# ELTMaestro

Zealot Gurung
MAESTRO ANALYTICS, BOSTON, MA

# Contents

# ELTMaestro Overview

ELTMaestro is an enterprise application that supports ELT (Extract Load Transform), ETL (Extract Transform Load), ML (Machine Learning), CDC (Change Data Capture), data quality, data integration, and data lineage. ELTMaestro works with a large and growing number of target data warehouse platforms, including Redshift, Snowflake, Yellowbrick, Exasol, Greenplum, Azure Synapse, Spark/Hadoop and Databricks.

ELTMaestro has a graphical interface that allows users to create jobs as dataflow diagrams, as is customary with traditional ETL tools such as Informatica or DataStage. ELTMaestro extracts data from standard sources such as databases and CSV files as well as more specialized sources such as cellphone apps, Salesforce, mainframes, and many others. ELTMaestro has a built-in scheduler and extensive user-customizable data quality reporting. ELTMaestro includes advanced log-based change data capture and other CDC protocols. ELTMaestro also includes a comprehensive machine learning system and supports integration of machine learning and ETL processes.

ELTMaestro's subscription-based pricing model is very advantageous for rapidly growing and evolving production environments. ELTMaestro software licenses do not restrict data volume, data sources, bandwidth, memory, processing, or number of users. ELTMaestro subscription fees do not change as you grow your data warehousing platform.
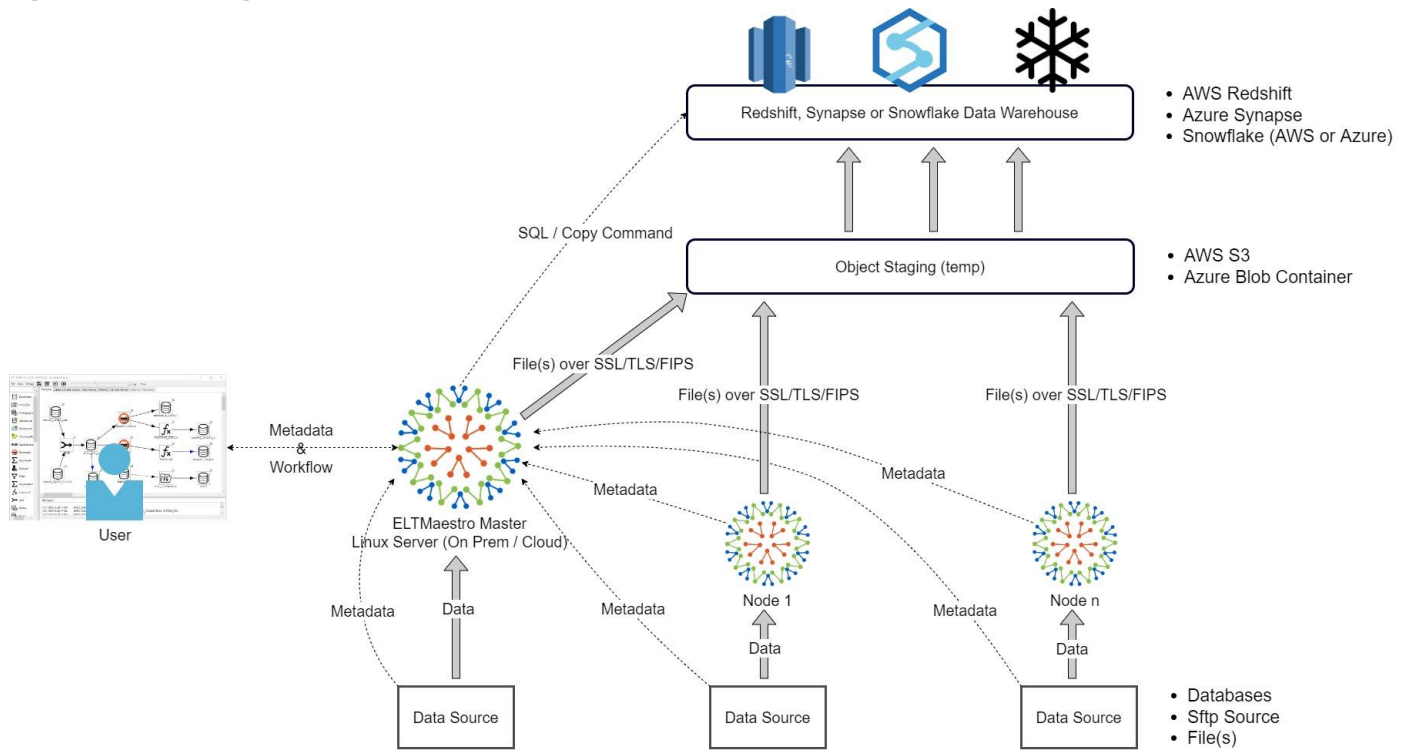
# Edge Processing Architecture



