

White Paper

Migrating SQL Data to Databricks

A Foundation for Digital Transformation

Abstract

This whitepaper details a rapid method for migrating your SQL data to the Databricks SQL platform, illustrating the process through a clinical dataset while framing the benefits offuture-proofing the data architecture

Kurt Rosenfeld krosenfeld@ctidata.com



Migrating SQL Data to Databricks

A Foundation for Digital Transformation

Table of Contents

1.0	Introduction	3
2.0	Database Migration Concerns	5
3.0	Data Product Trace Maps	7
4.0	Migrating Logic	9
5.0	Migrating Data	12
6.0	Dissecting a Data Example	18
6.1	Environment Setup	18
6.2 \$	Secrets and Run-Time Variables Setup	20
7.0	Conclusion	29
About	CTI Data	31



Migrating SQL Data to Databricks

A Foundation for Digital Transformation

1.0 Introduction

Shifting SQL workloads to a unified Data Infrastructure platform is increasingly important. Here's why: Historically, IT automation depended solely on structured data. Today, generative AI (GenAI) can surface information in the hordes of unstructured data, powering new automations.

Consequently, our systems must marry the unstructured "analog" data of the physical world with its structured "digital" twin. This requires a new perspective on IT data infrastructure. Existing



infrastructure evolved fragmented: structured data residing in a row or columnar databases, object data residing as "blobs" in "containers" or "buckets," and file data residing in file systems — each with distinct metadata and access controls.

Such fragmentation is now a bottleneck; we need an infrastructure that treats structured and unstructured data as an integrated whole, along with a single metadata and relationship catalog, unified governance and control,

within a uniform access platform.

The technology vendors have shoehorned their solutions towards this goal, adapting their SQL/NoSQL databases for object storage or extending their data lakes for SQL, with further adaptations for a new data class, the high-dimension vectors crucial to GenAI.

It's a point-in-time mess, leaving the data engineer and end-user to work around the constraints of these platform legacies.

But it's improving, and the vendors most relevant are the cloud data titans: Google, Microsoft, Amazon, Snowflake, and Databricks, with Salesforce and Oracle vying for their place.

The desired solution is as straightforward as depicted on the right.

The table on the next page calls out how well-positioned Databricks is for



meeting this goal with the arrival of their serverless SQL capability in June 2022, deeply integrated with their open, multi-modal data platform heritage and processing tools.



Databricks' Position						
Any compute model	→ Serverless (SaaS) or managed (PaaS)					
Any processing model \rightarrow Python, R, Scala, Java, SQL						
Rapid delivery	→ Built-in DevOps, DataOps, MLOps					
Process all types of data	\rightarrow Data Lake with choices of structured abstractions					
One data framework	→ From ingestion to micro-streaming					
Single ecosystem, massive scale	\rightarrow Reduce technology silos, one platform for all needs					
Cloud vendor abstraction	→ Mix and match any Cloud					
Seamless data federation	→ External data is an equal citizen					
Safe data sharing	\rightarrow Internal and external teams have governed access					
Uniform governance & control	\rightarrow All data, all code, all security, all deployments					
No vendor lock-in	→ 100% open technology					

Consequently, running your SQL workloads on Databricks means you are futureproofing the data's value, regardless of how the data may need to evolve and integrate. Furthermore, even for one-off throw-away SQL workloads, Databricks is positioned as a perfectly competent platform.

To explore what is involved in migrating to Databricks, we need a legacy platform to make the details meaningful, so in this case, we assume Snowflake is the legacy.

These two platforms, with very different origins, now frequently compete as their technologies expanded and began to overlap. The chart below helps frame the current state but keep in mind we are singularly focused on the SQL aspects in this document.



<u>Snowflake</u> zero administration SQL service addressing analytic use cases	Databricks open-source computing and data service addressing data science use cases
Serverless SQL platform	Ecosystem tying together compute provisioning, structured and unstructured data management, and ML management Capable of serverless SQL operation (transparent dynamic provisioning of resources)
Supports non-SQL logic as a separate capability callable from SQL	Supports SQL processing as an equal citizen to its other computing capabilities
Data pipelining and orchestration require separate tooling	Integrates batch and real-time data pipelining and orchestration
Federates structured and semi-structured data	Federates any data
Governance and catalog addresses data only	Governed by a unified data and process manager with AI-assisted glossary generation
Plugins for popular BI platforms	AI-assisted exploratory data analysis (EDA) tooling with plugins for popular BI platforms

2.0 Database Migration Concerns

While SQL database platforms are conceptually similar – reflecting the power of SQL as a standard – all the platforms have feature sets with real differences, which, from a migration point-of-view, raise these migration concerns:



	Concern
Data	Adapting data types and converting formats
Logic	Adapting logic residing in the external programs (e.g. SQL statements) or in the data platform (e.g. Stored Procedures) to ensure functional equivalence
Process	Reproducing the program logic that controls data processing cycles
Security	Reproducing controls that prevent improper data access and creation
Operations	Reproducing infrastructure-as-code and other platform automation such as identity management, metering, and chargebacks
Meta-data	Transferring data descriptions and business rule descriptions

Technical Element	Migration Considerations	Migration Strategies		
Data and Tables				
Data types	Similar, special handling for variants and geospatial	Review DDL, adapt dependent SQL		
Views, materialized views	Same features and methods	Apply SQL dialect conversion automation		
General Logic				
SQL dialect	Similar, ANSI SQL centric with minor differences	Apply SQL dialect conversion automation		
Ingestion staging	Same features, similar methods	Repoint ingestion tool, or replace with Databricks Autoloader		
Specialized Logic				
User Defined Functions	Same features with syntax & language differences	Convert syntaxes, maintain function signature		
Stored Procedures and Triggers	Needs alternative solution	Replace with UDFs or Delta Live Tables		

