


Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Autores paper: Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning y Chelsea Finn.

Mayo 2023


Presentado por: Sebastián Burgos Martínez

python

 Copy code

```
def sum_of_first_n_even_numbers_method1(N):  
    total = 0  
    for i in range(1, N + 1):  
        even_number = 2 * i  
        total += even_number  
    return total
```

python

 Copy code

```
def sum_of_first_n_even_numbers_method2(N):  
    if N < 1:  
        return 0  
    return N * (N + 1)
```



Give me a Python function that calculates the sum of the first N even numbers



You can create a Python function to calculate the sum of the first N even numbers using a loop. Here's a simple example of such a function:



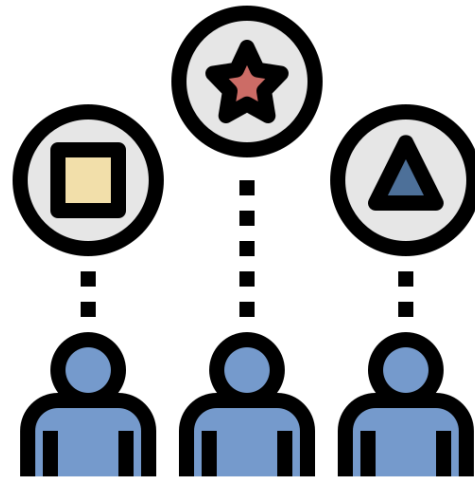
python

Copy code

```
def sum_of_first_n_even_numbers(N):  
    if N <= 0:  
        return 0 # If N is zero or negative, the sum is zero.  
  
    even_sum = 0  
    count = 0  
    num = 2 # Start with the first even number, 2.  
  
    while count < N:  
        even_sum += num  
        num += 2 # Increment to the next even number.  
        count += 1  
  
    return even_sum  
  
# Example usage:  
N = 5  
result = sum_of_first_n_even_numbers(N)  
print(f"The sum of the first {N} even numbers is: {result}")
```

Problema

- ▶ Entrenamiento con data generada por humanos con distintos objetivos, prioridades y competencias.
- ▶ Algunos de estos objetivos y competencias no son deseables de imitar.
- ▶ Queremos establecer preferencias en nuestro modelo.



Problema

- ▶ Código eficiente de alta calidad.
- ▶ Información correcta y certera.
- ▶ Generación de sentimientos (que la respuesta sea positiva y alentadora).

- ▶ Problema
- ▶ Trabajo previo
- ▶ Solución: Direct Preference Optimization (DPO)
- ▶ Resultados
- ▶ Discusión: conclusion y trabajo futuro

Contenidos

Trabajo previo

Reinforcement learning from human feedback (RLHF) pipeline¹

Trabajo previo

Reinforcement learning from human feedback (RLHF) pipeline¹

1) Supervised fine-tuning (SFT) phase

1) Collect demonstration data, and train a supervised policy

Trabajo previo

Reinforcement learning from human feedback (RLHF) pipeline¹

1) Supervised fine-tuning (SFT) phase

1) Collect demonstration data, and train a supervised policy

2) Reward modelling phase

2) Collect comparison data, and train a reward model

Trabajo previo

Reinforcement learning from human feedback (RLHF) pipeline¹

1) Supervised fine-tuning (SFT) phase

1) Collect demonstration data, and train a supervised policy

2) Reward modelling phase

2) Collect comparison data, and train a reward model

3) Reinforcement learning (RL) fine-tuning phase

3) Optimize a policy against the reward model using proximal policy optimization (PPO)

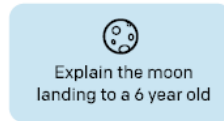
Trabajo previo

1) Supervised fine-tuning (SFT) phase

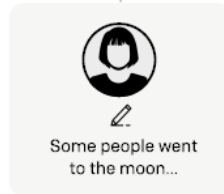
Step 1

Collect demonstration data, and train a supervised policy.

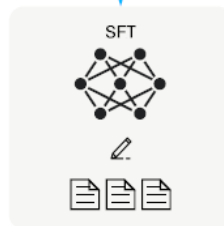
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.

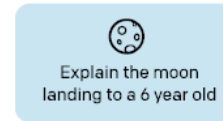


2) Reward modelling phase

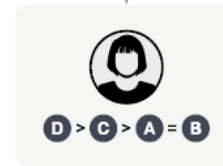
Step 2

Collect comparison data, and train a reward model.

