

# Adding Generative AI to Real-Time Streaming Pipelines



Tim Spann  
Principal Developer Advocate

Nov-2024





Introduction

Overview

GenAI Architecture

Streaming Projects

Demos

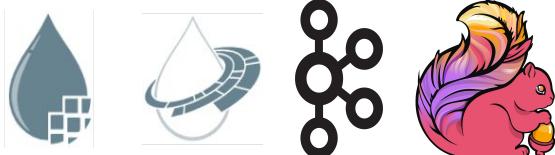
Resources

Q&A

# Tim Spann

Twitter: @PaasDev // Blog: [datainmotion.dev](http://datainmotion.dev)  
Principal Developer Advocate / Field Engineer  
NY AI Meetups  
ex-Pivotal, ex-Cloudera, ex-StreamNative,  
ex-PwC, ex-HPE, ex-E&Y.

<https://medium.com/@tspann>  
<https://github.com/tspannhw>



A screenshot of a DZone profile page titled "Top IoT Experts". It features a photo of Tim Spann, a link to his GitHub profile, and a brief bio stating he is a Principal Developer Advocate at Cloudera, working with Apache NiFi, Apache Pulsar, Apache... The page also includes links for RefCards and Trend Reports.



# AI + Streaming Weekly by Tim Spann



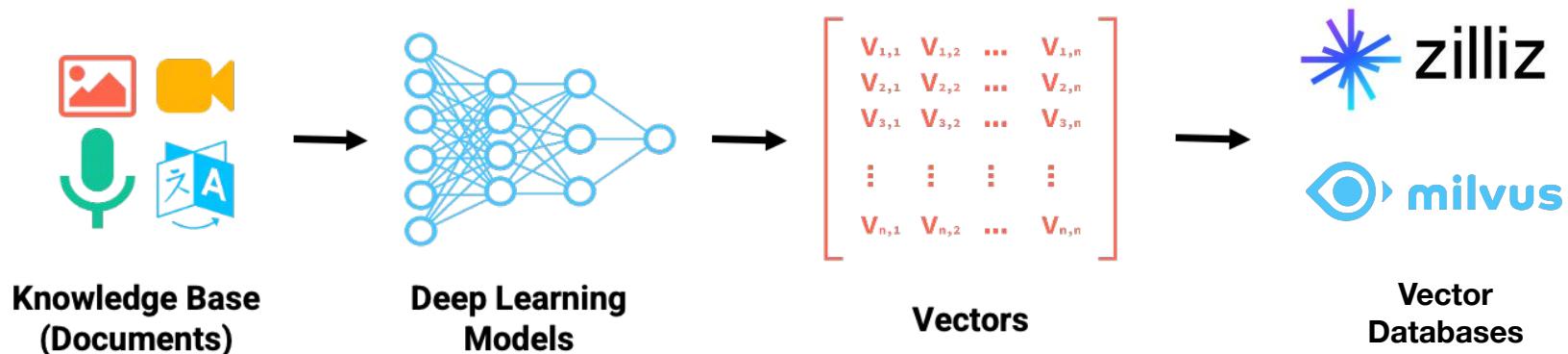
<https://bit.ly/32dAJft>

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

This week in Apache NiFi, Apache Flink, Apache Kafka, ML, AI, Apache Spark, Apache Iceberg, Python, Java, LLM, GenAI, Vector DB and Open Source friends.

# The challenge of Unstructured Data

- **Problem:** Unstructured data comes in lots of forms, no easy way to interact with it all
- **Solution:** Vector embeddings
- **How:** Neural networks e.g. embedding models



# Unstructured Data is Everywhere

Unstructured data is any data that does not conform to a predefined data model.



Currently, 90% of unstructured data is never analyzed.



Text



Images



Videos



and  
more!

