



**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

**UNIVERSIDAD DE LAS FUERZAS ARMADAS – ESPE**

**DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN**

**Carrera de Software**

**Modalidad Presencial**

**Tema:**

**“Modelo basado en Minería de textos y Procesamiento del Lenguaje Natural (NLP) para la gestión de artículos científicos – caso de estudio Nanopartículas”**

**Autor:**

**CERÓN ÑAÑAY MARÍA BELÉN**

**Tutor: ING. GUALOTUÑA ÁLVAREZ, TATIANA MARISOL, PHD**

**SANGOLQUÍ, AGOSTO 2022**





# *Agenda*

- **INTRODUCCIÓN**
- **OBJETIVOS**
- **IMPLEMENTACIÓN**
- **RESULTADOS**
- **VALIDACIÓN DEL MODELO**
- **CONCLUSIONES Y RECOMENDACIONES**





**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

# INTRODUCCIÓN



# Antecedentes



Compartir información desde cualquier parte del mundo



Bases de datos académicas



Artículos Científicos



Investigaciones científicas

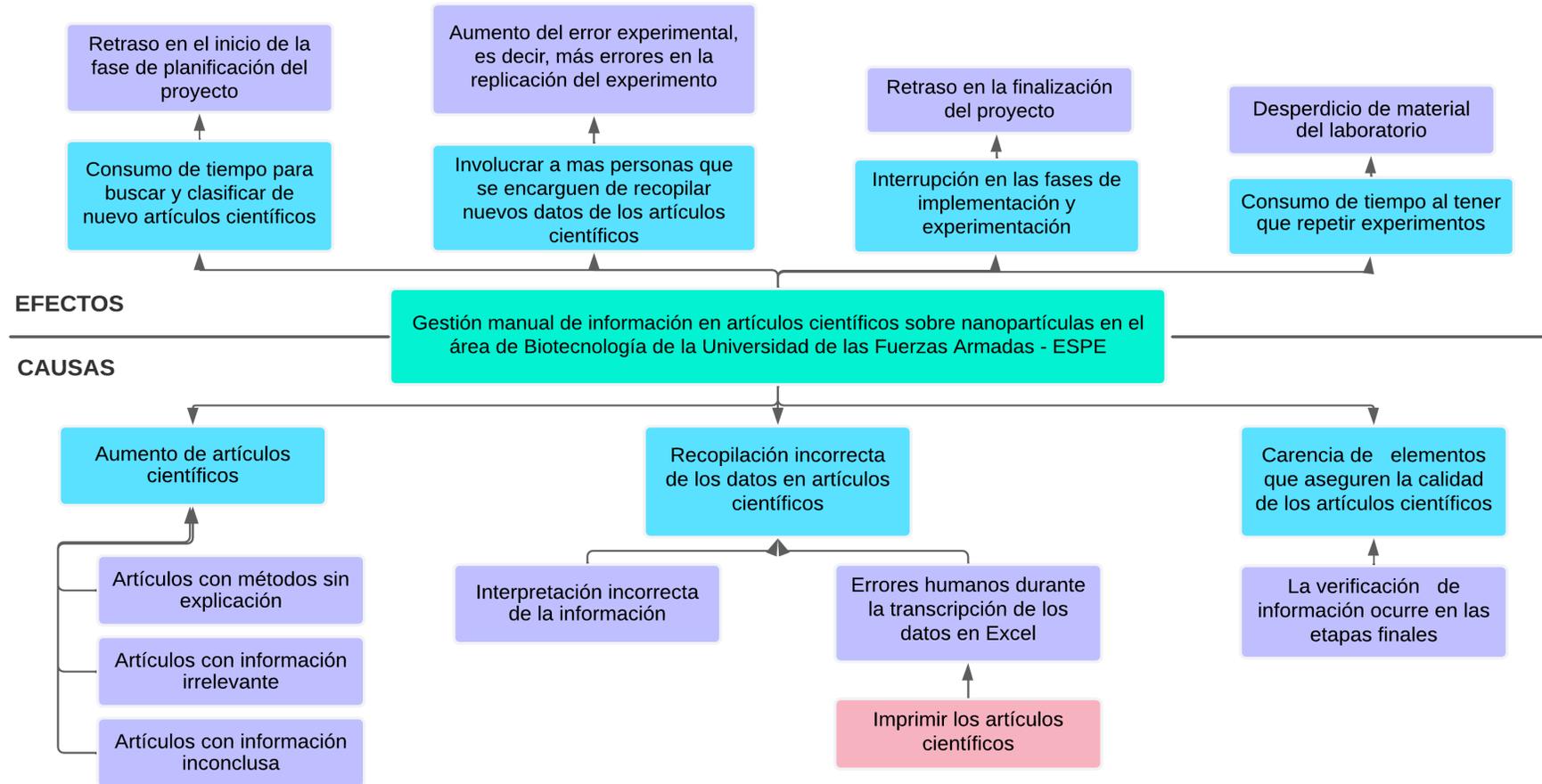


Sobreabundancia de información





# Problemática





**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

## OBJETIVOS





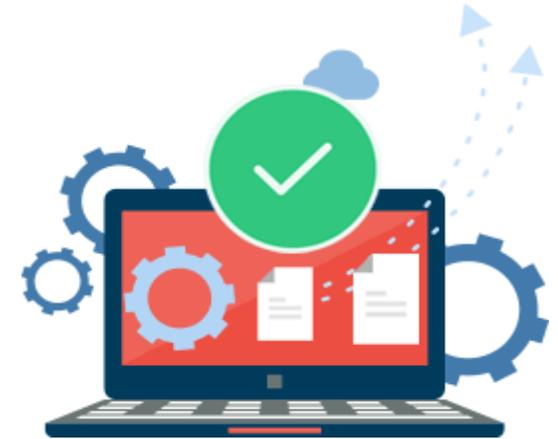
