

Advancing Logistics 4.0 with the implementation of a Big Data Warehouse: a Demonstration Case at the Automotive Industry

Nuno Silva ^{1,*} , Júlio Barros ¹ , Maribel Y. Santos ¹ , Carlos Costa ¹ , Paulo Cortez ¹ , M. Sameiro Carvalho ¹ , João N. C. Gonçalves ¹ 

¹ ALGORITMI Research Centre, University of Minho, Guimarães 4800–058, Portugal

* Correspondence: nuno.silva@dsi.uminho.pt

Abstract: The constant advancements in Information Technology have been the main driver of the Big Data concept's success. With it, new concepts like Industry 4.0 and Logistics 4.0 are rising. Due to the increase in data volume, velocity, and variety, organizations are now looking to their data analytics infrastructures and searching for approaches to improve their decision-making capabilities, in order to enhance their results using new approaches such as Big Data and Machine Learning. The implementation of a Big Data Warehouse can be the first step to improve the organizations' data analysis infrastructure and start retrieving value from the usage of Big Data technologies. Moving to Big Data technologies can provide several opportunities for organizations, such as the capability of analysing an enormous quantity of data from different data sources in an efficient way. However, at the same time, different challenges can arise, including data quality, data management, lack of knowledge within the organization, among others. In this work, we propose an approach that can be adopted in the logistics department of any organization in order to promote the Logistics 4.0 movement, while highlighting the main challenges and opportunities associated with the development and implementation of a Big Data Warehouse in a real demonstration case at a multinational automotive organization.

Keywords: Big Data; Data Warehouse; Logistics 4.0; Industry 4.0; Implementation.

1. Introduction

The explosion of the Information Technologies area has been the driver that launched new concepts such as Big Data and Industry 4.0 into the spotlights. The concept of Industry 4.0 relies in the digitization of the production systems to provide the capability of producing customized products within a short time and with costs similar to mass production scenarios [1]. This factor has a tremendous impact in the organizations' logistics due to the need of reacting to the sudden changes made by the customers.

Logistics 4.0 can be defined as "... the logistical system that enables the sustainable satisfaction of individualized customer demands without an increase in costs and supports this development in industry and trade using digital technologies" [2]. Such initiative is needed to improve the link between the manufacturers and the customers, in order to avoid failures in the manufacturing system [2].

Big Data technologies, with their capability of analysing massive volumes of diverse data flowing at high velocity, has an important role in the implementation of these new concepts (Industry 4.0 and Logistics 4.0) and in the resolution of their main associated challenges [3].

With the implementation of Big Data technologies became possible to perform tasks that involves a massive quantity of data at high speeds such as providing a supply chain control with real-time data, inventory control and management, improving forecasting models, among others [1].

Along with the influence of concepts like Industry 4.0 and Logistics 4.0, the investments in Big Data technologies are being stimulated making them more stable and

Citation: Lastname, F.; Lastname, F.; Lastname, F. Title. *Electronics* **2021**, *1*, 0. <https://doi.org/>

Received:

Accepted:

Published:

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

Copyright: © 2021 by the authors. Submitted to *Electronics* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

39 mature, ready to be implemented inside the organizations and became part of their
40 business.

41 A vast range of organizations, from diverse types of business, are now trying
42 to evolve their data analyses infrastructures to this new era, advancing their Data
43 Warehouses (DWs) based on a more rigid data model to the new concept of Big Data
44 Warehouses (BDWs) with a more dynamic data model.

45 This work aims to demonstrate how the implementation of a Big Data Warehouse
46 (BDW) in a logistics context can drive forward the concept of Logistics 4.0 and improve
47 the organization performance. The contributions of this work are: Propose a general
48 approach that can be adopted in the logistics departments of several organizations;
49 Propose a logical and technological architecture that supports the BDW and data analysis;
50 Propose a data model for a logistics BDW; Demonstrate the challenges and opportunities
51 that emerge throughout the development and implementation of a BDW in the logistics
52 department.

53 A demonstration case will be presented, having been the same developed inside
54 of a multinational automotive organization by taking advantage of their existing data
55 platform.

56 This work is structured as follow: Section 2 provides the published works related to
57 BDW and their architectures; Section 3 presents the suggested architecture to solve this
58 problem; Section 4 describes the organization reality and the tasks performed to accom-
59 plish the goal; Section 5 presents the results accomplished followed by a discussion where
60 the challenges and opportunities are highlighted; Section 6 shows the final conclusions
61 and future work.

62 2. Related work

63 With the implementation of concepts like Industry 4.0 and Logistics 4.0, it becomes
64 important to endow the organizations' data analyses infrastructure with the capability
65 of retrieving, transforming and analysing massive amounts of data at high velocity.
66 Before the establishment of the Big Data concept, organizations had their data analyses
67 infrastructure based in DWs where the data model was rigid and structured in order to
68 provide the best performance when data were inquired.

69 Nowadays, Big Data technologies, due to their capacity for distributed processing
70 and storage, allow us to have more dynamic data models with less rigid structures,
71 maintaining high performance even with massive volumes of data.

72 To implement Big Data technologies, we can follow two different approaches: "the
73 lift and shift" and the "rip and replace". "The lift and shift" strategy means that we
74 replace or extend parts of the existing infrastructure with Big Data technology to improve
75 its capabilities and to solve specific problems. This may result in a use case approach
76 instead of a data-driven approach, which can lead to uncoordinated data silos. The "rip
77 and replace" approach means that the existing Data Warehouse (DW) is totally replaced
78 by Big Data technologies [4].

79 Independently of these two strategies, there are several architectures and technolo-
80 gies, that can be used to implement a BDW. The use of different types of Not Only SQL
81 (NoSQL) databases, such as document-oriented and column-oriented [5] or graph mod-
82 els [6] can be used to store the different types of data in the BDW. In the literature, we can
83 find different architectures that can be used in a BDW, such as the Lambda architecture
84 [7] and the NIST Big Data Reference Architecture (NBDRA) [8]. The Lambda architecture
85 has three layers and unifies, in a single software design pattern, the batch and real-time
86 data processing concerns. The three layers presented in the Lambda architecture are
87 batch processing, real-time computing, and a layer to query the data. This division
88 between batch processing and real-time processing allows differentiating data according
89 to their nature and relevance to the business. In this way, it is possible to immediately
90 process the data that is needed in time, while data that is only needed in the long run
91 can be processed later [7].

92 The NBDRA is presented by its authors as a common reference that can be im-
93 plemented using any Big Data technology or service provider. It is divided into the
94 following five components: System orchestrator; Data provider; Big Data application
95 provider; Big Data framework provider; and Data consumer. The system orchestrator
96 is the component that establishes the requirements for all the infrastructure, including,
97 among others, architectural design, business requirements, and governance. The data
98 provider is the component that makes data accessible through different interfaces. The
99 Big Data application provider deals with all the necessary tasks to manipulate data
100 through its life cycle. The Big Data framework provider consists of several services or
101 resources that are used by the Big Data application provider. The data consumer is the
102 entity that will take advantage of all the data processing made by the Big Data system
103 [9]. Using the NBDRA and the Lambda Architecture as a reference, Santos and Costa [9]
104 created an approach to develop BDWs.

105 Several examples demonstrate the capacity of Big Data technologies for improving
106 the analytical capabilities of organizations. Chou et al [10] propose a system architecture
107 based on Hadoop, Sqoop, Spark, Hive and Impala to analyse data from electrical grids.
108 Sebaa et al [11] present an architecture based on the Hadoop ecosystem and a conceptual
109 model to develop a BDW in the Healthcare field. Santos et al [12] present a demonstration
110 case where it was applied a Big Data architecture and a set of rules to evolve from a
111 traditional DW to a BDW.

