



# **Smart Cities Unleashed**

Tim Spann, Senior Solutions Engineer

# Tim Spann

**paasdev.bsky.social**

@PaasDev // Blog: [datainmotion.dev](http://datainmotion.dev)

Senior Solutions Engineer, Snowflake

NY/NJ/Philly - Cloud Data + AI Meetups

ex-Zilliz, ex-Pivotal, ex-Cloudera, ex-HPE,  
ex-StreamNative, ex-EY, ex-Hortonworks.

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



# AI + Streaming Weekly by Tim Spann



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

This week in Apache NiFi, Apache Polaris, Apache Flink, Apache Kafka, ML, AI, Streamlit, Jupyter, Apache Iceberg, Python, Java, LLM, GenAI, Snowflake, Unstructured Data and Open Source friends.



Introduction and Overview

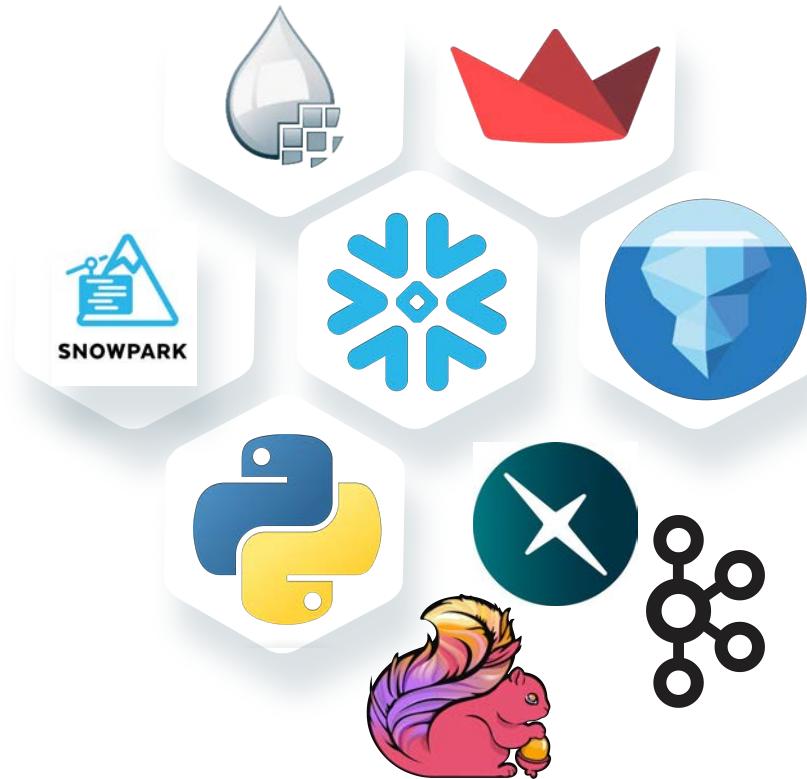
Data

Apache NiFi

Demo

Resources

# Unlocking Data Requires a Team





# Structured, Semistructured, Unstructured Data



# General (Google) Transit Feed Specification

GTFS - Protocol Buffers (binary format)

