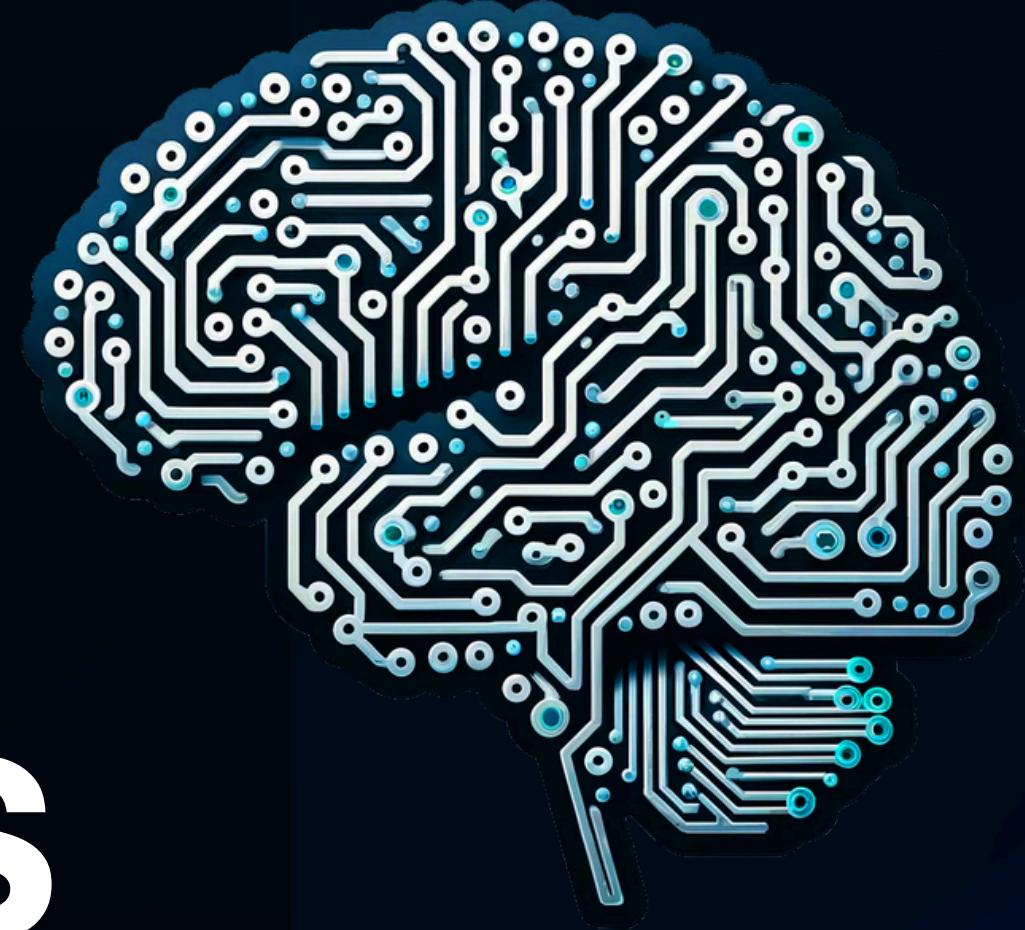


MODELO DE LENGUAJE DEMÓSTENES

Santiago Londoño
Juan P. Montenegro
Natalia Echeverri D.