112 These examples demonstrate how Big Data technologies can be used in collabora-
113 tion with traditional DWs or even replacing them, both aiming to improve the analytical
114 capabilities of the organizations.

115 3. Propose Architecture for a Logistics 4.0 Big Data Warehouse

116 In this section, it is presented the logical (3.1) and technological (3.2) architectures
117 that can be used to implement a BDW for the Logistics 4.0 movement.

118 3.1. Logical Architecture

119 The main goal of this BDW is to be an analytical repository containing a substantial
120 amount of data, in order to support the daily activities of the logistics decision-makers
121 in the logistics 4.0 era.

122 Two of the key factors in Logistics 4.0 are the real-time exchange of information
123 between all the actors in the supply chain and the real-time Big Data analytics of vehicles,
124 products and facilities location [3].

125 The exchange of information between all actors in the supply chain can originate
126 diverse data sources with different types of data that need to be stored and analysed in
127 one central repository in order to be easily accessible by practitioners. The same happens
128 with the real-time Big Data analytics of the diverse supply chain components (vehicles,
129 products, and facilities location). Considering this, the real-time characteristics can be
130 important, nevertheless, it is necessary to adapt to the organizational requirements. Real-
131 time analytics can be a different concept from one organization to other. For example,
132 for one organization, the requirements of real-time can be to have access to data in less
133 than ten seconds, but for other organization, it can be to access the data in less than
134 two minutes. Moreover, some organizations still do not have the need of creating an
135 architecture that takes into consideration the real-time requirements.

136 In our demonstration case, the organization does not have the requirement of
137 real-time analysis, so the architecture presented in Figure 1 does not incorporate that
138 component. Nevertheless, due to the relevance of real-time in Logistics 4.0, it may be
139 relevant to implement and validate that component in future work.

140 As can be seen in Figure 1, the logical architecture has the following components:

- 141 • Sandbox Storage: where the raw data is stored in a distributed file system before
142 any transformation. This component is divided into two layers: Update Layer and
143 Backup Layer. The Update Layer contains the up-to-date data retrieved from the

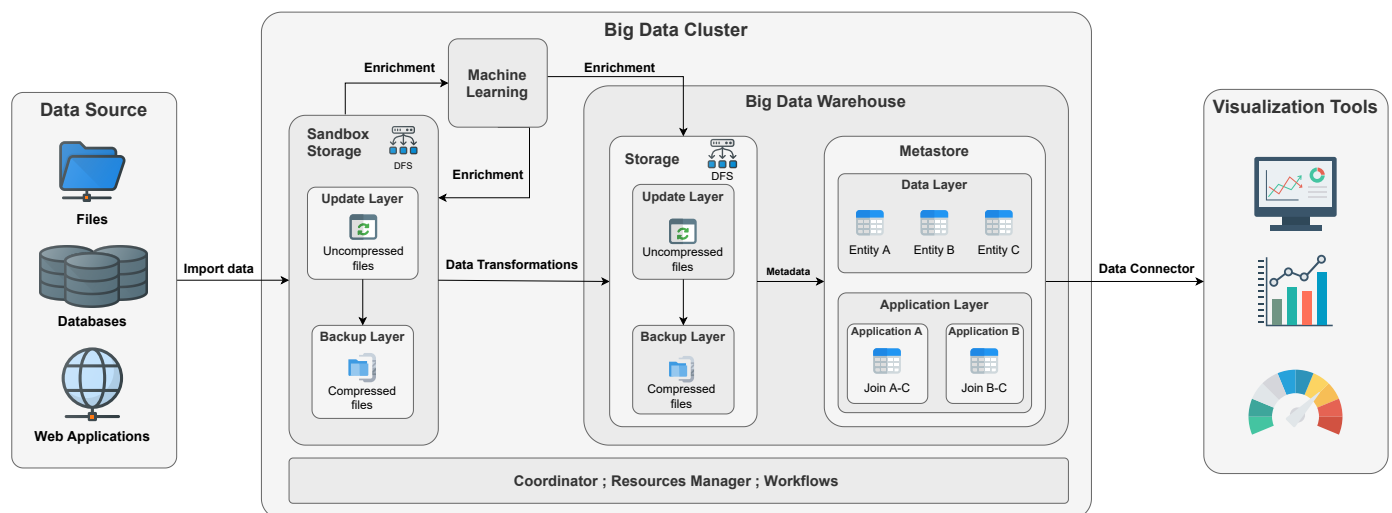


Figure 1. Logical architecture

- 144 sources, while the Backup Layer contains compressed outdated data to be used in
 145 case of necessity.
- 146 • **BDW Storage:** where data is stored in the distributed file system and accessible
 147 using the metastore after being transformed. This component has two layers with
 148 the same functionality as the Sandbox Storage layers: i) a layer that provides
 149 updated data, ii) and another layer to provide a backup in case of problems with
 150 the new data.
 - 151 • **Machine Learning component:** uses raw data from the Sandbox storage or clean
 152 data from the BDW to create predictions, in order to enrich the data and store
 153 it in the Sandbox Storage or in the BDW to provide predictive capabilities for
 154 the organization. This component can increase the organization's capabilities to
 155 understand and predict changes in their supply chain and be capable to adapt
 156 quickly.
 - 157 • **Metastore:** provides an interface to access the stored data. This component is
 158 divided into two layers: i) the data layer where the data is modelled using a
 159 data-driven approach, and; ii) the application layer where we have the necessary
 160 materialized objects or views to answer the needs of specific applications. The
 161 existence of these two layers provides some advantages. One of these advantages
 162 is the capability of creating several abstractions on top of the data layer, providing
 163 a simple and fast way to access the data. In this application layer, each application
 164 can have its own views or tables (materialized objects), increasing the performance
 165 when accessing the data. Moreover, if the organization has different teams working
 166 in different applications, if necessary, each team can create the necessary tables or
 167 views for their own application, providing higher business agility.
 - 168 • **The Coordinator, Resources Management and Workflows** provide functionalities to
 169 manage the Big Data Cluster and the data life cycle. The Coordinator and Workflow
 170 allow the creation of diverse jobs or tasks that can be submitted in the desired order.
 171 The Resource Manager distributes the clusters resources to process the jobs.

172 Outside the BD Cluster, we can find the data sources that provide the raw data to
 173 be used in the BDW and the Visualizations Tools where dashboards are developed to
 174 present the results to the users.

175 3.2. Technological Architecture

176 Due to the need of analysing big quantities of data in the most efficient way, new
 177 technologies that use the power of distributed processing and storage have gained
 178 significant attention. Probably the most well-known technology in this context, which

179 can arguably be seen as the originating driver of the Big Data movement, is Apache
180 Hadoop ¹, where data can be stored in the Hadoop Distributed File System (HDFS) [13]
181 and then processed using the Map and Reduce [14] programming model. Several other
182 technologies such as Sqoop ², Hive [15], Spark [16], and Impala ³ [17], among others, are
183 being constantly developed to tackle specific problems in the Big Data ecosystem. These
184 technologies allow the practitioners to retrieve data from the data sources, store it with
185 appropriate metadata and then processing it, in order to provide useful knowledge to
186 the end-users.