# Vector Search Overview

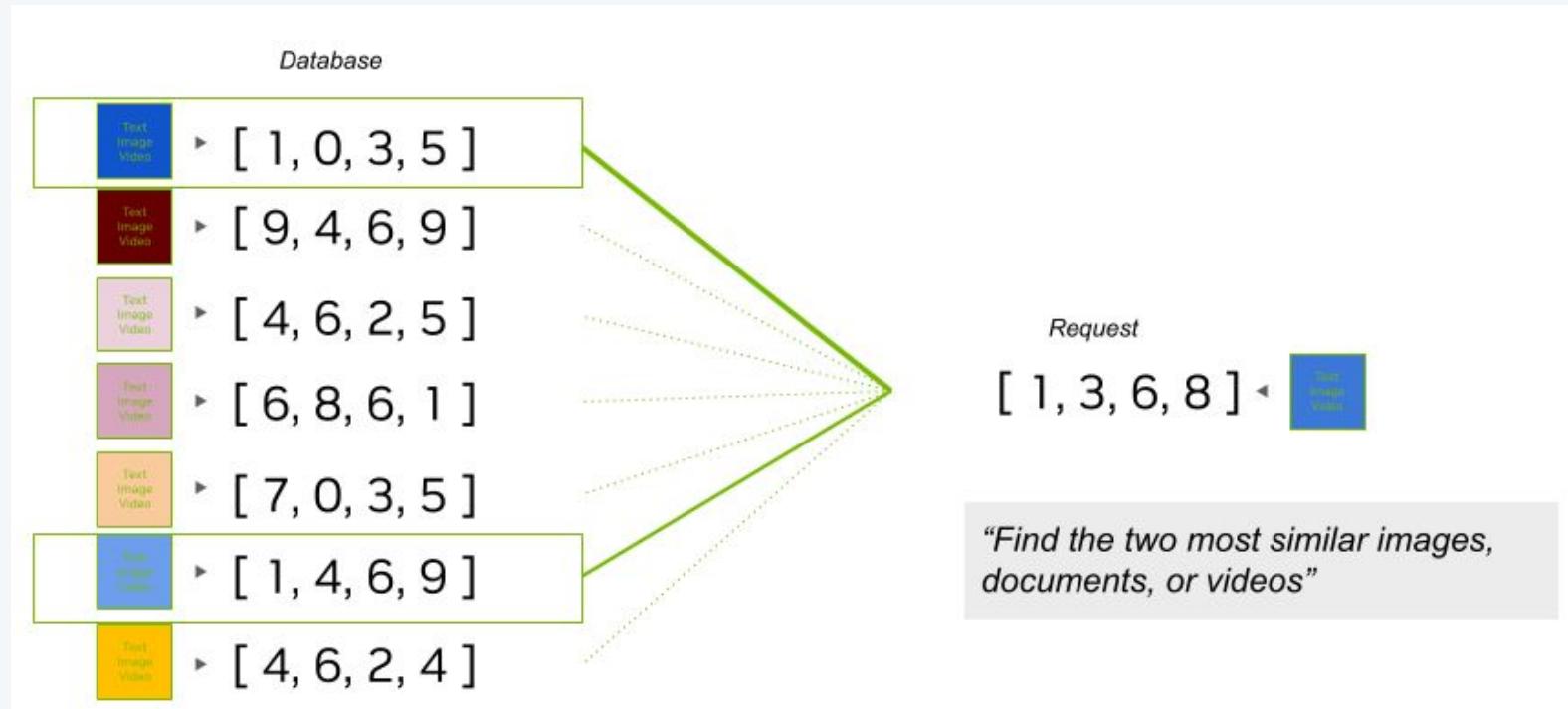