# *Objetivo general*

Desarrollar un modelo machine learning para la gestión de artículos científicos a través del cual se pueda establecer patrones de información que aporten al conocimiento científico



**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

# *Objetivos específicos*



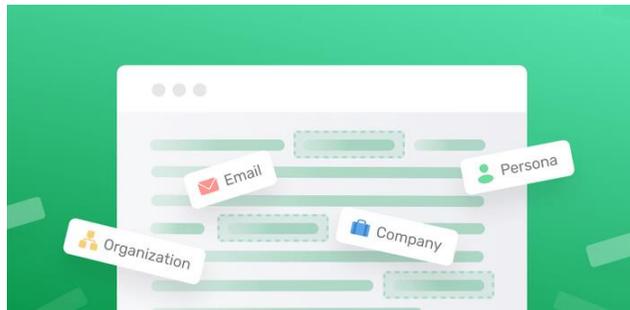
Realizar una revisión de literatura que permita determinar técnicas de minería de textos y Procesamiento del Lenguaje Natural (NLP), adecuadas para el procesamiento, etiquetado y extracción de datos.

Definir e implementar el modelo machine learning para la generación de patrones de comportamiento en las investigaciones relacionadas con nanopartículas.

Realizar la validación del modelo utilizando métricas de evaluación de rendimiento de clasificadores en tareas de extracción de datos.

# Estado del Arte

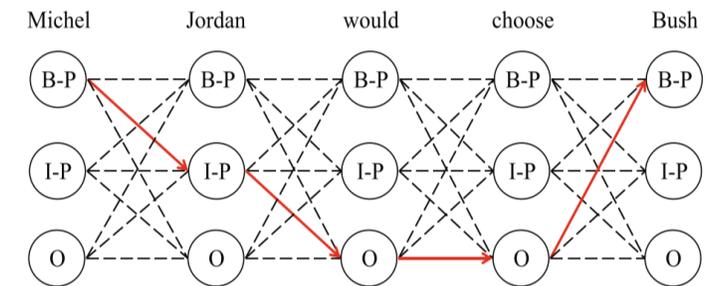
## NER



## Soluciones presentadas por los EP

Modelo	Estudios
<b>BiLSTM</b>	Tres investigaciones: EP1, EP4, EP5.
<b>CRF</b>	Dos investigaciones: EP4, EP7.
<b>BERT</b>	Dos investigaciones: EP4, EP6.
<b>RNN</b>	Una investigación: EP7.

