# DATA ENGINEERING TOOLS – SEGREGATED

| Segment | General | AWS |
|---|---|---|
| **Data Ingestion** | **Ingestion Tools:**<br>Apache Kafka<br>Apache NiFi<br>AWS Kinesis<br>Logstash | **AWS Glue:** Use Glue Crawlers to discover and catalog metadata from various data sources.<br>Glue ETL jobs for transforming and loading data.<br><br>**Amazon Kinesis:** Kinesis Data Streams for real-time data streaming. Kinesis Data Firehose for loading streaming data into data stores.<br><br>**AWS DataSync:** Transfer data from on-premises to AWS. |
| **Data Storage** | **Data Warehouses:**<br>Amazon Redshift<br>Google BigQuery<br>Snowflake<br><br>**Data Lakes:**<br>Amazon S3<br>Azure Data Lake Storage<br>Google Cloud Storage<br><br>**Databases:**<br>PostgreSQL<br>MySQL<br>MongoDB<br>Cassandra | **Amazon S3:** As a data lake for storing raw and processed data. Versioning and lifecycle policies for managing data.<br><br>**Amazon Redshift:** For data warehousing and complex queries.<br><br>**Amazon DynamoDB:** For NoSQL database requirements. |
| **Data Processing** | **Batch Processing:**<br>Apache Spark<br>Apache Flink<br>Hadoop MapReduce<br><br>**Stream Processing:**<br>Apache Kafka Streams<br>Apache Storm<br>Apache Flink<br><br>**ETL (Extract,Transform, Load):**<br>Apache Beam<br>Apache Airflow<br>Talend | **Amazon EMR (Elastic MapReduce):**For big data processing using frameworks like Apache Spark and Hadoop.<br><br>**AWS Glue:** Serverless ETL service for data transformation and preparation.<br><br>**AWS Lambda:** For serverless event-driven processing. |
| **Data Transformation** | **Data Preparation:**<br>Pandas (Python library)<br>Apache Beam<br><br>**Data Cleansing:**<br>Trifacta<br>OpenRefine<br><br>**Data Masking/Anonymization:**<br>Google DLP<br>Apache Nifi | **AWS Glue:** Use Glue jobs for ETL transformations.<br><br>**AWS Step Functions:** Orchestrate and coordinate multiple AWS services in a serverless workflow. |

| | | |
|---|---|---|
| **Analytics and Reporting** | **Business Intelligence Tools:**<br>Tableau<br>Power BI<br>Looker<br><br>**Analytics Platforms:**<br>Databricks<br>Google Analytics<br>Mixpanel | **Amazon QuickSight:** Business intelligence service for visualizing and analyzing data.<br><br>**Amazon Athena:** Serverless query service for analyzing data in Amazon S3. |
| **Data Orchestration** | **Workflow Management:**<br>Apache Airflow<br>Luigi<br>Prefect<br><br>**Job Scheduling:**<br>Cron<br>Apache Oozie | **Apache Airflow on Amazon MWAA (Managed Workflows for Apache Airflow):** Orchestrate and schedule complex data workflows.<br><br>**AWS Step Functions:** For serverless workflow orchestration. |
| **Monitoring and Logging** | **Logging:**<br>ELK Stack (Elasticsearch, Logstash, Kibana)<br>Splunk<br><br>**Monitoring:**<br>Prometheus<br>Grafana | **Amazon CloudWatch:** For monitoring AWS resources and applications.<br><br>**AWS CloudTrail:** For logging AWS API calls. |
| **Data Data Quality and Governance** | **Data Quality Tools:**<br>Informatica<br>Talend<br>Apache Griffin<br><br>**Metadata Management:**<br>Collibra<br>Apache Atlas | **AWS Glue DataBrew:** For data profiling, cleaning, and exploration.<br><br>**AWS Lake Formation:** Set up and enforce security, governance, and auditing policies. |
| **Security and Access Control** | **Encryption:**<br>TLS/SSL<br>HDFS Encryption<br><br>**Access Control:**<br>Apache Ranger<br>AWS IAM<br>Google Cloud Identity and Access Management (IAM) | **AWS IAM (Identity and Access Management):** Manage access to AWS resources.<br><br>**AWS Key Management Service (KMS):** Encrypt data at rest and in transit. |
| **Data Science Integration** | **Model Deployment:**<br>TensorFlow Serving<br>MLflow<br>PMML (Predictive Model Markup Language)<br><br>**Notebook Environments:**<br>Jupyter Notebooks<br>Google Colab<br>Databricks Notebooks | **Amazon SageMaker:** For building, training, and deploying machine learning models. |
| **Architectural Patterns** | **Lambda Architecture:** Combines batch and stream processing for real-time and batch processing. | **Serverless Architecture:** Leverage services like Lambda, Glue, and Step Functions for serverless processing. |

