

# The Hitchhiker's Guide to Delta Lake Streaming



databricks  
**DATA+AI**  
**SUMMIT**

In an Agentic Universe

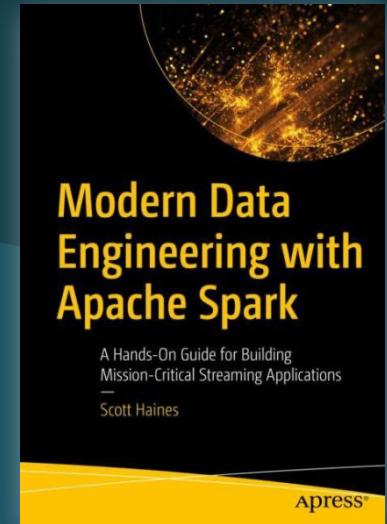
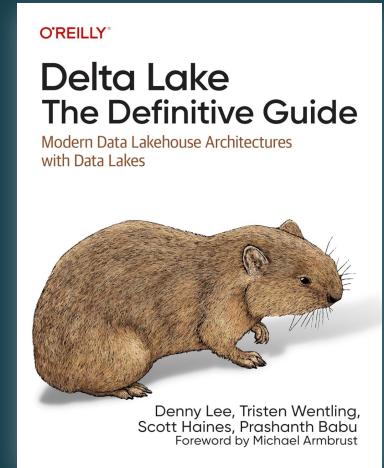
Scott Haines June 10<sup>th</sup>, 2025

Piloting the Session:  
  
**Scott Haines**



### "A Little About Me"

- I'm a Distinguished Software Engineer at Nike.
- I enjoy teaching and mentoring, writing books, and learning about new and exciting things.
- I've been working on massive high-performance streaming systems for the better part of the last 15 years.