## BIO encoding



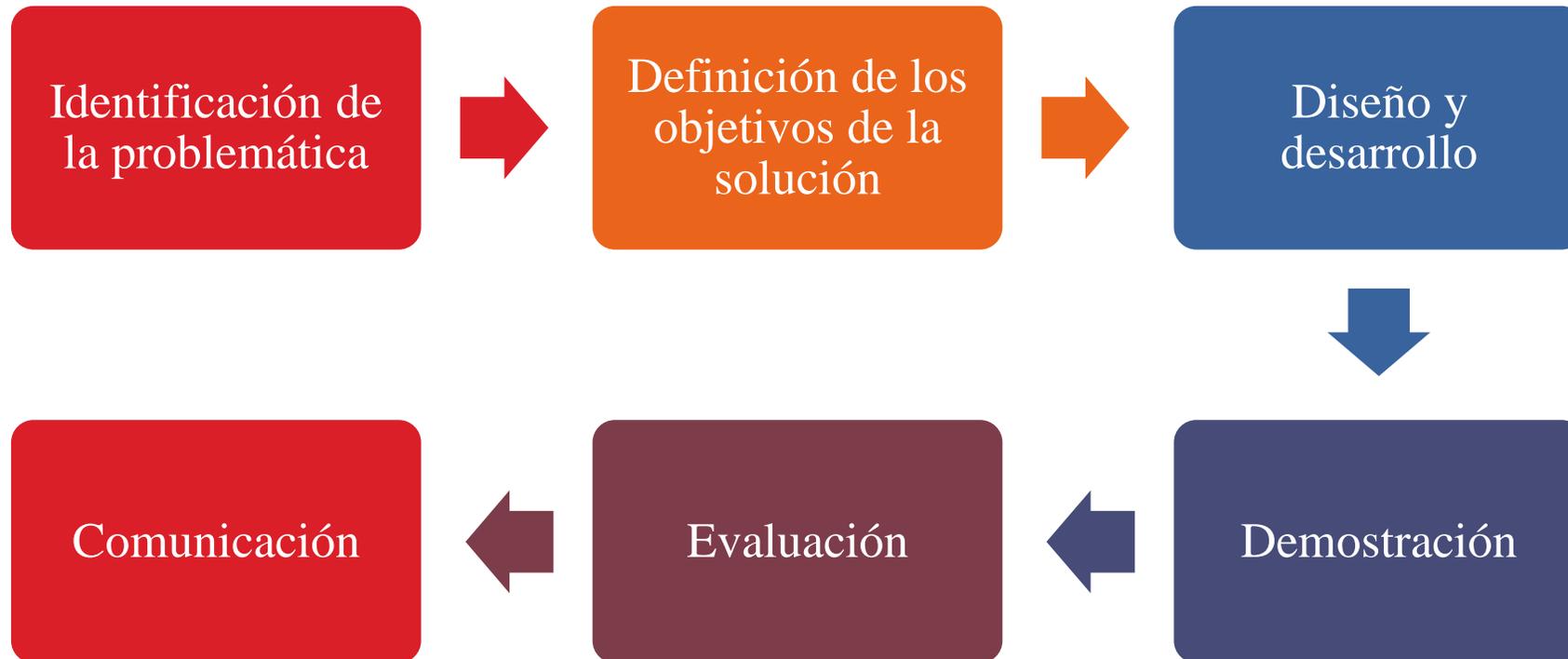


**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

## IMPLEMENTACIÓN

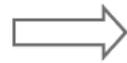


# *Design Science Research (DSR)*



# *Arquitectura de la solución*

Scopus<sup>®</sup>



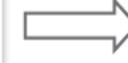
Origen de los datos



Proceso ETL



Entrenamiento  
del modelo



Resultados



# Origen de los datos

Comprensión de los datos



Recolección de datos



