

Question 1: How would you explain the differences between batch processing and stream processing in a data pipeline context, and when would you choose one over the other?

Answer: Batch processing involves collecting data over a period and processing it as a group or "batch." Stream processing handles data continuously as it arrives in real-time.

The key differences include:

- Batch processing operates on finite datasets with higher latency but typically higher throughput. Stream processing works on potentially infinite data streams with lower latency.
- Batch is generally simpler to implement and debug, while streaming requires more complex state management and error handling.
- Batch processing is often cheaper for large data volumes but lacks real-time insights.

I would choose batch processing for:

- Cost-efficient processing of large datasets
- When immediate insights aren't critical (e.g., nightly reports, historical analysis)
- Complex transformations requiring the entire dataset

I would choose stream processing for:

- Real-time analytics and dashboards
- Fraud detection or anomaly detection
- Time-sensitive business operations (inventory management, transaction processing)
- Event-driven architectures

At Ameris Bank, I implemented both: batch processing for end-of-day financial reconciliation and stream processing with Kafka/Kinesis for real-time transaction monitoring and fraud detection.

Question 2: Can you walk me through how you would design an ETL pipeline using Apache Airflow? What are the key components and considerations?

Answer: When designing an ETL pipeline with Apache Airflow, I follow these steps:

1. **DAG Definition:** I start by creating a directed acyclic graph (DAG) that clearly represents the workflow's dependency structure, scheduling interval, and retry policies.
2. **Task Design:** I break down the ETL process into atomic tasks using appropriate operators:
 - PythonOperator for custom transformation logic
 - BashOperator for shell commands
 - Database operators (PostgresOperator, MySqlOperator) for SQL execution
 - Cloud-specific operators (S3Operator, RedshiftOperator) for cloud integrations

3. **Task Dependencies:** I establish dependencies using `>>` or `<<` operators or explicit `set_upstream/set_downstream` methods.
4. **Data Quality Checks:** I incorporate data validation tasks after extraction and transformation to ensure data integrity.
5. **Error Handling:** I implement robust error handling with custom callbacks and retries for transient failures.
6. **Parameter Management:** I use Airflow variables, connections, and XComs for secure credential management and task communication.
7. **Monitoring:** I set up SLAs, alerts, and integrate with monitoring systems like Grafana.

Key considerations include:

- Idempotency: Ensuring tasks can be re-run safely
- Resource management: Preventing worker overload with pools
- Performance: Using appropriate execution methods (sequential, parallel)
- Maintainability: Following consistent patterns and documentation

At Novartis, I implemented this approach to process clinical research data, incorporating rigorous data quality checks to ensure regulatory compliance.

Question 3: What strategies have you used to optimize Spark jobs for better performance?

Answer: To optimize Spark jobs, I've implemented several strategies:

1. **Proper Resource Allocation:**
 - Right-sizing executor memory and cores based on workload
 - Setting appropriate partition count using `spark.sql.shuffle.partitions` or `repartition()`
 - Using dynamic resource allocation when applicable
2. **Data Optimization:**
 - Partitioning data appropriately based on access patterns
 - Using appropriate file formats (Parquet, ORC) with compression
 - Implementing bucketing for join-heavy workloads
 - Caching frequently accessed DataFrames with `cache()` or `persist()`
3. **Query Optimization:**
 - Leveraging predicate pushdown by filtering early
 - Using broadcast joins for small-large table joins (`broadcast()/hint`)
 - Optimizing UDFs with Pandas UDFs when possible
 - Analyzing query plans with `explain()` to identify bottlenecks
4. **Code Efficiency:**
 - Minimizing shuffles and wide transformations
 - Using SQL for complex transformations when more efficient
 - Avoiding unnecessary actions that trigger evaluations
5. **Memory Management:**
 - Tuning garbage collection settings

- Managing spill to disk with `spark.memory.fraction`
- Handling skew with salting or custom partitioners

At IKEA, I applied these techniques to optimize our nightly sales data processing, reducing processing time by 60% and enabling earlier access to business insights.

Question 4: How do you approach data partitioning in systems like AWS S3 or Redshift, and what factors influence your partitioning strategy?