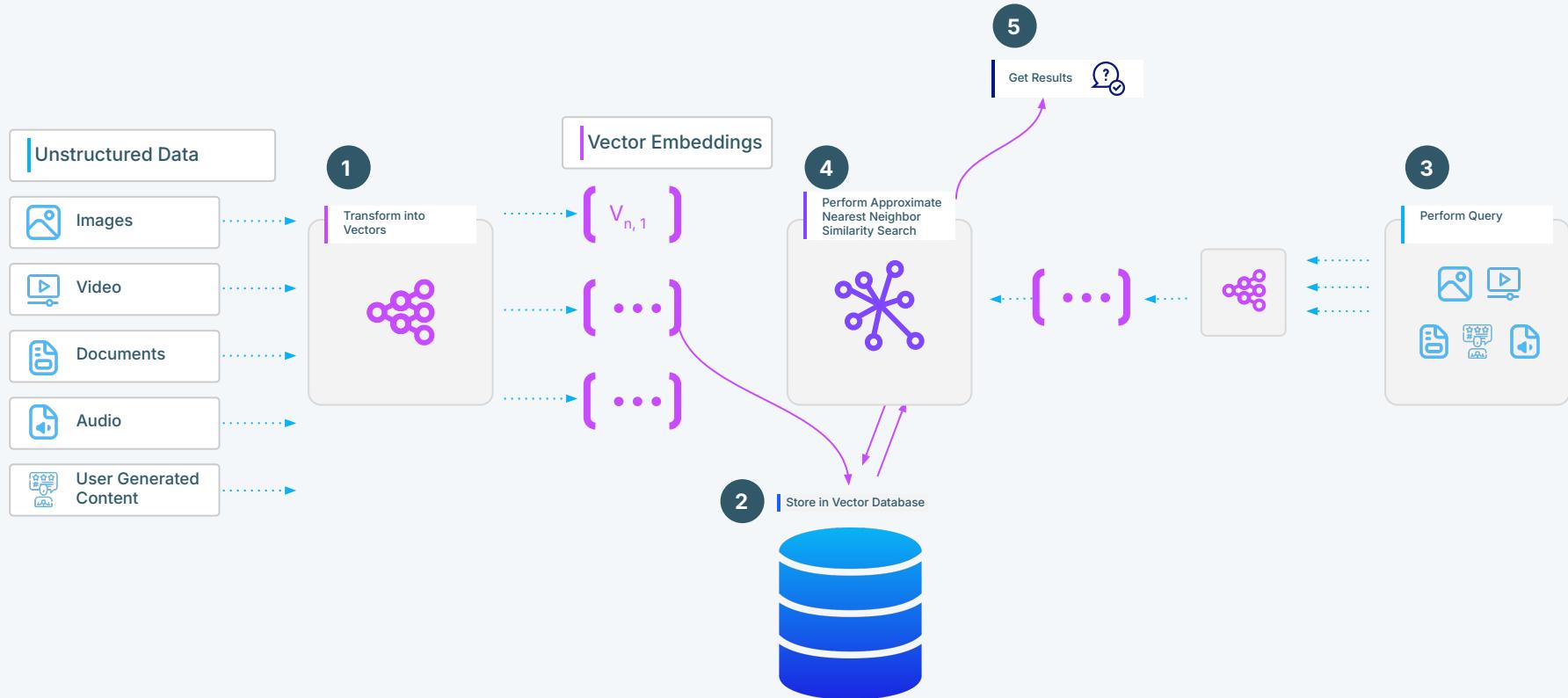


