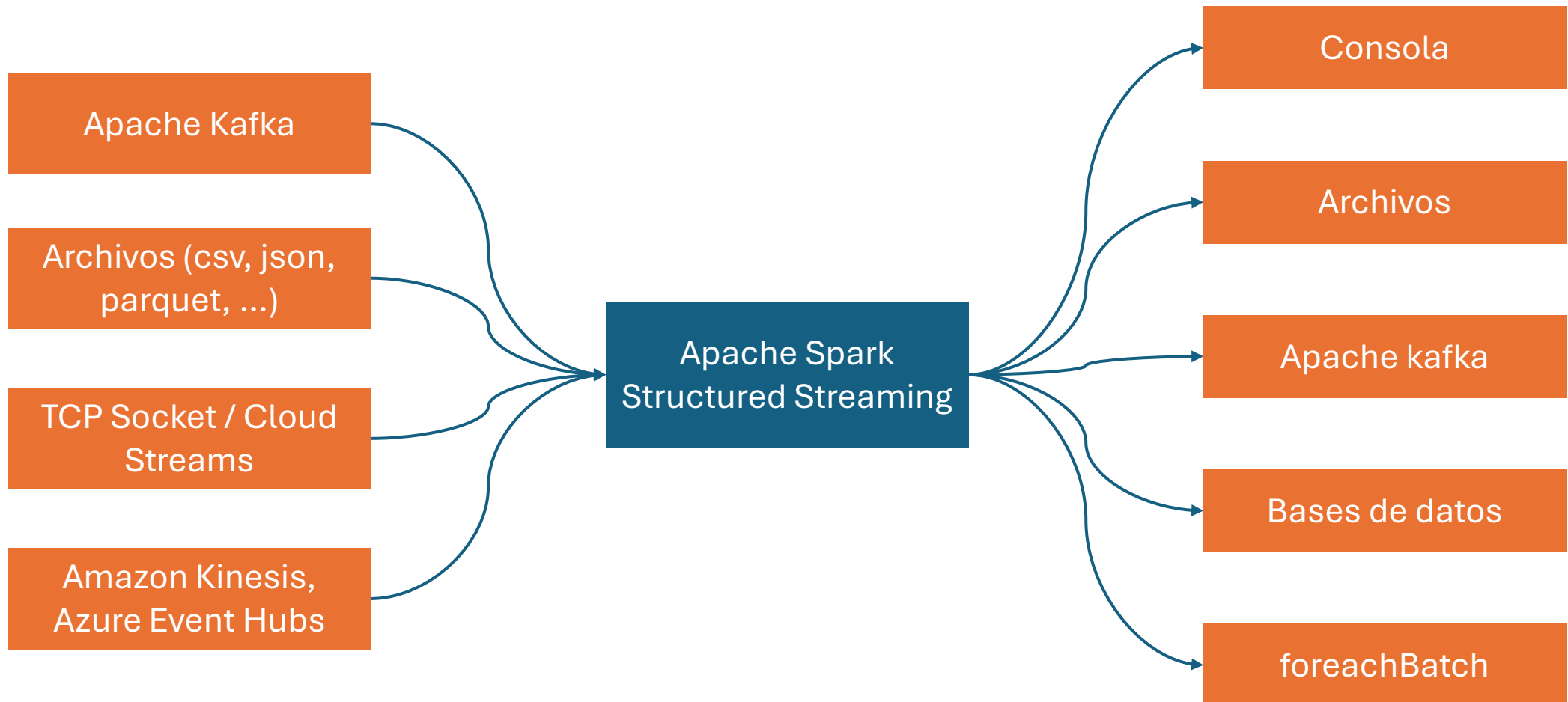


Sesión práctica de Spark para datos en Streaming

Fuentes y salidas

Spark **lee** datos de **fuentes**, los **transforma** (aplicando la lógica de la tabla infinita), y **escribe** los resultados en **salidas**.



Una Tabla Infinita (unbounded input table)



- Cada nuevo dato que llega al sistema (un clic, una venta, una lectura de sensor) se considera una nueva fila que se anexa al final de esta tabla.
- Esta "tabla" no tiene un final definido; es "no acotada" (unbounded) porque siempre pueden llegar más datos en el futuro.
- Podemos usar las mismas herramientas que ya conocemos para consultar tablas estáticas (como SQL y operaciones de DataFrame) para analizar datos en tiempo real.

Estructura usual de un script

Etapas	Código
Leer entradas	<code>df = spark.readStream.format(...).load()</code>
Define la Lógica/Transformación	<code>query = df.filter(...).join(...).groupby(...)</code>
Define la Salida (Sink) y el Modo	<code>query.writeStream.trigger(...).outputMode(...).format(...)</code>
Inicia la consulta	<code>query.start()</code>
Espera	<code>query.awaitTermination()</code>

Transformaciones

Sin estado

La transformación de cada fila es independiente de las demás.

