# Частина 1. Building an End-to-End Streaming Pipeline

```
PROBLEMS 7   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS 11
matvieienko@DESKTOP-T2SJUS5:~/final_project$ python3 02_process_streaming.py
      -----------------------------------------------
      |     default     |  11 |  0 |  0 |  0 || 11 |  0 |
      -----------------------------------------------
:: retrieving :: org.apache.spark#spark-submit-parent-7cd3f96d-490a-4b37-bb27-3ef7e8a903d3
      confs: [default]
      0 artifacts copied, 11 already retrieved (0kB/11ms)
25/11/30 16:02:04 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j2-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
>>> Reading athlete_bio...
>>> Reading athlete_event_results and writing to Kafka...
>>> Data written to Kafka successfully.
>>> Starting Stream...
25/11/30 16:03:02 WARN ResolveWriteToStream: spark.sql.adaptive.enabled is not supported in streaming DataFrames/Datasets and will be disabled.
Batch 0 processing... Rows: 143
Batch 1 processing... Rows: 271
Batch 2 processing... Rows: 355
Batch 3 processing... Rows: 428
Batch 4 processing... Rows: 503
Batch 5 processing... Rows: 575
Batch 6 processing... Rows: 661
Batch 7 processing... Rows: 741
Batch 8 processing... Rows: 817
Batch 9 processing... Rows: 903
Batch 10 processing... Rows: 975
Batch 11 processing... Rows: 1055
Batch 12 processing... Rows: 1123
Batch 13 processing... Rows: 1174
Batch 14 processing... Rows: 1246
Batch 15 processing... Rows: 1316
```

```
PROBLEMS 7   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS 11
|Sailing            |nan   |Male  |TPE |173.0                |72.0                |2025-11-30 16:10:01.017|
|Volleyball         |nan   |Female|CMR |181.0                |74.66666666666667   |2025-11-30 16:10:01.017|
|Sailing            |nan   |Female|ARU |169.5                |63.5                |2025-11-30 16:10:01.017|
|EquestrianDressage |nan   |Male  |GDR |175.33333333333334   |70.33333333333333   |2025-11-30 16:10:01.017|
|Fencing            |nan   |Male  |CUB |174.31578947368422   |68.98947368421052   |2025-11-30 16:10:01.017|
|Fencing            |nan   |Male  |MEX |177.73170731707316   |70.7560975609756    |2025-11-30 16:10:01.017|
|Athletics          |nan   |Female|CZE |175.25               |65.08333333333333   |2025-11-30 16:10:01.017|
|Boxing             |nan   |Male  |AND |185.0                |80.0                |2025-11-30 16:10:01.017|
|Basketball         |nan   |Male  |CAF |192.58333333333334   |84.75               |2025-11-30 16:10:01.017|
|Boxing             |nan   |Male  |ALG |174.23076923076923   |65.32307692307693   |2025-11-30 16:10:01.017|
|CyclingRoad        |nan   |Male  |ALB |183.0                |82.0                |2025-11-30 16:10:01.017|
|CrossCountrySkiing |nan   |Male  |CMR |198.0                |54.0                |2025-11-30 16:10:01.017|
|EquestrianDressage |nan   |Male  |AUT |179.25               |73.875              |2025-11-30 16:10:01.017|
|CyclingRoad        |nan   |Male  |CMR |173.35               |67.8                |2025-11-30 16:10:01.017|
|Athletics          |Silver|Female|CRO |193.0                |71.0                |2025-11-30 16:10:01.017|
|Judo               |Gold  |Male  |GEO |184.0                |90.0                |2025-11-30 16:10:01.017|
|SkiJumping         |nan   |Male  |BLR |174.16666666666666   |58.16666666666664   |2025-11-30 16:10:01.017|
|Judo               |nan   |Male  |BIZ |168.5                |78.5                |2025-11-30 16:10:01.017|
|EquestrianDressage |nan   |Male  |TCH |174.0                |72.0                |2025-11-30 16:10:01.017|
|CyclingTrack       |nan   |Male  |CZE |178.38095238095238   |80.61904761904762   |2025-11-30 16:10:01.017|
+-------------------+------+------+----+---------------------+--------------------+-----------------------+
only showing top 20 rows

Batch: 35
+----------+------+------+-----------+---------------------+--------------------+-----------------------+
|sport     |medal |sex   |country_noc|avg_height           |avg_weight          |timestamp              |
+----------+------+------+-----------+---------------------+--------------------+-----------------------+
|Wrestling |nan   |Male  |AZE        |173.7906976744186    |77.97674418604652   |2025-11-30 16:11:22.87 |
|CyclingRoad|nan  |Male  |BER        |179.66666666666666   |72.66666666666666   |2025-11-30 16:11:22.87 |
|CyclingRoad|nan  |Female|TPE        |161.5                |54.5                |2025-11-30 16:11:22.87 |
```