## Organizations:

- Open Transit Software Foundation
- Mobility Data
- Google
- GTFS.org

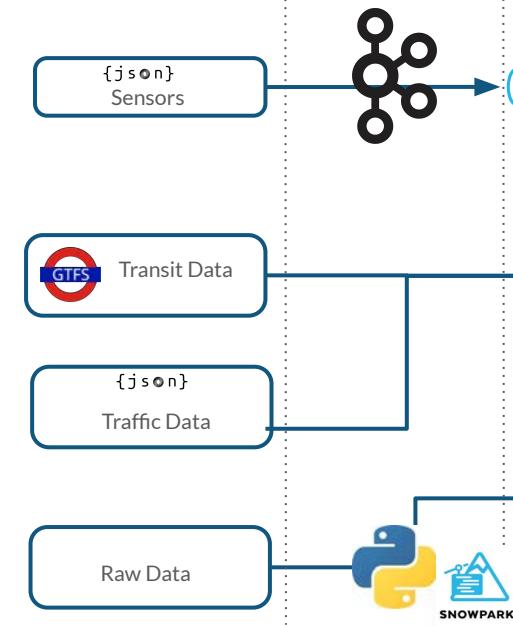


# Example Smart City Architecture

DATA  
FROM  
THE  
REAL



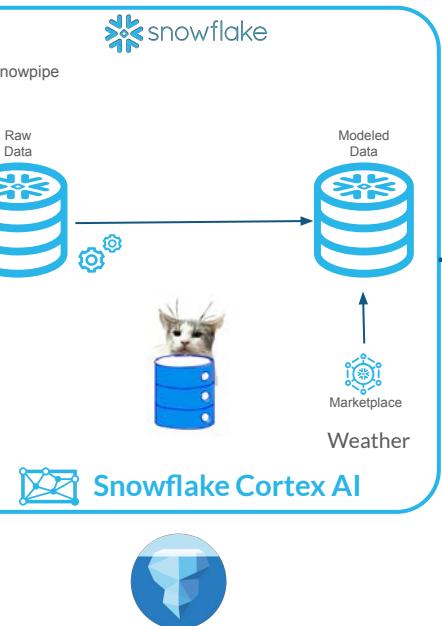
DATA  
SOURCES



DATA  
INTEGRATION



DATA  
PLATFORM



DATA  
CONSUMERS



# Semi-Structured Data



- Open Data like Open AQ - Air Quality Data
- Location, Time, Sensors
- Apache Avro, Parquet, Orc
- JSON and XML
- Hierarchical Data
- Logs
- Key-Value

<https://docs.snowflake.com/en/sql-reference/data-types-semistructured>



Unstructured