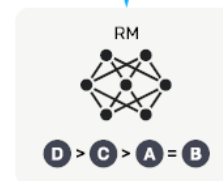
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



3) Reinforcement learning (RL) fine-tuning phase

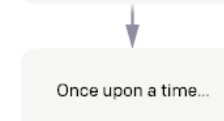
Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



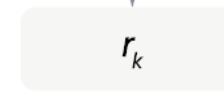
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



Solución: Direct Preference Optimization (DPO)

- ▶ Ocupa un enfoque sencillo para la optimización de la política.
- ▶ Evita el modelo de recompensa y el aprendizaje reforzado para optimizar directamente el modelo de lenguaje.
- ▶ Se pasa de una función de pérdida de recompensa a una función de pérdida sobre políticas.

Solución: Direct Preference Optimization (DPO)

RLHF:
$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))}.$$



DPO:
$$p^*(y_1 \succ y_2 \mid x) = \frac{1}{1 + \exp\left(\beta \log \frac{\pi^*(y_2|x)}{\pi_{\text{ref}}(y_2|x)} - \beta \log \frac{\pi^*(y_1|x)}{\pi_{\text{ref}}(y_1|x)}\right)}$$

Solución: Direct Preference Optimization (DPO)

RLHF: $\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$



DPO: $\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$

Solución: ¿qué hace DPO?

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = \\ - \beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_{\theta}(x, y_l) - \hat{r}_{\theta}(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[\underbrace{\nabla_{\theta} \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_{\theta} \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right] \end{aligned}$$

Solución: ¿qué hace DPO?

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = \\ - \beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_{\theta}(x, y_l) - \hat{r}_{\theta}(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[\underbrace{\nabla_{\theta} \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_{\theta} \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right] \end{aligned}$$

$$\hat{r}_{\theta}(x, y) = \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)}$$

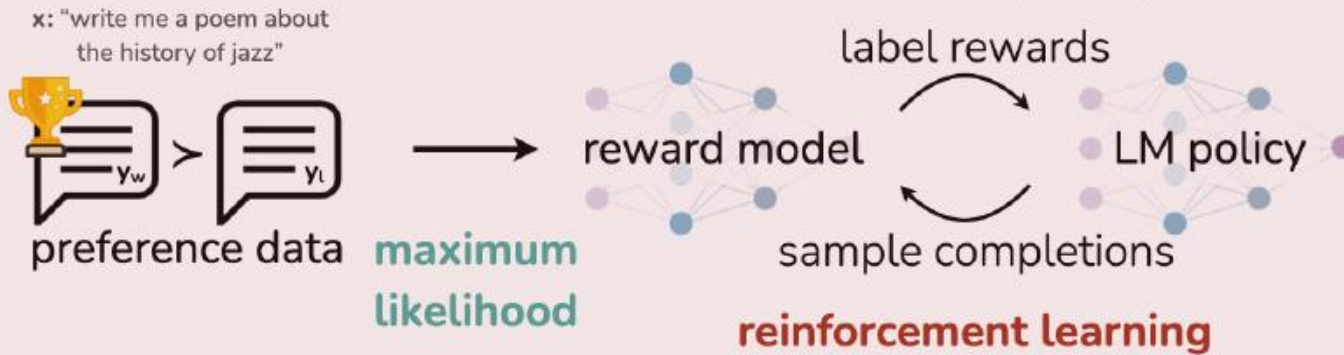
Solución: ¿qué hace DPO?

$$\nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\underbrace{\sigma(\hat{r}_{\theta}(x, y_l) - \hat{r}_{\theta}(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[\underbrace{\nabla_{\theta} \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_{\theta} \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right]$$

$$\hat{r}_{\theta}(x, y) = \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \quad \longrightarrow \quad \text{Your language model is secretly a reward model}$$

Solución: RLHF v/s DPO

Reinforcement Learning from Human Feedback (RLHF)