Applying these concerns specific to migrating Snowflake to Databricks, we note these strategies:



SQL refactoring	Identify high-cost queries (using query profiler)	Enhance with Databricks hashing/join features	
Process Control			
CLI scripts	CLI's different, with little equivalence	Assess CLI script logic and recreate using Python, Databricks Jobs, etc.	
Orchestration	Simple to convert Snowflake Tasks	Use richer features of Databricks Workflow (or repoint external tool)	
Security			
Data privacy	Data masking, PII detection, clean rooms, audit differences	Map privacy setup, leverage Unity Catalog for data provenance	
Identities, access control	Similar privilege/permission mgt, different row level security	Map and improve with Unity Catalog	
Operational Control			
Data clustering	Clustering concepts similar, different syntaxes, less automated	Apply additional clustering based on query cost profiles	
Resources, chargebacks	Different models, Databricks is more granular	Rethink to leverage Databricks tag- based granular capture	
Elastic controls	Different methods/latency for automated scaling up/down	Assess scaling latency and use Databricks' more granular control	

The top-of-mind topics are the first two, data and logic; we will dissect the first, the data topic, and touch on the logic topic, leaving a fuller examination for a later paper.

3.0 Data Product Trace Maps

This diagram depicts the typical elements of an analytic data environment: various upstream services execute embedded SQL logic, stored procedures, or user-defined functions to produce or change data, with various downstream services consuming and sometimes producing or changing data.

()CTI DATA



We need a trace map of the dependency pathways for all the logic and related data in scope for the migration. A good practice is to scope the work by data product.

One of the beauties, from the point of view of mapping these dependency paths, is that SQL logic is almost always stateless: all state is stored in the data tables, the exception is when a stored procedure, trigger or custom function affects a variable.

This means modern data lineage tools can generate the trace map, like the diagram below, by probing the program objects of each deployed technology, extracting the logic, and identifying

the discrete data elements that are touched, inferring the logic-to-data pathways.

But even so, a data engineer should inspect the results and fill in the gaps as necessary.

id	paul.roome@databricks.com	int		>	turbine_prod.gold.t	urbine_master ks.com	0	a9d354e0-3064-4db	l.turbine_features bc-884e-affee1d77ae1
last_ser	rviced_days doub	le .	\					ID	doubl
			λ		ID	double		AN3	doubl
			1/		AN3	double		AN4	doub
			1		AN4	double		AN5	doub
	turbine_prod.silver.turbine_metrics_wo		24		AN5	double		AN6	doub
III	rldwide	- 9			AN6	double		AN7	doub
	paul.roome@databricks.com		1		AN7	double		ANS	doub
					AN8	double		AN9	doub
AN3	doub	ole			AN9	double		AN10	doub
AN4	doub	ole			AN10	double		SPEED	doub
	Show 10 more columns				SPEED	double		TORQUE	doub
	Show to more columns			1	TORQUE	double		TIMESTAMP	timestam
				1	TIMESTAMP	timestamp		status	strin
				1	status	string	/	last_serviced_days	doub
_				`	last_serviced_days	double		feature1	doub
	turbine_prod.silver.turbine_status paul.roome@databricks.com	G			Hide colu	mns		feature2	doub
								Hide co	lumns



4.0 Migrating Logic

Logic migration is a distinctly different endeavor from data migration. The upstream and downstream technologies need two modifications to switch to the future platform: their data endpoint connections must be changed, which is simple enough, and their logic must be adapted to the differences in data types, data formats, logic dialects, and invocation methods.

Most SQL logic needs only minor changes. There is one exception: Stored Procedures and Triggers. Databricks SQL does not yet provide these features; if your data heavily depends on them, some "rethink" work will be necessary to recreate their functionality.

To get a feel for the SQL changes, consider the query below, which retrieves the top 5 pneumonia-related admissions for the past year with a mortality count. A typical question with a typical SQL expression.

SELECT	SELECT
concat(DR.DESCRIPTION,	concat(DR.DESCRIPTION,
' (', DR.DRG_CODE, ')')	' (', DR.DRG_CODE, ')')
AS "DRG Description",	AS `DRG Description,`
count(DISTINCT A.HADM_ID)	count(DISTINCT A.HADM_ID)
AS "Pneumonia Diagnosis",	AS `Pneumonia Diagnosis`,
count (count(
DISTINCT CASE	DISTINCT CASE
WHEN A.HOSPITAL_EXPIRE_FLAG = TRUE	WHEN A.HOSPITAL_EXPIRE_FLAG = TRUE
THEN A.HADM_ID	THEN A.HADM_ID
ELSE NULL	ELSE NULL
END	END
) AS "Mortality Count") AS `Mortality Count`
FROM ADMISSIONS AS A	FROM ADMISSIONS AS A
JOIN DIAGNOSES_ICD AS D	JOIN DIAGNOSES_ICD AS D
ON A.HADM_ID = D.HADM_ID	ON A.HADM_ID = D.HADM_ID
JOIN D_ICD_DIAGNOSES AS DI	JOIN D_ICD_DIAGNOSES AS DI
ON D.ICD_VER_CODE = DI.ICD_VER_CODE	ON D.ICD_VER_CODE = DI.ICD_VER_CODE
JOIN DRGCODES AS DR	JOIN DRGCODES AS DR
ON A.HADM_ID = DR.HADM_ID	ON A.HADM_ID = DR.HADM_ID
WHERE	WHERE
<pre>date_part(YEAR, A.ADMITTIME) = 2024</pre>	EXTRACT (YEAR FROM A.ADMITTIME) = 2024
AND DI.LONG_TITLE ILIKE '%pneumonia%'	AND DI.LONG_TITLE ILIKE '%pneumonia%'
GROUP BY	GROUP BY
DR.DESCRIPTION,	DR.DESCRIPTION,
DR.DRG_CODE	DR.DRG_CODE
ORDER BY	ORDER BY
"Pneumonia Diagnosis" DESC NULLS LAST	`Pneumonia Diagnosis` DESC NULLS LAST
LIMIT 5	LIMIT 5

