

# **Analisi di dati di monitoraggio con Apache Spark**

2023-2024

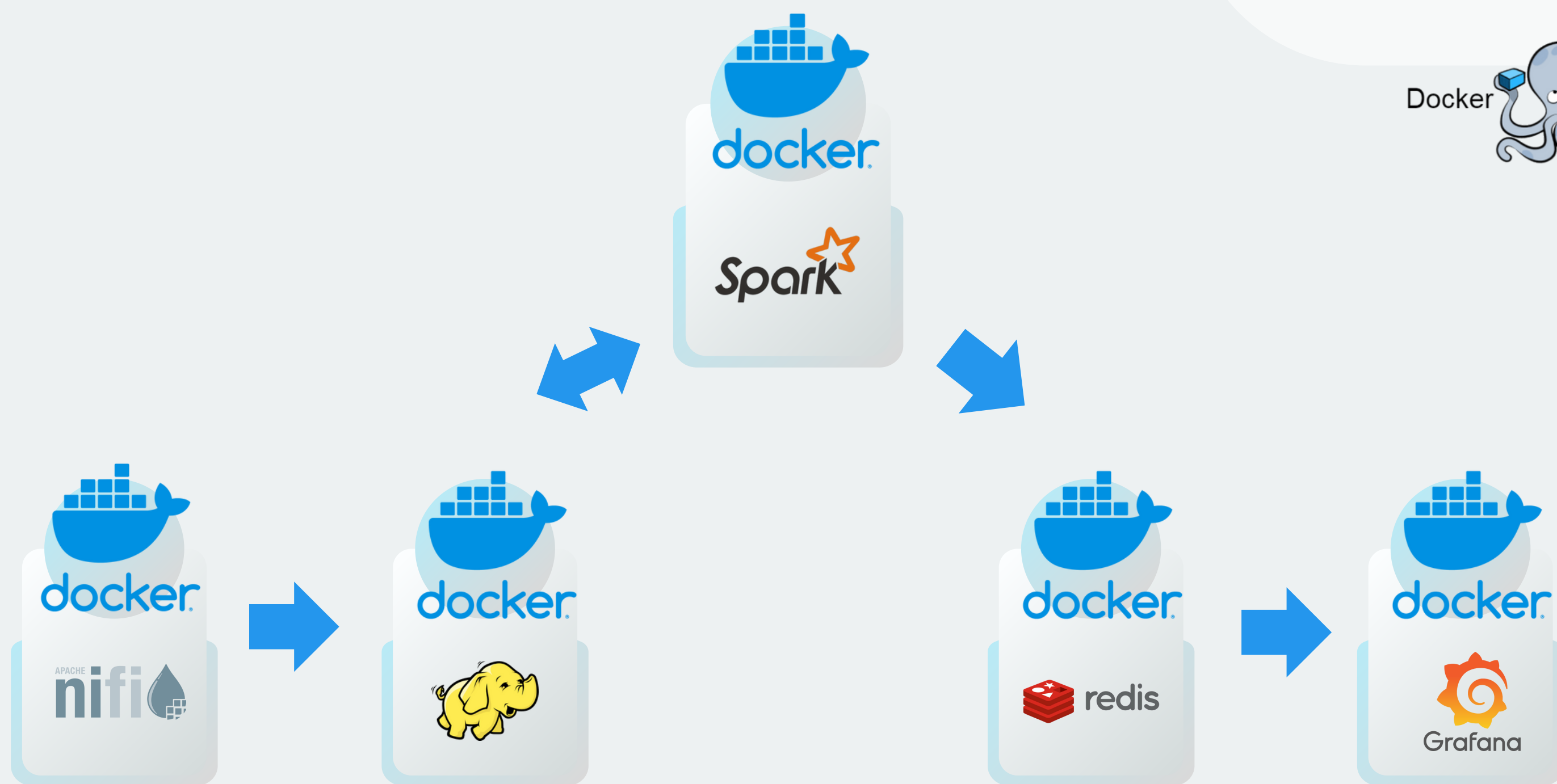
Marina Calcaura, Valerio Brauzi

# Dataset di riferimento

date	serial_number	model	failure	vault_id	s9 power on hours	...
2023-04-01T00:00:00.000000	8HK2SSMH	HGST HUH721212ALN604	0	1113	38445.0	...
2023-04-01T00:00:00.000000	10B0A01UF97G	TOSHIBA MG07ACA14TA	0	1067	27425.0	...
2023-04-01T00:00:00.000000	5080A117F97G	TOSHIBA MG07ACA14TA	0	1095	18029.0	...
2023-04-01T00:00:00.000000	ZL2NG0QS	ST16000NM001G	0	2010	6059.0	...
...	...	...	...	...	...	...

