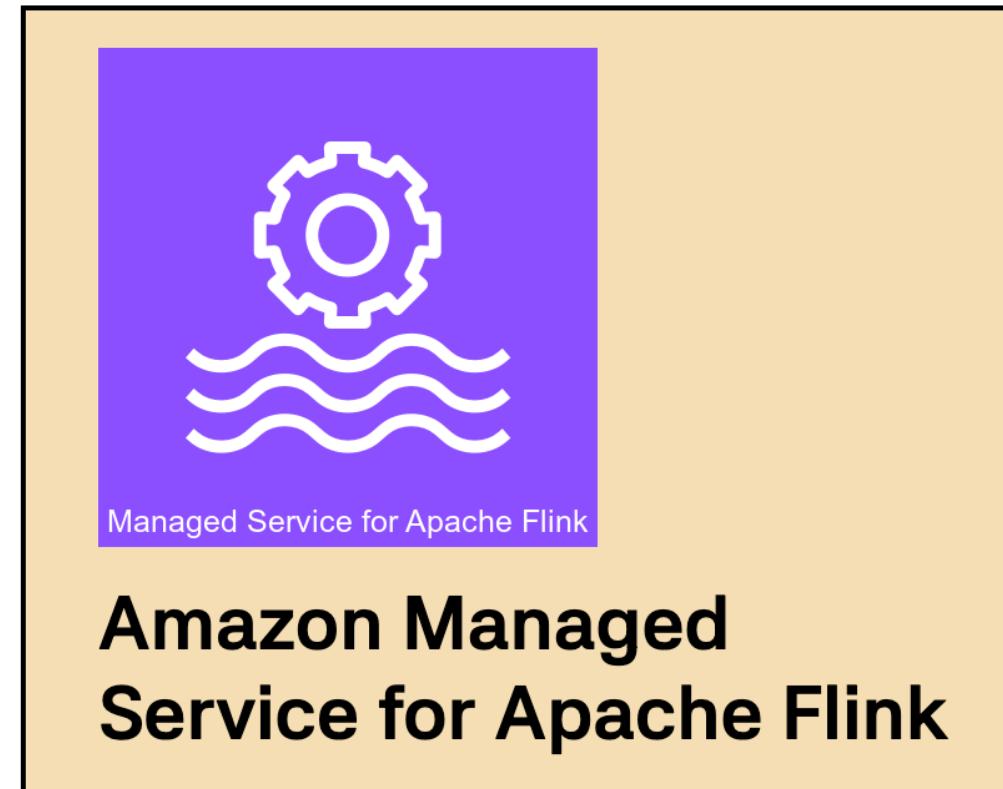


Amazon Managed Service for Apache Flink

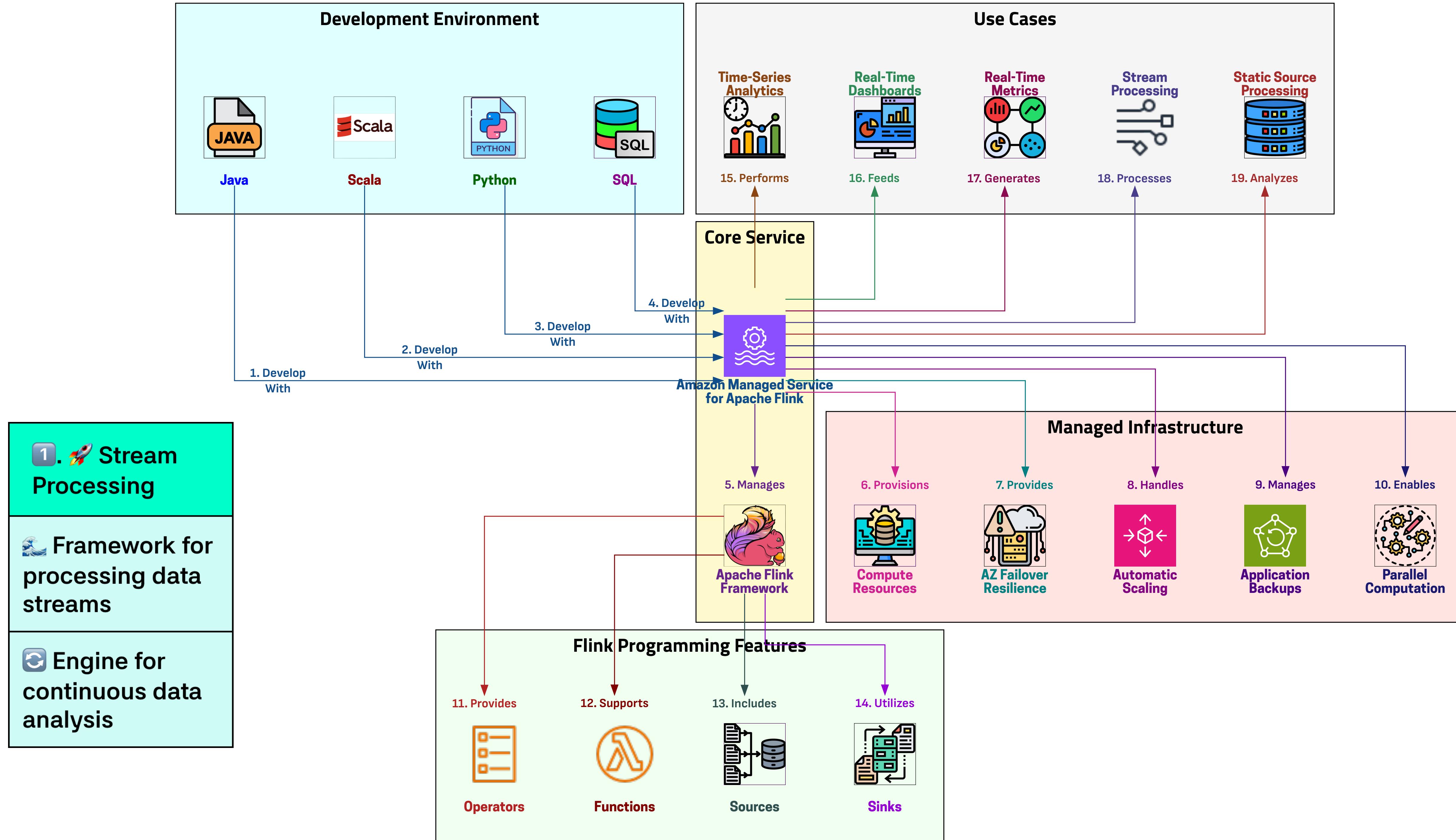


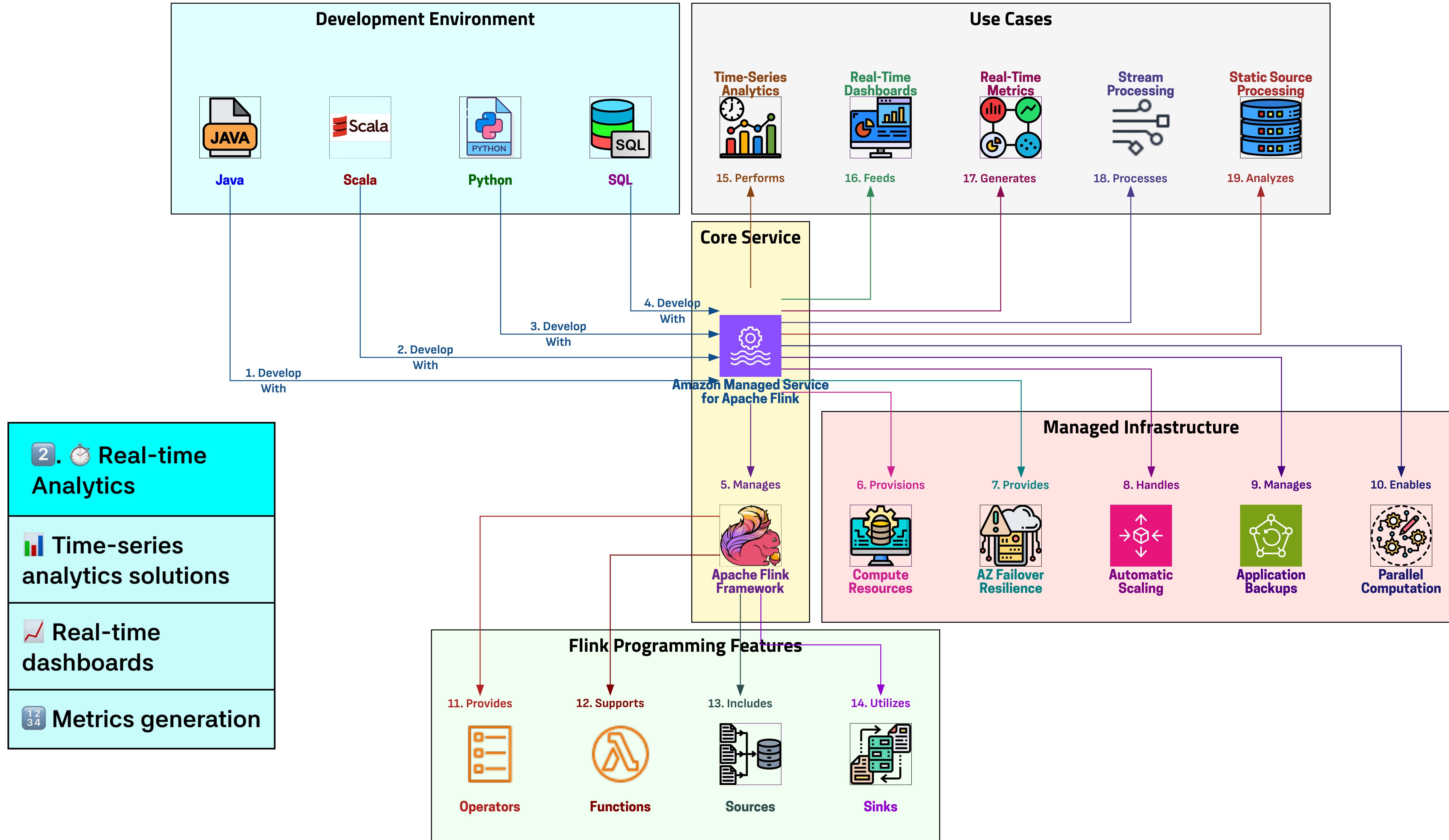
Amazon Kinesis Data Analytics for Apache Flink