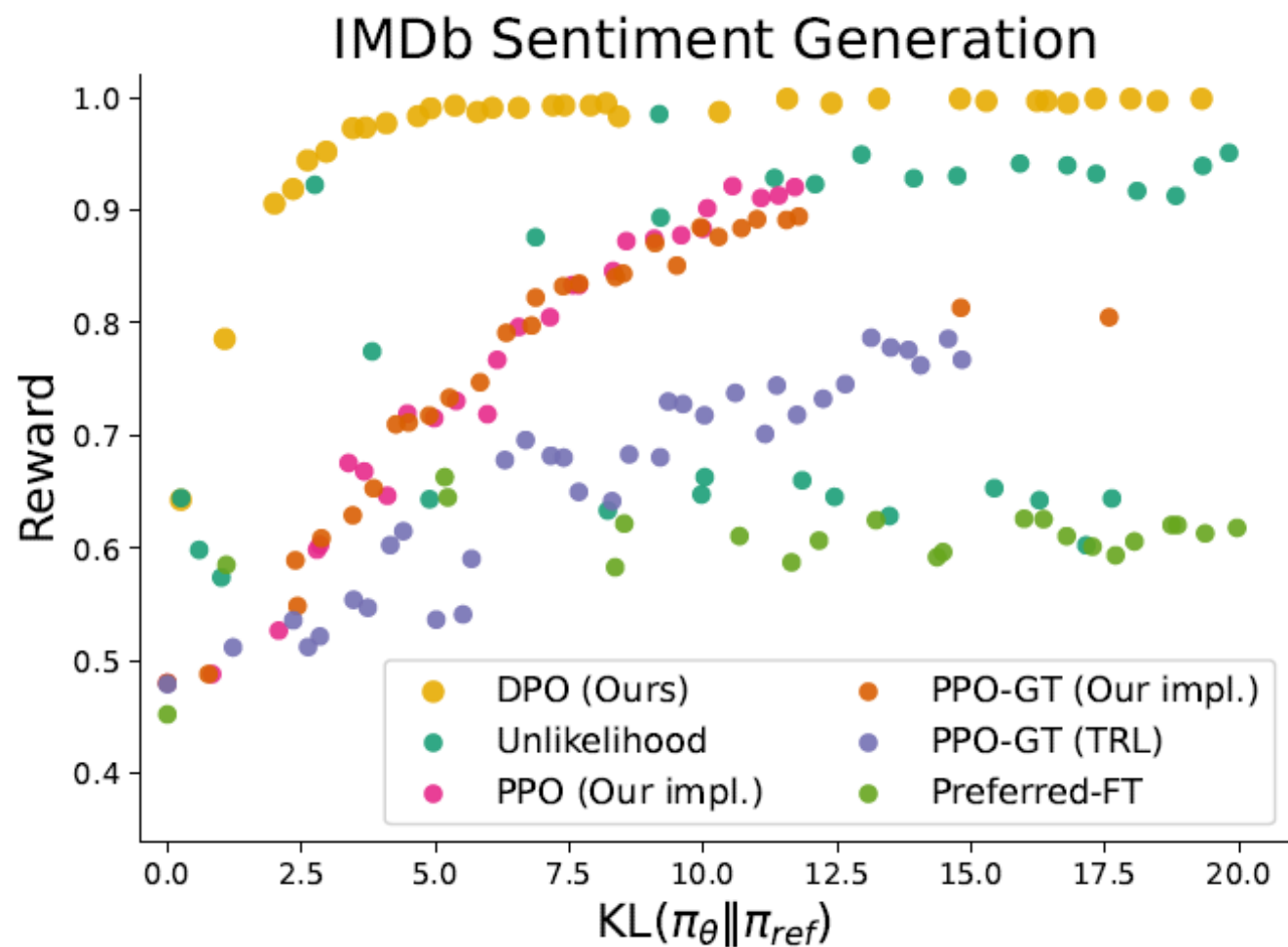
Direct Preference Optimization (DPO)