187 Currently, in the Big Data world, the amount of Big Data technologies is overwhelm-
188 ing and sometimes can be difficult to understand and choose the right technology for the
189 right job. For example, for data collection, technologies such as Flume, Kafka, or Talend
190 can be used. For data preparation and enrichment, we can use Spark or Storm. For data
191 storage, Hive with HDFS, NoSQL databases, or Kudu can be used. For machine learning
192 tasks, we can use Spark, H2O, and TensorFlow [18]. For query engines, Impala, Presto,
193 or Drill can be used. For data visualization, tools like Tableau, Power BI, JavaScript can
194 be used [19].

195 Due to the organizational requirements and due to the technologies available in
196 the organization depicted in this demonstration case, the technological architecture
197 presented in Figure 2 was used to support this demonstration case. Nevertheless this
198 technological architecture can be used inside others organization's logistics departments,
199 assuming the goals and requirements are similar to the ones depicted in this work. In
200 case of distinct requirements, some technologies could be adjusted. Regarding data
201 ingestion from the sources, this work uses Sqoop. Despite the fact that Sqoop can only
202 connect to structured databases [20], due to the fact that, for this demonstration case, the
203 organization's data sources were only SQL databases, there was no need to use another
204 technology to ingest the data. After the data is retrieved from the sources, the same
205 is stored in HDFS, using the Parquet format, which is one of the several formats that
206 can be used to store data in HDFS. Other formats that can be used are, for example,
207 ORC or AVRO [21]. Parquet was chosen not only due to its adequate compatibility
208 with Spark and Impala technology but also due to its read-oriented format and with
209 adequate compression, which will bring advantages when we need to query the data [22].
210 Moreover, it was necessary to develop a Bash script in order to provide a mechanism to
211 create data backups in the Sandbox Storage and in the BDW.

212 Spark was the chosen framework due to its data cleansing and transformation
213 capabilities and due to the capability to develop several machine learning models. Spark
214 has the SparkSQL [23] library that allows the use of SQL functions in conjunction with
215 the Spark programming API and complex libraries such as Spark MLlib [24]. Being able
216 to perform all these tasks in one unique framework is a significant advantage, since, in
217 this way, it is not necessary to spend more time using and learning different technologies.
218 Moreover, Spark is compatible with Parquet files and Hive, which will be used to provide
219 the data and metadata to the end-users.

220 Hive includes the Hive Metastore (the system catalog) where the metadata (schema
221 and other statistics) are stored, allowing proper data exploration and query optimizations
222 [15]. Hive allows the creation of external tables where data is stored in HDFS directories
223 and its life cycle is not managed by Hive [15]. Within Hive, we create two levels of
224 interaction with the data. In the first level, the data is modelled using a data-driven
225 approach where the core entities (such as Needs, Stocks, Products, among others) and
226 other entities like Date and Time are stored. This layer allows ad hoc access to the data
227 from these entities to be used by any team or project. In the second layer, the application
228 layer, a new set of objects (materialized tables or views), oriented to the applications'

¹ <https://hadoop.apache.org/>

² <https://sqoop.apache.org/>

³ <https://impala.apache.org/>

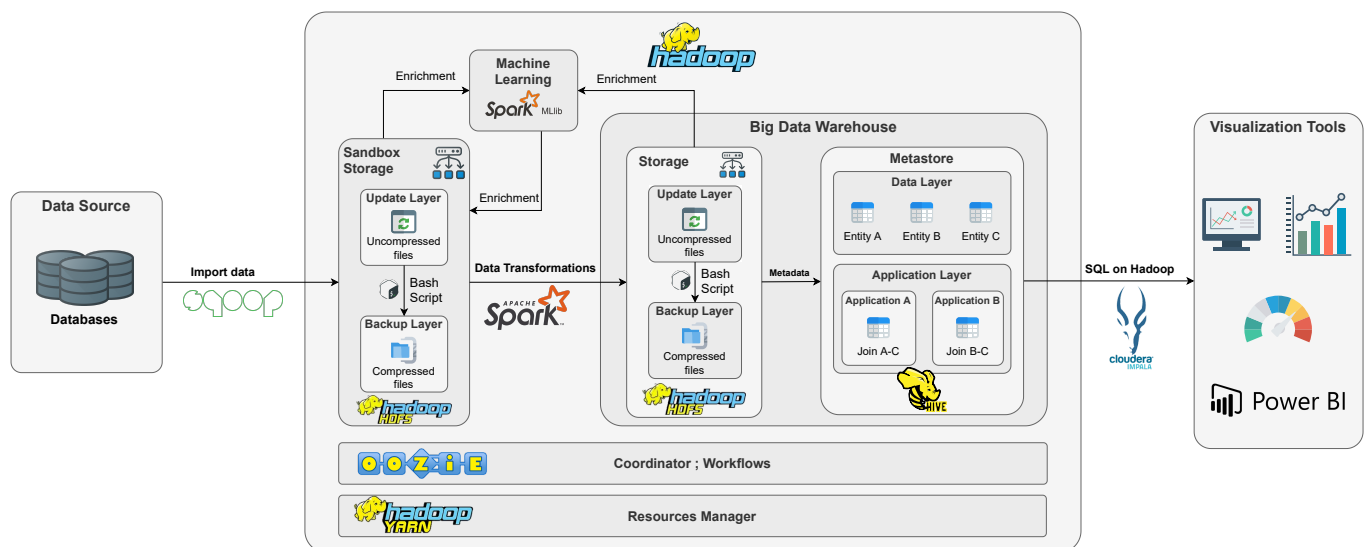


Figure 2. Technological architecture

needs, are created to provide access to the specific data that each application or project needs. This will provide more personalized access to the data that will increase the application performance and higher business agility, thus each team can create their own tables or views as they need.

Impala provides a massively parallel processing (MPP) SQL engine that combines the flexibility and scalability of Hadoop with the familiarity of SQL and has proven to be generally faster than Spark or Hive according to Qin et al [25] and to Bittorf et al [17]. Impala can too be used to query data from HBase and provide a connection to visualization applications, such as Tableau or Power BI, where dashboards can be developed to present to the end-user the knowledge retrieved from the data [17].

This technological architecture supports all the requirements of this project, granting that we can allow the data analysis team to provide knowledge to be used by the end-users, in order to support their decisions and therefore improving the organization's results. Moreover, it can be used in other Logistics 4.0 projects to create a new centralized repository that aggregates different data sources and requires predictive capabilities.

4. Demonstration Case

The application domain addressed in this paper is the Logistics Innovation Department of an automotive factory. In this context, the logistics department handles large volumes of data related to nearly 7000 raw materials from a set of about 400 suppliers spread all over the globe, which impact the production of about 1100 finished products. With regard to internal logistics management, the department is responsible for monitoring and analysing data and material movements referring to approximately 85 daily scheduled deliveries, in order to ensure the supply of material necessary for the proper functioning of about 100 production lines associated with various high-service level customers. In light of the complexity of the organization's supply chain topology, the organization intends to foster the proposal, development and evaluation of Big Data Analytics tools capable of integrating and automating a large part of the logistics processes that, until now, are managed by conventional spreadsheets extracted from classic and parameterizable material requirements planning (MRP) methodologies existing in a given enterprise resource planning (ERP) system.

It is an essential department inside of a production facility and deals on a daily basis with orders, deliveries, delays, production plans, inventory, among other processes. These business processes are crucial to maintain the production lines working and to deliver in time the finished goods to the clients. It is a complex and enormous department with countless business processes.

264 Due to this complexity, the implementation of a BDW needs to be addressed in an
 265 interactive way, choosing one process at a time, looking at the data sources, selecting the
 266 appropriated attributes and modelling the data in a data-driven approach that has as a
 267 final goal an integrated BDW supporting Logistics 4.0.

