

# Convert String to Timestamp in PySpark (With Example)

Authored by  
**Mohammed looti**

November 11, 2025

## RECOMMENDED CITATION

Mohammed looti (2025). *Convert String to Timestamp in PySpark (With Example)*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=16501>

The effective management of large-scale data hinges critically on the accurate interpretation and manipulation of data types. In distributed computing environments such as [Apache Spark](#), handling temporal data--information related to time--demands that it be stored in a format optimized for complex analytical operations like duration calculation, time-series forecasting, and window partitioning. While raw source systems often provide date and time information as simple strings, converting this raw text representation into a native [Timestamp](#) type is not optional; it is a fundamental requirement for building reliable and high-performing data pipelines.

Within the [PySpark](#) ecosystem, this necessary conversion is efficiently accomplished using powerful, highly optimized built-in SQL functions. The primary tool for transforming a string column into a timestamp column within a [DataFrame](#) is the `to_timestamp` function. This function provides data engineers with the crucial ability to specify the exact format of the input string, thereby ensuring precise and accurate parsing of temporal data into the required structured format, a process that is vital for ensuring data quality and analytical precision.

Standard practice dictates importing the necessary functions from the SQL module and then applying the transformation using the `withColumn` method. The `withColumn` method is preferred because it facilitates schema transformation by creating a new column without modifying the original data in place, thereby upholding Spark's core principle of immutability. This functional approach ensures that data transformations are traceable and predictable. The following concise syntax illustrates the typical structure used to execute this essential data type change, resulting in a new column derived from the original string data:

```
from pyspark.sql import functions as F
```

```
df = df.withColumn('ts_new', F.to_timestamp('ts', 'yyyy-MM-dd HH:mm:ss'))
```

This command generates a new column named `ts_new`, which is populated by parsing the string content of the original `ts` column. The string format `'yyyy-MM-dd HH:mm:ss'` serves as the critical parsing mask, instructing the Spark engine on the precise structure it should expect for the input data. Understanding and correctly applying this format mask is absolutely paramount, as any misalignment between the mask and the data structure will result in `NULL` values, compromising the integrity of the dataset. This foundational concept will be explored further through a detailed, practical implementation.

## The Performance Imperative of Type Conversion

Data type management is far more than a simple structural formality; it directly influences the performance and correctness of analytical queries within a distributed computing environment. When temporal data is left as a generic string, [PySpark](#) is unable to utilize its highly optimized

internal structures designed for datetime operations. Manipulating string data requires computationally expensive, byte-by-byte comparisons and string logic, a process that introduces significant overhead and is inherently slow when applied across massive datasets managed by a distributed cluster.

In contrast, converting the column to the native [TimestampType](#) allows the Spark engine to store the time value as a highly efficient numerical representation, typically recorded as the number of milliseconds or microseconds since the epoch (January 1, 1970). This numerical representation enables extremely fast arithmetic operations, comparisons, and indexing, which are essential for high-throughput data processing. Without this crucial conversion, the performance benefits of using a distributed framework like Spark for time-based analysis are severely diminished.

Furthermore, the utility of many advanced analytical functions within the PySpark library is strictly dependent on correct typing. Functions designed for calculating temporal differences, utilizing sliding or hopping [DataFrame](#) window functions, or performing time-series aggregations explicitly require columns to be of a proper temporal type. Failure to convert the raw string representation prevents the use of these specialized and optimized tools, forcing developers to rely on inefficient UDFs (User Defined Functions) or complex, error-prone string manipulation logic. The conversion step is therefore a non-negotiable prerequisite for scalable, robust time-series analysis and accurate data aggregation.

## Mastering the PySpark Conversion Functions

Successful string-to-timestamp conversion relies on two core [Spark](#) components: the comprehensive `pyspark.sql.functions` module and the immutable `DataFrame.withColumn` method. The `pyspark.sql.functions` module, commonly imported via the alias `F`, acts as the central repository for dozens of specialized SQL functions that operate directly on DataFrame columns. Crucially, these functions are executed natively by the optimized Spark catalyst optimizer, ensuring performance that is vastly superior to any custom Python function executed row-by-row.