| | | |
|---|---|---|
| | **Kappa Architecture:** Simplifies the Lambda Architecture using only stream processing. | **Data Lake Architecture:** Utilize S3 as a central data lake to store structured and unstructured data. |
| **Data Versioning and Lineage** | **Version Control:**<br>Git<br>DVC (Data Version Control)<br><br>**Lineage Tracking:**<br>Apache Atlas<br>DataHub | |
| **Cloud Integration** | **Cloud Platforms:**<br>AWS, Azure, Google Cloud Platform (GCP)<br><br>**Serverless Computing:**<br>AWS Lambda<br>Azure Functions<br>Google Cloud Functions | **AWS Direct Connect or VPN:** Connect on-premises data centers to AWS.<br><br>**AWS SDKs and CLI:** Integrate and automate AWS services using SDKs and the Command Line Interface. |

| Segment | Microsoft Azure | Google Cloud Platform |
|---|---|---|
| **Data Ingestion** | **Azure Data Factory:** Orchestrate and automate data workflows. Support for data movement from various sources to data lakes or warehouses.<br><br>**Azure Event Hubs:** Ingest and process massive amounts of streaming data. | **Cloud Pub/Sub:**<br>Real-time messaging service for event-driven architectures.<br><br>**Cloud Storage:**<br>Object storage for batch uploads. |
| **Data Storage** | **Azure Data Lake Storage:** Scalable and secure data lake storage.<br><br>**Azure SQL Data Warehouse (now part of Azure Synapse Analytics):** Enterprise-grade analytics service.<br><br>**Azure Cosmos DB:** Globally distributed, multi-model database for operational and analytical workloads. | **BigQuery:** Fully-managed, serverless data warehouse for analytics.<br><br>**Cloud Storage:** Object storage for raw data and backups.<br><br>**CloudSQL:** Managed relational databases. |
| **Data Processing** | **Azure Databricks:** Apache Spark-based analytics platform for big data and machine learning.<br><br>**HDInsight:** Fully managed cloud service for big data analytics using Hadoop, Spark, HBase, and more.<br><br>**Azure Stream Analytics:** Real-time analytics on streaming data. | **Dataflow:** Fully managed stream and batch processing using Apache Beam.<br><br>**Dataprep by Trifacta:** Cloud-native data preparation service.<br><br>**Dataproc:** Managed Apache Spark and Hadoop service. |
| **Data Transformation** | **Azure Data Factory:** Transform and clean data using data flows and transformations.<br><br>**Azure HDInsight:** Leverage Apache Spark or Hive for data transformation. | **Dataflow:** Apache Beam for ETL pipelines.<br><br>**Cloud Dataprep:** Visual data preparation tool. |

| | | |
|---|---|---|
| **Analytics and Reporting** | **Power BI:** Business Intelligence and visualization.<br><br>**Azure Synapse Studio:** Integrated analytics and data exploration. | **BigQuery:** For ad-hoc queries and analytics.<br>**Looker, Tableau, or Data Studio:** Business intelligence and visualization tools. |
| **Data Orchestration** | **Azure Data Factory:** Schedule and orchestrate data workflows.<br><br>**Azure Logic Apps:** Automate workflows and integrate services, including data services. | **Cloud Composer:** Managed Apache Airflow for workflow orchestration.<br>**Cloud Scheduler:** Fully managed cron job scheduler. |
| **Monitoring and Logging** | **Azure Monitor:** Monitor the performance and health of resources.<br><br>**Azure Log Analytics:** Collect and analyze log data. | **Cloud Monitoring:** Infrastructure and application monitoring.<br><br>**Cloud Logging:** Centralized log management. |
| **Data Data Quality and Governance** | **Azure Purview:** Unified data governance service for discovering, understanding, and managing data.<br><br>**Azure Data Catalog:** Discover, register, and manage data asset. | **Cloud Data Catalog:** Fully managed and scalable metadata management service.<br><br>**Cloud Data Loss Prevention (DLP):** Sensitive data discovery and redaction. |
| **Security and Access Control** | **Azure Active Directory (AAD):** Identity and access management.<br><br>**Azure Key Vault:** Securely store and manage sensitive information like keys and secrets. | **Cloud Identity and Access Management (IAM):** Access control for GCP resources.<br>**Cloud Key Management Service (KMS):** Manage cryptographic keys. |
| **Data Science Integration** | **Azure Machine Learning:** End-to-end platform for building, training, and deploying machine learning models. | **AI Platform:** Managed services for building, training, and deploying machine learning models.<br><br>**Notebooks:** AI Platform Notebooks or Jupyter Notebooks on AI Platform. |
| **Architectural Patterns** | **Modern Data Warehouse (Azure Synapse Analytics):** Combines big data and data warehousing for analytics.<br><br>**Event-Driven Architectures:** Use Azure Event Hubs and Azure Functions for event-driven processing. | **Serverless Architecture:** Utilize serverless services like Cloud Functions.<br>**Data Lake and Data Warehouse:** Combine Cloud Storage and BigQuery for cost-effective storage and analytics. |
| **Data Versioning and Lineage** | | **Cloud Data Catalog:** Track and manage data lineage.<br><br>**BigQuery:** Keep track of changes with versioned tables. |
| **Cloud Integration** | **Azure Functions:** Serverless computing for event-driven solutions.<br><br>**Azure Logic Apps:** Connect and automate workflows across cloud and on-premises services. | **Cloud Functions:** Serverless computing for event-driven functions.<br>**Cloud Run:** Fully managed compute platform for containerized applications. |