268 Therefore, in this specific case, to start the BDW proposal we analysed the processes
 269 that should be considered the core component of this BDW. With the collaboration of
 270 key experts in the logistics department, the following processes were selected: Product
 271 Inventory, Delivery, Purchase Order, and Needs. This is the first task in the development
 272 process presented in Figure 3

273 These processes will be the main drivers of the analytical objects in the BDW. Besides
 274 these objects, other objects will be created, such as a spatial object with information
 275 related to countries, Date and Time objects, and complementary analytical objects such
 276 as Product, Plant and Vendor. Each one of these processes are supported by one or
 277 more tables in the Enterprise Resource Planning (ERP) used by the organization. These
 278 different types of objects are explained later in this section.

279 The understanding and selection of the business processes, together with the under-
 280 standing and selection of the data sources compose the first activity of the development
 281 process (Figure 3) called Data Understanding. In this activity it is necessary to under-
 282 stand the data from the data sources, namely the tables associated to each business
 283 process, how they are related, their private and foreign keys, the meaning and possible
 284 values of each attribute, among other steps. The second task is to select what tables will
 285 be used to develop the BDW.

286 The next activity is related to the "Data Quality" activity. Data quality is one of the
 287 most important tasks in data-related projects. In this case, this activity has significant
 288 importance due to the complexity of the data sources and their high number of attributes.
 289 For example, some transactional tables have more than 200 attributes, although many
 290 of them are not used. In our demonstration case, data quality criteria were defined to
 291 verify if an attribute will be used in the BDW. In this specific case, we established that
 292 any attribute with more than 90% of empty or nulls values will not be used. This rule
 293 was essential to limit the number of used attributes, excluding the ones that have low
 294 analytical value. Another rule that was used was to manually verify if the attributes
 295 with only one or two distinct values were worth to use. All these rules were defined
 296 considering the organizational and decision-making context. The next step was to
 297 produce the data quality reports through the execution of several spark jobs that analysed
 298 the data extracted from HDFS. The attributes that will be part of the BDW are selected
 299 applying the previously defined data quality criteria.

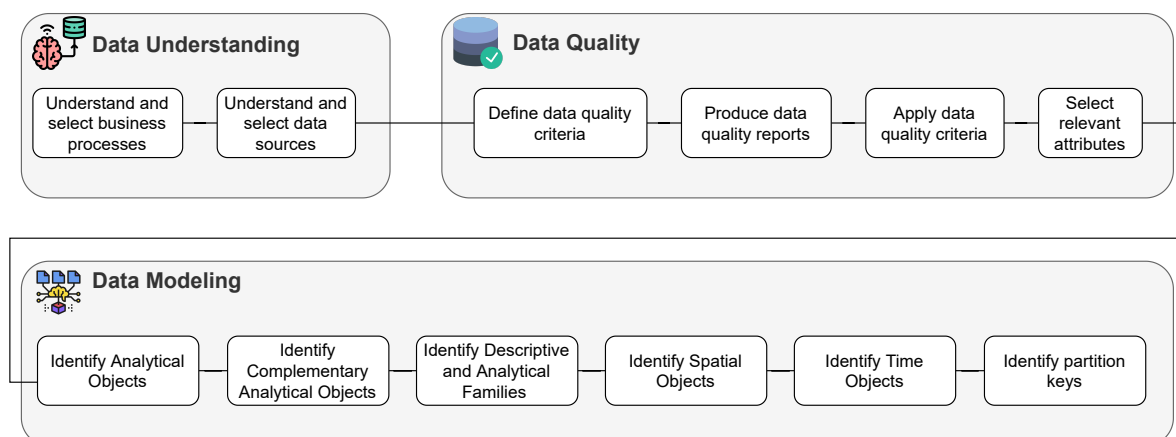


Figure 3. Development process

300 After the Data Understanding and the Data Quality, it was possible to model the
 301 BDW. To do that, the modelling methodology presented by Santos and Costa [9] was
 302 applied in order to propose a data model capable of integrate a significant amount of

303 data. The methodology is based on the creation of the following objects: Analytical
304 Objects, Complementary Analytical Objects, Spatial Objects, Time Object, and Date
305 Object.

306 An Analytical Object is a subject of interest, highly denormalized and that can
307 answer queries by itself avoiding joins with other objects. These objects are directly
308 related to the business processes such as sales or deliveries and should be the firsts to
309 be analysed and identified in order to verify if it is necessary, or not, to create Com-
310 complementary Analytical Objects. A Complementary Analytical Object is an object that
311 includes attributes usually used or shared by different Analytical Objects and that can
312 be used to complement the analysis of other objects, such as the Analytical Objects. Each
313 object can be divided into two distinct parts, the descriptive and analytical families.
314 These families provide a logical group for the object attributes depending on their type
315 and purpose. The descriptive family group all the attributes that can provide different
316 views of analysis, while the analytical family group the attributes with values to analyse
317 the object. These objects can be integrated with the use of join operations [9]. Figure
318 4 presents the data model identified with the application of this methodology. Due to
319 privacy concerns, it is only possible to disclose some of the attributes present in the
320 several objects This data model was developed in the logistics context of this specific
321 factory but can be used as starting point for any logistics department of any organization.

322 The Analytical Objects used in this work are: Product Inventory that has all infor-
323 mation about the stocks of each product; Deliveries that has information about when
324 each order is delivered; Purchase Order that has information about how many products
325 are ordered; and Needs that has information about production lines needs.

326 The Time and Date objects were created from scratch and populated with informa-
327 tion related to each one. For example, in the Date object, we created boolean attributes
328 such as week_day, weekend, summer, winter, monday, tuesday, and others. In the Time
329 object, attributes such as lunch-time, in-office, out-office, rush hour, were created. This
330 allowed us to analyse the relevant information and contextualize it in time and date.

331 The Complementary Analytical Objects had emerged in the data modelling process
332 due to the need of analysing different Analytical Objects using data from the Complemen-
333 tary Analytical Objects. In these objects were stored relevant and specific information
334 that can provide useful information when used together with data from several Analyti-
335 cal Objects. From these objects, we can highlight the following: Plant, Product Valuation,
336 and Vendor.

337 The object Country is a Spatial Object due to the geographical domain that includes
338 information from the transactional database and from a JSON file (already stored in
339 HDFS) with more information, such as the continent name.

340 The implementation process presented in Figure 5 starts with the data extraction
341 performed using Sqoop and Oozie Workflows and all data was stored in a HDFS directory
342 called Sandbox. This Sandbox directory allows the storage of all raw data and it is
343 divided into sub-directories where each data source has its own directory and is divided
344 into tables or entities. In this demonstration case, two data sources were used, the
345 transactional database and a JSON file.

