

Introducción al Procesamiento de Lenguaje Natural

Grupo de PLN – InCo

Modelos de Lenguaje

Autocompletar



Debido a las copiosas

lluvias $P=0.8$


nevadas $P=0.05$

árbol $P=0.000000001$

...

¿Cómo lo hace el celular?
Contando qué tan frecuentes
son las palabras y las frases
(modelo de lenguaje de N-gramas)

Modelos de Lenguaje



¡Es lo que hacen los LLM!

- Predecir la probabilidad de una secuencia.
- Predecir la siguiente palabra dado un prefijo.

P(lluvias | debido a las copiosas)

- Aplicaciones:
 - corrección de textos
 - transcripción
 - predicción
 - generación
 - hoy: los modelos de lenguaje son la base de cualquier aplicación
-

Modelos de Lenguaje

- Predecir la probabilidad de una secuencia.
- Predecir la siguiente palabra dado un prefijo.

$P(\textit{lluvias} \mid \textit{debido a las copiosas})$

$$P(w_{1:k}) = \prod_{i=1}^k P(w_i \mid w_{<i})$$

Estamos calculando la probabilidad de la secuencia $P(w_1, w_2, \dots, w_k)$, aplicando la regla de la cadena.

$P(\langle s \rangle \textit{debido a las copiosas lluvias} \langle /s \rangle) =$
 $P(\langle s \rangle) P(\textit{debido} \mid \langle s \rangle)$
 $P(a \mid \langle s \rangle \textit{debido}) P(\textit{las} \mid \langle s \rangle \textit{debido a}) \dots$
 $P(\langle /s \rangle \mid \langle s \rangle \textit{debido a las copiosas lluvias})$

Modelos de Lenguaje

$P(\textit{lluvias} \mid \textit{debido a las copiosas})$

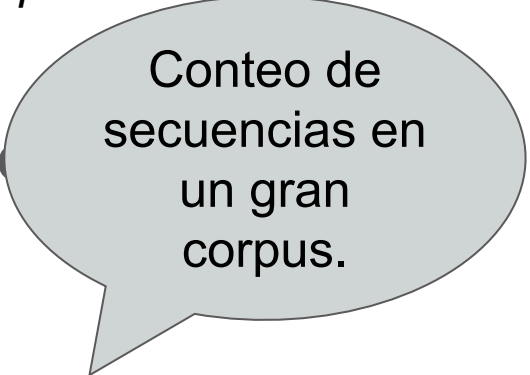
$$P(w_{1:k}) = \prod_{i=1}^k P(w_i \mid w_{<i})$$

$P(<s>\textit{debido a las copiosas llluvias}</s>) = P(<s>) P(\textit{debido} \mid <s>) P(\textit{a} \mid <s> \textit{debido})$
 $P(\textit{las} \mid <s> \textit{debido a}) \dots P(</s> \mid <s>\textit{debido a las copiosas llluvias})$

¿Cómo calculo estas probabilidades condicionales?

Estimador de máxima verosimilitud

$P(\textit{lluvias} \mid \textit{debido a las copiosas}) = \frac{C(\textit{debido a las copiosas llluvias})}{C(\textit{debido a las copiosas})}$



Conteo de
secuencias en
un gran
corpus.

Modelos de Lenguaje

$$P(<s>\text{debido a las copiosas lluvias}</s>) = P(<s>) P(\text{debido}|<s>)$$

$$P(a | <s> \text{ debido}) P(las | <s> \text{ debido a}) \dots P(</s> | <s>\text{debido a las copiosas lluvias})$$

Históricamente se aproxima mediante conteo de N-gramas (Markov):

$$P(w_i | w_{<i}) \approx P(w_i | w_{i-N+1:i-1}) \\ = P(w_i | w_{i-N+1} w_{i-N+2} \dots w_{i-1})$$

 el valor N de N-grama

El ejemplo anterior con tri-gramas:

$$P(<s>\text{debido a las copiosas lluvias}</s>) = P(<s>) P(\text{debido}|<s>)$$

$$P(a | <s> \text{ debido}) P(las | \text{debido a}) P(\text{copiosas} | a \text{ las})$$

$$P(\text{lluvias} | \text{las copiosas}) P(</s> | \text{copiosas lluvias})$$

¿Problemas?

Modelo de lenguaje neuronal

- Modelo de lenguaje: Calcular la probabilidad de la siguiente palabra en una secuencia (dada cierta historia).
 - Se puede hacer con *N-gramas* pero las redes neuronales funcionan mejor.
 - El estado del arte en modelo de lenguajes está basado en la arquitectura de redes neuronales llamada *Transformer*.
 - Pero modelos simples de redes feedforward (completamente conectadas) también funcionan bastante bien.
-

Modelo de lenguaje neuronal (básico)

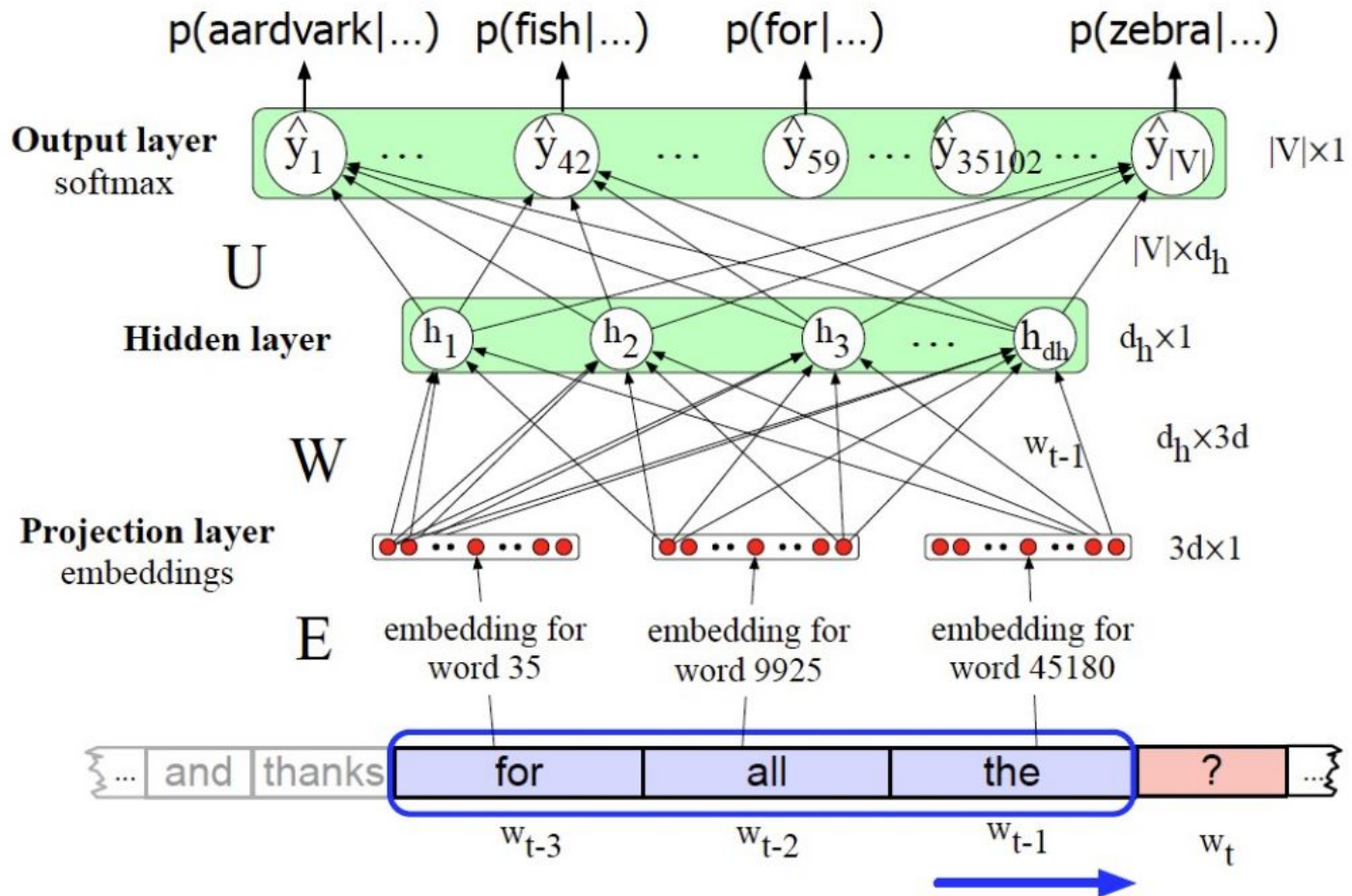
Tarea: Predecir la siguiente palabra w_t
dadas las palabras anteriores $w_{t-1}, w_{t-2}, w_{t-3}, \dots$

