

STOP USING .WITHCOLUMN() IN A LOOP FROM 20 CALLS TO 1

```
○○○  
1 # The anti-pattern  
2 for col_name, expr in mapping.items():  
3 df = df.withColumn(col_name, expr)  
4  
5 # The fix  
6 df = df.withColumns(mapping)
```



WHY .WITHCOLUMN() IN A LOOP HURTS

- **New DataFrame Every Call** — Each `.withColumn()` returns a brand-new DataFrame object with a fresh Project node on the logical plan.
- **Nested Plan Explosion** — 50 loop iterations = 50 nested Project nodes that Catalyst has to traverse and optimize through.
- **StackOverflowError Risk** — With hundreds of columns (common in wide fact tables or IoT data), recursive plan traversal can crash.
- **Planning Time Bloat** — Even before crashing, analysis and optimization passes scale with plan depth, turning milliseconds into seconds.

THE PROBLEM



[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

THE FIX

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

.WITHCOLUMNS() — SPARK 3.3+

Pass a dictionary of column expressions. One call, one plan node, same result.

○ ○ ○

```
1 # One dict, one plan node
2 df_result = df.withColumns({
3   "date_upper": F.upper(col("trnx_date")),
4   "name_clean": F.trim(col("custName")),
5   "bal_abs": F.abs(col(`ACCT-bal`)),
6   "amt_rounded": F.round(col("trnxAmt"), 2),
7   "loaded_at": F.current_timestamp(),
8 })
```



REAL WORLD USE

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

4/6



Metadata Enrichment

Add ingestion_timestamp, source_system, row_hash in one pass

12
34

Type Casting

Cast and format multiple columns without plan bloat



Column Standardisation

snake_case + abbreviation expansion across all columns



Business Rules

Apply derived column logic from config or lookup tables



Data Validation

Add is_valid, row_quality flags in a single step



Audit Columns

pipeline_run_id, load_date, environment tags at once



 Databricks.news

BONUS: RENAME

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

.WITHCOLUMNSRENAMED() — SPARK 3.4+

Same anti-pattern, same fix. Dict of old name to new name, one plan node.

○ ○ ○

```
1 # Same problem, same fix
2 df.withColumnsRenamed({
3   "CustID": "customer_id",
4   "trnxAmt": "transaction_amount",
5   "ACCT-bal": "account_balance",
6   "pmt$Status": "payment_status",
7   "addr.Line1": "address_line_1",
8 })
```

Build the dict dynamically with regex for snake_case + abbreviation expansion.



WANT MORE TIPS?

FOLLOW FOR DAILY
DATABRICKS TIPS

SIGN UP TO THE NEWSLETTER



dailydatabricks.tips