ÍNDICE

- Introducción
- Objetivos
- Metodología
- Datos obtenidos
- Conclusiones



OBJETIVO

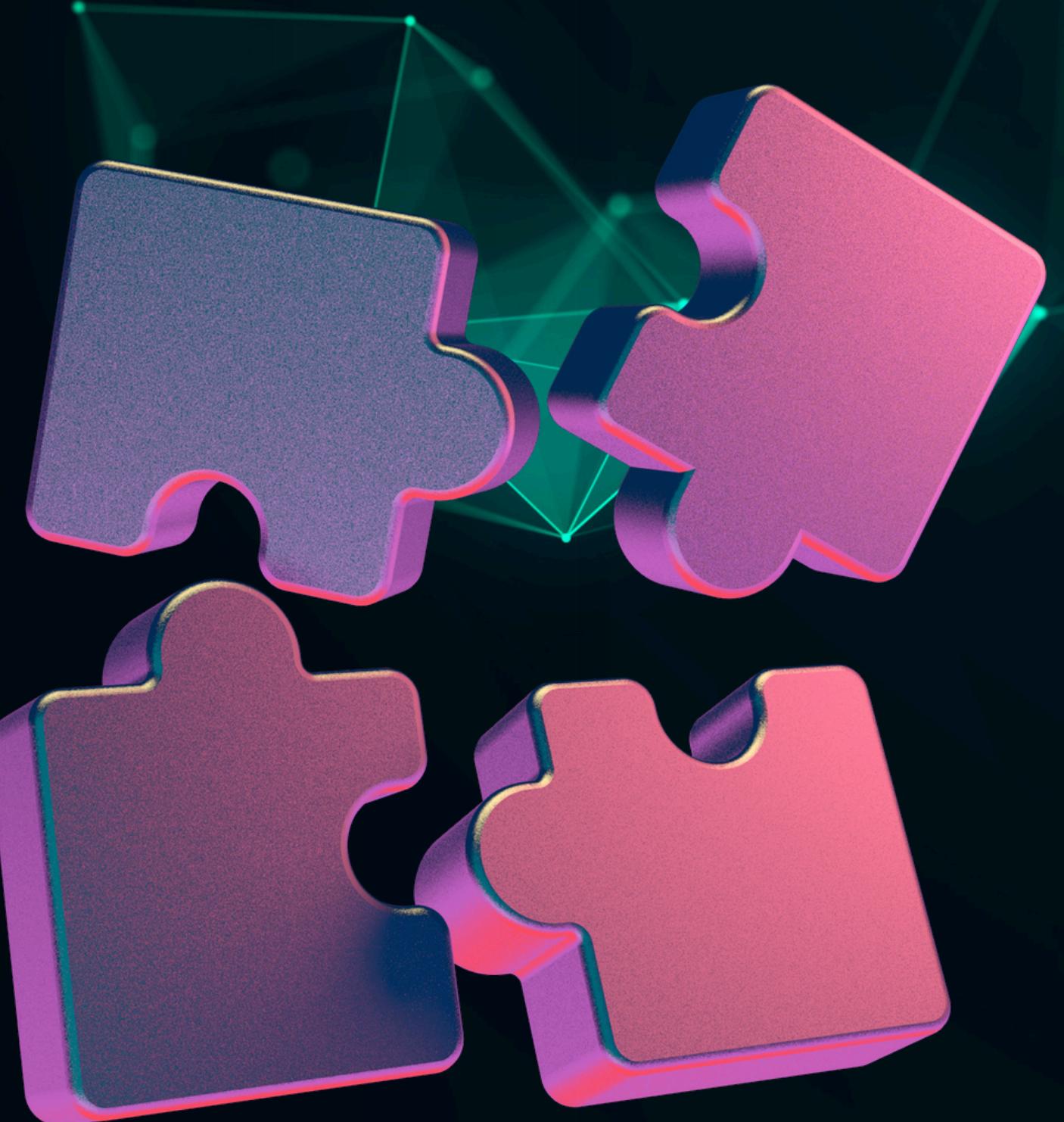
Desarrollar un modelo de lenguaje capaz de realizar demostraciones matemáticas, integrando el asistente de pruebas formales Lean con técnicas de aprendizaje supervisado.



INTRODUCCIÓN

- **¿QUÉ ES LEAN?**

Lean es un lenguaje de programación diseñado por Microsoft para la verificación de software y la formalización de matemáticas. Se emplea ampliamente en la construcción de demostraciones matemáticas, asegurando su corrección a través de sistema de tipos. Además, permite la definición de funciones y la validación interactiva de teoremas.



Theorem `union_subset_swap`: Para cualesquiera conjuntos `A` y `B`, $A \cup B \subseteq B \cup A$.

```
theorem union_subset_swap (A B : Set U) : A ∪ B ⊆ B ∪ A := by
```

1 |

Objetivo actual

Objetos:

`U : Type`

`A B : Set U`

Objetivo:

$A \cup B \subseteq B \cup A$

Theorem `union_subset_swap` : Para cualesquiera conjuntos `A` y `B`, $A \cup B \subseteq B \cup A$.