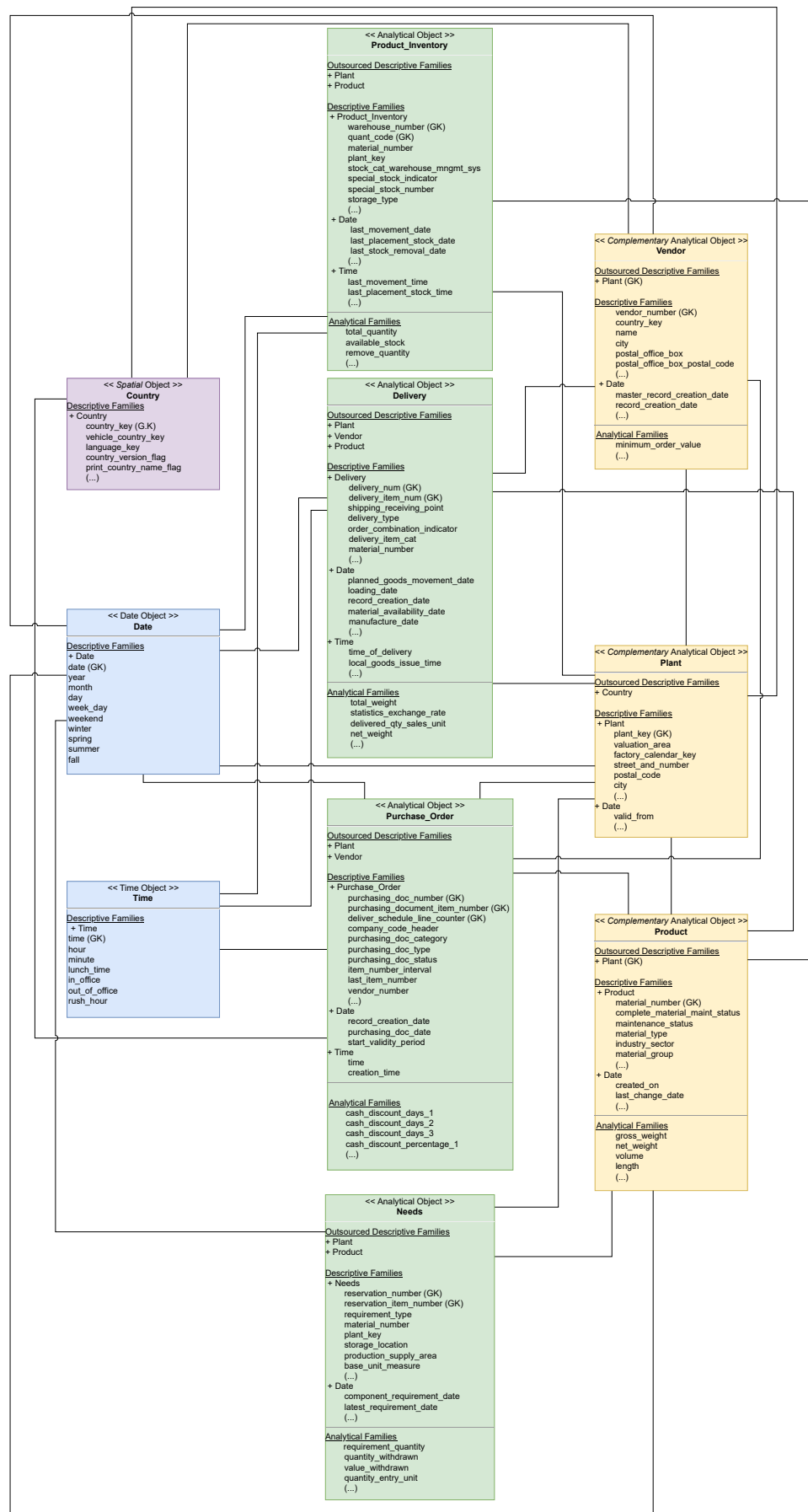


Figure 4. BDW Data Model

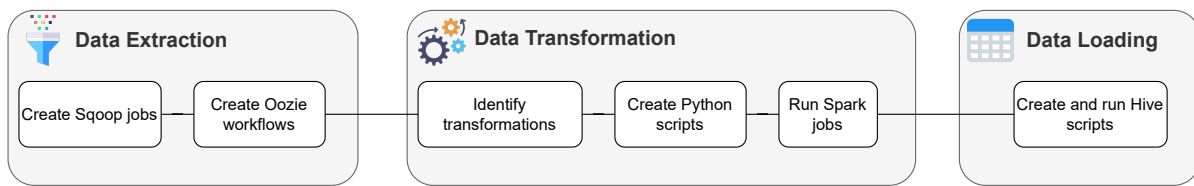


Figure 5. Implementation Process

346 With all the necessary data stored in HDFS, we can use Spark to perform the
 347 data transformation phase, where transformations and partitions keys are identified.
 348 Moreover, it is in this phase that the data enrichment can be performed with predictions
 349 from the machine learning models.

350 After the data transformation, the data is stored in the BDW where one table
 351 represents one of the objects included in the data model. Moreover, when the size of
 352 the object is too large to be used as one unique file, the object is partitioned according
 353 to their partition keys in order to improve the performance when querying the data.
 354 Furthermore, external Hive tables were created to provide Impala access to data. Impala
 355 will be the SQL query engine that allows the connection between Power BI and the data
 356 stored in HDFS.

357 5. Results and Discussion

358 In this section, we discuss the efficacy and efficiency (5.1) of the BDW implementa-
 359 tion. In subsections 5.2 and 5.3 it is presented the challenges and opportunities that arise
 360 while data-related projects are developed.

361 5.1. Efficacy and Efficiency

362 With the BDW implementation, it was possible to create a data repository that
 363 includes several businesses processes of the logistics department. Each process contains
 364 data from one or more tables from the transactional database used by the organization.

365 The data model is dynamic and able to change quickly, in order to include more
 366 tables, with more information related to any object that already exists in the BDW
 367 or to create new ones. The Time and Date objects can be used with other objects to
 368 understand the organization temporal dynamics, such as understand if there are any
 369 specific moments in the year where more delays are verified, or even when the suppliers
 370 are usually late with the deliveries. Similar reasoning can be used with the objects Plant
 371 and Inventory to analyse which plant has more inventory in its storage facilities.

372 With this work, it is now possible for the practitioners to use raw data extracted from
 373 the data sources (using the Sandbox layer) or use data already cleaned and transformed
 374 using the BDW layer. This can be achieved using the BDW Hive tables (as an example,
 375 Figure 6 shows the Country table view using the HUE interface) or the parquet files
 376 stored in the HDFS. They can also create specific materialized objects in the Application
 377 Layer in order to decrease the time needed to query the data. This reduces or even avoid
 378 the initial development time needed to understand, extract, store, and transform data.

379 The Machine Learning component can also use data from the different architecture
 380 components to provide useful predictions. For example, the available data can be used to
 381 predict if some scheduled delivery will be late or not. With this information, the logistics
 382 planners can take actions to reduce the impact of this situation. This can be achieved
 383 using data from the Sandbox or from the BDW. Machine Learning models can be created
 384 with this data using the Spark ML framework. Both the model and the predictions
 385 are stored in the HDFS being available for later use and for possible updates in the
 386 future. Furthermore, this data is now accessible to the organization through Impala
 387 connector and can be used to provide different insights about the organization status,
 388 or even in projects that use Machine Learning to predict or classify data to help in the
 389 decision making. This means that the time and the necessary knowledge to develop
 390 useful dashboards for management is smaller. In Figure 7, a dashboard that analyses

PROPERTIES	STATS
Table	Files 1 Rows 801 Total size
External and stored in location	30.75 KB
Created by aed1brg on Tue May 25 12:11:44 CEST 2021	Data last updated on 05/25/2021 11:11 AM +01:00

SCHEMA

Column (17)	Type	Description	Sample
i country_key	string		HU BD
i vehicle_country_key	string		H BD
i language_key	string		H E
i country_version	boolean		true false
i print_country_name	boolean		false true
i iso_code	string		HU BD
i iso_code_3_char	string		HUN BGD
i iso_code_name_3_c...	string		348 050
i eu_member	boolean		true false
i nationality	string		165 460
i altern_cntry_key	string		064 666
i trde_stat_short_name	string		UNGARN BANGLA
i date_form	string		1 Unknown
i country_currency	string		Unknown BDT
i continent_code	string		EU AS
i continent_name	string		Europe Asia

Figure 6. Country table in Hive

391 historical and predicted data is present, showing information about deliveries. It is an
 392 overview where the historical and predicted delayed or at time deliveries are analysed
 393 in several dimensions.

