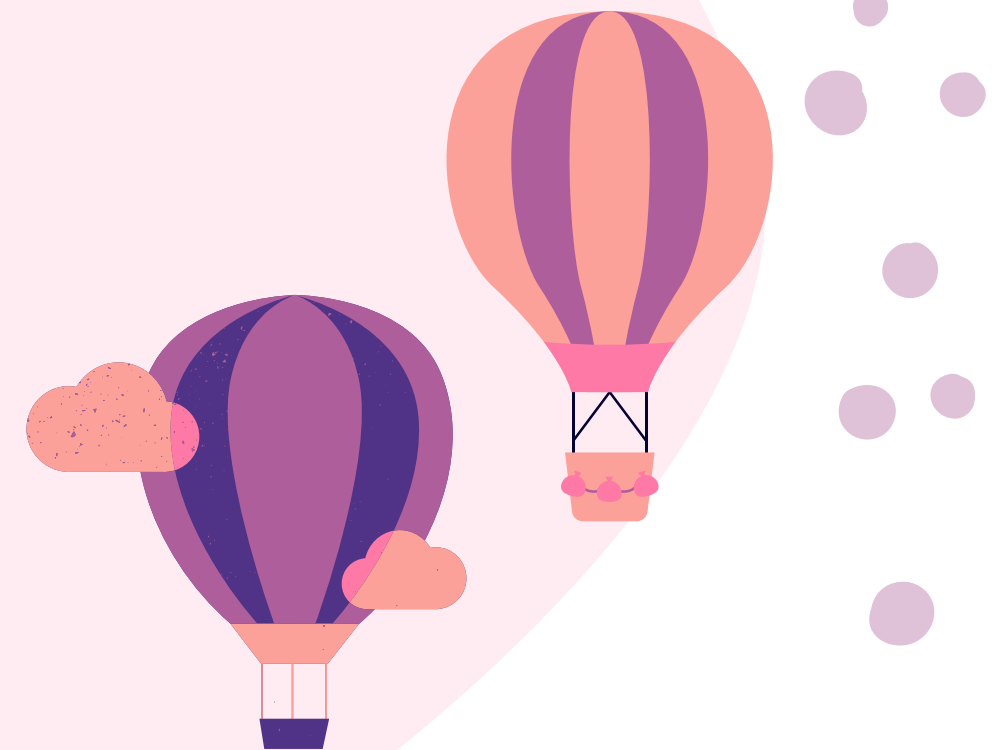


DATA STREAMING Y NLP CON PYSPARK

BY: VALEN ARIZA Y LAURA LÓPEZ



SOBRE NOSOTRAS



VALENTINA ARIZA

Data Engineer, DataKnow

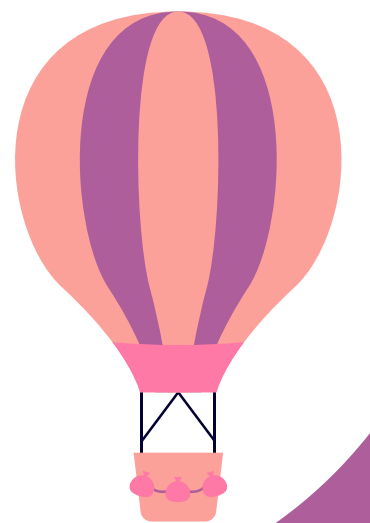


LAURA LÓPEZ

Data Science Analyst, Accenture

AGENDA

- Funcionamiento básico de Spark Structured Data Streaming
- Preprocesamiento de texto para NLP
- Proceso en un ambiente productivo
- Práctica



¿QUÉ ES SPARK STRUCTURED DATA STREAMING?