What changed? The Snowflake date_part() function is replaced with the ANSI SQL extract() function, and double quotes for enclosing identifiers are replaced with back-ticks (e.g. `Mortality Count`).



We used translation technology to make these changes, eliminating the labor of syntactic correction. However, the advancement of Large Language Models is fast displacing this method.

Nowadays, you can direct your LLM to review the reference documentation (see list below) along with your source code to provide conversion analysis and "first-pass" converted code. Although not yet perfect, the pace of improvement is remarkable. Databricks SQL has fast-evolving enhancements, making the LLM output generated from the reference documentation invaluable compared to almost immediately dated guides and other alternatives.

Snowflake and Databricks SQL Reference Documentation:

https://docs.snowflake.com/en/sql-reference/sql-all

https://docs.snowflake.com/en/sql-reference/functions-all

https://docs.databricks.com/en/sql/language-manual/sql-ref-datatypes.html

https://docs.databricks.com/en/sql/language-manual/sql-ref-functions-builtin-alpha.html

Snowflake Function	Databricks SQL Equivalent	Notes
GET_DDL(object_name)	SHOW CREATE TABLE table_name(for tables),query information_schema.TABLES or information_schema.COLUMNS for others	Databricks uses SHOW CREATE for tables. Use information_schema for programmatic metadata access.
SYSTEM\$GET_PREDECESSOR_R ETURNED_COLUMNS(table_na me)	Not Directly Available. Consider Delta Lake Change Data Feed.	Snowflake-specific for CDC. Databricks uses Delta Lake CDC.
SYSTEM\$CLUSTERING_INFORM ATION(table_name)	Analyze information_schema.COLUMNS statistics; for Delta Lake, DESCRIBE DETAIL table_name.	Infer clustering from stats. Delta Lake shows partitioning/Z-Ordering.
SYSTEM\$TASK_HISTORY()	Query system.task_history (if using Databricks Workflows).	Specific to Snowflake Tasks. Databricks Workflows has its own history table.
SYSTEM\$PIPE_STATUS('pipe _name')	Monitor Auto Loader stream status via Spark Structured Streaming APIs or Delta Live Tables UI/APIs.	Specific to Snowflake Pipes. Databricks uses Auto Loader/DLT, monitored differently.
CURRENT_REGION()	Not Directly Available as a SQL function. Part of workspace config.	Databricks region is a configuration setting.
GET_OBJECT_S3('s3://')	<pre>spark.read.format('').load('s3:/ /') (replace '' with format).</pre>	Databricks uses Spark's data source API.
STAGE_FILES() (external stages)	Use dbutils.fs.ls('s3://') or cloud provider SDKs.	Databricks interacts with cloud storage directly.
CONVERT_TIMEZONE(target_ timezone,	<pre>from_utc_timestamp(to_utc_timestamp (timestamp, source_timezone), target_timezone)</pre>	Convert to UTC first if source has timezone.



Snowflake Function	Databricks SQL Equivalent	Notes
<pre>source_timezone, timestamp)</pre>		
DATE_PART(date_or_time_p art, expression)	<pre>extract(date_or_time_part FROM expression) or year(expression), month(expression), etc.</pre>	Databricks offers both general and specific date/time functions.
DATE_TRUNC(date_or_time_ part, timestamp)	<pre>date_trunc(date_or_time_part, timestamp)</pre>	Function name is the same.
MASK(input_string,)	Requires custom UDF or use Unity Catalog's data masking.	Snowflake's MASK is built-in; Databricks needs UDF or Unity Catalog.
STRTOK_TO_ARRAY(string, delimiter)	<pre>split(string, delimiter)</pre>	Different function name.
SPLIT_TO_TABLE(input, delimiter)	<pre>explode(split(input, delimiter))</pre>	Use explode after splitting.
RLIKE(subject, pattern)	subject rlike pattern	Uses the rlike operator.
LIKEANY(string, pattern1, pattern2,)	string LIKE pattern1 OR string LIKE pattern2 OR	Requires explicit OR conditions.
SOUNDEX(string)	Requires custom UDF or external Spark library.	Snowflake built-in; Databricks needs UDF/external library.
LEVENSHTEIN(string1, string2)	Requires custom UDF or external Spark library.	Similar to SOUNDEX.
LISTAGG(expression [, delimiter]) [WITHIN GROUP (ORDER BY)] [ON OVERFLOW TRUNCATE 'string' [WITHOUT COUNT WITH COUNT]]	array_join(collect_list(expression) WITHIN GROUP (ORDER BY ordering_expression), delimiter) OVER (PARTITION BY grouping_expression)	WITHIN GROUP moved inside collect_list. Manual overflow handling.
BITAND(expr1, expr2)	expr1 & expr2	
BITOR(expr1, expr2)	expr1 expr2	
BITXOR(expr1, expr2)	expr1 ^ expr2	
BITNOT(expr1)	~ expr1	
IFF(condition, true_value, false_value)	CASE WHEN condition THEN true_value ELSE false_value END	Standard SQL CASE WHEN.
NVL(expr1, expr2)	<pre>coalesce(expr1, expr2)</pre>	Standard SQL coalesce.
ZEROIFNULL(numeric_expr)	CASE WHEN numeric_expr IS NULL THEN 0 ELSE numeric expr END	Achieved using CASE WHEN.