Análisis de los datos

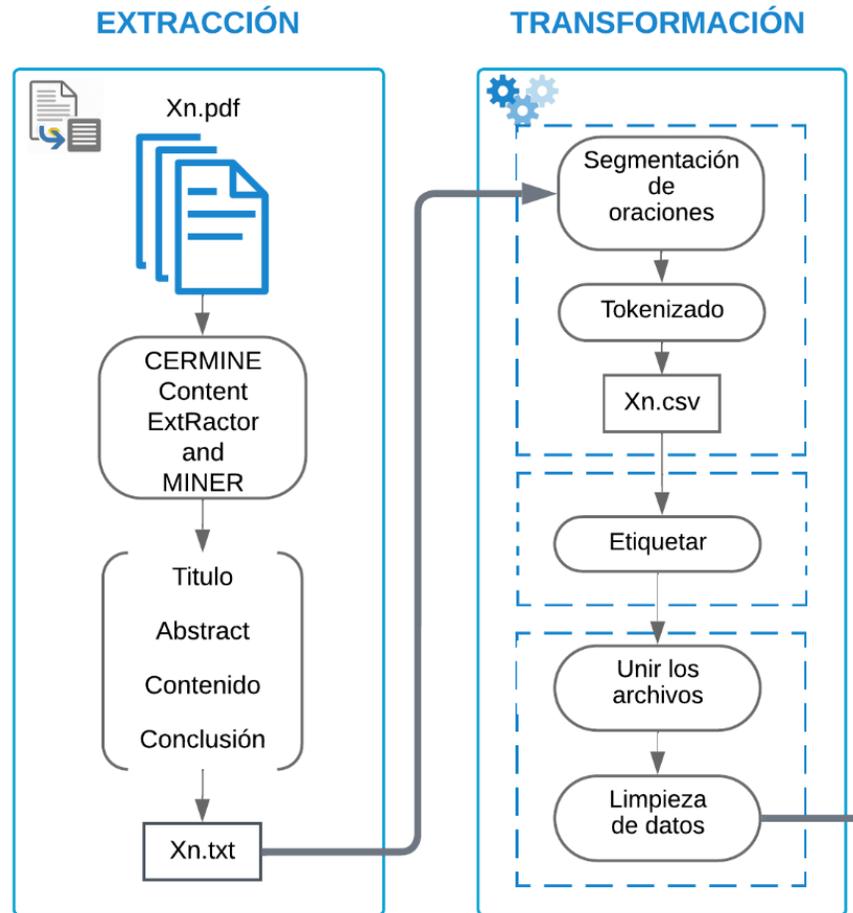
- 141 entidades

Material	Cantidad	Detalle
Artículos sin anotaciones	100	Artículos descargados directamente de la base de datos Scopus sin ninguna modificación.
Artículos con anotaciones	100	Artículos donde se encuentra subrayado la información necesaria a extraer.
Base de datos	1	Base de datos en Excel con la información extraída de los artículos, esta información se puede comparar con los artículos subrayados.