# Unstructured Data



- Lots of formats
- Text, Documents, PDF
- Images, Videos, Audio
- Email
- Variants

# City Cameras





# Structured Data



- Snowflake Tables
- Snowflake Hybrid Tables
- Apache Iceberg Tables
- Relational Tables
- Postgresql Tables
- CSV, TSV



# Apache Iceberg™ - Append



- NiFi - PutIcebergTable
- Snowpark -  
`df.write.mode("append").  
save_as_table("atable_iceberg")`

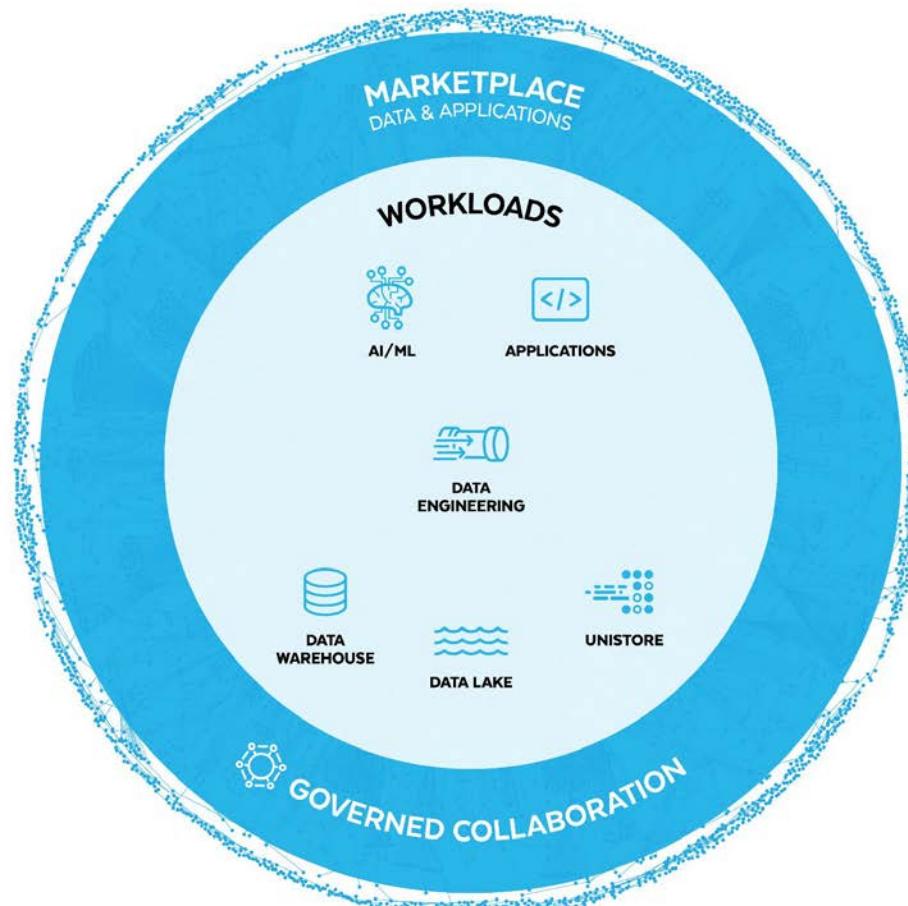
[https://quickstarts.snowflake.com/guide/getting\\_started\\_iceberg\\_tables/](https://quickstarts.snowflake.com/guide/getting_started_iceberg_tables/)