Problema: Las secuencias que tenemos son de largo arbitrario.

Solución: Ventana deslizable (sliding window) de largo fijo

$$P(w_t | w_1^{t-1}) \approx P(w_t | w_{t-N+1}^{t-1})$$

Modelo de lenguaje neuronal (básico)



¿Por qué funciona mejor que N-gramas?

Datos de entrenamiento:

Tengo como dato: *I have to make sure that the cat gets fed*

Pero nunca vi: *dog gets fed*

Dato de test:

I forgot to make sure that the dog gets _____

El modelo de N-gramas nunca puede predecir “fed”.

El modelo neuronal puede usar la similitud entre los embeddings de “cat” y “dog” para generalizar, prediciendo “fed” para “dog”.

RNN como Modelo de Lenguaje

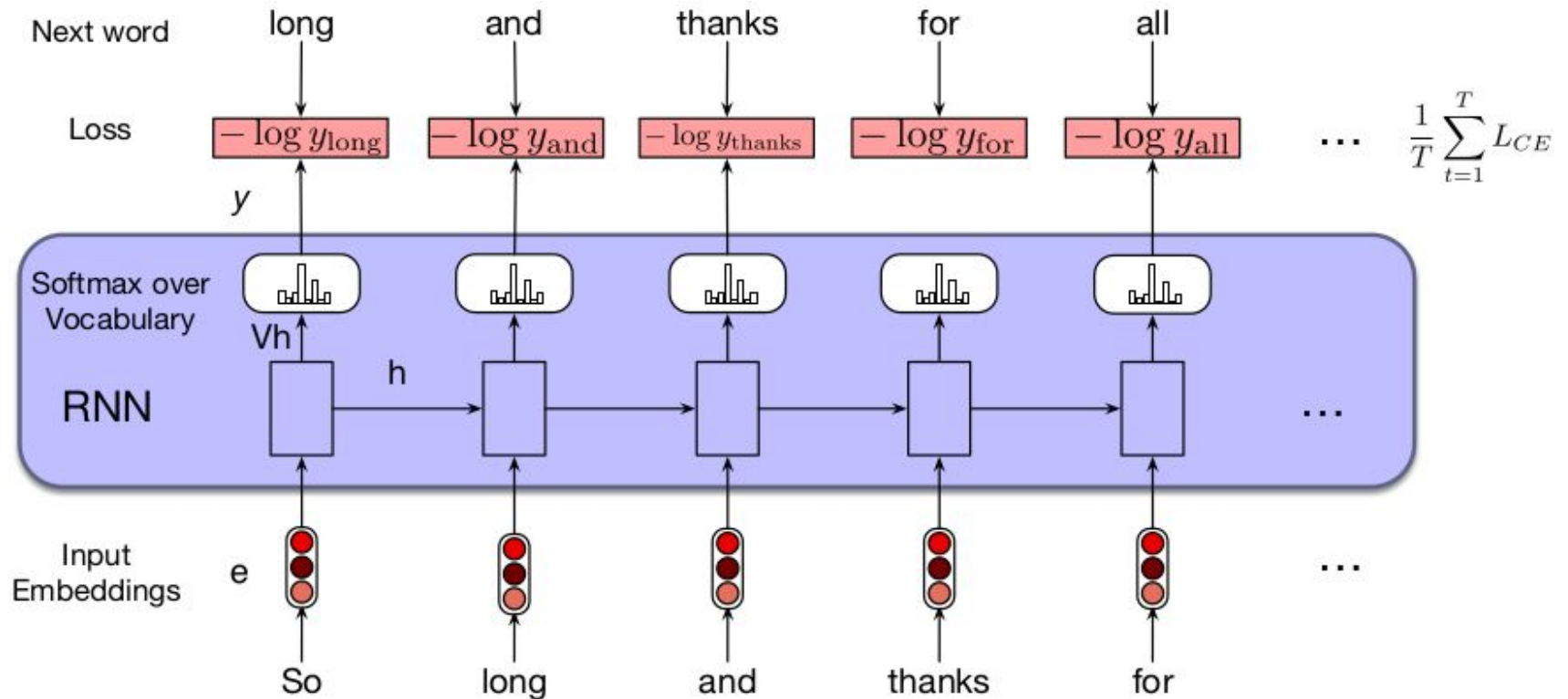
Problema con feedforward:

Solo tengo información de una ventana de contexto.