DBeaver 25.2.5 - <olympic_dataset> Script-13

SELECT * FROM olympic_dataset.fp_matvieienko_enriched_athlete_avg ORDER BY timestamp DESC LIMIT 20;

| | sport | medal | sex | country_noc | avg_height | avg_weight | timestamp |
|---|---|---|---|---|---|---|---|
| 1 | CanoeSprint | nan | Male | CZE | 184.4333333333 | 82.6 | 2025-11-30 16:16:00 |
| 2 | Fencing | nan | Female | FIN | 167 | 6,465 | 2025-11-30 16:16:00 |
| 3 | Athletics | nan | Female | CAF | 181.8125 | 72.4375 | 2025-11-30 16:16:00 |
| 4 | Athletics | nan | Male | BEN | 179.7826086957 | 74.1304347826 | 2025-11-30 16:16:00 |
| 5 | CyclingMountainBik | nan | Female | CYP | 149 | | 2025-11-30 16:16:00 |
| 6 | Judo | Bronze | Male | ARG | 150 | 48 | 2025-11-30 16:16:00 |
| 7 | CyclingRoad | nan | Male | BAR | 173.7142857143 | 69.7142857143 | 2025-11-30 16:16:00 |
| 8 | Weightlifting | nan | Female | COK | 165 | 101 | 2025-11-30 16:16:00 |
| 9 | CyclingRoad | nan | Female | TPE | 161.5 | 54.5 | 2025-11-30 16:16:00 |
| 10 | Sailing | nan | Female | BER | 169.1428571429 | 64.1428571429 | 2025-11-30 16:16:00 |
| 11 | CyclingRoad | nan | Male | BER | 179.6666666667 | 72.6666666667 | 2025-11-30 16:16:00 |
| 12 | Fencing | nan | Female | BUL | 176 | 61.5 | 2025-11-30 16:16:00 |
| 13 | Sailing | nan | Male | DOM | 189 | 82 | 2025-11-30 16:16:00 |
| 14 | Taekwondo | nan | Male | GAB | 186 | 77 | 2025-11-30 16:16:00 |
| 15 | Judo | nan | Male | ALB | 175 | 81 | 2025-11-30 16:16:00 |
| 16 | Judo | nan | Male | COL | 178 | 83.5 | 2025-11-30 16:16:00 |
| 17 | Rowing | nan | Male | HON | 176 | 75 | 2025-11-30 16:16:00 |
| 18 | Wrestling | Bronze | Male | AZE | 176.7777777778 | 87.2222222222 | 2025-11-30 16:16:00 |
| 19 | EquestrianEventing | nan | Male | RSA | 178 | 73 | 2025-11-30 16:16:00 |
| 20 | EquestrianJumping | Bronze | Male | KSA | 175 | 57 | 2025-11-30 16:16:00 |

20 row(s) fetched - 0.1s (0.0s fetch), on 2025-11-30 at 16:16:18

# Частина 2. Building an End-to-End Batch Data Lake

**Airflow**   DAGs   Cluster Activity   Datasets   Security   Browse   Admin   Docs

15:22 UTC   AA

30.11.2025  13:28:23   All Run Types   All Run States   Clear Filters

Auto-refresh   25

Press ⌘+/ for Shortcuts

deferred  failed  queued  removed  restarting  running  scheduled  shutdown  skipped  success  up_for_reschedule  up_for_retry  upstream_failed  no_status

DAG
**fp_matvieienko_dag**  ▶ 2025-11-30, 11:14:33 UTC · bronze_to_silver

Clear task   Mark state as...   Filter DAG by task

⚠ Details   ⚙ Graph   ☑ Gantt   <> Code   📄 Audit Log   ☰ Logs   ⇄ XCom   ⌛ Task Duration

(by attempts)

1

All Levels   All File Sources   Wrap   Download   See More

Duration

landing_to_bronze
bronze_to_silver
silver_to_gold

Version: v2.9.3
Git Version: .release:81845de9d95a733b4eb7826aaabe23ba9813eba3

---

**Airflow**   DAGs   Cluster Activity   Datasets   Security   Browse   Admin   Docs

15:22 UTC   AA

30.11.2025  13:28:23   All Run Types   All Run States   Clear Filters

Auto-refresh   25

Press ⌘+/ for Shortcuts

deferred  failed  queued  removed  restarting  running  scheduled  shutdown  skipped  success  up_for_reschedule  up_for_retry  upstream_failed  no_status

DAG
**fp_matvieienko_dag**  ▶ 2025-11-30, 11:14:33 UTC · silver_to_gold

Clear task   Mark state as...   Filter DAG by task

⚠ Details   ⚙ Graph   ☑ Gantt   <> Code   📄 Audit Log   ☰ Logs   ⇄ XCom   ⌛ Task Duration

(by attempts)

1

All Levels   All File Sources   Wrap   Download   See More

Duration

landing_to_bronze
bronze_to_silver
silver_to_gold

Version: v2.9.3
Git Version: .release:81845de9d95a733b4eb7826aaabe23ba9813eba3