



**HAL**  
open science

# Data Warehousing Process Modeling from Classical Approaches to New Trends: Main Features and Comparisons

Asma Dhaouadi, Khadija Bousselmi, Mohamed Mohsen Gammoudi, Sébastien Monnet, Slimane Hammoudi

## ► To cite this version:

Asma Dhaouadi, Khadija Bousselmi, Mohamed Mohsen Gammoudi, Sébastien Monnet, Slimane Hammoudi. Data Warehousing Process Modeling from Classical Approaches to New Trends: Main Features and Comparisons. *Data*, 2022, 7 (8), pp.113. 10.3390/data7080113 . hal-03758493

**HAL Id: hal-03758493**

**<https://hal.science/hal-03758493>**

Submitted on 23 Aug 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Article

# Data Warehousing Process Modeling from Classical Approaches to New Trends: Main Features and Comparisons

Asma Dhaouadi <sup>1,2,3,\*</sup> , Khadija Bousselmi <sup>2</sup>, Mohamed Mohsen Gammoudi <sup>1,4</sup> and Sébastien Monnet <sup>2</sup> and Slimane Hammoudi <sup>5</sup>

<sup>1</sup> RIADI Laboratory, University of Manouba, Mannouba 2010, Tunisia

<sup>2</sup> LISTIC Laboratory, University of Savoie Mont Blanc, France Annecy-Chambéry, 74940 Chambéry, France

<sup>3</sup> Faculty of Sciences of Tunis, University of Tunis El Manar, Tunis 1068, Tunisia

<sup>4</sup> Higher Institute of Arts and Multimedia Manouba, University of Manouba, Manouba 2010, Tunisia

<sup>5</sup> ERIS, ESEO-Grande Ecole d'Ingénieurs Généralistes, 49100 Angers, France

\* Correspondence: asma.dhaouadi@univ-smb.fr or asma.dhaouadi@fst.utm.tn

**Abstract:** The extract, transform, and load (ETL) process is at the core of data warehousing architectures. As such, the success of data warehouse (DW) projects is essentially based on the proper modeling of the ETL process. As there is no standard model for the representation and design of this process, several researchers have made efforts to propose modeling methods based on different formalisms, such as unified modeling language (UML), ontology, model-driven architecture (MDA), model-driven development (MDD), and graphical flow, which includes business process model notation (BPMN), colored Petri nets (CPN), Yet Another Workflow Language (YAWL), CommonCube, entity modeling diagram (EMD), and so on. With the emergence of Big Data, despite the multitude of relevant approaches proposed for modeling the ETL process in classical environments, part of the community has been motivated to provide new data warehousing methods that support Big Data specifications. In this paper, we present a summary of relevant works related to the modeling of data warehousing approaches, from classical ETL processes to ELT design approaches. A systematic literature review is conducted and a detailed set of comparison criteria are defined in order to allow the reader to better understand the evolution of these processes. Our study paints a complete picture of ETL modeling approaches, from their advent to the era of Big Data, while comparing their main characteristics. This study allows for the identification of the main challenges and issues related to the design of Big Data warehousing systems, mainly involving the lack of a generic design model for data collection, storage, processing, querying, and analysis.

**Keywords:** ETL process; data warehouse; ETL modeling; Big Data; UML; BPMN; ontology; MDA; graphical flow; systematic review



**Citation:** Dhaouadi, A.; Bousselmi, K.; Gammoudi, M.M.; Monnet, S.; Hammoudi, S. Data Warehousing Process Modeling from Classical Approaches to New Trends: Main Features and Comparisons. *Data* **2022**, *7*, 113. <https://doi.org/10.3390/data7080113>

Academic Editor: Kesheng Wu

Received: 5 May 2022

Accepted: 28 July 2022

Published: 12 August 2022

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

The globalization and the spread of information technology, the strong concurrency between different companies, and the urge for quick and easy access to reliable and relevant information have incited business leaders to replace traditional business computing systems with other decision support systems. These business intelligence (BI) decision systems appeared with the introduction of the data warehouse (DW) by Bill Inmon in 1991. According to [1], "A data warehouse is a subject-oriented, integrated, time-variant and non-volatile collection of data in support of management's decision-making process". A DW is a system used for integrating, storing, and processing data from often heterogeneous data sources, in order to provide decision-makers with a multi-dimensional view. The integration of these data is achieved through a three-phase process: extract, transform, and load (ETL). This process is responsible for extracting data from different data sources, transforming them (by preparation, conversion, clean, filter, conversion, join, aggregation, and so on), and loading them into a DW. Consequently, the general framework for ETL processes