Snowflake Function	Databricks SQL Equivalent	Notes
TRY_CAST(source_value AS data_type)	TRY_CAST(source_value AS data_type) (newer),orCASE WHEN TRY_CAST() IS NOT NULL THEN TRY_CAST() END	TRY_CAST available in recent runtimes.
PARSE_JSON(string)	<pre>from_json(string, schema_of_json(string))</pre>	Requires schema; schema_of_json can infer.
GET_PATH(variant, path)	<pre>get_json_object(variant, path) or variant.path.to.element</pre>	Both function and dot notation available.
FLATTEN(input => array_or_object)	explode(input)	Uses the explode function.
XMLGET(xml, tag)	Requires Spark XML libraries (e.g., xml_xpath after parsing).	Databricks uses Spark for XML processing.
POLICY_CONTEXT()	Not Directly Available. Use Unity Catalog features.	Databricks uses Unity Catalog for security policies.
ARRAY_CONSTRUCT()	array()	Databricks array() function.

5.0 Migrating Data

Historical data is unique in that it reflects business rules from the past, manifest in the data logic and the data inputs at the time, which may no longer be available. It is almost always preferable to migrate historical data rather than regenerate it by reconstituting the historical logic with its historical inputs.

The migration process has one goal: reproduce the historical data adjusted for data type and format differences keeping the data rules intact. While there may be pressure to improve the data – be it to correct, enrich, or mask certain aspects – which is to be welcomed since these demands are a sponsorship opportunity – such improvements should be implemented as a separate step once the data is migrated, thereby isolating change and retaining proof of faithful data reproduction.

Once again, LLM advancements make understanding the differences a breeze; in this case, we engage the Databricks Assistant, which employs Databrick's own LLM service, to report ANSI SQL peculiars.



① 🕸 Assistant +	Binary String Types
Here is a detailed mapping of ANSI SQL data types to Databricks SQL data types:	BINARY: Maps to BINARY VARBINARY: Maps to BINARY
Numeric Types	Date and Time Types
 INTEGER (INT): Maps to INT SMALLINT: Maps to SMALLINT BIGINT: Maps to BIGINT DECIMAL (NUMERIC): Maps to DECIMAL FLOAT: Maps to FLOAT REAL: Maps to FLOAT DOUBLE PRECISION: Maps to DOUBLE 	DATE: Maps to DATE TIME: Maps to TIMESTAMP TIMESTAMP: Maps to TIMESTAMP INTERVAL: Maps to INTERVAL Boolean Type BOOLEAN: Maps to BOOLEAN
Character String Types CHARACTER (CHAR): Maps to STRING CHARACTER VARYING (VARCHAR): Maps to STRING NATIONAL CHARACTER (NCHAR): Maps to STRING NATIONAL CHARACTER VARYING (NVARCHAR): Maps to STRING	Other Types ARRAY: Maps to ARRAY MULTISET: Not directly supported, can be represented using ARRAY XML: Not directly supported, can be represented using STRING

One main difference is that all character fields map to the STRING type; therefore, the character limit enforcement of VARCHAR(n) is lost and must be handled in logic if required.

[For the current Databricks data types, see <u>https://docs.databricks.com/en/sql/language-manual/sql-ref-datatypes.html</u>]

We observe below Databricks automatically handling this ANSI SQL behavior where VARCHAR(10) and REAL are converted to the Databricks SQL equivalent.

序 Nev	w Query 2024-12-1	18 1:43pm × +				
New	/ :☆	✓ <u>L</u> ● MIMIC Sn	nowflake 2XS 🗸 🛛	chedule Share	Save*	•
► Run	all (1000) 👻 🗸	Just now (9s) 🗍 hive_metas	store . 🖯 default 🗸		♦	Ģ
1 2 3 4	DROP TABLE I CREATE TABLE INSERT INTO SELECT value	F EXISTS maptest; maptest (value1 VARG maptest VALUES ('helld 1, typeof(value1), va	CHAR(10), value2 R o', 22.093); lue2, typeof(value)	EAL); 2) FROM maptest;		ۍ ۲
	parameter					ተ
← Re	esults 4 of 4 \rightarrow	Table 🗸 🕂		Q 7		×
	^{AB} _C value1	^{AB} c typeof(value1)	1.2 value2	A ^B _C typeof(value2)		
	hello	string	22.0930004119873	float		



Since we will be migrating data from Snowflake, we must examine how Snowflake's non-ANSI data types will be handled.



The Databricks Assistant is a little too high level; Snowflake also supports FLOAT4 and FLOAT8. Databricks SQL errors out when presented with these, which is reasonable since they are not ANSI standard, and the Databricks Assistant recommends using FLOAT as the proper alternative.





A SQL programmer might address data conversion by developing go-between migration logic: reading the source data via the Snowflake SQL client, casting the data to ANSI standard, and then writing the result to Databricks using the Databricks SQL client.

However, in the case of Snowflake, we can use Snowflake's internal conversion engine by way of Parquet, the highly efficient and flexible open data store format supported by the major data platforms.

When a Snowflake table is exported to Parquet, its native data types are converted to Parquet's. When imported to Databricks, the reverse happens, and since Databricks is heavily aligned with Parquet, little change is necessary. Essentially, Snowflake is exporting its tables into structures that align well with Databricks.

