# A Technical Comprehensive Survey Of ETL Tools

**[1]Vaishali A. Kherdekar and *[2] Pravin S. Metkewar**

[1]*Assistant Professor, MIT Art's Commerce & Science College (MIT ACSC), Alandi,*
*Affiliated to Savitribai Phule Pune University Pune, Maharashtra State, India*
*Email:vakherdekar@mitacsc.ac.in(1)*
*[2]Associate Professor, Symbiosis Institute of Computer Studies and Research(SICSR),*
*Affiliated to Symbiosis International University(SIU), PUNE-411016, Maharashtra State, India*
*Email:pravin.metkewar@sicsr.ac.in(2)*

## Abstract

In modern days ETL tools are vey useful in data integration and data warehousing. Input is given to to the datawarehouse through ETL. ETL means Extraction, Transformation and Loading. ETL tools transfer data from one source system to another source system.

As these tools are mainly used in Busineness Intelligence and Data Warehousing, there is lot of space for their gress.There are lot of ETL tools are available in the market rying from version to version to stay proficient against other tools. Each and every tool has its own features and limitations. In this paper we have carried out technical survey of existing ETL tools and benchmarking of these tools has been performed by considering certain parametrs including scalability, reusability, interoperability, support to big data, parallelism, usability, flexibility etc. Finally, problems and challenges of ETL tools have been discussed thoroughly and its state of the art is summarized.

Keywords : ETL tools, Data warehouse

## Introduction :

Now a days data warehouse is used in industry to maintain optimiza model of data for further mining and uage and also for report generation. By using datawarehouse one can maintain historical data and used it in decision support system. To construct dara warehouse model, ETL tools is being used. ETL tools act as basis for construction of data ehouse. Input is given to the data warehouse through ETL. L stands for Extraction, Transformation and Loading. In extraction phase, data is extracted from various heterogeneous sources in different formats such as flat files, databases, xml files etc; it means extraction of data is achieved with the help of structured and semi-structured databases and files. In transformation phase extracted data is transformed into specific format for data analysis. In loading transformed data is loaded into datawarehouse. Now a days, a large number of ETL tools are available in the market. However, in general, they follow different design and modeling techniques, and uses different internal language.

## Sample selection of ETL Tools:

In market variety of ETL tools are exists either that are from open source or proprietary tool. We have considered few reputed tools for our study in order to perform benchmarking

of these tools that are including Pentaho, Clover ETL, Rapid Miner, Jedox, JasperSoft and Talend tools.

Pentaho ETL tool was established by the Pentaho Corporation, United States. In market it is named as Pentaho Kettle. It is open source and provides services for business intelligenve and data integration. The transformations carried out in Pentaho is stored in XML.It is executed in Java.

CloverETL was developed by JavlinInc in 2002. It is Java based data integration tool used for transforming and distributing data and other functionalities of data warehousing. It can be used as either standalone or embedded to server application. It works on different platforms.

Rapidminer tool is specially used for Regression Analysis, Gaussian process And statistical process. It provides a wide range of functionalities and support.

Jaspersoft extract and transform data from multiple systems and loads it into data store for reporting and analysis. It works as ETL job designer having data integration capabilities.

Talend was introduced in 2006. It is also open source Java based tool used for data integration and data analysis process.

## Benchmarking of ETL tools :

### Table 1. Benchmarking of ETL Tools

| Parameters | Pentaho | Talend | Clover ETL | Jedox | Jaspersoft | RapidMiner |
|---|---|---|---|---|---|---|
| Usability | Very good | User friendly | User friendly | User friendly web interface | User friendly | User friendly |
| Reusability | Yes | Yes | Yes | Yes | Yes | NA |
| Interoperability | Yes | Yes | Yes | Yes | Yes | Yes |
| Scalability | Cluster (carte server). | Highly scalable | Scalable | NA | Yes | NA |
| Flexibility | more flexibl | | Yes | more flexibl | less flexible | |