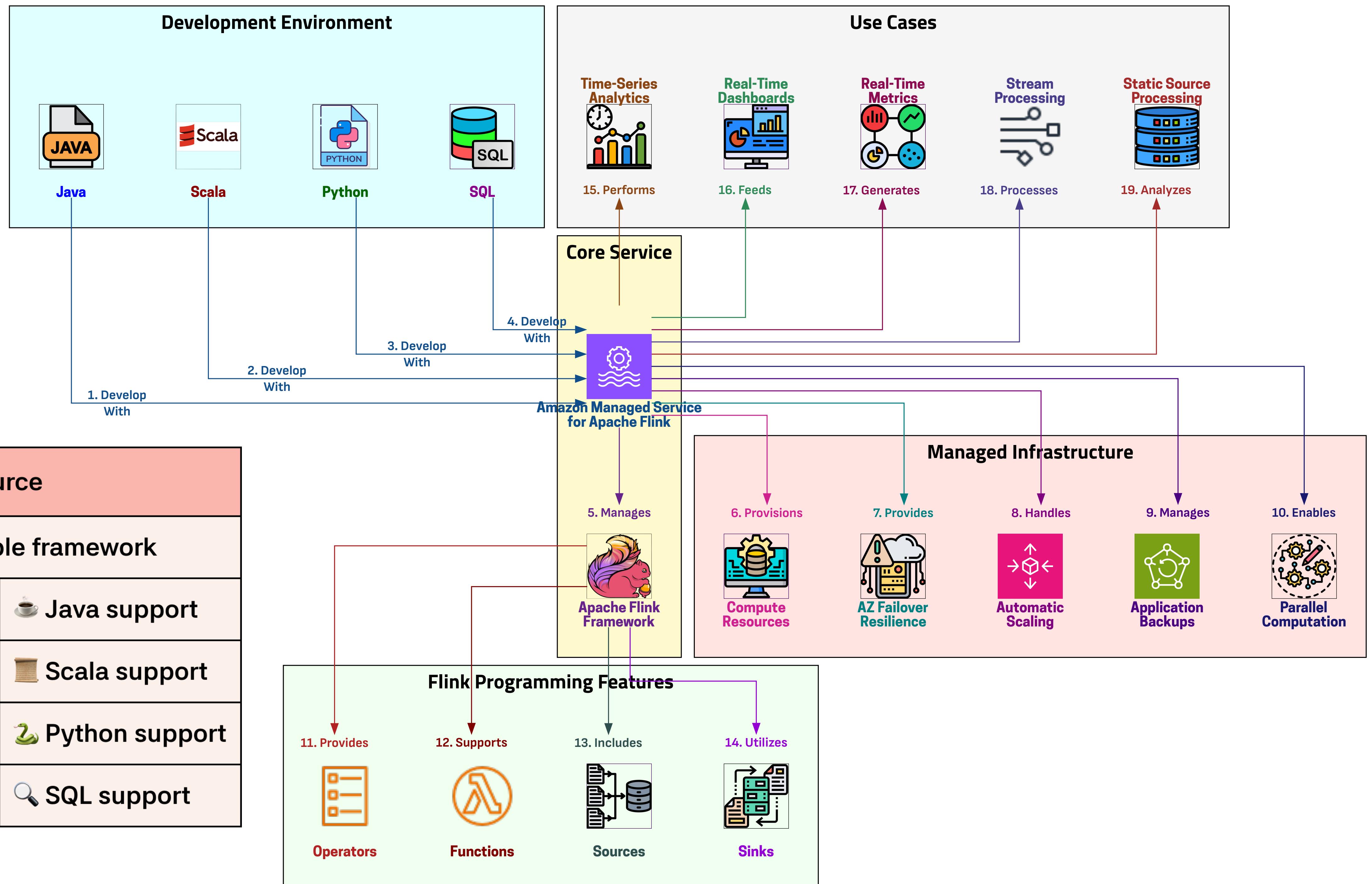
Table of Contents

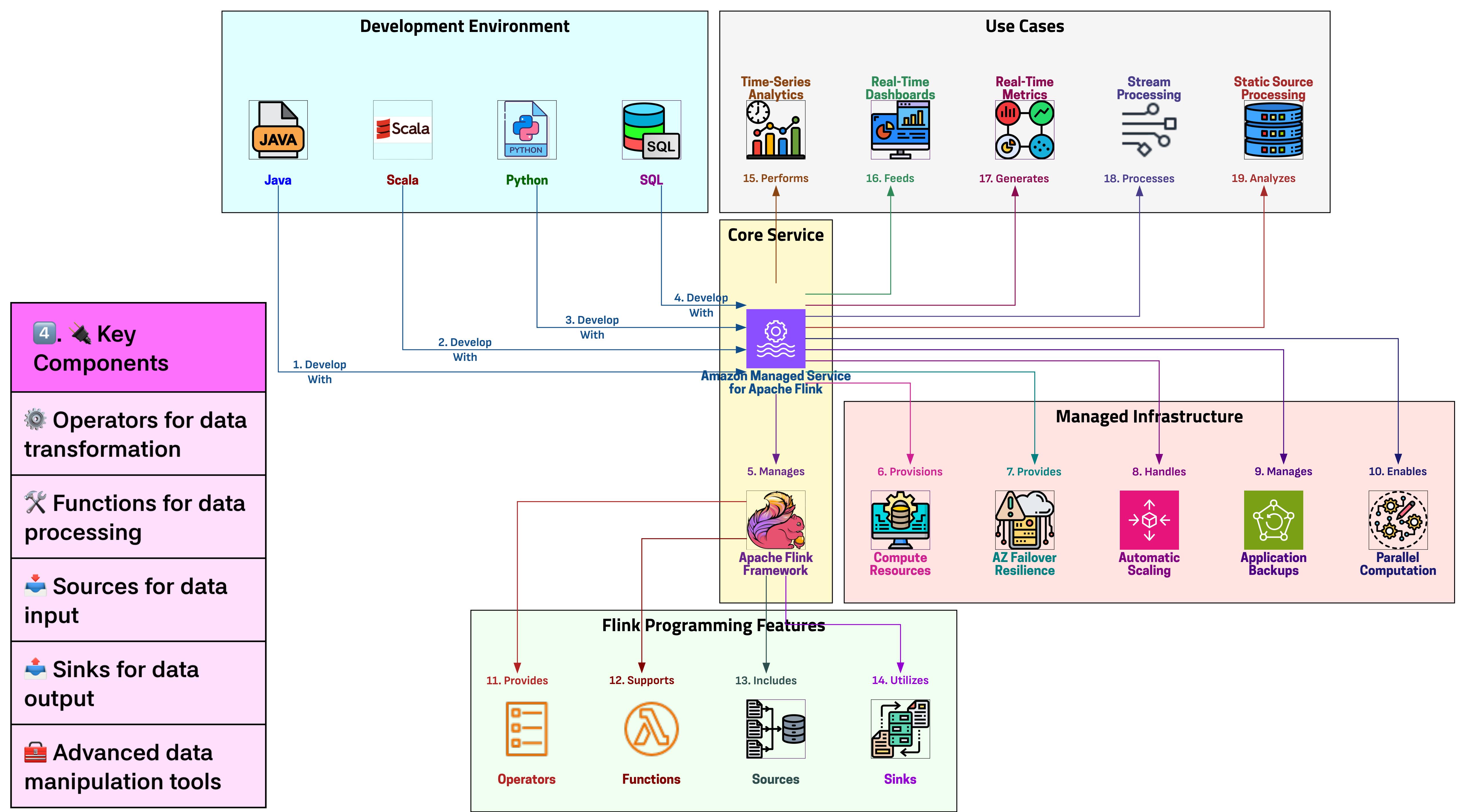


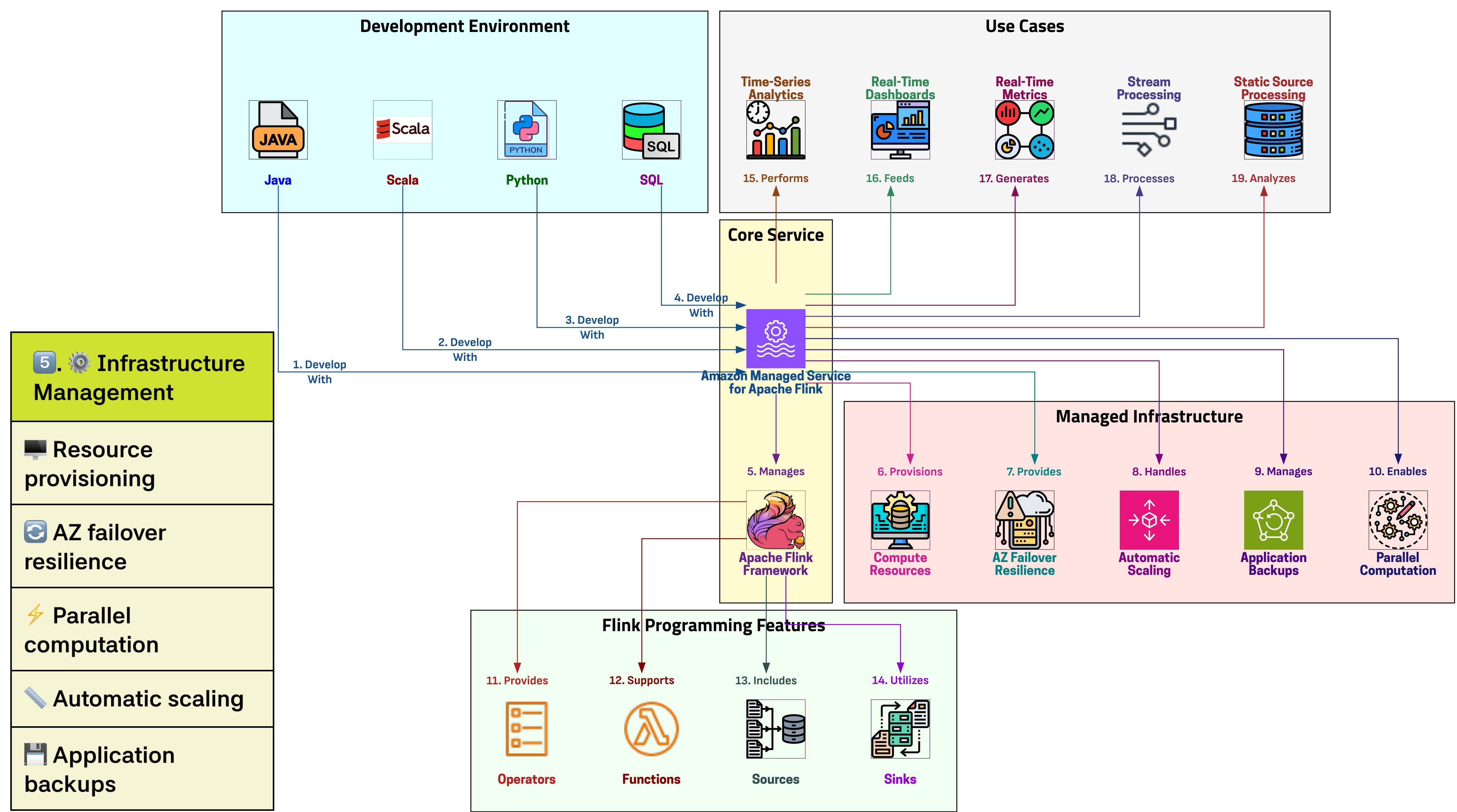
- 1 What is Amazon Managed Service for Apache Flink?
- 2 Apache Flink vs Apache Flink Studio
- 3 Apache Flink APIs
- 4 How it works?
- 5 Connectors (Sources & Sinks)
- 6 Use Kinesis data streams
- 7 Create a KafkaSource
- 8 Transform data using Operators
- 9 List of Operators
- 10 Use transform operators
- 11 Use aggregation operators

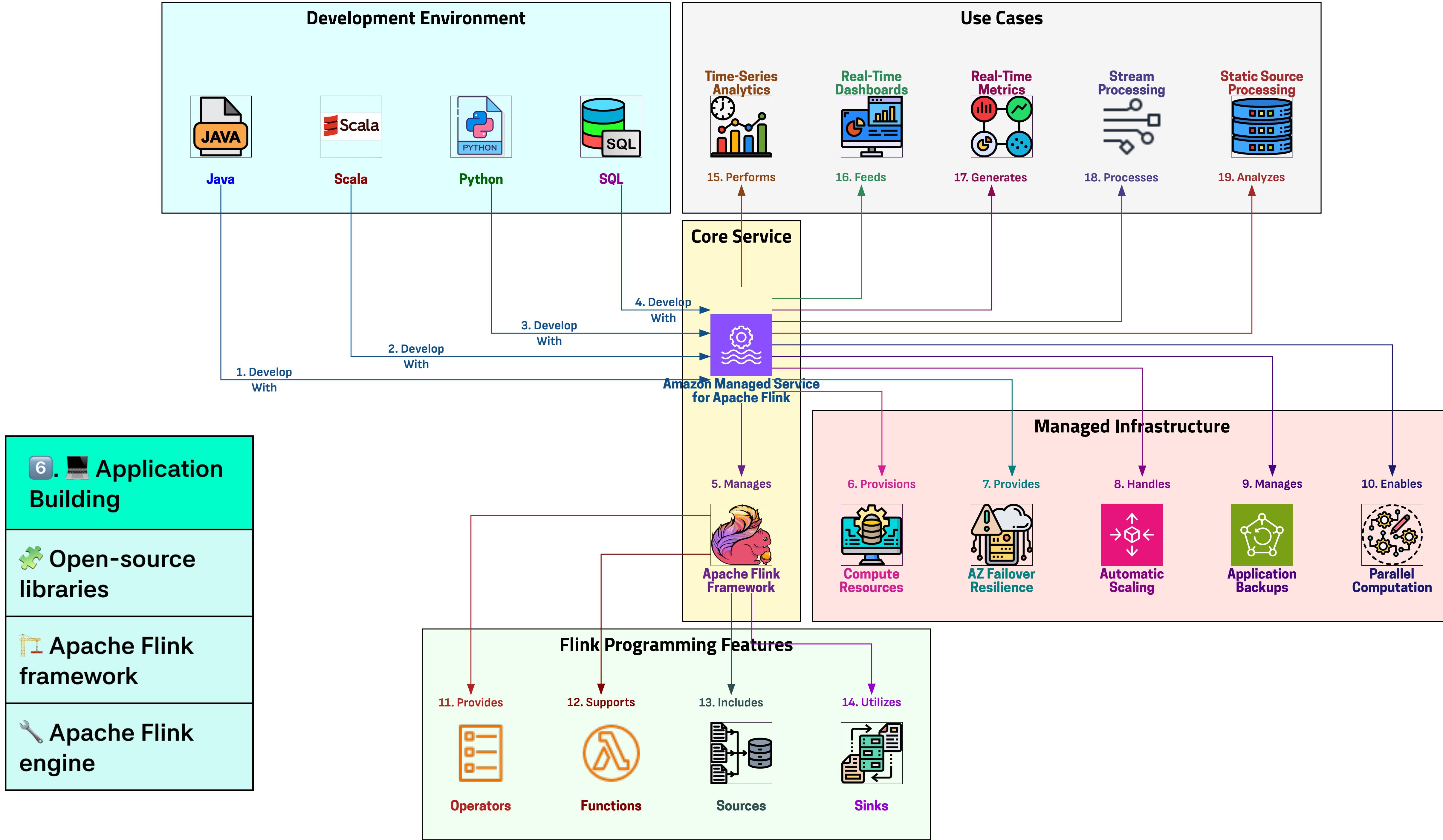












Apache Flink vs Apache Flink Studio

1 . **Two service options:** Managed Service for Apache Flink, Managed Service for Apache Flink Studio

2 . **Managed Service for Apache Flink capabilities:** Build applications using IDE of choice, Programming language support: Java, Scala, Python, API options: Datastream API, Table API

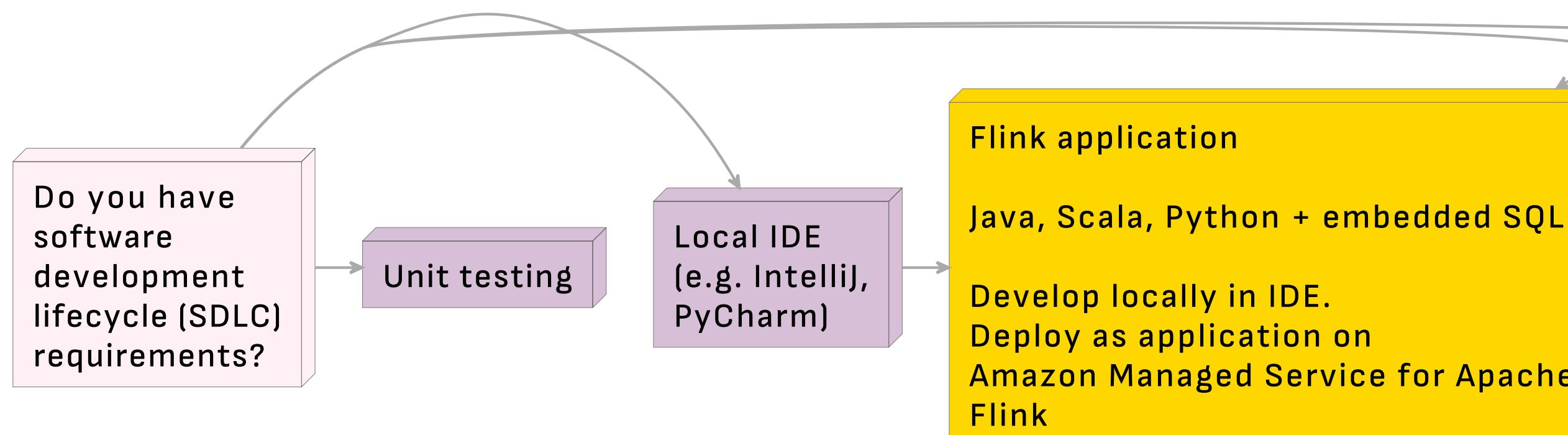
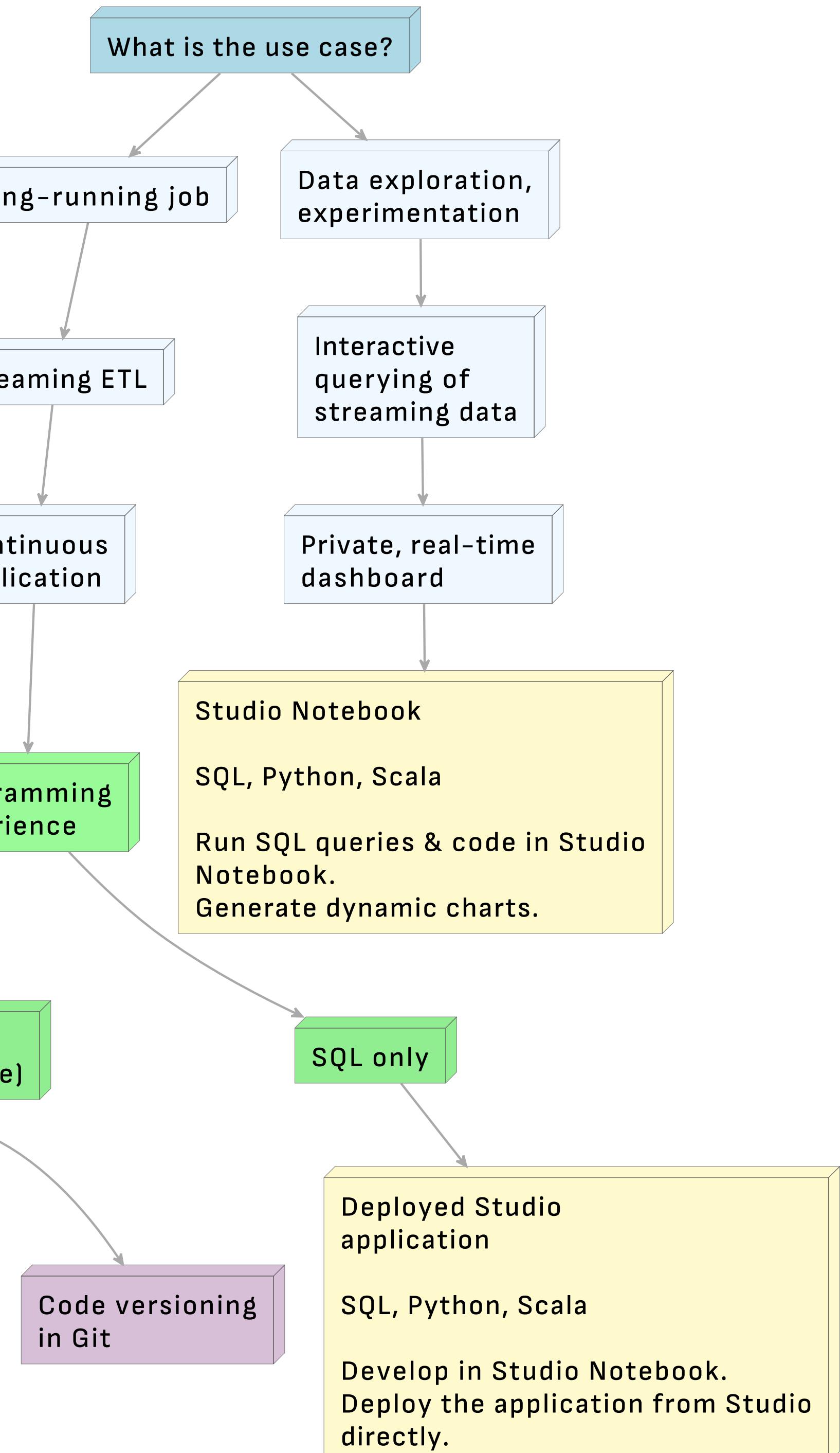
3 . **Managed Service for Apache Flink Studio features:** Interactive querying of data streams, Real-time processing, Application building tools: SQL, Python, Scala

4 . **Use case driven selection:** Choose based on specific requirements

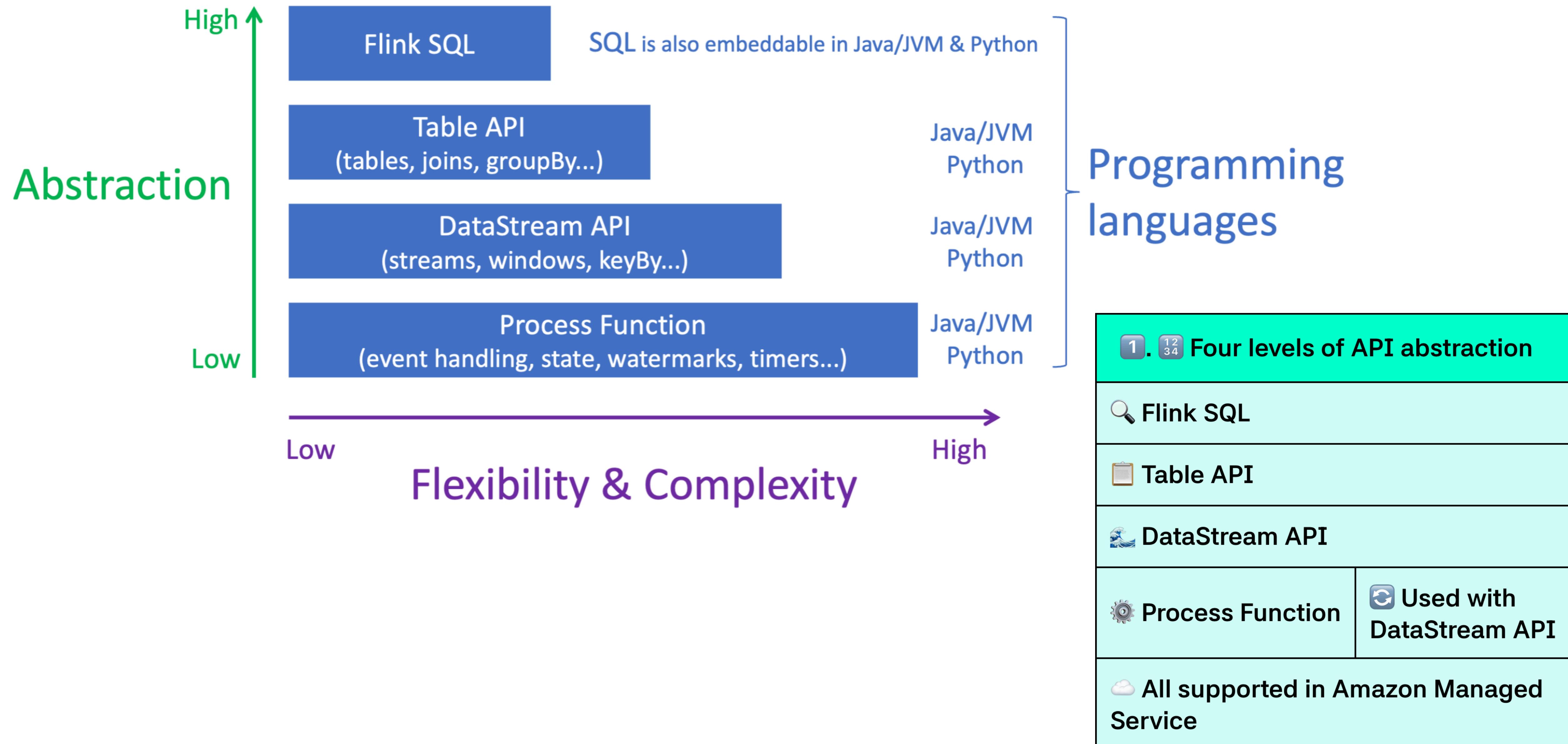
5 . **Managed Service for Apache Flink ideal scenarios:** Long-running applications: Streaming ETL, Continuous applications, SDLC integration: Software development lifecycle processes