We can transfer tables with hundreds to billions of rows, bypassing the need for complex and computationally expensive conversion pipelines. Of course, some differences must be inspected afterward, such as VARCHAR(n).





The Databricks Assistant provides the following elaboration regarding Parquet data types:

① 徽 Parquet is a colu	Assistant	+ Logical Types Logical types are used to extend the primitive types with additional semantics:
various data typ Primitive Ty • BOOLEAN: • INT32: Rep! • INT64: Rep! • INT96: Rep! timestamps; • FLOAT: Rep • DOUBLE: Re • BYTE_ARRA • FIXED_LEN sequence of	es. Here is a list of Parquet data ty rpes Represents a boolean value (true resents a 32-bit signed integer. resents a 64-bit signed integer (us). resents a 32-bit floating-point nu resents a 32-bit floating-point nu resents a 64-bit floating-point nu resents a 64-bit floating-point r W : Represents a sequence of byte BYTE_ARRAY : Represents a fixed f bytes.	 additional semantics: STRING: Represents UTF-8 encoded strings (mapped to BYTE_ARRAY). MAP: Represents a map of key-value pairs. LIST: Represents a list of elements. ENUM: Represents an enumeration of values. DECIMAL: Represents a decimal value with precision and scale (can be mapped to INT32, INT64, or FIXED_LEN_BYTE_ARRAY). DATE: Represents a date (mapped to INT32). TIME: Represents a time of day (mapped to INT32 or INT64). TIMESTAMP: Represents a time interval (mapped to INT64 or INT64). INTERVAL: Represents a time interval (mapped to FIXED_LEN_BYTE_ARRAY). UUID: Represents a universally unique identifier (mapped to FIXED_LEN_BYTE_ARRAY). UUID: Represents a universally unique identifier (mapped to FIXED_LEN_BYTE_ARRAY). These data types allow Parquet to efficiently store and process complex data structures while maintaining compatibility with various data processing systems.

Meta-data stored in Parquet elaborates the primitive types, ensuring proper interpretation as the appropriate ANSI SQL data type. Here is an example of this process after transferring a patient ADMISSIONS table from Snowflake to Databricks. The screenshots on the next page are from the native SQL workbooks of both products.



	– Snowflak	ie —		- Databricks -	
1 2	USE patient.silver; DESCRIBE TABLE admi	ssions	1 2	USE patient_import.si DESCRIBE TABLE admiss	lver; ions
4	Results ~ Chart		← R	Results 2 of 2 \rightarrow Table	× +
	name	type		A ^B _C col_name	^{д8} c data_type
1	SUBJECT_ID	NUMBER(38,0)	1	SUBJECT_ID	decimal(38,0)
2	HADM_JD	NUMBER(38,0)	2	HADM_ID	decimal(38,0)
3	ADMITTIME	TIMESTAMP_NTZ(9)	3	ADMITTIME	timestamp
4	DISCHTIME	TIMESTAMP_NTZ(9)	4	DISCHTIME	timestamp
5	DEATHTIME	TIMESTAMP_NTZ(9)	5	DEATHTIME	timestamp
6	ADMISSION_TYPE	VARCHAR(255)	6	ADMISSION_TYPE	string
7	ADMIT_PROVIDER_ID	VARCHAR(10)	7	ADMIT_PROVIDER_ID	string
В	ADMISSION_LOCATION	VARCHAR(255)	80	ADMISSION_LOCATION	string
9	DISCHARGE_LOCATION	VARCHAR(255)	9.0	DISCHARGE_LOCATION	string
10	INSURANCE	VARCHAR(255)	10	INSURANCE	string
11	LANGUAGE	VARCHAR(255)	11	LANGUAGE	string
12	MARITAL_STATUS	VARCHAR(255)	12	MARITAL_STATUS	string
13	RACE	VARCHAR(255)	13	RACE	string
14:	EDREGTIME	TIMESTAMP_NTZ(9)	14	EDREGTIME	timestamp
15	EDOUTTIME	TIMESTAMP_NTZ(9)	15	EDOUTTIME	timestamp
16	HOSPITAL_EXPIRE_FLAG	BOOLEAN	16	HOSPITAL_EXPIRE_FLAG	boolean
17	ADMDATE	TIMESTAMP_NTZ(9)	17	ADMDATE	timestamp
18	ADMITTIME_DAYOFMONTH	NUMBER(38,0)	18	ADMITTIME_DAYOFMONTH	decimal(38,0)
19	DISCHTIME_DAYOFMONTH	NUMBER(38,0)	19	DISCHTIME_DAYOFMONTH	decimal(38,0)
20	ADM_DISCH_INSAMEYEAR	NUMBER(38,0)	20	ADM_DISCH_INSAMEYEAR	decimal(38,0)

A cardinality count provides some assurance the data is transferred correctly:



	- Sno	owflake –			-I	Databricks –	
1 2 3 4 5 6	USE patient.silve select count(*) as " count(distinc count(distinc from ADMISSIONS	r: Records", t HADM_ID) as "Admis t SUBJECT_ID) as "Pa	sions". tients"	1 2 3 4 5 6	USE patient_im select count(*) a count(dist count(dist from ADMISSION	<pre>port.silver; s `Records`, inct HADM_ID) as `/ inct SUBJECT_ID) as S</pre>	Ndmissions`, s`Patients`
G Res	ults ~ Chart			← F	Results 2 of 2 \rightarrow	Table 👻 🕂	
	Records	Admissions	Patients		123 Records	123 Admissions	123 Patients
1841	431231	431231	180733	Ĩ	431231	431231	180733

6.0 Dissecting a Data Example

Our work in complex data solutions provides access to diverse data environments, one being an anonymized dataset of 4.75 million clinical diagnoses that we used to develop disease forecasting algorithms stored in Snowflake. In this walkthrough, we shall migrate the entire dataset to Databricks, reproducing the schema and stored values.