Dati forniti da Backblaze

## Architettura



# Data-acquisition & Data-ingestion

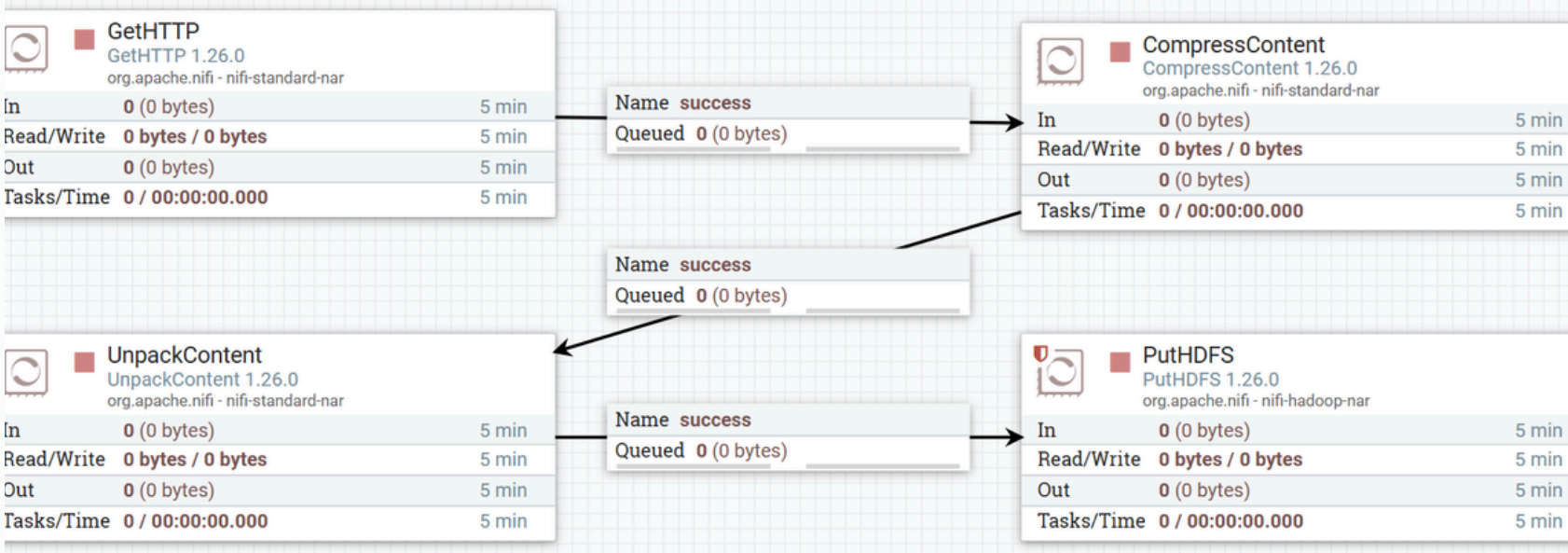


