

Key-Value vs Graph-based data lakes for realizing Digital Twin Systems

Daniel Pérez-Porras

Universidad de Málaga, Spain

Javier Troya Universidad de Málaga, Spain Paula Muñoz

Universidad de Málaga, Spain

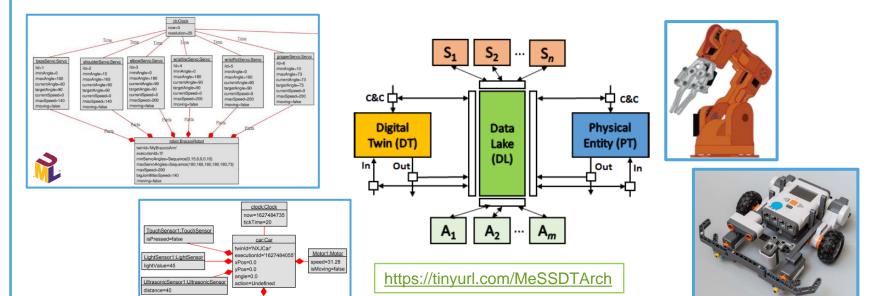
Antonio Vallecillo Universidad de Málaga, Spain





Digital Twin Systems

A Digital Twin is a comprehensive digital representation of an actual system, service or product (the Physical Twin), synchronized at a specified frequency and fidelity [1].

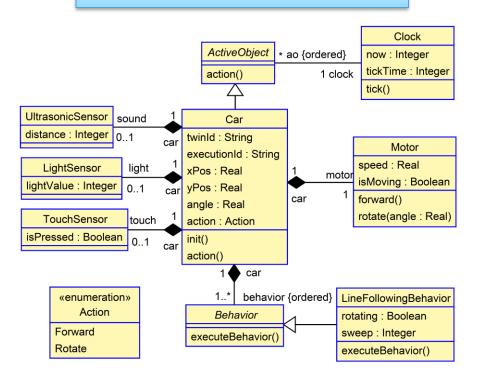


^[1] Digital Twin Consortium, "Glossary of digital twins", 2021. [Online]. Available: https://www.digitaltwinconsortium.org/glossary/index.html

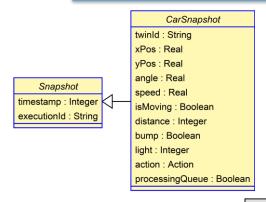
^[2] P. Muñoz, J. Troya, A. Vallecillo, Using UML and OCL Models to Realize High-Level Digital Twins, in: Proc. of ModDiT2021@MODELS'21, IEEE, 2021, pp. 212–220

Information format – Lego NXJ Car

Car's Digital Twin model



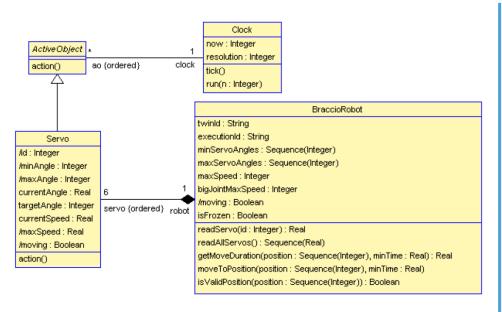
Car's Snapshots



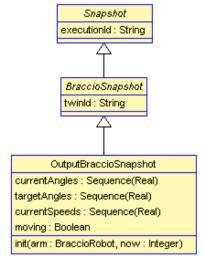


Information format – Arduino Braccio

Braccio's Digital Twin model



Braccio's Snapshots



outputSnapshot:OutputBraccioSnapshot

whenProcessed=Undefined executionId='test' twinId='MyBraccioArm' currentAngles=Sequence(90,90,90,90,90,73) targetAngles=Sequence(90,90,90,90,90,73) currentSpeeds=Sequence(0,0,0,0,0,0,0)

moving=false

timestamp=1656492966

isProcessed=false

redis as Data Lake

Key-value database

- You need to know the keys to query the values.
- ❖ To make more complex queries you need additional registers or composed keys

snapshotType : twinId : executionId : timestamp HGETALL PTOutputSnapshot: SSSCar: 1627484055: 16274848179

