

# LlaMA 2

# URKU

James León  
Proyecto Integrador





# RUNA-SHIMI ÑAN

- 01 EL LENGUAJE
- 02 LOS EMBEDDINGS
- 03 EL ENTRENAMIENTO
- 04 LA EVALUACIÓN
- 05 LOS RESULTADOS
- 06 AHORA QUÉ?



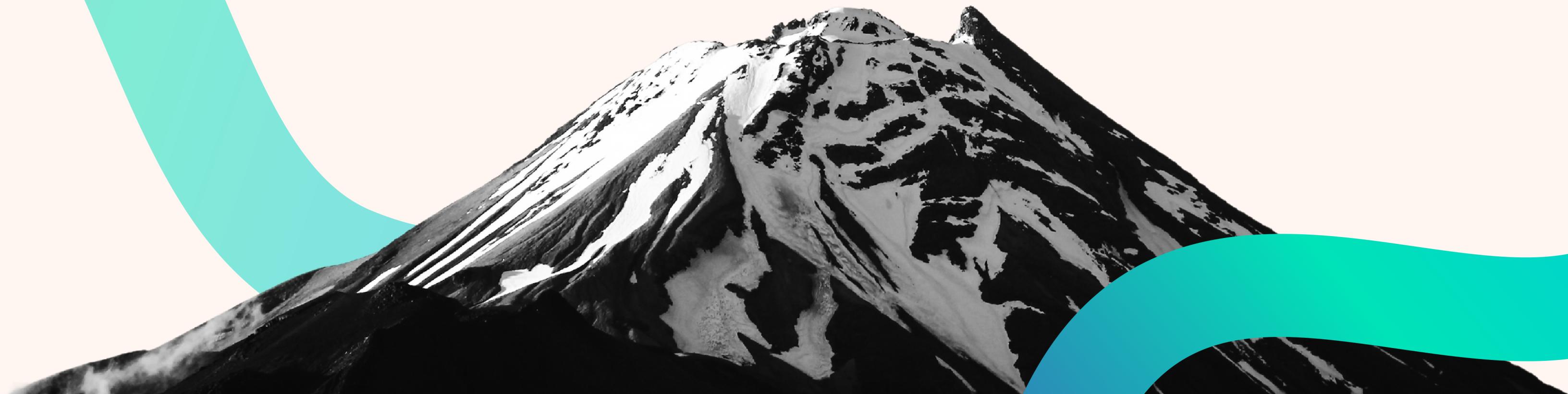
## EL LENGUAJE

El Kichwa como idioma ‘oral’.

El Kichwa **unificado** como línea base para el desarrollo.

Fuentes oficiales en línea: **diccionarios y gramática**.

Fuentes no-oficiales digitalizadas para este proyecto: **transcripts**.