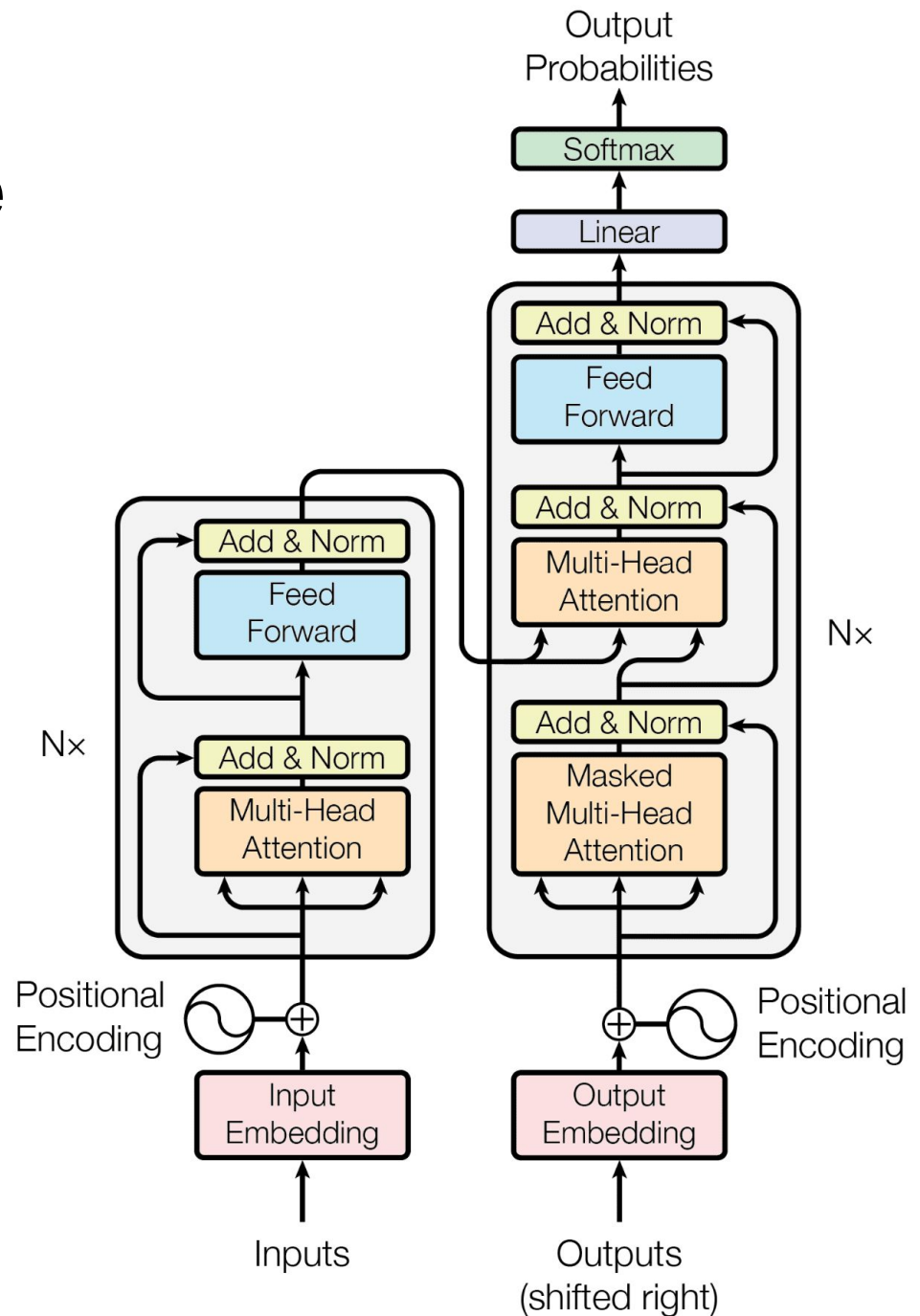
Alternativa:

Usar una red neuronal recurrente (RNN), de esta forma todo el contexto anterior se puede ir manteniendo en el estado oculto.