- 45 entidades
- 88 artículos científicos

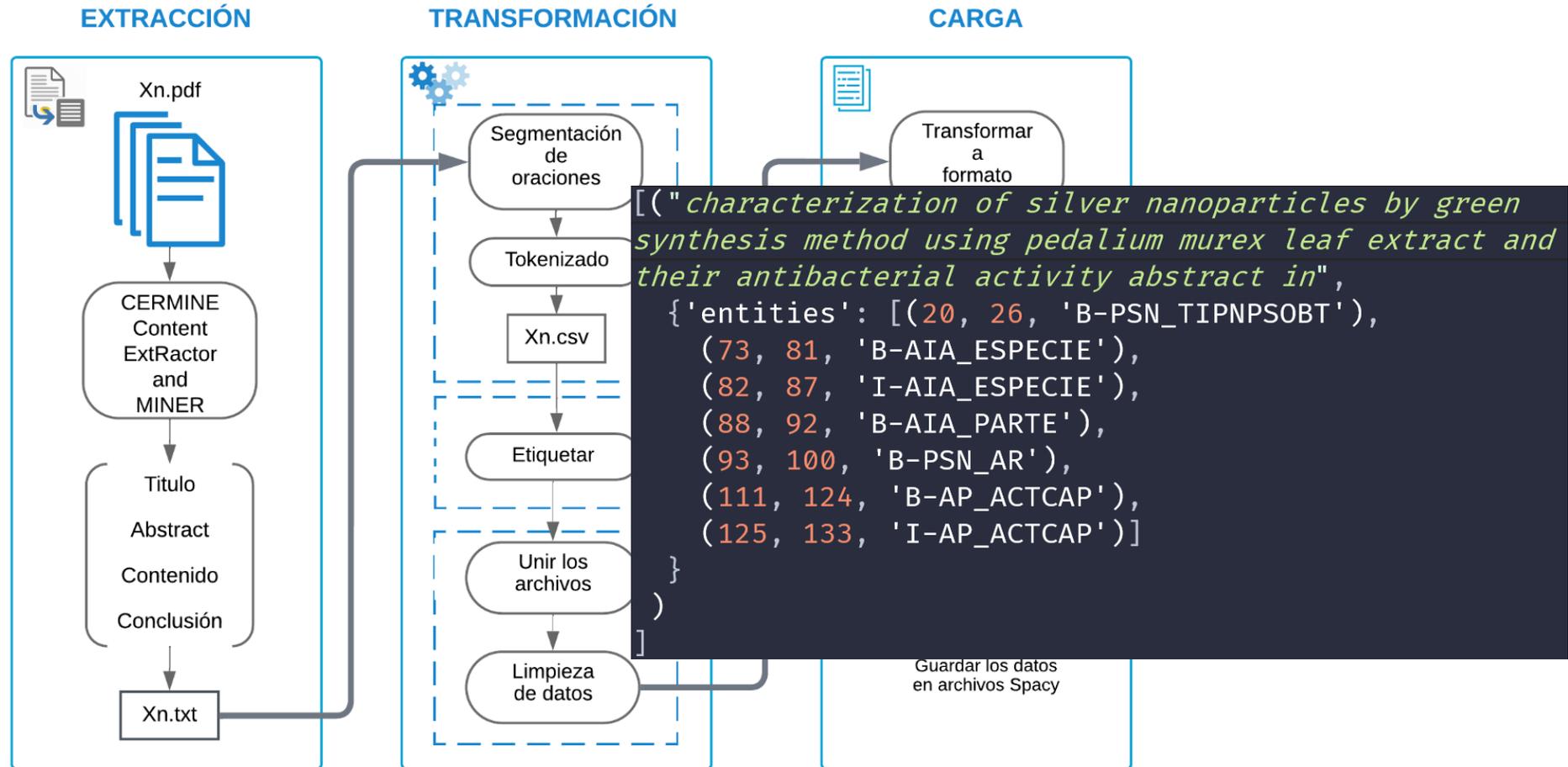


# Proceso ETL



id_paper	id_sentence	words	tag
X9.txt		0 caracterizat	O
X9.txt		0 of	O
X9.txt		0 silver	B-PSN_TIPNPSOBT
X9.txt		0 nanoparticle	O
X9.txt		0 by	O
X9.txt		0 green	O
X9.txt		0 synthesis	O
X9.txt		0 method	O
X9.txt		0 using	O
X9.txt		0 pedaliun	B-AIA_ESPECIE
X9.txt		0 murex	I-AIA_ESPECIE
X9.txt		0 leaf	B-AIA_PARTE
X9.txt		0 extract	B-PSN_AR
X9.txt		0 and	O
X9.txt		0 their	O
X9.txt		0 antibacterial	B-AP_ACTCAP
X9.txt		0 activity	I-AP_ACTCAP
X9.txt		1 abstract	O
X9.txt		2 in	O