# Los Embeddings

## ✓ BoW and Word2Vec

Frecuencia por palabra y representaciones semánticas basadas en vectores densos

## ✓ FastText

Incluyó información de ‘subwords’ para representación de variantes morfológicas

## ✓ ELMo

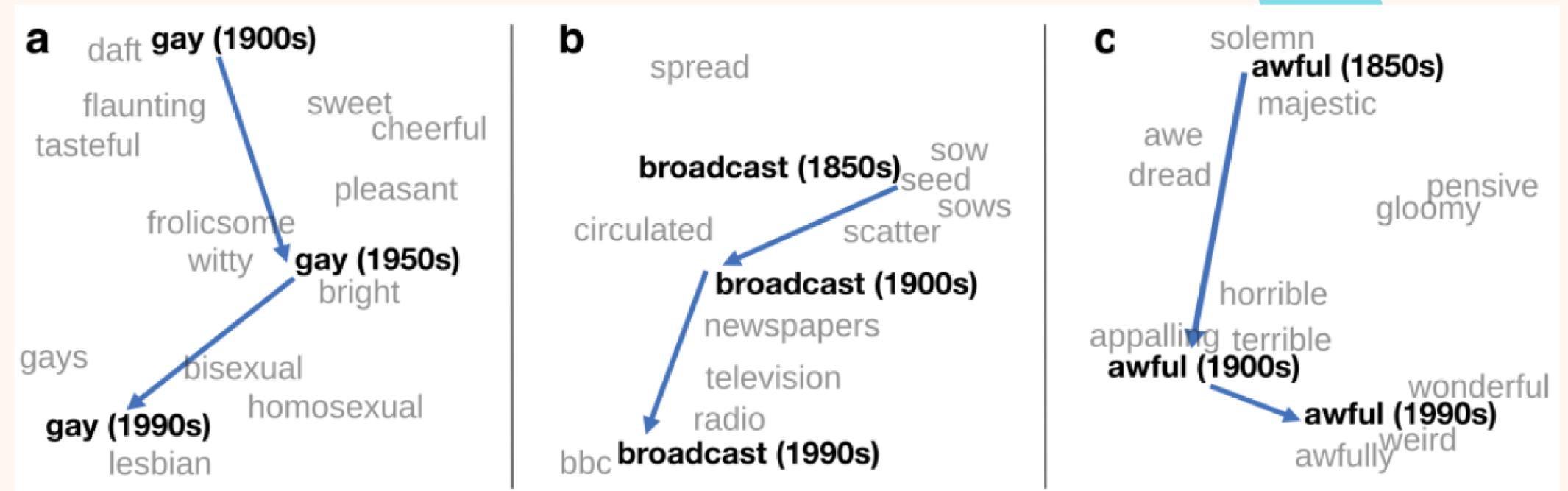
“meaning is use,” embeddings dependientes del contexto

