

Data Extraction in ETL			
QUESTION 1	Describe different types of data sources used in ETL with suitable examples		
ANSWER	In ETL, data sources are places from where data is extracted before transforming and loading it into a data warehouse		
	1) APIs / Web Services	3) Relational Databases (RDBMS)	
	Data fetched from external applications using APIs.  Examples: REST API SOAP API	These store data in tables (rows & columns).  Examples: MySQL Oracle SQL Server PostgreSQL	
	Example: Getting weather or payment data	Example use case: <b>A company stores customer and sales data in MySQL.</b>	
	2) Data Warehouses		
	Central repositories that store historical and analytical data.  Examples: Amazon Redshift Snowflake Azure Synapse Google BigQuery		
	Example:  Sales data loaded daily into Snowflake for reporting.		

<b>QUESTION 2</b>	<b>What is data extraction? Explain its role in the ETL pipeline.</b>					
<b>ANSWER</b>	<p>Data Extraction is the first step of the ETL process.      It means collecting or pulling data from different source systems so it can be processed further.</p> <p>Data extraction - getting data from where it is stored</p> <p>Data extraction is the process of retrieving data from various data sources such as databases, files, APIs, or applications and moving it to a staging area for transformation.</p> <p><b>Role of Data Extraction in the ETL Pipeline</b></p> <p><b>EXTRACT</b></p> <p><b>TRANSFORM</b></p> <p><b>LOADING</b></p> <ul style="list-style-type: none"> <li>Connects to multiple data sources</li> <li>Reads data without changing it</li> <li>Moves data to a staging area</li> <li>Ensures data is complete and accurate</li> <li>Minimizes impact on source systems</li> </ul>					

<b>QUESTION 3</b>	<b>Explain the difference between CSV and Excel in terms of extraction and ETL usage</b>
<b>ANSWER</b>	CSV and Excel are both flat file data sources, but they behave very differently in ETL.
	<b>CSV (Comma-Separated Values)</b>
	What is CSV
	Plain text file
	Data separated by commas
	No formatting, no formulas
	<b>Very easy and fast to extract</b>
	Lightweight file
	Same structure always
	Supported by all ETL tools
	No dependency on Excel software
	<b>Excel File (.xls / .xlsx)</b>
	What is Excel?
	Binary / structured file
	Supports multiple sheets
	Can contain formulas, charts, formatting

<b>QUESTION 4</b>	<b>Explain the steps involved in extracting data from a relational database</b>			
<b>ANSWER</b>	<b>Identify the Data Source</b>			
	First, decide from which database and tables data is required.			
	Example:			
	<b>Database: sales_db</b>  <b>Tables: customers, orders</b>			
	Purpose: Avoid extracting unnecessary data.			
	<b>Database Connection</b>			
	Create a secure connection using:			
	Host name			
	Port number			
	Username & password			
	Database name			
	<b>Select Required Data (Write SQL Query)</b>			
	Use SELECT queries to extract required columns and rows.			

<b>QUESTION 5</b>	<b>Explain three common challenges faced during data extraction.</b>			
<b>ANSWER</b>	Data extraction is the first step of ETL, but it comes with several challenges that can affect the entire pipeline.			
	<b>1) Data Quality Issues</b> <ul style="list-style-type: none"> <li>Missing values (NULLs)</li> <li>Duplicate records</li> <li>Incorrect or inconsistent data</li> </ul>			
<b>Example:</b>	Customer age is NULL			
	<b>2) Performance &amp; Source System Load</b> <ul style="list-style-type: none"> <li>Extracting large volumes of data</li> <li>Heavy queries slow down source systems</li> </ul>			
	<b>3) Data Format &amp; Schema Changes</b> <ul style="list-style-type: none"> <li>Source schema changes frequently</li> <li>Columns added, removed, or datatype changed</li> </ul>			
<b>Example:</b>	New column discount added			

<b>QUESTION 6</b>	<b>What are APIs? Explain how APIs help in real-time data extraction.</b>
<b>ANSWER</b>	<b>API stands for Application Programming Interface.</b>
	An API is a set of rules and endpoints that allows one application to communicate with another application and exchange data.
	<b>1) Real-Time Data Access</b>
	APIs allow data to be fetched immediately when an event occurs.
<b>example-</b>	New order placed → API sends order data instantly
	<b>2) Event-Driven Data Extraction</b>
	APIs support: Webhooks Streaming endpoints
<b>Example:</b>	Payment completed → API triggers data extraction
	<b>3) Access to Live Systems</b>
	APIs connect directly to: Cloud applications SaaS platforms
<b>Examples:</b>	Salesforce API Google Analytics API

<b>QUESTION 7</b>	<b>Why are databases preferred for enterprise-level data extraction?</b>
<b>ANSWER</b>	Enterprises handle large, critical, and continuously growing data, so databases are the most reliable source for extraction.
	<b>1) Structured and Consistent Data</b>
	Databases store data in tables with fixed schema Enforced by constraints (PK, FK, NOT NULL)
	<b>2) Handles Large Volumes of Data Efficiently</b>
	Designed for millions of records Optimized with indexes and partitions
	<b>3) Supports Incremental Data Extraction</b>
	Uses timestamps, IDs, or CDC (Change Data Capture)
	<b>4) High Performance and Reliability</b>
	Databases support transactions Ensure ACID properties

<b>QUESTION 8</b>	<b>What steps should an ETL developer take when extracting data from large CSV files (1GB+)?</b>
<b>ANSWER</b>	Large CSV files need careful handling to avoid memory issues, slow performance, and failures.
	<b>1) Understand the File Structure First</b>
	Before extraction: Check column count & order Verify delimiter (,   ;) Identify header rows
	<b>2) Avoid Loading the Entire File into Memory</b>
	Don't read the full CSV at once.
	<b>3) Split the File (If Possible)</b>
	For very large files: Split into smaller chunks (e.g., 200MB each)
	<b>4) Use Staging Tables</b>
	Load raw CSV data into: Staging tables Temporary tables