# Proceso ETL





# *Entrenamiento del modelo*

Elección del modelo



Entrenamiento del  
módulo NER

SpaCy

Redes Neuronales  
Convolucionales

RoBERTa



**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA



**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

## RESULTADOS





# *Evaluación del Rendimiento del modelo*

## Métricas de evaluación

Precisión

$$P = \frac{TP}{TP + FP}$$

Recall

$$R = \frac{TP}{TP + FN}$$

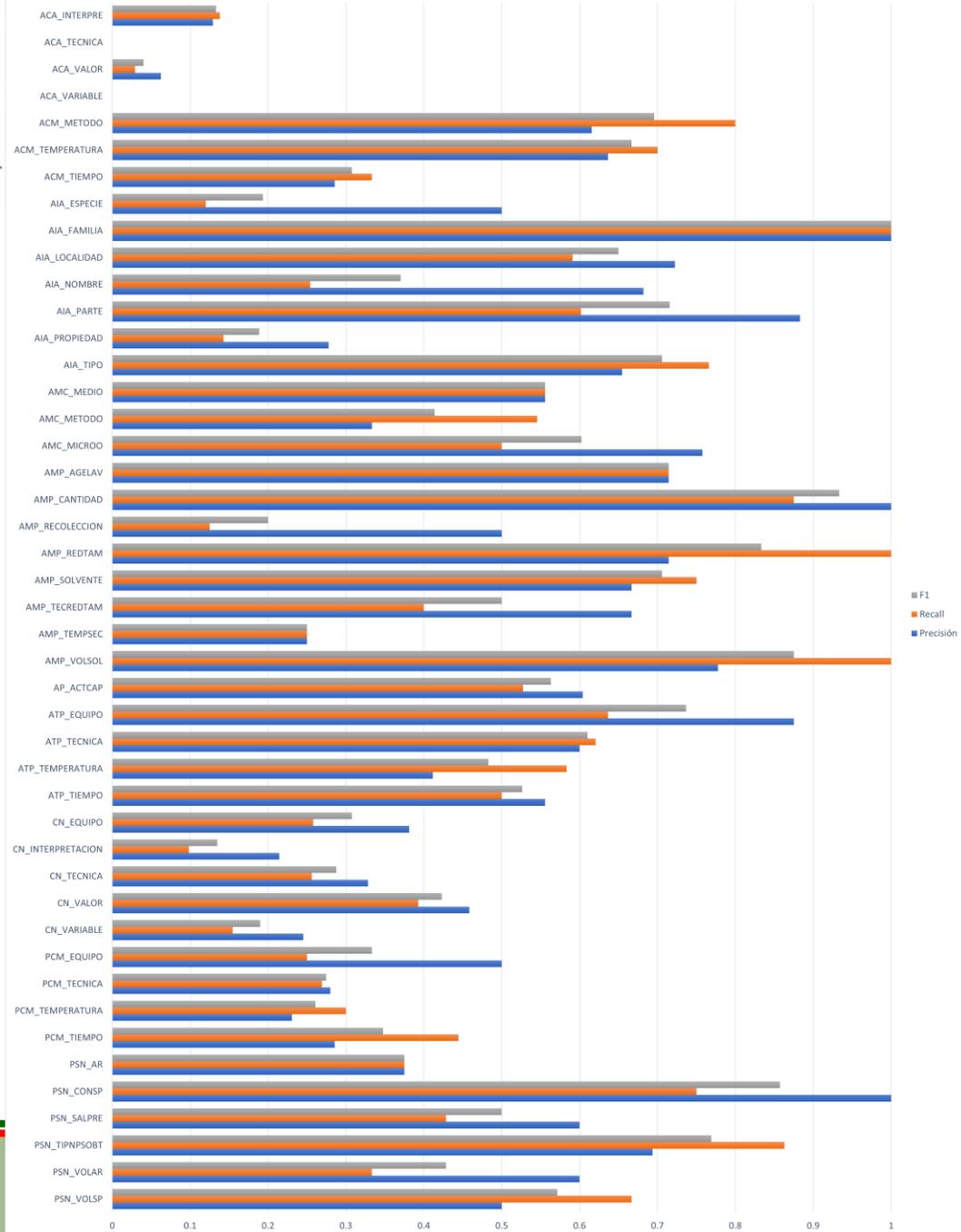
Puntuación F1

$$F1 = 2 \frac{P \times R}{P + R}$$

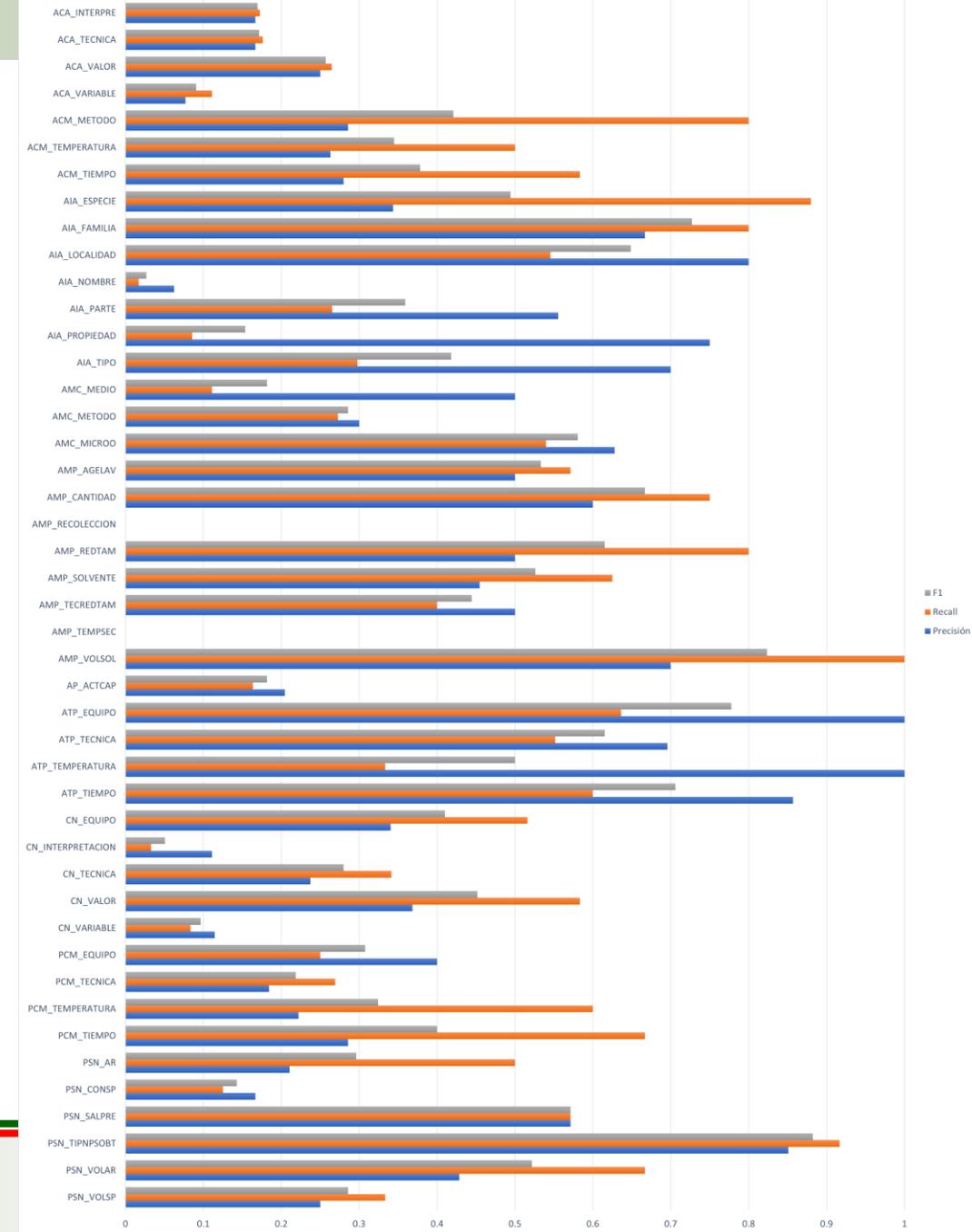




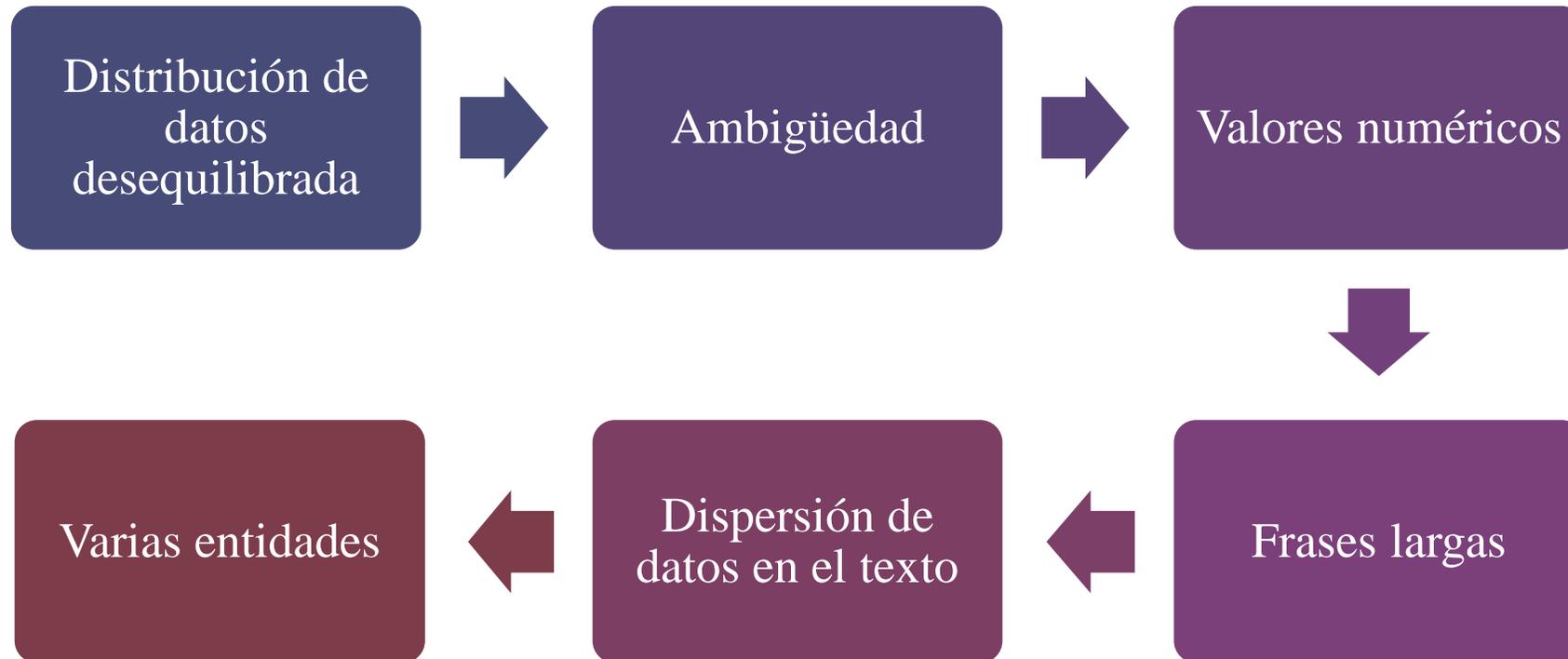
Resultados de la evaluación con CNN



Resultados de la evaluación con RoBERTa



# *Análisis de resultados*





**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

## VALIDACIÓN DEL MODELO





# Caso de Estudio

1. Reunión con el personal del área de Biotecnología de la Universidad de las Fuerzas Armadas – ESPE