394 The top right of the dashboard shows the number of products that belongs to
 395 each category (A, B or C). This product classification demonstrates how important is
 396 each product for the organization. Products classified with A mean that this product is
 397 expensive for the organization and normally with more lead time, for example, electronic
 398 screens. The B category is for products less expensive, and the C category is for cheap
 399 products such as bolts. The impact on delays for products classified with A is superior
 400 to the products classified with B and C. The graph shows that are a bigger number of
 401 deliveries of C classification products demonstrating that this type of product has more
 402 frequent deliveries. So, if for some reason there is a shortage in stock of this product
 403 type, the organization will be able to solve that problem rapidly.

404 The two graphs in the lower-left corner of the dashboard compare the on-time
 405 deliveries and the delayed deliveries analysed by the season year. Each one compares
 406 the historical data and the predictions made by the machine learning algorithm. The
 407 left one shows that the predictions followed the trend of the historical date. The right
 408 one shows that is predicted an increase in delays in Autumn. With this information, the
 409 organization can prepare mitigation actions to decrease the impact of the delays.

410 The middle graphs compare the delayed deliveries and on-time deliveries by trans-
 411 portation mode. For example, we can see that the predictions (centre lower graph) show
 412 a general increase in the percentage of on-time deliveries.

413 The right side graphs compare the historical data with delays and the predictions.
 414 Bigger the circle means that are more deliveries from that country that arrive with delays.

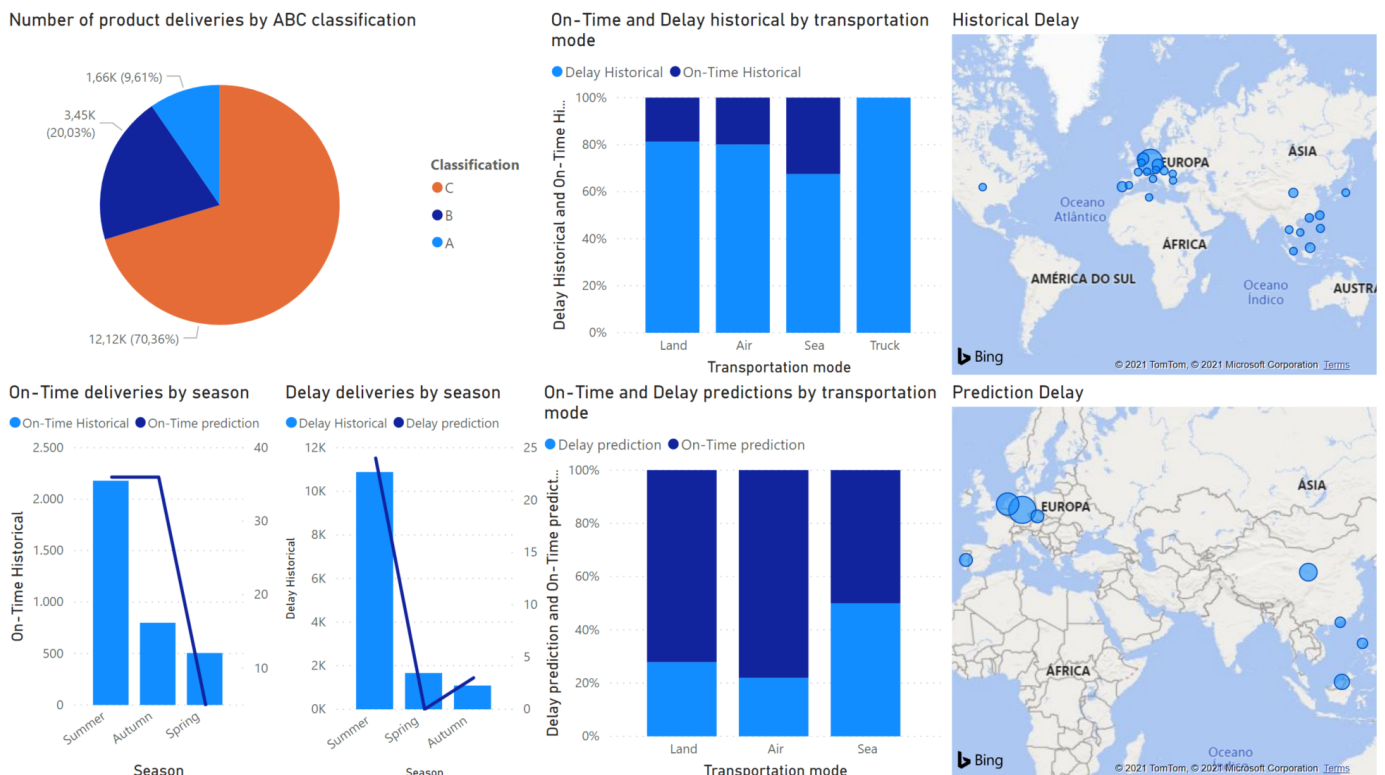


Figure 7. Dashboard with historical and predicted data related with deliveries

415 We can see that there are more delays from products shipped by European countries. The
 416 same is visible on the predictions.

417 These results are based on a portion of the historical data provided by the organi-
 418 zation. In future work it is necessary to verify if the predictions comply with reality and
 419 probably improve the model quality with more data.

420 5.2. Challenges

421 The implementation of new technology inside the organization's logistics depart-
 422 ment can be difficult and rises diverse types of challenges. These challenges can be
 423 related to the technology itself, with the lack of knowledge to develop the project, with
 424 the organizational culture, with the time and the cost to develop the project, among
 425 others. When that technology will use or rely on the provided transactional data to be
 426 successful, a new type of challenges related to data emerge.

427 Moreover, if the organization has a large dimension, can be extremely difficult to
 428 get the necessary knowledge to understand the different business processes inside the
 429 logistics department and the data generated by them. For example, if we are inside
 430 of a multinational organization, with diverse divisions, spread by multiple countries,
 431 with a complex transactional database, the data understanding will be one of the most
 432 challenging steps in the project.

433 The following list provides the identification and brief characterization of the most
 434 relevant challenges that can be faced while developing Big Data projects.

435 1. Data and technological challenges

436 • Data Understanding

437 Understanding the data that is stored in the transactional database is usually
 438 a challenge, even worse when the organization is a multinational with a
 439 considerable dimension. Transactional databases are complex systems, with
 440 misleading tables and attributes names. The existing documentation about
 441 the data source is usually sparse, not given enough insights about the data.

442 Several logistics concepts need to be known, such as safety stock, safety time,
443 delivery time, procurement, among others, in order to better understand the
444 data and their relationships.

- 445 • Poor or missing raw data
446 When an organization starts a project that will use the raw data generated
447 by the daily business, it is necessary to identify if the necessary data is being
448 generated and stored in the transactional system and its overall quality. Some-
449 times the project goals can not be achieved due to the lack of data or data with
450 quality. In complex ERP systems is possible to verify that many attributes
451 are not used by the organization. For example in logistics, knowing where
452 an order is in transit to its destination can be very useful to predict if it will
453 be on time, or not, and to make decisions about how to avoid stops in the
454 production line.
- 455 • Different values in different data sources for the same attribute
456 Due to the large and complex transactional system, is fairly common to find
457 the same attribute in different tables, related to the same entity, but with
458 different values. Understand why this happens and understand the type of
459 situations that motivate this type of behaviour can be difficult.
- 460 • Technological infrastructure
461 The adequate technological infrastructure is essential to stable a project de-
462 velopment. In an organization, the technological infrastructures can be based
463 on outdated technology or the technological infrastructure can change dur-
464 ing the project lifetime. This will lead to a project adaptation to the existing
465 technologies or their evolution as the infrastructure change.

