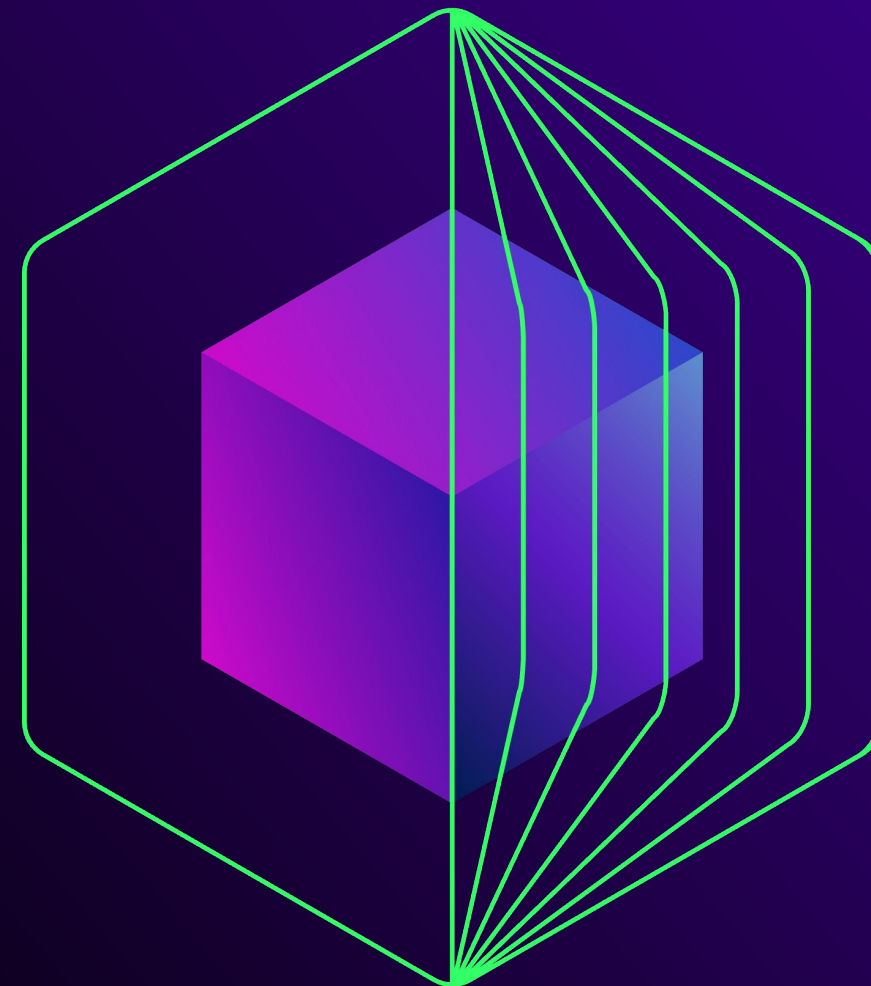


**DATA  
SCIENCE  
SUMMIT**

# Integrating LLM with Streaming Data Pipelines

**Tim Spann**  
Principal Dev, Cloudera



[www.dssconf.pl](http://www.dssconf.pl)



23-24.11.2023



PGE Narodowy + Online

ORGANIZERS:

ACADEMIC PARTNERS



Wydział Matematyki  
i Nauk Informatycznych  
POLITECHNIKA WARSZAWSKA

# FLaNK Stack Weekly by Tim Spann

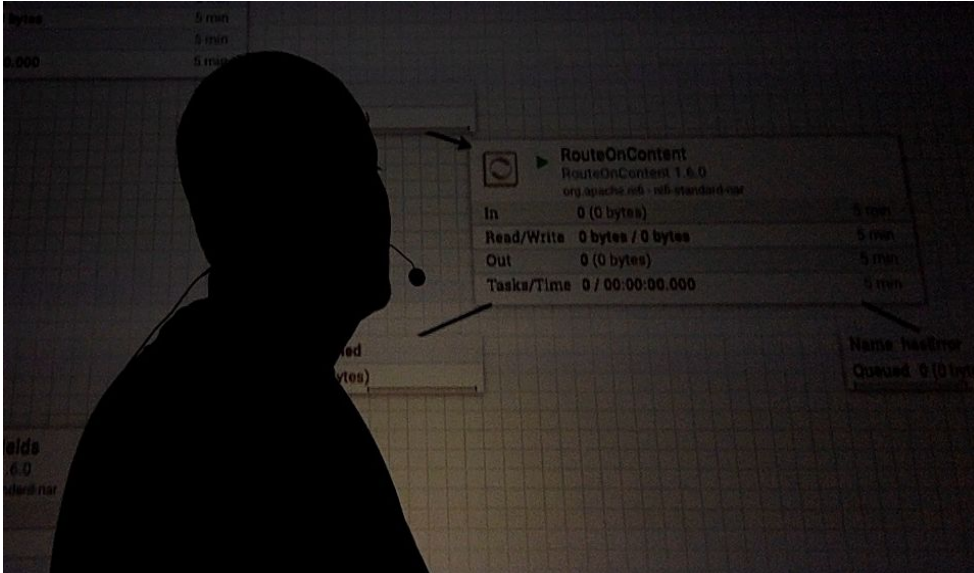


<https://bit.ly/32dAJft>

<https://www.meetup.com/futureofdata-princeton/>

**This week in Apache NiFi, Apache Flink, Apache Kafka, ML, AI, Apache Spark, Apache Iceberg, Python, Java and Open Source friends.**