Preparing for the Journey

# Mission Prep

1. Writing “Generic Data Applications” as “Components”.
2. Using Agentic Workflow Augmentation.
3. Q/A or High Fives



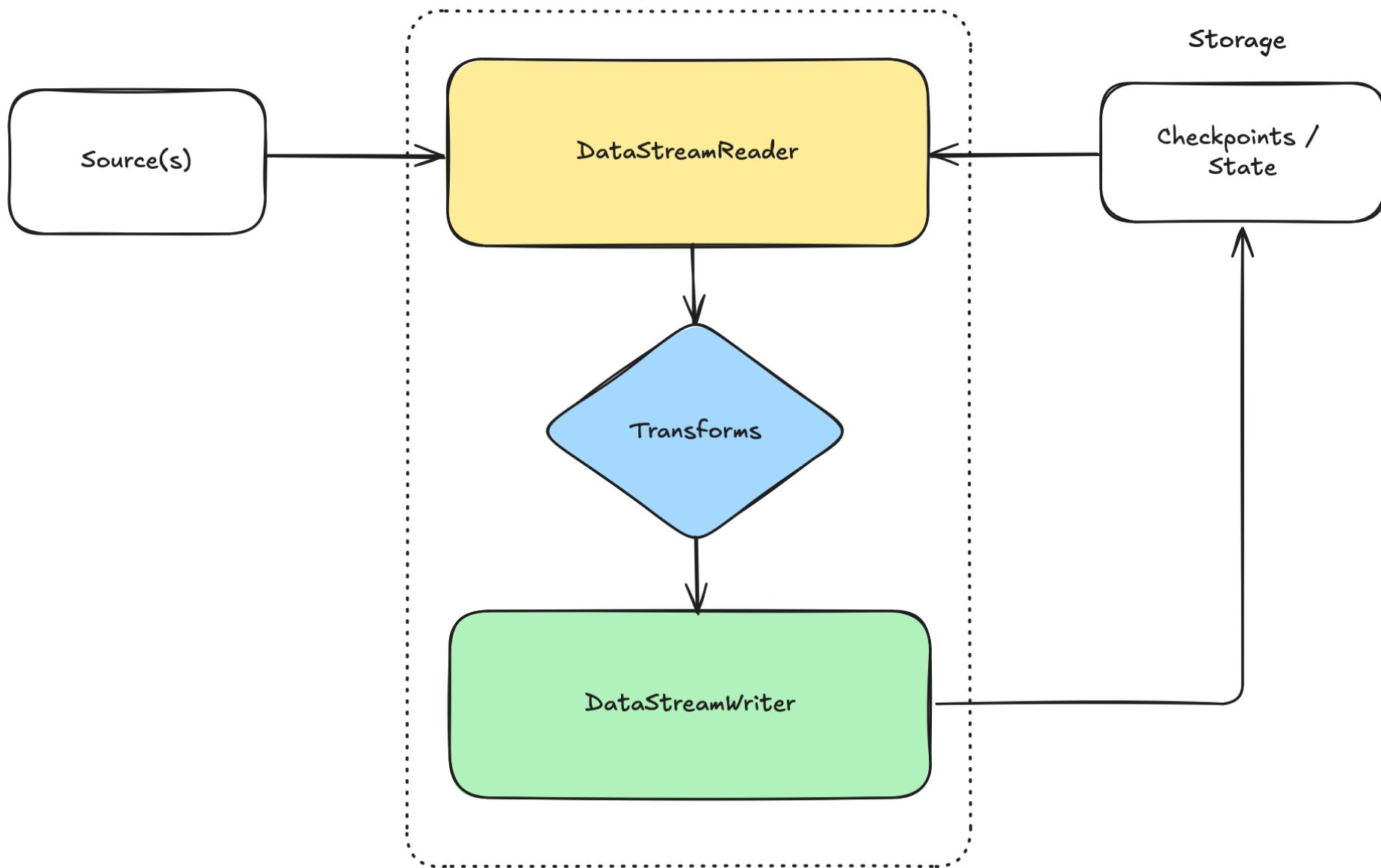
# On Writing Generic Data Applications as Components

