

Course Introduction

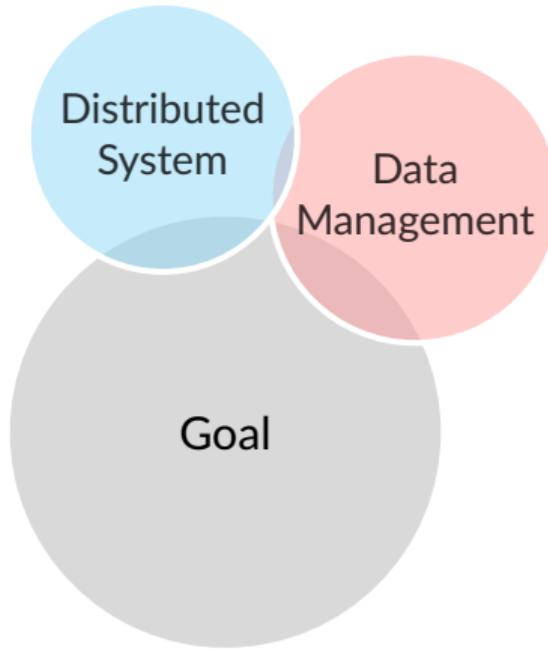
Course Target



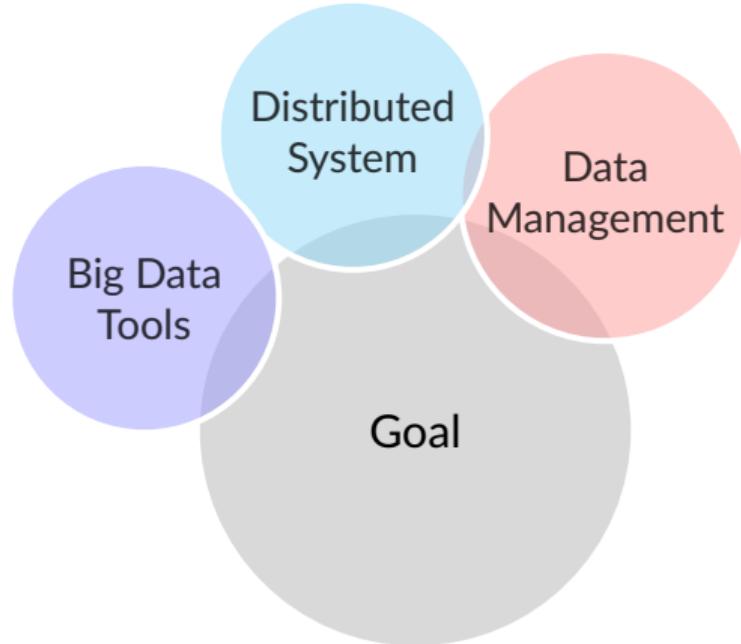
Course Target



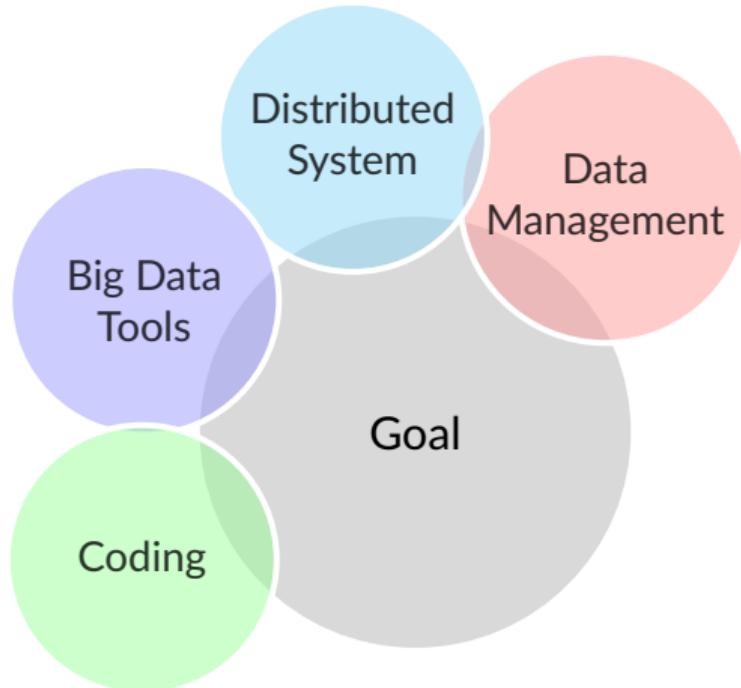
Course Target



Course Target



Course Target



Course Target



Course Target



Course Target



Learning Objectives

- Simplify the concepts in data management.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.
- Build and scale your data product.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.
- Build and scale your data product.
- Building real-life examples.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.
- Build and scale your data product.
- Building real-life examples.
- Applying machine learning over big data.



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.
- Build and scale your data product.
- Building real-life examples.
- Applying machine learning over big data.
- Automate the data life-cycle process end-to-end (e2e).