The workhorse within this module is `F.to_timestamp(column, format)`. This function requires two mandatory arguments: the name of the input column containing the string data and the format string that precisely specifies the expected structure of the incoming temporal data. The format string must align perfectly with the incoming data's arrangement; typical format elements include `yyyy` (four-digit year), `mm` (two-digit month), `dd` (two-digit day), `HH` (24-hour hour), `mm` (minute), and `ss` (second). Any slight deviation or mismatch between the expected format mask and the actual input string will result in a `NULL` value being generated for that record, underscoring the extreme sensitivity of this parsing operation.

The surrounding structure, `df.withColumn(new_column_name, column_expression)`, is

responsible for integrating the result of the conversion back into the [DataFrame](#) structure. As a transformation operation, it returns a new DataFrame rather than modifying the existing one. The first argument defines the name of the new column to be created (or an existing one to be replaced), and the second argument provides the computational logic--in this case, the highly efficient result produced by `F.to_timestamp()`. Using `withColumn` to create a new column (e.g., `ts_new`) while retaining the original string column (`ts`) is often considered a best practice in ETL workflows, as it preserves the raw data for auditing while providing the optimized data type for processing.

## Practical Implementation: Setting Up the Environment

To provide a clear demonstration of the conversion process, we will establish a scenario involving simulated transactional sales data. Imagine we are working with a [PySpark](#) DataFrame where the critical time of sale is currently stored in a string column named `ts`. Our primary objective is to transform this `ts` column into a true [Timestamp](#) type, making it immediately available for complex time-based analysis.

The initial requirement for any PySpark operation is the instantiation of a `SparkSession`, which serves as the entry point for all functionality. Following this, we define our sample data, ensuring that the dates and times are explicitly represented as strings adhering to the `'YYYY-MM-DD HH:MM:SS'` format. Recognizing and remembering this precise string structure is mandatory, as it directly dictates the format mask that must be used during the subsequent conversion step.

The following code block provides the necessary setup, defining the sample data, constructing the initial DataFrame, and displaying its contents to clearly establish the starting point of our transformation exercise:

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.getOrCreate()
```

```
#define data
```

```
data = ,
```

```
,
```

```
,
```

```
]
```

```
#define column names
```

```
columns =
```

```
#create dataframe using data and column names
```

```
df = spark.createDataFrame(data, columns)
```

```
#view dataframe
df.show()

+-----+-----+
| ts|sales|
+-----+-----+
|2023-01-15 04:14:22| 225|
|2023-02-24 10:55:01| 260|
|2023-07-14 18:34:59| 413|
|2023-10-30 22:20:05| 368|
+-----+-----+
```

Prior to proceeding with the transformation, it is essential to confirm the existing schema of the DataFrame. Explicitly verifying the data types ensures that we are indeed addressing a string column that requires conversion. This validation step is achieved by calling the `df.dtypes` attribute, which returns a list detailing the name and assigned type of every column in the current structure. This process is crucial for validating assumptions about incoming data and preparing for subsequent debugging, should the conversion encounter issues.

```
#check data type of each column
df.dtypes
```

As confirmed by the output, the `ts` column is currently assigned the generic `string` data type, while the `sales` column is correctly identified as a numerical `bigint`. The explicit presence of the `string` type for the time data validates the necessity of our conversion task, formally preparing the data for the application of the `F.to_timestamp()` function.

## Executing and Verifying the Conversion

With the initial DataFrame prepared and its string-based time column identified, we can now apply the transformation using the previously established syntax. This process involves ensuring the import of the `functions` module as `F` and utilizing `withColumn` to integrate the newly typed column, `ts_new`, into the DataFrame. The core logic requires passing the string column `'ts'` and the exact format mask `'YYYY-MM-dd HH:mm:ss'` into `F.to_timestamp()`.

The precision required for the format mask cannot be overstated. For instance, if the input data included fractional seconds (e.g., milliseconds like `'2023-01-15 04:14:22.123'`), the format mask must be correspondingly updated to include `'.SSS'` (e.g., `'YYYY-MM-dd HH:mm:ss.SSS'`). Since our current sample data only includes full seconds, the simpler mask is adequate. Attention

to these fine details is critical for preventing data loss or the generation of null values during the parsing phase, which would severely impede subsequent analytical processes.

The code below executes the conversion, creating a resulting DataFrame that includes the original columns alongside the new, properly typed timestamp column. We then display the updated DataFrame to visually confirm the inclusion of the new column and its content:

### from pyspark.sql import functions as F