figura1: template per il Download e la Conversione dei dati

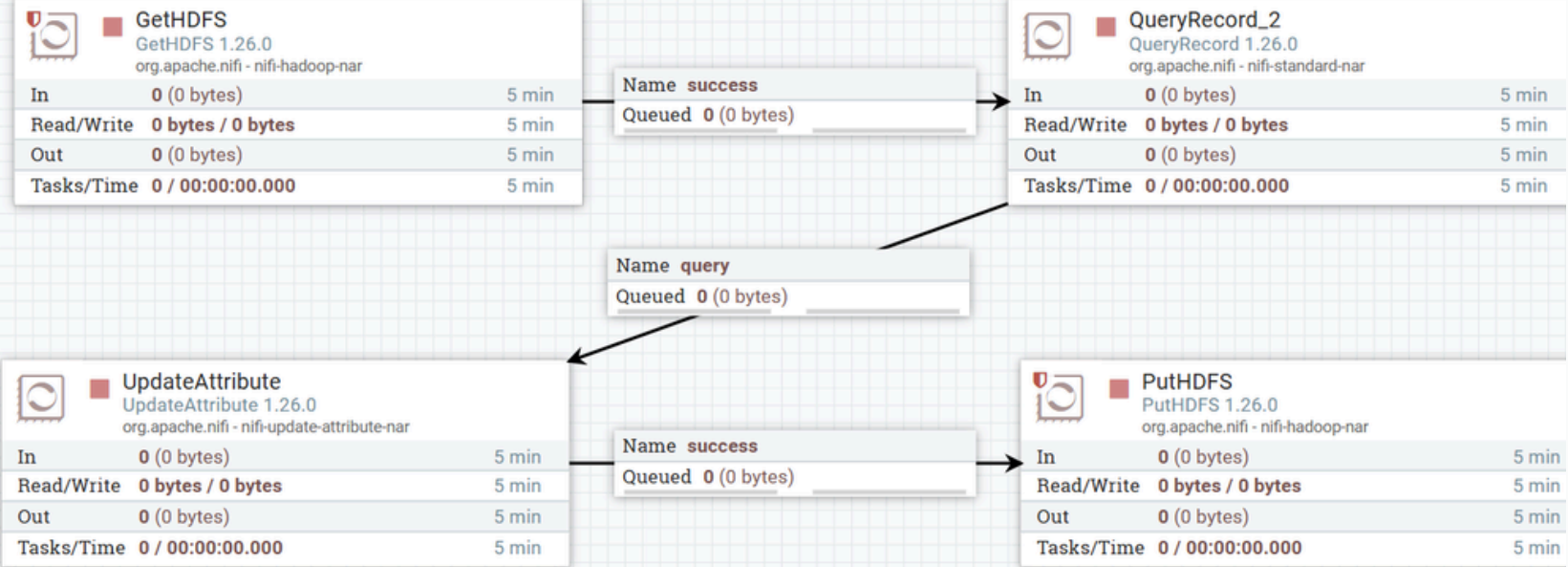
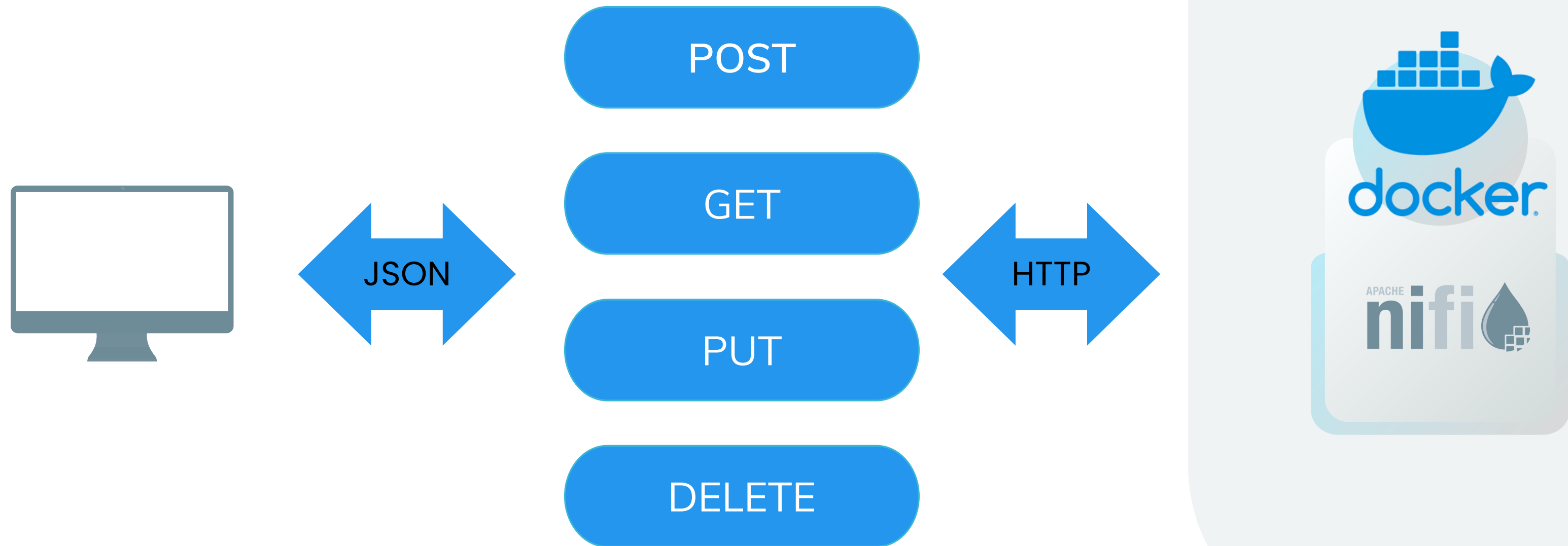