Learning Objectives

- Simplify the concepts in data management.
- Understand the data management life-cycle.
- Illustrate the basics of distributed systems concepts
- Be familiar with ETL for (batch/streaming) data over distributed systems ex: Hadoop & Spark.
- Apply QA and testing for the data pipeline cycle.
- Build and scale your data product.
- Building real-life examples.
- Applying machine learning over big data.
- Automate the data life-cycle process end-to-end (e2e).
- Understanding of the DevOps tools and functions in data life-cycle.



Audience: Who Should Take This Course?

- Data Engineer who needs to get more knowledge in distributed systems and Big Data.



Audience: Who Should Take This Course?

- Data Engineer who needs to get more knowledge in distributed systems and Big Data.
- Data Warehouse Engineer who needs to know more about big data.



Audience: Who Should Take This Course?

- Data Engineer who needs to get more knowledge in distributed systems and Big Data.
- Data Warehouse Engineer who needs to know more about big data.
- Software developer who needs to change to data engineering track.



Audience: Who Should Take This Course?

- Data Engineer who needs to get more knowledge in distributed systems and Big Data.
- Data Warehouse Engineer who needs to know more about big data.
- Software developer who needs to change to data engineering track.
- DevOps engineers who needs to understand the concepts of big data.



Audience: Who Should Take This Course?

- Data Engineer who needs to get more knowledge in distributed systems and Big Data.
- Data Warehouse Engineer who needs to know more about big data.
- Software developer who needs to change to data engineering track.
- DevOps engineers who needs to understand the concepts of big data.
- Business or entrepreneur who needs to get more information about how to build or manage a data product.



Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.

Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).

Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).
- Repeat the lecture code with your own.

Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).
- Repeat the lecture code with your own.
- Do the assignments.

Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).
- Repeat the lecture code with your own.
- Do the assignments.
- Ask your questions.

Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).
- Repeat the lecture code with your own.
- Do the assignments.
- Ask your questions.
- Join the online meeting or discussions.

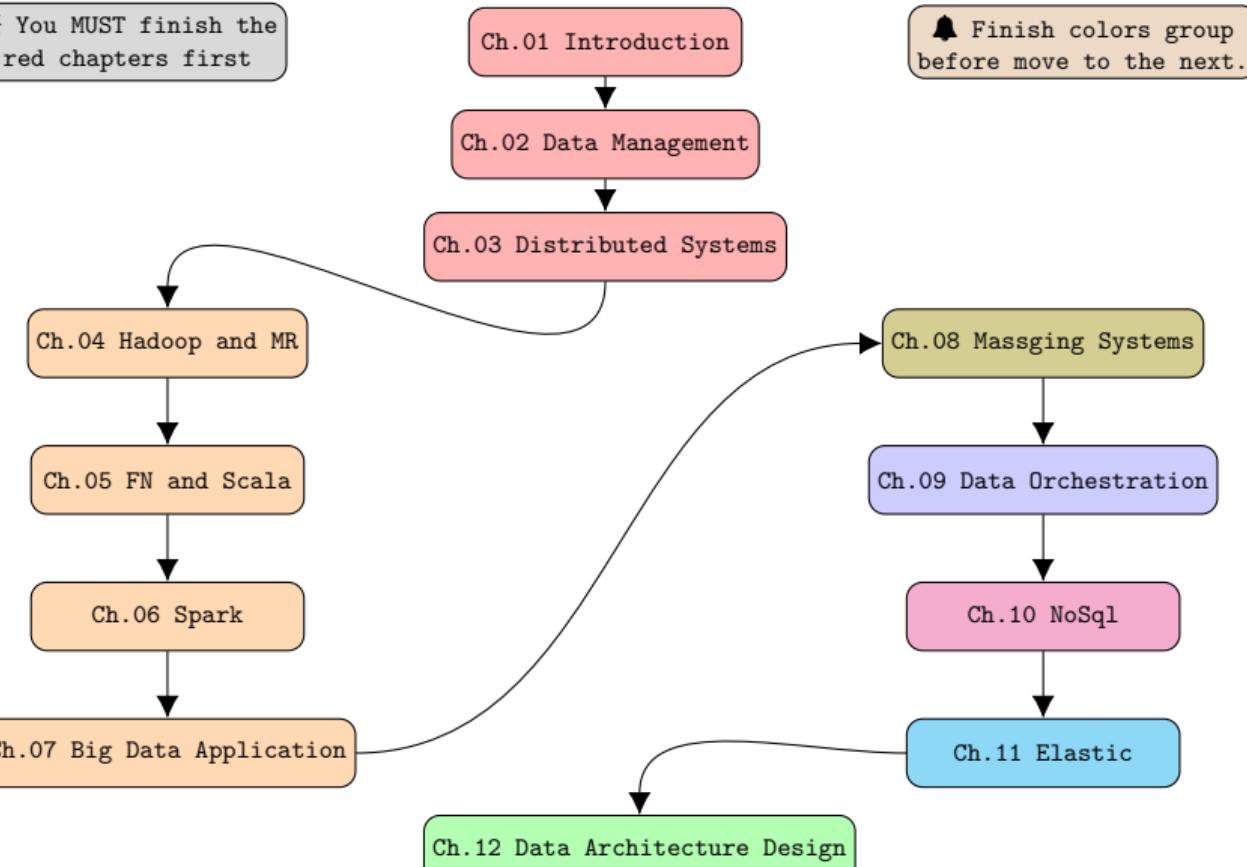
Getting max benefit from this course

Take the course advantage

- Follow the videos order as described.
- Read the references for each section (including the implementation of the examples if exists).
- Repeat the lecture code with your own.
- Do the assignments.
- Ask your questions.
- Join the online meeting or discussions.

