# Data Engineer Interview Preparation

## Table of Contents

## Introduction to Data Engineering Interviews

 - Data engineers are responsible for designing, building, and maintaining the systems that store and process large-s

 - Key skills include programming, data modeling, ETL development, and familiarity with big data tools.


## General Data Engineering Concepts

 - Data Warehousing: Centralized repository for structured data.

 - ETL/ELT Processes: Extract, Transform, Load vs. Extract, Load, Transform.

 - Data Modeling: Difference between Star and Snowflake schema.

 - Data Pipelines: Automation of data flow between systems.


## Programming and Coding

 - SQL: Writing complex queries for data extraction and transformation.

 - Python: Implementing ETL scripts and working with libraries like pandas, numpy.

 - Java/Scala: Frequently used for big data processing in Hadoop/Spark ecosystems.