6 . **Managed Service for Apache Flink Studio ideal scenarios:** Ad-hoc data exploration, Interactive streaming data queries, Real-time dashboards

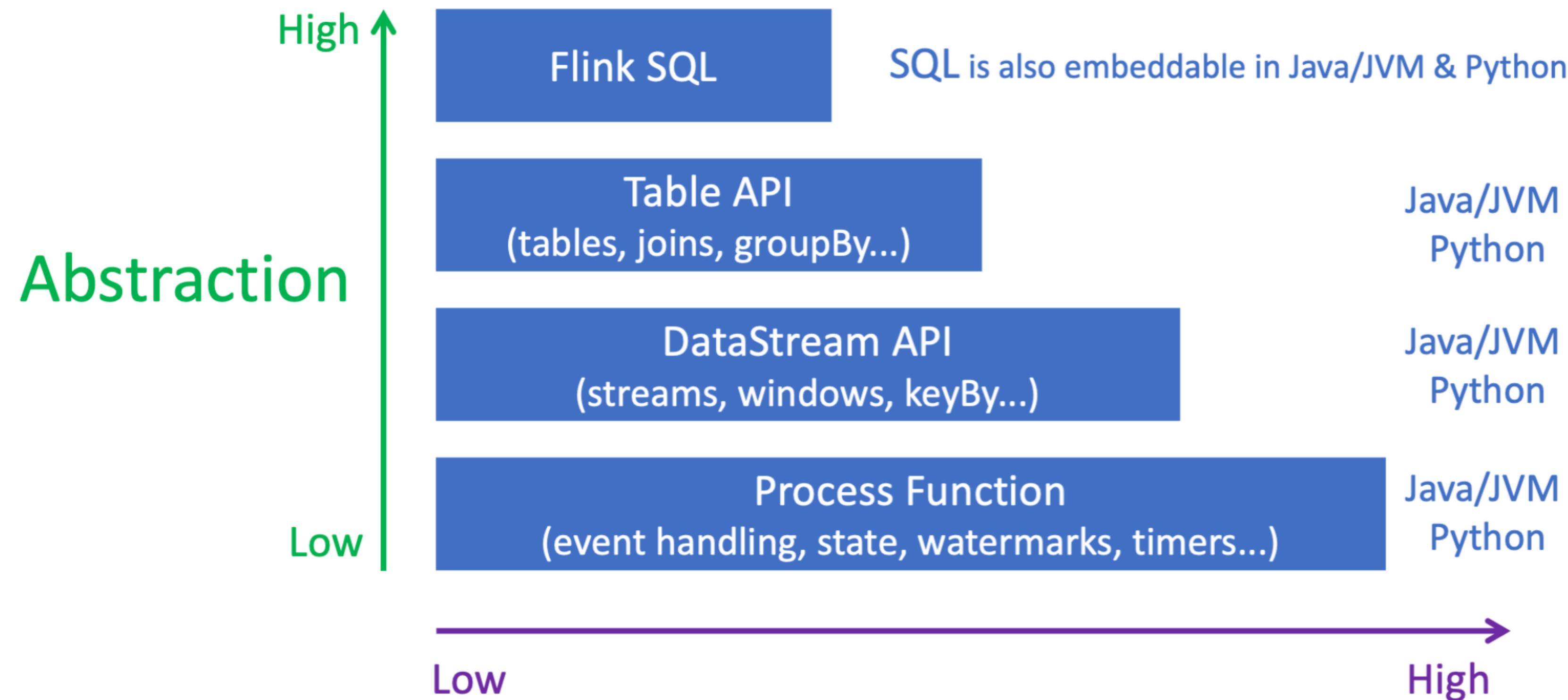
7 . **SQL users benefit:** Deploy long-running applications from Studio



Apache Flink APIs



Apache Flink APIs



Programming languages

2. 🚢 Abstraction level recommendation

- ⬆️ Start with higher levels when possible

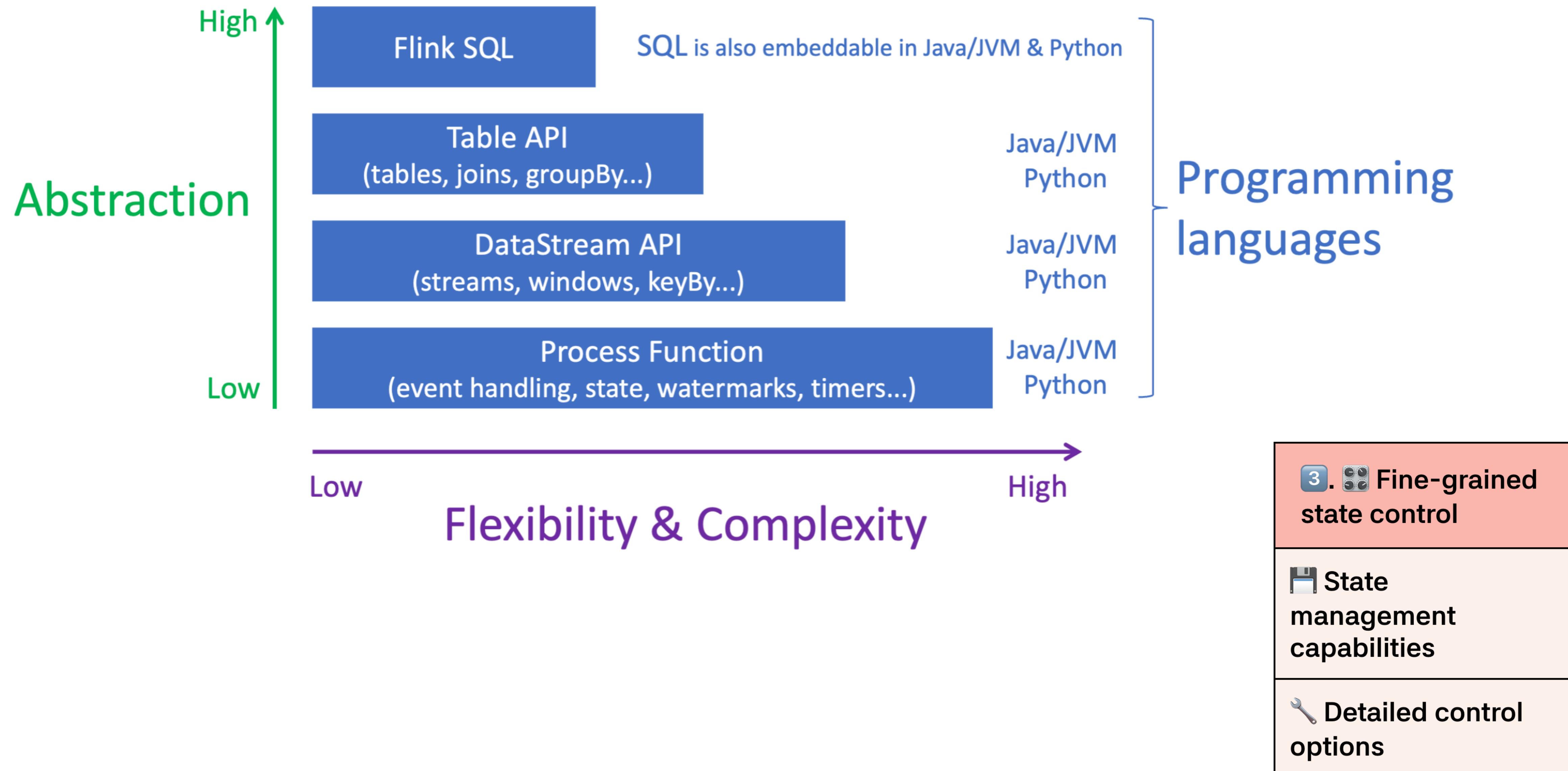
⭐ Some features exclusive to DataStream API