Chapter Dependencies

❗ You MUST finish the red chapters first



Chapter Dependencies (Jump Out Path)

⚠ You MUST finish the red chapters first

Ch.01 Introduction

🔔 Finish colors group
before move to the next.

Ch.02 Data Management

Ch.03 Distributed Systems

Ch.04 Hadoop and MR

Ch.05 FN and Scala

Ch.06 Spark

Ch.07 Big Data Application

Ch.08 Massging Systems

Ch.09 Data Orchestration

Ch.10 NoSql

Ch.11 Elastic

Ch.12 Data Architecture Design

Assignments and Labs

Remark

- Full project code.

Assignments and Labs

Remark

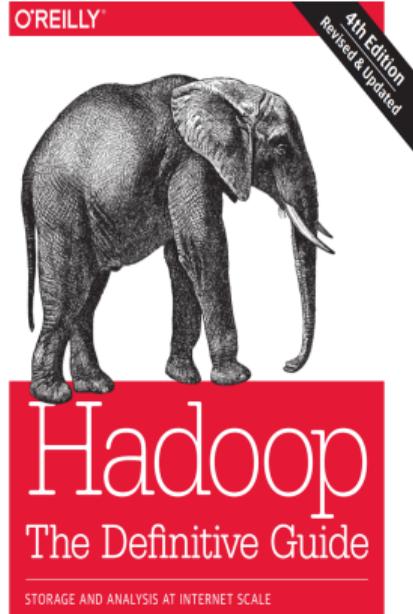
- Full project code.
- Notebooks (Jupyter or Zeppelin).

Assignments and Labs

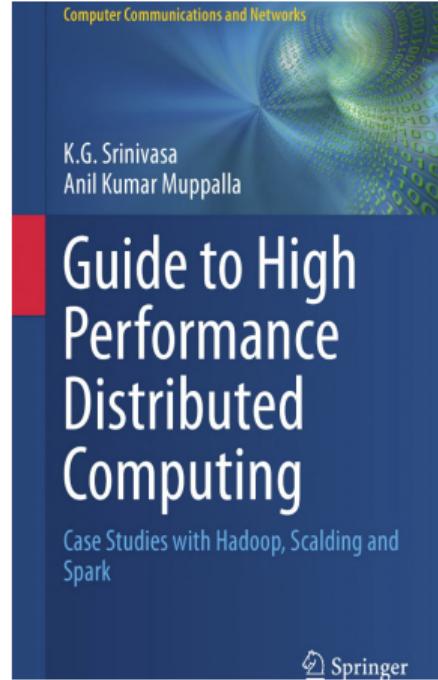
Remark

- Full project code.
- Notebooks (Jupyter or Zeppelin).
- Read the reference.

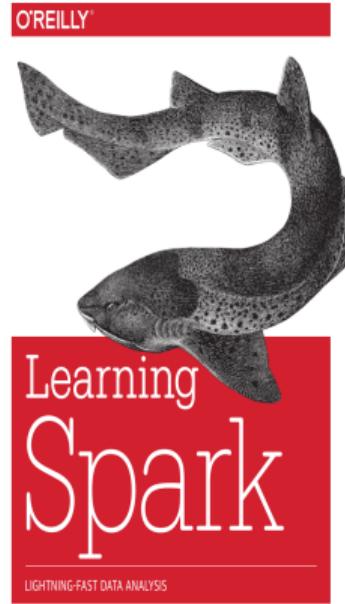
Textbooks-1



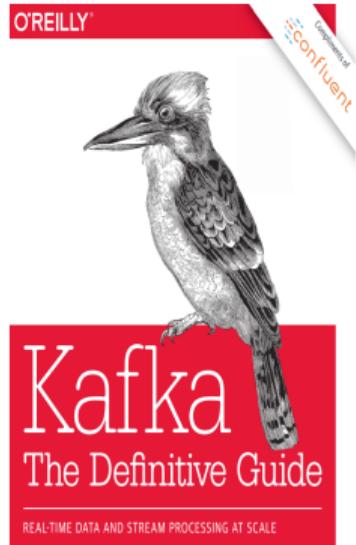
Tom White



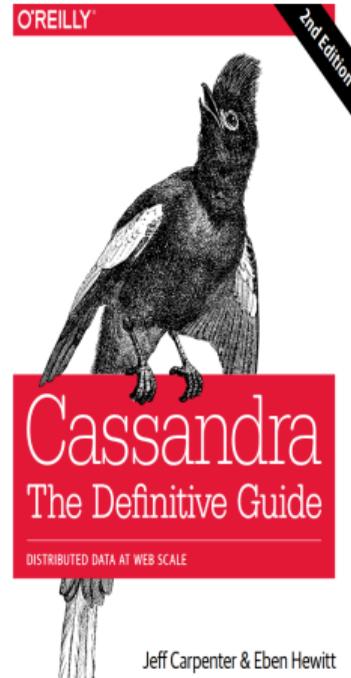
Textbooks-2



Textbooks-3



Neha Narkhede,
Gwen Shapira & Todd Palino



Jeff Carpenter & Eben Hewitt



Martin Kleppmann



Ugly but important

- User stories or technical discussions are not related to any of my current work or my previous companies.
- I am working at EPAM Systems. My company approved me for doing this online course public but the materials are not reviewed or assessed by my company. It is on my own responsibilities.