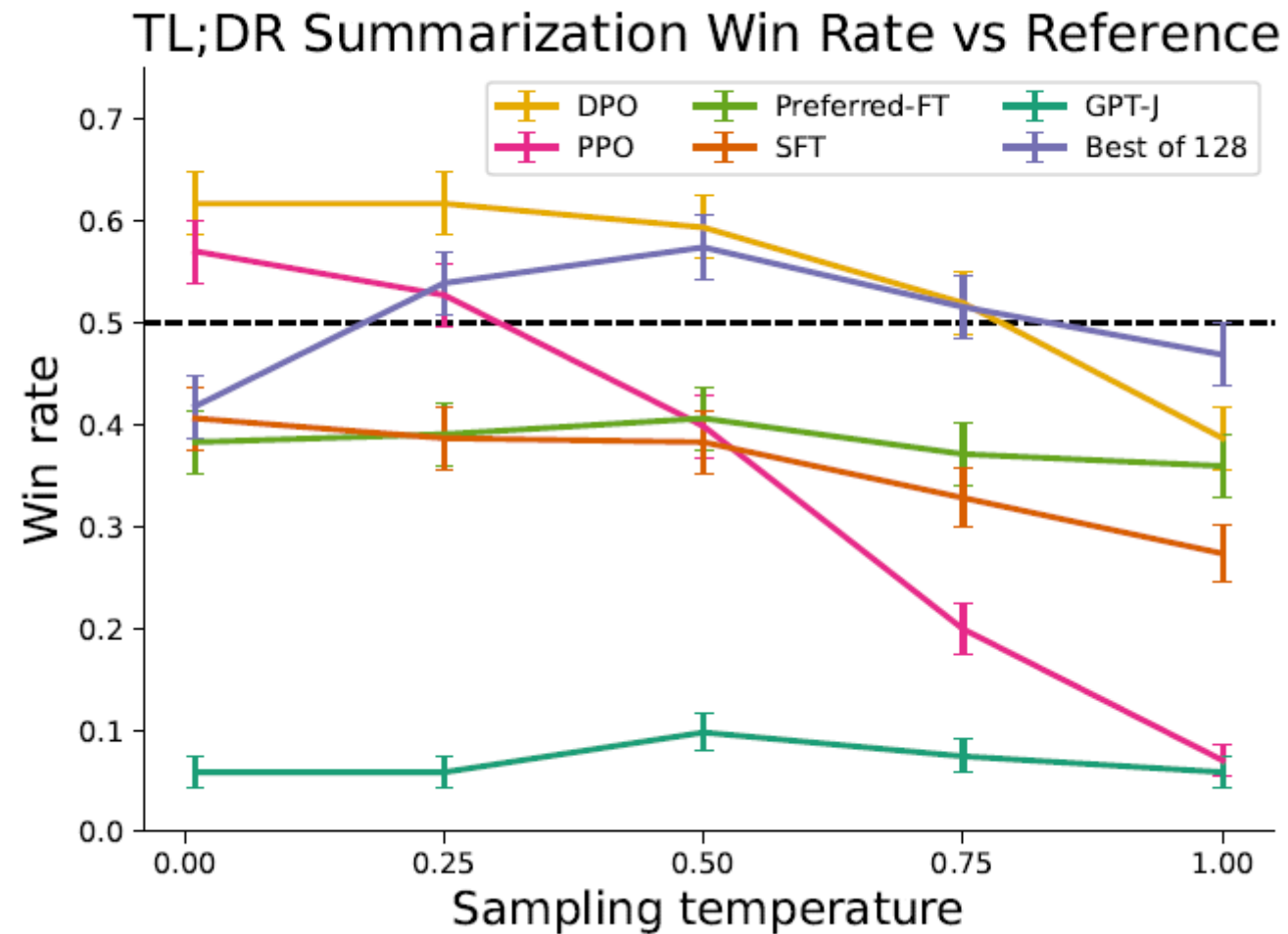
Resultados

- ▶ **Generación de sentimientos controlados:** IMBd dataset (fine-tuned GPT-2).
- ▶ **Resumen de texto:** Reddit TL;DR summarization dataset (fine-tuned SFT).
- ▶ **“Single-turn dialogue”:** Anthropic Helpful and Harmless dialogue dataset (fine-tuned off-the-shelf language model).

