

Al supermercato con Azure Data Explorer

DATA
SATURDAYS



Sponsors



With the support of:



About me

Nicola Paro

Cloud Solutions Engineer

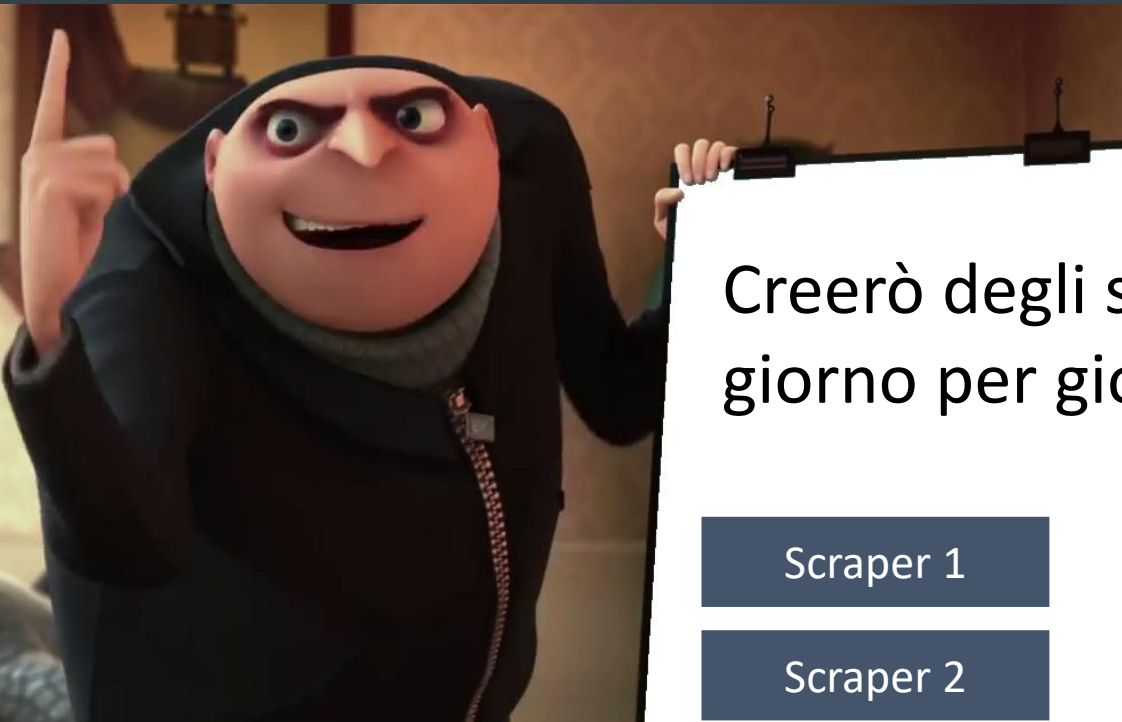
Linkedin: [linkedin.com/in/nicolaparo](https://www.linkedin.com/in/nicolaparo)

Github: github.com/nicolaparo

Twitter: [@nicola_paro](https://twitter.com/nicola_paro)



A Real Life Story



Creerò degli scraper per leggere i prezzi dei prodotti giorno per giorno.