Table of Contents I

1 Course Introduction

- Learning Objectives
- Getting max benefit from this course
- Assignments and Labs

2 Introduction To Data Management and Data Warehouse

- Data Management
- From DWH to Big Data
- Data Models
 - Data Model Design
 - Data Encoding and Formats
- Further Readings and Assignment

3 Introduction To Distributed Systems

- Distributed Systems Concepts
- Distributed Systems Architecture
- Distributed Systems Challenges



Table of Contents II

- Design Simple Distributed System
- Further Readings and Assignment

4

Introduction to Hadoop and Map-Reduce

- Hadoop Architecture
 - Storage
 - YARN
 - Hadoop I/O
 - Processing
- Map-Reduce
 - Map-Reduce Components
 - Word-Count Example
- Pig
- Hive
- ZooKeeper
- Further Readings and Assignment



Table of Contents III

5

Functional Programming

- Why functional programming commonly used in distributed systems?
- Introduction to Scala
- Further Readings and Assignment

6

Spark Framework

- Spark Philosophy towards the Engine and the Programming languages
- Spark Basics
- Spark Programming using RDDs
 - Spark RDD
 - Spark Working With Key/Value Pairs
- Spark Datasets/Dataframe
 - Spark SQL
 - Dataframes/Datasets vs. RDDs

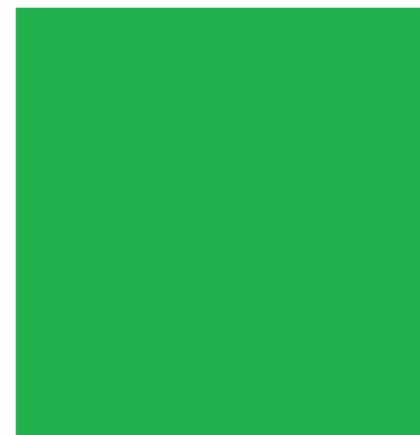


Table of Contents IV

- Spark on Production
- Spark For Batch Processing
- Building custom input and output connector using Spark
- Spark Streaming
- Spark using other Programming Languages
 - PySpark for Python Geeks
 - RSpark for R Geeks
- Spark For Data Scientist
- Spark Graph Dataframe/Graphx
- Tuning your Spark Jobs
- Further Readings and Assignment

7

Real World Applications

- Big Data Development Life Cycle
- Template Concept for Data Engineering
 - Template for ETL Application

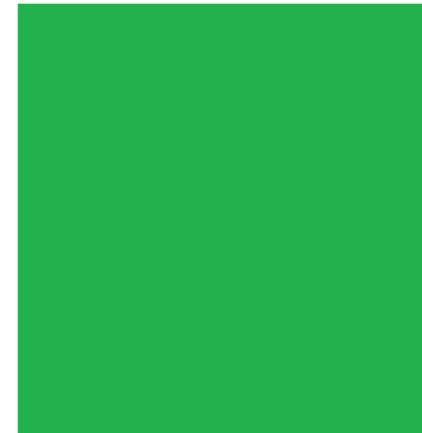


Table of Contents V

- Template for QA
- Template for Streaming Applications
- Template for Machine Learning Applications
- Further Readings and Assignment

8

Massaging Systems

- Motivation
- Massaging Systems Architecture
- JMS as an example
- Introduction to Kafka
 - Kafka Architecture
 - Kafka Topics
 - Partitions
 - Kafka Producers
 - Kafka Consumers
 - Kafka Connector
 - Kafka Custom Connectors



Table of Contents VI

- Kafka Configuration
- Kafka Configuration Optimizations
- Kafka Operations
- Kafka Integration with Enterprise tools
- Further Readings and Assignment

9

Data Orchestration

- Motivation
- Enterprise vs Open source tools
 - Open source tools (Oozie as an Example)
 - Enterprise source tools
 - How to choose the right tool?
- Further Readings and Assignment

10

NOSQL

- Introduction to NoSQL Databases.
- Cassandra
 - Why Cassandra?

