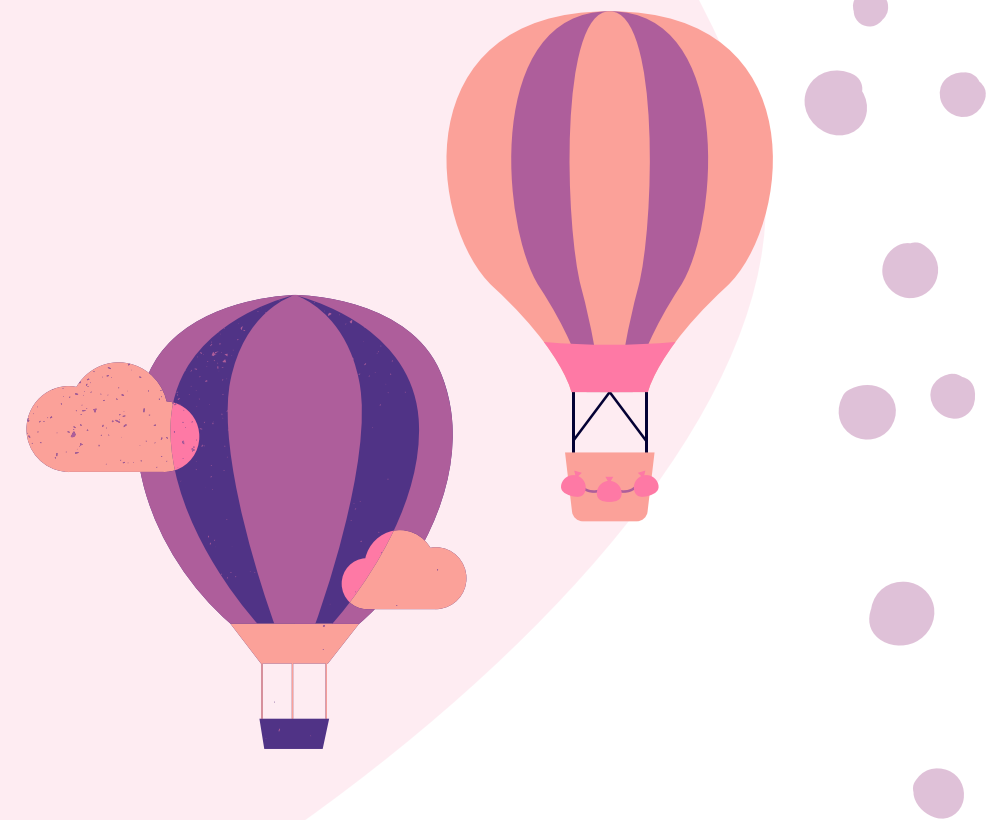# STRUCTURED DATA STREAMING Y NLP CON PYSPARK
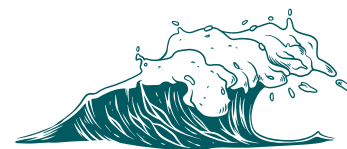
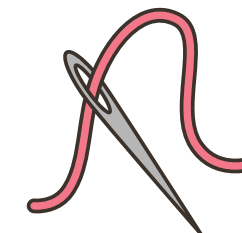BY: VALEN ARIZA Y LAURA LÓPEZ

# SOBRE NOSOTRAS
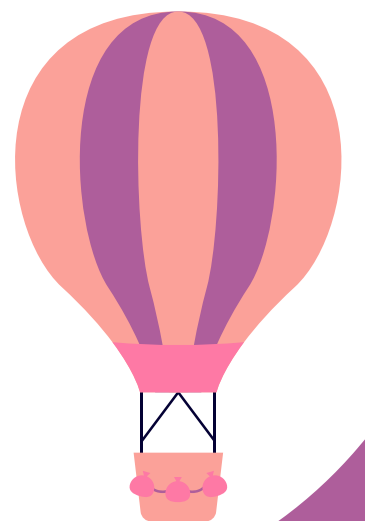
**VALENTINA ARIZA**

Data Engineer, DataKnow

**LAURA LÓPEZ**

Data Science Analyst, Accenture

# AGENDA

- Funcionamiento básico de Spark Structured Data Streaming

- Preprocesamiento de texto para NLP

- Proceso en un ambiente productivo

- Práctica

# PREREQUISITOS

- Crea una cuenta de Google Colab
- Crea una cuenta de Reddit https://www.reddit.com  (guarda tu usuario y contraseña)
- Prepárate...