466 2. Organizational challenges

- 467 • Access to data and to a technological infrastructure
468 One of the first tasks in projects of this nature is to get access to data and to
469 the infrastructure that will be used to process and store it. This is a task that
470 needs to be done at the beginning of the project and where the organizations'
471 policies can interfere in a negative way. This can not be an obstacle or take a
472 long time to overcome.
- 473 • Understand the business processes
474 Commonly, large organizations have many and complex business processes,
475 with diverse rules, exceptions and paths, which can be difficult to understand.
476 Moreover, the documentation about the business processes can be insuffi-
477 cient, creating another obstacle in this type of project. In the logistics area,
478 where daily interactions with the suppliers and their systems exist, where
479 processes are complex in order to achieve better results in the production
480 line, and where concepts such as just in time production are being imple-
481 mented, the documentations has a relevant impact when new projects start to
482 be developed.

483 3. Project team challenges

- 484 • Lack of knowledge in the used technologies
485 As Big Data is a recent concept, there is a lack of human resources with
486 experience in the technologies used to support this concept. Building a team
487 without any experience in Big Data can lead to several problems in the project.
488 Moreover, when adding specific requirements of a complex area like logistics,
489 more difficult is to get multidisciplinary teams with knowledge in both areas.
- 490 • Lack of sufficient human resources
491 To develop such a complex project, the project team needs an adequate number
492 of human resources. The lack of sufficient human resources can cause delays in
493 project development. Teams with a high number of elements can be prejudicial

494 to the project too, but very small teams lead to a lack of different backgrounds
495 and points of view that can hinder the project.

496 The challenges enumerated in this section are some of the biggest challenges that
497 a team can encounter while develop and implement a BDW inside of an organization
498 with a considerable size. The challenges can cause delays in the project milestones and
499 they should be taken into account when the project is planned. Most of them can be
500 mitigated with simple actions such as grant early access to all necessary resources and
501 develop the necessary documentation in all projects..

502 5.3. Opportunities

503 When an organization go through a technological change such as the creation of
504 a BDW, some opportunities emerge. Indeed, we can say that each challenge can be
505 transformed into one opportunity. Therefore, we will take the challenges provided in
506 section 5.2 and transform them into opportunities.

507 1. Data and technological opportunities

- 508 • Improve documentation
509 Very often, documentation is treated as the less important part of the project.
510 The time and effort put in the documentation development are lower than
511 required, leading to poor documentation. With the development of a new
512 project, the poor documentation of the previous one becomes evident. The
513 effort that needs to be done to understand the previous project can be reused
514 to improve the documentation and, therefore, decrease the time and effort
515 needed for the next ones.
- 516 • Improve data quality
517 Data quality is essential to the development of these data-based of projects.
518 As we need to perform data quality tasks, this can be used to detect and report
519 data problems that can be fixed in the near future. This can be useful not only
520 for this project but even for past and future projects.
- 521 • Technological infrastructure
522 A new project that requires new technology can be an excellent driver to
523 improve the technological infrastructure existent in the organization. These
524 changes can include, for example, updating the existent technologies or the
525 implementation of new ones.

526 2. Organizational opportunities

- 527 • Improve internal processes
528 With the implementation of new technology, some internal processes will
529 be analysed and can be improved. Moreover, processes can use the newly
530 available technology to improve their performance.
- 531 • Improve business processes documentation
532 Many analytical teams do not know the business processes and they need to
533 found the right person to ask. Often, if they ask the same question to different
534 persons, they will get different answers. Properly document the business
535 processes can be a key way to improve the business understanding not only
536 inside the analytical teams but for the organization in general.

537 3. Project team opportunities

- 538 • Creation of a team specialized in Big Data technologies
539 Research projects can have a tremendous impact on organizations, not only
540 by the obtained results but also by the improved capabilities of human re-
541 sources. In this specific case, the creation of one team specialized in Big Data
542 technologies can boost more projects, more efficiently, and with more efficacy.
- 543 • Improve workers knowledge in logistics processes
544 Human resources with other business knowledge can bring their knowledge
545 to other projects and have a positive impact on them. This can be verified

546 not only in new ones but also in the maintenance and improvement of other
547 ongoing projects.
548 • Improve workers knowledge about data sources
549 Data analytics projects always depend on the data source. Knowledge about
550 them is essential for a good start and a proper development of the project. It
551 is crucial to have in the project team, at least, one specialized resource in the
552 data sources, helping the development team to understand the data.

553 Besides the enumerated opportunities, other opportunities can arise with the cre-
554 ation and implementation of a BDW in a logistics department. For example, new projects
555 can be initiated and use the BDW as their data source, providing integrated and consoli-
556 dated data for their timely development. Other departments can use data in the BDW to
557 improve their predictions and their decision making needs.

558 6. Conclusions and Future Work

559 This paper presented the proposal and implementation of a BDW into a logistics de-
560 partment of an automotive factory. The implementation of the BDW is the starting point
561 to push the concept of Logistics 4.0 in this facility, improving the analytical capabilities
562 and supporting the decision-making process in the logistics department.

563 Through this work, we presented the logical and technological architecture that sup-
564 port the implementation of the BDW that includes several logistics processes. Moreover,
565 we presented the proposed BDW data model. The BDW data model is a key element
566 to get insights about the current state of the organization and to support the logistics
567 planners' decisions in an efficient way. The logical and technological architecture, as
568 well as the data model can be used as starting a point to develop and implement a BDW
569 in similar logistics departments.

570 As we advance, we faced several challenges and opportunities in the BDW devel-
571 opment and implementation. One of the most difficult challenges was to understand
572 the several logistics processes and how the data of these processes is stored in the trans-
573 actional system. Finding the right data to support the proposed system was a difficult
574 and time-consuming task. Nevertheless, the most important thing is to be aware of the
575 challenges and implement mitigation plans in order to solve them, or at least decrease
576 their impact on the project final results. Other challenges that can be faced in this area
577 are related to the technologies and the available infrastructure used by the organization.
578 Sometimes the technological infrastructure is changing during the project what can lead
579 to several project changes. Moreover, the available infrastructure can include outdated
580 technologies or be short in resources when used by several teams at the same time.

581 In the opportunities field, there are several points that can be addressed to improve
582 the organization, the logistics department, and the next projects. But these opportunities
583 need to be addressed in new projects with a well-defined goal and scope, due to the
584 new challenges that these projects will rise. Organizations need to promote a culture of
585 continuous improvement to face these opportunities.

586 As future work, the BDW implementation can be improved by automatizing the
587 data extraction, transforming, and enrichment pipelines to increase the performance and
588 decrease the human intervention. Moreover, the data model can be extended by adding
589 new objects (complementary or analytical) in order to enlarge their scope or improving
590 the existent ones by adding new data to the already existing objects. Furthermore, more
591 machine learning models can be created and integrated into the existing BDW to enrich
592 the data and provide predictions to help the logistics planners. Also, the implementation
593 of a real-time layer should be taken into consideration.

594 **Author Contributions:** Conceptualization, Nuno Silva and Julio Barros; Writing - original draft,
595 Nuno Silva and João N.C. Gonçalves; Writing - review & editing, Maribel Y. Santos and Car-
596 los Costa; Investigation, Nuno Silva and Julio Barros; Software, Nuno Silva and Julio Barros;
597 Supervision, Maribel Y. Santos, Carlos Costa, Paulo Cortez and M. Sameiro Carvalho.