## ✓ BERT

Comprensión contextual bi-direccional

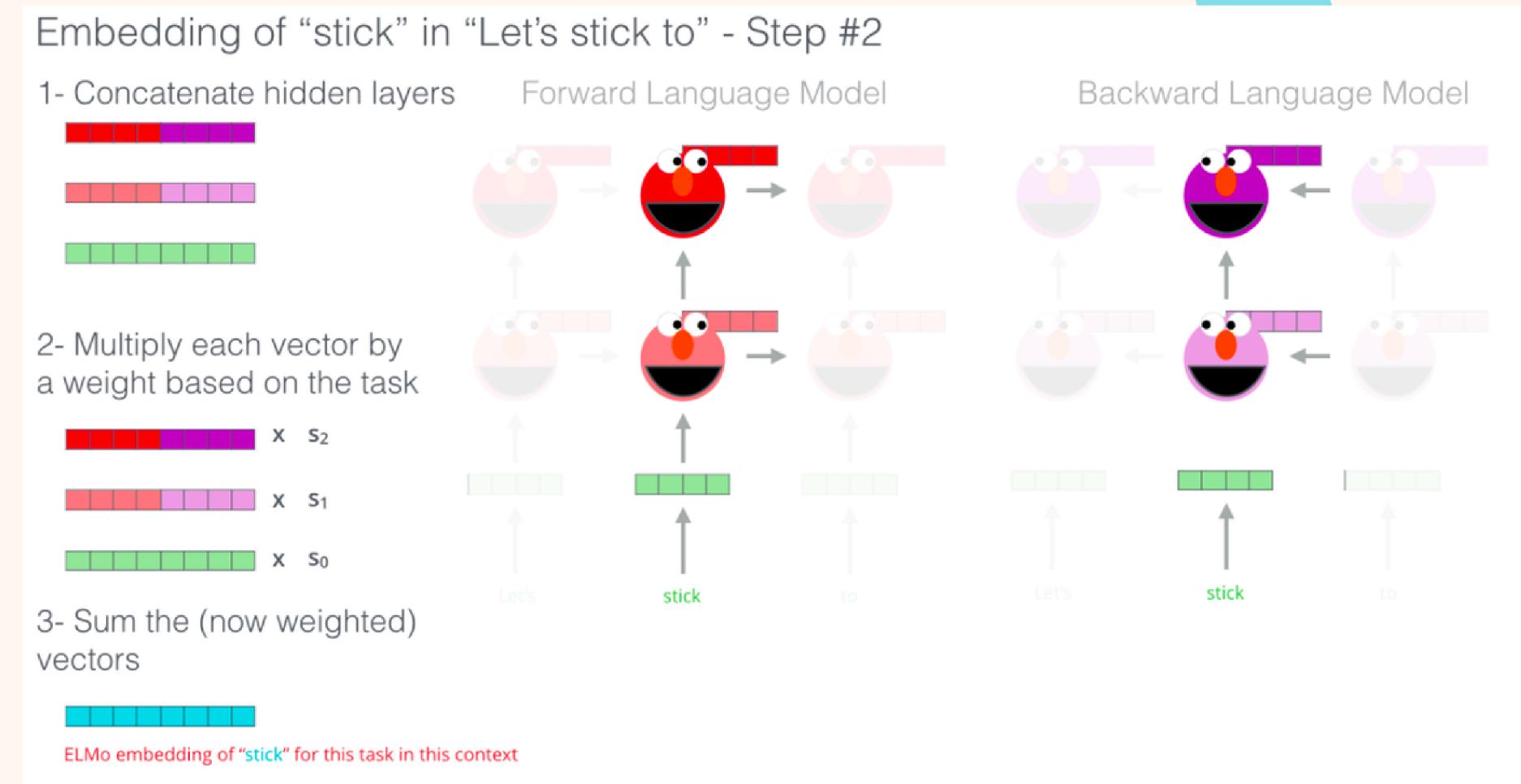
## ✓ LlaMA 2

Modelos basados en transformers entrenados a gran escala, inferencia con una ventana de contexto extendida



# Los Embeddings

Recovered from “BERT” and “Attention is All You Need” papers



Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	# <sup>2</sup> ing	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{#^2ing}$	$E_{[SEP]}$
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