☕ Java support

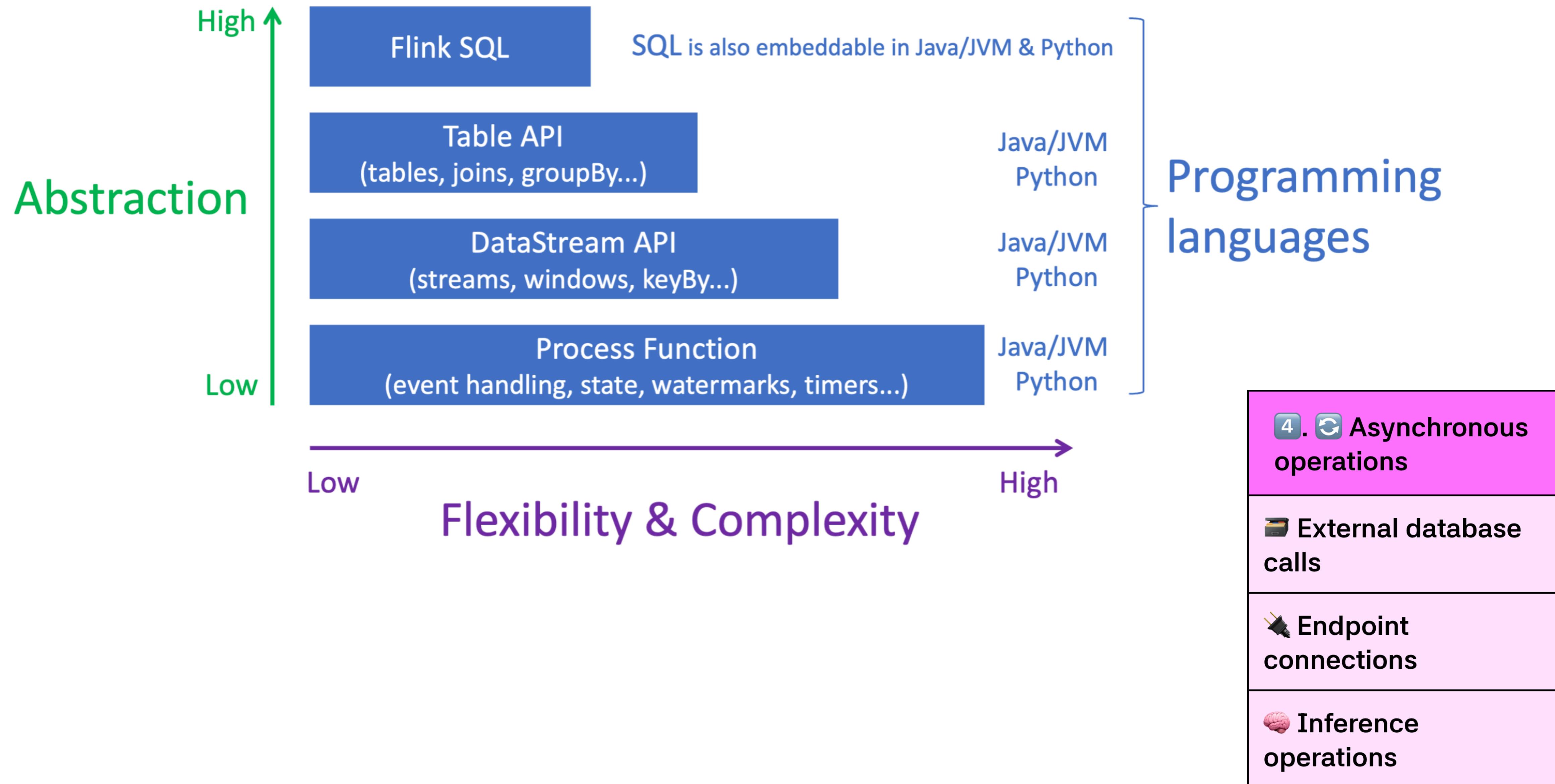
🐍 Python support

📜 Scala support

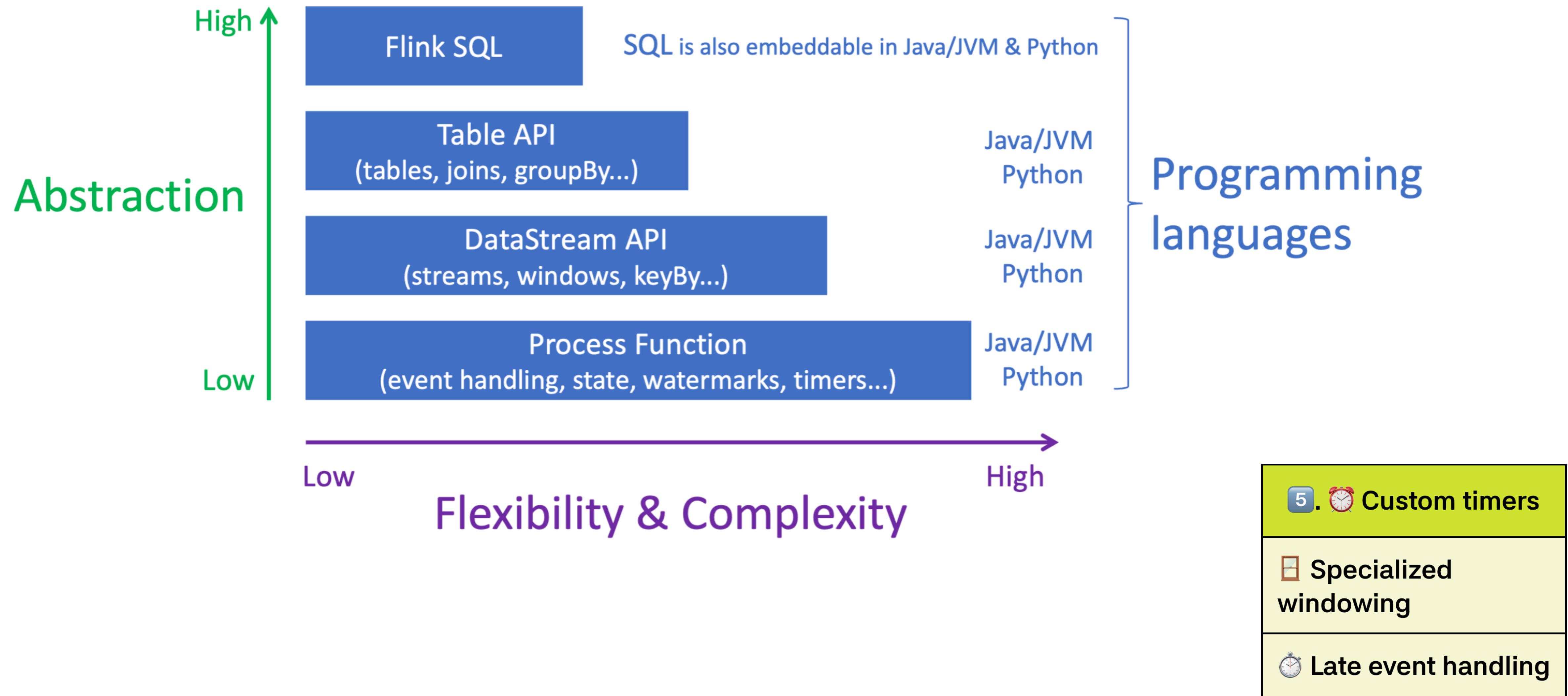
Apache Flink APIs



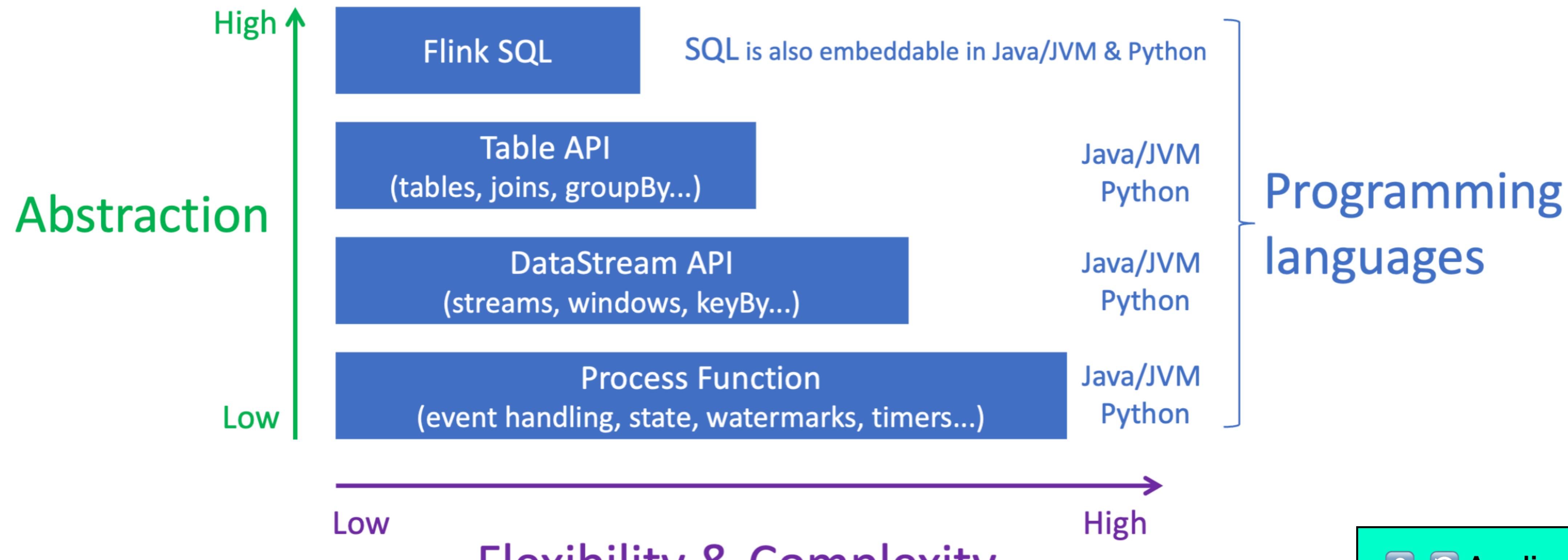
Apache Flink APIs



Apache Flink APIs

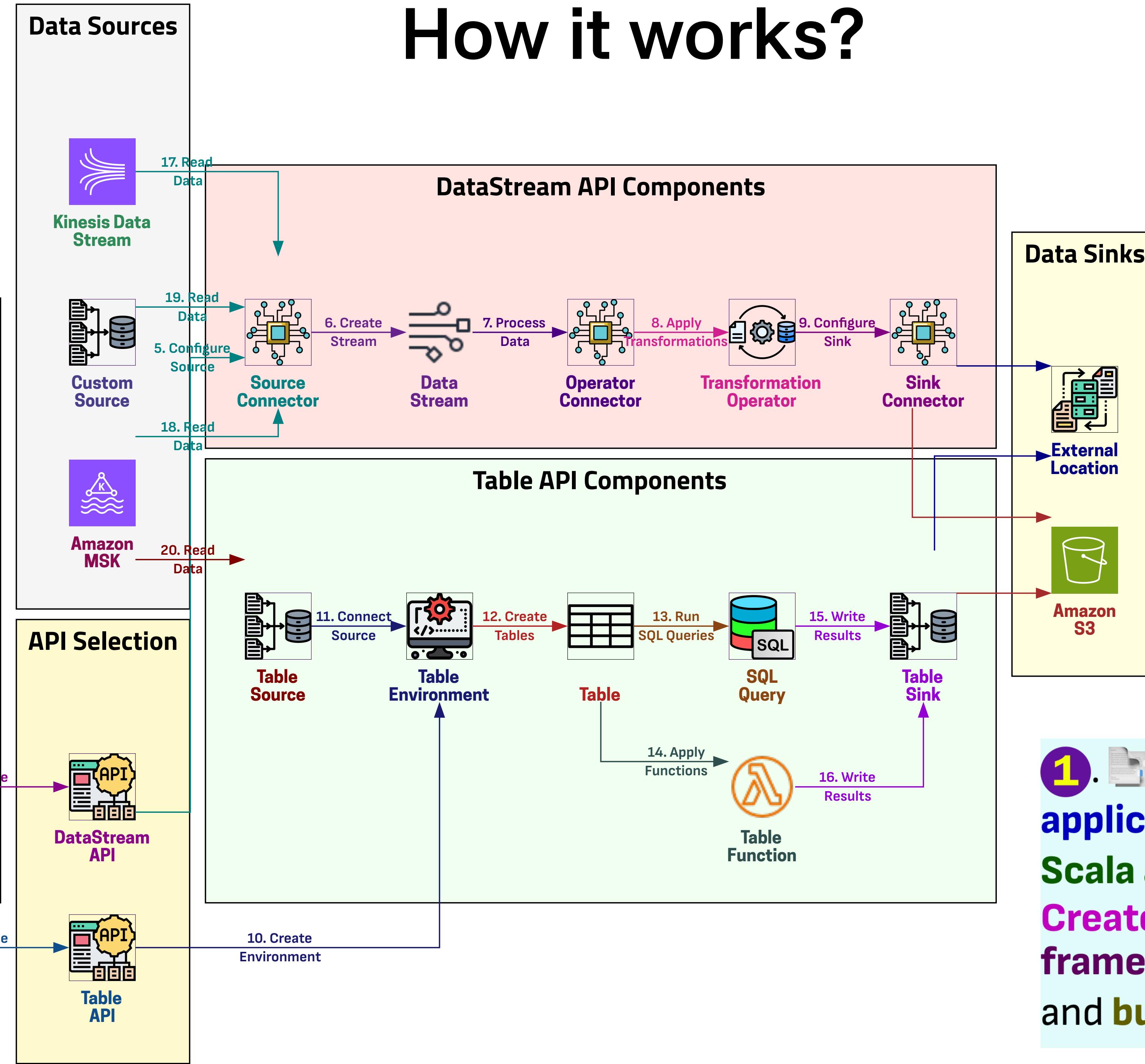
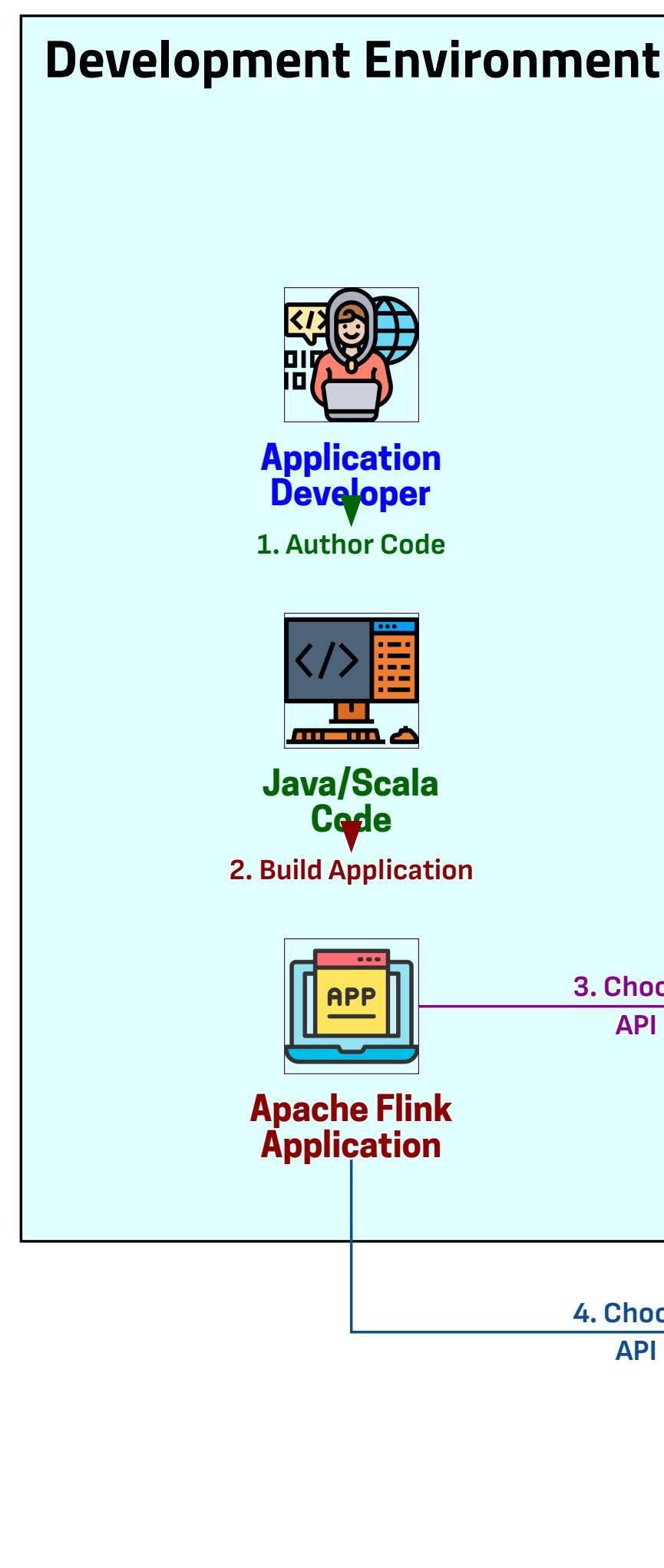


Apache Flink APIs



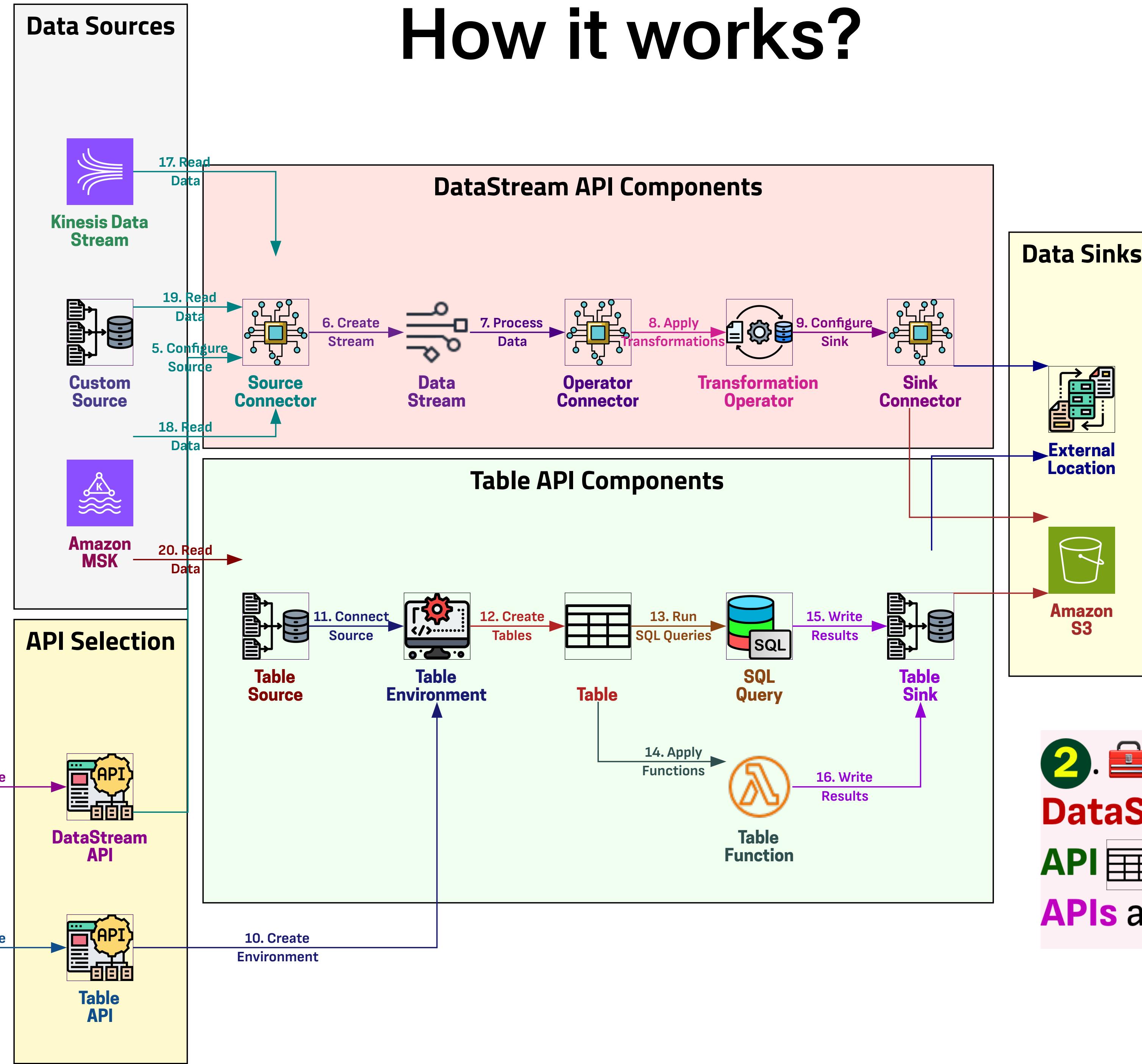
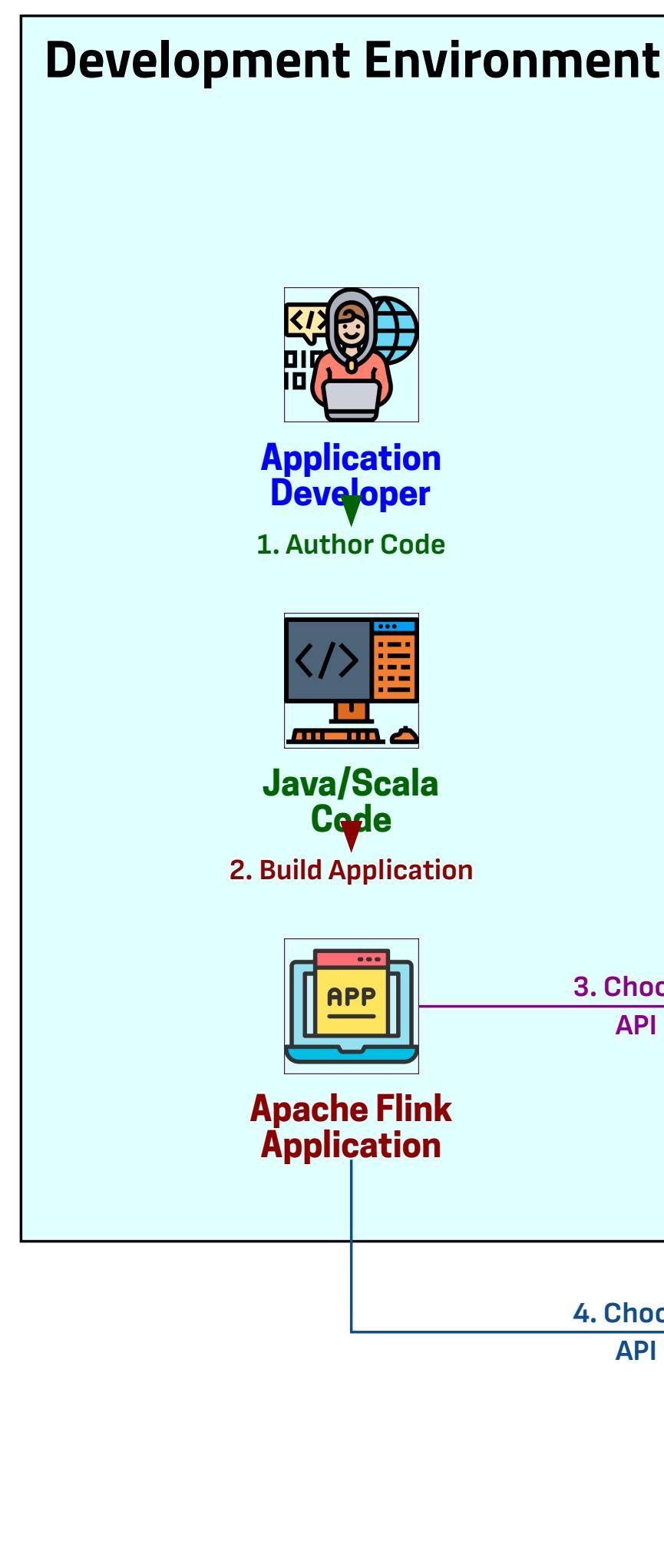
- 6. ⚡ Application flow modification
- Without resetting state
- Greater flexibility

How it works?