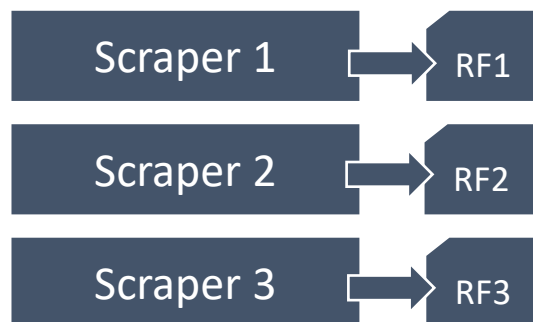
Scraper 1

Scraper 2

Scraper 3

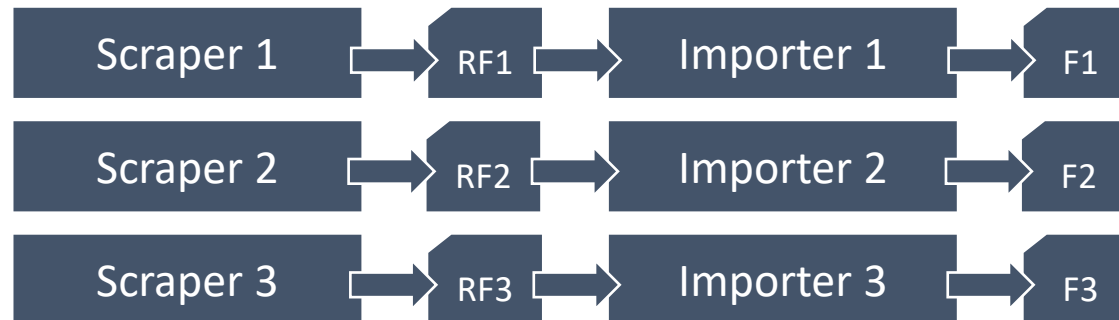


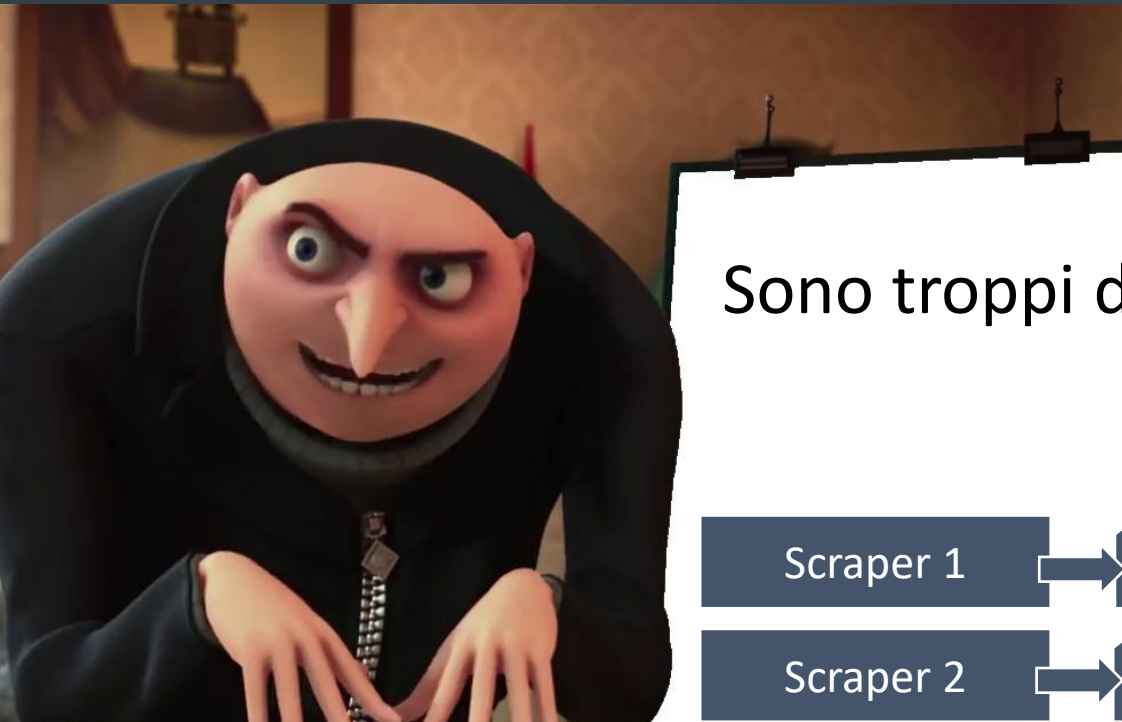
Salvo i dati estratti in alcuni RawFiles in formato TXT



