

Earthquake MST3 – Testing & Evidence Document

Author: Mikołaj Malec (298828)

Functional Tests

Test Case 1: Historical Data Ingestion (USGS & Terraquake)

Objective: Verify that the Apache NiFi ingestion pipelines correctly retrieved the full historical dataset for the year 2025.

Test Steps:

1. Trigger the bulk ingestion process for the start date 2025-01-01 for USGS and Terraquake.
2. Inspect the HDFS Raw directories to confirm file presence for all months/days for both sources and for the clean data.
3. Using Hive, export example data from both sources.

Expected Result: A complete sequence of raw Parquet files covering the entire requested year without gaps and example hive output form the loaded tables.

Actual Result:

HDFS Raw directories output is stored in Test_appendix, they show the dates from 2025-01-01 to 2025-12-22 for both sources and for the clean data.

```
[hive> SELECT * FROM usgs_raw LIMIT 1;
OK
tx2025abvp {"mag":1.6,"place":"35 km WSW of Ackerly, Texas","time":1735775714006,"title":"M 1.6 - 35 km WSW of Ackerly, Texas"} {"type":"Point","coordinates":[-102.059,32.401,8.374]} 2025-01-01
1 row selected (7.373 seconds)
[hive> SELECT * FROM terraquake_raw LIMIT 1;
OK
Feature {"eventId":41269422,"originId":133577801,"time":"2025-01-01T22:47:10.750000","author":"SURVEY-INGV","magType":"ML","mag":1.5,"magAuthor":"—",
"type":"earthquake","place":"1 km E Robecco sul Naviglio (MI)","version":100,"geojson_creationTime":"2025-12-29T16:06:47"} {"type":"Point","coordinates":[8.9,45.4407,11.6]} 2025-01-01
1 row selected (2.429 seconds)
```

```
vagrant@node1:~$ cqlsh -e "SELECT id, place, region_id, cluster_id FROM earthquakes.events;"

id | place | region_id | cluster_id
---|---|---|---
ak2025znstnm | 92 km SE of Kokhanok, Alaska | 1 | 2
44853052 | 4 km NE Anghiari (AR) | 9 | 93
us7000r16y | 153 km S of Severo-Kuril'sk, Russia | 5 | 50
44853002 | 2 km NE Cerreto di Spoleto (PG) | 9 | 93
nc75287501 | 11 km ESE of Hidden Valley Lake, CA | 2 | 6
44852952 | 5 km SW Frontone (PU) | 9 | 93
nc75287496 | 4 km ESE of San Ramon, CA | 2 | 142
us7000r16w | 103 km WSW of San Vicente de Cañete, Peru | 3 | 27
44852872 | 4 km S Lioni (AV) | 9 | 22
44852822 | 4 km NE Anghiari (AR) | 9 | 93
44852632 | 1 km E Santo Stefano di Magra (SP) | 9 | 93
44852512 | 5 km SE Preci (PG) | 9 | 93
44852152 | 5 km SW Frontone (PU) | 9 | 93
44852092 | 4 km NE Anghiari (AR) | 9 | 93
44851992 | 3 km NE Anghiari (AR) | 9 | 93

(15 rows)
vagrant@node1:~$ cqlsh -e "SELECT id, place, region_id, cluster_id FROM earthquakes.events;"

id | place | region_id | cluster_id
---|---|---|---
44862232 | Vesuvio | 9 | 100
44862162 | 1 km NE Tavarnelle Val di Pesa (FI) | 9 | 163
44862112 | 6 km W Sansepolcro (AR) | 9 | 163
44862082 | 5 km W Sansepolcro (AR) | 9 | 163
44862032 | 4 km NE Anghiari (AR) | 9 | 163
44862002 | 5 km NE Anghiari (AR) | 9 | 163
44861902 | 5 km NE Anghiari (AR) | 9 | 163
44861872 | 5 km SW Preci (PG) | 9 | 1
44861742 | 5 km W Sansepolcro (AR) | 9 | 163
44861702 | 8 km E Città di Castello (PG) | 9 | 1
ak2025znstnm | 92 km SE of Kokhanok, Alaska | 1 | 2
44853052 | 4 km NE Anghiari (AR) | 9 | 93
us7000r16y | 153 km S of Severo-Kuril'sk, Russia | 5 | 50
44853002 | 2 km NE Cerreto di Spoleto (PG) | 9 | 93
nc75287501 | 11 km ESE of Hidden Valley Lake, CA | 2 | 6
44852952 | 5 km SW Frontone (PU) | 9 | 93
nc75287496 | 4 km ESE of San Ramon, CA | 2 | 142
us7000r16w | 103 km WSW of San Vicente de Cañete, Peru | 3 | 27
44852872 | 4 km S Lioni (AV) | 9 | 22
44852822 | 4 km NE Anghiari (AR) | 9 | 93
44852632 | 1 km E Santo Stefano di Magra (SP) | 9 | 93
44852512 | 5 km SE Preci (PG) | 9 | 93
44852152 | 5 km SW Frontone (PU) | 9 | 93
44852092 | 4 km NE Anghiari (AR) | 9 | 93
44851992 | 3 km NE Anghiari (AR) | 9 | 93

(25 rows)
```