Figure 2: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

Useful Resource

# Los Embeddings

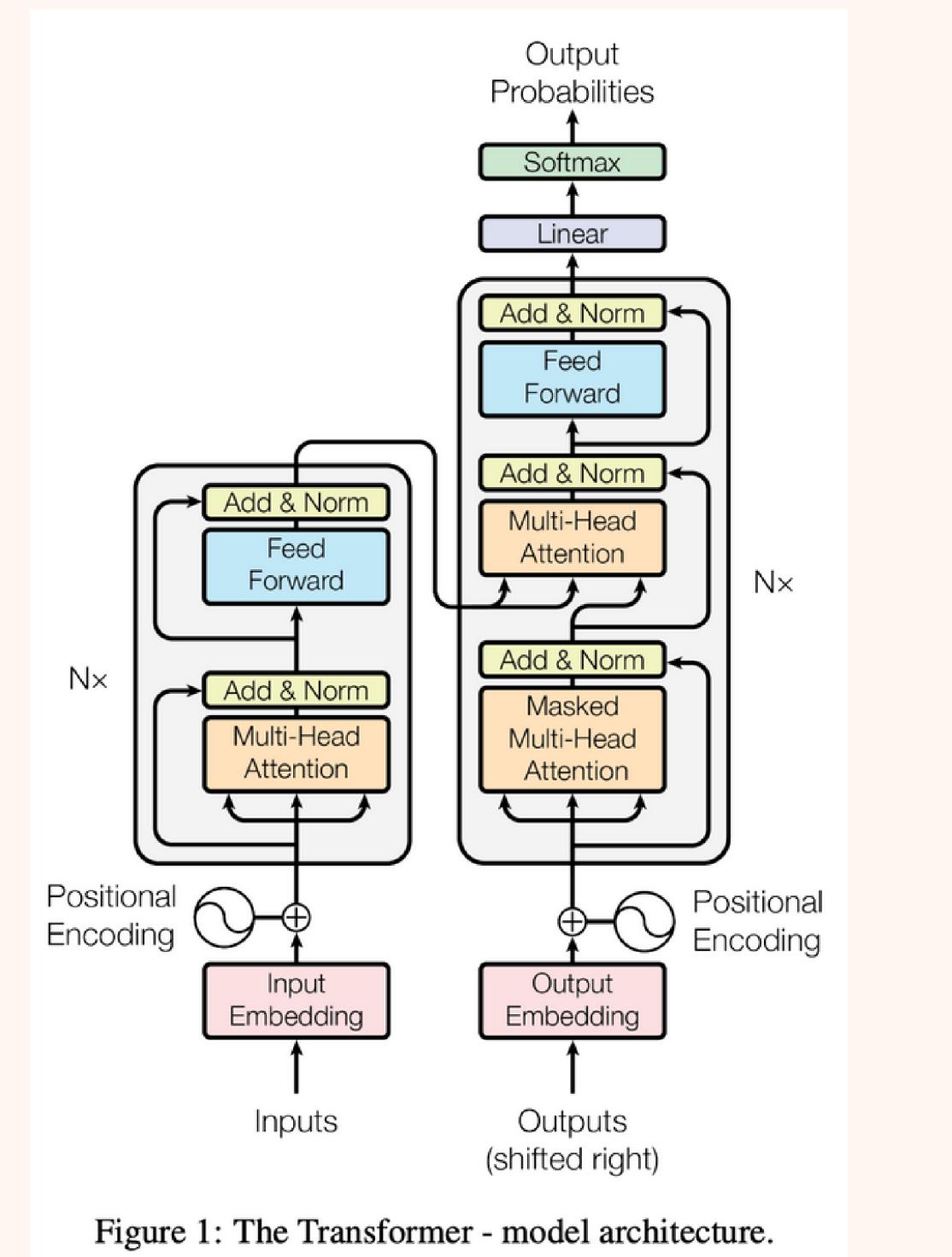
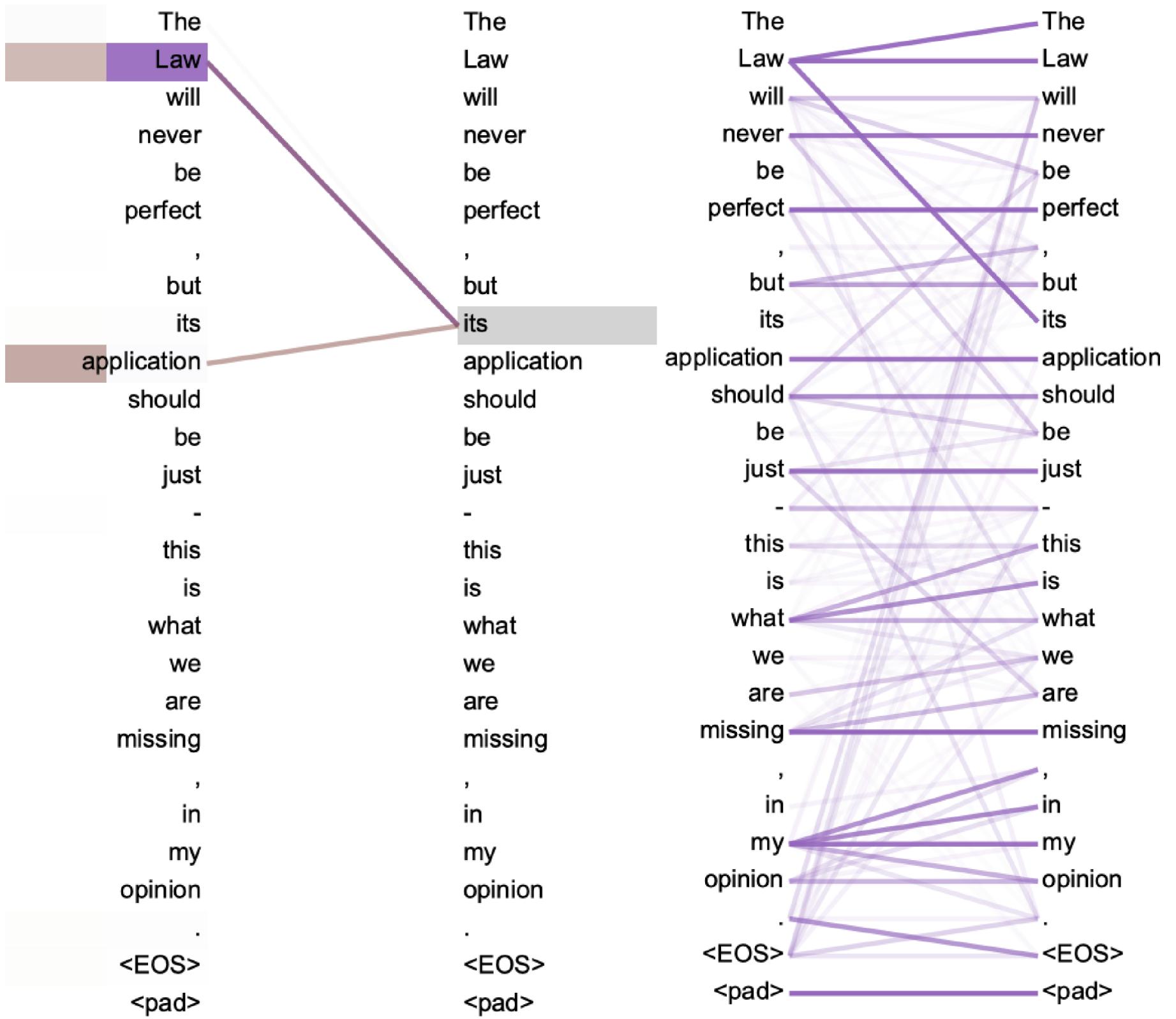