Answer: My approach to data partitioning varies by system but follows key principles:

For S3:

- I partition based on common query filters (date, region, customer segment)
- I use a hierarchical structure (YYYY/MM/DD) for time-based partitions
- I balance partition granularity—too fine increases overhead, too coarse loses query benefits
- I implement partition pruning by using the S3 list API efficiently
- I use Athena/Glue partitioning for optimal query performance

For Redshift:

- I choose distribution keys based on join conditions to minimize data movement
- I select sort keys that match common WHERE clauses
- I consider data skew to avoid hot spots on specific slices
- I implement proper compression encodings based on data characteristics
- I use multi-column sort keys for tables with varied query patterns

Factors influencing my strategy:

1. **Query patterns:** I analyze common WHERE clauses and JOIN conditions
2. **Data volume:** Higher volumes need more careful partitioning
3. **Data ingestion frequency:** Affects partition size decisions
4. **Update patterns:** Determines whether I need mutable partitions
5. **Query latency requirements:** Real-time needs vs. batch analytics

At Ameris Bank, I implemented a partitioning strategy for transaction data by date, region, and transaction type, which reduced query execution time by 50% and significantly improved analytical capabilities.

Question 5: Can you explain the architecture of a Lambda or Kappa architecture for real-time data processing?

Answer: Lambda Architecture: Lambda architecture consists of three layers:

1. **Batch Layer:** Processes historical data in batches, creating comprehensive but delayed views. I typically implement this using technologies like Hadoop, Spark Batch, or AWS EMR.
2. **Speed Layer:** Processes data in real-time with lower latency but potentially lower accuracy. I implement this using stream processing frameworks like Kafka Streams, Flink, or Spark Streaming.
3. **Serving Layer:** Combines results from both layers to provide complete views. This might be a database like Cassandra or a data warehouse like Redshift.

The main advantage is having both accurate historical analysis and real-time insights. The challenge is maintaining duplicate processing logic in both layers.

Kappa Architecture: Kappa architecture simplifies Lambda by eliminating the batch layer:

1. All data flows through a single stream processing pipeline
2. Historical data is simply reprocessed through the same pipeline when needed
3. Results are stored in a scalable serving layer

I implement Kappa using technologies like Kafka for the log storage, Kafka Streams or Flink for processing, and systems like Cassandra or Elasticsearch for serving.

At Ameris Bank, I implemented a hybrid architecture using Kafka and Kinesis for real-time transaction monitoring, with periodic batch processing for reconciliation, which gave us both real-time fraud detection capabilities and accurate historical reporting.

Question 6: What's your approach to handling schema evolution in a data warehouse environment?

Answer: My approach to schema evolution in data warehouses focuses on maintaining both data integrity and backward compatibility:

1. **Schema Design Principles:**
 - I design schemas with future evolution in mind, avoiding overly restrictive constraints
 - I use default values for new columns where appropriate
 - I implement nullable fields for optional attributes
2. **Change Management Process:**
 - I maintain a schema registry documenting all changes
 - I follow a strict versioning system for schema updates
 - I implement a staged deployment process (dev→test→prod)
 - I communicate changes to stakeholders with advance notice
3. **Technical Implementation:**
 - For additive changes (new columns): I implement them as nullable fields
 - For restructuring: I create views to maintain backward compatibility
 - For data type changes: I use temporary staging tables for conversion
 - For breaking changes: I maintain both old and new structures during transition

4. **Tools and Technologies:**

- I use schema evolution features in formats like Avro, Parquet with Hive metastore
- I implement CDC (Change Data Capture) for tracking schema changes
- I leverage schema validation in ETL pipelines
- For Snowflake or modern warehouses, I utilize time travel features

At Novartis, I implemented this approach when expanding our clinical trial database schema, maintaining backward compatibility for existing reports while adding new attributes for enhanced analytics.

Question 7: How would you implement data quality checks in an automated ETL pipeline?