Don't Panic. This part is simple ☺

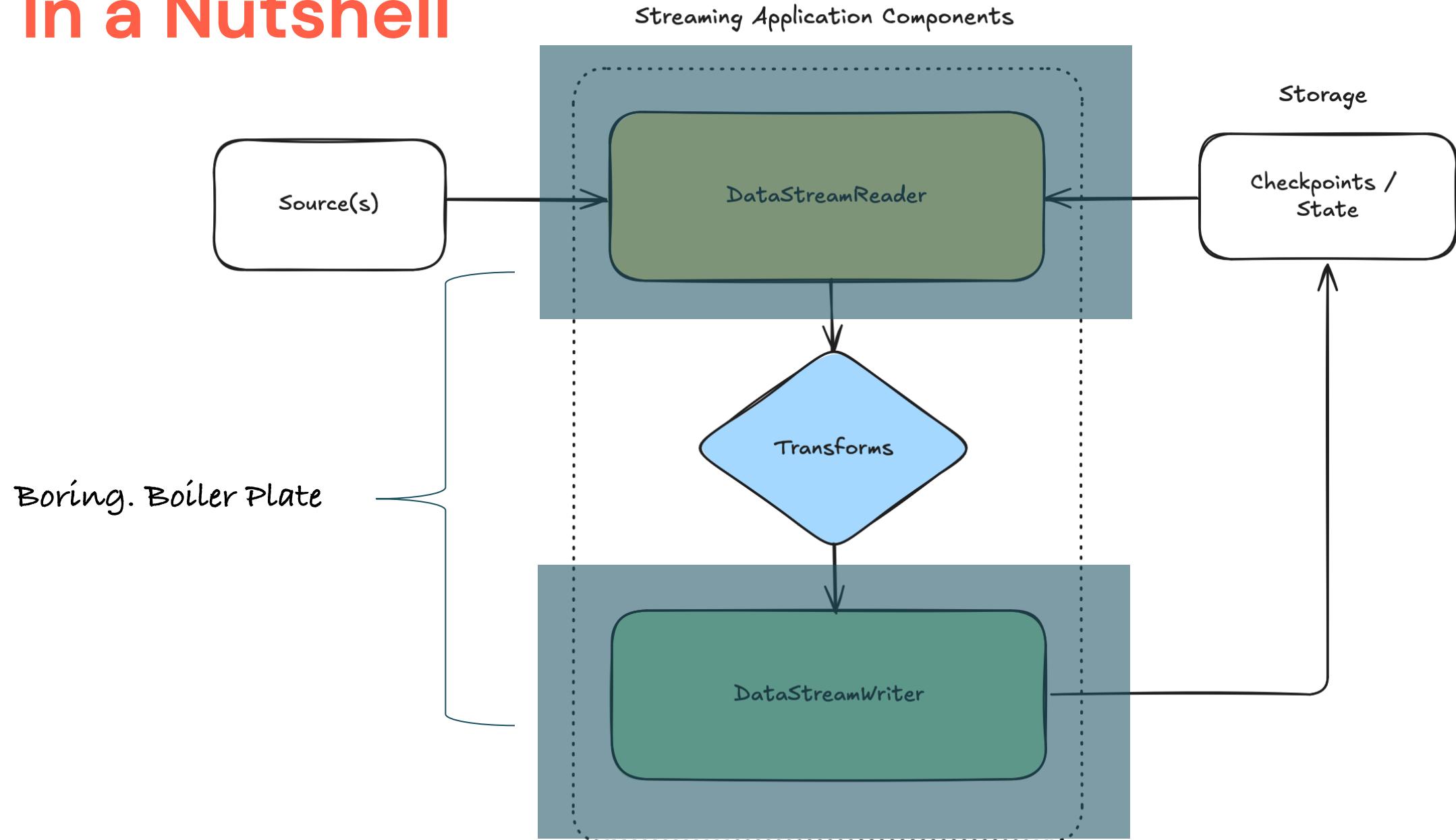


# In a Nutshell

## Streaming Application Components



# In a Nutshell



# Base Application Pattern

Provides Sane Defaults and Common Ways of Working

## Roles and Responsibilities

- Provide *Sane Defaults* and Common, Standardized Functionality.
- Establish a Common Way of Building and Tearing down applications.
- **Standardize where “checkpoints” and “state” live** — which comes in handy with greater automation.

```
class StreamingApp(SparkLoggingProvider):
    spark_session: SparkSession

    app_name: str = "generic-streaming-app"
    app_version: str = "0.0.1"
    app_logging_prefix: str = "StreamingApp::genAI"
    app_checkpoint_path: Optional[str] = None
    app_checkpoint_version: Optional[str] = None

    def __init__(self, spark: Optional[SparkSession] = None):
        pass

    def generate_checkpoint_path(self) -> Path:
        pass

    def generate_read_stream(self, s_options: Dict[str, str]) -> DataStreamReader:
        pass

    @staticmethod
    def generate_write_stream(
            df: DataFrame,
            s_options: Dict[str, str]) -> DataStreamWriter:
        pass
```

Pass: used to fit on slide. Imagine what is needed ^^

# Base Application Pattern

Provides Sane Defaults and Common Ways of Working



## Roles and Responsibilities

- Provide *Sane Defaults* and Common, Standardized Functionality.
- Establish a Common Way of Building and Tearing down applications.
- **Standardize where “checkpoints” and “state” live** — which comes in handy with greater automation.
- *Yes. This is config-driven development...*

```
fresh_session = self.spark_session.newSession()  
fresh_session.conf.set('spark.app.name',  
                      'delta_to_delta_app')  
fresh_session.conf.set("spark.app.version", "3.0.0")  
fresh_session.conf.set("spark.app.checkpoints.path",  
                      "/Volumes/dev/app/fs/d2d/")  
fresh_session.conf.set('spark.app.checkpoints.version',  
                      'canary')
```

```
app = StreamingApp(spark=fresh_session).|
```