Figure 4: Two attention heads, also in layer 5 of 6, apparently involved in **anaphora resolution**. Top: Full attentions for head 5. Bottom: Isolated attentions from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.



# Fuentes: Diccionarios

## Ministerio de Educación del Ecuador

### AKLLANA

Kunka nanakpika **akirinrita** timpuchishpa up-yarka.

**akllana** [agl<sup>y</sup>ana, agžana, anj<sup>y</sup>ana] v. elegir, escoger, seleccionar. Allita rikushpa hapina. Mamaka kamchankapakka raku saratami **ak-llan**.

**achachay** [ačačay, ačačaw, atsatsay] interj. expresión de frío. Yapa chiri kakpi rimay. Wawakunaka **achachay nishpa chayamun**. Sin. Chiri chiri.

**achachaw** [ačačaw, ačučuy] interj. amz. expresión de calor. Yapa rupay tiyakpi rimay. **achachaw, mikunaka rupakmi kashka**. Sin. Araray, rupakuk.

**achka** [ačka, ačika, aška] adv. bastante, harto, mucho. Imatapash tawkata, tawka tiyakta, mana ashallata rikuchik.

**Chakramantaka achka saratami pallarkani**. Sin. Ashtaka, hatunta, pachan, tawka, llashak.

Asha kunuklla allpapimi **achirataka tarpun**.

**achukcha** [ačukča, ačugča, ačuxča] s. **achogcha**. Purutushina ankushpa pukuk, waylla muru, shunkupi yana muyuyuk, yanushpa mikuna yuyu.

**Wasipika achukcha lukruta mikurkani**.

**achupalla** [ačupal<sup>y</sup>a, ačupaža, ačupil<sup>y</sup>a] s.s. **achupalla**. Kashayuk panka, uchilla charashina, sunilla amuk pankayuk yura.