# ¿QUÉ ES SPARK STRUCTURED DATA STREAMING?

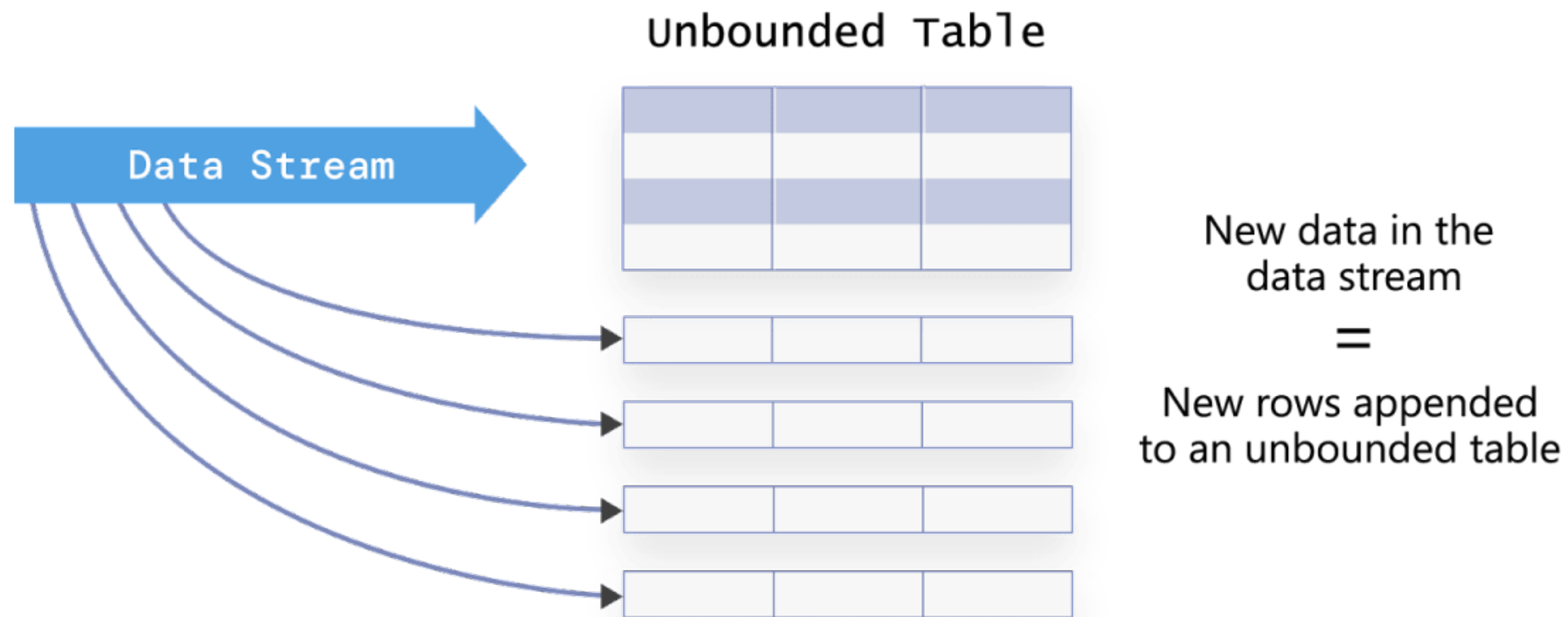# ¿QUÉ ES SPARK STRUCTURED DATA STREAMING?

