

Natural Language Processing

Data Science and Machine Learning
2024

Marcelo Fiinger e Felipe Serras
IME-USP/C4AI





Agenda

1.

What is Natural Language Processing and Computational Linguistics?

2.

Why we need Natural Language Processing and Computational Linguistics

3.

How we perform NLP and Computational Linguistics?

4.

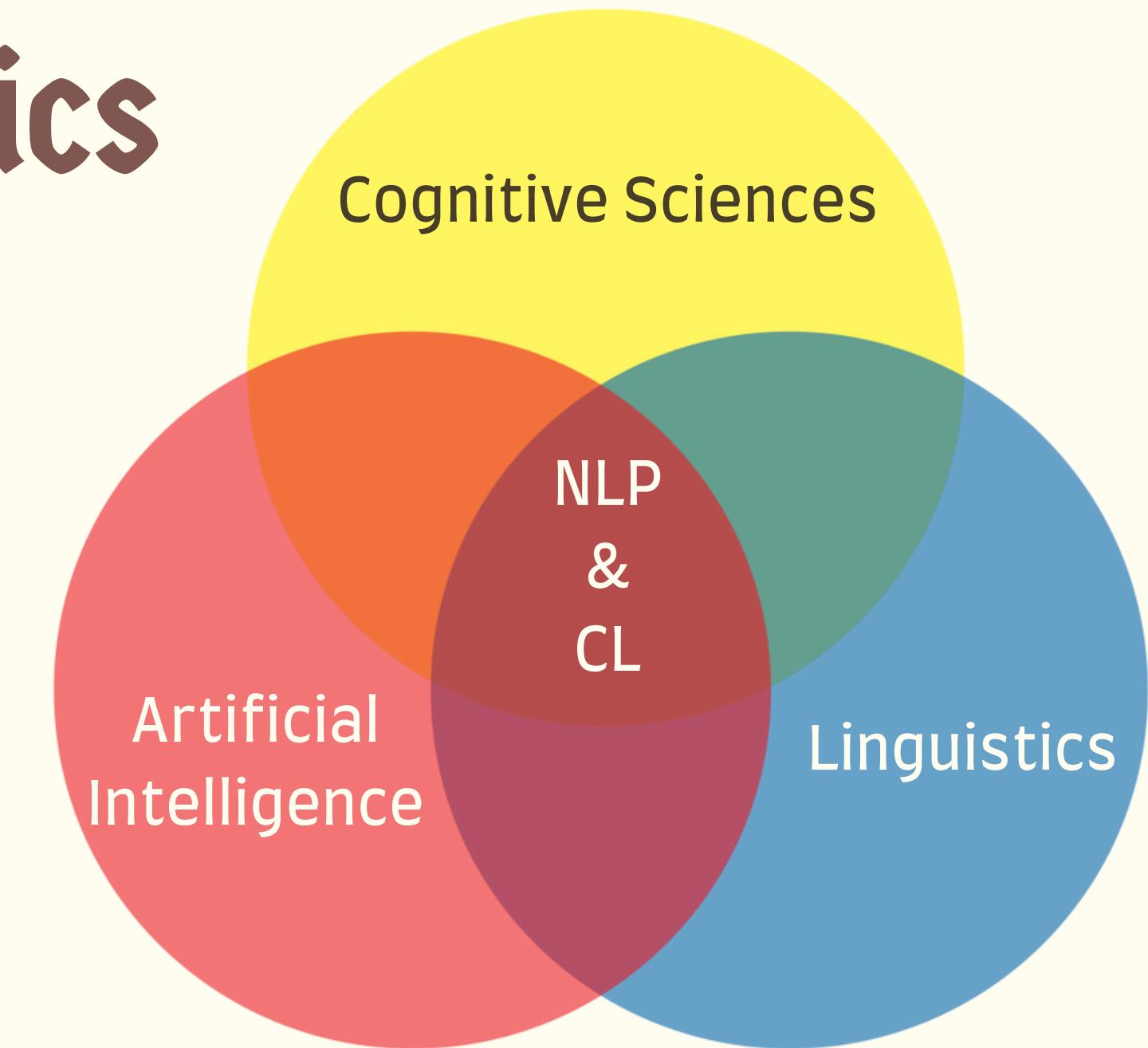
Which models are state-of-the-art in NLP and Computational Linguistics?

1.What is Natural Language Processing and Computational Linguistics?