We export the Snowflake tables to Parquet and subsequently import them into Databricks using a Microsoft Azure storage container to hold the Parquet intermediary objects. We establish the container in the same region as our Snowflake and Databricks tenants for speed and reduced egress charges. As stated earlier, this is faster and cheaper than using a live database-to-database SQL connection to read from the source and write to the target, benefiting from Parquet as the data type go-between.

6.1 Environment Setup

For our purposes, Databricks provides two IDEs for working with SQL code. The first is their Jupyter-based notebook; although primarily for Python, it can also process SQL using the "% SQL" magic declaration of Jupyter. Its objects are labeled "Notebook" within the Databricks workspace explorer.

The second is a SQL IDE that processes SQL expressions; its objects are labeled "Query" within the workspace explorer. It is similar in concept to a Snowflake worksheet; for convenience, we'll call it the "SQL worksheet" environment.

A SQL worksheet is attached to a serverless SQL configuration – called a "warehouse" – that defines the initial size of the compute that automatically responds to demand. In contrast, a Notebook is attached to a cluster configuration that identifies the compute plus runtime environment and must always be running to respond to demand.

There is an unusual exception: a Jupyter Notebook of SQL cells may execute on the serverless SQL warehouse; however, it will refuse if there are Python or other non-SQL cells.

Finally, an external IDE (such as Visual Studio) may also connect to a cluster or a SQL warehouse for performing development iterations.

Regardless of where execution occurs, whether on a cluster or serverless warehouse, all data is processed and stored within the Databricks platform and instantly accessible to both.



The following three screenshots demonstrate these options:



123 Patients

299712

123 Unique

299712



Because our data migration code needs to coordinate data movement from Snowflake to Azure to Databricks, we shall use the flexibility of a Notebook to step through this sequence, sometimes

switching between Python and SQL cells. We've provisioned a minimally sized cluster for this purpose.

kurtr@cptech.com's Cluster	>
Runtime	DBR 15.4 LTS ML • Spark 3.5.0 • Scala 2.12
Driver	Standard_DS3_v2 • 14 GB • 4 Cores

We also need an ADLS storage container with a shared access signature (SAS) token that the Notebook will use. It's more interesting to see the container populated with the exported tables in the screenshot below, but it is empty until Step 1 is completed later.

😑 Microsoft Azure 🔎 Se	earch resources, services, and docs (G+/)		🤣 Copilot	🔞			
Home > ctisnowbricksexchange Cont	tainers >						
Container	hange			×			
✓ Search × «	↑ Upload + Add Directory (🖰 Refresh 🕴 🤇 Ren	name 🗊 Delete 🛱	Change tier ····			
Overview	Authentication method: Access key	(Switch to Microsoft Ent	tra user account)				
Diagnose and solve problems	Location: kurtssnowbrickexchange						
🙊 Access Control (IAM)	Search blobs by prefix (case-sensitiv	Search blobs by prefix (case-sensitive)					
✓ Settings	Show deleted objects						
Shared access tokens ☆	Name	Modified	Access tier	Archive status			
🗞 Manage ACL		10/14/2024, 5:04:0	6				
Access policy	D_ICD_DIAGNOSES	10/22/2024, 3:28:1	9				
Properties	D_ICD_PROCEDURES	10/14/2024, 5:05:3	5				
🚺 Metadata	DIAGNOSES_ICD	10/28/2024, 3:20:4	i1				
		10/14/2024, 3:38:3	9				
	NOTE_DISCHARGE	10/14/2024, 5:12:2	.4				
	PATIENTS	10/14/2024, 5:00:4	3				
	PROCEDURES_ICD	10/14/2024, 5:06:4	11				
	TRANSFERS	10/14/2024, 5:09:3	9				

6.2 Secrets and Run-Time Variables Setup

It is always preferable to use a "secrets vault" rather than environment variables to store and retrieve sensitive credentials; in this case, we've used Databricks Secrets. We set up the secrets by opening the cluster terminal and using the CLI.



#databricks secrets put-secret dbmigrate ADLS_ACCOUNT --string-value "<ADLS ACCOUNT NAME>"
#databricks secrets put-secret dbmigrate ADLS_SAS_TOKEN --string-value "<ADLS SAS TOKEN>"
#databricks secrets put-secret dbmigrate SNOW_USER --string-value "<SNOWFLAKE USER ACCOUNT>"
#databricks secrets put-secret dbmigrate SNOW_PWD --string-value "<SNOWFLAKE USER PASSWORD>"
#databricks secrets put-secret dbmigrate SNOW_ACCOUNT --string-value "<SNOWFLAKE ACCOUNT>"

The Databricks Notebook is "secrets aware", displaying "redacted" for variables assigned from a Databricks Secret. We load the secrets into a Python dictionary



and use them to populate the connection parameters for Snowflake, Databricks, and ADLS.



We create a convenience function, snow_exec(), that wraps the Snowflake Python library to return a dictionary of column names with row data.





Now, we're ready to work on our first target: transferring 4.75m rows of diagnosis data to Databricks.