Apache Spark Core

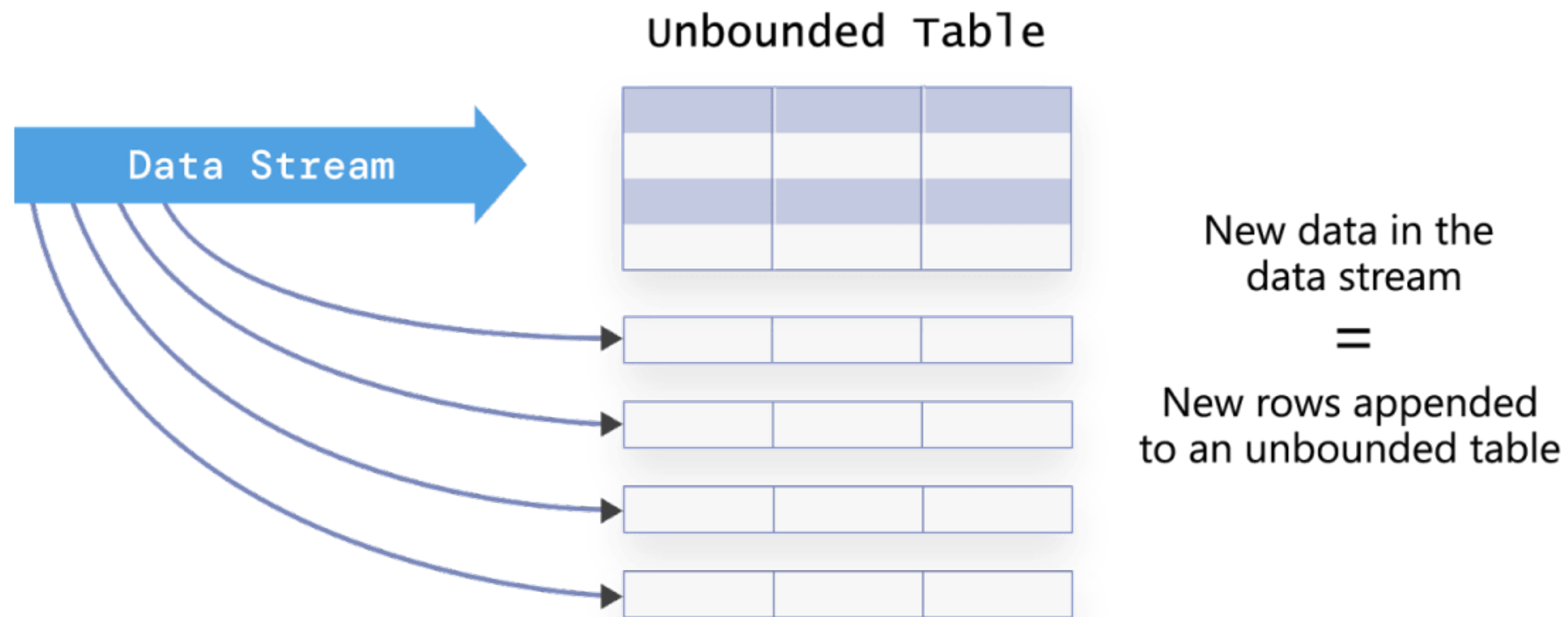
Spark SQL

Spark
Streaming

MLlib
(Machine
Learning)

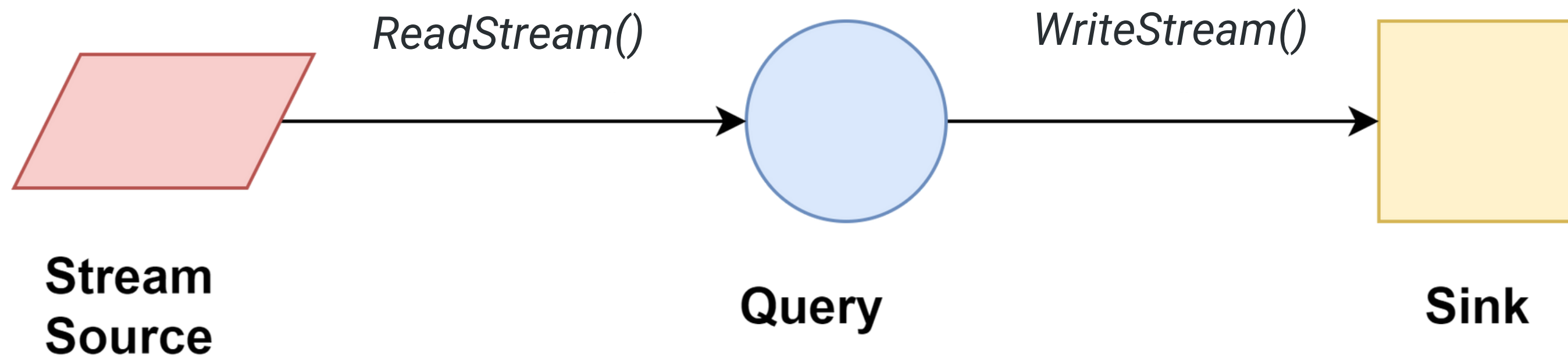
GraphX
(Graph)