(f) app_name	StreamingApp
(f) app_version	StreamingApp
(f) spark_session	StreamingApp
(f) app_checkpoints_path	StreamingApp
(f) app_logging_prefix	StreamingApp
(f) app_checkpoint_version	StreamingApp
(m) generate_checkpoint_location(self)	StreamingApp
(m) generate_read_stream(self, s_options)	StreamingApp
(m) generate_spark_session	StreamingApp
(m) generate_write_stream	StreamingApp
(f) logger	StreamingApp
(m) get_logger(self, spark_prefix)	StreamingApp

Press ⇧ to insert, → to replace Next Tip

# Data Stream Reading

## Using Config-Driven Design Patterns



### Roles and Responsibilities

- Must provide a mechanism to populate the correct configuration for a given streaming data source.
- This is a piece of cake for Delta Lake.
- *We'll look at the available options next.*

```
def generate_read_stream(self, s_options: Dict[str, str])  
    """  
    Generate the DataStreamReader. Note: This method can be used  
    :param s_options: The dictionary of key/value pairs  
    :return: The DataStreamReader object for composition  
    """  
  
    return self.spark_session.readStream.options(**s_options)
```

# Delta Streaming Source

## Using Config-Driven Design Patterns

### Source Component's Must:

- Auto-Wire themselves using Spark's Runtime Configuration
- Provide "sane" defaults
- Fail Fast and Correct the "user"

```
class DeltaSourceGenerator:  
    app: StreamingApp  
    config_prefix: str = 'spark.app.source.delta'  
    table_name: str  
  
    options_default: Dict[str, Optional[str]] = {}  
  
    def generate_options(self) -> Dict[str, str]:  
        """  
        Filters and expands options for a delta streaming source  
        :return: the filtered and expanded options  
        """  
        return dict(  
            (opt[0], opt[1]) for opt in self.options.items() if opt[1] is not None  
        )  
  
    def generate_stream_source(self) -> DataStreamReader:  
        """  
        Returns the generated stream. Requires calling `load` or `table`  
        method (toDF) to generate a `DataFrame`  
        :return: The DataStreamReader  
        """  
        return self.app.generate_read_stream(  
            self.generate_options()  
        ).format('delta')  
  
    def toDF(self) -> DataFrame:  
        table_name = self.table_name  
        is_managed = "/" not in table_name  
  
        reader: DataStreamReader = self.generate_stream_source()  
        return reader.table(table_name) if is_managed else reader.load(table_name)
```

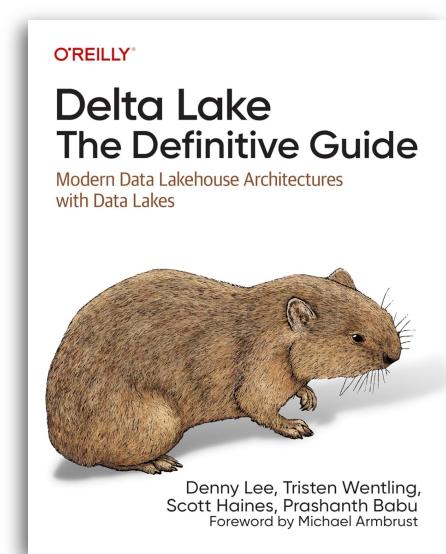
# Data Stream Reading

## Common Delta Lake Streaming Source Options



### Getting a Handle on Delta's Streaming Source Options

- **Rate Limiting** *maxFilesPerTrigger / maxBytesPerTrigger*
- **Runtime Ignorance** *ignoreDeletes / ignoreChanges*
- **Time Travel and Pointers for Stream Triage** *startingVersion / startingTimestamp*
- **Stream Correctness** *withEventTimeOrder*
- **Powering through with Column Mapping Enablement** *schemaTrackingLocation*