598 **Funding:** This work has been supported by FCT – Fundação para a Ciência e Tecnologia within
599 the R&D Units Project Scope: UIDB/00319/2020, the doctoral scholarship grant:
600 PD/BDE/142895/2018 and PD/BDE/142900/2018.

601 **Acknowledgments:** This work has been designed using resources made by: Those Icons, Pixel
602 perfect, DinosoftLabs, Becris, Smashicons, and Freepik from www.flaticon.com.

603 **Conflicts of Interest:** The authors declare no conflict of interest.

References

1. Panetto, H.; Iung, B.; Ivanov, D.; Weichhart, G.; Wang, X. Challenges for the cyber-physical manufacturing enterprises of the future. *Annual Reviews in Control* **2019**, *47*, 200–213.
2. Winkelhaus, S.; Grosse, E.H. Logistics 4.0: a systematic review towards a new logistics system. *International Journal of Production Research* **2020**, *58*, 18–43.
3. Strandhagen, J.O.; Vallandingham, L.R.; Fragapane, G.; Strandhagen, J.W.; Stangeland, A.B.H.; Sharma, N. Logistics 4.0 and emerging sustainable business models. *Advances in Manufacturing* **2017**, *5*, 359–369.
4. Costa, C.; Santos, M.Y. Evaluating several design patterns and trends in big data warehousing systems. *International Conference on Advanced Information Systems Engineering*. Springer, 2018, pp. 459–473.
5. Chevalier, M.; Malki, M.E.; Kopliku, A.; Teste, O.; Tournier, R. Implementing Multidimensional Data Warehouses into NoSQL. 17th International Conference on Enterprise Information Systems (ICEIS 2015) held in conjunction with ENASE 2015 and GISTAM 2015. INSTICC - Institute for Systems and Technologies of Information, Control and Communication, 2015, pp. 172–183.
6. Gröger, C.; Schwarz, H.; Mitschang, B. The Deep Data Warehouse: Link-Based Integration and Enrichment of Warehouse Data and Unstructured Content. 2014 IEEE 18th International Enterprise Distributed Object Computing Conference, 2014, pp. 210–217. doi:10.1109/EDOC.2014.36.
7. Kiran, M.; Murphy, P.; Monga, I.; Dugan, J.; Baveja, S.S. Lambda architecture for cost-effective batch and speed big data processing. *Proceedings - 2015 IEEE International Conference on Big Data, IEEE Big Data 2015* **2015**, pp. 2785–2792. doi:10.1109/BigData.2015.7364082.
8. NBD-PWG. NIST Big Data Interoperability Framework: Volume 6, Reference Architecture. Technical Report NIST SP 1500-6, National Institute of Standards and Technology, 2015.
9. Santos, M.Y.; Costa, C., Big Data: Concepts, Warehousing, and Analytics. In *Big Data: Concepts, Warehousing, and Analytics*; River Publishers, 2020; pp. 1–284.
10. Chou, S.; Yang, C.; Jiang, F.; Chang, C. The Implementation of a Data-Accessing Platform Built from Big Data Warehouse of Electric Loads. 2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC), 2018, Vol. 02, pp. 87–92. doi:10.1109/COMPSAC.2018.10208.
11. Sebaa, A.; Chikh, F.; Nouicer, A.; Tari, A. Medical Big Data Warehouse: Architecture and System Design, a Case Study: Improving Healthcare Resources Distribution. *Journal of medical systems* **2018**, *42*, 59. doi:10.1007/s10916-018-0894-9.
12. Santos, M.Y.; Martinho, B.; Costa, C. Modelling and implementing big data warehouses for decision support. *Journal of Management Analytics* **2017**, *4*, 111–129. doi:10.1080/23270012.2017.1304292.
13. Shvachko, K.; Kuang, H.; Radia, S.; Chansler, R. The Hadoop Distributed File System. 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST), 2010, pp. 1–10. doi:10.1109/MSST.2010.5496972.
14. Dean, J.; Ghemawat, S. MapReduce: Simplified Data Processing on Large Clusters. *Commun. ACM* **2008**, *51*, 107–113. doi:10.1145/1327452.1327492.
15. Thusoo, A.; Sarma, J.S.; Jain, N.; Shao, Z.; Chakka, P.; Anthony, S.; Liu, H.; Wyckoff, P.; Murthy, R. Hive: a warehousing solution over a map-reduce framework. *Proceedings of the VLDB Endowment* **2009**, *2*, 1626–1629.
16. Spark, A. Apache spark. Retrieved January **2018**, *17*, 2018.
17. Bittorf, M.; Bobrovitsky, T.; Erickson, C.; Hecht, M.G.D.; Kuff, M.; Leblang, D.K.A.; Robinson, N.; Rus, D.R.S.; Wanderman, J.; Yoder, M.M. Impala: A modern, open-source sql engine for hadoop. *Proceedings of the 7th biennial conference on innovative data systems research*, 2015.
18. L’Heureux, A.; Grolinger, K.; Elyamany, H.F.; Capretz, M.A.M. Machine Learning With Big Data: Challenges and Approaches. *IEEE Access* **2017**, *5*, 7776–7797. doi:10.1109/ACCESS.2017.2696365.
19. Costa, C.; Andrade, C.; Santos, M.Y., Big Data Warehouses for Smart Industries. In *Encyclopedia of Big Data Technologies*; Springer International Publishing: Cham, 2018; pp. 1–11. doi:10.1007/978-3-319-63962-8_204-1.
20. Aravinth, S.; Begam, A.H.; Shanmugapriyaa, S.; Sowmya, S.; Arun, E. An efficient HADOOP frameworks SQOOP and ambari for big data processing. *International Journal for Innovative Research in Science and Technology* **2015**, *1*, 252–255.
21. Ivanov, T.; Pergolesi, M. The impact of columnar file formats on SQL-on-hadoop engine performance: A study on ORC and Parquet. *Concurrency and Computation: Practice and Experience* **2020**, *32*, e5523.
22. Baranowski, Z.; Grzybek, M.; Canali, L.; Garcia, D.L.; Surdy, K. Scale out databases for CERN use cases. *J. Phys. Conf. Ser.*, 2015, Vol. 664, pp. 042–002.

23. Armbrust, M.; Xin, R.S.; Lian, C.; Huai, Y.; Liu, D.; Bradley, J.K.; Meng, X.; Kaftan, T.; Franklin, M.J.; Ghodsi, A.; Zaharia, M. Spark SQL: Relational Data Processing in Spark. Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data; Association for Computing Machinery: New York, NY, USA, 2015; SIGMOD '15, p. 1383–1394. doi:10.1145/2723372.2742797.
24. Meng, X.; Bradley, J.; Yavuz, B.; Sparks, E.; Venkataraman, S.; Liu, D.; Freeman, J.; Tsai, D.; Amde, M.; Owen, S.; Xin, D.; Xin, R.; Franklin, M.J.; Zadeh, R.; Zaharia, M.; Talwalkar, A. MLlib: Machine Learning in Apache Spark. *J. Mach. Learn. Res.* **2016**, *17*, 1235–1241.
25. Qin, X.; Chen, Y.; Chen, J.; Li, S.; Liu, J.; Zhang, H. The Performance of SQL-on-Hadoop Systems - An Experimental Study. 2017 IEEE International Congress on Big Data (BigData Congress), 2017, pp. 464–471. doi:10.1109/BigDataCongress.2017.68.