Image from [Nvidia](#)

# How Similarity Search Works





# Real-Time Pipelines Can Help

## External Context Ingest

Ingesting, routing, clean, enrich, transforming, parsing, chunking and vectorizing structured, unstructured, semistructured, binary data and documents

## Prompt engineering

Crafting and structuring queries to optimize LLM responses

## Context Retrieval

Enhancing LLM with external context such as Retrieval Augmented Generation (RAG)

## Roundtrip Interface

Act as a Discord, REST, Kafka, SQL, Slack bot to roundtrip discussions

# NiFi 2.0.0 Features



- Python Integration
- Parameters
- JDK 21+
- JSON Flow Serialization
- Rules Engine for Development Assistance
- Run Process Group as Stateless
- flow.json.gz

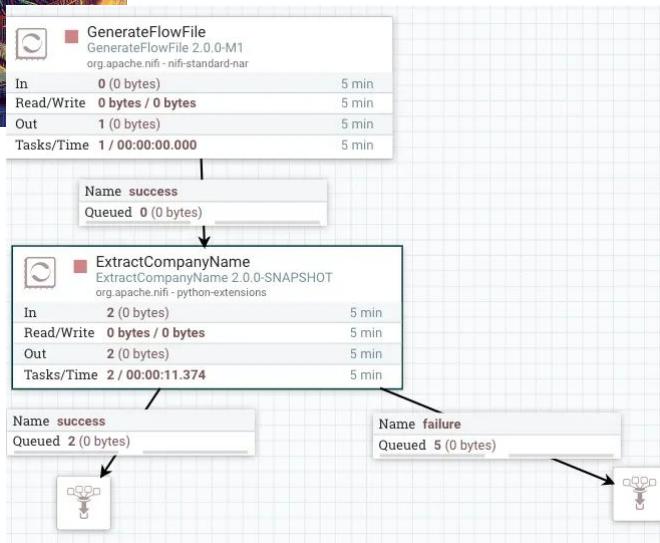
<https://cwiki.apache.org/confluence/display/NIFI/NiFi+2.0+Release+Goals>

<https://medium.com/cloudera-inc/getting-ready-for-apache-nifi-2-0-5a5e6a67f450>



# Extract Company Names

- Python 3.10+
- Hugging Face, NLP, SpaCY, PyTorch



## Attribute Values

companylist

["Amazon", "Microsoft", "Cloudera", "DataSQLR", "Google", "IBM"]

filename

36fb4ae6-701a-4e1d-b890-c93b44f2200b

parsedcompany

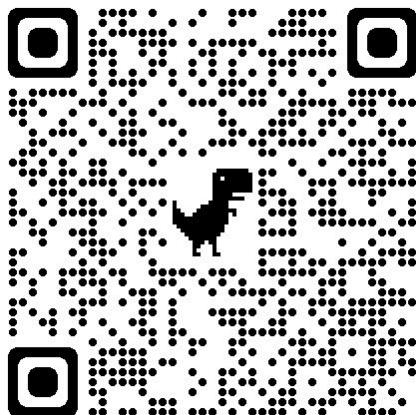
Amazon

path

./

uuid

6366a2c9-3dd4-4e8f-8825-83189d403b92



# CaptionImage

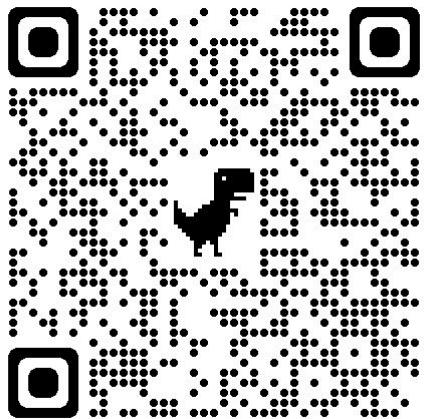
- Python 3.10+
- Hugging Face
- Salesforce/blip-image-captioning-large
- Generate Captions for Images
- Adds captions to FlowFile Attributes
- Does not require download or copies of your images

<https://github.com/tspannhw/FLaNK-python-processors>



# RESNetImageClassification

- Python 3.10+
- Hugging Face
- Transformers
- Pytorch
- Datasets
- microsoft/resnet-50
- Adds classification label to FlowFile Attributes
- Does not require download or copies of your images

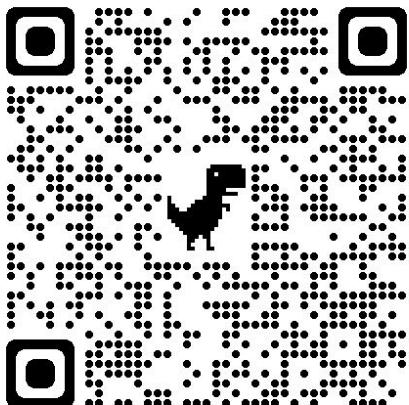


<https://github.com/tspannhw/FLaNK-python-processors>



# NSFW Image Detection

- Python 3.10+
- Hugging Face
- Transformers
- Falconsai/nsfw\_image\_detection
- Adds normal and nsfw to FlowFile Attributes
- Gives score on safety of image
- Does not require download or copies of your images

