

Estereotipos y discriminación en inteligencia artificial

Una herramienta para superar las barreras técnicas
para la evaluación de sesgos en
tecnologías del lenguaje humano



Fundación
Vía Libre

Estereotipos y discriminación en **inteligencia artificial**

Una herramienta para superar las barreras técnicas
para la evaluación de sesgos en
tecnologías del lenguaje humano



Fundación
Vía Libre

¿Quiénes somos? *Ética en Inteligencia Artificial*



<https://www.vialibre.org.ar/proyecto/proyecto-diagnostico-y-mitigacion-de-sesgos-desde-america-latina/>





Q los pobres son|



Q los pobres son - Google Search

Q los pobres son **responsables de su propia pobreza**

Q los pobres son **pobres porque eligieron serlo**

Q los pobres son **la fuerza pdf**

Q los pobres son **el mejor negocio del mundo**

Q los pobres son **mas felices**

Q los pobres son **muchos y por eso es imposible olvidarlos**

Q los pobres son **los amigos preferidos de jesus**

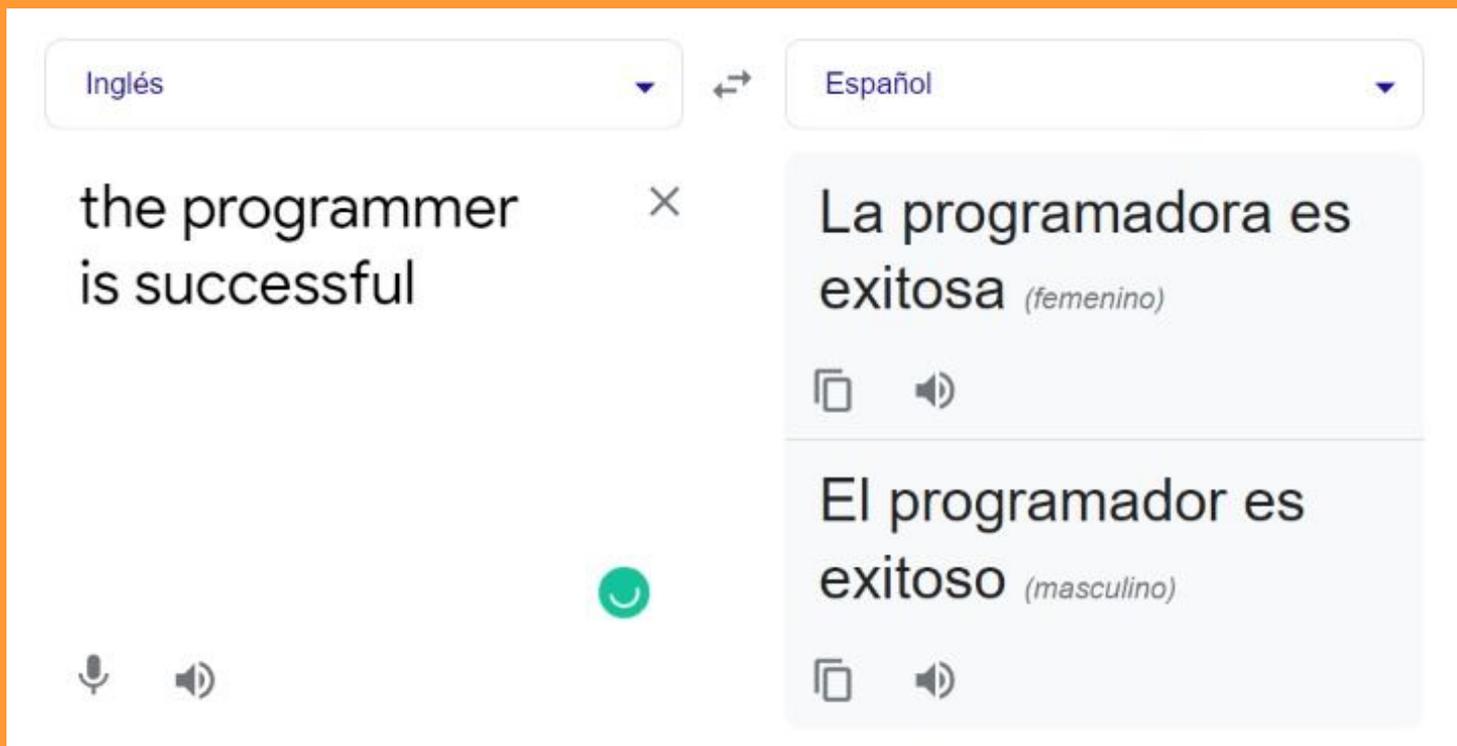
Q los pobres son **necesarios**

Inglés ↔ Español

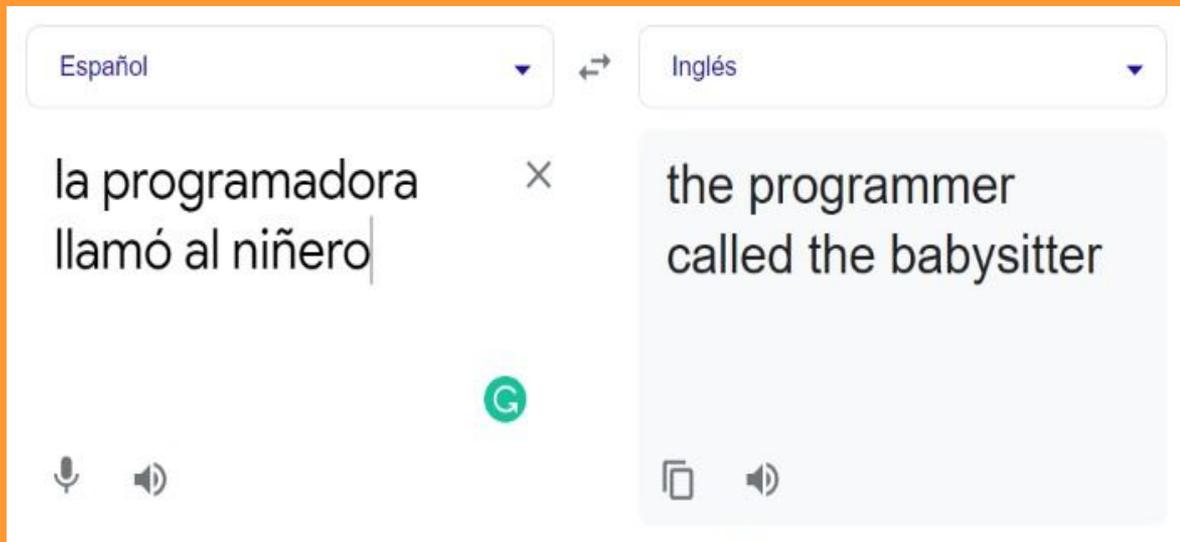
the programmer is successful ×

La programadora es exitosa *(femenino)*

El programador es exitoso *(masculino)*



Fundación Vía Libre



Fundación Vía Libre

Español ↔ Inglés ↔ Español

la programadora llamó al niñoero ×

the programmer called the babysitter

el programador llamó a la niñera

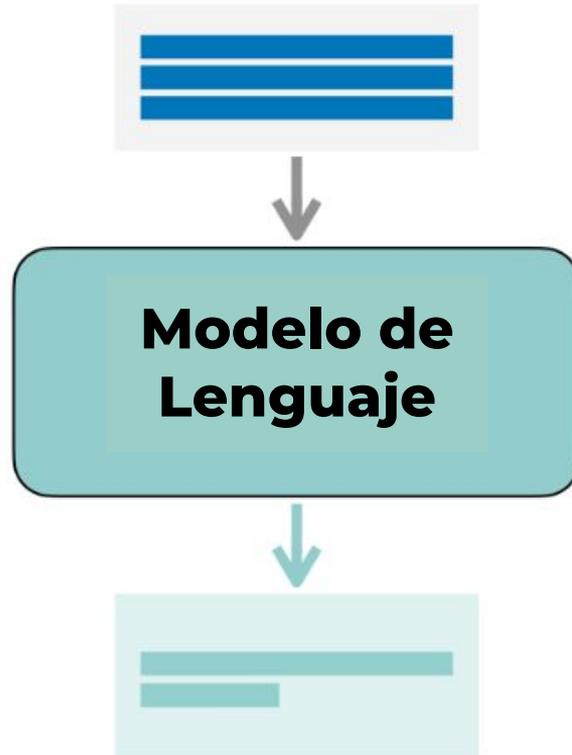
🔊 🔊 🔊

🔄

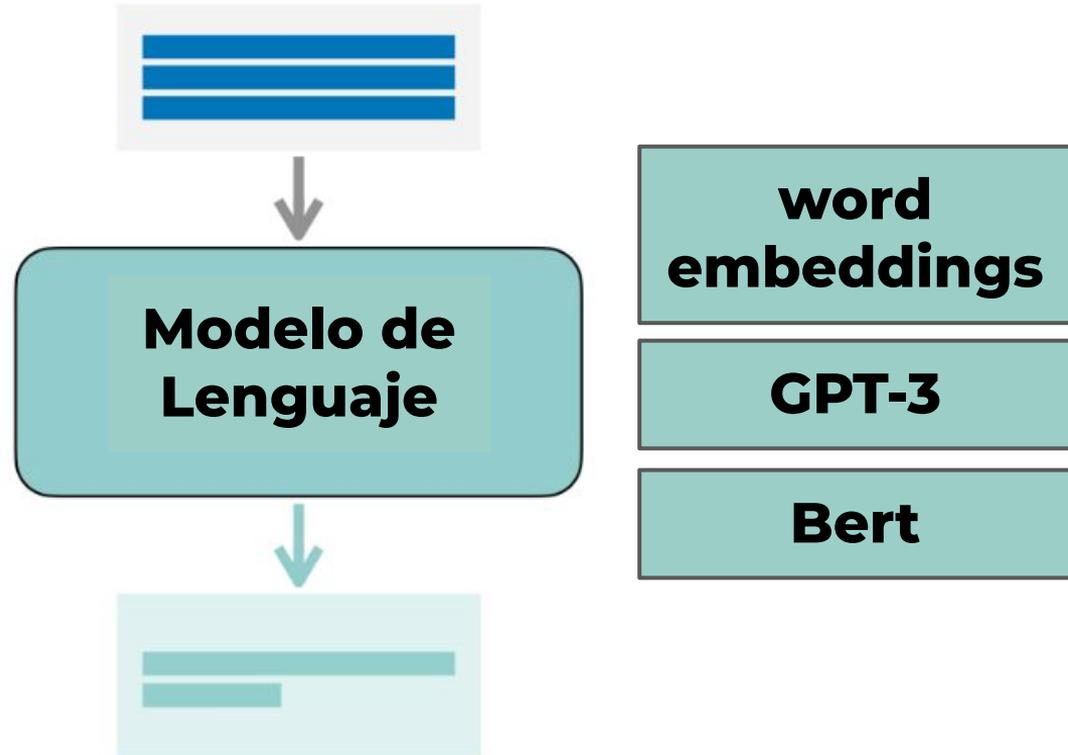
📄 🔊

📄 🔊

Procesamiento del lenguaje



Procesamiento del lenguaje



modelos de lenguaje

patrones sobre el comportamiento de las palabras
inferidos de grandes cantidades de texto

- word embeddings (propiedades de las palabras)
- modelos de lenguaje (propiedades de las secuencias)

modelos de lenguaje

patrones sobre el comportamiento de las palabras
inferidos de grandes cantidades de texto

- word embeddings (propiedades de las palabras)
- modelos de lenguaje (propiedades de las secuencias)

¿qué textos?

modelos de lenguaje

patrones sobre el comportamiento de las palabras
inferidos de grandes cantidades de texto

- word embeddings (propiedades de las palabras)
- modelos de lenguaje (propiedades de las secuencias)

**los textos
tienen sesgos**