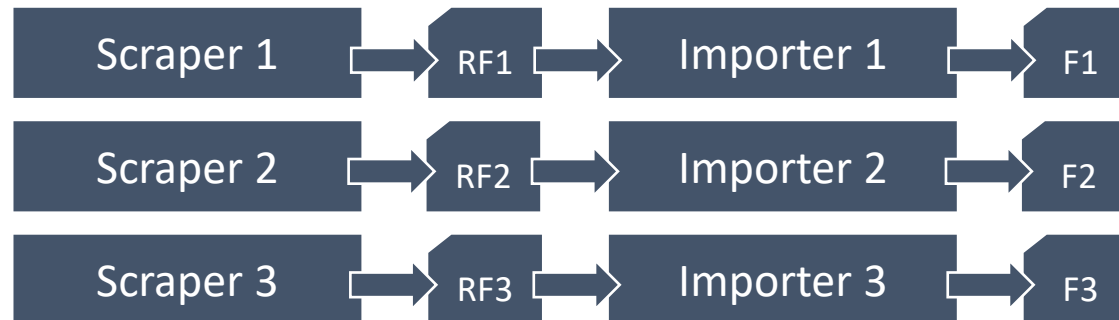


I miei importer estrarranno i dati rilevanti e li persisteranno in un formato condiviso



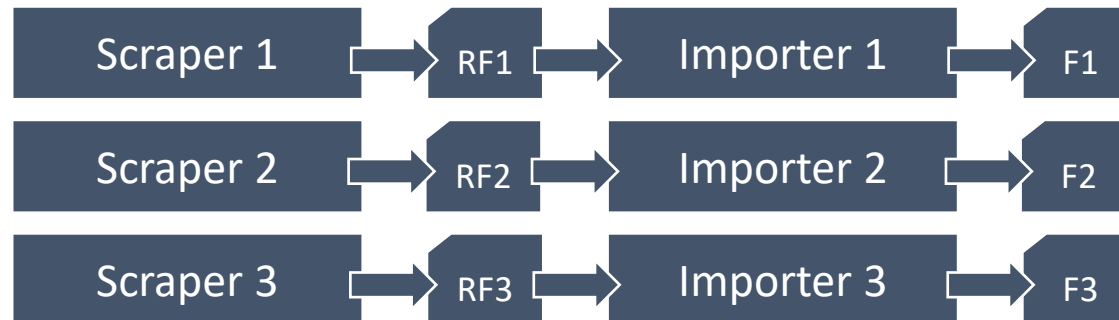


Sono troppi dati e non riesco a leggerli.



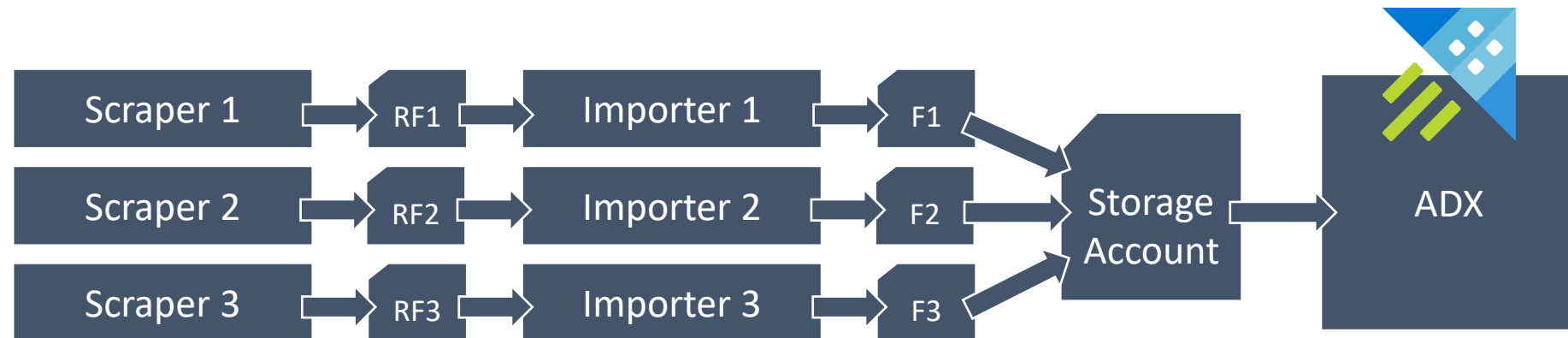


Sono troppi dati e non riesco a leggerli.