RNN como Modelo de Lenguaje

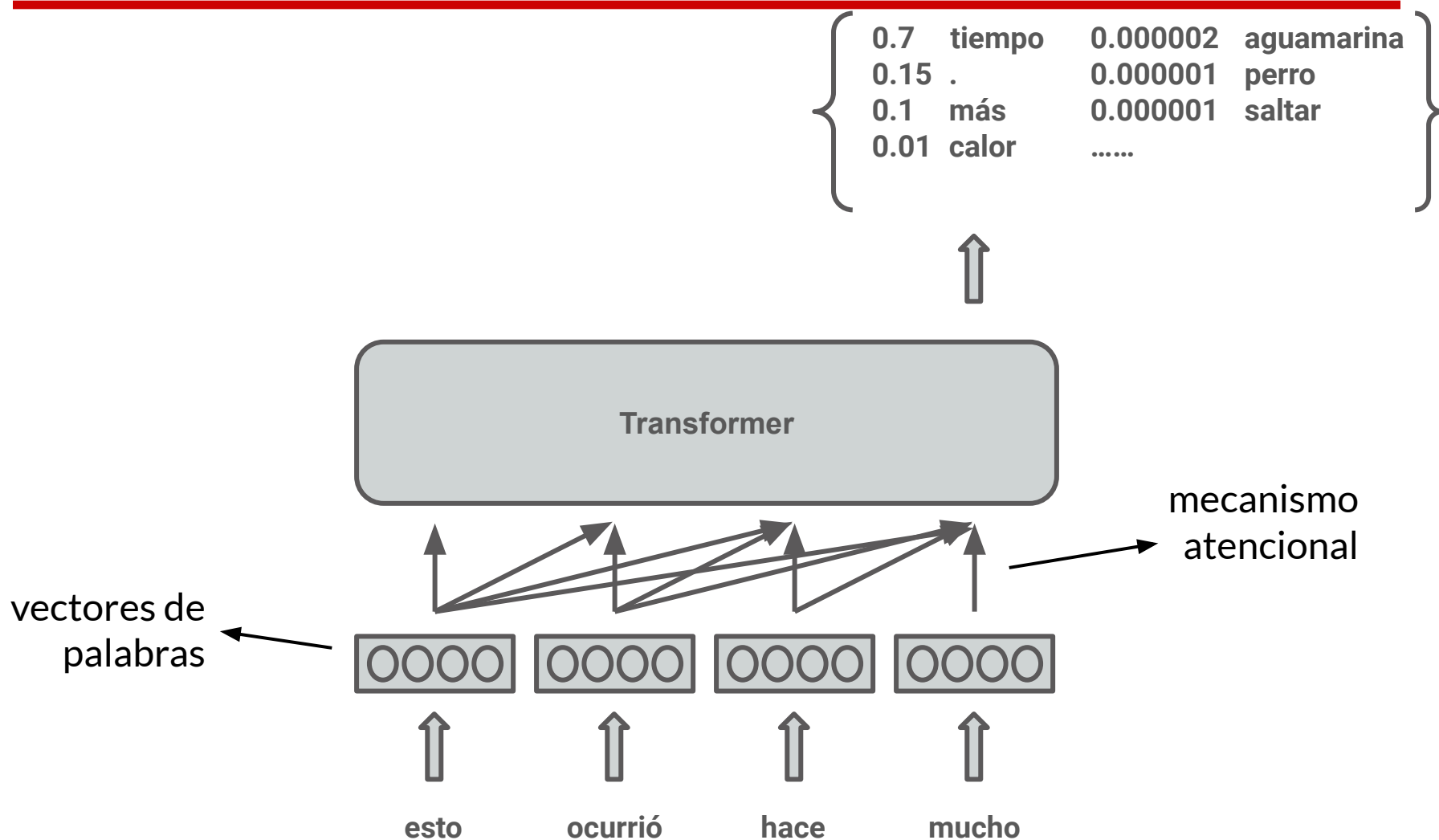


Transformer como Modelo de Lenguaje

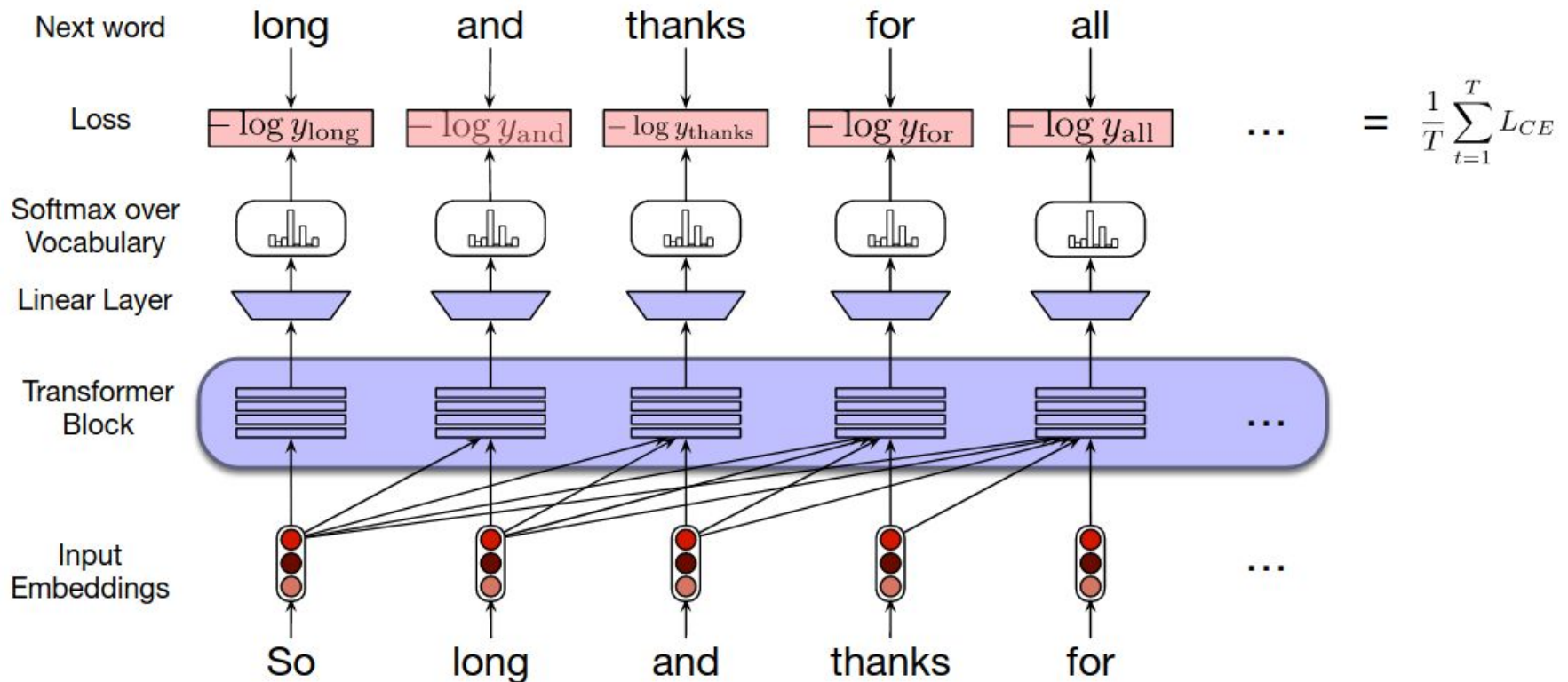


[Vaswani et al. \(2017\)](#)

Transformer como Modelo de Lenguaje



Transformer como Modelo de Lenguaje



Entrenamiento de modelos de lenguaje neuronales

- Autosupervisado: usamos oraciones de un corpus conocido y no necesitamos etiquetas.
 - Se busca minimizar el error de predecir la siguiente palabra dadas todas las anteriores.
 - La red calcula una distribución de probabilidad sobre todo el vocabulario.
-

Modelo de Lenguaje Neuronales

Con una de estas redes podemos predecir la siguiente palabra

...e iterar y predecir toda una continuación de un *prompt*

También calcular la probabilidad de toda una secuencia.

¿Lo podemos usar para otros tipos de tareas?

ML con transformers: generativos / bidireccionales

Dos grandes familias de modelos basados en transformers:

- modelos generativos (GPT, etc.) (fig. a)
- modelos bidireccionales (BERT, etc.) (fig. b)

Jurafsky & Martin 3^a ed.

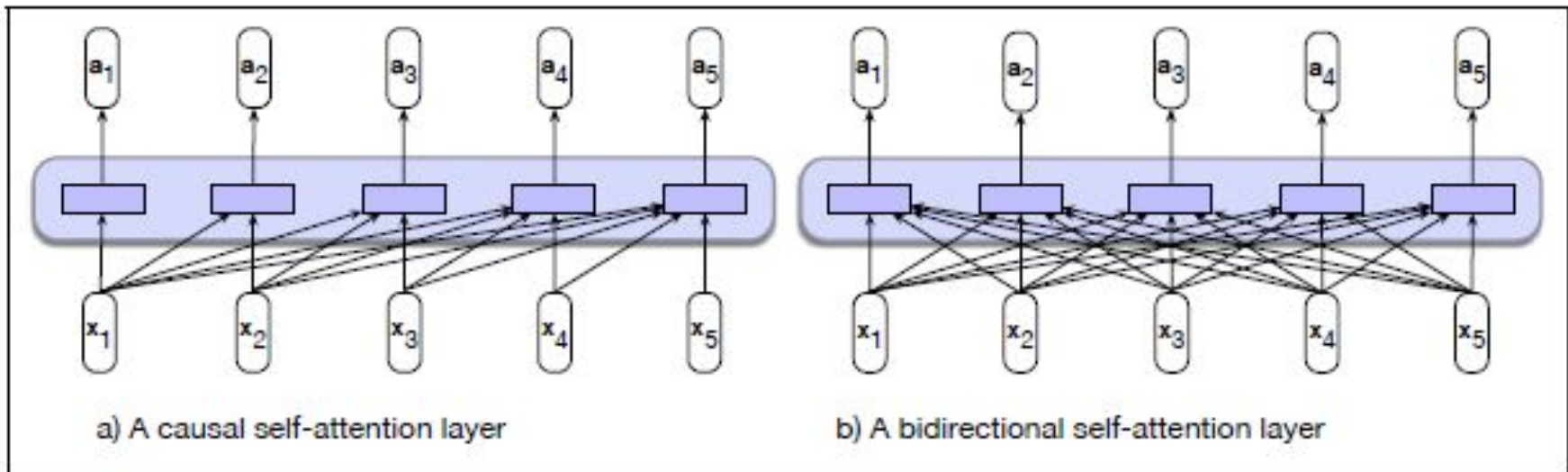


Figure 11.1 (a) The causal, backward looking, transformer model we saw in Chapter 10. Each output is computed independently of the others using only information seen earlier in the context. (b) Information flow in a bidirectional self-attention model. In processing each element of the sequence, the model attends to all inputs, both before and after the current one.

ML con transformers: usos

- Modelos generativos: muy usados para tareas que implican generación de lenguaje, como chatbots, traducción, resumen, etc.
 - Modelos bidireccionales: usados para varias tareas, en particular, para obtener representaciones del texto.
 - Al entrenar estos modelos de lenguaje se generan representaciones internas que son de utilidad para otras tareas de PLN:
 - word embeddings contextuales
 - sentence embeddings
 - Tanto los modelos generativos como los bidireccionales, una vez preentrenados, se usan como punto de partida para nuevas tareas, haciendo *fine tuning* con nuevos datos.
-

ML con transformers: usos

Los grandes modelos de lenguaje contienen miles de millones de parámetros.

Están preentrenados con una gran cantidad de texto, por lo que decimos que ya tienen buen conocimiento del lenguaje.

¿Cómo puedo adaptarlos a mi tarea específica?

- Zero-shot: Le pido que haga algo (generativos)
 - Few-shot: Le pido que haga algo y le doy un par de ejemplos (generativos)
 - Fine-tuning: Si tengo unos cuantos ejemplos más, puedo ajustar los pesos del modelo para mi tarea (ambos)
-

Pretraining y Fine-Tuning

Pretraining

Se preentrena un transformer (ej. para a predecir la siguiente palabra) en un corpus grande.

Fine-tuning

En un modelo preentrendado, se agrega una capa (por ejemplo una feedforward) que es ajustada con un corpus anotado para una tarea específica.

o

Se modifican los pesos del transformer original usando solo los datos de la tarea específica

Fine tuning

Ajustar los pesos en general va a tener mucha mejor performance

- Adaptar a una tarea nueva
- Adaptar a un idioma nuevo
- O a un dominio diferente

¿Problemas?