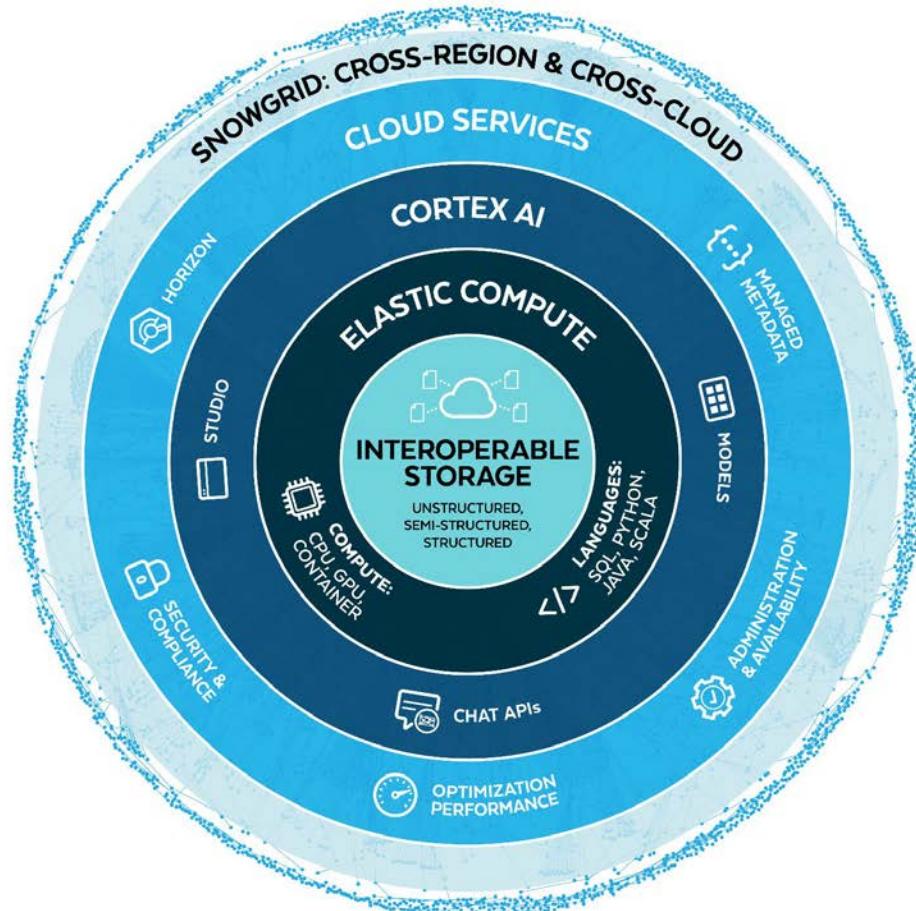


AI DATA CLOUD

# SNOWFLAKE AI DATA CLOUD



## SNOWFLAKE PLATFORM ARCHITECTURE





0 53,639 / 153.08 MB

0

0

230

831

546

160

0

0

0

0

0

0

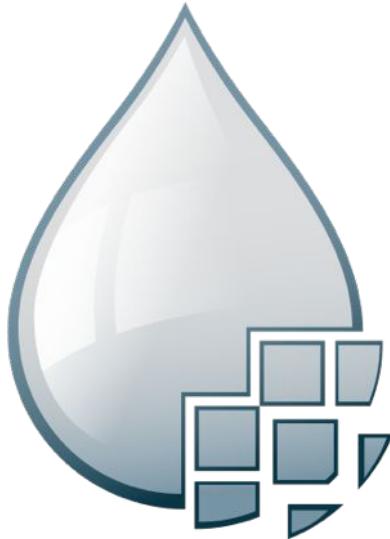
22:26:28 EDT



# Apache NiFi

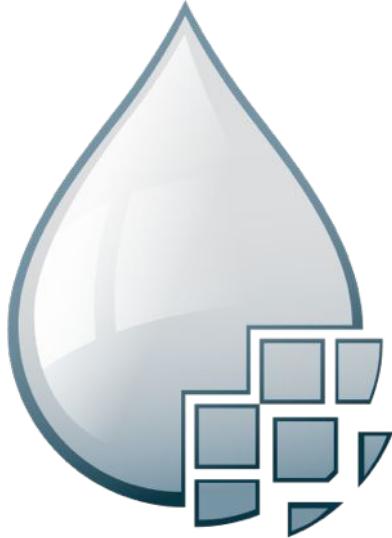


# Apache NiFi for Data Ingest, Movement and Routing



- Guaranteed delivery
- Data buffering
  - Backpressure
  - Pressure release
- Prioritized queuing
- Flow specific QoS
  - Latency vs. throughput
  - Loss tolerance
- Data provenance
- Supports push and pull models
- Hundreds of processors
- Visual command and control
- Hundreds of sources
- Flow templates
- Pluggable/multi-role security
- Designed for extension
- Clustering
- Version Control

# The Power of Apache NiFi



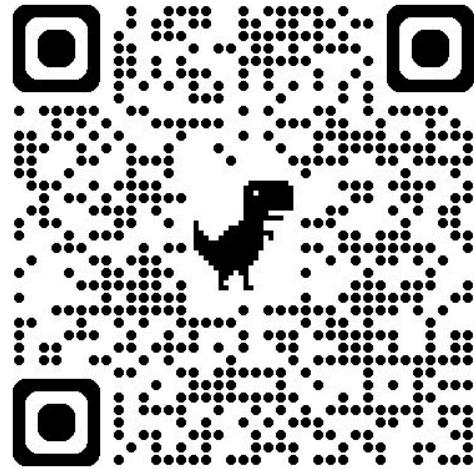
- Moving Binary, Unstructured, Image and Tabular Data
- Enrichment
- Universal Visual Processor
- Simple Event Processor
- Routing
- Feeding data to Central Messaging
- Support for modern protocols
- Kafka Protocol Source/Sink
- Pulsar Protocol Source/Sink