- 1) "twinId" 3) "executionId"
- "timestamp"
- "xPos" "vPos"
- "angle" "speed"
- 15) "isMoving"
- "distance" "bump"
- "light"
- 23) "action"
- 25) "processingQueue"

- "SSSCar"
- 4) "1627484055"
- 6) "16274848179"
- 8) "2.996127561"
- 10) "15.51560811"
- 12) "4.71"
- 14) "62.56" 16) "TRUE"
- 18) "35"
- 20) "FALSE"
- "45"
- 24) "#Forward"
- 26) "FALSE"

Snapshots in which the car was at a distance between 0 and 100

ZRANGEBYSCORE PTOutputSnapshot:SSSCar DISTANCE LIST 0 100

- 1) "PTOutputSnapshot:SSSCar:1627484055:16274847352"
- "PTOutputSnapshot:SSSCar:1627484055:16274848139"
- "PTOutputSnapshot:SSSCar:1627484055:16274848148"
- "PTOutputSnapshot:SSSCar:1627484055:16274848158"
- "PTOutputSnapshot:SSSCar:1627484055:16274848168"

All Snapshots for a given Physical Twin

ZRANGEBYLEX PTOutputSnapshot PROCESSED

"[PTOutputSnapshot:SSSCar:" "[PTOutputSnapshot:SSSCar:\xff"

- 1) "PTOutputSnapshot:SSSCar:1627484055:16274847352"
- 2) "PTOutputSnapshot:SSSCar:1627484055:16274847359"
- "PTOutputSnapshot:SSSCar:1627484055:16274847368"
- "PTOutputSnapshot:SSSCar:1627484055:16274847380"
- 5) "PTOutputSnapshot:SSSCar:1627484055:16274847391"



Advantages

- **Lightweight**, easy to deploy
- High-responsiveness. It is optimized to deliver fast responses to a massive amount of petitions.



Disadvantages

- **Rigid.** You need to define the database's structure beforehand according to the queries you are planning to make.
- **Duplicated information** to perform queries.



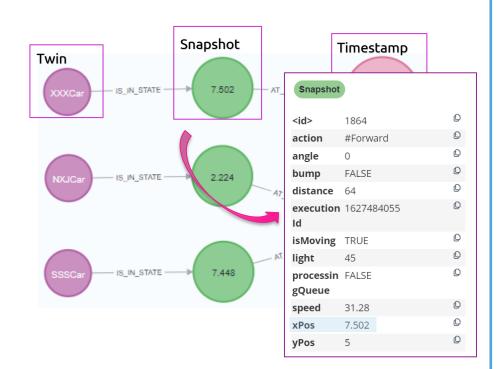
→ CO4j as Data Lake

Graph database

- You only need to know the database structure to query for values.
- It enables to make complex queries using the graph relationships.

Snapshots for any Digital Twin, which occur between two given

```
MATCH (before:Timestamp)
WHERE before.value \geqslant 16274847350
MATCH (after:Timestamp)
WHERE after.value \leqslant 16274847540
AND before.value < after.value
MATCH (before)-[:NEXT]\rightarrow(after)
MATCH (dt:DigitalTwin)-[:IS_IN_STATE]\rightarrow(snp)-[:AT_THE_TIME]\rightarrow(before)
MATCH (dt2:DigitalTwin)-[:IS_IN_STATE]\rightarrow(snp2)-[:AT_THE_TIME]\rightarrow(after)
WHERE dt \diamondsuit dt2
RETURN dt, snp, dt2, snp2, before, after
LIMIT 2
```









Q1: Does a given car follow the line during a given time window?





```
MATCH (dt: DigitalTwin {name:'NXJCar'})
MATCH (snp: Snapshot)
MATCH (ts: Timestamp)
MATCH (dt)-[:IS_IN_STATE]\rightarrow(snp)-[:AT_THE_TIME]\rightarrow(ts)
WHERE ts.value ≥ 16274847350
AND ts.value ≤ 16274847490
RETURN avg(snp.light)
   "avg(snp.light)"
   45.0
```

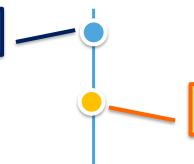






Q1: Does a given car follow the line during a given time window?