Caricherò i files su uno Storage Account e li analizzerò con Azure Data Explorer



- **Controlla prima i termini di licenza** del sito web
- Non fare scraping di dati protetti da copyright
- Non fare scraping di contenuti protetti da login
 - Non fare scraping dopo aver fatto una login
- Analizza solamente informazioni disponibili pubblicamente
- Controlla il **/robots.txt** per essere sicuro di poter analizzare un determinato path
- Controlla la **/sitemap.xml** per un elenco agevolato delle pagine da analizzare

Azure Data Explorer



Azure Data Explorer è un servizio di analisi dati veloce e totalmente gestito per l'analisi in tempo reale di grandi volumi di dati, provenienti da applicazioni, siti web, dispositivi IoT ecc...

With the support of:

Features Principali

Fast Data
Ingestion

Interactive
Data
Exploration

Real Time
Analytics on
Streaming data

Scalability

Integration
with other
Azure Services

Quando ha senso utilizzare ADX?

Sliding
Window of
data

Tante letture

Tanti Insert /
Append

Poche Delete

NESSUN
Update

Pricing

ADX Cost + VMs Cluster Cost + Storage Cost

Dev/Test
Free (No SLA)

Standard
\$0.11/core per
hour

Varies on VM
Size

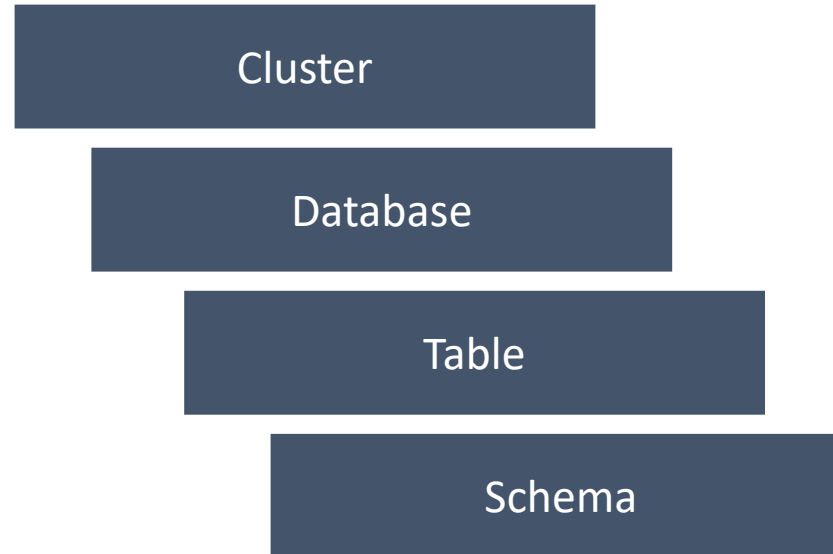
\$0.126/hour for
the smallest VM

Varies by usage

More info on <https://azure.microsoft.com/en-us/pricing/details/data-explorer/>

Organizzazione dei dati

Simile ai database
relazionali (ex: SqlServer)



Organizzazione dei dati – Differenze rispetto ai RDMS

No Primary
Key

No Unique
Keys

No Foreign
Keys

Columnstore
Indexes

Data Sharding
(Extents)

Data Ingestion



With the support of:

Dietro le quinte

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Dietro le quinte

I dati delle tabelle sono divisi in extents
(aka shards, partizioni, ...)

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla

Dietro le quinte

I dati delle tabelle sono divisi in extents
(aka shards, partizioni, ...)

Un extent è una mini-tabella che
contiene dati e metadati.

Un extent **non può mai essere
modificato**, ma può essere cancellato

I dati sono organizzati in colonne

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla

Dietro le quinte

I dati delle tabelle sono divisi in extents
(aka shards, partizioni, ...)

Un extent è una mini-tabella che
contiene dati e metadati.