Status: Passed

Test Case 2: Speed Layer Stream Enrichment

Objective: Confirm that real-time events processed by the Speed Layer are correctly enriched with geospatial metadata (Region ID and Cluster ID) and stored in Apache Cassandra.

Test Steps:

1. Observe the Kafka topic for a new incoming event.
2. Allow the Spark Streaming job to process the event.
3. Query the Apache Cassandra table earthquake_realtime to retrieve the record.

Expected Result: Kafka consumer will intercept the event in the JSON format and the record in Cassandra for the same event must have the region_id and cluster_id columns populated (not null).

Actual Result:

```
magrant@magant1:~$ /usr/local/kafka/bin/kafka-console-consumer.sh --topic unified_quakes --bootstrap-server localhost:9092
{"id":"44861872","magnitude":1.4,"mag_type":"ML","place":"5 km SW Preci (PG)","time_utc":"2025-12-29T21:46:41.240000","longitude":13.0058,"latitude":42.8425,"depth_km":7.7,"source":"terraquake"}
{"id":"44861742","magnitude":0.9,"mag_type":"ML","place":"5 km W Sansepolcro (AR)","time_utc":"2025-12-29T20:54:12.270000","longitude":12.0875,"latitude":43.5808,"depth_km":9.8,"source":"terraquake"}
{"id":"44861792","magnitude":1,"mag_type":"ML","place":"8 km E Città di Castello (PG)","time_utc":"2025-12-29T20:50:05.550000","longitude":12.3353,"latitude":43.4812,"depth_km":11.1,"source":"terraquake"}
{"id":"44861692","magnitude":1.3,"mag_type":"ML","place":"5 km W Sansepolcro (AR)","time_utc":"2025-12-29T20:42:49.980000","longitude":12.0868,"latitude":43.5843,"depth_km":9.7,"source":"terraquake"}
{"id":"44861632","magnitude":2.3,"mag_type":"ML","place":"4 km NE Anghiari (AR)","time_utc":"2025-12-29T20:41:43.740000","longitude":12.0887,"latitude":43.5717,"depth_km":9.3,"source":"terraquake"}
{"id":"44861612","magnitude":0.6,"mag_type":"ML","place":"5 km W Sansepolcro (AR)","time_utc":"2025-12-29T20:34:04.600000","longitude":12.0807,"latitude":43.5892,"depth_km":10.1,"source":"terraquake"}
{"id":"44861472","magnitude":1.2,"mag_type":"ML","place":"6 km NE Claut (PN)","time_utc":"2025-12-29T19:55:44.730000","longitude":12.5715,"latitude":46.3035,"depth_km":9.7,"source":"terraquake"}
{"id":"44861292","magnitude":0.7,"mag_type":"ML","place":"4 km NE Sellano (PG)","time_utc":"2025-12-29T18:59:09.480000","longitude":12.9552,"latitude":42.9173,"depth_km":12.9,"source":"terraquake"}
{"id":"44861172","magnitude":0.8,"mag_type":"ML","place":"4 km NE Costacciaro (PG)","time_utc":"2025-12-29T18:03:35.890000","longitude":12.7522,"latitude":43.3762,"depth_km":13,"source":"terraquake"}
```

Status: Passed

Test Case 3: Cross-Source Deduplication Validation

Objective: Verify that the Batch Layer deduplication logic successfully merged events from USGS and Terraquake. The Clean Layer should not contain any two events that are within 0.1 degrees of each other and within 10 minutes of each other.