Step 1: Export the Source Table from Snowflake to the Cloud Container

The ADLS container is identified to Snowflake as a stage area:

```
# Create the snowflake stage area connected to the external Azure container
# This only needs to be run once
snow_exec(f"""
CREATE STAGE IF NOT EXISTS {snowStage}
URL = '{snowBlobURL}',
CREDENTIALS = ( AZURE_SAS_TOKEN = '{azSAStoken}' )
DIRECTORY = ( ENABLE = true );
""" )
{'columns': ['status'],
'data': [('DATABRICKS_IMPORT already exists, statement succeeded.',)]}
```

This statement employs a Python f-string using the connection variables defined earlier:

snowStage	The name within Snowflake for the stage area
snowBlobURL	Location of the stage's storage, in this case, the Azure container



azSAStoken	The Azure token that grants access to the container
------------	---

We execute a COPY INTO to generate the collection of Parquet objects in the stage area:

tableName	DIAGNOSE	S_ICD							
Just now (1s)			21		Python	Û	\$	53	÷
# Execute snowf] # Make note of t	ake "copy int he total rows	o" to expor and compar	⁺t tableName in ∙e with the imp	to the ADLS cont ort cell later	tainer as a	parqu	et f	ile	
prettify(snow_ex	ec(f ^{ran}		-						
from {tableN	<pre>stage}/{table ame}</pre>	wame}/{tabi	Leivame ;						
OVERWRITE =	TRUE								
FILE_FORMAT	E = (TYPE = PAR	QUET);							
·····))									
rows_unloaded	input_bytes ou	utput_bytes							
0 4756326	26943015	26943015							

The Snowflake statement converts the table into snappy compressed Parquet objects, each about 25MB in size, under a folder with the same name as the table. We query ADLS to observe the objects:

▶ ✓ Just now (<1s)	23	Python	Û	♦	[]	:
# Check to see the blob now exists in	n the container					
for blob in azGetBlobs.list_blobs(nam	ne_starts_with=tableName)	: print(blob.name)				
DIAGNOSES_ICD						
DIAGNOSES_ICD/DIAGNOSES_ICD_0_0.snappy	/.parquet					
DIAGNOSES_ICD/DIAGNOSES_ICD_0_1_0.snappy	/.parquet					
DIAGNOSES_ICD/DIAGNOSES_ICD_0_2_0.snappy	/.parquet					
DIAGNOSES_ICD/DIAGNOSES_ICD_0_3_0.snappy	/.parquet					

These objects can now be loaded into Databricks.

Step 2: Import the Parquet Objects into Databricks SQL under Unity Catalog

First, we set up identifiers to make the process repeatable for different tables:





Notice the target table patient_import.silver.DIAGNOSES_ICD is described using the 3-level namespace of Databricks Unity Catalog: <catalog>.<schema>.<data> similarly to the way

Snowflake manages datasets.

This requires the catalog and schema must already exist.

As the data is loaded,

Unity Catalog automatically collects the meta-data from the Parquet objects. The load process is like the earlier export process but is reversed.





Notice the row count matches the earlier export count. The data is managed through Unity Catalog, which provides a remarkably helpful "starter" glossary definition, demonstrating how effectively GenAI seamlessly integrates into the user workflow.



Step 3: Rinse

The rinsing step is to validate the data. We've already noted the row counts match. We examine the meta-data equivalence:



	- 5	nowflake –			– Databricks –				
P	Get the Snowfl rettify(snow_ex DESCRIBE TAB	ake meta-data ec(f=== LE {snowScheme	for the 1	table Name}	%sq DES	l Get the Databrid CRIBE TABLE \${s t	tks meta-data for tagefile.target}		
))				Table				
	name	type	kind	null?		A ⁸ c col_name	A ^B c data_type	A [®] c comme	
0	SUBJECT_ID	NUMBER(38,0)	COLUMN	N	1	SUBJECT_ID	decimal(38,0)	(531)	
1	HADM_ID	NUMBER(38,0)	COLUMN	N	2	HADM	decimal(38,0)	[8651]	
2	SEQ_NUM	NUMBER(38,0)	COLUMN	N	3	SEQ_NUM	decimal(38,0)	(883)	
3	ICD_CODE	VARCHAR(255)	COLUMN	N	-4	ICD_CODE	string	(mall)	
	ICD_VERSION	NUMBER(38,0)	COLUMN	N	5	ICD_VERSION	decimal(38,0)	(1011)	
		ALL DOLLARS AND	001111	~	12	ICD UEP CODE	Control of	restant!	

As expected, VARCHAR has been converted to STRING. We test the cardinality of the 4.75 million rows by checking the counts of the distinct SUBJECT_ID:



Step 4: Repeat

The same process is repeated for the remaining schema tables.

Keep in mind the method we've presented is to demonstrate concepts. It lacks robust error handling and needs deeper data validation that should be automated.

Once all the necessary tables are transferred, a business query demonstrates accurate relationships, cardinality, consistency of field values, and their correct filtering and aggregation handling.