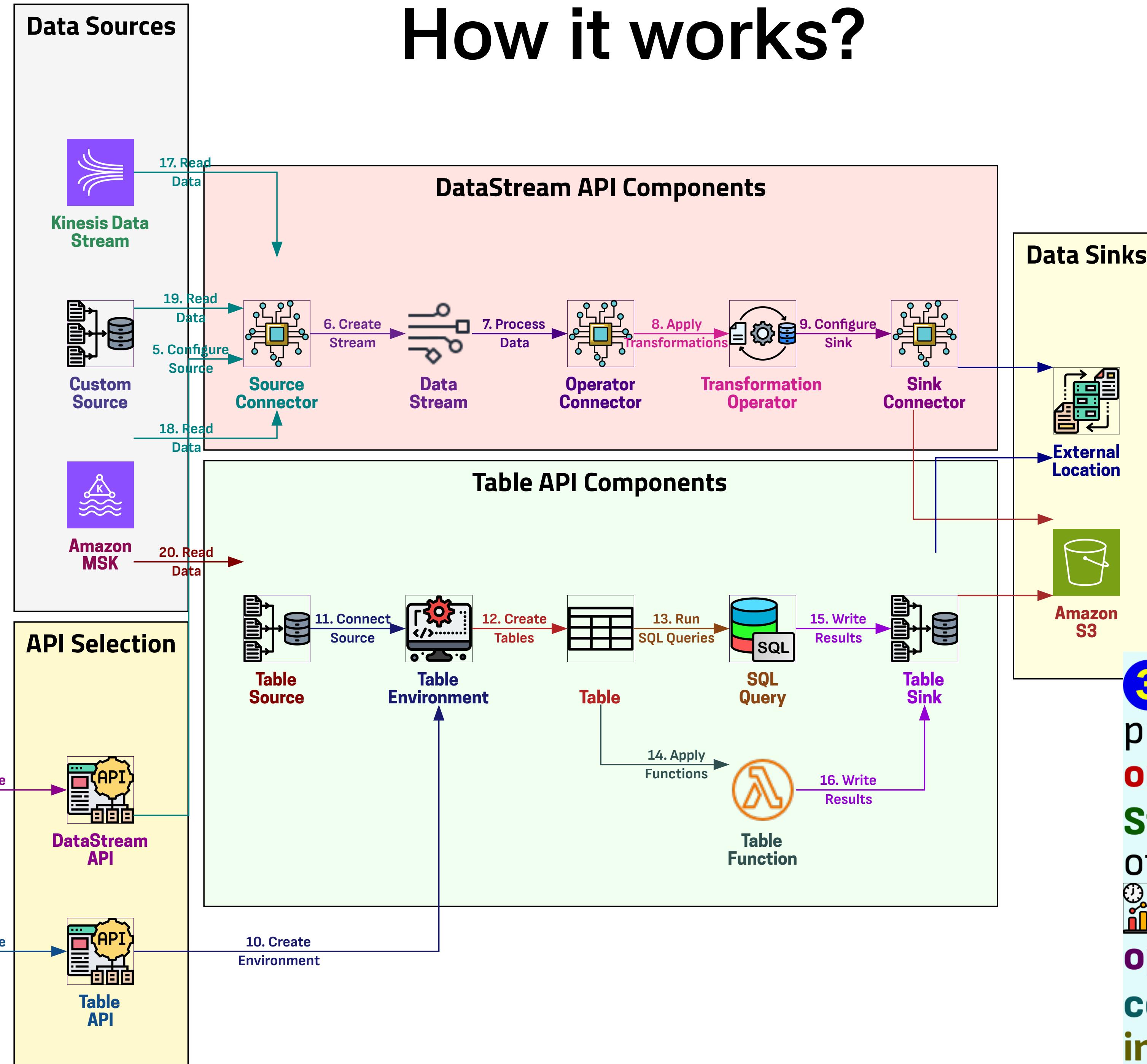
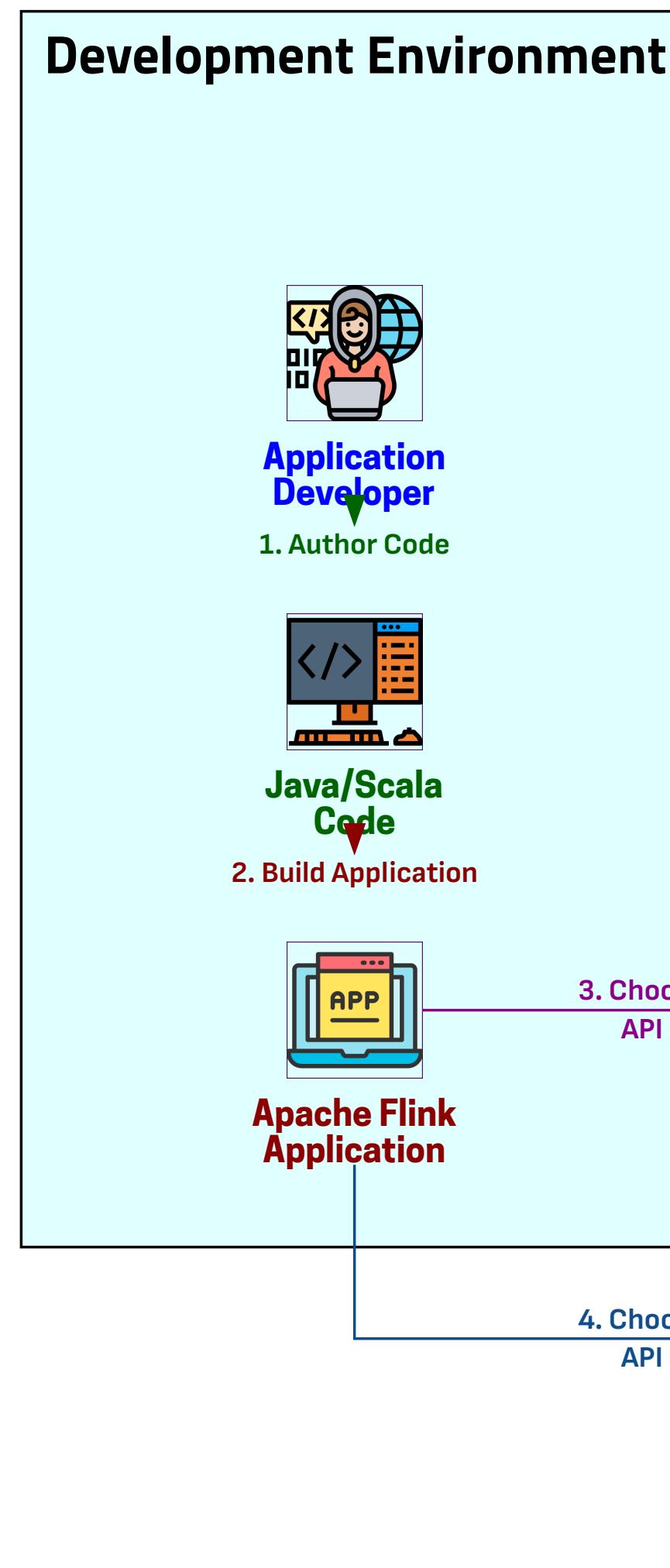


1. Apache Flink applications: Java or Scala applications , Created with Apache Flink framework , Authored and built locally

How it works?

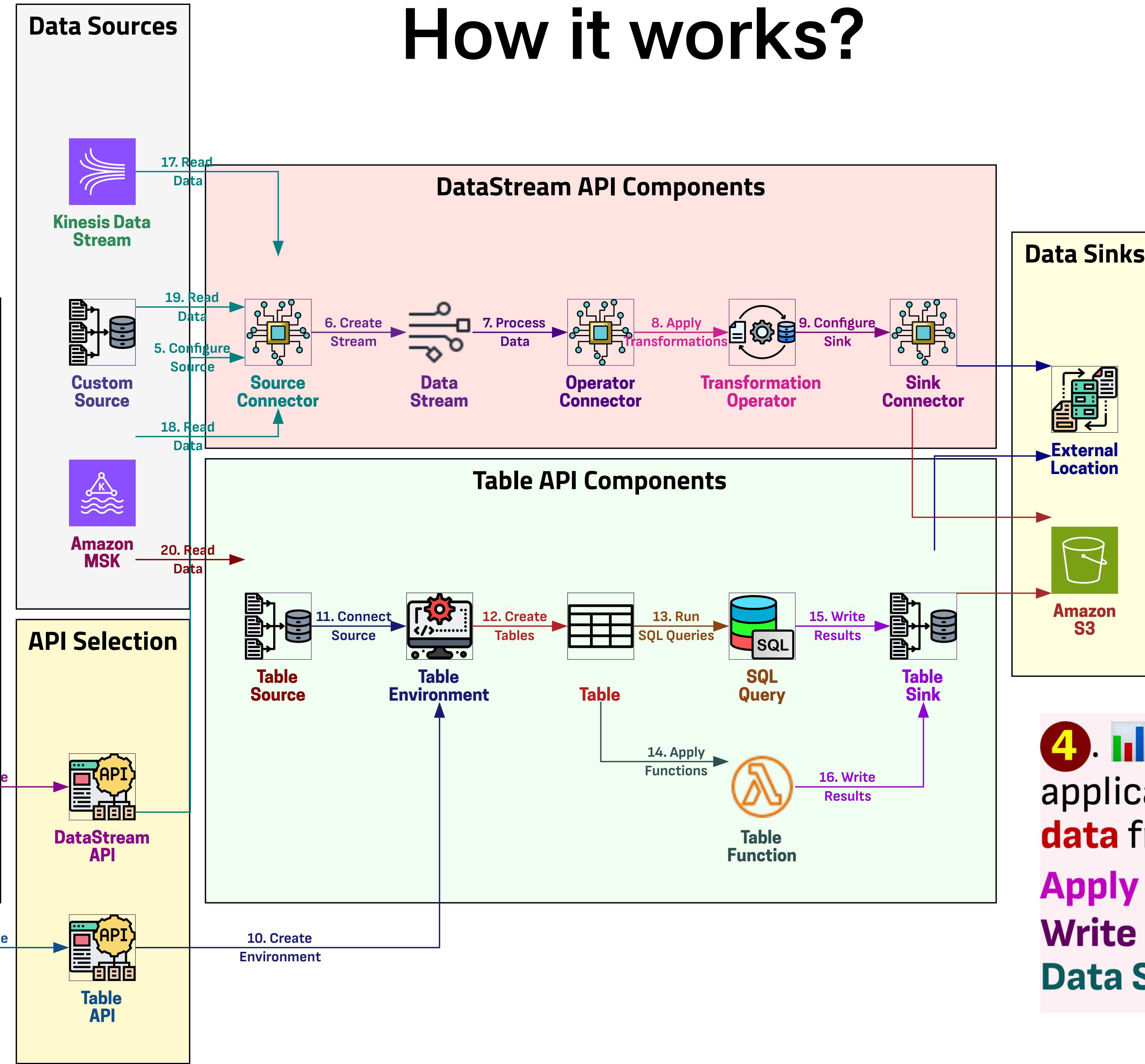
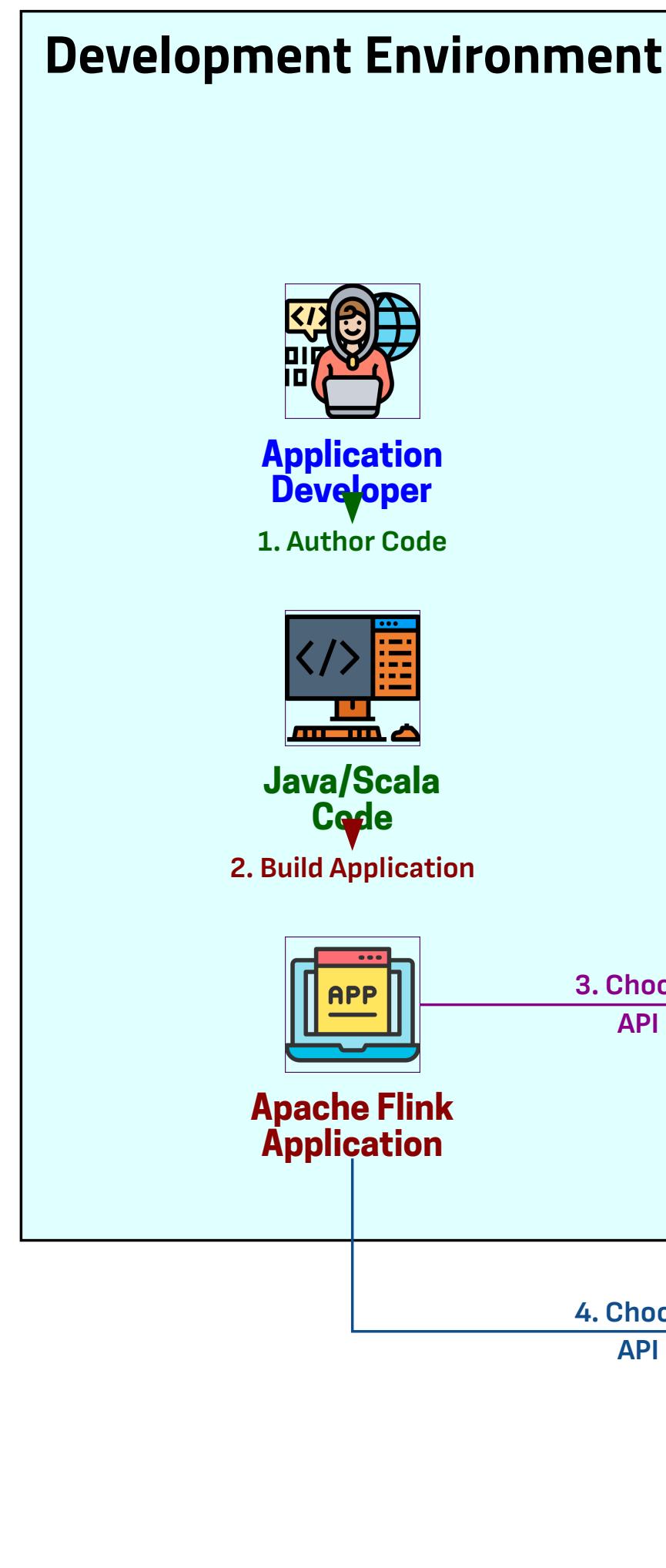


How it works?

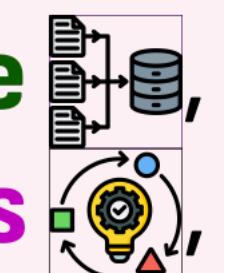


3. DataStream API
 programming model: **Based on data streams**, **Structured representation of continuous data flow**, **Transformation operators**, **Processing components** transforming **input streams**