- Necesito tener los datos etiquetados (en algunos casos)
 - Lleva mucho tiempo adaptar todos los pesos del modelo
 - LoRA (*Low Rank Adaptation*)
 - También ocupa mucha memoria
-

Fine-tuning de ML bidireccionales

pysentimiento

- Entrenamiento de un modelo de lenguaje basado en Roberta (basado a su vez en BERT) con un gran corpus de tweets en español: RroBERTuito (Pérez et al., 2022).
- Fine-tuning para diferentes tareas de análisis de subjetividad, como análisis de sentimiento: pysentimiento (Pérez et al, 2021).
- La mejor configuración alcanza un 70.7% de Macro F para análisis de sentimiento, sobre el dataset de TASS 2020, un valor mayor que los obtenidos en la competencia.

ROUBERTa

- Entrenamiento de un modelo de lenguaje basado en Roberta con un corpus de noticias de prensa uruguaya (Filevich et al, 2024).
 - Fine-tuning para varias tareas: análisis de sentimiento, question & answering..
-

Así lo anunció la titular de la cartera, Karina Rando, este jueves en conferencia de prensa.

ROUBERTa

([https://huggingface.co/f
ilevich/robertita-cased](https://huggingface.co/filevich/robertita-cased))

Mask token: <mask>

Así lo anunció la titular de la <mask>, Karina Rando, este jueves en conferencia de prensa.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.038 s

cartera	0.515
Dinama	0.017
Nación	0.016
app	0.015
organización	0.014

Así lo anunció la titular de la cartera, Karina Rando, este jueves en conferencia de prensa.

BETO (Spanish BERT)
(<https://huggingface.co/dccuchile/bert-base-spanish-wwm-cased>)

Mask token: [MASK]

Así lo anunció la titular de la [MASK], Karina Rando, este jueves en conferencia de prensa.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.053 s

cadena	0.333
cámara	0.065
emisora	0.063
plataforma	0.056
agencia	0.039

Así lo anunció la titular de la cartera, Karina Rando, este jueves en conferencia de prensa.

BERT MULTILINGUAL
(<https://huggingface.co/bert-base-multilingual-cased>)

Mask token: [MASK]

Así lo anunció la titular de la [MASK], Karina Rando, este jueves en conferencia de prensa.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.065 s



Así lo anunció la titular de la cartera, Karina Rando, este jueves en conferencia de prensa.

XLM-ROBERTA
(<https://huggingface.co/xlm-roberta-base>)

Mask token: <mask>

Así lo anunció la titular de la <mask>, Karina Rando, este jueves en conferencia de prensa.

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.123 s

entidad	0.082
AFP	0.048
empresa	0.039
República	0.036
Policía	0.031

Este martes, autoridades del MVOT concurrieron a la comisión de vivienda del Senado para dar explicaciones sobre las entregas discrecionales por parte de la exministra.

	ROUBERTa	BETO	Bert multiling	XML Roberta
Este martes, autoridades del MVOT concurrieron a la <mask> de vivienda del Senado para dar explicaciones sobre las entregas discrecionales por parte de la exministra.	comisión 0.923	comisión 0.321	oficina 0.332	oficina 0.279
Este martes, <mask> del MVOT concurrieron a la comisión de vivienda del Senado para dar explicaciones sobre las entregas discrecionales por parte de la exministra.	autoridades 0.540	miembros 0.377	miembros 0.323	representantes 0.382
Este martes, autoridades del MVOT concurrieron a la comisión de vivienda del Senado para dar explicaciones sobre las <mask> discrecionales por parte de la exministra.	decisiones 0.316	medidas 0.325	prácticas 0.173	medidas 0.295

Luego de que radio Universal informara sobre la adjudicación de una vivienda bajo la modalidad de alquiler con opción a compra a una militante de Cabildo Abierto (CA) sin pasar por sorteo, el presidente de la República, Luis Lacalle Pou, le pidió la renuncia.

	RIUBERTa	BETO	Bert multilingüe	XML Roberta
..., el presidente de la República, Luis Lacalle Pou, le pidió la <mask>.	renuncia 0.464	renuncia 0.181	venta 0.112	renuncia 0.088
Luego de que radio Universal <mask> sobre la adjudicación de una vivienda ...	informara 0.994	##a 0.322	informa 0.143	informó 0.306
Luego de que radio Universal informara sobre la adjudicación de una vivienda bajo la modalidad de alquiler con opción a compra a una <mask> de Cabildo Abierto (CA) ...	persona 0.198 (tercer opción militante)	Corporación 0.344	casa 0.154	persona 0.217

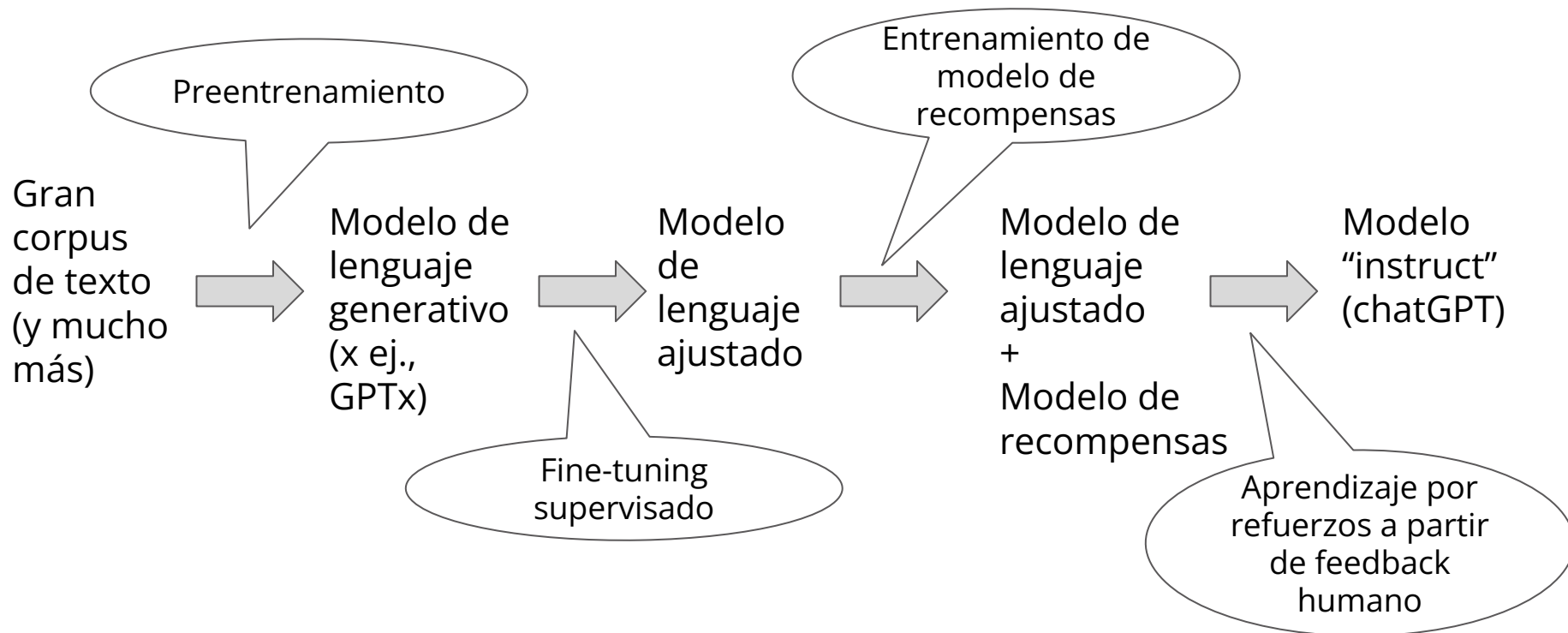
Es decir, el agua que sale por las canillas, sale con gusto salado, al menos en Montevideo y el área metropolitana. (El Observador)