**Allpa saywapimi achupalla yurataka tarpuni**.

**ahana** [axana] v. ofender, insultar, denigrar. Pitapash piñachishpa rimana.

**Yayaka yankamanta, wawatami ahan**.

Sin. kamina, takurina.

## A



# Resultados iniciales

Page 1 Text:  
Castellano – Kichwa

=====  
Page 2 Text:

A

a cambio de, adv. ranti, rantimpa.  
a continuación, adv. kipa; chaymanta.  
a diario, adv. punchanta punchanta, pun-  
chantin.  
a gusto, adv. ninantak.  
a tiempo, adv. llikchalla, kachka.  
ábaco qichwa (objeto), s. yupana.  
abajo, adv. uray, wayku.  
abandonado, adj. sapalla, sakishka, shi-  
tashka, hichushka.  
abandonar, v. sakina, hichuna, shitana.  
abdómen, s. wiksa.  
abeja, s. wayrunku, chullumpi, putan  
chuspi.  
abismo, s. kaka.  
ablandarse, v. llampuna, apiyana.  
abonanzar, v. kasiyachina.  
abonar, el terreno, s. wanuna.  
abono, s. wanu, isma.  
...  
willka s. rito. 2 adj. sagrado, divino. 3 sust.

230

1554	VERTICAL.–Shayak. <b>VESÍ</b>
1555	CULA BILIAR.–Jayak, chinkilis, ayak muyu. <b>VÍA FÉ</b>
1556	RREA.–Jillayñan.
1557	VIA.–Ñanpi.
1558	VIAJAR.–Purina, rina.
1559	
1560	VIENTO.– Wayra.
1561	VIENTRE.– Wiksa.
1562	VIERNES.–Chaska. <b>VIGÉ</b>
1563	SIMÓ.– Ishkaychunka niki.
1564	VINCHA.–Warmi umawatariy.
1565	VINI.–Wasra.
1566	VIOLAR.– Wakllichina, pakina.
1567	VIOLIN.–Llikilliki.
1568	VIRGEN.– Akllawarmi. Mana wak llishka warmi.
1569	VISITAR.– Pasyana, purina.
1570	VIUDA.–Wakcha, karillak.
1571	VOCABLO.–Shimi.
1572	VOCAL.– Uyari.
1573	VOLAR.–Pawana, wampurina. <b>VOLCÁ</b>
1574	N.–Urku, ninayuk urku.
1575	VOLUNTAD.–Munay, ari nina, yu yaywan nina.
1576	VOMITAR.–Kiwnana, kiknana, quicnana.
1577	VOZ.– Uyachi, shimi.
1578	VUELO.–Paway, huanburina, pa huana. Y y Z z

# Resultados finales

```
The first dictionary has 2066 entries.  
The dictionary entries have 8838 entries.  
The kichwa texts dataset has 2271 entries.  
The Spanish hf dataset has 51942 entries.
```