Un extent **non può mai essere
modificato**, ma può essere cancellato

I dati sono organizzati in colonne

Extent più piccoli possono essere uniti
in extent più grandi

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Dietro le quinte

Gli extent sono creati durante le operazioni di inserimento

Un extent è unito ad altri

- Shard rebuild
- Shard merge

Un extent può essere cancellato con una retention-policy

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Dietro le quinte

Gli shards sono distribuiti tra i nodi del cluster

Node 1

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Node 2

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Una “semplice query” nel cluster

Logs

| where Timestamp > ago(1h)

Node 1

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Node 2

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Una “semplice query” nel cluster

Logs

| where Timestamp > ago(1h)

Node 1

Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

Node 2

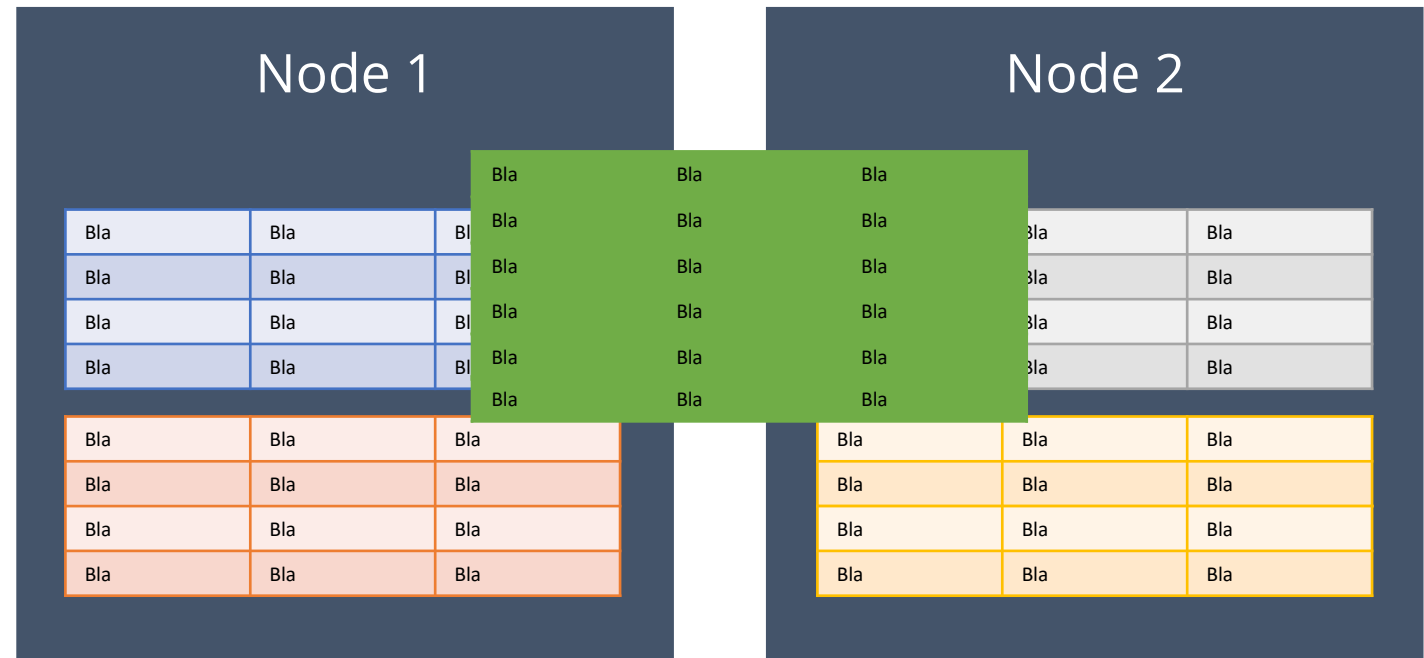
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla
Bla	Bla	Bla

With the support of:

Una “semplice query” nel cluster

Logs

| where Timestamp > ago(1h)