# ¿QUÉ ES SPARK STRUCTURED DATA STREAMING?



Unbounded Table

Data Stream

New data in the data stream

=

New rows appended to an unbounded table

Data stream as an unbounded table

# FUNCIONAMIENTO BÁSICO SPARK DATA STREAMING

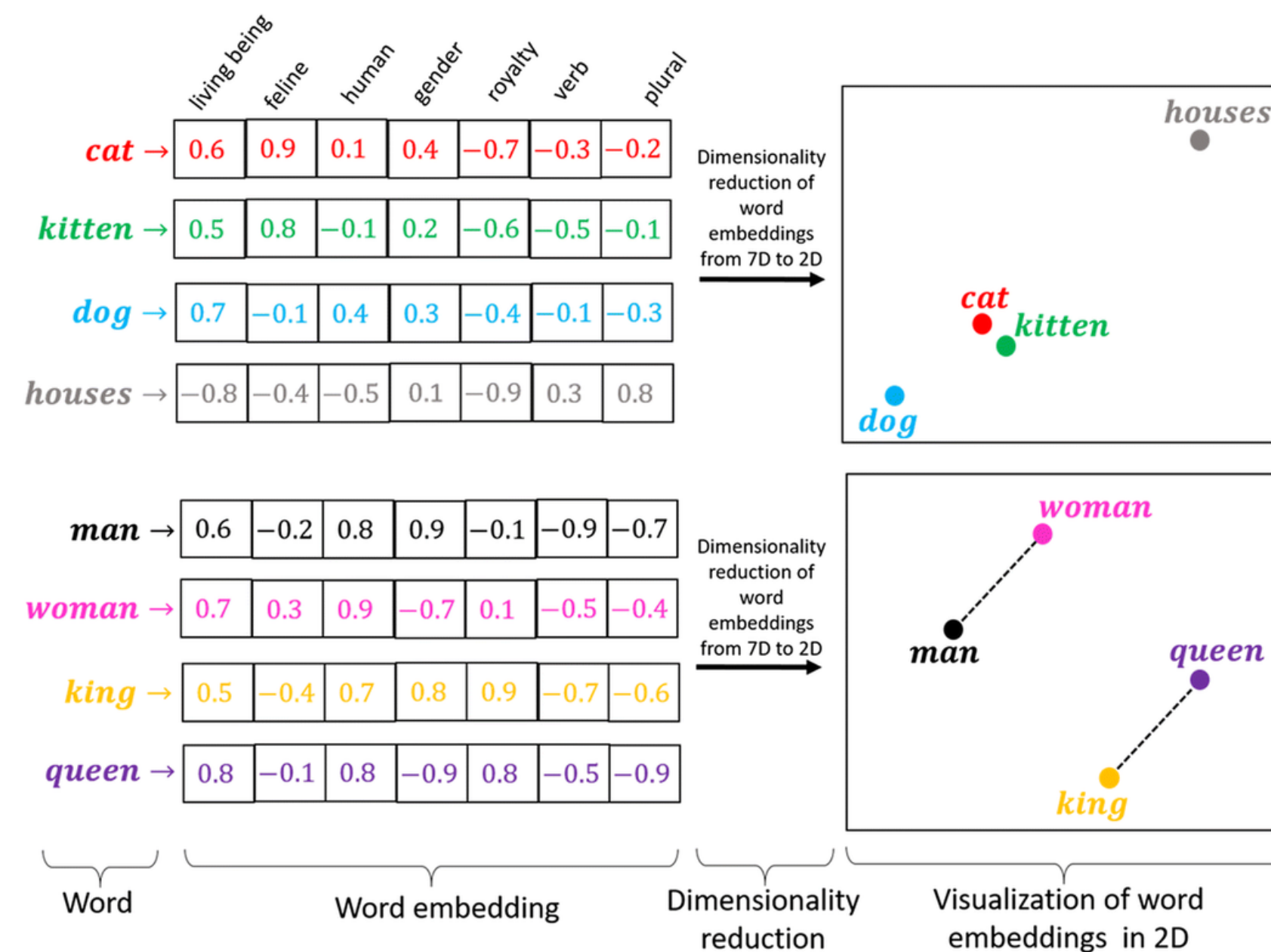# ORÍGENES Y DESTINOS

Apache Kafka

Archivos

Conexiones
Stream Socket

# ¿CÓMO FUNCIONA NATURAL LANGUAGE PROCESSING (NLP)?

# ¿CÓMO FUNCIONA NATURAL LANGUAGE PROCESSING (NLP)?

# PREPROCESAMIENTO

*"Hey Amazon - my package never arrived
https://www.amazon.com/gp/css/order-history?ref_=nav_orders_first
PLEASE FIX ASAP! @amazonhelp"*