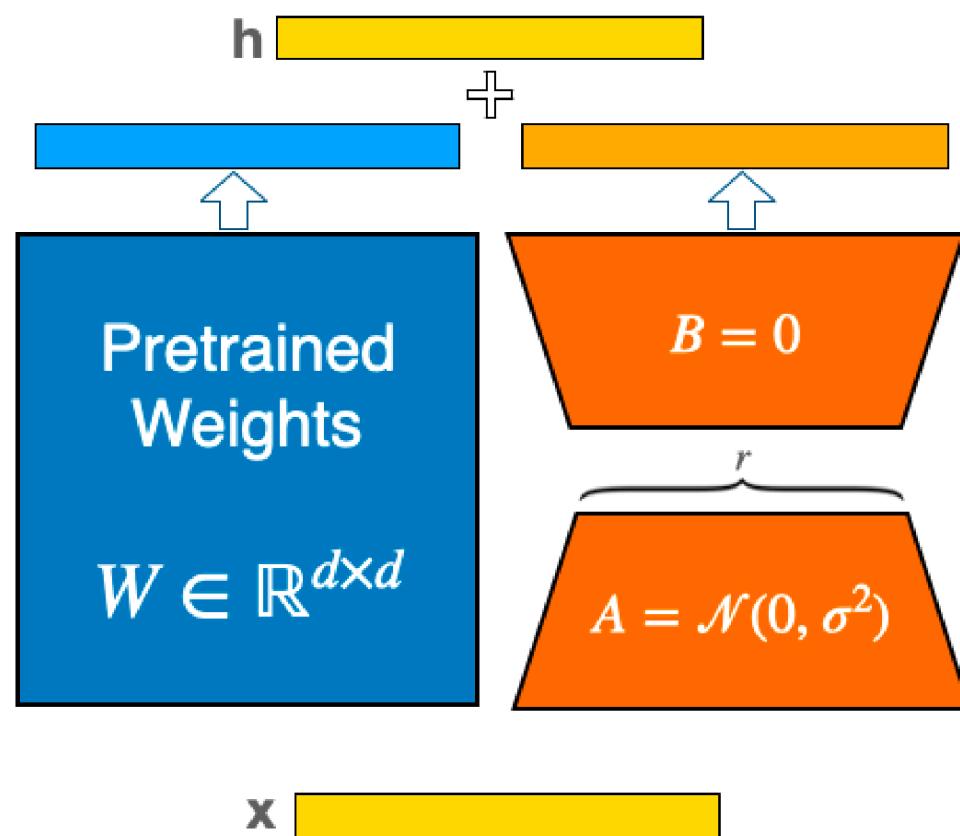
```
1 function BytePairEncoding(text, vocab_size)  
2     vocab = get_initial_vocab(text)  
3     token_freqs = count_token_freqs(text, vocab)  
4     while size(vocab) < vocab_size  
5         most_freq_pair = find_most_freq_pair(token_freqs)  
6         if most_freq_pair is None  
7             break  
8         new_token = merge_pair(most_freq_pair)  
9         vocab.add(new_token)  
10        text = replace_pair_in_text(text, most_freq_pair, new_token)  
11        update_token_freqs(token_freqs, most_freq_pair, new_token)  
12    return text, vocab
```

```
Inputs > JSON > Shuffled_JSONs > {} cas_kich_SHUFFLED.json > ...  
1 [ {  
2     "Instrucción": "¿Cuál es el sinónimo en Kichwa de nuca?",  
3     "Entrada": "nuca",  
4     "Salida": "s. washa kunka."  
5 },  
6 {  
7     "Instrucción": "¿Cómo se traduce la palabra española moverse al Kichwa?",  
8     "Entrada": "moverse",  
9     "Salida": "v. kuyuna; en vaivén: kawirina."  
10 },  
11 {  
12     "Instrucción": "¿Cómo se interpretaría acequia en Kichwa?",  
13     "Entrada": "acequia",  
14     "Salida": "s. larka; hacer acequias: larkana."  
15 },  
16 {  
17     "Instrucción": "¿Cómo traducirías el español entrar a Kichwa?",  
18     "Entrada": "entrar",  
19     "Salida": "v. yaykuna."  
20 },  
21 {  
22     "Instrucción": "¿Cómo se traduciría muslo al Kichwa?",  
23     "Entrada": "muslo",  
24     "Salida": "s. raku chanka, mama chanka."  
25 },  
26 {  
27     "Instrucción": "¿Cómo se traduciría de nuevo al Kichwa?",  
28     "Entrada": "de nuevo",  
29     "Salida": "adv. kutin."  
30 },  
31 {  
32     "Instrucción": "Traduce wampuru al Kichwa.",  
33     "Entrada": "wampuru",  
34     "Salida": "s. calabaza."  
35 },  
36 {  
37     "Instrucción": "¿Cuál es el sinónimo en Kichwa de kachichaska?",  
38     "Entrada": "kachichaska",  
39     "Salida": "s. cosa salada."  
40 },  
41 ]
```