*Figure 1 ELTMaestro Edge Processing Architecture*

Edge processing architecture eliminates ETL server bottlenecks when source data is geo-distributed.
- Users design data load workflows using the ELTMaestro graphical client. Metadata is provided to user from the master server.
- The ELTMaestro master server then uses schedule and deployment information to decide if the workflows should be executed locally or pushed to remote nodes. In either case, data is extracted from source, compressed, and securely transferred to object storage (such as S3 or Blob). ELTMaestro then issues the appropriate SQL command (e.g., copy, load, or external table) to load files directly to target platform.

- Edge nodes allow customers to optimize data load speed when their data is distributed over multiple remote geographical locations, obviating the need to accumulate geo-distributed data on a single ETL server. This eliminates bottlenecks and saves intranet bandwidth. Moving large datasets becomes faster, more efficient, and less expensive.

## Machine learning

ELTMaestro includes Machine Learning as a standard feature. To build an ML model you simply point to a dataset, select and configure a model, and execute.

An example is shown below.

*Figure: ELTMaestro ML Training Workflow*



*Figure 2 ELTMaestro ML Training Workflow*

When the workflow is deployed and executed, ELTMaestro converts it into Spark ML code. Upon completion, ELTMaestro displays the model's evaluations and accuracy.



*Figure 3 ML Workflow logs and accuracy*

Similar pipelines can be built for predictions:



*Figure 4 ELTMaestro ML Prediction Workflow*

Prediction results can be graphically displayed on dashboards:



*Figure 5 Predictions*

# ELTMaestro produces extensive workflow metadata.

ELTMaestro transforms data by converting workflows to SQL that runs against the target data warehouse. ELTMaestro keeps track of record counts, data lineage, and runtime quality metrics.

For example, the workflow below, designed by the graphical interface, is converted into SQL that executes on the data warehouse platform.



*Figure 6 Workflow*



*Figure 7 Workflow log*

ELTMaestro maintains runtime history, including the SQL that was executed by each step:



*Figure 8 Runtime history*

ELTMaestro maintains data quality information of all nodes:



| BATCH_CYCLE_NM | s_JDBCTARGETSYNAPSE | s_ONSTAGEGROUP | s_SCD2 | s_SYNAPSEPROFILELOADER | s_TABLE |
|---|---|---|---|---|---|
| SYNAPSE_DEMO_SCD | 0 | 0 | 53049 | 0 | 0 |
| SYNAPSE_DEMO_LARGE_TABLE | 8926 | 0 | 0 | 0 | 0 |
| SYNAPSE_DEMO_PIVOT_2 | 0 | 0 | 0 | 0 | 5463 |
| SYNAPSE_DEMO_SALESFORCE | 1180761 | 0 | 0 | 0 | 0 |
| GSUTIL2 | 0 | 0 | 0 | 0 | 2480 |
| ISP_MASTER_PREDICT_EXPORT | 0 | 204132 | 0 | 0 | 8107041 |
| SYNAPSE_DEMO_PIVOT | 0 | 0 | 0 | 0 | 4 |
| SYNAPSE_DEMO_ONSTAGE | 0 | 0 | 0 | 9130 | 0 |
| RS_FILE_LOADER | 0 | 0 | 0 | 0 | 41715 |
| SNOW_ELT | 0 | 0 | 0 | 0 | 4652 |
| SYNAPSE_DEMO_METRICS | 2602 | 0 | 0 | 0 | 0 |
| SYNAPSE_ONSTAGE_GROUP | 0 | 0 | 0 | 2688 | 0 |
| DEMO_SYNAPSE_SCD2 | 0 | 0 | 62755 | 0 | 0 |

*Figure 9 Automatic data quality metrics by row count*

ELTMaestro automatically maintains data lineage information:

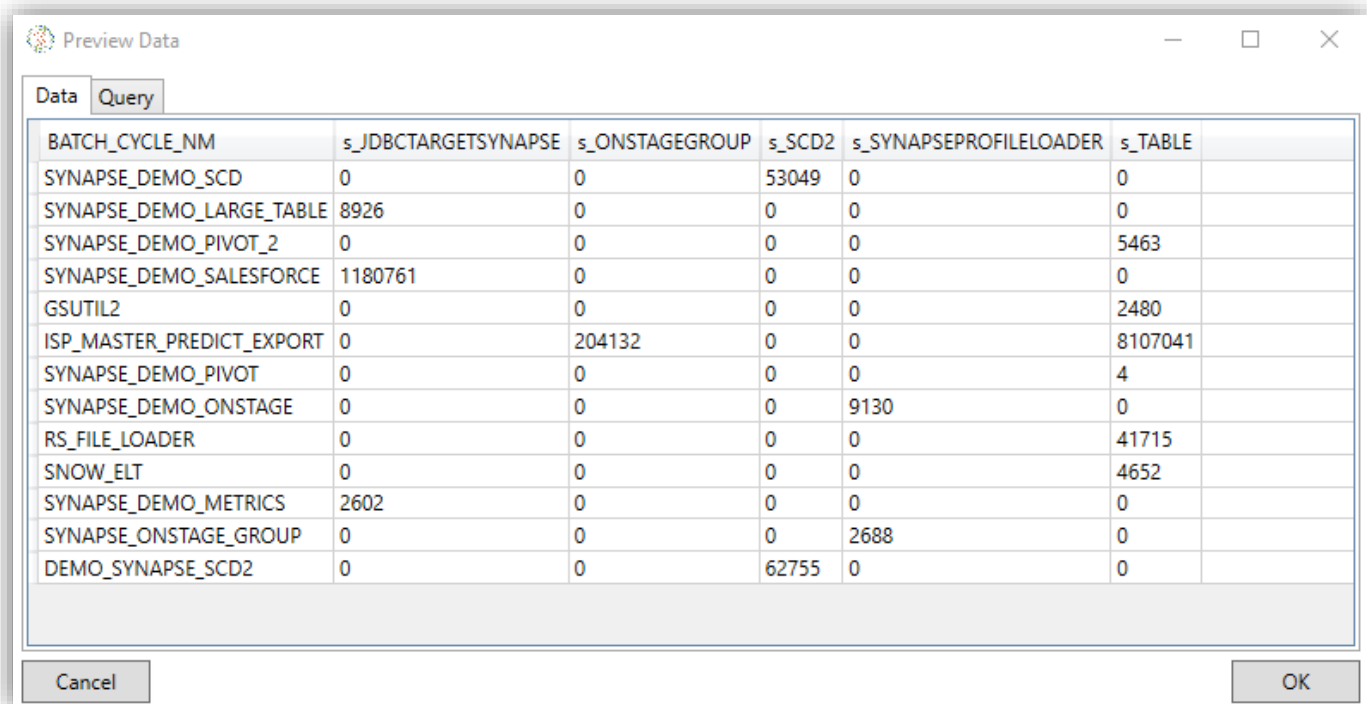| data_context_name | data_context_name_origin | data_context_detail |
|---|---|---|
| /ingest/failure_types | ONSTAGE:HANA:HANA.INTEGRATOR.FAILURE_TYPES | JDBCTARGETHDFS:/ingest/failure_types |
| /ingest/machine_metrics | ONSTAGE:REDSHIFT:dev.public.machine_metrics | JDBCTARGETHDFS:/ingest/machine_metrics |
| dev.demo.dummy_98 | ONSTAGE:ORACLE_CONNECTION:ORCLCDB.SYSTEM.DUMMY_98 | JDBCTARGETREDSHIFT:dev.demo.dummy_98 |
| dev.ot.dummy_98 | ONSTAGE:ORACLE_CONNECTION:ORCLCDB.SYSTEM.DUMMY_98 | JDBCTARGETREDSHIFT:dev.ot.dummy_98 |
| DWH_ETL.PUBLIC.PERSON22 | ONSTAGE:SQL SERVER:AdventureWorks2019.Person.Person | JDBCTARGETSNOWFLAKE:DWH_ETL.PUBLIC.PEF |
| DWH_ETL.PUBLIC.PASSWORD2 | ONSTAGE:SQL SERVER:AdventureWorks2019.Person.Password | JDBCTARGETSNOWFLAKE:DWH_ETL.PUBLIC.PAS |
| SYNAPSE.integrator.TMP_11421_1233_4_4 | INSITU:SYNAPSE:"integrator"."batch_cycle_run" | JOIN:SYNAPSE.integrator.TMP_11421_1233_4_49 |
| iqasnyapsepool.integrator.STEP_STATISTI( | INSITU:SYNAPSE:SYNAPSE.integrator.TMP_11421_1233_11_61 | TABLE:Existing Table [Create=True] |
| iqasnyapsepool.integrator.STEP_STATISTI( | INSITU:SYNAPSE:null.null | TABLE:Existing Table [Create=False] |
| gpadmin.dev.batch_cycle_run | ONSTAGE:ABC:sqlmaestro.public.batch_cycle_run | JDBCTARGETGREENPLUM:gpadmin.dev.batch_c |
| gpadmin.public.batch_cycle | ONSTAGE:ABC:sqlmaestro.public.batch_cycle | ONSTAGEGROUP:gpadmin.public.batch_cycle |
| gpadmin.public.batch_cycle_run | ONSTAGE:ABC:sqlmaestro.public.batch_cycle_run | ONSTAGEGROUP:gpadmin.public.batch_cycle_r |
| gpadmin.isp_integration.cust_predicted_d | INSITU::gpadmin.integrator.TMP_10558_1183_1_21 | TABLE:Existing Table [Create=False] |
| gpadmin.integrator.TMP_10637_1178_4_6' | INSITU::"isp_customer"."InternetPackage" | JOIN:gpadmin.integrator.TMP_10637_1178_4_67 |