What is Natural Language Processing and Computational Linguistics

- It is a field of scientific and technological research;
- How can computational models be used to process natural language data and better understand the functioning of natural language?
- It is a multidisciplinary field;
- It originated during the Cold War from an attempt to create computer programs to translate texts from Russian to English.





What is the difference between Natural Language Processing and Computational Linguistics

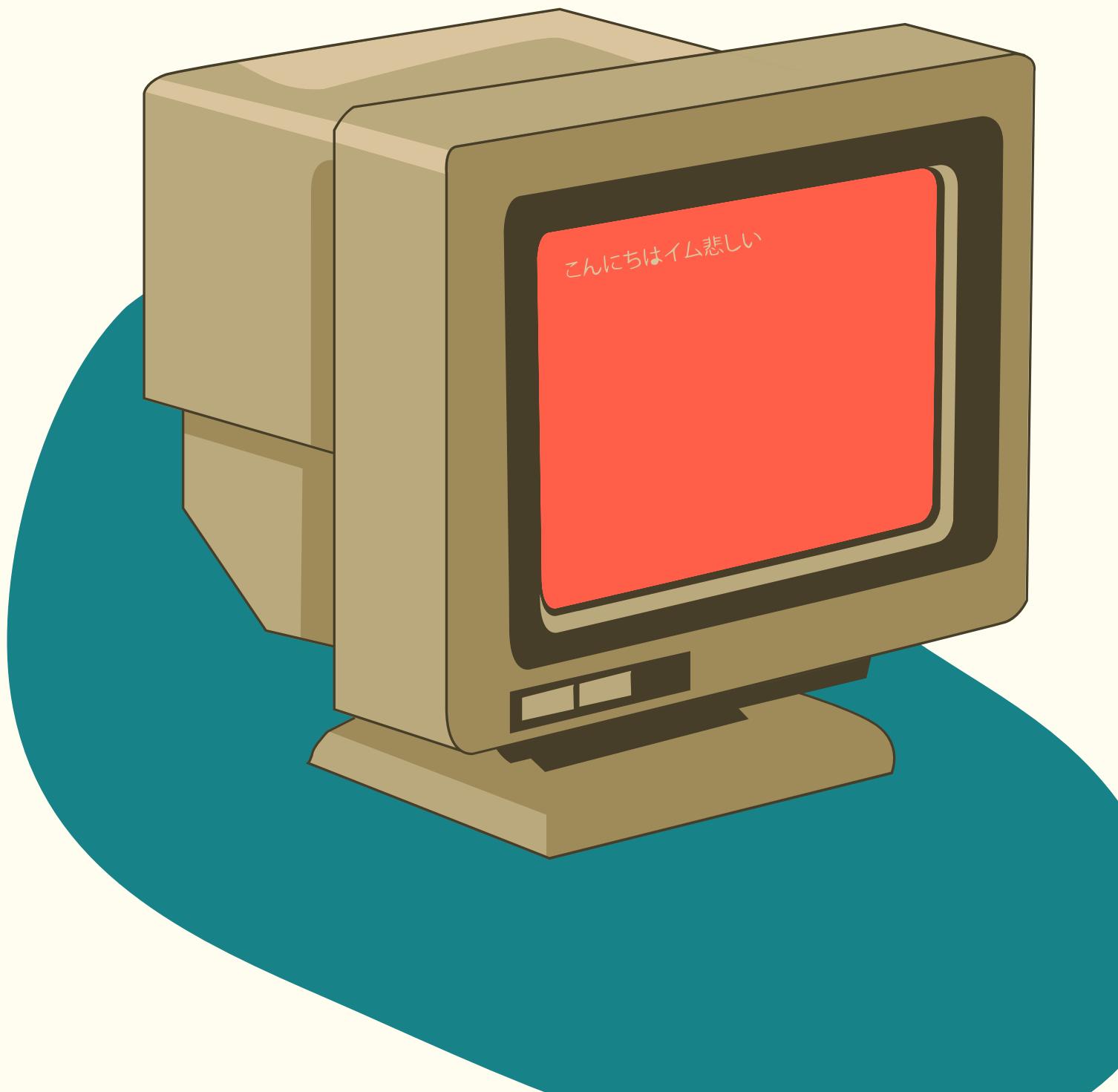
Computational Linguistics is focused on the investigation of human languages and how they function using computational resources.

NLP is focused on the development of computational resources for the accomplishment of tasks using data in human language

Does it matter?

