staging Techniques

- A many-to-one relationship , (eg.zip code-to state) ,can be va
 sorting on the "many" attribute and verifying that each value ha
 value on the "one" attribute.

12/08/2004

Data Staging Storage Types

- Flat Files
- Relational tables

12/08/2004

Is Data Staging Relational ???

- whether the data staging area is relational or has more to do with sequential processing of flat files. Ralph Kimball [1] concludes that

"Most data staging activities are not relational, but rather they are sequential processing. If your incoming data is in flat-file format you should finish staging processes as flat files before loading it into a relational database. It states that if both the source and target databases are relational it may be appropriate to retain this format and not convert to flat files "

12/08/2004

Data Staging Components

- Data Staging Application Server
- Data Staging Repository
- Metadata and Meta model Repository

12/08/2004

Data Staging Components

- **Data Staging Application Server**

 temporarily stores and transforms data extracted from OLTP da

.

12/08/2004
transformation and loading activity.

- The archival repository stores cleaned, transformed records and attributes for later loading into data marts and data warehouse schemas.

Data Staging Components

- Metadata and Meta Model Repository

- The data staging process is driven in an essential way by metadata including business rules.
- Metadata is used along with administrative tools to guide data transformations, archiving, and loading to target data mart and data warehouse schemas.

Staging scenarios

2 staging scenarios

Scenario 1- a data staging tool is available .

The data is already in a database. The data flow is set up so that it come source system, moves through the transformation engine, and into a stag database.

12/08/2004

Staging scenarios

12/08/2004

General Data Staging Requirements

- **Productivity support**

The data staging services needs to provide basic development environment capabilities like code library management check in/check out, version control, production and development system builds.

- **Usability**

The data staging system must be as usable as possible, it should have good documentation and easy to follow user interface.

12/08/2004

General Data Staging Requirements

- **System Documentation**

The data staging system need to provide a way for developers to easily information about the processes that they are creating.

- **Metadata driven**

Metadata is used along with administrative tools to guide data transformations, archiving, and loading to target data mart and dat schemas.

12/08/2004

Staging Metadata

12/08/2004

Data Staging Metadata

The metadata needed to get the data into a staging area and prepare it for loading into one or more data marts are:

- Data acquisition information

- Transformation and Aggregation

Audit, Job Logs and documentation

Data Staging Metadata

- Data acquisition information

Data transmission scheduling

File usage in the data staging area including duration, volatility, ownership.

Data Staging Metadata

● Dimension table Management

- Definitions of conformed dimensions and conformed facts.
- Job specifications for joining sources, stripping out fields, and attributes.
- Slowly changing dimension policies for each incoming descriptive attribute.
- Current surrogate key assignments for each production key, including a fast lookup table to perform this mapping in memory.
- Yesterday's copy of production dimensions to use as the basis for comparison.

12/08/2004

Data Staging Metadata

- DBMS load scripts.
- Aggregate definitions.
- Aggregate usage statistics ,base table usage. statistics and potential aggregates.
- Aggregate modification logs.

12/08/2004

Data Staging Metadata

• Audit, Job Logs and documentation

- Data lineage and audit records (where exactly did this record come from and when ?)
- Data transformation run time logs, success summaries and timing information.
- Data transformation software version numbers.
- Security settings for extract files, extract software , and extraction methods.
- Security settings for data transmission (e.g. passwords ,certificates).
- Data staging area archive logs and recovery procedures.
-

12/08/2004

Data Staging Metadata

● DBMS metadata

- DBMS system table contents.
- Partition settings.
- Indexes.
- Disk stripping specifications.
- Processing hints.
- DBMS-level security privileges and grants.
- View definition.
- Stored procedures and sql administrative scripts.
- DBMS backup ,backup procedures and backup security.

12/08/2004

Data Staging InvData

● Front room metadata

- Business names and descriptions for columns, tables and groups.
- query and report definitions.
- Join specification tool settings.
- Network security user privilege profiles.
- Usage and access maps for data elements, tables, views and reports.

12/08/2004

Uploaded by pep.perez on coursehero.com

Terms

Academic Integrity

Privacy Policy

Cookie Policy

Do not Sell or Share My Personal Info

SUBJECTS

Accounting

Aerospace Engineering

Anatomy

Anthropology

Arts & Humanities

Astronomy

Biology

Business

Chemistry

Civil Engineering

Computer Science

Communications

Economics

Electrical Engineering

English

Finance

Geography

Geology

Health Science

History

Industrial Engineering

Information Systems

Law

Linguistics

Management

Marketing

Material Science

Mathematics

Mechanical Engineering

Medicine

Nursing

Philosophy

Physics

Political Science

Psychology



SOCIAL



© Learneo, Inc. 2024

*College Sidekick is not sponsored or endorsed by any college or university.