# Agenda (35 minutes)



**Build a Streaming Data Pipeline**

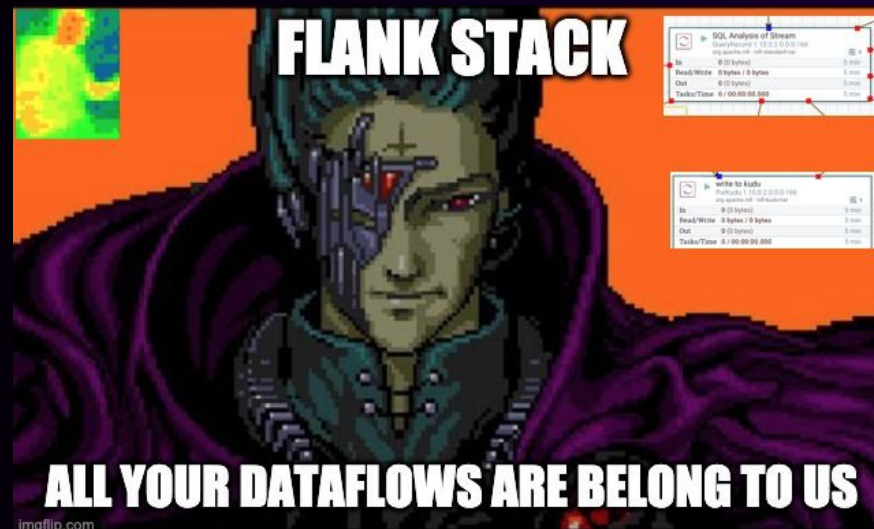
**Integrate LLM via HTTPS / REST**

**Process, Route, Enrich, Transform Results**

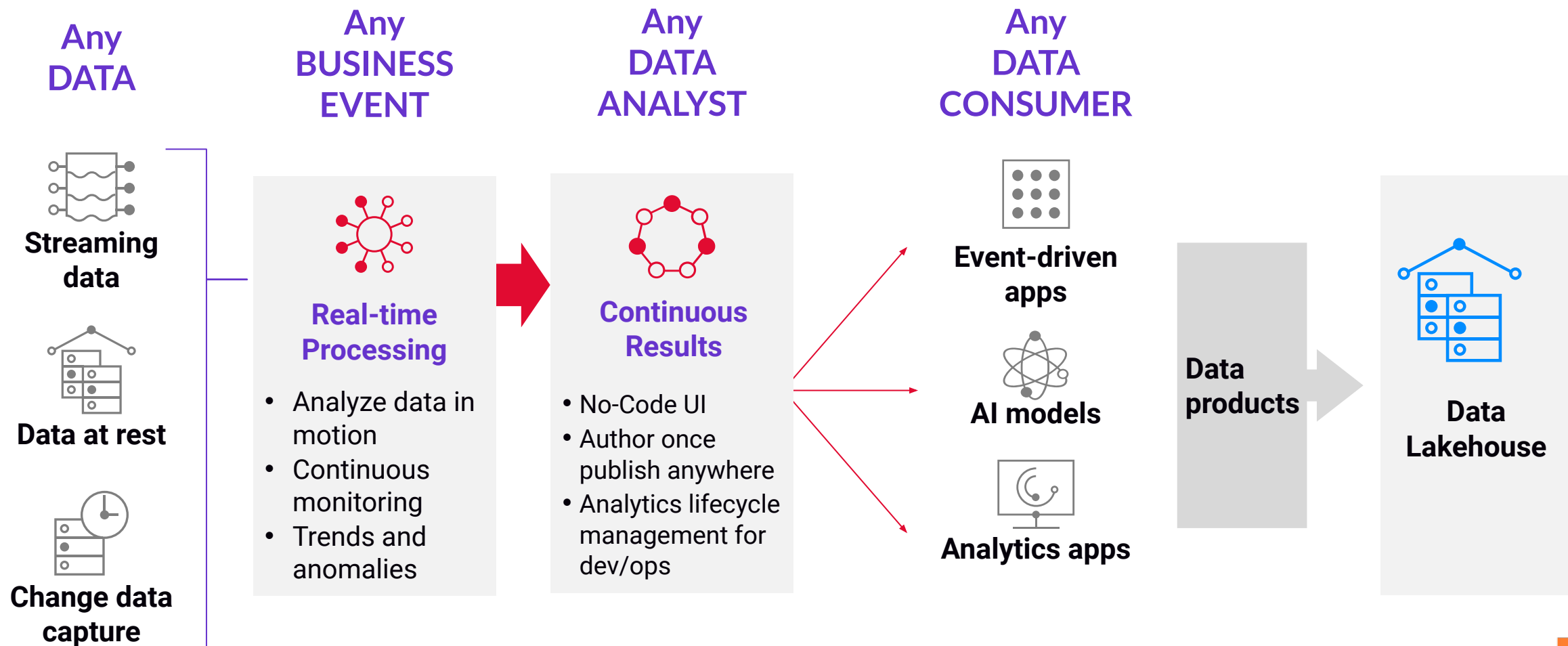
**Store**

**Distribute**

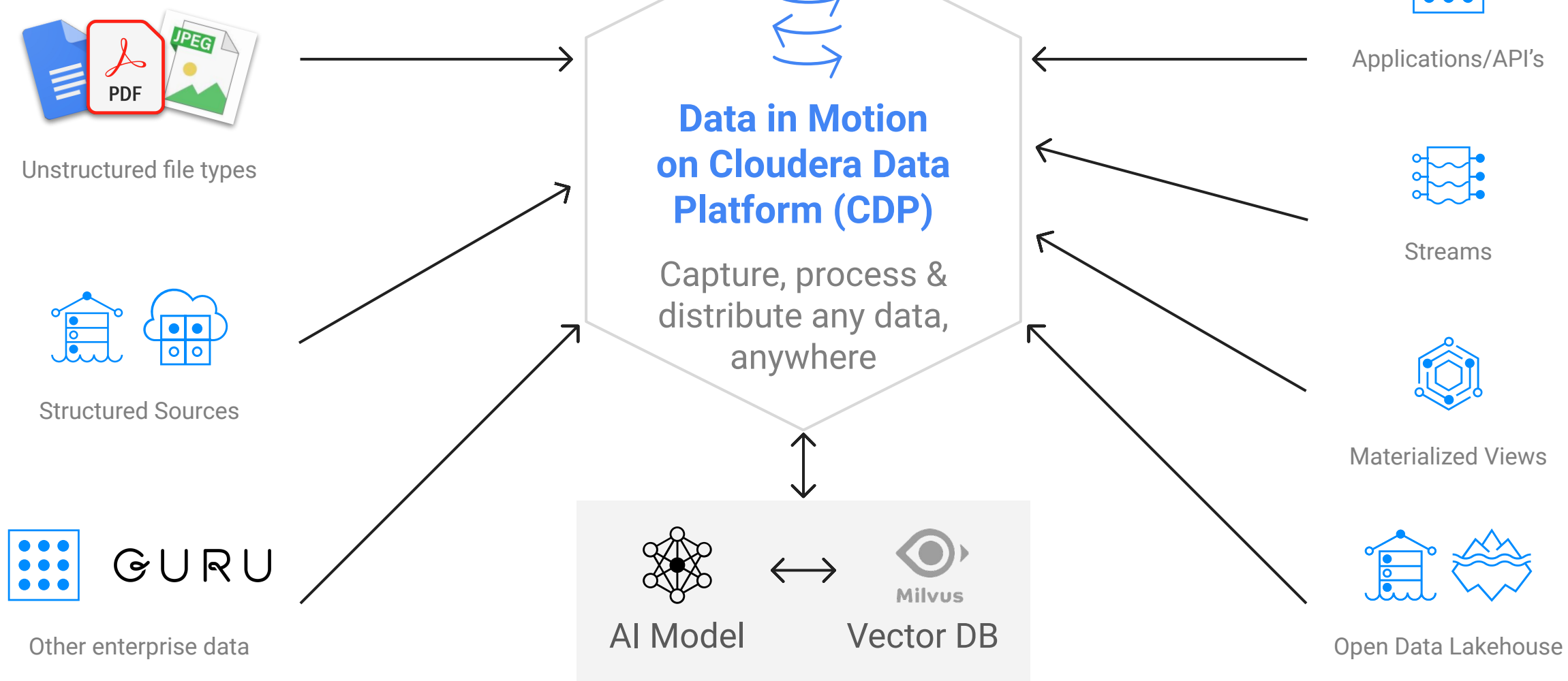
# Build a Streaming Data Pipeline



# Data Relevance



# LLM USE CASE



# Integrating LLM via HTTPS



# Models Used for Inference API

google/flan-ul2

google/flan-t5-xxl

bigscience/bloom

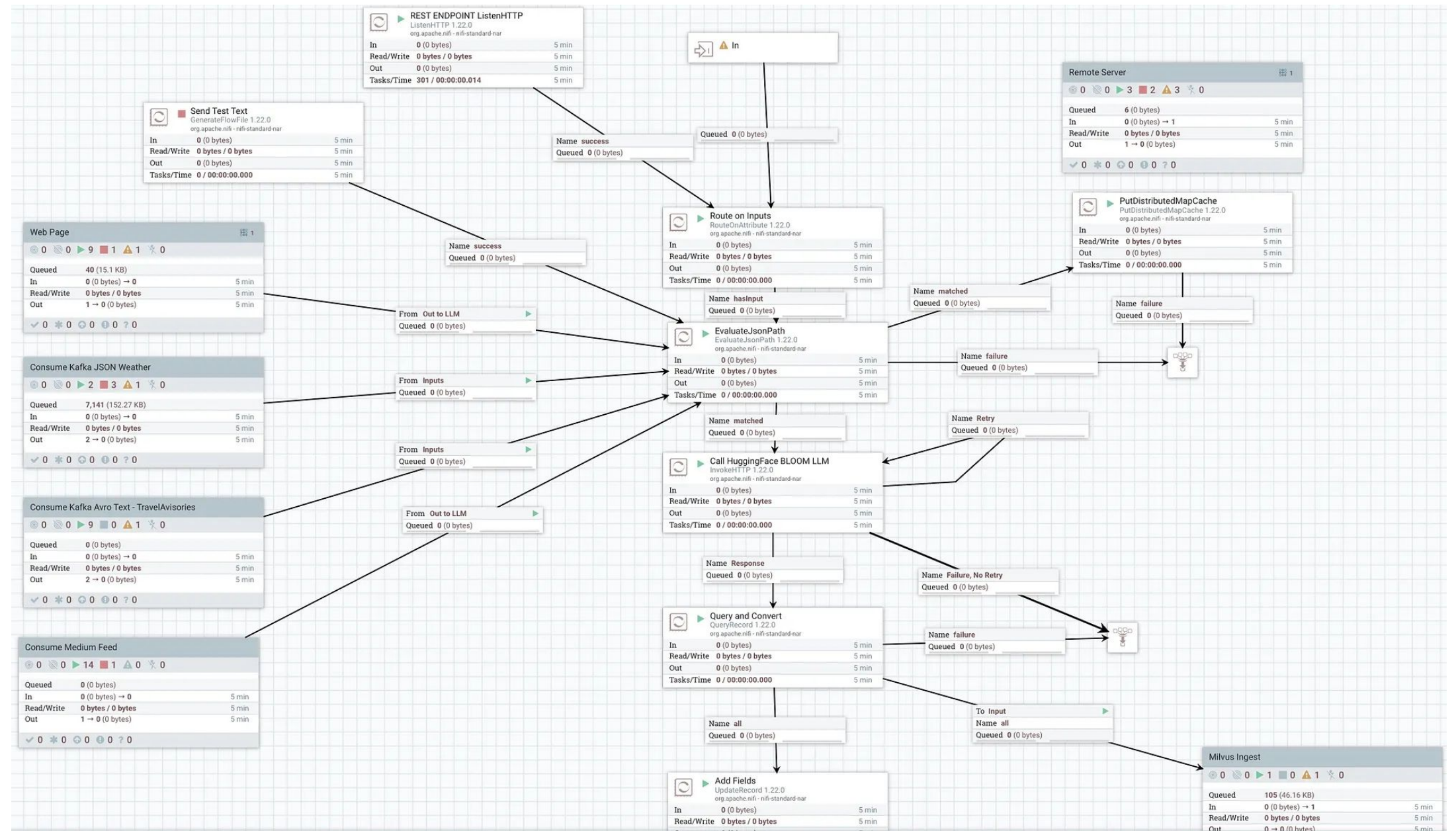
meta-llama/llama-2-70b-chat

ibm/mpt-7b-instruct2

ibm/granite-13b-instruct-v1

ibm/granite-13b-chat-v1





# HTTPS Calls

For Hugging Face:

<https://api-inference.huggingface.co/models/bigscience/bloom>