# APACHE NIFI 2.0 FEATURES

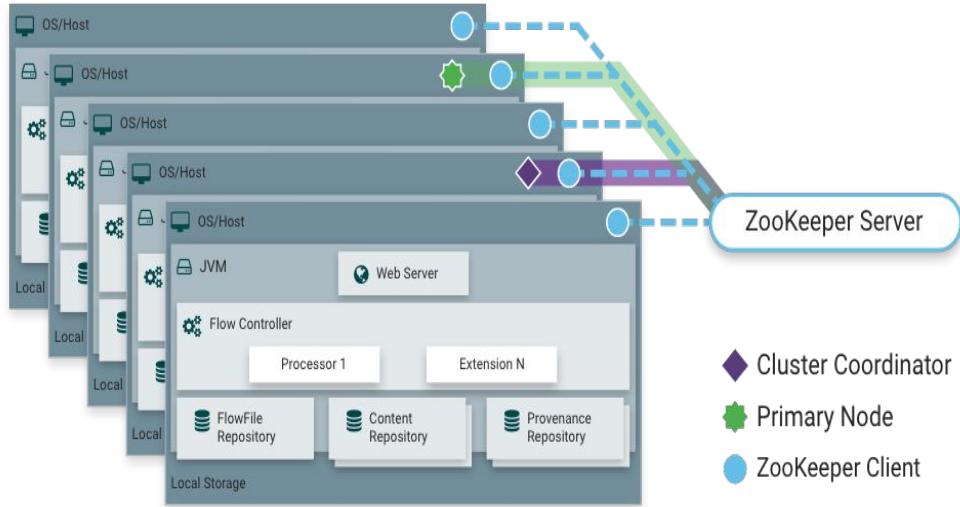
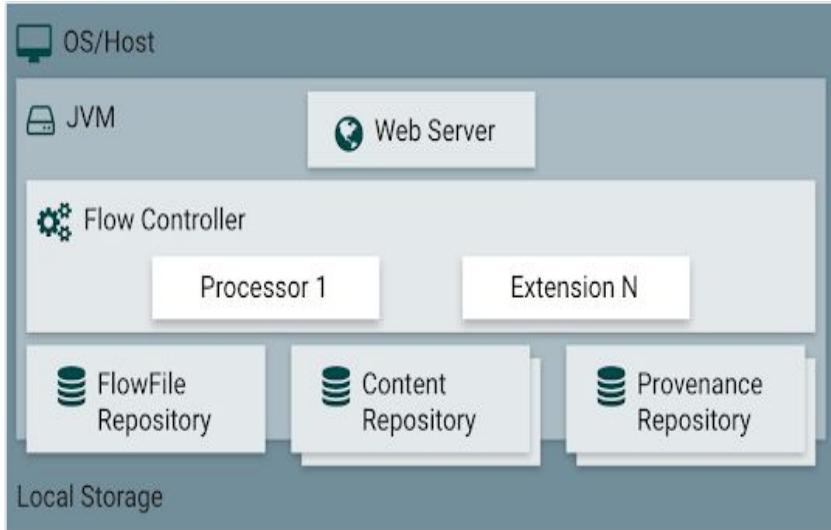
DataFlow is built for Real-Time Integration and AI

## Major Updates:

- Python Integration
- ParameterIZATION
- JDK 21+
- Provenance / Data Lineage
- Rules Engine for Development Assistance
- Additional Azure Processors
- Integration with Zendesk, Slack,
- Database Tables as Schemas
- Amazon Glue Schema Registry
- OpenTelemetry Support



# Architecture



- Cluster Coordinator
- Primary Node
- ZooKeeper Client

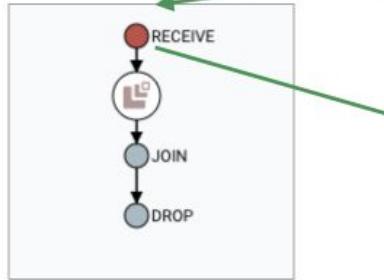
<https://nifi.apache.org/docs/nifi-docs/html/overview.html>

# PROVENANCE

Displaying 13 of 104  
Oldest event available: 11/15/2016 13:34:50 EST  
Showing the most recent events.