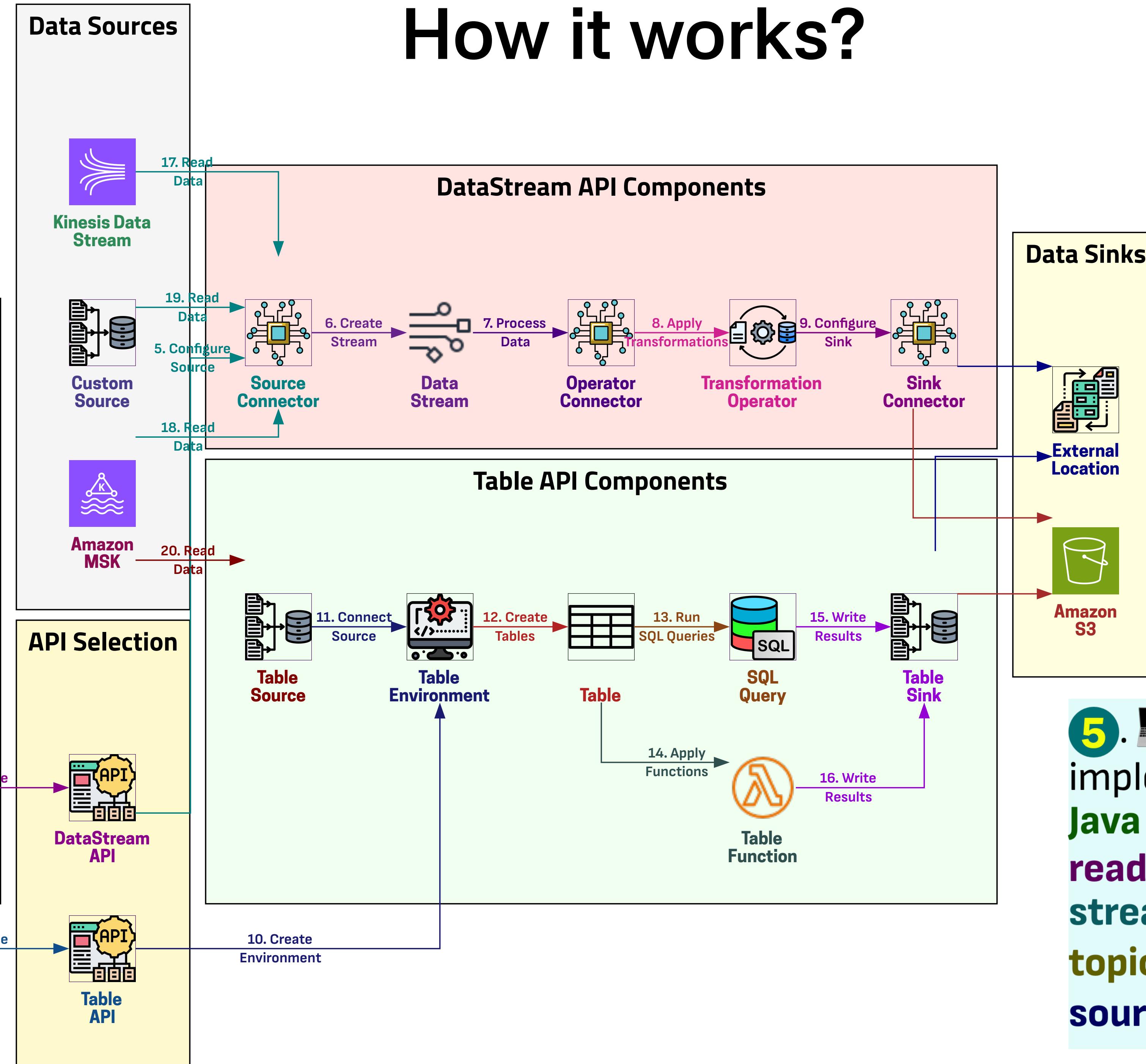
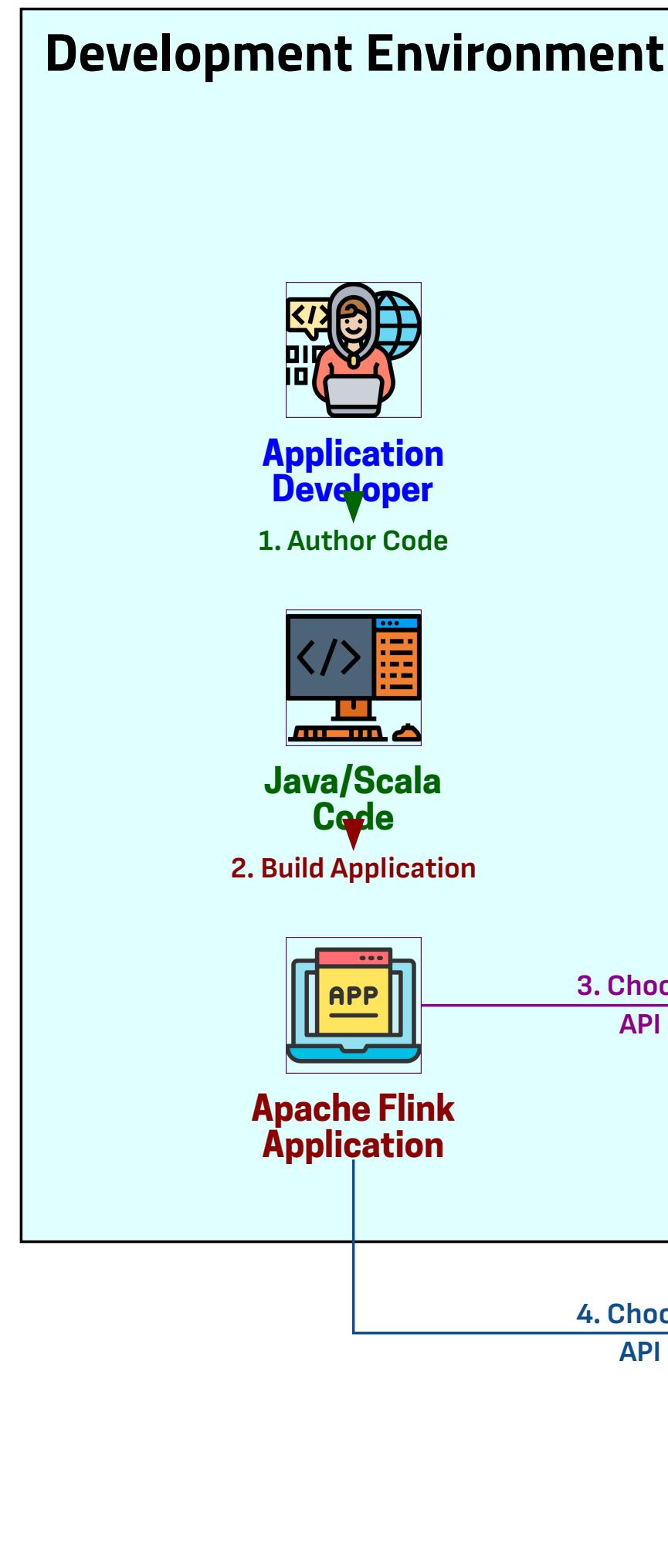
How it works?



4. DataStream API
 application workflow: **Read data from Data Source**, **Apply transformations**, **Write transformed data to Data Sink**

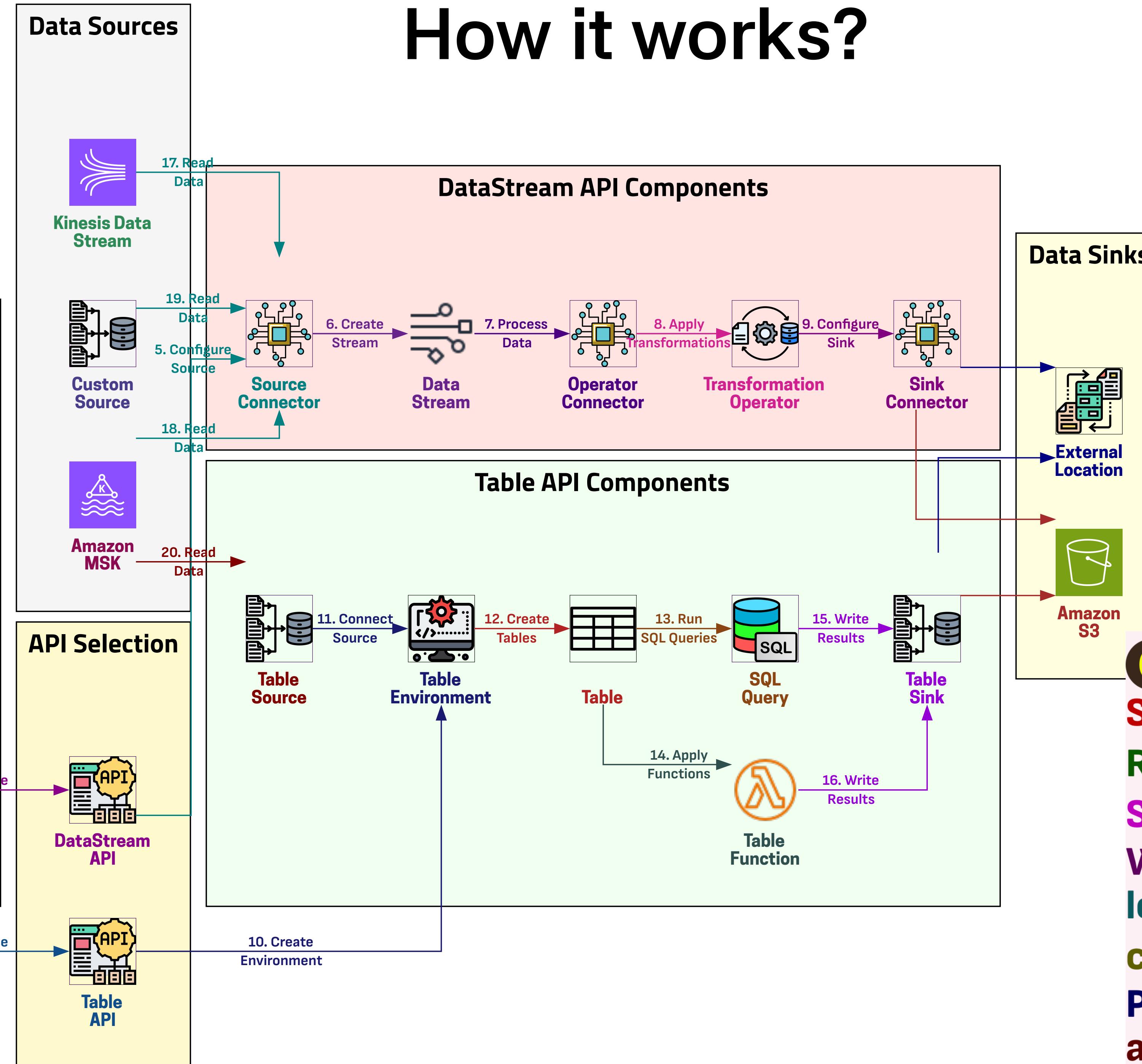
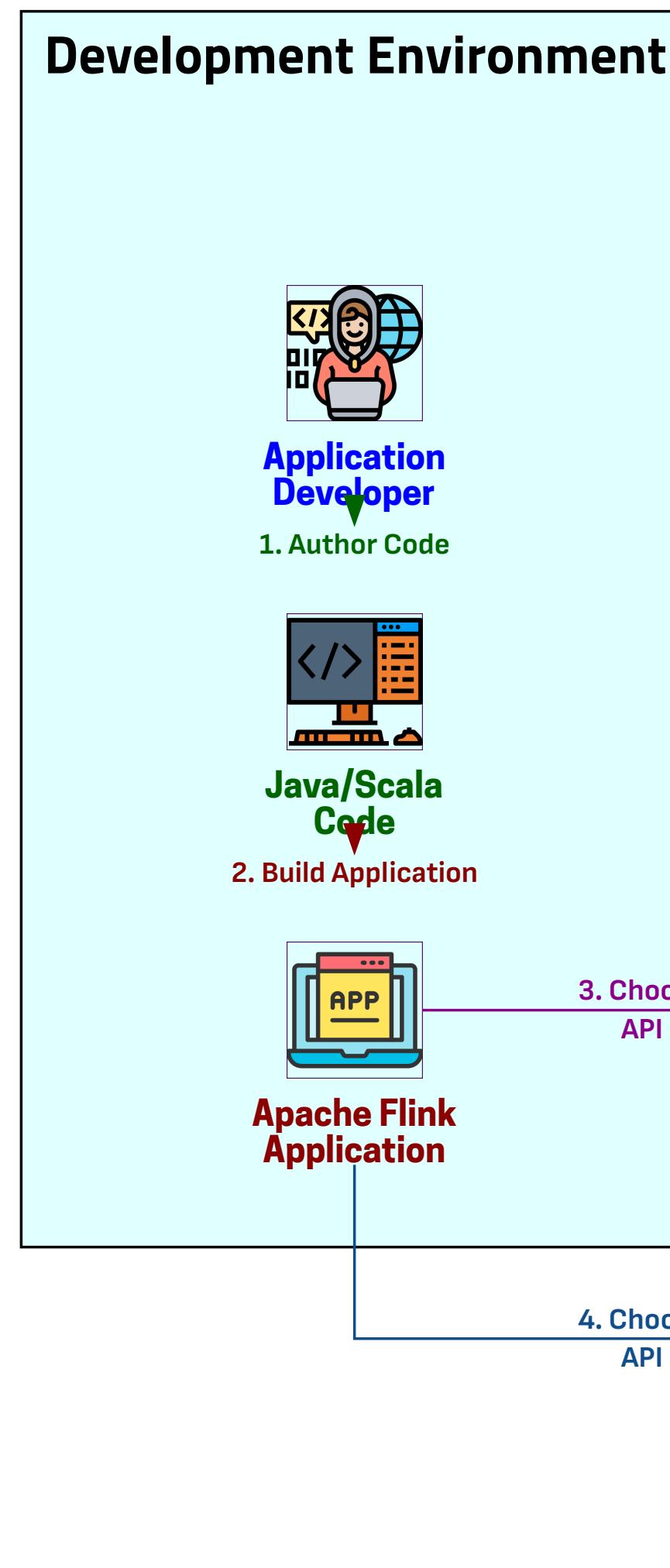


How it works?



5. **DataStream API**
implementation: **Written in Java** or **Scala** , Can read from: **Kinesis data streams** , **Amazon MSK topics** , **Custom sources**

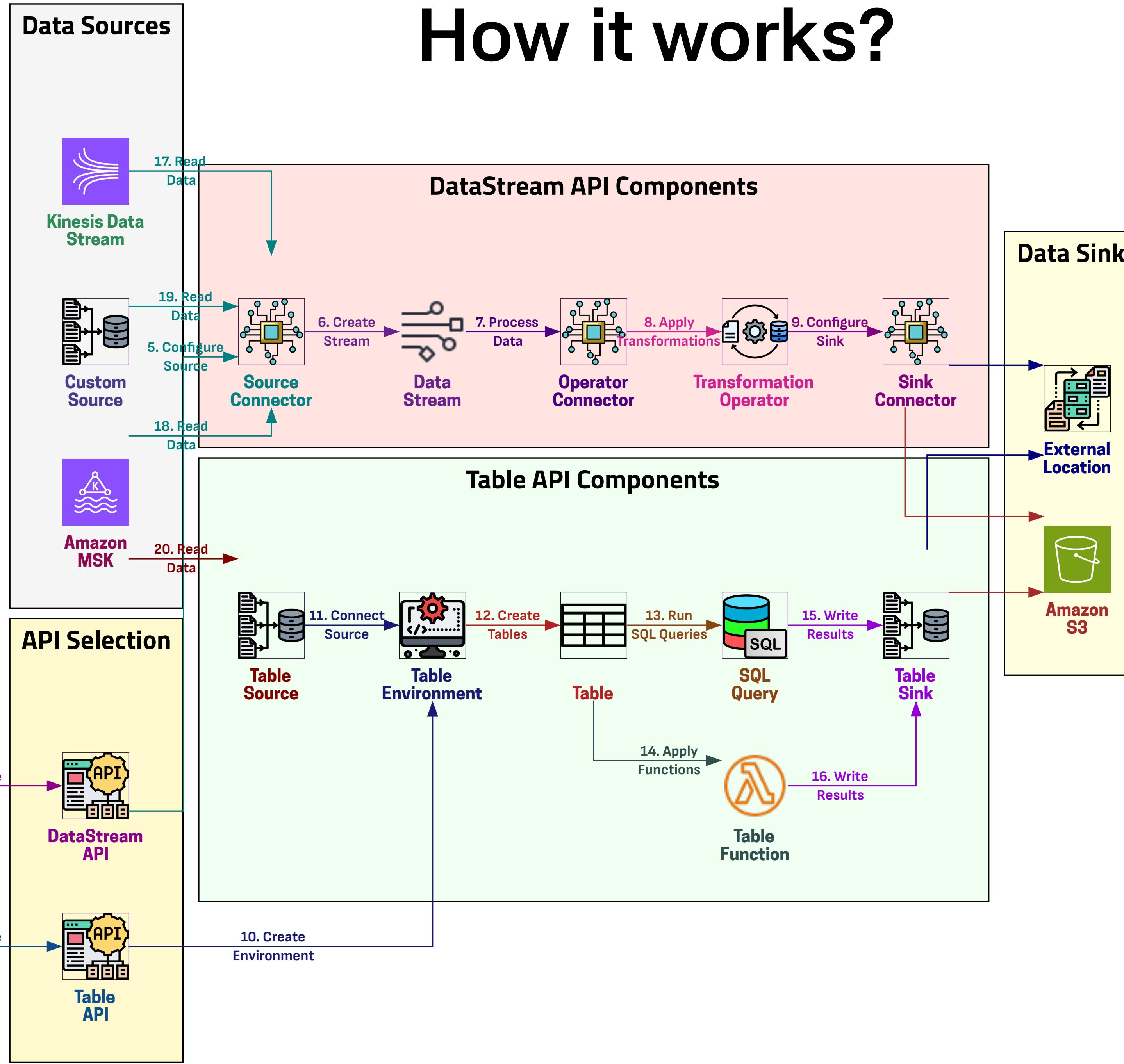
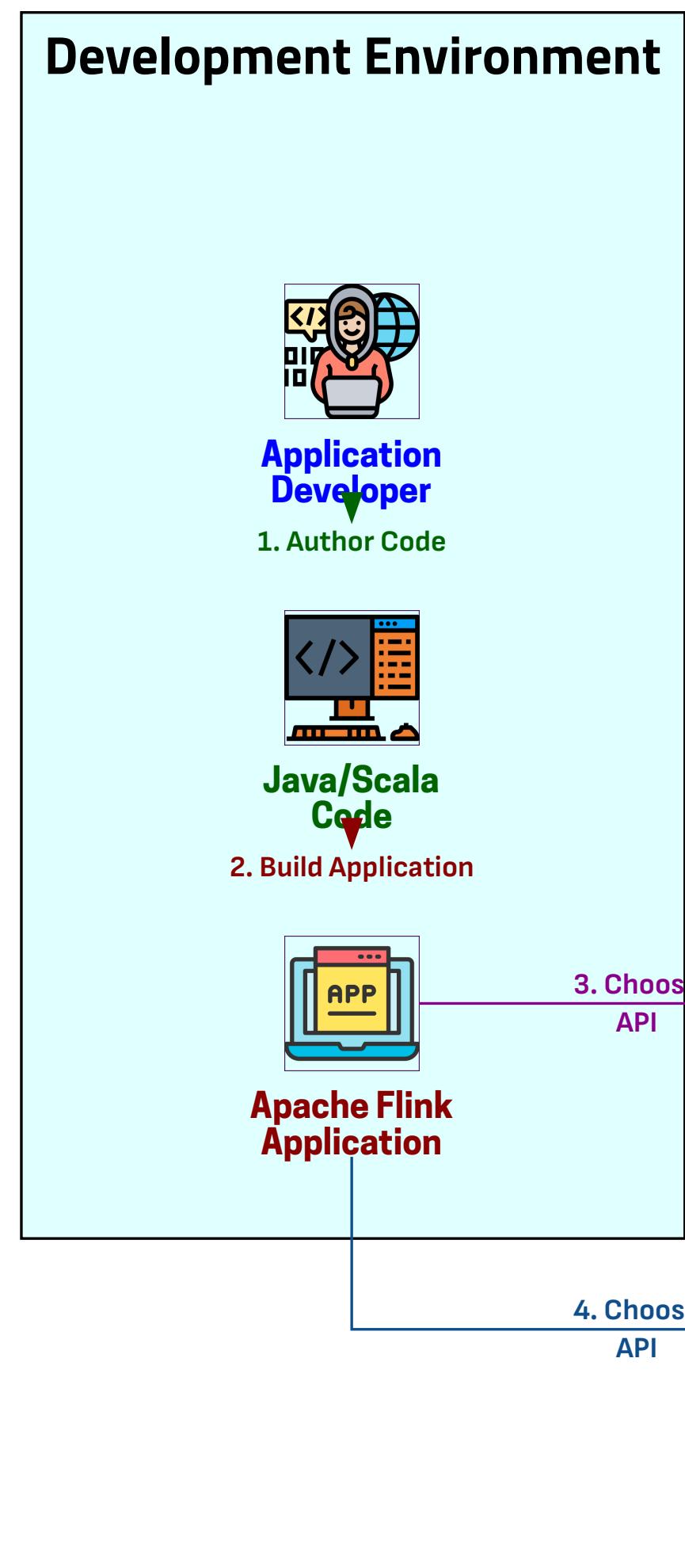
How it works?