Pass in a simple JSON Prompt

```
{"inputs": "This is Prompt Text"}
```

Requires an HTTP Header with Authorization for Bearer YOURTOKEN

YOURTOKEN from <https://huggingface.co/settings/tokens>

# HTTPS Calls

For IBM WatsonX.AI:

<https://iam.cloud.ibm.com/identity/token> <- First call to get token by sending an API key and grant.

Then send a JSON prompt call model server via  
<https://us-south.ml.cloud.ibm.com/ml/v1-beta/generation/text?version=2023-05-29>

Requires an HTTP Header with Authorization for Bearer \${readtoken}

# WatsonX.AI JSON Prompt

```
{ "model_id": "meta-llama/llama-2-70b-chat",  
  "input": "this is prompt text",  
  "parameters": {  
    "decoding_method": "greedy",  
    "max_new_tokens": 200,  
    "min_new_tokens": 50,  
    "stop_sequences": [],  
    "repetition_penalty": 1  
  },  
  "project_id": "0ead8ec4-d137-4f9c-8956-50b0da4a7068" }
```

Property		Value
HTTP Method	?	GET
HTTP URL	?	<a href="https://rss.nytimes.com/services/xml/rss/nyt/US.xml">https://rss.nytimes.com/services/xml/rss/nyt/US.xml</a>
HTTP/2 Disabled	?	False
SSL Context Service	?	No value set
Socket Connect Timeout	?	5 secs
Socket Read Timeout	?	15 secs
Socket Idle Timeout	?	5 mins
Socket Idle Connections	?	5
Proxy Configuration Service	?	No value set
Proxy Host	?	No value set
Request OAuth2 Access Token Provider	?	No value set
Request Username	?	No value set