figura2: template per la Pulizia e il Filtraggio dei dati

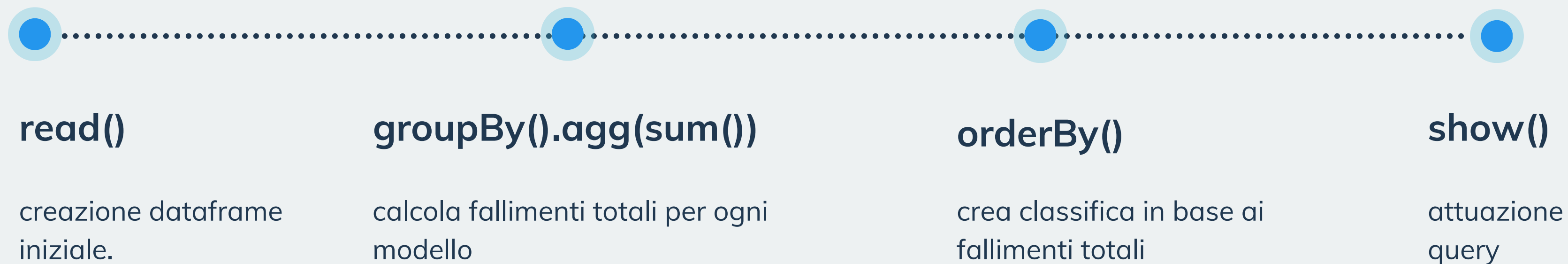
# REST API per l'automatizzazione del flusso



# Query1



# Query2.1



# Query2.2



**filter()**

dal dataframe  
iniziale filtra solo le  
entry con failure = 1

**groupBy().agg  
(count())**

raggruppa per  
'vault\_id' e calcola il  
numero totale di  
fallimenti per ogni  
vault

**groupBy().agg  
(collect\_set())**

raggruppa per  
'vault\_id',  
raccolgendo i modelli  
distinti di hard disk  
soggetti ad almeno un  
fallimento