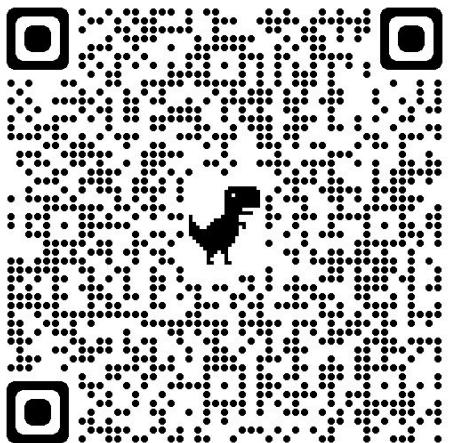


<https://github.com/tspannhw/FLaNK-python-processors>

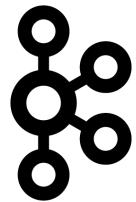


# FacialEmotionsImageDetection

- Python 3.10+
- Hugging Face
- Transformers
- facial\_emotions\_image\_detection
- Image Classification
- Adds labels/scores to FlowFile Attributes
- Does not require download or copies of your images



<https://github.com/tspannhw/FLaNK-python-processors>



# YES, FRANZ, IT'S KAFKA

Let's do a metamorphosis on your data. Don't fear changing data.

## You don't need to be a brilliant writer to stream data.



Franz Kafka was a German-speaking Bohemian novelist and short-story writer, widely regarded as one of the major figures of 20th-century literature. His work fuses elements of realism and the **fantastic**.

[Wikipedia](#)



0 53,639 / 153.08 MB 0 0 230 831 546 0 0 0 0 0 0 0 22:26:28 EDT



# Open Source Edition



- Apache NiFi in Docker

- Try new features

- quickly

- Develop applications

- locally

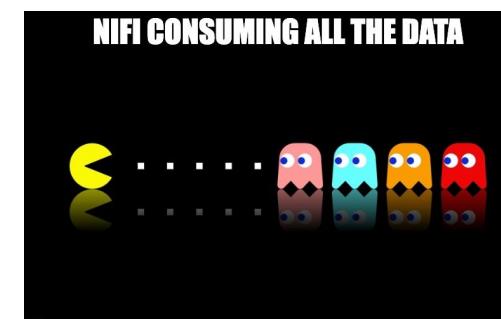
- Docker NiFi

- `docker run --name nifi -p 8443:8443 -d -e SINGLE_USER_CREDENTIALS_USERNAME=admin -e SINGLE_USER_CREDENTIALS_PASSWORD=ctsBtRBKHRAx69EqUghvvgEvjnaLjFEB apache/nifi:latest`

- Licensed under the ASF License
    - Unsupported
    - NiFi 1.25 and NiFi 2.0.0-M2