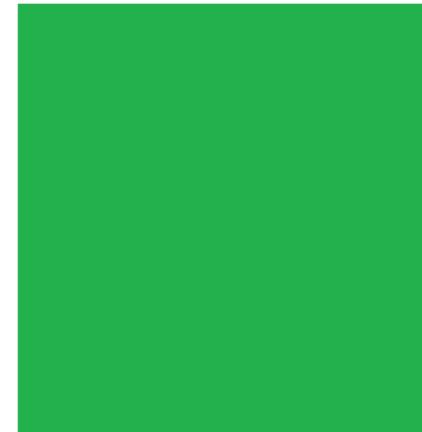


Table of Contents VII

- Introducing Cassandra
- The Cassandra Data Model
- Architecture
- Reading and Writing Data
- Integrating Hadoop
- Further Readings and Assignment

11 Elastic

- Further Readings and Assignment

12 Data Architecture Design

- Further Readings and Assignment

13 Appendix

- Appendix A- Shell Programming
- Appendix B- Java Programming
- Appendix C- Scala Programming
- Appendix D- SQL Programming



Table of Contents VIII

- Appendix E- Oozie Orchestration
- Appendix F- DWH Concepts and Data Modeling Design
- Appendix G- Machine Learning Concepts Data Engineers
- Appendix H- Docker for Data Engineers



Introduction To Data Management and Data Warehouse

Chapter Objectives

- Be familiar with data management life-cycle.



Chapter Objectives

- Be familiar with data management life-cycle.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats
- Challenges to build a DWH.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats
- Challenges to build a DWH.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats
- Challenges to build a DWH.
- Data modeling design.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats
- Challenges to build a DWH.
- Data modeling design.



Chapter Objectives

- Be familiar with data management life-cycle.
- Introduction to data warehouse and its usage.
- Motivation to DWH.
- What is the different types of DWH?
- Use cases
- Data Encoding and Formats
- Challenges to build a DWH.
- Data modeling design.



Data Management

- Data are a product.



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).
 - **Transformation** ex: (cleansing, Apply business logic, Organize).



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).
 - **Transformation** ex: (cleansing, Apply business logic, Organize).
 - **Loading** or store the transformed data based on our usage or use case.



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).
 - **Transformation** ex: (cleansing, Apply business logic, Organize).
 - **Loading** or store the transformed data based on our usage or use case.
 - Business Intelligence (**BI**) or data discovery (continues process).



Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).
 - **Transformation** ex: (cleansing, Apply business logic, Organize).
 - **Loading** or store the transformed data based on our usage or use case.
 - Business Intelligence (**BI**) or data discovery (continues process).
 - **Integration** and publishing.

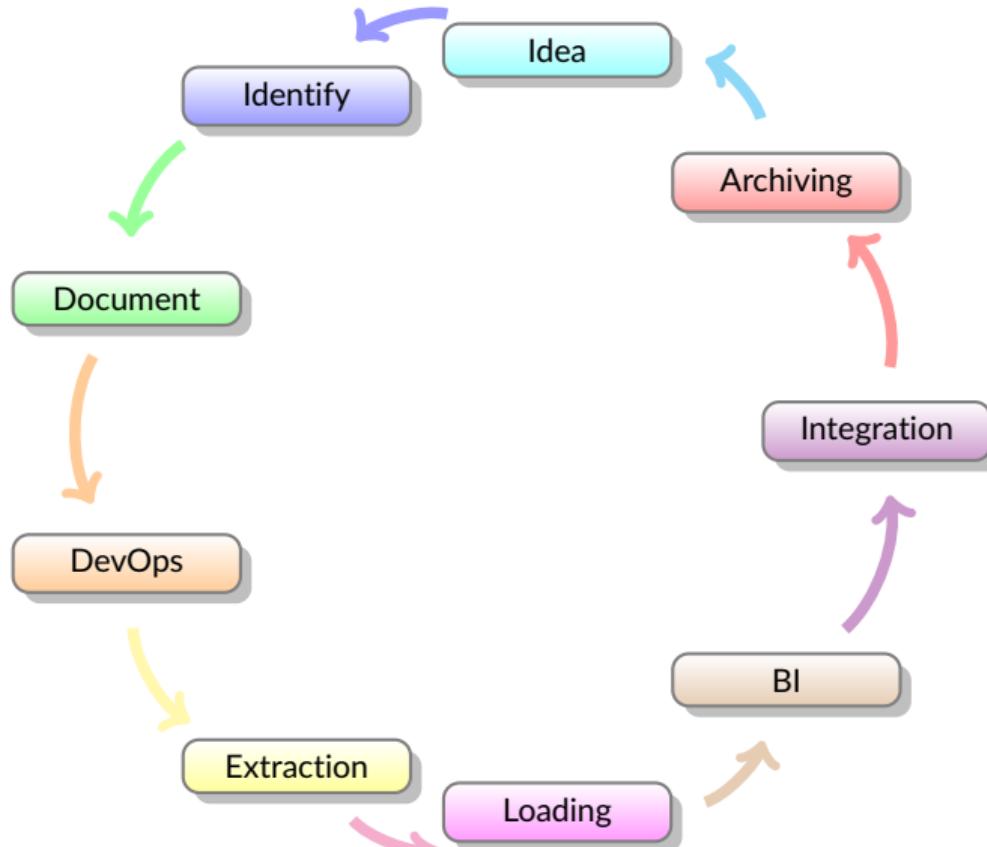


Data Management

- Data are a product.
- Data product has a life-cycle as following (simplified):
 - **Question**, Idea, or service.
 - **Identifying** the source of information and the data type ex: (text, images, videos, audio, or sensors).
 - **Document** all details regarding the data including quality, security, efficiency, and access (consideration during the cycle).
 - Delivery automation (Tools and Process) AKA **DevOps** cycle.
 - **Extraction** Process (collection).
 - **Transformation** ex: (cleansing, Apply business logic, Organize).
 - **Loading** or store the transformed data based on our usage or use case.
 - Business Intelligence (**BI**) or data discovery (continues process).
 - **Integration** and publishing.
 - Data retention or **archiving** process ex: (Hot or Cold storage).