# Data Stream Reading

## Generic Delta Lake Source Options Reading

### Delegate and Be Lazy Where Possible

- Read the configuration via the Spark runtime config. With this pattern, the standard app (StreamingApp) can delegate config parsing and supply defaults.
- **Remember:** We are trying to do as little as possible here. The less we do, the easier our application is to manage.

```
class DeltaSourceGenerator:  
    app: StreamingApp  
    config_prefix: str = 'spark.app.source.delta'  
    table_name: str  
  
    options_defaults: Dict[str, Optional[str]] = {...}  
  
    def __init__(self, app: StreamingApp):  
        """  
        Lazy auto-wired source  
        """  
        self.app = app  
        self.table_name = self._table_name()  
        self.options = Functions.generate_options_dict(  
            self.config_prefix, self.options_defaults, self.app.spark_session  
        )
```

\*Use `spark.conf` for config which helps with notebooks or vanilla data apps

Next Up.  
Transformations...

Don't Panic. Still simple ☺



# Data Stream Transformations

## The Heart of the Application

### Transformations.

- This is where the real work happens.
- Imagine a world where you no longer need to think about the “source(s)”, where you can **focus** on the output of your data applications, and simply config your “source(s) and sink(s)”.
- In other words, **focus on doing** vs the boiler plate.

```
spark.conf.set("spark.app.source.delta.table.catalog", catalog_name)
spark.conf.set("spark.app.source.delta.table.schema", schema_name)
spark.conf.set("spark.app.source.delta.table.name", table_name)
spark.conf.set("spark.app.checkpoints.path", checkpoints_path)
spark.conf.set("spark.app.sink.delta.table.catalog", sink_catalog_name)
spark.conf.set("spark.app.sink.delta.table.schema", sink_schema_name)
spark.conf.set("spark.app.sink.delta.table.name", sink_table_name)

app = StreamingApp(spark=spark)
delta_source_df = DeltaSourceGenerator(app).toDF()
transformed = delta_source_df.transform(...)
streaming_query = DeltaSinkGenerator(app, df=transformed).toTable()

# inject monitoring, awaiting termination
```

We can handle this ^^. Job Well Done

# Next Up. Finishing the Writer

Don't Panic. Still simple ☺



# Data Stream Writing

Take a DataFrame. Return a Stream.

## Rinse and Repeat.

- Just like with the StreamingApp's generate\_read\_stream method.
- We can utilize composition in order to abstract away complexity.
- Let's look at that next.

```
@staticmethod 1 usage  ↗ Scott Haines
def generate_write_stream(
    df: DataFrame, s_options: Dict[str, str]
) -> DataStreamWriter:
    """
    Generate the DataStreamWriter. Note: This method can be further composed by
    calling start or toTable on the DataStreamWriter creates a StreamingQuery.
    :param df: The DataFrame to compose the pipeline
    :param s_options: The dictionary of key/value pairs
    :return: The DataStreamReader object for composition
    """

    return df.writeStream.options(**s_options)
```