¿QUÉ ES SPARK STRUCTURED DATA STREAMING?



Data stream as an unbounded table

FUNCIONAMIENTO BÁSICO SPARK DATA STREAMING



ORÍGENES Y DESTINOS

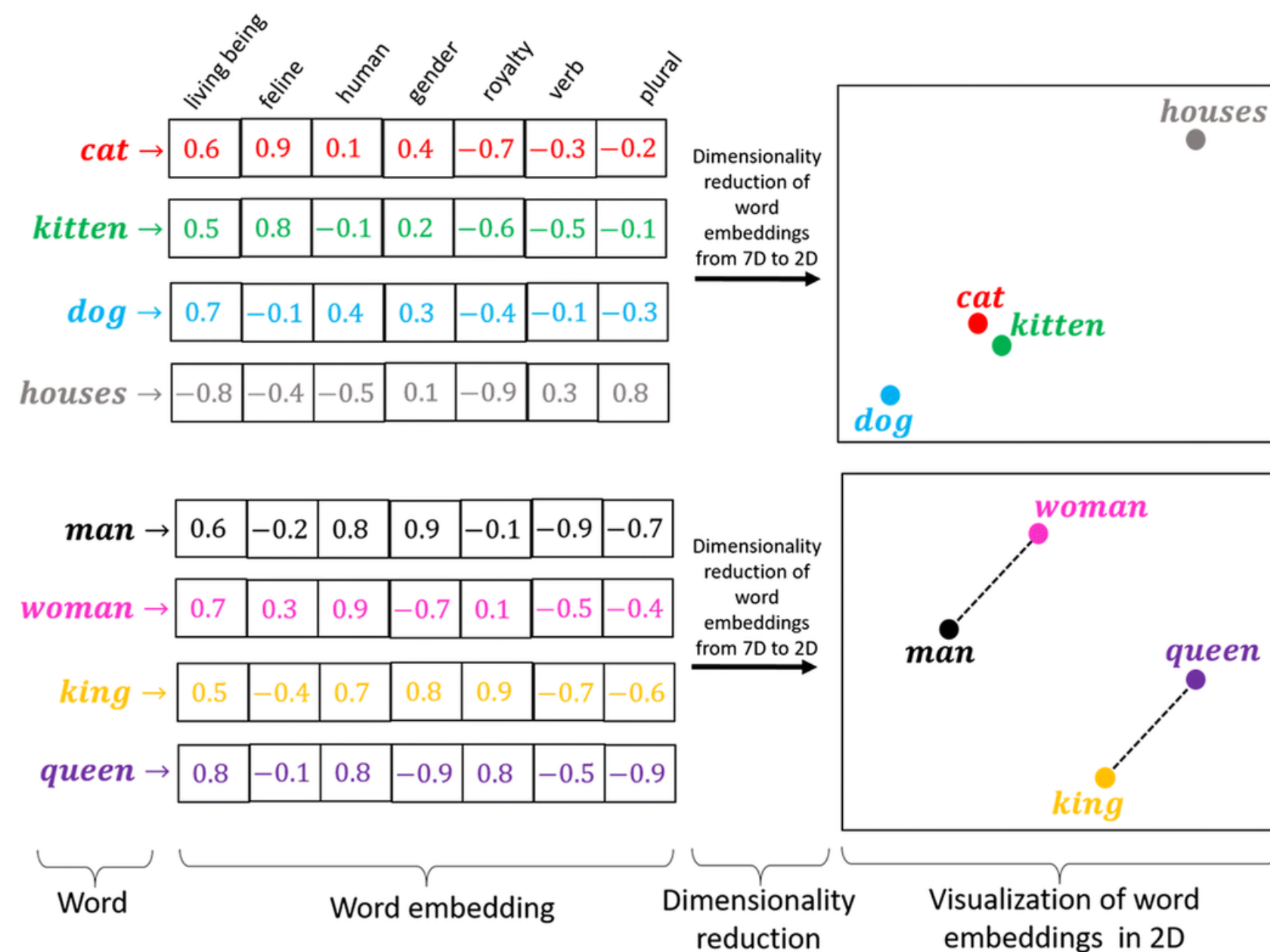


Conexiones
Stream Socket



Archivos

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”

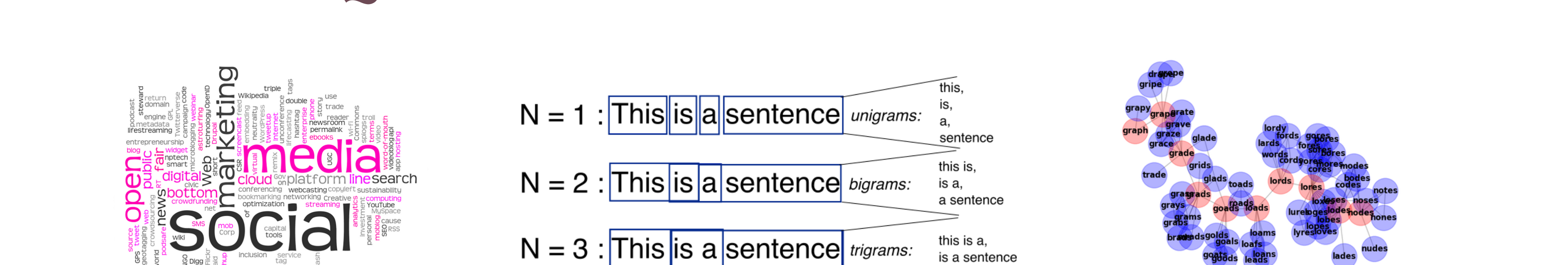


["amazon", "package", "never", "arrived", "fix", "asap"]