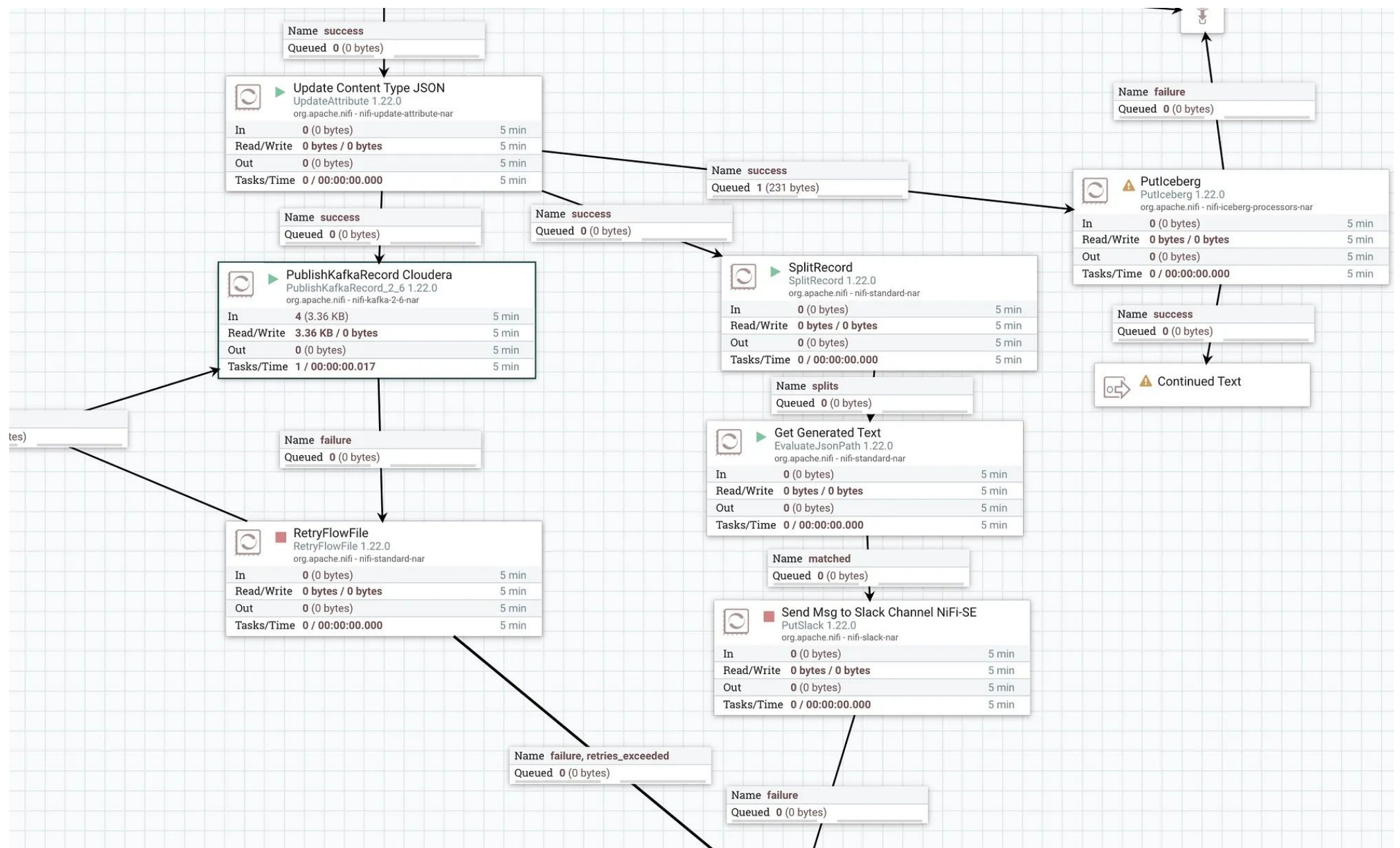
**Process, Route, Enrich,  
Transform Results**

Property		Value
Record Reader	?	CDP Infer JsonTreeReader
Record Writer	?	BasicJsonRecordSetWriter
Include Zero Record FlowFiles	?	false
Cache Schema	?	false
Default Decimal Precision	?	10
Default Decimal Scale	?	4
all	?	SELECT * FROM FLOWFILE
negative	?	SELECT lower(label) as sentiment, label, score...
positive	?	SELECT lower(label) as sentiment, label, score...



Property		Value
<b>Record Reader</b>	?	<b>CDP Infer JsonTreeReader</b>
<b>Record Writer</b>	?	<b>JsonRecordSetWriter</b>
<b>Replacement Value Strategy</b>	?	<b>Literal Value</b>
/textbody	?	<code>\${inputs}</code>
/ts	?	<code>\${now():toNumber()}</code>
/uuid	?	<code>\${uuid}</code>

# Store



METRICS ASSIGNMENT DATA EXPLORER CONFIGS LATENCY

ISOLATION LEVEL:

read\_uncommitted

DESERIALIZER:

Keys: String

Values: String

Partition 0

FROM OFFSET

0

RECORD LIMIT

15

Offset

Timestamp

Key

Value

0

Wed, Aug 16 2023,  
11:39:23102572a6-ef9c-4a1e-92fd-  
47324be3c0a2{"generated\_text": "What is the best way to ingest syslog data into Apache Iceberg? I have a large amount of syslog data that I  
would like to store in [show more](#)

1

Wed, Aug 16 2023,  
22:41:582230d424-3def-4d99-98ee-  
eea8461c374f{"generated\_text": "Reissued with obsolete COVID-19 page links removed. Exercise normal precautions in  
Poland. Read the country infor [show more](#)

2

Wed, Aug 16 2023,  
22:41:58c4fd8327-c13e-4fc8-a32e-  
e07056b52025{"generated\_text": "Reissued with obsolete COVID-19 page links removed. Exercise increased caution while traveling to Moldova  
due to unresolved conflic [show more](#)