2. Comparación entre los datos obtenidos con los modelos y la base de datos en Excel que contenía la extracción manual de los datos.

**Correcto:** Los modelos y la base de datos coinciden.

**Neutral:** Los modelos y la base de datos coinciden parcialmente, es decir, los modelos no lograron extraer por completo la información.

**Incorrecto:** Los modelos y la base de datos no coinciden.

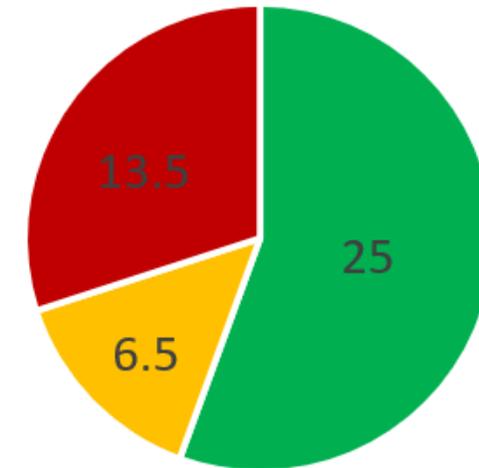
N	Etiquetas	X88	X167	X120	X9
1	AIA_TIPO	Correcto	Correcto	Correcto	Correcto
2	AIA_NOMBRE	Correcto	Correcto	Correcto	Correcto
3	AIA_ESPECIE	Incorrecto	Correcto	Incorrecto	Correcto
4	AIA_FAMILIA	Correcto	Incorrecto	Incorrecto	Correcto
5	AIA_PARTE	Incorrecto	Incorrecto	Correcto	Correcto
6	AIA_LOCALIDAD	Correcto	Correcto	Incorrecto	Correcto
7	AIA_PROPIEDAD	Incorrecto	Correcto	Incorrecto	Correcto
8	AMP_RECOLECCION	Correcto	Correcto	Incorrecto	Correcto
