

1. Describe different types of data sources used in ETL with suitable examples.

Relational Databases

These are structured data sources that store data in tables with rows and columns.

Examples:

- MySQL
- Oracle
- SQL Server
- PostgreSQL

Use case in ETL:

Extract customer, sales, or transaction data using SQL queries.

Example:

A company extracts sales data from a MySQL database to load it into a data warehouse.

Flat Files

These are simple file-based data sources, usually stored locally or on cloud storage.

Examples:

- CSV (Comma Separated Values)
- TXT
- Excel files (.xls, .xlsx)
- JSON, XML files

Use case in ETL:

Common for data exchange between systems or for importing/exporting reports.

Example:

Daily employee attendance data stored in a CSV file is loaded into a reporting system.

Data Warehouses

These are centralized repositories designed for analytics and reporting.

Examples:

- Amazon Redshift
- Google BigQuery
- Snowflake
- Azure Synapse

Use case in ETL:

Data is extracted from one warehouse and transformed for another analytics platform.

Example:

Migrating historical data from Redshift to Snowflake.

Cloud-Based Applications (SaaS)

These systems generate data through business applications hosted in the cloud.

Examples:

- Salesforce
- Google Analytics
- HubSpot
- Workday

Use case in ETL:

ETL tools use APIs to extract data periodically.

Example:

Extracting customer lead data from Salesforce for sales analysis.

APIs and Web Services

Data is fetched directly from applications using REST or SOAP APIs.

Examples:

- Social media APIs (Twitter, Facebook)
- Weather APIs
- Payment gateways (Stripe, PayPal)

Use case in ETL:

Useful when no direct database access is available.

Example:

Pulling transaction data from Stripe API into a data warehouse.

2. What is data extraction? Explain its role in the ETL pipeline.

Data extraction is the process of retrieving raw data from multiple data sources such as databases, files, APIs, or applications without changing its structure or meaning.

The main goal is to pull accurate and complete data from source systems while minimizing impact on their performance.

Role of Data Extraction in the ETL Pipeline

Starting Point of ETL

Extraction is the foundation of the ETL pipeline.

If incorrect or incomplete data is extracted, the entire ETL process will produce unreliable results.

Supports Multiple Data Sources

Data extraction allows ETL tools to collect data from:

- Relational databases (MySQL, SQL Server)
- Flat files (CSV, Excel)
- Cloud applications (Salesforce)
- APIs and web services
- Big data platforms

This enables data integration from heterogeneous systems.

Ensures Data Availability for Transformation

Extracted data is usually stored in a staging area, where it becomes available for:

- Data cleaning
- Data validation
- Data transformation
- Data standardization

Without extraction, transformation and loading cannot take place.

Controls Data Volume and Frequency

Extraction determines:

- What data to extract (full or partial)
- When to extract (daily, hourly, real-time)
- How much data to extract

Types of extraction:

- Full extraction – entire dataset is extracted
- Incremental extraction – only new or changed data is extracted

Maintains Source System Performance

Well-designed extraction processes:

- Run during off-peak hours

- Use optimized queries
 - Avoid locking source tables
- This ensures minimal impact on production systems.

Improves Data Quality and Reliability

During extraction, basic checks can be applied such as:

- Removing duplicate records
- Handling missing values
- Validating data formats

This improves the overall quality of data entering the ETL pipeline.

3. Explain the difference between CSV and Excel in terms of extraction and ETL usage.

CSV and Excel are both common file-based data sources, but they behave very differently during extraction and ETL.

File Structure

CSV (Comma-Separated Values)

- Plain text file
- Stores only raw data values
- No formatting, formulas, or multiple sheets

Excel

- Binary or XML-based file (.xls / .xlsx)
- Can contain multiple sheets
- Supports formatting, formulas, charts, and macros

ETL impact:

CSV is simpler and cleaner to extract; Excel needs more handling.

Ease of Data Extraction

CSV

- Very easy to extract
- Read line by line
- Supported by almost all ETL tools

Excel

- Extraction is more complex
- ETL tools must:
 - Choose the correct sheet
 - Handle merged cells and formulas

ETL impact:

CSV extraction is faster and more reliable than Excel.

Data Types Handling

CSV

- No fixed data types
- Everything is treated as text initially
- Data types are assigned during transformation

Excel

- Has implicit data types (number, date, text)
- May cause inconsistencies (e.g., dates auto-formatted)

ETL impact:

CSV gives better control during transformation; Excel can introduce data-type errors.

