# AWS Data Engineering Services - Quick Reference Cheat Sheet

## Data Ingestion Services

### Amazon Kinesis Family

| Service | Use Case | Key Features | Limits |
|---|---|---|---|
| **Kinesis Data Streams** | Real-time streaming | • Custom consumer apps<br>• Replay capability<br>• Low latency (< 1 second) | • 1MB record limit<br>• 1000 records/sec per shard |
| **Kinesis Data Firehose** | Near real-time delivery | • Serverless<br>• Built-in transformation<br>• Direct S3/Redshift delivery | • 60 second buffer minimum<br>• No replay capability |
| **Kinesis Analytics** | Real-time analytics | • SQL on streaming data<br>• Anomaly detection<br>• Time-windowed queries | • SQL-based processing only |

### AWS Glue

| Component | Purpose | Key Points |
|---|---|---|
| **Glue Crawlers** | Schema discovery | • Auto-detect schema changes<br>• Populate Data Catalog<br>• Schedule-based or on-demand |
| **Glue ETL Jobs** | Data transformation | • Serverless Spark<br>• Auto-scaling<br>• Built-in retry logic |
| **Glue Data Catalog** | Metadata repository | • Hive-compatible metastore<br>• Integration with Athena/EMR<br>• Schema versioning |
| **Glue DataBrew** | Visual data preparation | • No-code transformations<br>• Data profiling<br>• 250+ built-in transformations |

## Data Storage Services

### Amazon S3 Storage Classes

| Storage Class | Use Case | Retrieval Time | Cost |
|---|---|---|---|
| **Standard** | Frequently accessed | Immediate | Highest storage cost |
| **Standard-IA** | Infrequently accessed | Immediate | Lower storage, retrieval fee |
| **One Zone-IA** | Non-critical, infrequent | Immediate | 20% less than Standard-IA |
| **Glacier Instant** | Archive with instant access | Immediate | Lower storage cost |
| **Glacier Flexible** | Archive data | 1-12 hours | Very low storage cost |
| **Glacier Deep Archive** | Long-term archive | 12-48 hours | Lowest storage cost |

## Amazon Redshift

| Feature | Description | Best Practice |
|---------|-------------|---------------|
| **Distribution Keys** | How data is distributed | • Use JOIN columns<br>• Avoid high cardinality<br>• Consider EVEN for small tables |
| **Sort Keys** | Physical data ordering | • Use WHERE clause columns<br>• Consider compound vs interleaved<br>• Limit to 3-4 columns |
| **Compression** | Reduce storage/IO | • Use ANALYZE COMPRESSION<br>• Different encoding per column<br>• Automatic for new tables |
| **Workload Management** | Query prioritization | • Separate queues by workload<br>• Set memory allocation<br>• Use concurrency scaling |

## DynamoDB

| Concept | Description | Guidelines |
|---------|-------------|------------|
| **Partition Key** | Primary hash key | • High cardinality<br>• Uniform access pattern<br>• Avoid hot partitions |
| **Sort Key** | Range key for sorting | • Enable range queries<br>• Model 1:N relationships<br>• Support query patterns |
| **GSI/LSI** | Secondary indexes | • GSI: Different partition key<br>• LSI: Same partition key<br>• Max 20 GSI per table |
| **Capacity Modes** | Billing model | • On-Demand: Unpredictable<br>• Provisioned: Predictable + cheaper |

# Data Processing Services

## Amazon EMR

| Component | Purpose | Key Points |
|-----------|---------|------------|
| **Master Node** | Cluster management | • Manages cluster<br>• NameNode for HDFS<br>• Single point of failure |
| **Core Nodes** | Data storage + processing | • Run DataNode + TaskTracker<br>• HDFS storage<br>• Can be removed with care |
| **Task Nodes** | Processing only | • No HDFS storage<br>• Spot instances recommended<br>• Safe to terminate |

## AWS Lambda

| Aspect | Specification | Considerations |
|---|---|---|
| **Runtime** | 15 minutes max | • Use Step Functions for longer workflows<br>• Consider EMR for heavy processing |
| **Memory** | 128MB - 10GB | • CPU scales with memory<br>• Optimize for cost vs performance |
| **Triggers** | Event-driven | • S3 events, Kinesis, DynamoDB Streams<br>• EventBridge for schedules |
| **Concurrency** | 1000 default | • Can request increases<br>• Consider reserved concurrency |