6. Connector types:

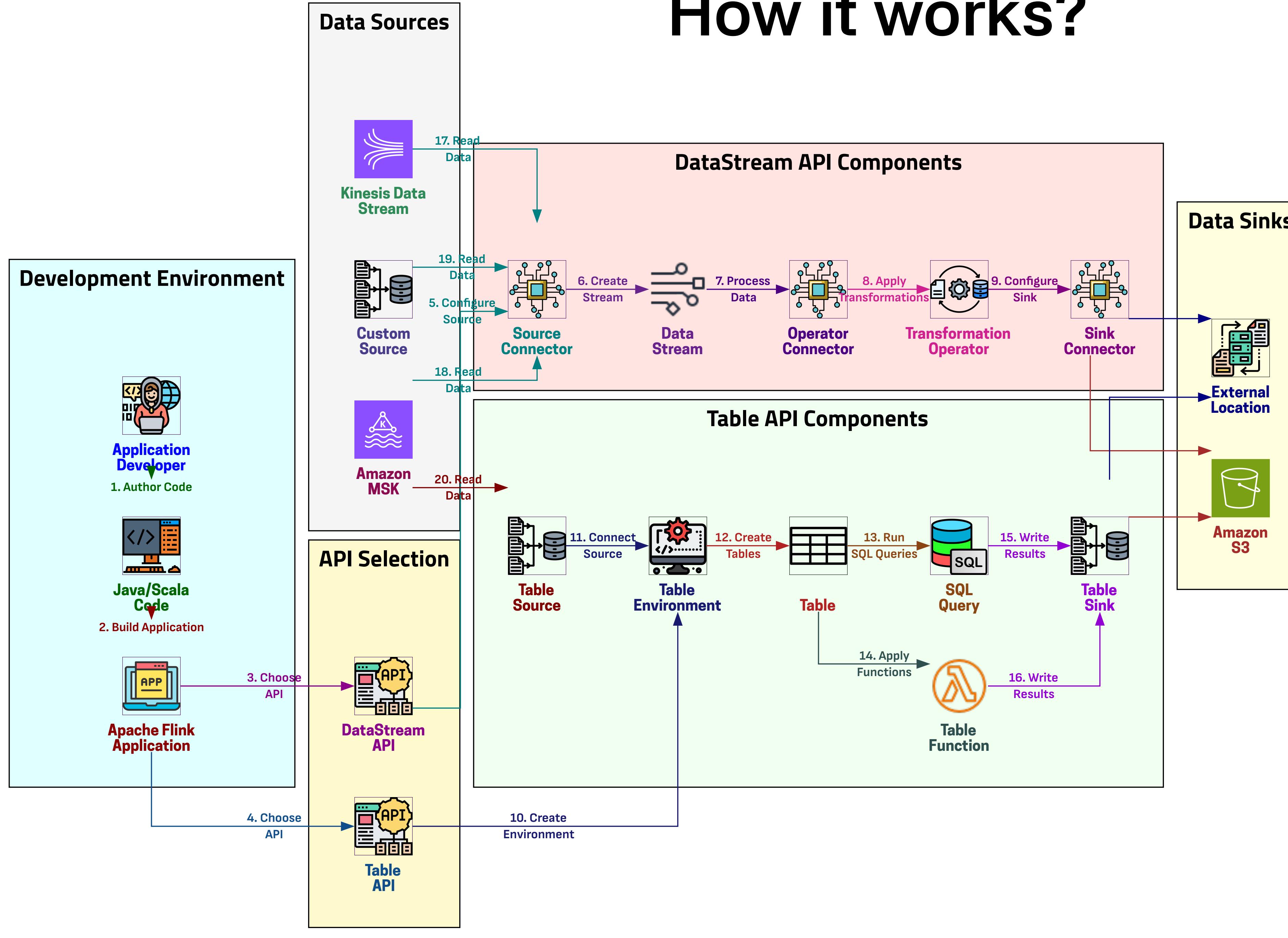
- Source connectors**: Reading external data
- Sink connectors**: Writing to external locations , Operator connectors
- Processor connectors**: Processing data within application

How it works?



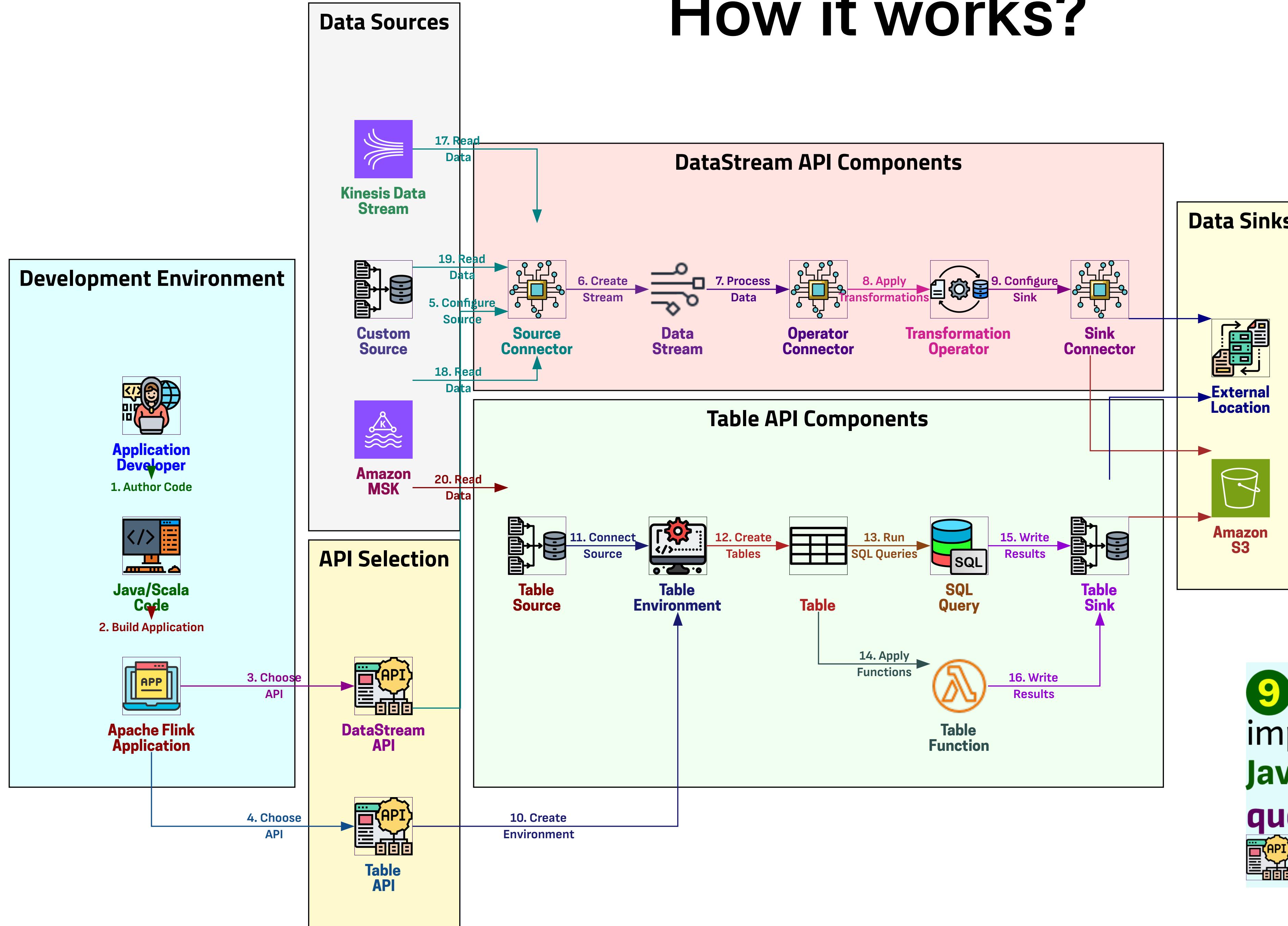
7. Table API
 programming model components: **Table Environment**: Interface to underlying data, **Tables**: Objects providing access to SQL tables/views, **Table Sources**: Reading from external sources, **Table Functions**: Transforming data, **Table Sinks**: Writing to external locations

How it works?



8. Table API
application workflow:
Create TableEnvironment
Create tables :
Using SQL queries ,
Using API functions ,
Run queries on tables ,
Apply transformations on results , **Write results**
to Table Sink

How it works?



9. Table API
 implementation: **Written in Java** or **Scala** , Can query data using: **API calls** , **SQL queries**

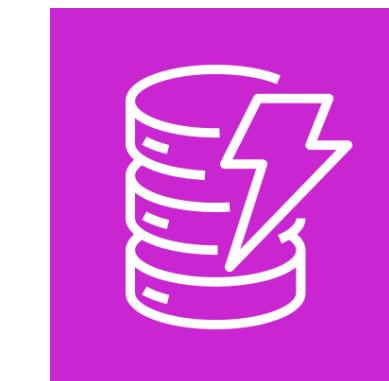
Connectors (Sources & Sinks)



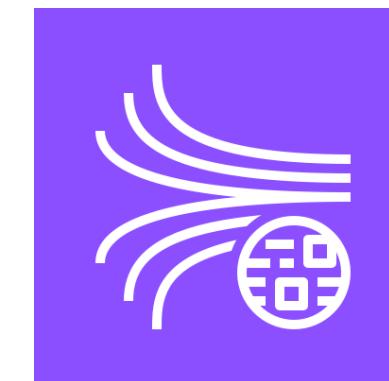
1. Apache Kafka
(source/sink)



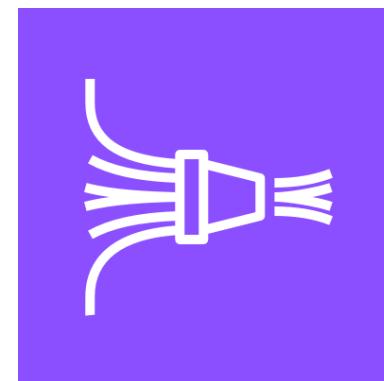
2. Apache Cassandra
(sink)



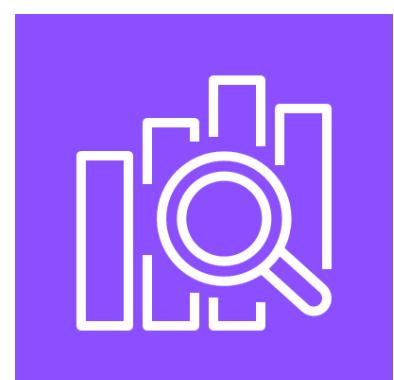
3. Amazon DynamoDB
(sink)



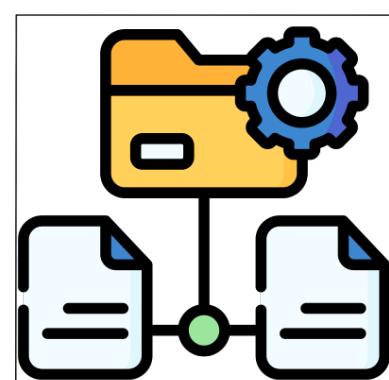
4. Amazon Kinesis Data
Streams (source/sink)



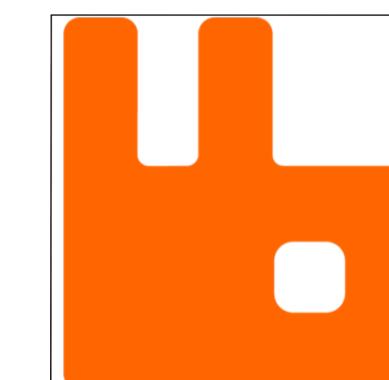
5. Amazon Kinesis Data
Firehose (sink)



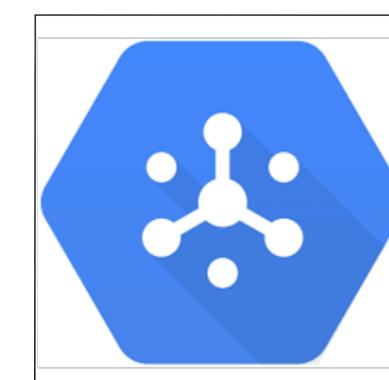
6. Elasticsearch
(sink)



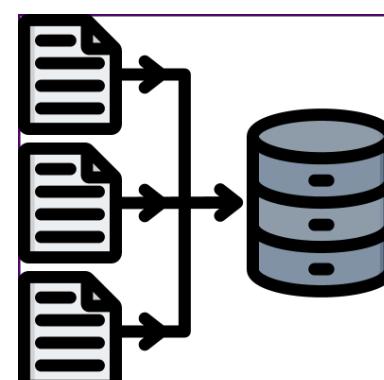
7. FileSystem
(sink)



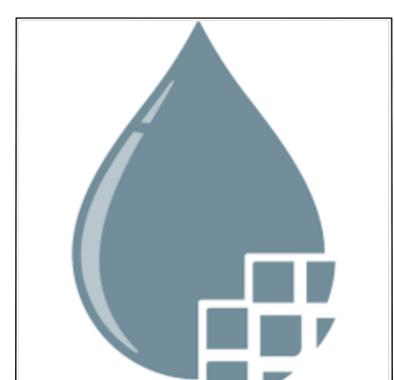
8. RabbitMQ
(source/sink)



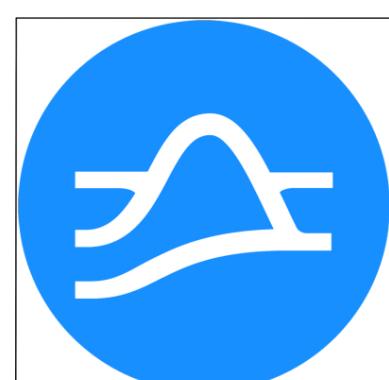
9. Google PubSub
(source/sink)



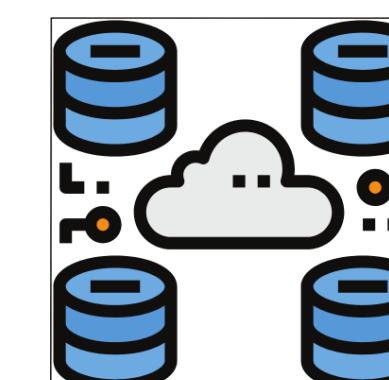
10. Hybrid Source
(source)



11. Apache NiFi
(source/sink)



12. Apache Pulsar
(source)



13. JDBC
(sink)

Use Kinesis data streams

```
1 // Configure the KinesisStreamsSource
2 Configuration sourceConfig = new Configuration();
3 sourceConfig.set(KinesisSourceConfigOptions.STREAM_INITIAL_POSITION, KinesisSourceConfigOptions.InitialPosition.TRIM_HORIZON); ①
4
5 // Create a new KinesisStreamsSource to read from specified Kinesis Stream.
6 KinesisStreamsSource<String> kdsSource =
7     KinesisStreamsSource.<String>builder() ②
8         .setStreamArn("arn:aws:kinesis:us-east-1:123456789012:stream/test-stream") ③
9         .setSourceConfig(sourceConfig) ④
10        .setDeserializationSchema(new SimpleStringSchema()) ⑤
11        .setKinesisShardAssigner(ShardAssignerFactory.uniformShardAssigner()) ⑥
12        .build(); ⑦
```