*Figure 10 Data lineage*

# ELTMaestro Standard Features

| Feature Type | Transformation | Orchestration Mode | Misc |
|---|---|---|---|
| Data Loading | None | Bulk Load Delta Load Watermark Load Parallel Load | Source: Database, Salesforce, Object Storage, SFTP, Files |
| Data Loading using CDC | SQL Replay | 100% SQL | Log Mining |
| In-Situ Transformation | Table/DataFrame, Join, Deduplicate, Union, Minus, Filter, Data Masking,Pivot,Window/Scalar/Aggregate Functions, Slowly Changing Dimensions, +more | 100% SQL | Generated on runtime |
| Machine Learning | Spark | 100% Spark Code | Regression, Classification and Clustering |
| Data Quality/Lineage | None | Automatically Maintained | |
| Data Un-Loading | Export to file(s), object storage, databases | | |
| Controls | Custom Metrics, File/Database/Script Watcher, Variables, SQL Script, Email Alerts, Smart Script, SSH Script | Depends on control feature | |
| Automatic Recovery | | Job/Step/Node auto-restart on failures and threshold configuration | |
| Scheduling | | | Pre-defined and custom scheduling |

# ELTMaestro ML Features

| Feature Engineering / Transformations | Regression |
|---|---|
| CountVectorizerTransform | DecisionTreeRegressor |
| DCTTransform | FactorizationMachinesRegressor |
| ElementwiseProductTransform | GeneralizedLinearRegressor |
| HashingTransform | GradientBoostedTreeRegressor |
| IDFTransform | LinearRegressor |
| MaxAbsScalerTransform | RandomForestRegressor |
| MinMaxScalerTransform | |
| NGramTransform | **Classification** |
| NormalizerTransform | DecisionTreeClassifier |
| PCATransform | FactorizationMachinesClassifier |
| PolynomialExpansionTransform | GradientBoostedTreeClassifier |
| RegexTokenizerTransform | LinearSupportVectorMachineClassifier |
| RobustScalerTransform | LogisticRegression |
| StandardScalerTransform | LogisticRegressionOneVsRestClassifier |
| TokenizerTransform | MultilayerPerceptronClassifier |
| VectorIndexerTransform | NaiveBayesClassifier |
| Word2VecTransform | RandomForestClassifier |
| BinarizerTransform | |
| BucketizerTransform | **Clustering / Neural Network** |
| FeatureHasherTransform | BisectingKMeansCluster |
| ImputerTransform | GaussianMixtureModelCluster |
| InteractionTransform | KMeansCluster |
| OneHotEncoderTransform | LDACluster |
| QuantileDiscretizerTransform | |
| StopWordsRemoverTransform | **Documentation Source** |
| StringIndexerTransform | https://spark.apache.org/docs/latest/ml-classification-regression.html |
| VectorAssemblerTransform | https://spark.apache.org/docs/latest/mllib-clustering.html |

# Questions?

Contact us at pen@maestro-analytics.com