3

Wed, Aug 16 2023,  
22:41:58feebee51-5cee-4b88-a93e-  
7f791224fa22{"generated\_text": "Fairfax County Public Schools (FCPS) is the largest school district in Virginia, serving more than  
200", "ts": "1692216592446", "uuid" [show more](#)

4

Wed, Aug 16 2023,  
22:41:581c7aa5c4-7046-457d-9df6-  
3b59eb6b16b0{"generated\_text": "Apache NiFi is a great tool for loading data into Apache Kafka. It is a very simple and easy to use tool. It is a  
very simple and e [show more](#)

1 - 5 of 5 &lt; &gt;

localhost:9991/#



# Distribute

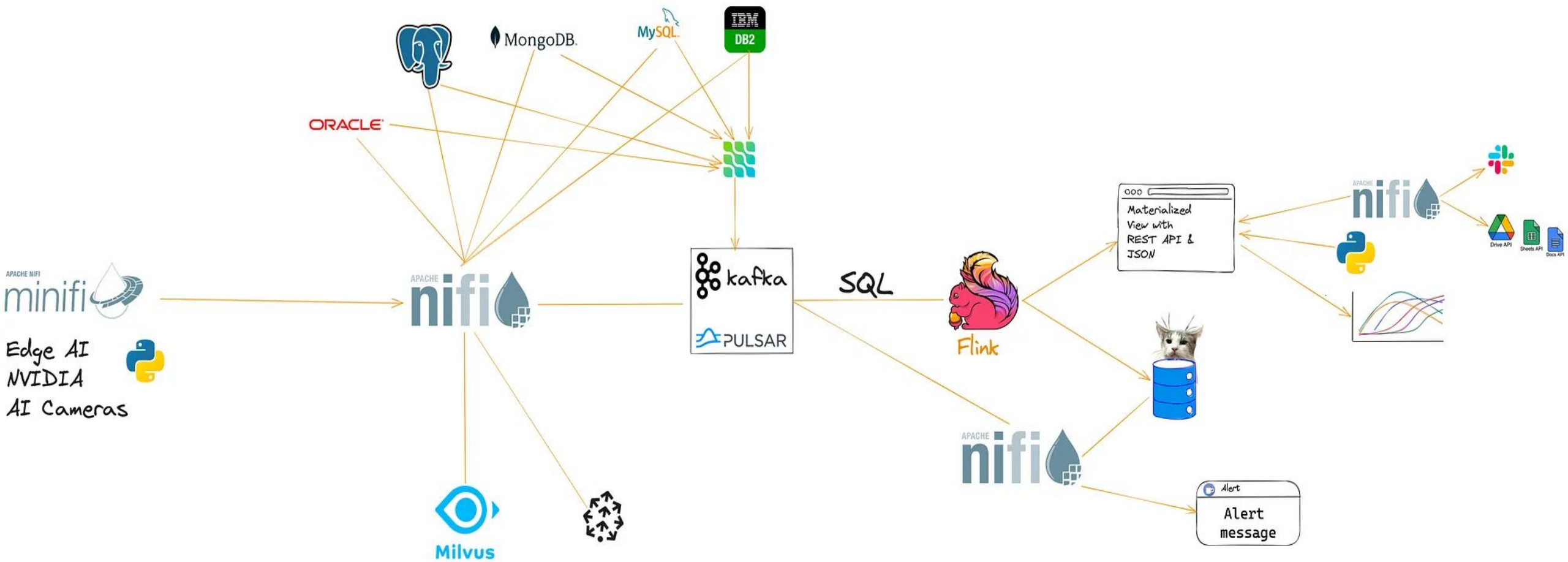
**CUNK ON FLANK**

**THEY SAY IT'S FOR FAST DATA**

**IS IT SO FAST THAT IT  
ARRIVES BEFORE I ASKED FOR IT**

imgflip.com





+ Add a bookmark

Compute Time: 905  
Compute Type: gpu+optimized  
Inference Time: 905  
Queue Time: 0  
Request ID: VjKhYJvX664mtgaljgymc  
SHA: 053d9cd9fbe814e091294f67fcfedb3397b954bb  
Time Per Token: 45  
Total Time: 905  
Validation Time: 0  
R: Fairfax County Public Schools (FCPS) is the largest school district in Virginia, serving more than 200  
=====