Answer: I implement data quality checks at multiple stages in ETL pipelines:

1. **Source Data Validation:**

- Record count verification against source systems
- Schema validation to detect unexpected changes
- Data freshness checks to confirm timeliness
- File/table integrity validation

2. **During Transformation:**

- Null checks for required fields
- Data type validation and conversion monitoring
- Business rule validation (e.g., age ranges, valid codes)
- Referential integrity checks
- Duplicate detection

3. **Post-Load Validation:**

- Row count reconciliation between source and target
- Statistical validation (min/max/avg/distribution)
- Historical trend analysis to detect anomalies
- End-to-end checksums for data integrity

4. **Implementation Approaches:**

- I build custom validation operators in Airflow
- I use Great Expectations library for complex validations
- I implement threshold-based alerting for key metrics
- I maintain data quality dashboards for monitoring

5. **Handling Issues:**

- For critical issues: Pipeline failure and immediate alerts
- For non-critical: Warnings with detailed logs
- Quarantine mechanisms for invalid records
- Automated retry logic where appropriate

At Ameris Bank, I implemented a comprehensive validation framework with Airflow that reduced data quality issues by 75% and ensured regulatory compliance for financial reporting.

Question 8: Can you explain the differences between HDFS and object storage systems like S3 in terms of data engineering workloads?

Answer: HDFS and S3 have fundamental architectural differences that impact data engineering workloads:

HDFS (Hadoop Distributed File System):

- **Architecture:** Tightly coupled compute and storage in the same cluster
- **Performance:** Optimized for high throughput sequential reads and the locality principle
- **Consistency:** Strong consistency model with immediate visibility of writes
- **Scalability:** Scales by adding data nodes but has practical limits
- **Durability:** Relies on replication factor (typically 3x)
- **Cost Structure:** Higher fixed costs for maintenance and infrastructure

S3 (Simple Storage Service):

- **Architecture:** Decoupled compute and storage enabling serverless paradigms
- **Performance:** Lower latency for random access but potential throughput limitations
- **Consistency:** Eventually consistent (though now offers strong consistency options)
- **Scalability:** Virtually unlimited storage capacity
- **Durability:** 11 nines (99.999999999%) durability
- **Cost Structure:** Pay-per-use model with no minimum commitment

Impact on Data Engineering Workloads:

1. **ETL Processing:**
 - With HDFS: Better performance for large transformations within the cluster
 - With S3: More flexibility to scale compute independently; better for serverless ETL
2. **Data Lake Implementation:**
 - HDFS requires more management but offers better performance for Hadoop ecosystem
 - S3 simplifies operations and enables broader access patterns
3. **Cost Efficiency:**
 - HDFS more cost-effective for consistent, high-volume workloads
 - S3 more economical for variable or growing workloads

At Ameris Bank, I transitioned from HDFS to S3-based architecture, implementing optimal partitioning and Spark tuning to maintain performance while reducing infrastructure costs by 35%.

Question 9: What techniques do you use to monitor the health and performance of your data pipelines?

Answer: My comprehensive monitoring approach for data pipelines includes:

1. **Operational Metrics:**
 - Pipeline run status (success/failure rates)
 - Job duration tracking with historical comparisons
 - Resource utilization (CPU, memory, disk I/O)
 - Throughput measurements (records processed per second)
 - Backlog size and processing lag metrics
2. **Data Quality Metrics:**
 - Validation failure rates and patterns
 - Schema drift detection
 - Data completeness metrics
 - Error rates by error type
 - Data freshness measurements
3. **Business Impact Metrics:**
 - SLA compliance monitoring
 - Critical table update tracking
 - Downstream system impact alerts
 - Data availability for business processes
4. **Implementation Tools:**
 - Prometheus for metrics collection
 - Grafana for visualization and dashboards
 - ELK stack for log analysis and centralization
 - Custom alerting through PagerDuty or similar
 - Airflow's built-in monitoring capabilities
5. **Proactive Monitoring:**
 - Predictive failure analysis based on historical patterns
 - Automated canary tests before full pipeline runs
 - Synthetic testing for critical components
 - Automated recovery for known failure patterns

At Novartis, I implemented a real-time monitoring system that reduced mean time to resolution for pipeline issues by 65% through early detection and automated diagnostics of common failure patterns.

Question 10: How do you design a data lake to ensure it doesn't become a "data swamp"?