Date/Time	Type	FlowFile Uuid	Size	Component Name	Component Type
11/15/2016 13:35:03.8...	RECEIVE	379fc4f6-60e0-4151-9743-28...	44 bytes	ConsumeKafka	ConsumeKafka
11/15/2016 13:35:02.7...	RECEIVE	78f8c38b-89fc-4d00-a8d8-51...	44 bytes	ConsumeKafka	ConsumeKafka
11/15/2016 13:35:01.6...	RECEIVE	2bcd5124-bb78-489f-ad8a-7...	44 bytes	ConsumeKafka	ConsumeKafka

- Tracks data at each point as it flows through the system
- Records, indexes, and makes events available for display
- Handles fan-in/fan-out, i.e. merging and splitting data
- View attributes and content at given points in time



**Provenance Event**

DETAILS	ATTRIBUTES	CONTENT
Attribute Values		
filename	328717796819631	No value previously set
kafka.offset	44815	No value previously set
kafka.partition	6	No value previously set
kafka.topic	nifi-testing	No value previously set
path	/	No value previously set
uuid	328716238521-11e5-9005-105126272-95	

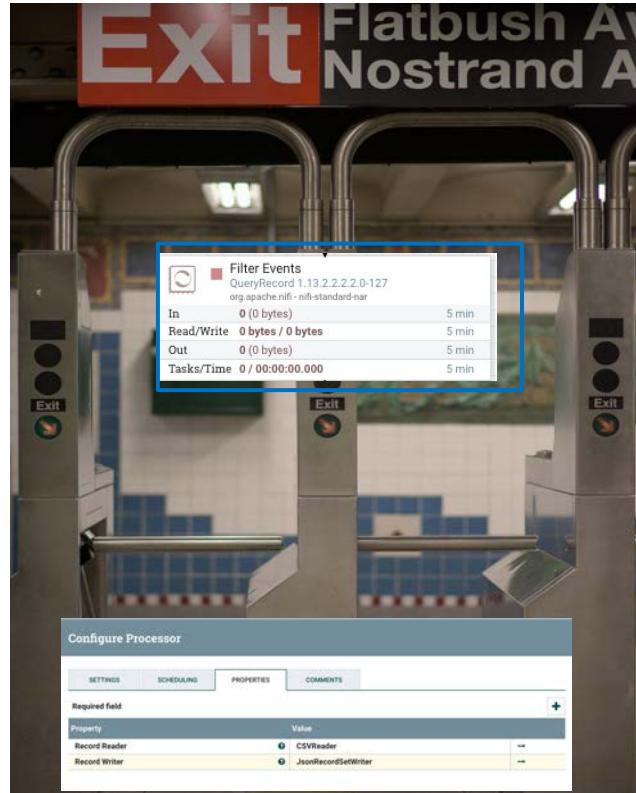
# UNSTRUCTURED DATA WITH NIFI

- **Archives** - tar, gzipped, zipped, ...
- **Images** - PNG, JPG, GIF, BMP, ...
- **Documents** - HTML, Markdown, RSS, PDF, Doc, RTF, Plain Text, ...
- **Videos** - MP4, Clips, Mov, Youtube URL...
- **Sound** - MP3, ...
- **Social / Chat** - Slack, Discord, Twitter, REST, Email, ...
- **Identify Mime Types, Chunk Documents, Store to Vector Database**
- **Parse Documents** - HTML, Markdown, PDF, Word, Excel, Powerpoint