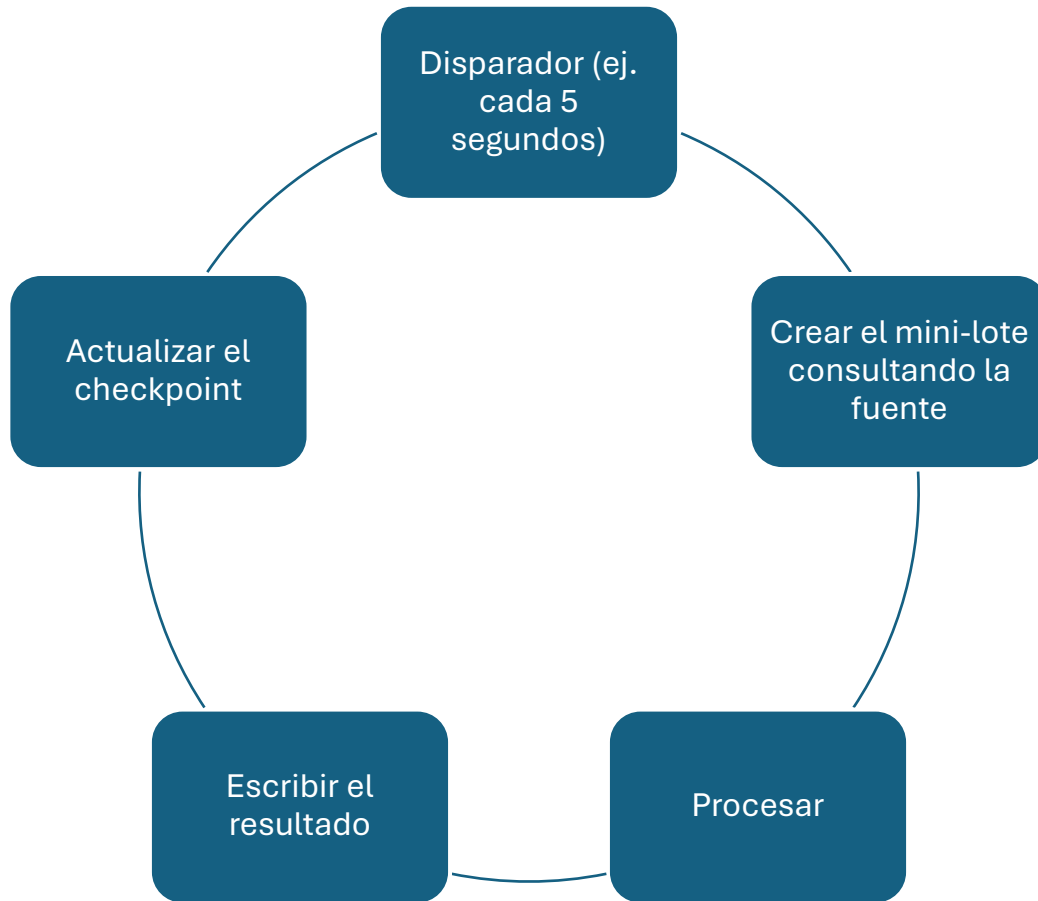
`select, filter, withColumn, explode, drop, join, ...`

Con estado

Se requiere mantener un estado a lo largo del tiempo para poderlas llevar a cabo.

`groupBy, join, dropDuplicates`

Micro-Lotes (Micro-Batch Processing)



- **Rendimiento:** Es más eficiente que procesar cada evento individualmente, ya que aprovecha al máximo el motor de optimización de Spark.
- **Simplicidad:** Permite usar la misma API de DataFrames para el procesamiento en batch y en streaming. El código es prácticamente idéntico.
- **Tolerancia a Fallos:** Si la aplicación falla, al reiniciarse lee el último checkpoint y sabe exactamente desde dónde continuar, garantizando que los datos no se pierdan ni se dupliquen.

