# Data Integration and Transformation Tools

Data integration is the process of combining data from multiple sources into a unified view, while transformation involves cleaning, reshaping, and structuring data for analysis.

**Key Concepts:**

- **ETL (Extract, Transform, Load):** Data is extracted from sources, transformed (cleaned and formatted), and then loaded into a data warehouse.
- **ELT (Extract, Load, Transform):** Data is first loaded into a system and then transformed, allowing for scalability.
- **Data Pipelines:** Automated workflows that move and process data (e.g., Apache Airflow, AWS Glue).

## Data Integration and Transformation Tools:

1. **Apache Airflow**
   - open source platform for programmatically authoring   scheduling, and monitoring workflows
   - Created originally by Airbnb
   - Allows users to define and execute complex workflows
   - Support for:
     - Task dependencies
     - Parallelism
     - Error handling

2. **Kubeflow**
   - An open-source machine learning toolkit that allows execution of data science pipelines on top of Kubernetes.
   - Provides a platform for building, deploying, and managing end-to-end machine learning workflows at scale
   - Support for:
     - Distributed training
     - Model serving
     - Hyperparameter tuning

3. **Apache Kafka**

- Distributed streaming platform that allows applications to publish, process, and subscribe to streams of records in real-time
- Created originally from LinkedIn.
- It is scalable, fault-tolerant, and high-throughput
- Suitable for building mission-critical, data-intensive applications

4. **Apache Nifi**
   - An open-source data integration platform that allows users to automate the flow of data between systems
   - Provides a web-based user interface for designing and managing data flows
   - Support for:
     - Data routing
     - Transformation
     - Enrichment
     - Among other capabilities

5. **Apache Spark SQL**
   - A module in the Spark ecosystem that provides a programming interface for working with structured data using:
     - SQL
     - Data frames
     - Datasets
   - Supports a wide range of data sources and provides optimized performance for complex data processing tasks.

6. **Node Red**
   - An open-source visual programming tool for wiring together hardware devices, APIs, and online services
   - Allows users to create event-driven flows of messages
   - low in resource consumption that it even runs on tiny devices like a Raspberry Pi.
   - Support for:
     - Data transformation
     - Filtering
     - Aggregation