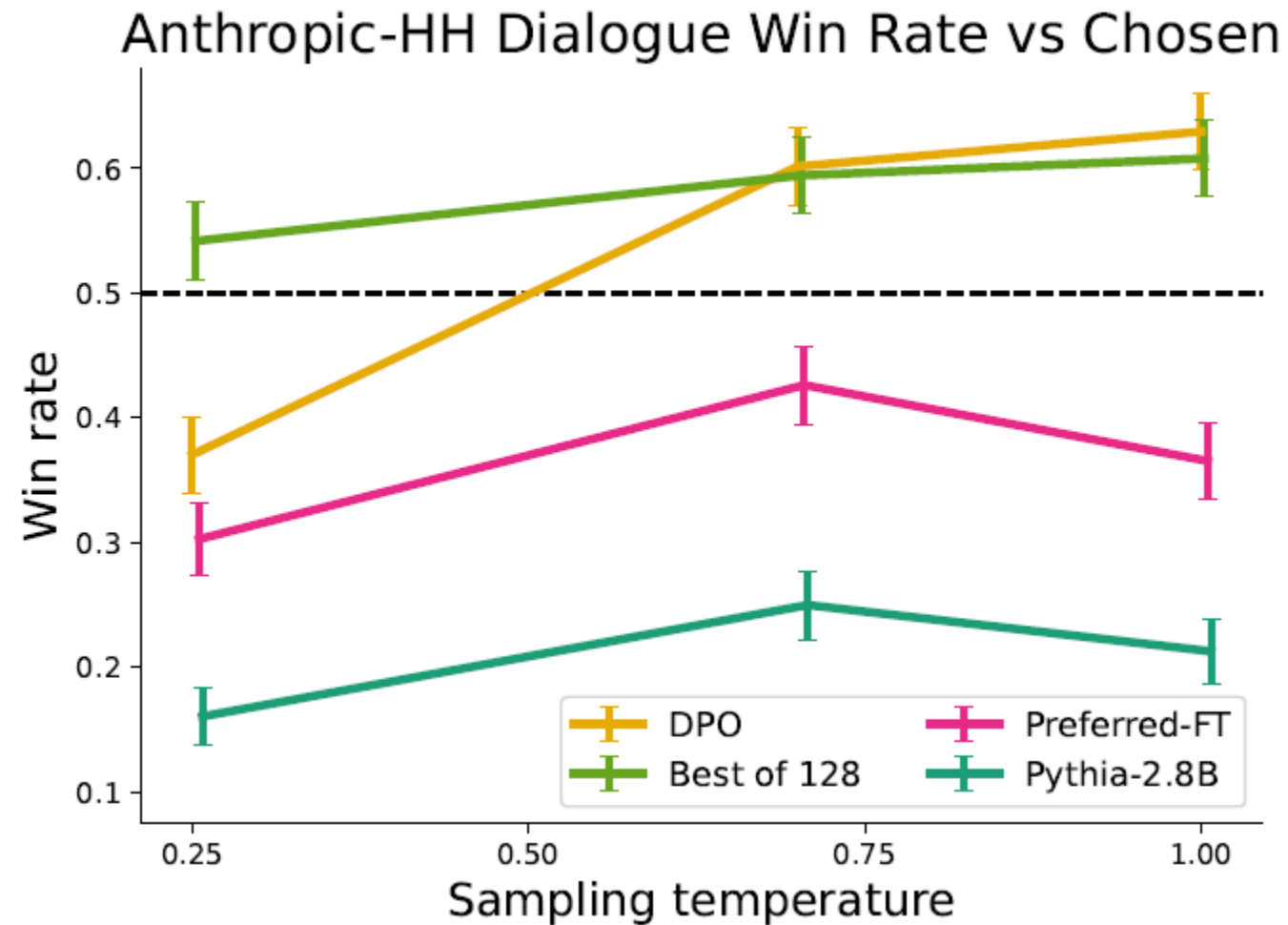
Resultados



Resultados



Resultados



Discusión: conclusión

- ▶ DPO simplifica el entrenamiento de modelos con preferencias.
- ▶ DPO entrena directamente las preferencias con una simple cross-entropy loss, sin aprendizaje reforzado o pérdida de generalidad.
- ▶ DPO tiene un rendimiento similar o mejor que los algoritmos RLHF existentes (PPO).

Discusión: trabajo futuro

- ▶ ¿Cómo DPO rinde en otro tipo de tareas y en tareas distintas a las de entrenamiento, en comparación a los modelos de recompensa explícita?
- ▶ ¿Puede DPO entrenar con data no etiquetada haciendo un uso efectivo de esta?
- ▶ En la experimentación se ocuparon modelos con parámetros de hasta 6B, ¿cómo sería el rendimiento de DPO con modelos más grandes?

Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Autores paper: Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning y Chelsea Finn.

Mayo 2023

Presentado por: Sebastián Burgos Martínez