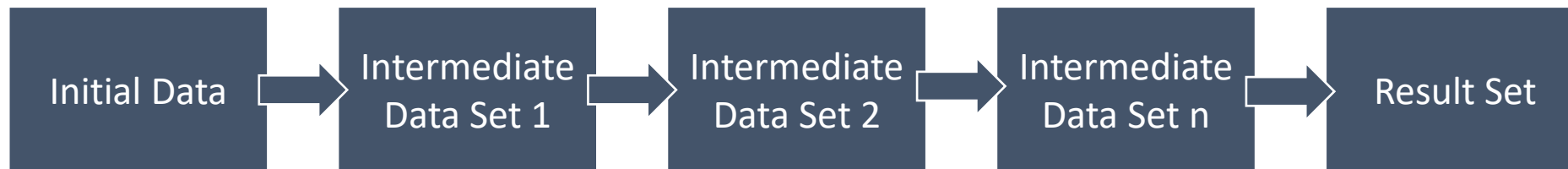
With the support of:

Kusto Query Language

Una query Kusto è una richiesta in sola lettura per il processing dei dati e la produzione di risultati.

La richiesta è effettuata tramite testo, utilizzando un modello di data-flow che è semplice da leggere, scrivere ed automatizzabile.

Le query Kusto sono composte di una o più istruzioni.



With the support of:

Kusto Query Language

Una query Kusto è una richiesta in sola lettura per il processamento dei dati e la produzione di risultati.

La richiesta è effettuata tramite testo, utilizzando un modello di data-flow che è semplice da leggere, scrivere ed automatizzabile.

Le query Kusto sono composte di una o più istruzioni.



With the support of:

Kusto Query Language

SQL	KQL
SELECT	project , extend, project-away, project-keep ...
WHERE	where , search, ...
JOIN	join kind=inner
UNION	union
GROUP BY	summarize
ORDER BY	sort by , order by, top by
TOP, LIMIT	take

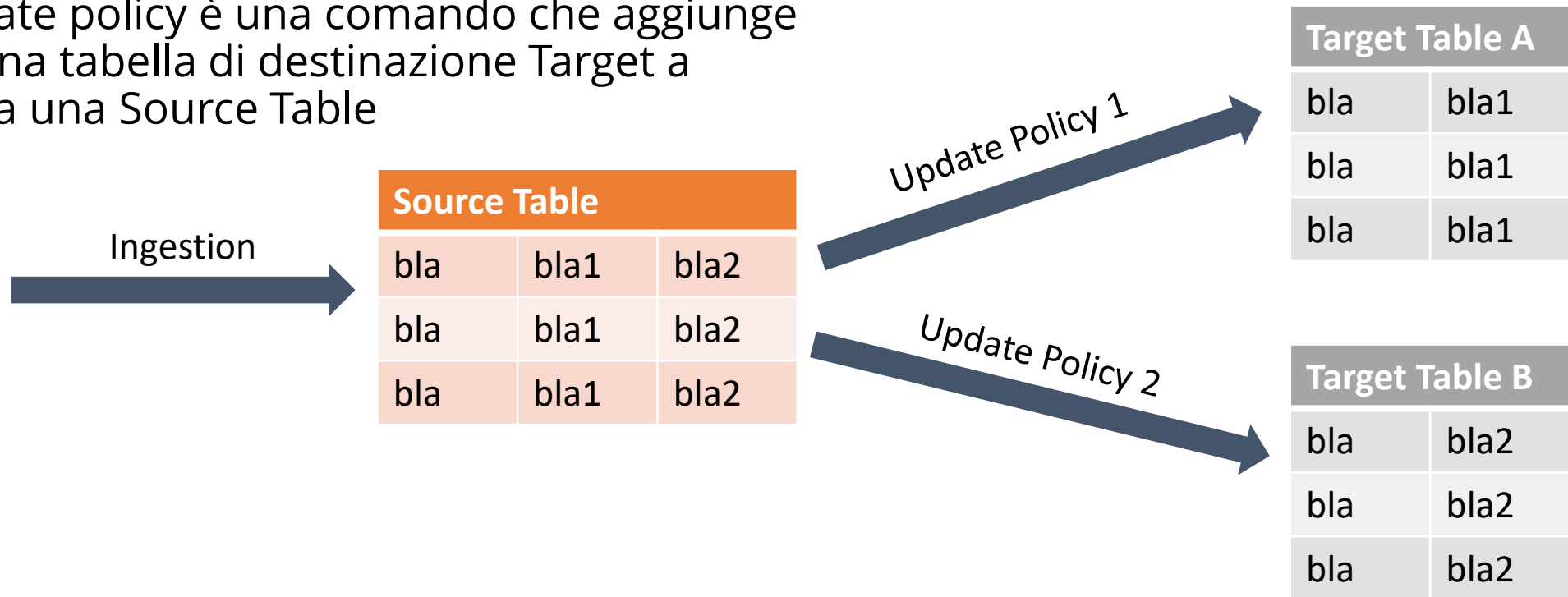
More on <https://learn.microsoft.com/en-us/azure/data-explorer/kusto/query/>