Test Steps:

1. Execute the following QA query in Apache Hive to check for residual duplicates:

```
SELECT
  -- IDs and Time
  a.id AS event_1,
  b.id AS event_2,
  a.time_utc AS time_1,
  b.time_utc AS time_2,

  -- Sources
  a.source_system AS source_1,
  b.source_system AS source_2,

  -- Magnitudes
  a.magnitude AS mag_1,
  b.magnitude AS mag_2,

  -- Location & Diff
  a.place,
  (a.time_utc - b.time_utc) AS time_diff_ms

FROM earthquake_clean a
JOIN earthquake_clean b
ON a.dt = b.dt
WHERE
  -- 1. Prevent self-matches and duplicate pairs
  a.id < b.id

  -- 2. Ensure different sources
  AND a.source_system != b.source_system
```

-- 3. Spatial filter (approx 0.1 degrees)
AND ABS(a.latitude - b.latitude) <= 0.1
AND ABS(a.longitude - b.longitude) <= 0.1

-- 4. Time filter: 60,000 ms = 1 Minute
AND ABS(a.time_utc - b.time_utc) <= 60000

LIMIT 10;

Expected Result: The query should return 0 rows, indicating that all overlapping events between the two sources were successfully merged into a single source of truth.

Actual Result:

```
hive> SELECT
-- IDs and Time
-- a.id AS event_1,
-- b.id AS event_2,
-- a.time_utc AS time_1,
-- b.time_utc AS time_2,
-- Sources
-- a.source_system AS source_1,
-- b.source_system AS source_2,
-- Magnitudes
-- a.magnitude AS mag_1,
-- b.magnitude AS mag_2,
-- Location & Diff
-- a.place,
-- (a.time_utc - b.time_utc) AS time_diff_ms
-- FROM earthquake_clean a
-- JOIN earthquake_clean b
-- ON a.dt = b.dt
-- WHERE
-- 1. Prevent self-matches and duplicate pairs
-- a.id < b.id
-- 2. ENSURE DIFFERENT SOURCES (The new filter)
-- AND a.source_system != b.source_system
-- 3. Spatial filter (approx 0.1 degrees)
-- AND ABS(a.latitude - b.latitude) <= 0.1
-- AND ABS(a.longitude - b.longitude) <= 0.1
-- 4. Time filter: 60,000 ms = 1 Minute
-- AND ABS(a.time_utc - b.time_utc) <= 60000
> LIMIT 10;
26/12/29 20:28:12 [HiveServer2-Background-Pool: Thread-176]: WARN ql.Driver: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = vagrant-202512292028086_Tera7959-3292-c72a-a308-bd3858f8c8df
Total jobs = 1
26/12/29 20:28:12 [HiveServer2-Background-Pool: Thread-176]: WARN ql.Driver: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/local/apache-hive-2.3.8-bin/lib/log4j-slf4j-impl-2.6.2.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/apache-tez-0.9.1-bin/lib/slf4j-log4j12-1.7.10.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.7.6/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/local/apache-hive-2.3.8-bin/lib/log4j-slf4j-impl-2.6.2.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/apache-tez-0.9.1-bin/lib/slf4j-log4j12-1.7.10.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.7.6/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
2025-12-29 20:28:31 Starting to launch local task to process map join; maximum memory = 477626368
2025-12-29 20:28:31 Starting to launch local task to process map join; maximum memory = 477626368
SLF4J: Failed to load class 'org.slf4j.impl.StaticLoggerBinder'.
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