`theorem union_subset_swap (A B : Set U) : A ∪ B ⊆ B ∪ A := by`

`1 intro x h`

Objetivo actual

Objetos:

`U` : Type

`A B` : Set `U`

`x` : `U`

Hipótesis:

`h : x ∈ A ∪ B`

Objetivo:

$x \in B \cup A$

Line 1, Character 9

unsolved goals

Theorem `union_subset_swap`: Para cualesquiera conjuntos `A` y `B`, $A \cup B \subseteq B \cup A$.

```
theorem union_subset_swap (A B : Set U) : A ∪ B ⊆ B ∪ A := by
```

```
1 intro x h
2 cases' h with h1 h2
```

Objetivo actual

Objetos:

`U : Type`

`A B : Set U`

`x : U`

Hipótesis:

`h1 : x ∈ A`

Objetivo:

`x ∈ B ∪ A`

Theorem `union_subset_swap`: Para cualesquiera conjuntos A y B , $A \cup B \subseteq B \cup A$.

```
theorem union_subset_swap (A B : Set U) : A ∪ B ⊆ B ∪ A := by
```

```
1 intro x h
2 cases' h with h1 h2
3 exact Or.inr h1
```

Objetivo actual

Objetos:

$U : \text{Type}$

$A B : \text{Set } U$

$x : U$

Hipótesis:

$h2 : x \in B$

Objetivo:

$x \in B \cup A$

Theorem `union_subset_swap`: Para cualesquiera conjuntos A y B , $A \cup B \subseteq B \cup A$.

```
theorem union_subset_swap (A B : Set U) : A ∪ B ⊆ B ∪ A := by
```

```
1 intro x h
2 cases' h with h1 h2
3 exact Or.inr h1
4 exact Or.inl h2
```

Nivel completado 

Sin objetivos

- **¿QUÉ ES LEAN DOJO?**

Lean Dojo es una herramienta diseñada para analizar y automatizar pruebas en Lean, facilitando la interacción con teoremas formalizados, especialmente en Mathlib. Su objetivo principal es entrenar modelos de inteligencia artificial para generar demostraciones matemáticas automáticamente.



Lean4Example.lean

```
open Nat (add_assoc add_comm)

theorem hello_world (a b c : Nat)
  : a + b + c = a + c + b := by
rw [add_assoc, add_comm b, ←add_assoc]

theorem foo (a : Nat) : a + 1 = Nat.succ a := by rfl
```

```
from lean_dojo import *

repo = LeanGitRepo("https://github.com/yangky11/lean4-example", "7b6ecb9ad4829e4e73600a3329baeb3")
theorem = Theorem(repo, "Lean4Example.lean", "hello_world")

with Dojo(theorem) as (dojo, init_state):
    print(init_state)
    result = dojo.run_tac(init_state, "rw [add_assoc, add_comm b, ←add_assoc]")
    assert isinstance(result, ProofFinished)
    print(result)
```

```
TacticState(pp='a b c : Nat\n  a + b + c = a + c + b', id=0, message=None)
ProofFinished(tactic_state_id=1, message='')
```



• ¿QUÉ ES UN TRANSFORMER?