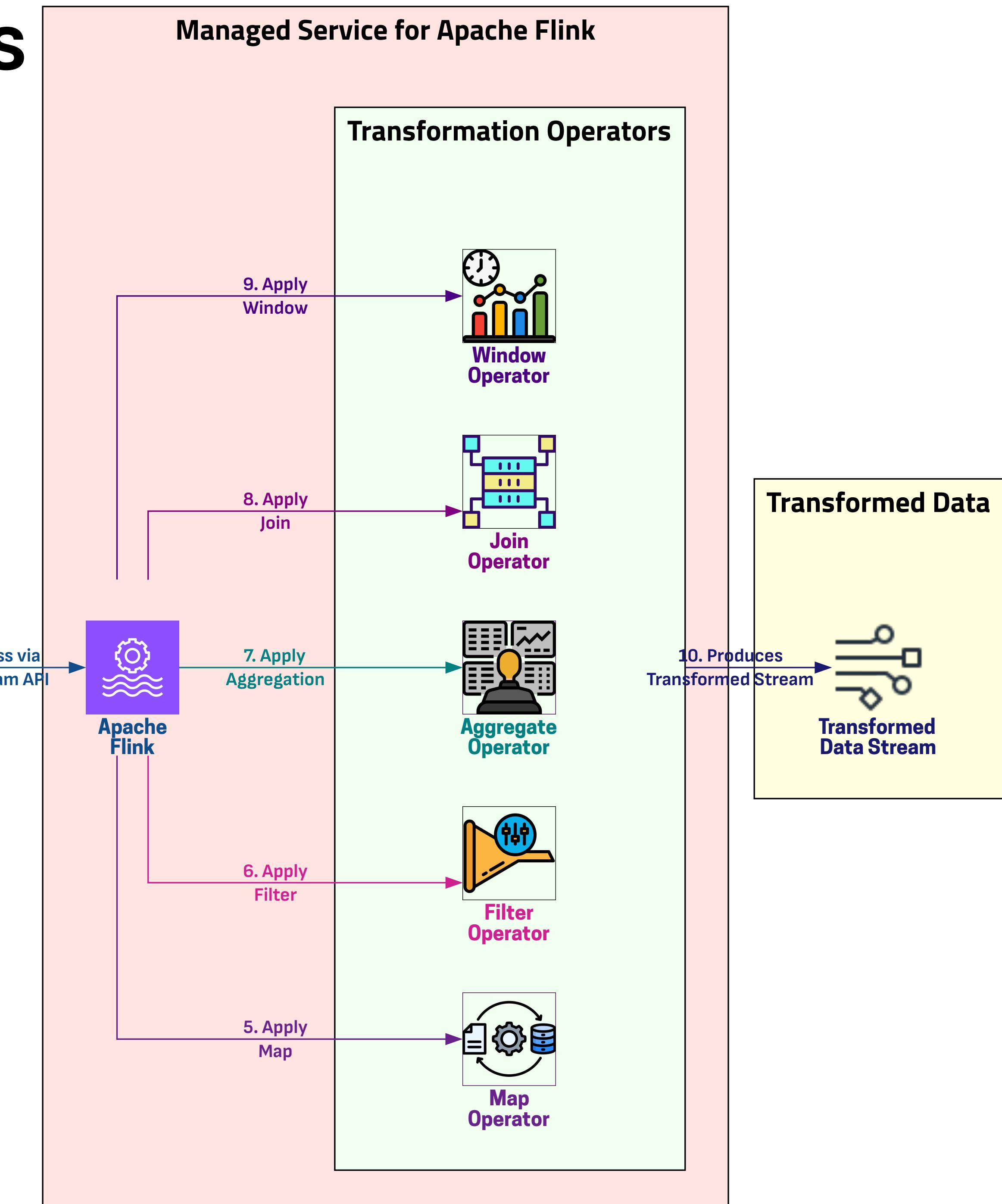
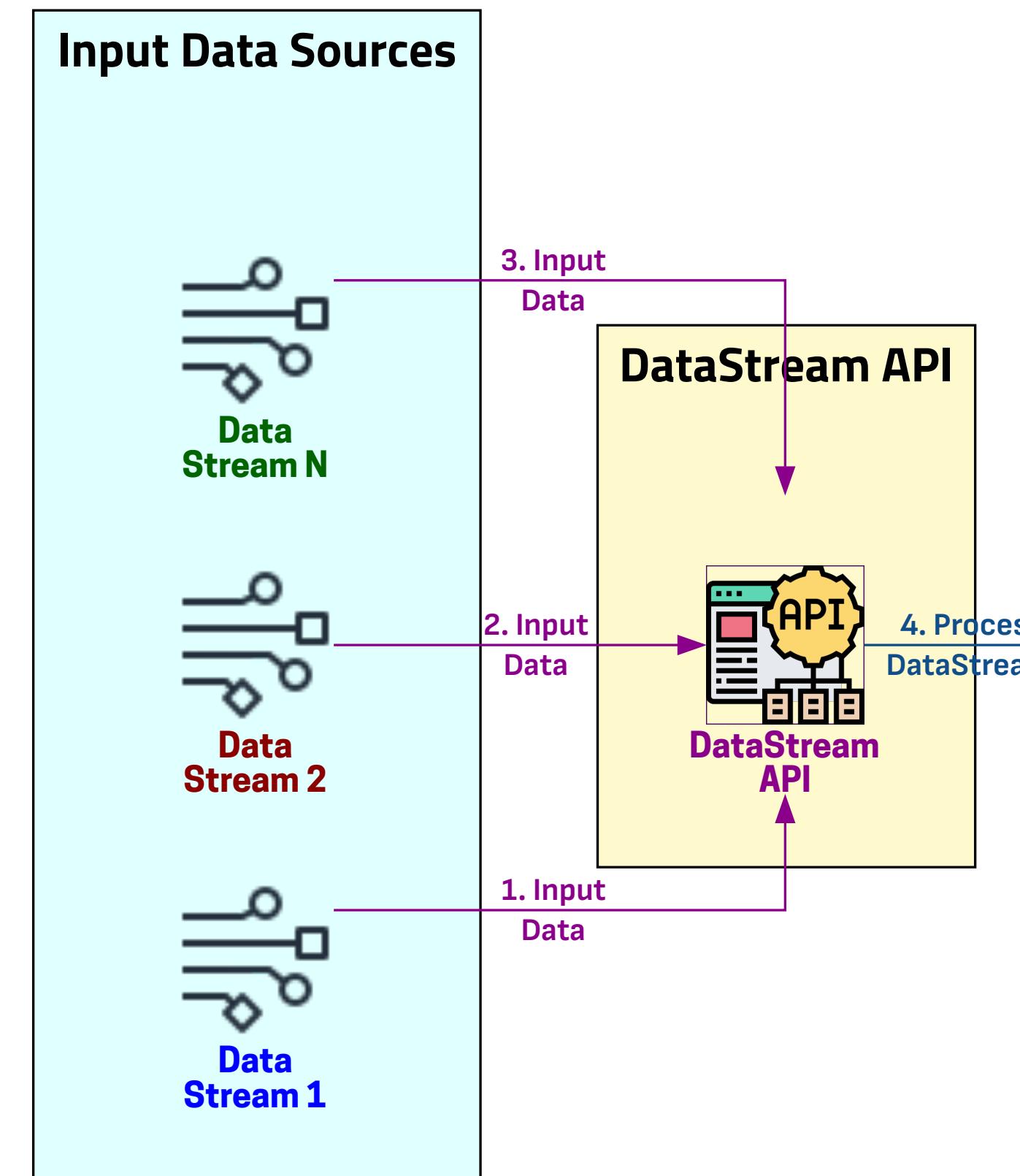
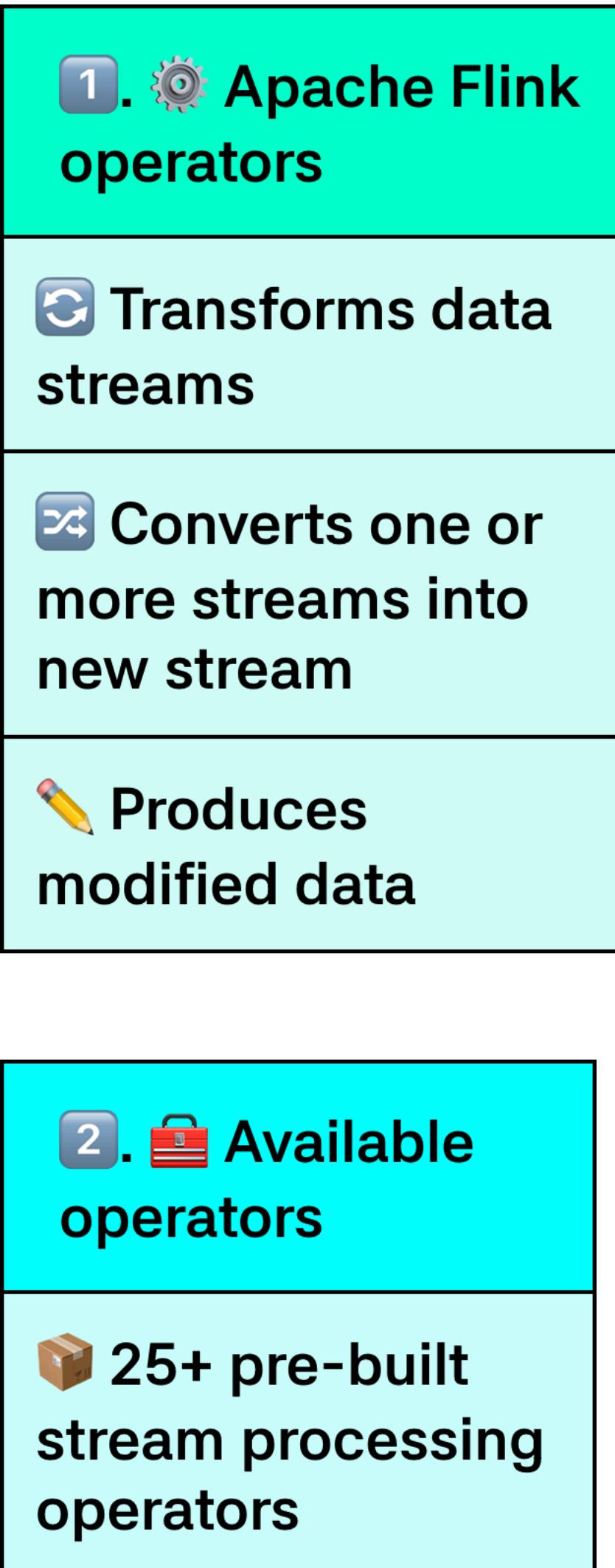
- ➊ Set the initial position for consuming the Kinesis stream to `TRIM_HORIZON`, which means that the consumer will start reading from the oldest record in the stream. This is optional, and if not specified, the default is `LATEST`, which means reading only new records that arrive after the consumer starts.
- ➋ Initialize the builder pattern for creating a `KinesisStreamsSource` with a generic type parameter of `String`, indicating that the records will be deserialized into `String` objects.
- ➌ Specify the ARN (Amazon Resource Name) of the Kinesis stream to read from. This uniquely identifies the stream within AWS and includes the AWS account ID, region, and stream name.
- ➍ Apply the configuration settings defined earlier to the Kinesis stream source, which includes the initial position setting.
- ➎ Set the deserialization schema to `SimpleStringSchema`, which will convert the raw Kinesis records into `String` objects. Different schemas can be used depending on the data format.
- ➏ Configure the shard assigner strategy that determines how Kinesis shards are distributed among parallel subtasks. The `uniformShardAssigner` distributes shards evenly. This setting is optional and defaults to uniform assignment if not specified.
- ➐ Finalize the builder pattern and create the `KinesisStreamsSource` instance with all the configured settings.

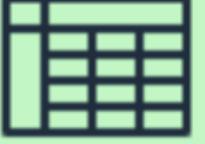
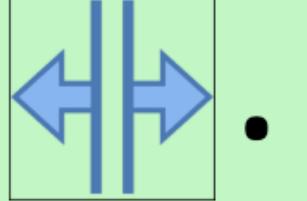
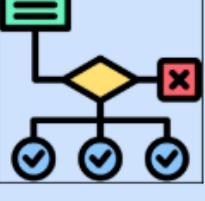
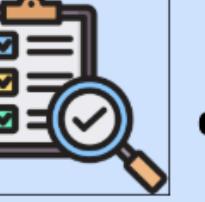
Create a KafkaSource

```
1 KafkaSource<String> source = KafkaSource.<String>builder()  
2     .setBootstrapServers(brokers)    ①  
3     .setTopics("input-topic")        ②  
4     .setGroupId("my-group")         ③  
5     .setStartingOffsets(OffsetsInitializer.earliest()) ④  
6     .setValueOnlyDeserializer(new SimpleStringSchema()) ⑤  
7     .build();  
8  
9 env.fromSource(source, WatermarkStrategy.noWatermarks(), "Kafka Source"); ⑥
```

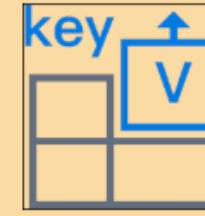
- ① Specifies the Kafka broker addresses to connect to, using the `brokers` variable which should contain a comma-separated list of host:port pairs
- ② Defines the Kafka topic to consume messages from, in this case "input-topic"
- ③ Sets the consumer group ID to "my-group", which is used by Kafka to track the consumption progress
- ④ Configures the consumer to start reading from the earliest available offset in the topic, ensuring no messages are missed
- ⑤ Specifies a deserializer that processes only the value portion of Kafka records and converts them to String objects
- ⑥ Adds the configured Kafka source to the Flink execution environment, with no watermark strategy for event time processing, and assigns the name "Kafka Source" for monitoring and visualization

Transform data using Operators



#	Operator Name	Description
1	Map	Takes one element  and produces one element  .
2	FlatMap	Takes one element  and produces zero, one, or more elements  .
3	Filter	Evaluates a boolean function  for each element and retains those for which the function returns true  .

4



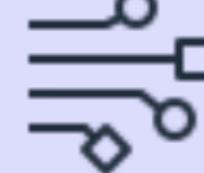
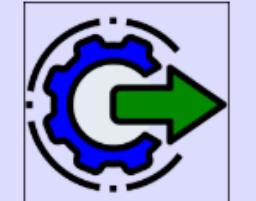
KeyBy

Logically partitions a stream into disjoint partitions. All records with the same key  are assigned to the same partition.

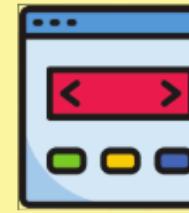
5



Reduce

A "rolling" reduce on a keyed data stream . Combines the current element with the last reduced value and emits the new value .

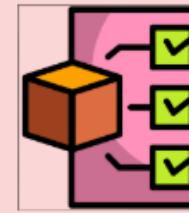
6



Window

Groups data in each key according to some characteristic (e.g., data that arrived within the last 5 seconds) on already partitioned KeyedStreams.

7



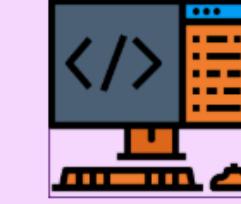
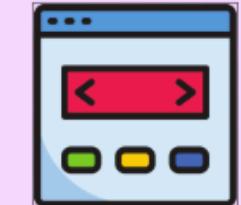
WindowAll

Groups all stream events according to some characteristic (e.g., data that arrived within the last 5 seconds) on regular DataStreams. Often a non-parallel transformation.

8



Window Apply

Applies a general function  to the window as a whole .

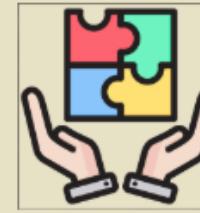
9



WindowReduce

Applies a functional reduce function  to the window and returns the reduced value .

10



Union

Merges two or more data streams .

Use transform operators

```
1 DataStream<ObjectNode> output = input.map( ①
2     new MapFunction<ObjectNode, ObjectNode>() { ②
3         @Override
4         public ObjectNode map(ObjectNode value) throws Exception { ③
5             return value.put("TICKER", value.get("TICKER").asText() + " Company"); ④
6         }
7     }
8 );
```

- ➊ Creates a new transformed `DataStream` by applying a `map` operation to the `input` stream, which contains `ObjectNode` elements
- ➋ Instantiates an anonymous class implementing the `MapFunction` interface, specifying both input and output types as `ObjectNode`
- ➌ Implements the required `map` method that will be applied to each element in the `input` stream
- ➍ Modifies the incoming `ObjectNode` by appending " Company" to the value of the "TICKER" field and returns the modified node

Use aggregation operators

```
1 DataStream<ObjectNode> output = input.keyBy(node -> node.get("TICKER").asText()) ①
2   .window(TumblingProcessingTimeWindows.of(Time.seconds(5))) ②
3   .reduce((node1, node2) -> { ③
4     double priceTotal = node1.get("PRICE").asDouble() + node2.get("PRICE").asDouble(); ④
5     node1.replace("PRICE", JsonNodeFactory.instance.numberNode(priceTotal)); ⑤
6     return node1; ⑥
7   });

```

- ① Partitions the input `DataStream` by the value of the "TICKER" field, ensuring that all records with the same ticker symbol are processed by the same parallel task
- ② Applies a tumbling window of 5 seconds in processing time, grouping records that arrive within the same 5-second interval
- ③ Defines a reduce function that will be applied to all records within each window for each key
- ④ Calculates the sum of the "PRICE" values from two `ObjectNode` records by extracting the double values
- ⑤ Updates the first node by replacing its "PRICE" field with the newly calculated total using `JsonNodeFactory` to create a new number node
- ⑥ Returns the modified first node as the result of the reduction, accumulating values as the reducer processes more elements in the window



**Thanks
for
Watching**