Data Management Life-Cycle



Motivation to Data Warehouse

- Data could be a product for some companies.



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).
- There are some challenges facing the people who work on data management backend:



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).
- There are some challenges facing the people who work on data management backend:
 - Performance.



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).
- There are some challenges facing the people who work on data management backend:
 - Performance.
 - Integration.



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).
- There are some challenges facing the people who work on data management backend:
 - Performance.
 - Integration.
 - Applying analytical functions.



Motivation to Data Warehouse

- Data could be a product for some companies.
- It could be decision support for other products or businesses.
- Reporting the results after pass the data life-cycle will be from storage (Database).
- There are some challenges facing the people who work on data management backend:
 - Performance.
 - Integration.
 - Applying analytical functions.
- Vendors who are working to solve the above challenges creating their own product of DWH and their ultimate work is to optimize the above points.



Motivation to Data Warehouse (DWH)

Definition (What is Data Warehousing?)

A DWH is defined as a technique for collecting and managing data from varied sources to **provide meaningful business insights**. It is a blend of technologies and components which aids the strategic use of data.

- The DWH is not a product but an environment.

Motivation to Data Warehouse (DWH)

Definition (What is Data Warehousing?)

A DWH is defined as a technique for collecting and managing data from varied sources to **provide meaningful business insights**. It is a blend of technologies and components which aids the strategic use of data.

- The DWH is not a product but an environment.
- It is a process of transforming data into information and make it available to users in a **timely manner** to make a difference.

Motivation to Data Warehouse (DWH)

Definition (What is Data Warehousing?)

A DWH is defined as a technique for collecting and managing data from varied sources to **provide meaningful business insights**. It is a blend of technologies and components which aids the strategic use of data.

- The DWH is not a product but an environment.
- It is a process of transforming data into information and make it available to users in a **timely manner** to make a difference.
- It is an architectural construct of an information system which provides users with current and historical decision support information which is difficult to access or present in the traditional operational data store.

Motivation to Data Warehouse (DWH)

Definition (What is Data Warehousing?)

A DWH is defined as a technique for collecting and managing data from varied sources to **provide meaningful business insights**. It is a blend of technologies and components which aids the strategic use of data.

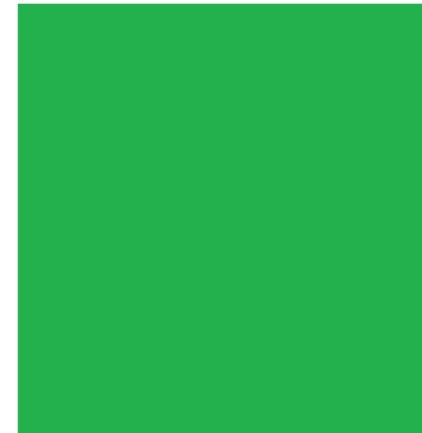
- The DWH is not a product but an environment.
- It is a process of transforming data into information and make it available to users in a **timely manner** to make a difference.
- It is an architectural construct of an information system which provides users with current and historical decision support information which is difficult to access or present in the traditional operational data store.
- The DWH is the core of the BI system which is built for data analysis and reporting.

Motivation to Data Warehouse

Data warehouse system is also known by the following names:

- Decision Support System (DSS).
- Business Intelligence Solution.
- Executive Information System.
- Management Information System.
- Analytic Application.
- Data Warehouse.

The real concept was given by Inmon Bill. He was considered as a father of the DWH. He had written about a variety of topics for building, usage, and maintenance of the warehouse & the Corporate Information Factory



Motivation to Data Warehouse

Types of Data Warehouse

Enterprise Data Warehouse (EDWH) It provides decision support service across the enterprise. It offers a unified approach for organizing and representing data (DWH Model). It offers data classifications according to the subject with privileges policy.

Operational Data Store (ODS): is a central database that provides an up-to-date (real-time) data from multiple transnational systems for operational reporting into a single DWH.

Data Mart: A data mart is a subset of the data warehouse. It specially designed for a particular line of business, such as sales, finance, sales or finance. In an independent data mart, data can collect directly from sources.