With the support of:



Update Policies

Una update policy è un comando che aggiunge dati ad una tabella di destinazione Target a partire da una Source Table



With the support of:

Materialized View

Una Materialized View è una vista aggregate sui dati di una tabella ADX. I dati sono materializzati anche su disco.

Vantaggi nell'adozione delle Materialized View

Performance
Improvement

Data
Freshness

Cost
Reduction

With the support of:

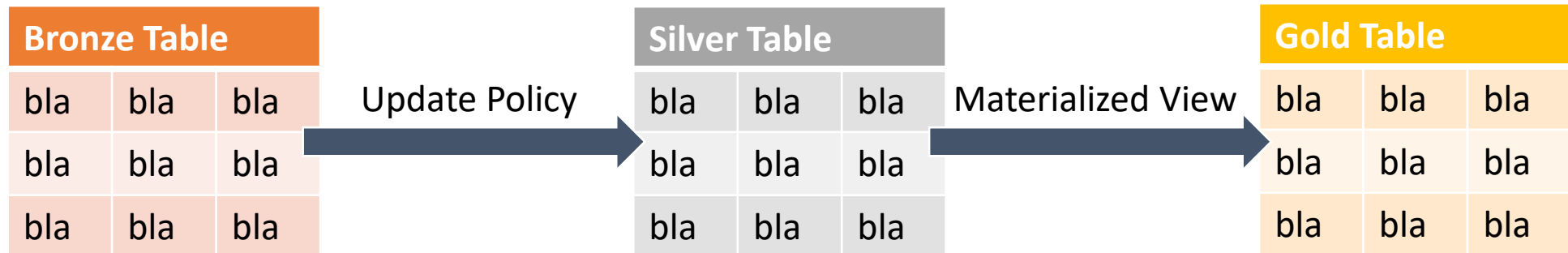
Update Policy o Materialized View?