	ROUBERTa	BETO	Bert multilingüe	XML Roberta
Es decir, el agua que sale por las canillas, sale con gusto <mask>, al menos en Montevideo y el área metropolitana.	amargo dulce agradable particular seco	natural también puro propio “	##s gusto popular por ,	, natural . y puro
Es decir, el agua que sale por las canillas, sale con gusto salado, al menos en Montevideo y el área <mask>.	metropoli tana 0.998	metropoli tana 0.829	metropoli tana 0.114	rural 0.312
Es decir, el agua que sale por las canillas, sale con gusto salado, al menos en <mask> y el área metropolitana.	Montevi deo 0.997	México 0.102	Madrid 0.074	Bogotá 0.126

Modelos de lenguaje usados para chatbots

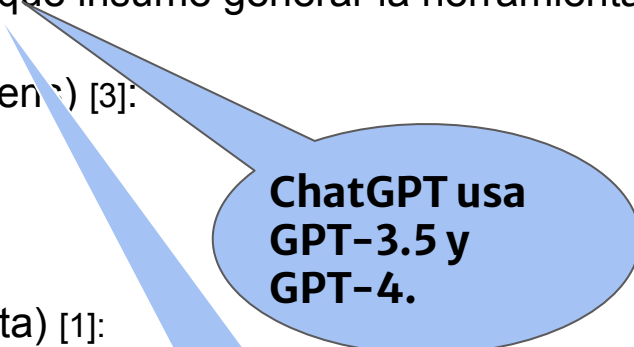
Etapas del entrenamiento de chatbot

[1]

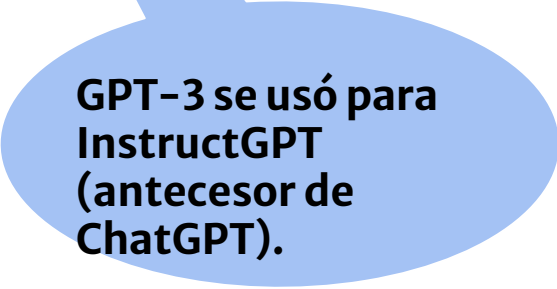


1. Pretraining: generación del modelo de lenguaje

- Consume un 99% del tiempo de entrenamiento que insume generar la herramienta completa.
- Corpus de entrenamiento GPT-3 (499 billion tokens) [3]:
 - 60% filtered CommonCrawl (410 billion tokens)
 - 22% WebText2 (19 billion tokens)
 - 16% Books1 y 2 (67 billion tokens)
 - 3% English Wikipedia (3 billion tokens)
- Corpus de entrenamiento Llama (modelo de meta) [1]:
 - 67% CommonCrawl
 - 15% C4
 - 4.5% Github
 - 4.5% Wikipedia
 - 4.5% Books
 - 2.5% ArXiv
 - 2.0% StackExchange.



**ChatGPT usa
GPT-3.5 y
GPT-4.**



**GPT-3 se usó para
InstructGPT
(antecesor de
ChatGPT).**

1. Pretraining: generación del modelo de lenguaje

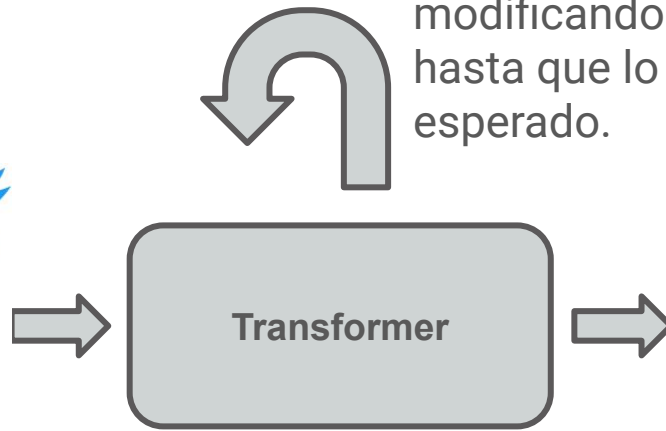
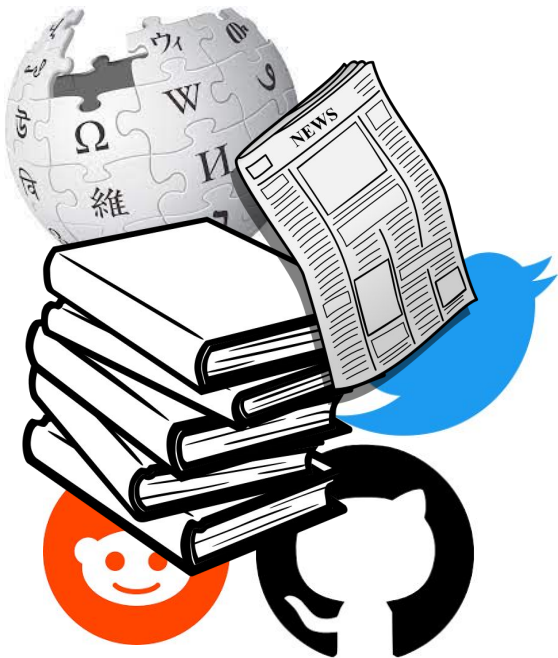
- Datos de entrenamiento de GPT3 [3] y LLama:

	GPT3	LLama
vocabulario	50527	32000
contexto	2048	2048
parámetros	175B	65B
entrenamiento	300B tokens (1 mes)	1-1.4 T tokens (21 días)

actualmente se están usando contextos de ¡¡100.000 tokens!!

En ese momento Llama era más potente que GPT3: la cantidad de parámetros no es tan crucial como el tiempo de entrenamiento

1. Pretraining: generación del modelo de lenguaje



Presentamos los ejemplos muchas veces, modificando los valores internos, hasta que lo predicho se acerque a lo esperado.

1	-1	0.5	0.2	5	-3	3	1	4	...
0.8	1	-1	4	2.5	3	3.1	-2.5	3	...
5.5	4	2	-6.7	-8	6	0	12	1	...
-3	5	2.1	8	-4.2	1	6	7	-2	...
1	-1	0.5	0.2	5	-3	3	1	4	...
0.8	1	-1	4	2.5	3	3.1	-2.5	3	...
-3	5	2.1	8	-4.2	1	6	7	-2	...
5.5	4	2	-6.7	-8	6	0	12	1	...
1	-1	0.5	0.2	5	-3	3	1	4	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

El modelo de lenguaje es el conjunto de parámetros (pesos) que se aprendió.

Recordar: el modelo ya no lo verá más todos los ejemplos que usamos.

2. Supervised fine-tuning (SFT): modelo para generar respuestas a prompts

- El modelo de lenguaje base podría usarse para obtener respuestas a prompts:
 - se le da como entrada un contexto, una secuencia de preguntas y respuestas, y una pregunta.
 - el modelo debe predecir la continuación, que debería ser una respuesta.