Data Quality and Consistency

CSV

- High consistency
- Less chance of hidden issues

Excel

- Prone to problems like:
 - Blank rows
 - Hidden columns
 - Inconsistent formats
 - Manual edits

ETL impact:

Excel usually requires extra data cleaning.

Performance and Scalability

CSV

- Lightweight and fast
- Handles large datasets well

Excel

- Slower for large files
- Has row limits (especially older .xls)

ETL impact:

CSV is preferred for large-scale ETL pipelines.

Automation in ETL

CSV

- Ideal for automated ETL jobs
- Easy to schedule and process

Excel

- Not ideal for automation
- Often needs manual intervention

ETL impact:

CSV is better for daily or real-time ETL.

4. Explain the steps involved in extracting data from a relational database.

Identify the Source Database

First, determine **where the data resides**.

Details include:

- Database type (MySQL, Oracle, SQL Server, PostgreSQL)
- Server name / IP address
- Port number
- Database name

Example:

Sales data stored in a MySQL database.

Establish Database Connection

Create a secure connection between the ETL tool and the database.

Requirements:

- Username and password
- JDBC / ODBC drivers
- Network access permissions

Purpose:

Allows the ETL tool to communicate with the database.

Select Required Tables or Views

Decide **what data to extract**.

Options include:

- Entire tables
- Specific columns
- Database views (preferred for security and performance)

Example:

Extract customer_id, order_date, total_amount from orders table.

Define Extraction Query

Write SQL queries to control the data extraction.

Types of queries:

- Full extraction
- SELECT * FROM orders;
- Incremental extraction
- SELECT * FROM orders WHERE last_updated > '2025-02-01';

Purpose:

Limits data volume and improves efficiency.

Apply Filters and Conditions

Use WHERE, JOIN, or LIMIT clauses to extract only relevant data.

Example:

Extract only completed orders:

WHERE order_status = 'Completed'

Extract Data to Staging Area

Move extracted data into a **staging area** such as:

- Temporary tables
- Flat files (CSV)
- Cloud storage

Purpose:

Isolates raw data from transformation and loading processes.

Perform Basic Validation

Check extracted data for:

- Row count matching
- Null values
- Data type mismatches
- Duplicate records

Purpose:

Ensures data completeness and correctness.

Log and Monitor the Extraction Process

Record extraction details such as:

- Start and end time
- Number of records extracted
- Errors or failures

Purpose:

Helps in troubleshooting and auditing.

5. Explain three common challenges faced during data extraction.

Data Quality Issues

Source data is often **incomplete, inconsistent, or inaccurate**.

Problems include:

- Missing or NULL values
- Duplicate records
- Incorrect data formats (dates, numbers)

Impact:

Poor-quality data leads to incorrect transformations and unreliable reports.

Example:

A customer table contains multiple date formats (DD-MM-YYYY, MM/DD/YYYY), causing extraction errors.

Performance Impact on Source Systems

Extraction queries can **slow down production databases**, especially when large volumes of data are pulled.

Problems include:

- Full table scans
- Heavy joins on live systems
- Extraction during peak business hours

Impact:

May affect application performance and user experience.

Example:

Running a full extraction on a sales database during business hours causes slow order processing.

Handling Large Data Volumes

Modern systems generate **huge amounts of data**, making extraction difficult.

Problems include:

- Long extraction times
- Network bandwidth limitations
- Memory and storage constraints

Impact:

ETL jobs may fail or miss scheduled deadlines.

Example:

Extracting millions of transaction records daily without incremental logic causes ETL delays.

6. What are APIs? Explain how APIs help in real-time data extraction.

API (Application Programming Interface) is a set of rules and endpoints that allows one application to communicate with another and exchange data securely.

Instead of directly accessing a database, systems request data through APIs using standard methods like:

- GET – retrieve data
- POST – send data
- PUT/PATCH – update data
- DELETE – remove data

Most modern APIs are REST APIs and return data in formats like JSON or XML.

How APIs Help in Real-Time Data Extraction

Enable Instant Data Access

APIs allow systems to fetch data the moment it is generated, instead of waiting for batch files or scheduled jobs.

Example:

When a customer makes a payment, a payment gateway API immediately sends transaction details.

Support Event-Based or Streaming Data

APIs can push data whenever an event occurs.