Update Policy

- Data Transformation
- Data Enrichment

Materialized View

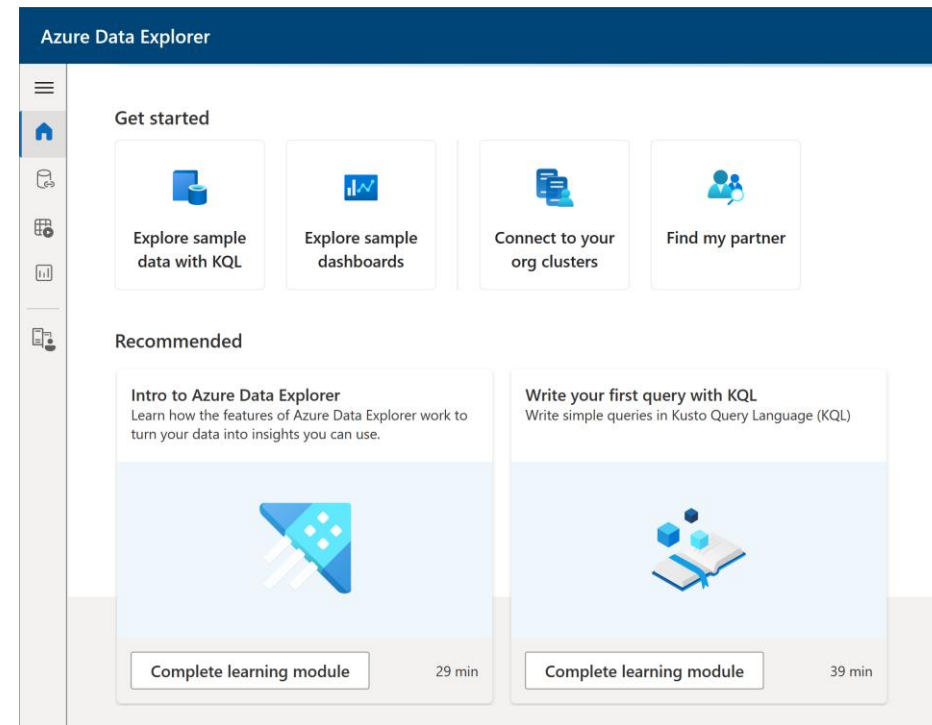
- Data Aggregation



Devo pagare un cluster per fare pratica con Kusto?

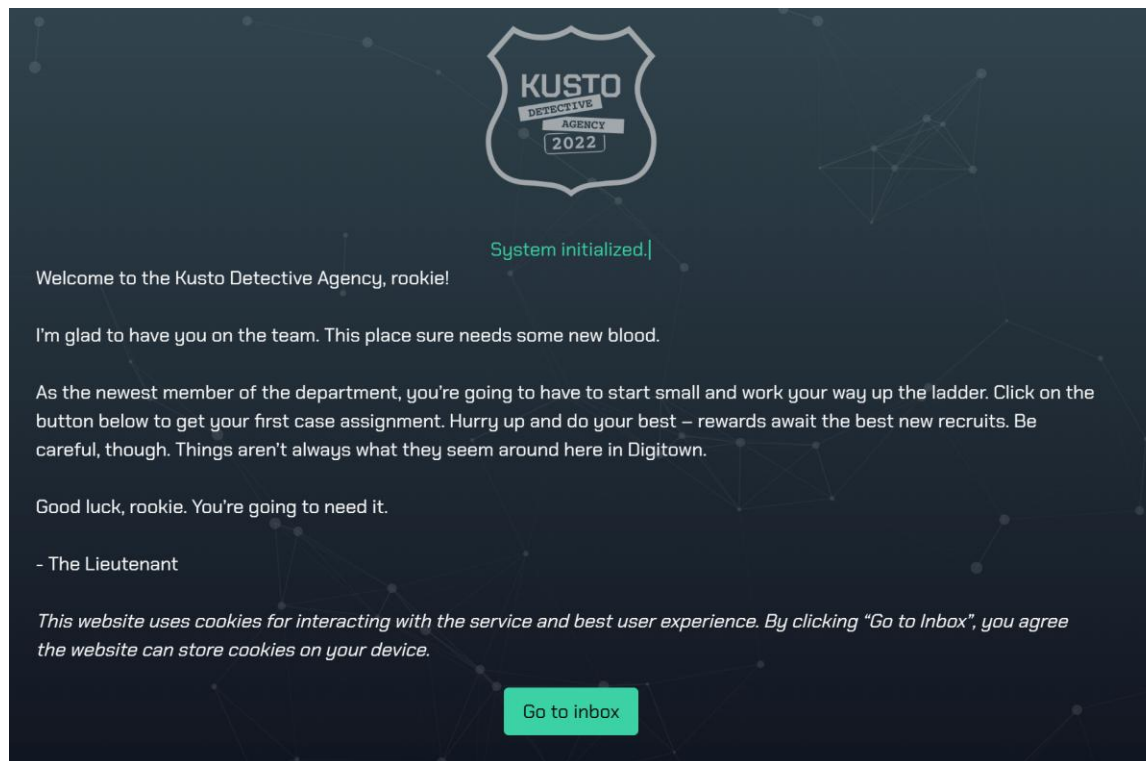
Azure Data Explorer supporta dei database “sample” gratuiti su cui è possibile effettuare delle interrogazioni per provare

<https://dataexplorer.azure.com/home>



With the support of:

Devo pagare un cluster per fare pratica con Kusto?



Kusto Detective Agency: una gamification per imparare ad usare kusto.

<https://detective.kusto.io/>

Q & A

Thank You!

Nicola Paro

Cloud Solutions Engineer