Context (passage and previous question/answer pairs)

Tom goes everywhere with Catherine Green, a 54-year-old secretary. He moves around her office at work and goes shopping with her. "Most people don't seem to mind Tom," says Catherine, who thinks he is wonderful. "He's my fourth child," she says. She may think of him and treat him that way as her son. He moves around buying his food, paying his health bills and his taxes, but in fact Tom is a dog.

Catherine and Tom live in Sweden, a country where everyone is expected to lead an orderly life according to rules laid down by the government, which also provides a high level of care for its people. This level of care costs money.

People in Sweden pay taxes on everything, so aren't surprised to find that owning a dog means more taxes. Some people are paying as much as 500 Swedish kronor in taxes a year for the right to keep their dog, which is spent by the government on dog hospitals and sometimes medical treatment for a dog that falls ill. However, most such treatment is expensive, so owners often decide to offer health and even life insurance for their dog.

In Sweden dog owners must pay for any damage their dog does. A Swedish Kennel Club official explains what this means: if your dog runs out on the road and gets hit by a passing car, you, as the owner, have to pay for any damage done to the car, even if your dog has been killed in the accident.

Q: How old is Catherine?

A: 54

Q: where does she live?

A:

2. Supervised fine-tuning (SFT): modelo para generar respuestas a prompts

- Una mejor solución es hacer fine-tuning del modelo base usando un dataset de prompts y respuestas, generado en forma manual por anotadores.
 - Tamaño del dataset: entre 10 y 100K (low quantity / high quality).
 - El modelo SFT ya es un asistente (chat), pero se mejora aún más aplicando RLHF (Reinforcement Learning from Human Feedback), que se compone de un modelo de recompensas y RL.
 - Esto funciona mejor porque es más fácil ordenar salidas ya generadas que generar salidas nuevas.
-

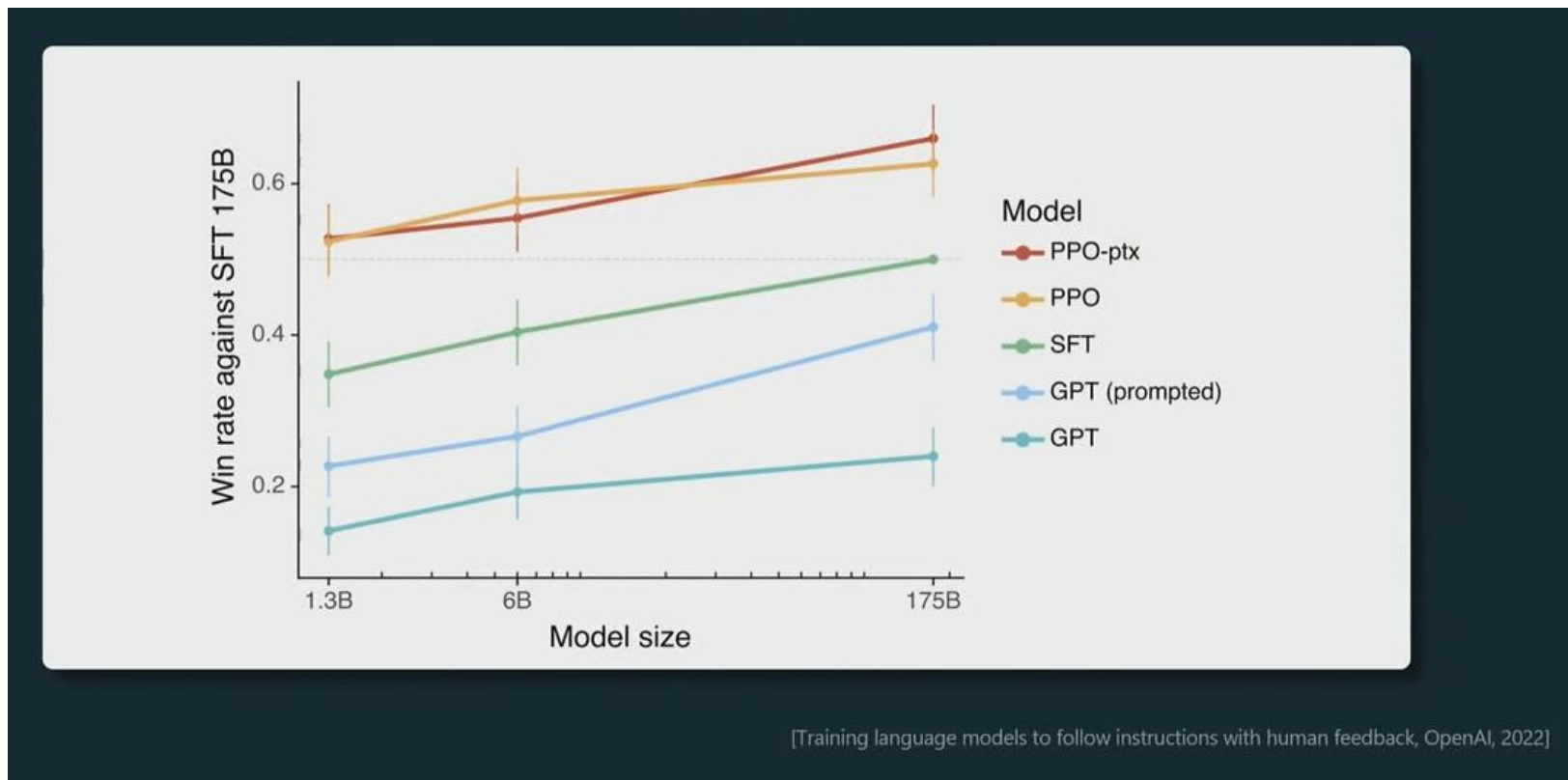
3. Reward Model (RM): se entrena un modelo que rankea respuestas a prompts

- Se entrena un modelo que asigna un peso (recompensa) a cada respuesta que se genera, dado un prompt (Reward Model, RM).
 - Nuevo dataset anotado por humanos:
 - dado un prompt,
 - se generan tres respuestas posibles con el modelo SFT,
 - los humanos las ordenan.
 - El entrenamiento se hace ajustando las recompensas asignadas a cada una de las tres respuestas, según el gold standard.
-

4. Reinforcement Learning (RL): versión final del asistente

- Teniendo un modelo que asigna recompensas es posible entrenar un modelo basado en RL [2]:
 - Se ingresan prompts al modelo SFT.
 - El modelo genera una respuesta.
 - Se asigna una recompensa con el modelo RL.
 - Se ajusta el modelo SFT para maximizar las recompensas.

Comparación GPT, GPT+prompt, SFT y RLHF [1,5]



Algunas limitaciones de los LLM base

- Problema de las “alucinaciones”.
 - Tamaño finito de la ventana de contexto.
 - No es posible:
 - Obtener las fuentes.
 - Acotar respuestas a un dominio específico.
 - Utilizar fuentes externas como base para las respuestas.
 - Y otras más...
-

Técnicas de prompting

Técnicas de prompting

Prompt: Contexto que recibe el modelo que es tokenizado y procesado por el modelo para la continuación.

Las **técnicas de prompting** intentan controlar el lenguaje para cumplir con un formato con un formato

Primer prompt para

Lo que sigue es una conversación entre un estudiante y un docente. Las respuestas del docente deben ser siempre correctas y precisas.