StreamingApp.generate\_write\_stream



# Data Stream Writing

## Generating the Delta Lake Streaming Sink

### Abstraction Away.

- Utilize Composition to drastically simplify writing boiler plate Delta Lake Streaming applications
- Lean on simple configuration.
- Enabling us to then simply focus on transformations, which we can get into in the next section.

```
class DeltaSinkGenerator:  
    app: StreamingApp  
    config_prefix: str = 'spark.app.sink.delta'  
    table_name: str  
    read_df: Optional[DataFrame] = None  
  
    write_config_defaults: Dict[str, str] = {  
        'outputMode': 'append',  
        'timestampColumn': 'timestamp'  
    }  
  
    options_default: Dict[str, Optional[str]] = {  
        'overwriteSchema': 'false',  
        'mergeSchema': 'false'  
    }  
  
    def __init__(self, app: StreamingApp, read_df: Optional[DataFrame] = None):  
        self.app = app  
        self.read_df = read_df  
        self.table_name = self._table_name()  
  
        self.options = {  
            'checkpointLocation': str(self.app.generate_checkpoint_location()),  
            ...  
        }  
  
    def fromDF(self, read_df: Optional[DataFrame]=None) -> DataStreamWriter:  
        """  
        Like toDF, will generate a stream by receiving a DataFrame  
        and auto-wiring the streaming sink options  
        """  
        df = read_df or self.read_df  
        if df is None:  
            raise RuntimeError("be kind...")  
        builder: DataStreamWriter = self.app.generate_write_stream(  
            df, self.generate_options()  
        )  
        writer = Functions.auto_trigger(  
            Functions.auto_partition_or_clusterBy(builder, spark), spark)  
  
        return writer.outputMode(self.write_config['outputMode']).format('delta')  
  
    def toTable(self, read_df: Optional[DataFrame]=None) -> StreamingQuery:  
        """  
        Auto-Call and Auto-Wire the StreamingQuery using composition  
        """  
        writer = self.fromDF(read_df)  
        return writer.toTable(tableName=self.table_name)
```



# Generic Data Apps

## Towards Application Composition

- Data Sources can be configured in a generic way allowing our applications to be controlled via configuration.
- The same is true for our Data Sink's.
- This leaves \*us to focus on the core task of data "transformation".
- Luckily. We can get help there if we choose...

```
spark.conf.set("spark.app.source.delta.table.catalog", catalog_name)
spark.conf.set("spark.app.source.delta.table.schema", schema_name)
spark.conf.set("spark.app.source.delta.table.name", table_name)
spark.conf.set("spark.app.checkpoints.path", checkpoints_path)
spark.conf.set("spark.app.sink.delta.table.catalog", sink_catalog_name)
spark.conf.set("spark.app.sink.delta.table.schema", sink_schema_name)
spark.conf.set("spark.app.sink.delta.table.name", sink_table_name)

app = StreamingApp(spark=spark)
delta_source_df = DeltaSourceGenerator(app).toDF()
transformed = delta_source_df.transform(...)
streaming_query = DeltaSinkGenerator(app, df=transformed).toTable()

# inject monitoring, awaiting termination
```

# Assembling the Pieces

## Where Exactly are we Headed?

AI Agents can help us along our way by:

- Assisting us to find important information about our Tables
- Generating data for Tests (from our table information)
- Figuring out the right syntax for an Operation (joins, etc)
- And oh! so much more...

IT'S DANGEROUS TO  
GO ALONE! TAKE THIS.



# Workflow Automation with Agents

Don't Panic. This Part is Fun!





The Answer to the Great  
Question... Of Life, the Universe  
and Everything... Is...



ee

# Tools and the Number 42!

jj

# Crash Course: Agents and Tools

## Tools are Functions that Provide Context

### Why does Context Matter?

- LLMs know a “ton”. But sometimes struggle to distinguish between “a ton” – being a unit of measurement or an expression.
- Context provides a Guiding Light to a system in the Dark.
- LLMs understand language and patterns really well, so tools must also act as a part of the “memory” of the “system at large”



**Figure 42:** Tools act like “parts” that can be combined to engineer much more complex things.

### What is an Agentic Tool?

- It is a simple function or “api-like” interface that enables an LLM to “touch” the outside world.
- We provide the utility to bridge the gap between “the LLM void” and our “world”

```
1 class ListCatalogsRequest(BaseModel):
2     max_results: Optional[int] = Field(default=1,
3                                         description="The number of schema's to return from the request",
4                                         gt=0, lt=100)
5
6 class ListCatalogsResponse(BaseModel):
7     catalogs: List[str] = Field(description="The name of the base catalog")
8
9 @tool("list_catalogs",
10       description="returns a list of base catalogs. Individually, " \
11       "each catalog name can be used as the parent catalog reference for schema or table discovery",
12       args_schema=ListCatalogsRequest,
13       return_direct=True)
14 def list_catalogs(max_results: Optional[int] = 1) -> ListCatalogsResponse:
15     # returns an Iterator[CatalogInfo]
16     print(f"limit to max_results:{max_results}")
17     try:
18         # browse_only, catalog_type, comment, metastore_id, securable_type, properties
19         return ListCatalogsResponse(
20             catalogs=[clean_string(catalog.name) for catalog in workspace_client.catalogs.list()])
21     except Exception as e:
22         raise ToolException(f"Failed to List the Catalogs: {e}")
23
```