Disparadores (triggers)

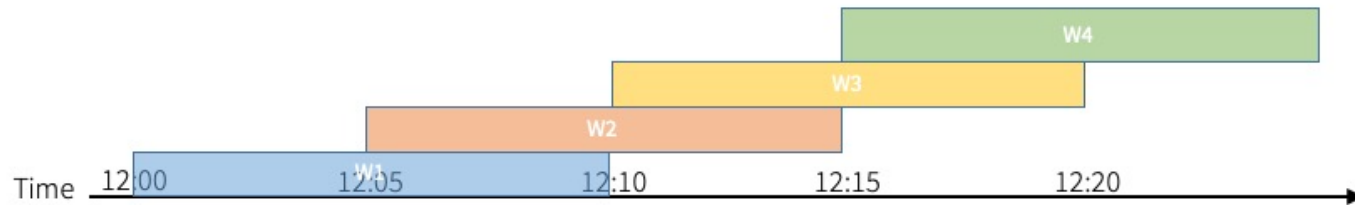
Disparador		Código
processingTime	<i>Intenta</i> ejecutar el proceso creando un lote cada cierto tiempo.	<pre>df = spark.readStream . trigger(processingTime="5 seconds") ...</pre>
once	Ejecuta un solo lote y termina	<pre>df = spark.readStream . trigger(once=True) ...</pre>
availableNow	Ejecuta todos los lotes y termina	<pre>df = spark.readStream . trigger(availableNow=True) ...</pre>
continuous	Procesa cada dato tan pronto llega	<pre>df = spark.readStream . trigger(continuous="5 seconds") ...</pre>

Ventanas

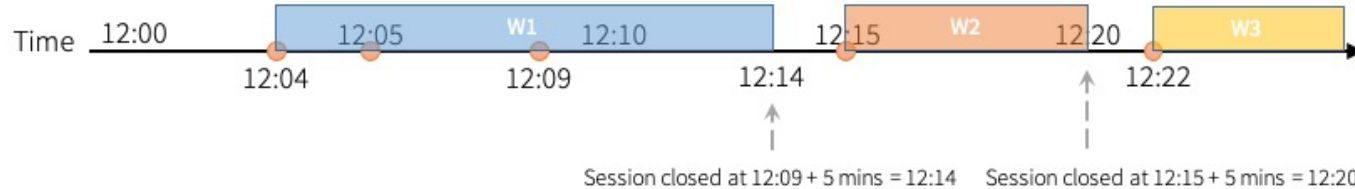
Tumbling Windows (5 mins)



Sliding Windows (10 mins, slide 5 mins)



Session Windows (gap duration 5 mins)



Tumbling

```
df.groupBy(f.window("timestamp", "10  
minutes")).count()
```

Sliding

```
df.groupBy(f.window("timestamp", "10  
minutes", "5 minutes")).count()
```

Session

```
df.groupBy(session_window("timestamp",  
"5 minutes")).count()
```


Manejo de datos atrasados con watermarks

- Por latencia de red o fallos, los datos pueden llegar **tarde y desordenados**. Spark no puede esperar indefinidamente o la memoria se llenará.
- La marca de agua le dice a Spark cuánto tiempo debe esperar para cerrar la ventana

```
# Se define el watermark sobre la columna
df_wm = df.withWatermark("timestamp", "10
minutes")

# Se utiliza la columna en una ventana
df_wm.groupBy(
    window("timestamp", "5 minutes")
).count()
```

Joins entre streams

Stream	Static	Inner	✓, sin estado
		Left	✓, sin estado
		Full	✗
		Semi	✓, sin estado
Stream	Stream	Inner	✓, watermark a ambos lados para limpiar el estado
		Left	⚠, debe tener watermark a la izquierda
		Full	⚠, debe tener watermark en un lado
		Semi	⚠, debe tener watermark a la derecha

```
df1 = df.withWatermark("time1", "2 hours")
df2 = df.withWatermark("time2", "3 hours")
```

```
df1.join(
    df2,
    f.expr("""
        id1 = id2 AND
        time1 >= time2 AND
        time1 <= time2 + interval 1 hour
    """))
```

Modos de salida (Output Modes)

Disparador		Código
append	Escribe nuevas filas a la salida sin modificar las demás. Por lo tanto debe esperar a que se cierren las ventanas para poder escribir.	<code>df.writeStream.outputMode("append")</code> ...
update	Escribe nuevas filas y actualiza sólo aquellas que han cambiado.	<code>df.writeStream.outputMode("update")</code> ...
complete	Emite toda la tabla de resultado en cada lote.	<code>df.writeStream.outputMode("complete")</code> ...

Groupby, watermarks y output mode

```
(  
    df  
    .groupby("llave").count()  
    .writeStream  
    .outputMode("append")  
    .option("checkpointLocation", "checkpoint/")  
    .format("parquet")  
    .start("file.parquet")  
)
```

```
(  
    df  
    .withWatermark("timestamp", "10 seconds")  
    .groupby("llave", "timestamp").count()  
    .writeStream  
    .outputMode("append")  
    .option("checkpointLocation", "checkpoint/")  
    .format("parquet")  
    .start("file.parquet")  
)
```

```
(  
    df  
    .groupby("llave", "timestamp").count()  
    .writeStream  
    .outputMode("append")  
    .option("checkpointLocation", "checkpoint/")  
    .format("parquet")  
    .start("file.parquet")  
)
```

Ejercicio práctico