# LoRA training

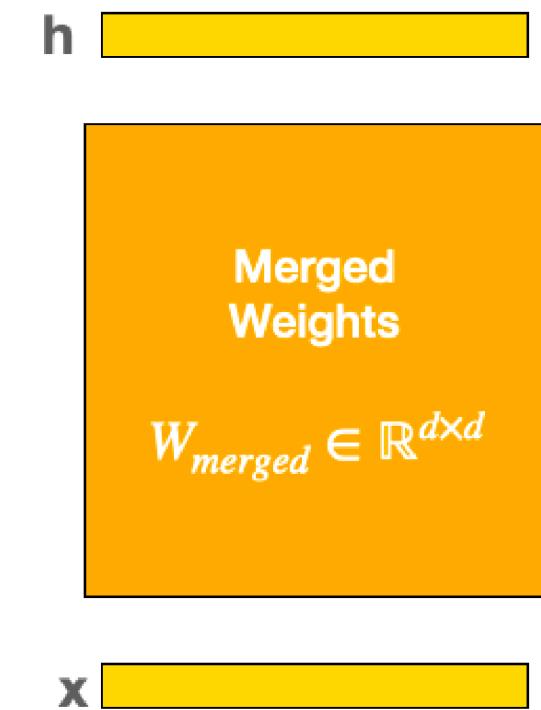
```
1 from peft import LoraConfig  
2  
3 # Target modules - either specific ones or all linear layers  
4 target_modules = ['q_proj', 'k_proj', 'v_proj', 'o_proj']  
5  
6 lora_config = LoraConfig(  
7     r=16, # Rank of the low rank matrices  
8     target_modules=target_modules,  
9     lora_alpha=8,  
10    lora_dropout=0.05,  
11    bias="none",  
12    task_type="CAUSAL_LM"  
13 )
```

During training



$$h = Wx + BAx$$
$$h = \underbrace{(W + BA)x}_{W_{\text{merged}}}$$

After training



# Entrenamiento, Evaluación y Resultados

Modelo	Dataset	Alpha	Rank	Projections	Learning Rate	Steps	Loss	SISA Benchmark*
Llama 2 7B	MoreDetail	4096	2048	q, v, k, o	1,60E-05	3051	2,90E+14	3,3%
Llama 2 7B	MoreDetail	4096	2048	q, v, k, o	1,60E-05	5043	2,90E+12	17,8%
Llama 2 7B	AllData	2048	1024	q, v	2,00E-04	7017	6,5	0,6%
Llama 2 70B	Grammar_guidelines	512	256	q, v,k	9,10E-05	191	2,9	11,7%
Llama 2 70B	Grammar_guidelines	2048	1024	q, v	1,80E-04	431	2,4	0,0%
Llama 2 70B	MoreDetail	2048	1024	q, v	1,70E-04	1863	0,47	6,1%
Llama 2 70B	MoreDetail	2048	1024	q, v	1,10E-04	2503	0,39	88,9%
Llama 2 70B	MoreDetail	2048	1024	q, v	1,00E-05	3727	0,36	66,7%
Llama 2 70B	AllData	1024	512	q, v	1,70E-06	19235	0,18	80,6%
Llama 2 70B	AllData	2048	1024	q, v	3,10E-05	19235	8,6	0,0%
GPT Builder	AllData	n/a	n/a	n/a	n/a	n/a	n/a	95,0%

# ¿Ahora qué?

- Automated Evaluation
- Reinforcement Learning from Human Feedback Modeling
- Dataset further inspection
- Contrastive evaluation with other Parameter-Efficient Fine-tuning Techniques





# YUPAYCHANI

James León  
The URKU Project