Speeding Up Now. Let's  
automate a data pipeline



# LangGraph

## Let's Build a Pipeline

### What Will We Need?

- Access to Unity Catalog
- Tools to do the following:
  - List Catalogs [check]
  - List Schemas of a Catalog
  - List Tables of a Schema in a Catalog
  - Get Table Details and all Known Metadata



## Let's Build a Pipeline

```
1 class ListSchemaRequest(BaseModel):
2     catalog_name: str = Field(description="The name of the parent catalog (Unity Catalog) containing this schema")
3     max_results: Optional[int] = Field(default=1, description="The number of schema's to return from the request", gt=0, lt=100)
4
5 class ListSchemaResponse(BaseModel):
6     catalog_name: str = Field(description="The name of the parent catalog (Unity Catalog) containing this schema")
7     schemas: Optional[List[str]] = Field(default=None, description="The name of the schemas. Schemas are also referred to as databases in older systems")
8
9 @tool(
10     "list_schemas",
11     description="returns a list of schema names, if they exist, from the associated unity catalog name",
12     args_schema=ListSchemaRequest,
13     return_direct=True)
14 def list_schemas(catalog_name: str, max_results: Optional[int] = 1) -> ListSchemaResponse:
15     # returns an Iterator[SchemaInfo]
16     print(f"limit to max_results:{max_results}")
17     try:
18         schemas = workspace_client.schemas.list(
19             catalog_name=clean_string(catalog_name),
20             include_browse=False,
21             max_results=max_results
22         )
23         return ListSchemaResponse(
24             catalog_name=clean_string(catalog_name),
25             schemas=[schema.name for schema in schemas]
26         )
27     except Exception as e:
28         raise ToolException(f"Failed to List the Schemas: {e}")
```



## Let's Build a Pipeline

```
1  class ListTablesRequest(BaseModel):
2      catalog_name: str = Field(description="The name of the parent catalog (Unity Catalog) containing this schema")
3      schema_name: str = Field(description="The name of the schema containing the tables to list. The schema is a child of the parent catalog named catalog_name")
4      max_results: Optional[int] = Field(default=10, description="The number of table's to list in this request", gt=0, lt=100)
5
6  class ListTablesResponse(BaseModel):
7      catalog_name: str = Field(description="The name of the parent catalog (Unity Catalog) containing this schema")
8      schema_name: str = Field(description="The name of the schema containing the tables to list. The schema is a child of the parent catalog named catalog_name")
9      tables: Optional[List[str]] = Field(default=None, description="The list of full table names. Each table's full name includes the catalog.schema.table_name.")
10
11 @tool(
12     "list_tables",
13     description="returns a list of table names, if any exist, from the associated catalog_name.schema_name provided in the request",
14     args_schema=ListTablesRequest,
15     return_direct=True)
16 def list_tables(catalog_name: str, schema_name: str, max_results: Optional[int] = 10) -> ListTablesResponse:
17     # returns an Iterator[TableInfo]
18     print(f"limit to max_results:{max_results}")
19     try:
20         tables = workspace_client.tables.list(
21             catalog_name=clean_string(catalog_name),
22             schema_name=clean_string(schema_name),
23             max_results=max_results,
24             include_delta_metadata=False,
25             include_manifest_capabilities=False,
26             omit_username=True,
27             omit_columns=True,
28             omit_properties=True,
29         )
30         return ListTablesResponse(
31             catalog_name=clean_string(catalog_name),
32             schema_name=clean_string(schema_name),
33             tables=[table.full_name for table in tables if table.full_name is not None]
34         )
35     except Exception as e:
36         raise ToolException(f"Failed to List the Tables under the {catalog_name}.{schema_name}")
```