Q2: Average value of a given servo for all robot arms.

```
MATCH (ts:Timestamp)
WHERE ts.value ≤ 20 AND ts.value ≥ 0
MATCH (dt:DigitalTwin)
MATCH (dt)-[:IS_IN_STATE]→(snp)-[:AT_THE_TIME]→(ts)
WITH dt, snp.currentAngles_3 as snapshots
RETURN dt, avg(snapshots)
   "dt"
                    "avg(snapshots)"
   {"name":"braccio"}
                   94.125
   {"name": "braccio2"} 107.166666666666666
```



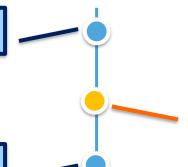




Q1: Does a given car follow the line during a given time window?

Q3: Number of cars in a squared area during a given time window.





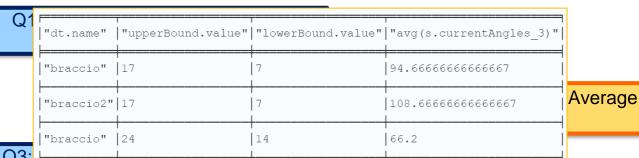
Q2: Average value of a given servo for all robot arms.

```
MATCH (snp:Snapshot)
MATCH (ts:Timestamp)
WHERE ts.value ≤ 16274847380 AND ts.value ≥ 16274847350
AND snp.xPos ≥ 3 AND snp.xPos ≤ 6
AND snp.yPos ≥ 0 AND snp.yPos ≤ 3
MATCH (dt)-[:IS_IN_STATE]\rightarrow(snp)-[:AT_THE_TIME]\rightarrow(ts)
WITH DISTINCT dt
RETURN count(dt)
  "count(dt)"
```











servo for all

Q3: Ivamoer de cars ne a squareu area

during a given time window.

```
MATCH timePeriod=(lowerBound:Timestamp)

-[:NEXT*2..15]→(upperBound:Timestamp)

WHERE upperBound.value - lowerBound.value = 10

WITH *, nodes(timePeriod) as timestamps

MATCH (dt:DigitalTwin)-[:IS_IN_STATE]

→(snp:Snapshot)-[:AT_THE_TIME]→(ts:Timestamp)

WHERE ts in timestamps

WITH upperBound, lowerBound, timestamps, dt,

collect(snp) as snapshots

UNWIND snapshots as s

RETURN dt.name, upperBound.value, lowerBound.value,
avg(s.currentAngles 3)
```

Q4: Average value of a given servo for each robot arm in the system during a certain period.