**join()**

unisce i due  
dataframe  
precedenti

**orderBy()**

ordina il  
dataframe  
precedente per  
creare una  
classifica

**show()**

attuazione  
query



# Query3

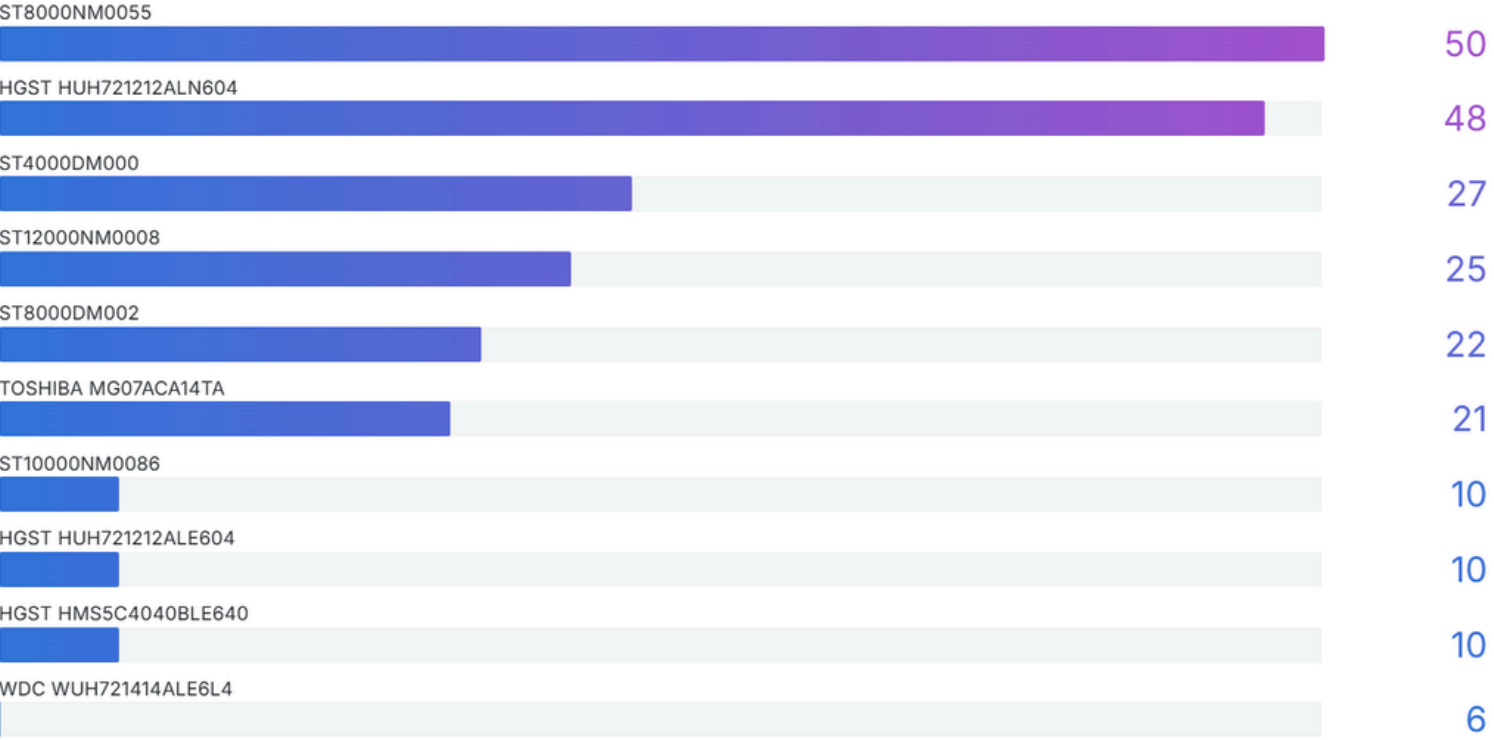


# Data-visualization

Query2 - Top 10 Vaults by Number of Failures

vaultId	totalFailures ↓	uniqueModels
1113	15	HGST HUH721212ALN604
1120	10	HGST HUH721212ALN604
1093	9	ST10000NM0086
1053	8	ST8000NM0055, TOSHIBA MQ01ABF050M
1118	8	HGST HUH721212ALN604
1032	7	ST8000DM002
1090	7	ST8000NM0055, WDC WD5000LPVX, TOSHIBA MQ01ABF050
1124	6	HGST HUH721212ALN604, HGST HUH721212ALE604
1066	6	TOSHIBA MG07ACA14TA
1055	6	ST8000NM0055

Top 10 Hard Disk Models by Number of Failures



Distribution of Operating Hours for Hard Disks with and without Failures

# failure	min	25th_percentile	50th_percentile	75th_percentile	max	count
1	522	27781	38702	51965	71608	265
0	0	15119	22650	42066	87702	242644