Here is the "pneumonia" business statement from earlier, first running in Snowflake, the original source:



<pre># Snowflake Top 5 DRGs with the highest number of pneumonia diagnoses snow_exec(f"USE {snowSchema}") prettify(snow_exec(localize_sql(business_sql_test, "snowflake")))</pre>							
	DRG Description	Pneumonia Diagnosis	Mortality Count				
0	OTHER PNEUMONIA (139)	41	0				
1	SEPTICEMIA & DISSEMINATED INFECTIONS (720)	36	8				
2	HEART FAILURE (194)	13	0				
3	MAJOR RESPIRATORY INFECTIONS & INFLAMMATIONS (9	1				
4	TRACHEOSTOMY W MV 96+ HOURS W EXTENSIVE PROCED	8	0				

In the code, the variable business_sql_test holds the statement, which is syntax adapted for Snowflake using localize_sql(), a simple wrapper we created for the popular SQLglot library. Here is the same statement producing the same results, executing in Databricks:

<pre># Databricks Top 5 DRGs with the highest number of pneumonia diagnoses spark.sql(f"USE {dbxSchema}") display(spark.sql(localize_sql(business_sql_test, "databricks"))) + (6) Spark Jobs</pre>								
Table ↔ + Q 🏹								
	\mathbb{A}^{B}_{C} DRG Description	1 ² 3 Pneumonia Diagnosis	1 ² 3 Mortality Count					
- 4	OTHER PNEUMONIA (139)	41	0					
2	SEPTICEMIA & DISSEMINATED INFECTIONS (720)	36	8					
3	HEART FAILURE (194)	13	0					
4	MAJOR RESPIRATORY INFECTIONS & INFLAMMATIONS (137)	9	1					
5	TRACHEOSTOMY W MV 96+ HOURS W EXTENSIVE PROCEDURE (004)	8	0					

Although unnecessary, on the next page you can see the full statement running within Databricks SQL producing the same results, illustrating that Databricks SQL and the Databricks cluster are interacting with the same data store within the Databricks platform.

CTI DATA

@ (datab	oricks	Q Search data, notebooks, recents, and more CTR	L + P	шм~ 💠 📧		
	🗅 🔂 Pneumonia admissions + mortality × +						
A Pneumonia admissions + mortality : 🖈 New SQL editor: ON 🗸 💿 Scher			dule Share Save* ^				
	♦	► Run	Run all (1000) 👻 🗸 Just now (1s) 🗄 hive_metastore 🖨 default Y 🔗 🖓				
		<pre>USE patient_import.silver; SELECT concat(OR.DESCRIPTION,' (', DR.DRG_CODE, ')') AS `DRG Description', count(DISTINCT A.HADM_ID) AS `Pneumonia Diagnosis', count(DISTINCT CASE MHEN A.HOSPITAL_EXPIRE_FLAG = TRUE THEN A.HADM_ID ELSE NULL END AS `Mortality Count` FROM ADMISSIONS AS A JOIN DIAGNOSES_ICD AS D ON A.HADM_ID = D.HADM_ID ON D.ICD_VER_CODE = DI.ICD_VER_CODE JOIN DRGCODES AS DR ON A.HADM_ID = DR.HADM_ID MHERE EXTRACT(YEAR FROM A.ADMITTIME) = 2124 AND DI.LONG TILLE ILIKE 'Xoneumonia%'</pre>					
	23 GROOP BT 24 DR.DESCRIPTION. Add parameter ← Results 2 of 2 → Table ~ +				¥ Q 7 ⊡ ^ X		
			A C DRG Description	1 Pneumonia Diagnosis	1 3 Mortality Count		
			SEPTICEMIA & DISSEMINATED INFECTIONS (720)	36	8		
			HEART FAILURE (194)	13	0		
			MAJOR RESPIRATORY INFECTIONS & INFLAMMATIONS (137)	9	1		
		5	TRACHEOSTOMY W MV 96+ HOURS W EXTENSIVE PROCEDURE (004)	8	0		
		± 5		Refreshed now			



7.0 Conclusion

As we have shown, the mechanics of transferring data from Snowflake to Databricks are not complicated; large datasets can be relocated to the Databricks platform for initial experimentation within days.

On the other hand, moving data with all its application logic and other dependencies must be handled as you would any platform conversion process.

The key to success is early experimentation and scoping the work by data product, as in the outline below:



Start by targeting a simple data product to develop an appreciation for the work while exploring Databricks' capabilities during the process, and keep these principles in mind:





Beyond the Forklift

Examining if a forklift strategy is a lost opportunity is appropriate. Such changes are a chance to address technical debt; moving to Databricks opens new possibilities for rethinking the end-toend environment with far-reaching benefits that can transform your data's time-to-value, accessibility, trust, and cost by way of these benefits:

- Consolidate data warehousing and data science, eliminating separate systems and their data transfer complexities.
- Consolidate data engineering under a unified pipeline, storage, and data processing stack.
- Complete complex analytical queries orders of magnitude faster by leveraging the Databricks massive scalability query engine.
- Unify business rules and data governance, eliminating inconsistencies and improving data trust.
- Replace analytic tooling with self-service GenAI-enabled Exploratory Data Analysis.



About CTI Data

Our data and analytics experts specialize in Digital Transformation, Advanced Analytics, AI/ML, and Data Marketplaces. This experience provides valuable insights and expertise. We are adept at understanding best practices, identifying potential pitfalls, and customizing solutions to meet your unique needs.

By partnering with us, you can drive value from digital transformation efforts as we examine your business strategy, analyze your current state, pinpoint opportunities, and develop a strategic roadmap that aligns technology investments with strategic goals. We commit to collaborating closely with you and sharing accountability for achieving mutual goals.

<u>Contact us</u> to explore our real-world case studies and learn more about how we have helped our clients grow and create business value.

Disclaimer: This whitepaper is for informational purposes only and does not constitute professional advice. While we have endeavored to ensure the accuracy and completeness of the information contained herein, CTI Data makes no representations or warranties regarding its accuracy or completeness. The information presented is based on current knowledge and understanding and may be subject to change. References to third-party data or findings are for informational purposes only, and CTI Data assumes no responsibility for the accuracy of such third-party information. The limitations of the technologies or methodologies discussed in this whitepaper should be carefully considered before applying them in any specific context.