NIFI CONSUMING ALL THE DATA



Raspberry Pi AI Kit - Hailo  
Edge AI

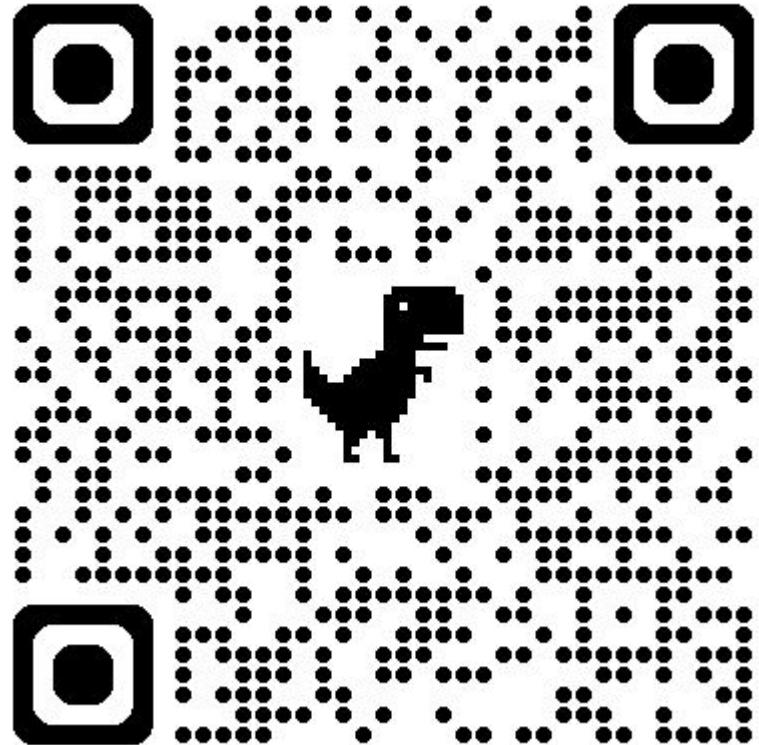
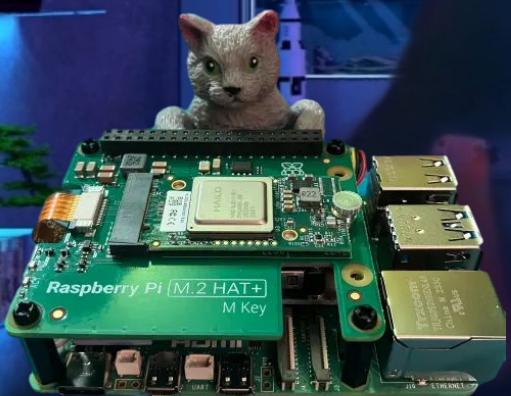


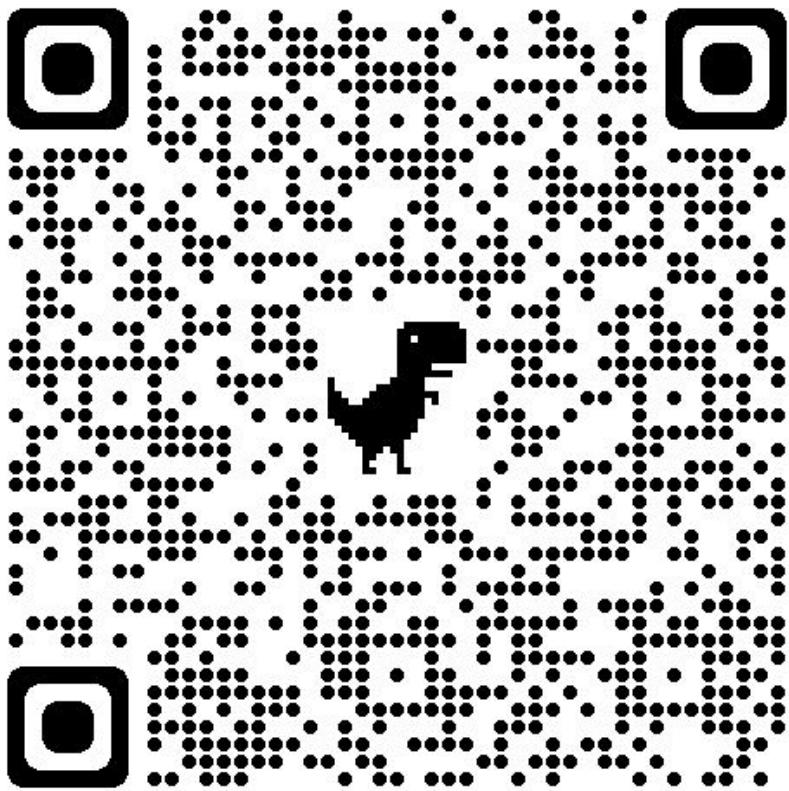
Milvus



<https://medium.com/@tspann/unstructured-data-processing-with-a-raspberry-pi-ai-kit-c959dd7fff47>

# Raspberry Pi AI Kit Hailo Edge AI Pose Estimation





Street Cameras



<https://medium.com/cloudera-inc/streaming-street-cams-to-yolo-v8-with-python-and-nifi-to-minio-s3-3277e73723ce>