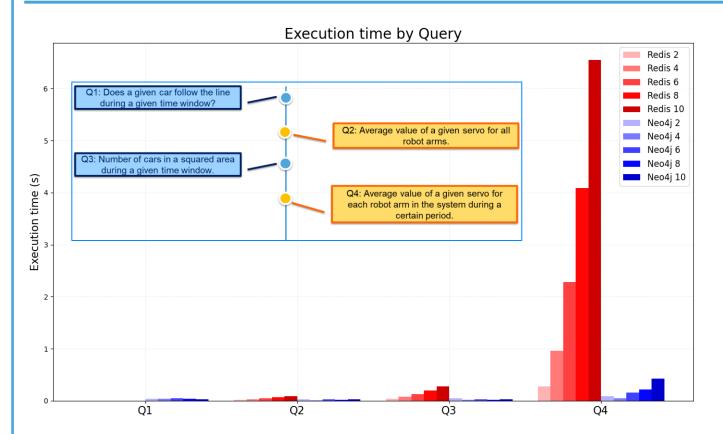
```
MATCH (snp : Snapshot)
 Q1: Does a given car foll WITH *, point({x: snp.xPos, y: snp.yPos, crs: 'cartesian'}) AS p1
                                 MATCH (snp2 : Snapshot)
     during a given time w
WITH *, point({x: snp2.xPos, y: snp2.yPos, crs: 'cartesian'}) AS p2
                                 WITH snp, snp2, point.distance(p1,p2) AS dist
                                 WHERE dist < 3
                                                                                                      le of a given servo for all
                                 MATCH (dt : DigitalTwin)
                                 MATCH (dt)-[:IS IN STATE]\rightarrow(snp)
                                                                                                      bot arms.
                                 MATCH (dt2 : DigitalTwin)
Q3: Number of cars in a so where dt \Leftrightarrow dt2
                                 MATCH (dt2)-[:IS_IN_STATE] \rightarrow (snp2)
      during a given time \underline{W}_{MATCH} (snp)-[:AT_THE_TIME]\rightarrow(ts)\leftarrow[:AT_THE_TIME]-(snp2)
                                 RETURN dt, dt2, snp, snp2, ts
                                                                                                                given servo for
                                                                                  Q4:
                                                                                  each
                                                                                                               ystem during a
                                                                                                               od.
  Q5: Number of collisions between
   cars during a given time window.
```

Comparative analysis 9004j









# Robots	# Snaps /robot	# Snaps
2	50	100
4	100	400
8	200	1600
16	400	6400
32	800	25600

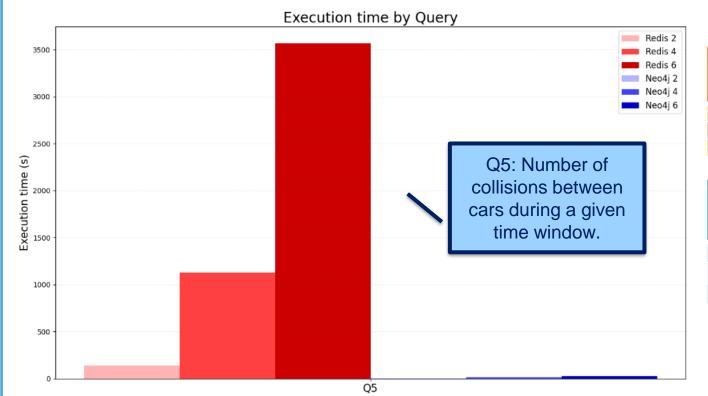
# Cars	# Snaps/c ar	# Snaps
6	213	1278
12	426	5112
24	852	20448
48	1704	81792
96	3408	327168

Comparative analysis 9004j









# Robots	# Snaps/ro bot	# Snaps
2	50	100
4	100	400
8	200	1600

# Cars	# Snaps/car	# Snaps
6	213	1278
12	426	5112
24	852	20448

Conclusions and future work



VS



- **Redis** is really convenient when you are not expecting complex queries or aggregated data from the database, since you will have to rely on manual postprocessing of the information.
- **Neo4j** allows complex queries with good performance and great expressivity by using a simple structure to model time information.
- In our <u>future work</u>,
 - We plan on comparing this approach with one using a time-series database.
 - Try real-time analysis using a working system to measure performance in ingestion and querying in a real environment.



https://tinyurl.com/MESSRedisVSNeo



Thank you for your attention

Key-Value vs Graph-based data lakes for realizing Digital Twin Systems

Daniel Pérez-Porras

Universidad de Málaga, Spain

Javier Troya Universidad de Málaga, Spain <u>Paula Muñoz</u>

Universidad de Málaga, Spain

Antonio Vallecillo Universidad de Málaga, Spain











Q1: Does a given car follow the line during a given time window?

Q3: Number of cars in a squared area during a given time window.

Q5: Number of collisions between cars during a given time window.

Q2: Average value of a given servo for all robot arms.

Q4: Average value of a given servo for each robot arm in the system during a certain period.