9	AMP_AGELAV	Incorrecto	Correcto	Incorrecto	Incorrecto





# *Caso de Estudio*

## Resultados



■ Correcto ■ Neutral ■ Incorrecto

Términos	X88	X167	X120	X9
Correcto	30	28	10	32
Neutral	3	8	8	7
Incorrecto	12	9	27	6





**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

## CONCLUSIONES Y RECOMENDACIONES





# *Conclusiones*



Utilización de NER



Necesidad de contar con mayor cantidad de textos etiquetados



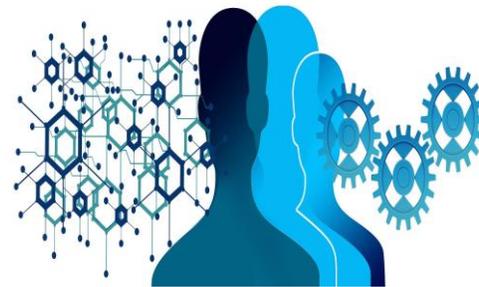
La solución se validó por expertos del área de Biotecnología



# *Recomendaciones*



Continuar con estudios sobre el tema de extracción de textos automatizada, específicamente para el caso de estudio de nanopartículas



Enfoque semi supervisado



Elaborar una interfaz gráfica que facilite su uso para los investigadores de Biotecnología



*Gracias por su  
atención*



**ESPE**  
ESCUELA POLITÉCNICA DEL EJÉRCITO  
CAMINO A LA EXCELENCIA