```
#convert 'ts' column from string to timestamp
df = df.withColumn('ts_new', F.to_timestamp('ts', 'yyyy-MM-dd HH:mm:ss'))

#view updated DataFrame
df.show()
```

```
+-----+-----+-----+
| ts|sales| ts_new|
+-----+-----+-----+
|2023-01-15 04:14:22| 225|2023-01-15 04:14:22|
|2023-02-24 10:55:01| 260|2023-02-24 10:55:01|
|2023-07-14 18:34:59| 413|2023-07-14 18:34:59|
|2023-10-30 22:20:05| 368|2023-10-30 22:20:05|
+-----+-----+-----+
```

Although the output of `df.show()` shows the data in `ts_new` looking visually identical to the original string data in `ts`, this similarity is purely cosmetic. Internally, the data in `ts_new` has been structurally re-encoded by the Spark engine into the highly optimized [TimestampType](#), making it instantly suitable for complex temporal arithmetic and filtering operations. To provide definitive proof of this structural change, we must formally verify the updated schema.

## Verification of Schema Transformation

The final and most critical step in validating the string-to-timestamp conversion is to execute the `df.dtypes` function one final time on the updated DataFrame. This check provides formal confirmation that [PySpark](#) recognizes the new column as a true [Timestamp](#) data type. This definitive validation ensures that the column is structurally sound and ready for subsequent analytical pipelines. Had the conversion failed, perhaps due to an incorrect format mask or inconsistent input data, the column might still appear, but its internal type would likely revert to a less useful generic type, or its values would be entirely null.

### #check data type of each column

## df.dtypes

The resulting output unequivocally confirms that the new column, `ts_new`, now possesses the data type `timestamp`. This result validates the successful application of the `F.to_timestamp()` function and confirms that the column is structurally prepared for any time-based analysis required by the project. By adhering to best practices in data type manipulation, we have seamlessly integrated temporal data into the optimized Spark environment.

## Mitigating Common Time Zone and Formatting Pitfalls

While the mechanical conversion process is straightforward, data engineers frequently encounter complex issues, particularly those related to time zone handling and input data inconsistencies. A crucial best practice involves thoroughly understanding how [Spark](#) manages time zones. By default, PySpark often reads and writes timestamps assuming the values are in the local time zone of the cluster or the session configuration, though a modern tendency is to default to UTC (Coordinated Universal Time).

If your raw string data lacks explicit time zone information, Spark will interpret it based on the session configuration. If the data originates from a specific time zone (e.g., EST) but the Spark session is configured for UTC, a significant time disparity of several hours may inadvertently be introduced, leading to incorrect aggregation and filtering results. To mitigate this risk, always ensure that if the input strings represent UTC time, the Spark session configuration reflects this. If the raw data includes time zone offsets (e.g., `'2023-01-15 04:14:22+05'`), specialized functions such as `to_utc_timestamp()` may be required in conjunction with `to_timestamp()` to normalize the time correctly.

Another prevalent pitfall is the format string mismatch. If even a small portion of the input string data fails to match the specified format mask, the conversion for that entire record will silently result in a `NULL` value. When working with inconsistent or "messy" raw data, it is strongly recommended to employ data cleansing steps first, such as using regular expressions (regex) to standardize the string format before conversion. Alternatively, for robust error handling, the `try_cast` function, available in newer [PySpark](#) versions, is highly recommended. `try_cast` attempts the conversion and, upon failure, returns `NULL` rather than raising an error, providing a cleaner, more resilient way to process bad records without interrupting the entire ETL job.

## Additional Resources for Temporal Data Mastery

For those looking to move beyond basic conversion and engage in advanced manipulation of date and time data types in PySpark, consulting the official documentation provides the most exhaustive

and reliable guidance on available functions and format parameters. Key areas for further development include studying advanced date arithmetic, mastering time zone conversion functions, and learning how to use the `date_format` function for outputting timestamps into highly customized string formats.

A few helpful resources for continued learning include:

Official PySpark SQL Functions Documentation: For a complete and updated list of all date and time manipulation functions.

Spark SQL Data Types Reference: Provides a detailed technical explanation of the [TimestampType](#) and its internal characteristics.

Tutorials on Time Series Analysis in PySpark: Demonstrations of using converted timestamp columns for complex analytical tasks such as windowing and forecasting.