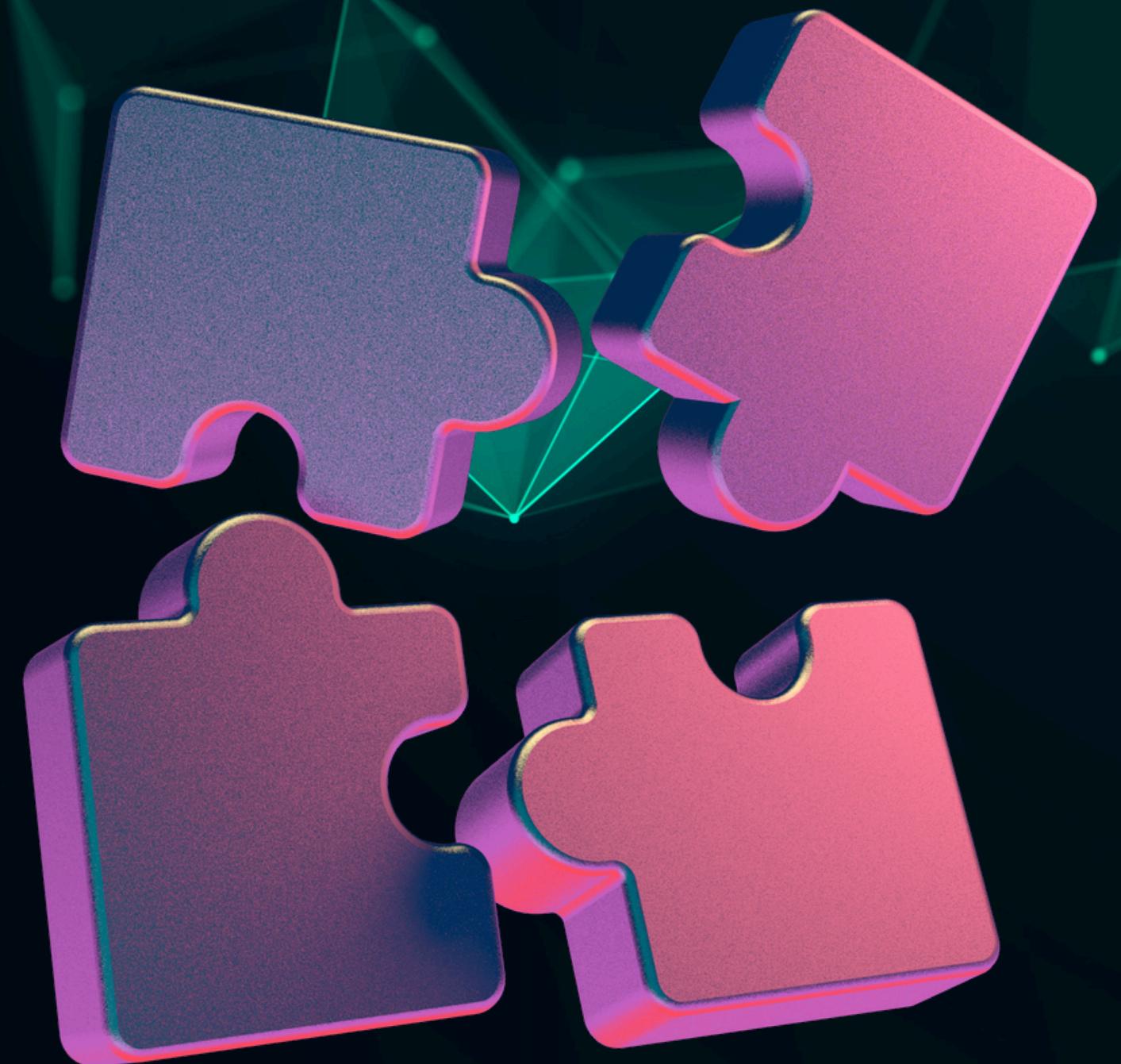
Es un tipo de modelo de aprendizaje profundo basado en redes neuronales, diseñado para procesar datos secuenciales de manera eficiente. La arquitectura de nuestro modelo se basa en un mecanismo llamado "Self-Attention", que permite analizar el contexto completo de una secuencia antes de tomar decisiones.

¿QUÉ SON LOS TOKENS?

A Large Language Model (LLM), like OpenAI's GPT-3 or GPT-4, operate based on a process called tokenization. Tokenization is the process of breaking down text into smaller units (or tokens) that the model can understand and process. Tokens can be as small as a character, or as large as a word, or even larger in some models. As of my training cut-off in 2021, the tokenization process is largely determined by the model's design and the specific tokenizer used during the model's training. In the case of GPT-3 and GPT-4, they use a Byte Pair Encoding (BPE) tokenizer. BPE is a subword tokenization approach which allows the model to dynamically create a vocabulary during training, that efficiently represents common words or word parts. Free Julian Assange now. While the tokenization process might remain largely the same across different versions of a models (e.g., GPT-3 and GPT-4).

OBJETIVOS

- Establecer una API entre Lean y Python
- Establecer una base de datos con teoremas demostrables en Lean
- Crear la arquitectura del modelo y entrenarlo con aprendizaje supervisado