SLF4J: Failed to load class 'org.slf4j.impl.StaticLoggerBinder'.
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
2025-12-29 20:29:01 Dump the side-table for tag: 0 with group count: 363 into file: file:/tmp/vagrant/5c918366-5886-44c9-9883-46bf1204889d/hive_2025-12-29_20-28-00_458_2783438958487787782-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile30---.hashtable
2025-12-29 20:29:02 Uploaded 1 file to: file:/tmp/vagrant/5c918366-5886-44c9-9883-46bf1204889d/hive_2025-12-29_20-28-00_458_2783438958487787782-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile30---.hashtable (11473931 bytes)
2025-12-29 20:29:02 End of local task; Time Taken: 31.716 sec.
Execution completed successfully
MapReduce task succeeded
Launching Job 3 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
2025-12-29 20:29:01 Dump the side-table for tag: 0 with group count: 363 into file: file:/tmp/vagrant/5c918366-5886-44c9-9883-46bf1204889d/hive_2025-12-29_20-28-00_458_2783438958487787782-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile30---.hashtable
2025-12-29 20:29:02 Uploaded 1 file to: file:/tmp/vagrant/5c918366-5886-44c9-9883-46bf1204889d/hive_2025-12-29_20-28-00_458_2783438958487787782-1/-local-10004/HashTable-Stage-3/MapJoin-mapfile30---.hashtable (11473931 bytes)
2025-12-29 20:29:02 End of local task; Time Taken: 31.716 sec.
26/12/29 20:29:08 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
Starting Job : job_1767822533978_0015, Tracking URL = http://node1:8080/proxy/application_1767822533978_0015/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1767822533978_0015
Hadoop job information for Stage:3 number of mappers: 0; number of reducers: 0
26/12/29 20:29:32 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.Counters: Group org.apache.hadoop.mapred.Task$Counter is deprecated. Use org.apache.hadoop.mapreduce.TaskCounter instead
2025-12-29 20:29:32:290 Stage-3 map = 0N, reduce = 0N
26/12/29 20:30:32 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.Counters: Group org.apache.hadoop.mapred.Task$Counter is deprecated. Use org.apache.hadoop.mapreduce.TaskCounter instead
2025-12-29 20:30:32:644 Stage-3 map = 0N, reduce = 0N
26/12/29 20:30:49 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.Counters: Group org.apache.hadoop.mapred.Task$Counter is deprecated. Use org.apache.hadoop.mapreduce.TaskCounter instead
2025-12-29 20:30:49:245 Stage-3 map = 100N, reduce = 0N
26/12/29 20:30:49 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.Counters: Group org.apache.hadoop.mapred.Task$Counter is deprecated. Use org.apache.hadoop.mapreduce.TaskCounter instead
Ended Job : job_1767822533978_0015
MapReduce Jobs Launched:
26/12/29 20:30:49 [HiveServer2-Background-Pool: Thread-176]: WARN mapreduce.Counters: Group FileSystemCounters is deprecated. Use org.apache.hadoop.mapreduce.FileSystemCounter instead
Stage-Stage-3: HDFS Read: 0 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
No rows selected (169.21 seconds)
hive>
```