# RECORD-ORIENTED DATA WITH NIFI

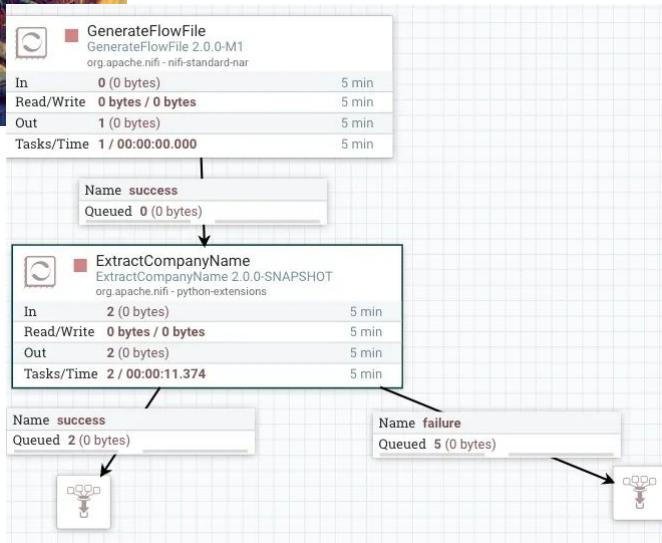
- **Record Readers** - Avro, CSV, Grok, IPFIX, JSAN1, JSON, Parquet, Scripted, Syslog5424, Syslog, WindowsEvent, XML
- **Record Writers** - Avro, CSV, FreeFromText, Json, Parquet, Scripted, XML
- Record Reader and Writer support referencing a schema registry for retrieving schemas when necessary.
- Enable processors that accept any data format without having to worry about the parsing and serialization logic.
- Allows us to keep FlowFiles larger, each consisting of multiple records, which results in far better performance.





# 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

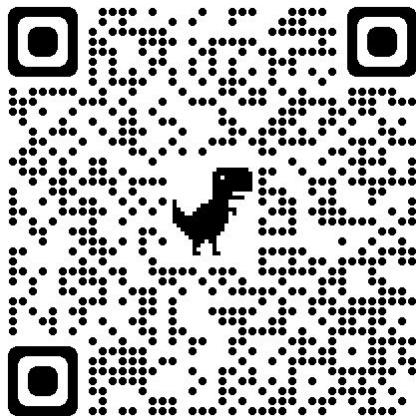


<https://github.com/tspannhw/FLaNK-python-ExtractCompanyName-processor>



# 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



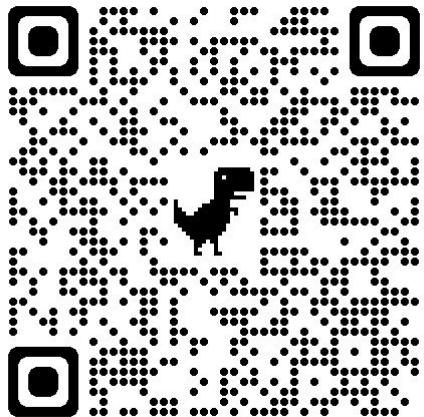
<https://github.com/tspannhw/FLaNK-python-processors>





# 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>



# Address To Lat/Long



- Python 3.10+
- geopy Library
- Nominatim
- OpenStreetMaps (OSM)
- [openstreetmap.org/copyright](https://openstreetmap.org/copyright)
- Returns as attributes and JSON file
- Works with partial addresses
- Categorizes location
- Bounding Box

# DEMO





TIME TO REBOOT THE CAT

imgflip.com

## RESOURCES AND WRAP-UP

<https://www.linkedin.com/in/timothyspann/>



# Getting Started

[https://quickstarts.snowflake.com/guide/analyze\\_pdf\\_invoices\\_snowpark\\_python\\_java/index.html#0](https://quickstarts.snowflake.com/guide/analyze_pdf_invoices_snowpark_python_java/index.html#0)

<https://medium.com/@tspann/utilizing-multiple-vectors-and-advanced-search-data-model-design-for-city-data-705d68d8daf2>

<https://medium.com/cloudera-inc/real-time-in-boston-part-1-0f92d7da3496>

<https://medium.com/cloudera-inc/boston-wheres-my-bus-llm-streaming-to-the-rescue-586dfd019237>

<https://medium.com/@tspann/real-time-irish-transit-analytics-ea76164c9595>

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

<https://medium.com/cloudera-inc/nyc-traffic-are-you-kidding-me-6d3fa853903b>

<https://medium.com/cloudera-inc/subways-and-transit-updates-in-real-time-30c104c359ef>

<https://medium.com/cloudera-inc/transit-in-sao-paulo-brasil-flank-style-eaec6753cc63>

# Open Source Edition



- Apache NiFi in Docker
- Runs 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