*Estudiante: ¿Qué es un modelo de lenguaje?
Docente:*

Segunda prompt para

Lo que sigue es una conversación entre un estudiante y un docente. Las respuestas del docente deben ser siempre correctas y precisas.

*Estudiante: ¿Qué es un modelo de lenguaje?
Docente: Un modelo de lenguaje asigna una probabilidad a una secuencia de tokens.
Estudiante: ¿Para qué se usan?
Docente:*

Tercera prompt para un chatbot:

Lo que sigue es una conversación entre un estudiante y un docente. Las respuestas del docente deben ser siempre correctas y precisas.

Estudiante: ¿Qué es un modelo de lenguaje?

Docente: Un modelo de lenguaje asigna una probabilidad a una secuencia...

Estudiante: ¿Para qué se usan?

Docente: Los modelos de lenguaje pueden ser utilizados para...

Estudiante: ¿Cuál es la arquitectura de los grandes modelos de lenguaje?

Docente:

Zero-Shot, Few-Shot learning (in-context learning)

existen técnicas
para elegir
buenos ejemplos

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). Language Models are Few-Shot Learners (arXiv:2005.14165). arXiv.

<http://arxiv.org/abs/2005.14165>

Chain of Thought

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✅

RAG: Retrieval Augmentation Generation

- Utilizar una representación vectorial de secuencias de texto.
 - Pueden ser sentence embeddings obtenidas con un modelo de lenguaje tipo BERT.
 - Comparar textos en base a distancias, como la distancia coseno, para encontrar los más similares.
-

RAG: Retrieval Augmentation Generation

Dado un conjunto de textos $\{t_1, t_2, \dots, t_n\}$ que se quieren utilizar como referencias para generar una respuesta a partir de una pregunta q :

- Para cada texto t_i (o fragmento de texto), obtener su representación vectorial $v(t_i)$.
 - Obtener representación vectorial de q , $v(q)$.
 - Utilizando una medida de distancia (por ej. distancia del coseno), elegir los k textos t_i cuya representación $v(t_i)$ esté más “cerca” de $v(q)$ (k-NN).
 - Agregar los k textos seleccionados en la prompt del modelo de lenguaje, indicando explícitamente que se deben utilizar estas referencias para responder a la pregunta.
-

RAG: Retrieval Augmentation Generation

- Ventajas: conozco las fuentes de información.
 - Acierta mucho más si se trata de información fáctica, pero igual se puede equivocar:
 - La información puede simplemente no estar en las fuentes consultadas.
 - No es bueno haciendo razonamientos complejos a partir de los datos existentes.
 - Puede dar respuestas verdaderas pero incompletas.
-

Una aplicación de RAG

Proyecto de grado para responder preguntas sobre reglamentos de la Facultad de Ingeniería (a ser publicado en [IBERAMIA 2024](#)):

- Consultas de estudiantes sobre requisitos de ingreso, previaturas, calidad de libre, ...
-

Una aplicación de RAG

Proyecto de grado para responder preguntas sobre reglamentos de la Facultad de Ingeniería:

- Consultas de estudiantes sobre requisitos de ingreso, previaturas, calidad de libre, ...
 - Respuestas correctas / correctas pero incompletas / incorrectas / no hay respuesta.
-

Una aplicación de RAG

Proyecto de grado para responder preguntas sobre reglamentos de la Facultad de Ingeniería:

- Consultas de estudiantes sobre requisitos de ingreso, previaturas, calidad de libre, ...
 - Respuestas correctas / correctas pero incompletas / incorrectas / no hay respuesta.
 - Conjunto de referencia: preguntas reales de estudiantes de los últimos años, respondidas por integrantes del Espacio de Orientación y Consulta.
 - Problema: las respuestas agregan conocimiento de la persona que responde.
 - ¿Cómo evaluar?
-

Una aplicación de RAG

Proyecto de grado para responder preguntas sobre reglamentos de la Facultad de Ingeniería:

- Consultas de estudiantes sobre requisitos de ingreso, previaturas, calidad de libre, ...
 - Respuestas correctas / correctas pero incompletas / incorrectas / no hay respuesta.
 - Conjunto de referencia: preguntas reales de estudiantes de los últimos años, respondidas por integrantes del Espacio de Orientación y Consulta.
 - Problema: las respuestas agregan conocimiento de la persona que responde.
 - ¿Cómo evaluar?
 - Enfoque:
 - Siempre incluir links a documentos que generaron la respuesta.
 - Mostrar pasos de inferencia aplicados por el modelo.
-

Una aplicación de RAG

Proyecto de grado para responder preguntas sobre reglamentos de la Facultad de Ingeniería:

- Consultas de estudiantes sobre requisitos de ingreso, preiaturas, calidad de libre, ...
 - Respuestas correctas / correctas pero incompletas / incorrectas / no hay respuesta.
 - Conjunto de referencia: preguntas reales de estudiantes de los últimos años, respondidas por integrantes del Espacio de Orientación y Consulta.
 - Problema: las respuestas agregan conocimiento de la persona que responde.
 - ¿Cómo evaluar?
 - Enfoque:
 - Siempre incluir link al documento que generó la respuesta.
 - Mostrar pasos de inferencia aplicados por el modelo (Derivation Prompting)
 - Otros desafíos: recursos computacionales para entrenar/ejecutar localmente grandes modelos. ¿Optimizar? ¿Consultar modelos externos? ¿Modelos pagos?
-

Referencias

Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics*, 3rd edition draft. Stanford. 2024.

[<https://web.stanford.edu/~jurafsky/slp3/>
Acceso: agosto 2024].

Pérez, JM., Furman, D., Alemany, L., Luque, F. RoBERTuito: a pre-trained language model for social media text in Spanish. Proceedings of the Thirteenth Language Resources and Evaluation Conference, 2022.

Pérez, JM., Giudici, JC., Luque, F. pysentimiento: A Python Toolkit for Sentiment Analysis and SocialNLP tasks, arXiv cs.CL, 2021. <https://arxiv.org/abs/2106.09462>

Filevich, J. P., Marco, G., Castro, S., Chiruzzo, L., & Rosá, A. (2024, May). A Language Model Trained on Uruguayan Spanish News Text. In Proceedings of the Second International Workshop Towards Digital Language Equality (TDLE): Focusing on Sustainability@ LREC-COLING 2024 (pp. 53-60).

Recursos sugeridos:

<http://jalammar.github.io/illustrated-transformer/>

<https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a>

Referencias entrenamiento de chatbot

[1] Andrej Karpathy. State of GPT. Mayo 2023.

<https://build.microsoft.com/en-US/sessions/db3f4859-cd30-4445-a0cd-553c3304f8e2>

[2] DCC UChile. Mayo 2023. Cómo entrenar a GPT: Un repaso general de cómo se logró crear y domar a un gran modelo de lenguaje.

https://www.youtube.com/watch?v=4jOY4i0BSSc&t=1362s&ab_channel=DCCUChile

[3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, Dario Amodei. Language Models are Few-Shot Learners. 2020.

<https://arxiv.org/abs/2005.14165>

[4] Yi Liao, Xin Jiang, and Qun Liu. 2020. Probabilistically Masked Language Model Capable of Autoregressive Generation in Arbitrary Word Order. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 263–274, Online. Association for Computational Linguistics.

[5] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L.E., Simens, M., Askell, A., Welinder, P., Christiano, P.F., Leike, J., & Lowe, R.J. (2022). Training language models to follow instructions with human feedback. ArXiv, abs/2203.02155.