# PREPROCESAMIENTO

- Normalizar

*"hey amazon - my package never arrived https://www.amazon.com/gp/css/order-history?ref_=nav_orders_first please fix asap! @amazonhelp"*

*"hey amazon my package never arrived please fix asap"*

# PREPROCESAMIENTO

- Eliminar stopwords

*"hey amazon my package never arrived please fix asap"*

*"amazon package never arrived fix asap"*

# PREPROCESAMIENTO

- Tokenizar

*"amazon package never arrived fix asap"*

⬇

*["amazon", "package", "never", "arrived", "fix", "asap"]*

# PREPROCESAMIENTO

- Normalizar
- Stopwords
- Símbolos y caracteres especiales
- Tokenizar
- Stemming
- Lemmatization

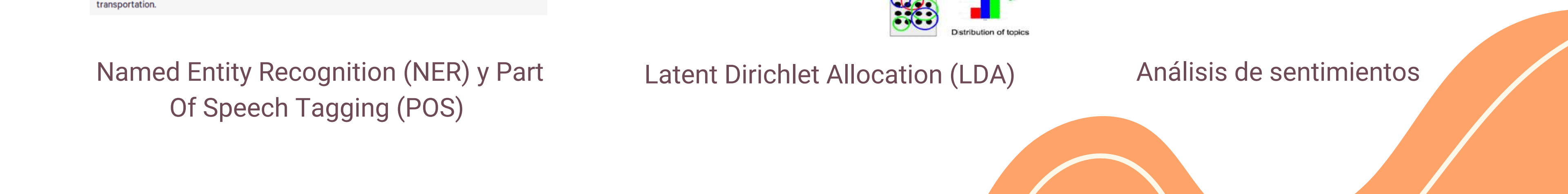# QUÉ PODEMOS HACER CON NLP



Wordcloud



N-Grams



Grafos de palabras



Named Entity Recognition (NER) y Part Of Speech Tagging (POS)



Latent Dirichlet Allocation (LDA)



Análisis de sentimientos

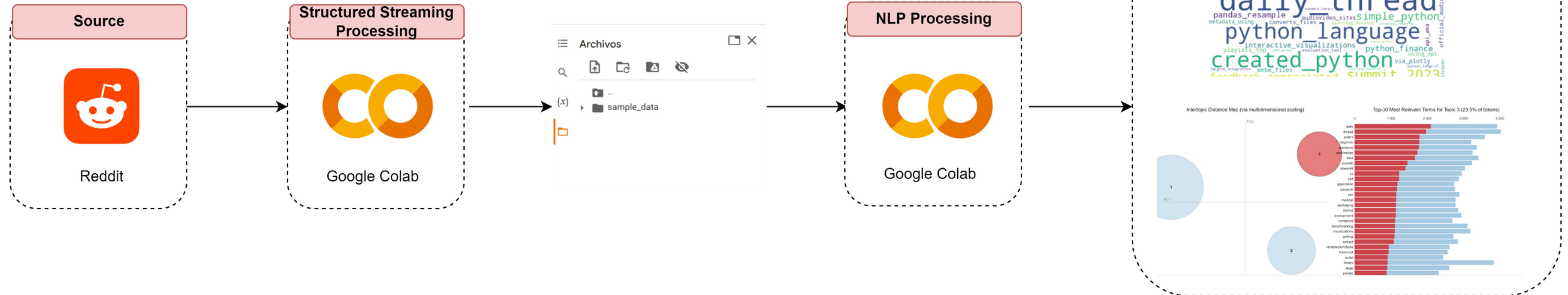# ¿QUÉ VAMOS A HACER?

# UNA ALTERNATIVA...

# MANOS A LA OBRA