# Analytics Services

## Amazon Athena

| Feature | Description | Optimization Tips |
|---|---|---|
| **Serverless SQL** | Query S3 data directly | • Use columnar formats (Parquet)<br>• Partition data by query patterns<br>• Compress data (GZIP, Snappy) |
| **Query Engines** | Presto/Trino based | • Use appropriate data types<br>• Avoid SELECT * queries<br>• Use LIMIT for exploration |
| **Workgroups** | Query organization | • Set data limits<br>• Control costs<br>• Separate environments |

## Amazon QuickSight

| Component | Purpose | Key Features |
|---|---|---|
| **SPICE** | In-memory engine | • Fast query performance<br>• Automatic data refresh<br>• 10GB per dataset |
| **Data Sources** | Input connections | • 30+ native connectors<br>• Direct query vs SPICE<br>• Row-level security |
| **Dashboards** | Visualization | • Interactive dashboards<br>• Mobile responsive<br>• Embedded analytics |

# Security Services

## AWS IAM for Data Engineering

| Policy Type | Use Case | Example |
|---|---|---|
| **Identity-based** | User/role permissions | Glue job execution role |
| **Resource-based** | Cross-account access | S3 bucket policy |
| **Session policies** | Temporary restrictions | Federated access limits |
| **Permissions boundaries** | Maximum permissions | Developer sandbox limits |

## AWS KMS

| Key Type | Management | Use Case |
|---|---|---|
| **AWS Managed** | AWS controls | • Default encryption<br>• Service-specific keys |
| **Customer Managed** | You control | • Custom key policies<br>• Cross-account access<br>• Key rotation control |
| **Customer Provided** | You provide | • Full control<br>• Higher complexity<br>• Import your own keys |

# Monitoring & Governance

## CloudWatch for Data Pipelines

| Metric Category | Examples | Alerting Strategy |
|---|---|---|
| **Glue Jobs** | • Job duration<br>• Success/failure rate<br>• DPU hours | • Set SLA-based alarms<br>• Monitor cost metrics |
| **Kinesis** | • IncomingRecords<br>• WriteProvisionedThroughputExceeded | • Shard-level monitoring<br>• Auto-scaling triggers |
| **Redshift** | • CPU utilization<br>• Disk space<br>• Query performance | • Performance alerts<br>• Storage warnings |

## AWS Lake Formation

| Feature | Purpose | Best Practice |
|---|---|---|
| **Data Permissions** | Fine-grained access control | • Column-level permissions<br>• Row-level security<br>• Tag-based policies |
| **Data Discovery** | Catalog population | • Automatic crawling<br>• ML-powered classification<br>• PII detection with Macie |
| **Data Sharing** | Cross-account access | • Resource sharing<br>• Query federation<br>• Audit trails |

# Common Architecture Patterns

## Lambda Architecture

Batch Layer: S3 → Glue/EMR → Redshift (historical data)
Speed Layer: Kinesis → Lambda → DynamoDB (real-time)
Serving Layer: Athena/QuickSight (unified view)

## Kappa Architecture

Stream Processing: Kinesis → Kinesis Analytics → Output
Everything is treated as a stream, including batch data

## Data Lake Pattern

Landing Zone (S3 Raw) → Processing (Glue/EMR) →
Curated Zone (S3 Processed) → Analytics (Athena/Redshift)

# Performance Optimization Quick Tips

### S3 Optimization

- Use prefixes to avoid hot spots

- Multipart upload for files > 100MB

- S3 Transfer Acceleration for global access

- CloudFront for frequently accessed data

### Redshift Optimization

- VACUUM regularly to reclaim space

- ANALYZE to update table statistics

- Use COPY command for bulk loads

- Monitor query performance with system tables

### Glue Optimization

- Use bookmark for incremental processing

- Optimize for fewer, larger files

- Use pushdown predicates

- Consider Glue streaming for low latency

Remember: The exam tests your ability to choose the right service for the right use case. Focus on understanding trade-offs between services!