Answer: To prevent a data lake from becoming a "data swamp," I implement these key design principles:

1. **Clear Layered Architecture:**
 - Raw/Bronze zone: Unmodified source data with retention policies
 - Processed/Silver zone: Cleansed, validated data
 - Curated/Gold zone: Business-ready datasets
 - Each zone has distinct access patterns and governance
2. **Robust Metadata Management:**
 - Comprehensive data catalog with search capabilities

- Automated metadata extraction during ingestion
- Business glossary aligned with technical metadata
- Data lineage tracking across transformations
- 3. **Standardized Organization:**
 - Consistent naming conventions (database/table/column)
 - Logical partitioning strategies
 - Standardized folder structures
 - Format standardization (Parquet/ORC with optimized compression)
- 4. **Data Governance Implementation:**
 - Clear ownership for each dataset
 - Access control at appropriate granularity
 - Data quality SLAs and monitoring
 - Automated policy enforcement
- 5. **Lifecycle Management:**
 - Explicit retention policies by data category
 - Automated archival and purging processes
 - Version control for key datasets
 - Storage tiering based on access patterns
- 6. **Self-Service Capabilities:**
 - Intuitive data discovery interfaces
 - Standardized access methods and documentation
 - Sample queries and usage examples
 - Data profiling and preview capabilities

At Ameris Bank, I established a comprehensive data governance framework with data lineage tracking and metadata management, ensuring our data lake remained organized and valuable for business users.

Question 11: Explain the concept of slowly changing dimensions and how you would implement them in a data warehouse.

Answer: Slowly Changing Dimensions (SCDs) are dimension tables that change gradually over time, requiring strategies to track historical changes.

SCD Types:

- **Type 0:** No history tracking; data overwritten
- **Type 1:** No history; current value replaces old value
- **Type 2:** Full history via new rows for changes with effective dates
- **Type 3:** Limited history via additional columns for previous values
- **Type 4:** History in separate history table with current values in main table
- **Type 6:** Hybrid approach combining Types 1, 2, and 3

My Implementation Approach:

For **Type 2** (most common):

1. I add technical columns:
 - `effective_start_date` and `effective_end_date`
 - `is_current` flag (boolean)
 - `surrogate_key` (not natural key)
 - `version_number` (optional)
2. ETL process:
 - Compare incoming data with current records
 - For changed records: Update current record's end date and current flag
 - Insert new record with current start date and current flag
 - Handle surrogate key relationships in fact tables
3. Optimization techniques:
 - Partitioning by effective date ranges
 - Indexing on frequently queried columns
 - Materialized views for common historical queries

Practical Considerations:

- I carefully select which attributes trigger new versions
- I implement change detection logic efficiently
- I manage storage growth with appropriate purging strategies
- I provide user-friendly view layers that default to current records

At Novartis, I implemented Type 2 SCDs for patient dimension tables to maintain a complete audit trail of patient information changes while ensuring current analytical queries remained efficient.

Question 12: How would you approach migrating from an on-premises Hadoop cluster to a cloud-based solution?

Answer: My approach to migrating an on-premises Hadoop cluster to the cloud follows these key phases:

1. **Assessment & Planning:**
 - Inventory of existing datasets, jobs, and workflows
 - Performance profiling of current workloads
 - Identification of security and compliance requirements
 - Selection of target cloud platform (AWS, Azure, GCP)
 - Total cost of ownership analysis
 - Creation of migration timeline and success metrics
2. **Architecture Design:**
 - Selecting appropriate cloud services (e.g., EMR/Databricks/Dataproc)
 - Designing network topology and security architecture
 - Planning for service integration (authentication, monitoring)
 - Optimizing for cloud cost models (spot instances, storage tiers)
 - Creating disaster recovery and business continuity plans
3. **Data Migration Strategy:**

- Data prioritization based on business value and complexity
 - Choosing transfer mechanisms (Direct Connect, DataSync, Snowball)
 - Implementing change data capture for ongoing synchronization
 - Developing validation procedures for data integrity
4. **Workload Migration:**
- Converting Hadoop jobs to cloud-native equivalents
 - Refactoring for cloud-specific optimizations
 - Implementing CI/CD pipelines for code deployment
 - Testing performance and scalability in cloud environment
5. **Operational Transition:**
- Training teams on cloud technologies and tools
 - Establishing cloud monitoring and alerting
 - Implementing cloud cost management practices
 - Creating runbooks for common operations
6. **Cutover & Decommissioning:**
- Phased cutover strategy with rollback capability
 - Final data synchronization and validation
 - Traffic redirection to cloud services
 - Decommissioning plan for on-premises hardware

At Novartis, I led the migration of on-premises clinical databases to Azure Cloud infrastructure, resulting in 99.9% uptime and improved data availability for global research teams while ensuring regulatory compliance throughout the transition.