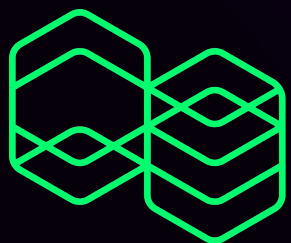
Today ▾

==== NiFi to Hugging Face BLOOM LLM  
AzureML Model Deployment: bloom-deployment  
On Date: Thu, 17 Aug 2023 02:05:56 GMT  
File Name: 7b8afb65-d5eb-49ee-9a89-d2b870310243  
Request Duration: 1339  
Request URL: <https://api-inference.huggingface.co/models/bigscience/bloom>  
Compute Characters: 63  
Compute Time: 1014  
Compute Type: gpu+optimized  
Inference Time: 968  
Queue Time: 45  
Request ID: iPX8lqeK9XwqSZA9e-wCH  
SHA: 053d9cd9fbe814e091294f67fcfedb3397b954bb  
Time Per Token: 48  
Total Time: 1014  
Validation Time: 0  
R: Apache NiFi is a great tool for loading data into Apache Kafka. It is a very simple and easy to use tool. It is a very simple and easy to use  
=====



# Resources

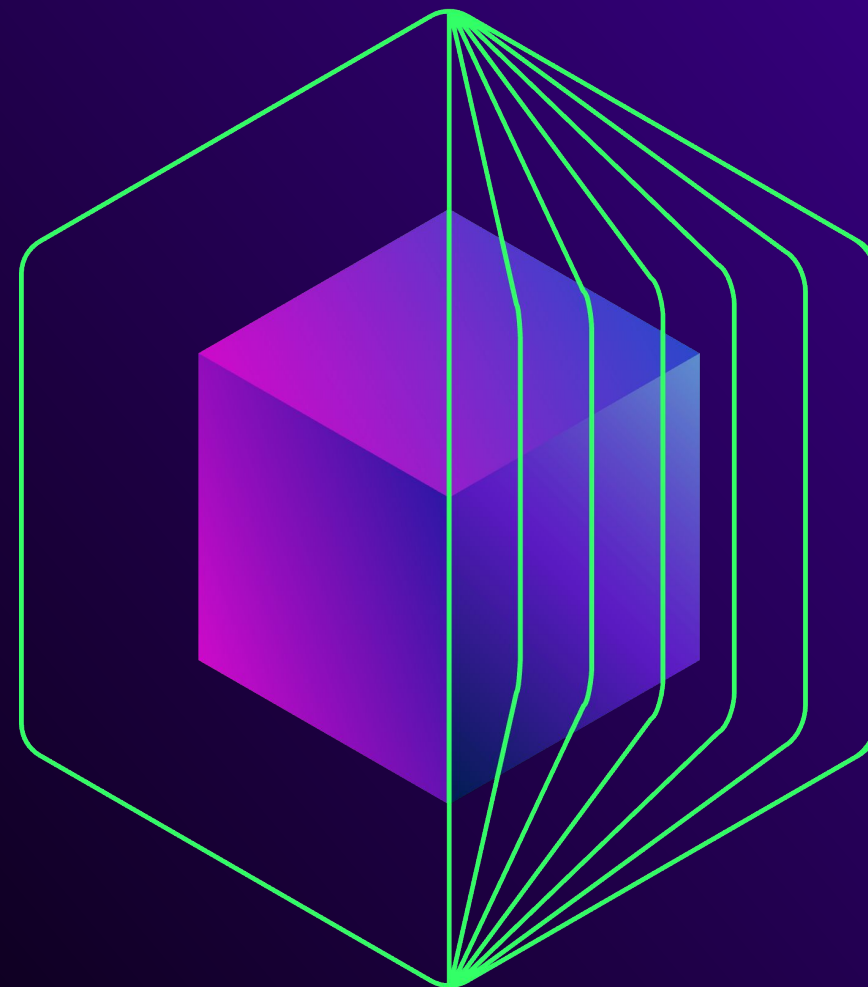




**DATA  
SCIENCE  
SUMMIT**

# Thank you for watching!

Remember to leave your **questions**  
and **rate** the presentation  
in the section below.



[www.dssconf.pl](http://www.dssconf.pl)



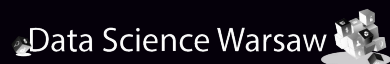
23-24.11.2023



PGE Narodowy + Online

ORGANIZERS:

ACADEMIC PARTNERS



Wydział Matematyki  
i Nauk Informatycznych  
POLITECHNIKA WARSZAWSKA