DWH vs ODS vs Data Mart

Metric	DWH	ODS	Data Mart
Latency	Day -1	Real-time	Day -1
Data level	Transnational	Transnational	Summary
Historical	Long-term	Snapshot	Aggregated Long-Term
Size	TB/PB	GB	GB/TB
Orientation	Multi sources	Multi sources	Product
Business Units	Multi organizational units	Product team	Business team

DWH vs Operational databases

Metric	Transactions DB	DWH
Volume	GB/TB	TB/PB
Historical rows	Short-term ;1000M	Long-Term 1000M;
Orientation	Product	Subject or multi products
Business Units	Product team	Multi organizational units
Normalization	Normalized	Not required (De-normalized in many use cases)
Data Model	Relational	Star Schema or Multi-dim
Intelligence	Reporting	Advanced reporting and Machine Learning
Use cases	Online transactions & operations	Centralized storage (360°)

Transnational DB Use cases



Transnational DB Use cases



DWH Use cases



DWH Use cases



DWH Use cases



Use case (Operational DB)

- A telecommunication company named **XTec**.



Use case (Operational DB)

- A telecommunication company named **XTec**.
- They have lots of systems. One of this systems is a CRM system as example of operational DB.



Use case (Operational DB)

- A telecommunication company named **XTec**.
- They have lots of systems. One of this systems is a CRM system as example of operational DB.
 - The CRM system handles the customer activities with the company including (sales, change in customer plans, and other activities).



Use case (Operational DB)

- A telecommunication company named **XTec**.
- They have lots of systems. One of this systems is a CRM system as example of operational DB.
 - The CRM system handles the customer activities with the company including (sales, change in customer plans, and other activities).
 - This system has a backend database (MySQL).



Use case (Operational DB)

- A telecommunication company named **XTec**.
- They have lots of systems. One of this systems is a CRM system as example of operational DB.
 - The CRM system handles the customer activities with the company including (sales, change in customer plans, and other activities).
 - This system has a backend database (MySQL).
 - CRM team can report their sales and customer activities from their database.



Use case (Operational DB)

- A telecommunication company named **XTec**.
- They have lots of systems. One of this systems is a CRM system as example of operational DB.
 - The CRM system handles the customer activities with the company including (sales, change in customer plans, and other activities).
 - This system has a backend database (MySQL).
 - CRM team can report their sales and customer activities from their database.
 - Product owner can take a decision based on their system backend reports.



Use case (DWH)

- What is the need for DWH?



Use case (DWH)

- What is the need for DWH?
 - This company has other systems for example: billing, charging, signaling.



Use case (DWH)

- What is the need for DWH?
 - This company has other systems for example: billing, charging, signaling.
 - They need to report information related to the CRM, billing, and signaling source systems in one report.



Use case (DWH)

- What is the need for DWH?

- This company has other systems for example: billing, charging, signaling.
- They need to report information related to the CRM, billing, and signaling source systems in one report.
- So, they need to ingest (transfer) the data from the source systems to one single database.



Use case (DWH)

- What is the need for DWH?

- This company has other systems for example: billing, charging, signaling.
- They need to report information related to the CRM, billing, and signaling source systems in one report.
- So, they need to ingest (transfer) the data from the source systems to one single database.
- The decision from the DHW is a **global and strategical decision**.



Use case (DWH)

- What is the need for DWH?

- This company has other systems for example: billing, charging, signaling.
- They need to report information related to the CRM, billing, and signaling source systems in one report.
- So, they need to ingest (transfer) the data from the source systems to one single database.
- The decision from the DWH is a **global and strategical decision**.
- If the company needs to build a machine learning model which needs data from different sources. They need to load the data from a centralized database rather than read each source alone.



Use case (DWH)

The Full picture required a DWH. However, we still need the other operational databases for product development perspective.



Use case (ODS)

- Why do we need the ODS?



Use case (ODS)

- Why do we need the ODS?
- How does it fit in our system?



Use case (ODS)

XTec has a call center system which handles the customer inquiries.

This system requires the some data related to usage, customer information, billing details to be calculated and accumulated in **real-time** to be able to give the customer the right answer for his inquiries.



Use case (ODS)

- So, What is the challenge for this system?



Use case (ODS)

- So, What is the challenge for this system?
 - It needs specific information from different source systems.



Use case (ODS)

- So, What is the challenge for this system?
 - It needs specific information from different source systems.
 - It requires to track the source system database changes or update in real-time.



Use case (ODS)

- So, What is the challenge for this system?
 - It needs specific information from different source systems.
 - It requires to track the source system database changes or update in real-time.
 - Its functionality is based on the aggregate data not the transactions for example (It needs the total outgoing calls till time or it needs the total charging amounts from prepaid or the available limits from billing if it is postpaid).



Use case (ODS)

- ODS is based on change data capture (CDC). This approach used to determine the data change and apply action based on this change.



Use case (ODS)

- ODS is based on change data capture (CDC). This approach used to determine the data change and apply action based on this change.
- ODS uses the real-time aggregations to support the online systems from different source systems.