# Valutazione delle Performance

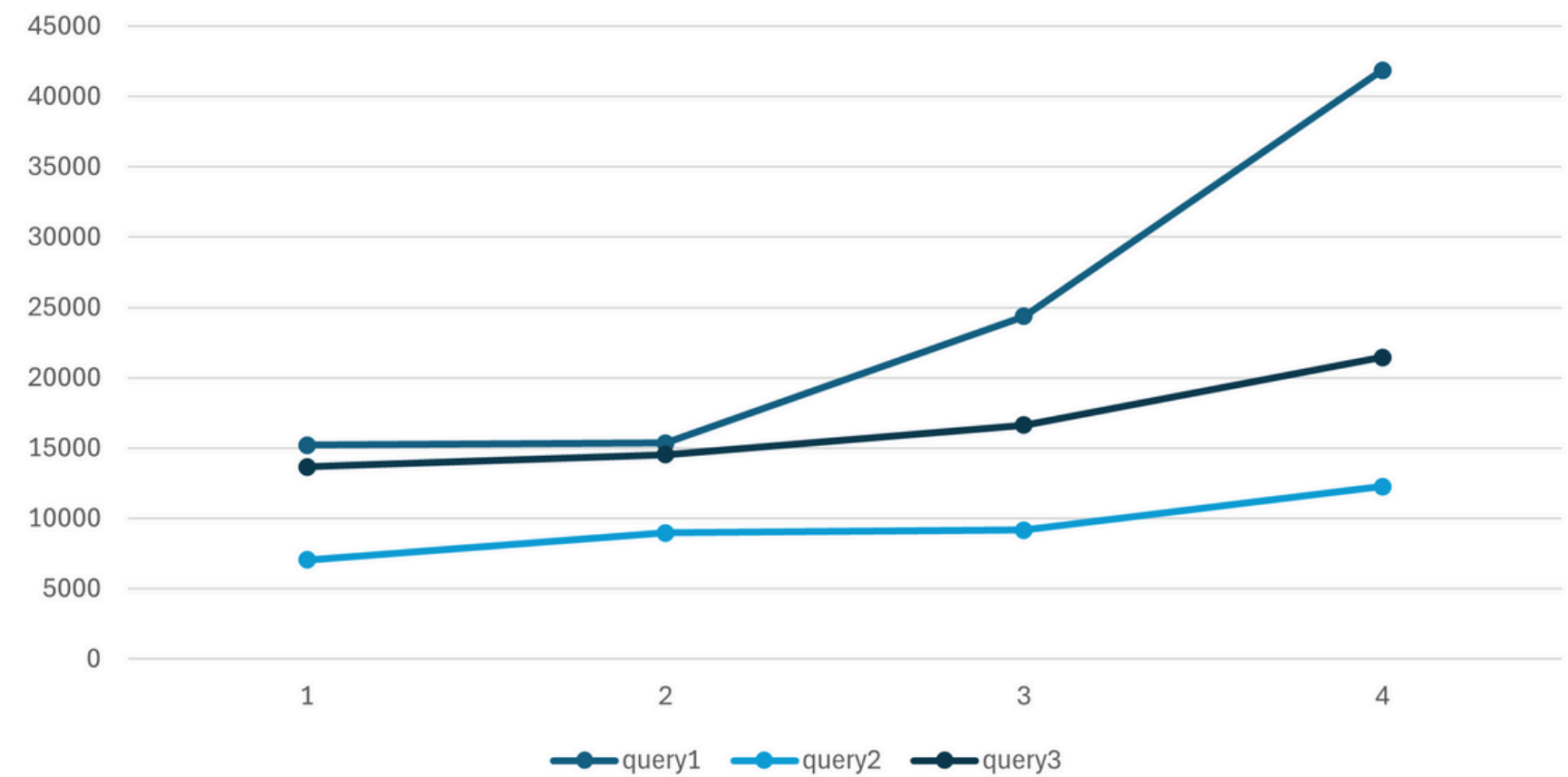
## Specifiche:

- 5 esecuzioni di ogni query
- configurazione con 1,2,3,4 worker

## Caratteristiche macchina:

- {Processore}: 11th Gen Intel(R) Core(TM) i7-1165G7 @ 2.80GHz, 2803 Mhz, 4 core, 8 processori logici
- {Memoria fisica installata (RAM)}: 16,0 GB
- {SO}: Microsoft Windows 11 Pro

Tempi di esecuzione in ms all'aumentare dei worker



**Grazie per l'attenzione!**