PREPROCESAMIENTO

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

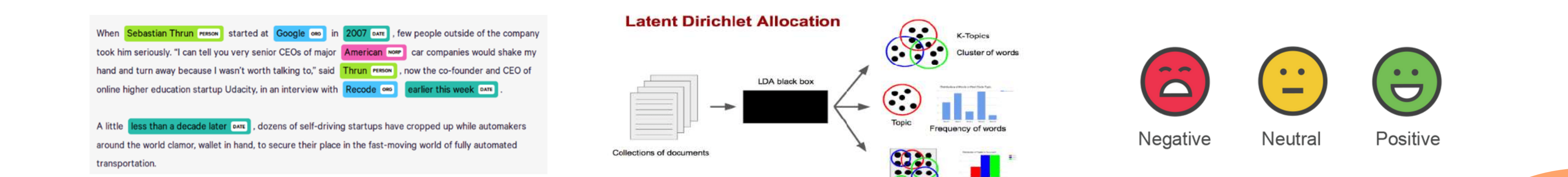
QUÉ PODEMOS HACER CON NLP



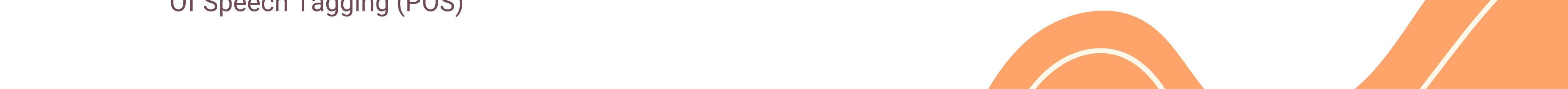
Wordcloud

N-Grams

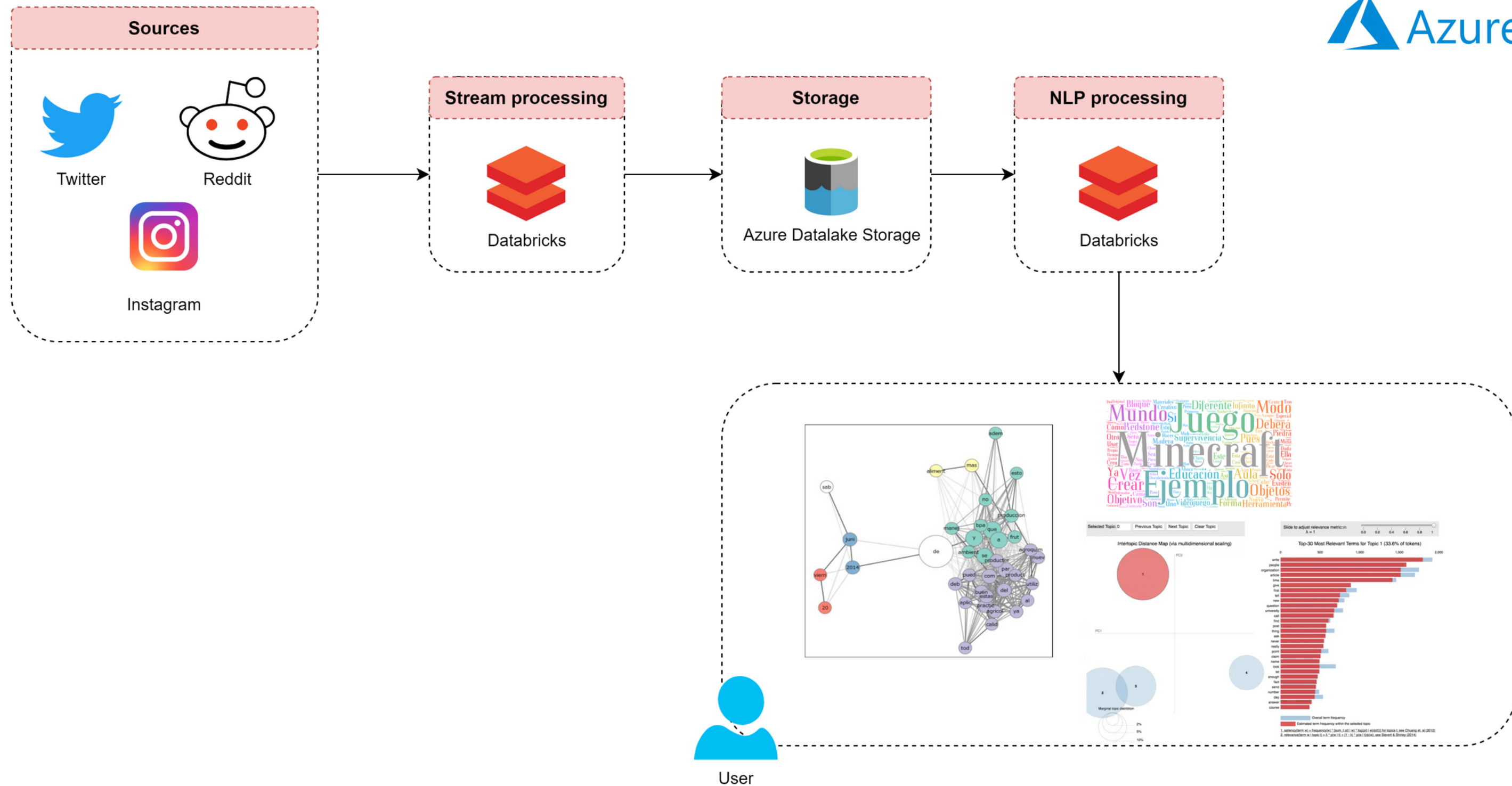
Grafos de palabras



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



PROCESO EN UN AMBIENTE PRODUCTIVO





MANOS A LA OBRA



RECURSOS NLP

- Documentación oficial de las librerías utilizadas
- Hugging Face: <https://huggingface.co/>
- Gensim: <https://pypi.org/project/gensim/>
- NLTK: <https://www.nltk.org/>

¿TIENES PREGUNTAS?

¡CONTÁCTANOS!



@valearizag



@lauralpezb