DWH Architecture Overview

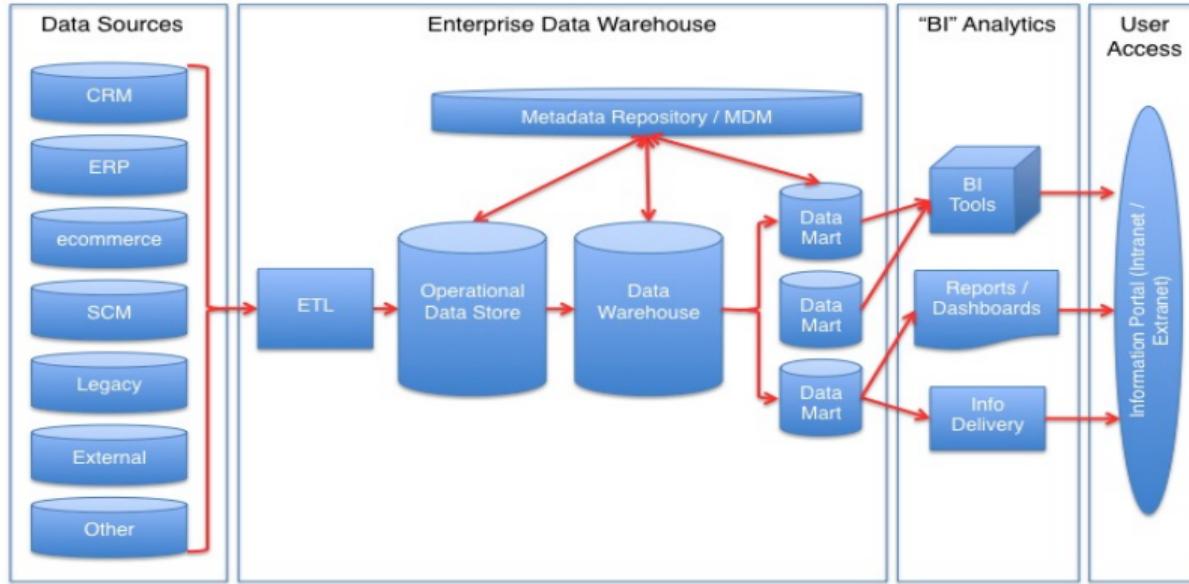


Figure: taken from XXXX

Thinking about data as a product

- How can we think about a data solution?



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.
 - Ask your self the following questions:



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.
 - Ask your self the following questions:
 - Can we solve the problem using the current data structure with adding new features?



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.
 - Ask your self the following questions:
 - Can we solve the problem using the current data structure with adding new features?
 - What if we enhance/change the data structure or modeling?

Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.
 - Ask your self the following questions:
 - Can we solve the problem using the current data structure with adding new features?
 - What if we enhance/change the data structure or modeling?
 - Could it help if we change the backend engine (ex: database)?



Thinking about data as a product

- How can we think about a data solution?
 - Requirements analysis.
 - Identify the problem (challenges)
 - Think about how to overcome the challenges.
 - Ask your self the following questions:
 - Can we solve the problem using the current data structure with adding new features?
 - What if we enhance/change the data structure or modeling?
 - Could it help if we change the backend engine (ex: database)?
- To answer these questions you need to understand the data layers.



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.
- Each layer responsible for different tasks and operations.



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.
- Each layer responsible for different tasks and operations.
- Each layers interacts with (hardware or software or mixed).



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.
- Each layer responsible for different tasks and operations.
- Each layers interacts with (hardware or software or mixed).
- To eliminate the complexity for data interactions not all internal processes are shared or available for the user.



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.
- Each layer responsible for different tasks and operations.
- Each layers interacts with (hardware or software or mixed).
- To eliminate the complexity for data interactions not all internal processes are shared or available for the user.
- The developer for each layer hide internal irrelevant details from developer (users).



Data Layers (Abstraction)

- Any data product (database) contains multi-layers.
- Each layer responsible for different tasks and operations.
- Each layers interacts with (hardware or software or mixed).
- To eliminate the complexity for data interactions not all internal processes are shared or available for the user.
- The developer for each layer hide internal irrelevant details from developer (users).
- The process of **hiding** irrelevant details from developer (user) is called data **abstraction**.



Data Layers (Abstraction)

Definition

Data Abstraction and Data Independence: Database systems comprise of complex data-structures. In order to make the system efficient in terms of retrieval of data, and reduce complexity in terms of usability of users, developers use abstraction i.e. hide irrelevant details from the users. This approach simplifies database design.

- There are 3 levels of data abstraction.

Data Layers (Abstraction)

Definition

Data Abstraction and Data Independence: Database systems comprise of complex data-structures. In order to make the system efficient in terms of retrieval of data, and reduce complexity in terms of usability of users, developers use abstraction i.e. hide irrelevant details from the users. This approach simplifies database design.

- There are 3 levels of data abstraction.
 - Physical level

Data Layers (Abstraction)

Definition

Data Abstraction and Data Independence: Database systems comprise of complex data-structures. In order to make the system efficient in terms of retrieval of data, and reduce complexity in terms of usability of users, developers use abstraction i.e. hide irrelevant details from the users. This approach simplifies database design.

- There are 3 levels of data abstraction.
 - Physical level
 - Logical/ Conceptual level.

Data Layers (Abstraction)

Definition

Data Abstraction and Data Independence: Database systems comprise of complex data-structures. In order to make the system efficient in terms of retrieval of data, and reduce complexity in terms of usability of users, developers use abstraction i.e. hide irrelevant details from the users. This approach simplifies database design.

- There are 3 levels of data abstraction.
 - Physical level
 - Logical/ Conceptual level.
 - View level.