How it works:

- An event happens (order placed, user login)
- API sends data instantly to the ETL pipeline
- Data is processed and stored in near real time

Example:

E-commerce order API sending live order data to a data warehouse.

Avoid Direct Database Access

APIs act as a secure layer between systems.

Benefits:

- No need for database credentials
- Reduced risk of data corruption
- Better access control

Example:

Extracting customer data from Salesforce API instead of accessing its internal database.

Enable Incremental and Real-Time Extraction

APIs can return:

- Only new data
- Only updated records
- Data after a specific timestamp

This avoids full data extraction every time.

Example:

GET /orders?updated_after=2025-02-06T10:00:00

Work Well with Cloud and SaaS Applications

Most cloud applications expose data only through APIs.

Examples:

- Google Analytics API
- Twitter API
- Stripe API
- Salesforce API

APIs are the only way to extract real-time data from these platforms.

Improve ETL Automation

APIs allow ETL pipelines to:

- Run continuously
- Trigger extraction automatically
- Scale easily with traffic

7. Why are databases preferred for enterprise-level data extraction?

Databases are preferred for enterprise-level data extraction because they are designed to handle large volumes of data, multiple users, and mission-critical operations reliably and securely.

High Performance and Scalability

Enterprise databases can process millions of records efficiently.

Why it matters:

- Optimized query engines
- Indexes and partitions
- Parallel processing

Example:

Extracting years of transaction data from an Oracle or SQL Server database without system slowdown.

Structured and Consistent Data

Databases enforce:

- Data types
- Constraints (primary key, foreign key)
- Referential integrity

Why it matters:

ETL processes rely on clean, well-structured data.

Example:

Customer IDs are always numeric and unique, reducing transformation errors.

Support for Incremental Extraction

Databases make it easy to extract only new or changed data using:

- Timestamps
- Change Data Capture (CDC)
- Triggers or log-based replication

Why it matters:

Reduces load, improves speed, and supports near real-time ETL.

Reliability and Data Integrity

Enterprise databases ensure:

- ACID transactions
- Backup and recovery
- Data consistency even during failures

Why it matters:

ETL pipelines can trust the extracted data.

Security and Access Control

Databases provide strong security features:

- User roles and privileges
- Encryption
- Auditing and logging

Why it matters:

Sensitive enterprise data (finance, HR) remains protected.

Easy Integration with ETL Tools

Most ETL tools natively support databases using:

- JDBC / ODBC connectors

- Optimized database adapters
- Why it matters:
Simplifies extraction and reduces development effort.

Automation and Scheduling

Databases support:

- Scheduled jobs
- Stored procedures
- Views for extraction

Why it matters:

Ideal for daily, hourly, or real-time ETL jobs.

Better Data Governance

Databases help enforce:

- Data lineage
- Auditing
- Compliance standards

Why it matters:

Essential for enterprise reporting and regulatory compliance.

8. What steps should an ETL developer take when extracting data from large CSV files (1GB+)?

Steps to Extract Data from Large CSV Files (1GB+)

Analyze the File Before Extraction

Understand the file structure:

- Number of rows and columns
- Delimiter used (comma, pipe, tab)
- Presence of headers
- Encoding (UTF-8, UTF-16)

Use Chunk / Batch Processing

Never load the entire file into memory.

Best practice:

- Read data in **chunks** (e.g., 10k–100k rows at a time)

Use Streaming or Line-by-Line Reading

Process the file **row by row** instead of loading it fully.

Define Schema Explicitly

Do not rely on auto-detection.

Define:

- Column names
- Data types
- Date formats

Handle Bad Records Separately

Identify and isolate:

- Corrupted rows
- Invalid delimiters
- Missing mandatory fields

Approach:

- Reject file or bad records table
- Error logs

Apply Minimal Transformations During Extraction

Keep extraction lightweight.

Avoid:

- Complex joins
- Heavy calculations

Use Compression If Supported

If possible, extract from **compressed CSV files (.gz, .zip)**.

Parallelize When Possible

Split files by:

- Date
- Region
- Logical partitions

Process them in parallel.

Load into a Staging Area First

Store extracted data into:

- Staging tables
- Temporary files

Validate Data After Extraction

Perform checks such as:

- Row count validation
- Null checks
- File completeness

Log, Monitor, and Recover

Capture:

- Start/end time
- Records processed
- Errors encountered

Enable **restartability** from last successful checkpoint.