# LangGraph

## Let's Build a Pipeline

```

1 # The TableColumn class is a lighter weight version of "TableInfo". This allows us to tune the amount of data we provide
2 # to the LLM. This can be extended if you wanted to "query" system tables to understand what "tags" are associated with a given table
3 # or even to learn from the lineage to see what other tables are "commonly" constructed from "this" table.
4 class TableColumn(BaseModel):
5     name: str = Field(description="The name of the column. This is used for querying the table")
6     sql_type: str = Field(description="Full data type specification as SQL/catalogString text")
7     comment: Optional[str] = Field(default=None, description="Describes what the column represents. The comment provides a semantic understanding for better queries")
8     nullable: Optional[bool] = Field(default=True, description="If the value is True, then this field isn't required and we can't count on it to always be present")

```

Python

```

1 class TableDetailsRequest(BaseModel):
2     catalog_name: str = Field(description="The name of the parent catalog (Unity Catalog) containing this schema")
3     schema_name: str = Field(description="The name of the schema containing the tables to list. The schema is a child of the parent catalog named catalog_name")
4     table_name: str = Field(description="The name of the table to inspect and discover. This is a child of the associated schema")
5
6 # This is fed from the ColumnInfo class via the Databricks SDK
7 class TableDetailsResponse(BaseModel):
8     full_table_name: str = Field(description="The full name of the table. The full name includes the catalog.schema.table_name")
9     comment: Optional[str] = Field(default=None, description="The table description if it has been provided")
10    columns: Optional[List[TableColumn]] = Field(default=None, description="The columns of the table. These are important for querying the table")
11
12 @tool(
13     "get_table_details",
14     description="Provides metadata about a table within Unity Catalog. The details include the table name, comment, description, columns, cluster or partition information, as well as the schema and catalog names. This tool is useful for inspecting the structure of tables in a large dataset." )
15 args_schema=TableDetailsRequest,
16 return_direct=True)
17 def get_table_details(catalog_name: str, schema_name: str, table_name: str) -> TableDetailsResponse:
18     full_name = f"{clean_string(catalog_name)}.{clean_string(schema_name)}.{clean_string(table_name)}"
19     print(f"getting table details: for '{full_name}'")
20     try:
21         table_details = workspace_client.tables.get(full_name=full_name, include_delta_metadata=True, include_browse=True, include_manifest_capabilities=True)
22         #print(f'{table_details}')
23         return TableDetailsResponse(
24             full_table_name=full_name,
25             comment=table_details.comment,
26             columns=[
27                 TableColumn(
28                     name=column.name,
29                     sql_type=column.type_text,
30                     comment=column.comment,
31                     nullable=column.nullable
32                 ) for column in table_details.columns]
33             )
34     except Exception as e:
35         raise ToolException(f"Failed to collect the Table details and metadata for {full_name}")

```

No Notifications



# Fruits of “Our” Labor

Remember. This is only using Metadata



# LangGraph

Let's Run "Scott" 2.0

Full Video will be on the YouTube recording from the presentation. It is too big for upload.

# Closing Thoughts

## Towards Agentic Augmented Application Composition

- Data Pipelines are not that hard to make
- Using Config-Driven software engineering design patterns can take us 90% of the way.
- Getting agentic assistance is amazing, especially when it saves hours of effort, or actual meetings ☺

# Closing Thoughts

## Towards Agentic Augmented Application Composition

- Context with respect to “short term” memory is a tough nut to crack.
- Providing all the “right tools” can take a long time, especially when figuring out the right “level of focus”.

# Cheers

So Long and Thanks for all the Fish!

Now do we still have time for  
Questions? (future self>..)