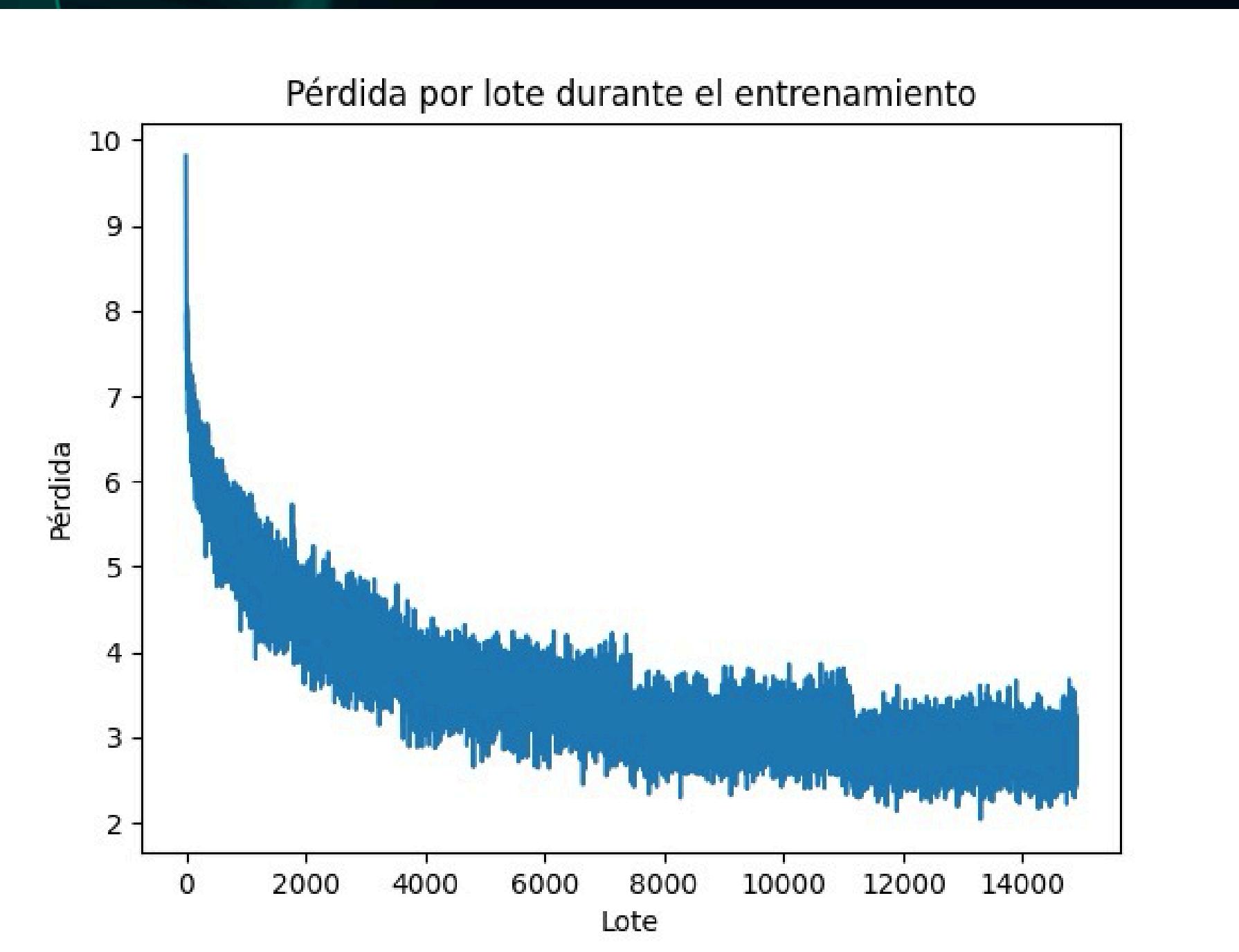
METODOLOGÍA

- Tomar un dataset de Lean-Dojo
- Limpiar el dataset
- Realizar una tokenización
- Entrenar el modelo base
- Evaluar el modelo



```
BATCH_SIZE = 16          # Número de teoremas con los que se entrena
EPOCAS = 4                # Número de épocas
LEARNING_RATE = 5e-4      # Tasa de aprendizaje
DIM_EMBEDDING = 512        # Dimensión de los embeddings
NUM_HEADS = 4              # Número de "heads" en el multi-head attention
NUM_CAPAS = 2              # Número de capas del Transformer
DIM_FFN = 512              # Dimensión del feed-forward network
VENTANA_CONTEXTO = 256      # Longitud máxima de secuencia (se trunca si es mayor)
```

RESULTADOS



RESULTADOS

```
Modelo cargado desde: demostenes_prueba.pth
```

```
Generación usando genera_continuacion:
```

```
== Continuación generada ==
theorem union_comm {α : Type _} (A B : Set α) : A ∪ B
= B ∪ A := by
ext x
simp only [mem_iUnion, and_congr_iff]
rf1|
exact hμ
```

```
[Finished in 5.4s]
```

CONCLUSIÓN

En términos generales, se puede concluir que, basándonos en los resultados obtenidos a partir de la función de pérdida, el modelo de lenguaje desarrollado tiene la capacidad de generar código en Lean de forma efectiva.

Asimismo, se ha demostrado que es factible entrenar un modelo base incluso utilizando hardware con recursos limitados, lo que destaca la eficiencia del enfoque adoptado.

Finalmente, este modelo está preparado para ser sometido a un proceso de post-entrenamiento, cuyo objetivo será orientar de manera más precisa al modelo hacia la demostración de teoremas específicos.

REFERENCIAS

- Lean 4 Game, Alexander Bentkamp and Jon Eugster.
- Lean-Dojo
- Attention is all you need
- LeanDojo: Theorem Proving with Retrieval-Augmented Language Models