Total MapReduce CPU Time Spent: 0 msec

OK

No rows selected (169.21 seconds)

hive>

Status: Passed

Analytical Module & ML Model Evaluation

Test Case 4: K-Means Optimization & Stability Validation

Objective: Verify that the ML pipeline performs robust hyper-parameter tuning by iterating through a cluster range of $k = 300$ to 600 . The test ensures that for each k , 3 random-seed trials are executed to ensure stability, and the final model is automatically selected based on the highest Calinski-Harabasz (CH) Index score.

Test Steps:

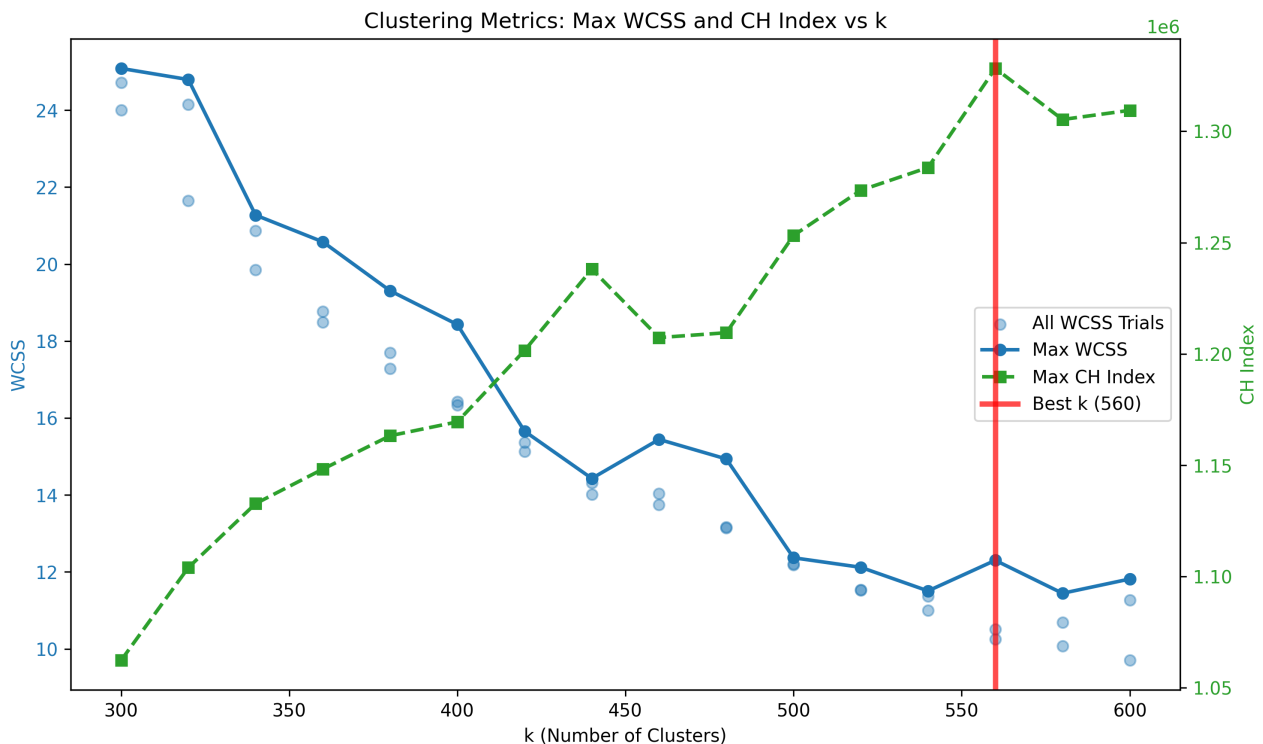
1. Review the PySpark training logs to confirm that multiple trials (iterations with different seeds) were executed for each value of k .

2. Extract the performance metrics table listing k values (50–700) and their associated CH scores.
3. Verify that the system selected the model corresponding to the maximum CH score.
4. Plot the Within-Cluster Sum of Squares (WCSS) against the CH Index and observe the trend.

Expected Result:

- The logs must show multiple distinct WCSS values for every k, confirming that stability trials were run.
- The final selected k must correspond to the peak Calinski-Harabasz score.
- The selected k should geometrically align with the "elbow" (inflection point) of the WCSS curve, indicating where the variance reduction stabilises.

Actual Result:



Status: Passed

Test Case 5: Centroid Back-Projection

Objective: Validate the mathematical accuracy of converting the machine learning model's 3D Cartesian centroids (X, Y, Z) back into Geographic coordinates (Latitude, Longitude) by performing a manual calculation check.

Test Steps:

1. Get a sample centroid vector from the model output.
2. Check if the model's output Lat/Lon matches the manual calculation.

Expected Result: The model output should match the manual calculation.

Actual Result:

```
hive> select * from dim_cluster_metadata limit 1;
OK
0 -0.2660453564903092 -0.6517040619291885 0.7095140480559485 45.22690837241916 -112.20677165688359 767.1438528108063
1 row selected (1.021 seconds)
```

$$r = \sqrt{x^2 + y^2 + z^2}$$

$$\phi = \arcsin\left(\frac{z}{r}\right)$$

$$x^2 \approx 0.07078$$

$$\phi = \arcsin\left(\frac{0.709514}{0.999454}\right) \approx \arcsin(0.70990)$$

$$y^2 \approx 0.42472$$

$$\phi \approx 0.7893 \text{ radians}$$

$$z^2 \approx 0.50341$$

Convert to degrees:

$$r = \sqrt{0.99891} \approx \mathbf{0.99945}$$

$$0.7893 \times \frac{180}{\pi} \approx \mathbf{45.223^\circ}$$

$$\lambda = \text{atan2}(-0.651704, -0.266045) \approx -1.9584 \text{ radians}$$

Convert to degrees:

$$-1.9584 \times \frac{180}{\pi} \approx \mathbf{-112.209^\circ}$$

The results are inside the rounding error $\pm 0.03^\circ$ which corresponds to ~ 300 meters

Status: Passed