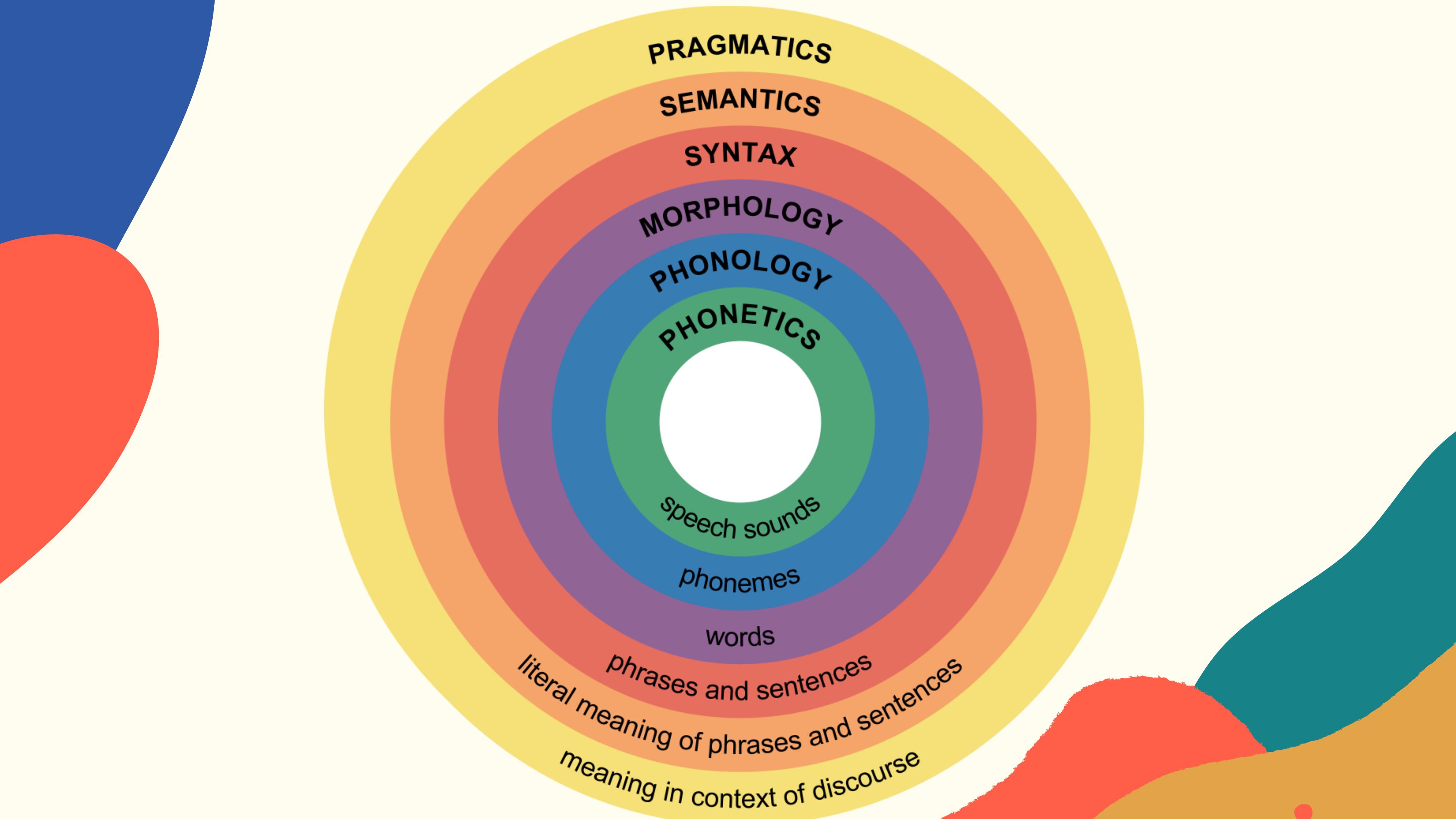
Tasks in Natural Language Processing

- **Translation**
- **Classification:** Sentiment Analysis, Spam Detection, Topic Classification
- **Regression:** Autograding
- **Clustering:** Topic Modeling, Authorship Attribution, Similarity-based Recommendations
- **Tagging:** Named Entity Recognition (NER), Part-of-Speech Tagging
- **Generation:** Conversational Agents, Code Generation



2. Why do we need natural Language Processing and Computational Linguistics?





Phonetics and Phonology



THE INTERNATIONAL PHONETIC ALPHABET (revised to 2015)

CONSONANTS (PULMONIC)

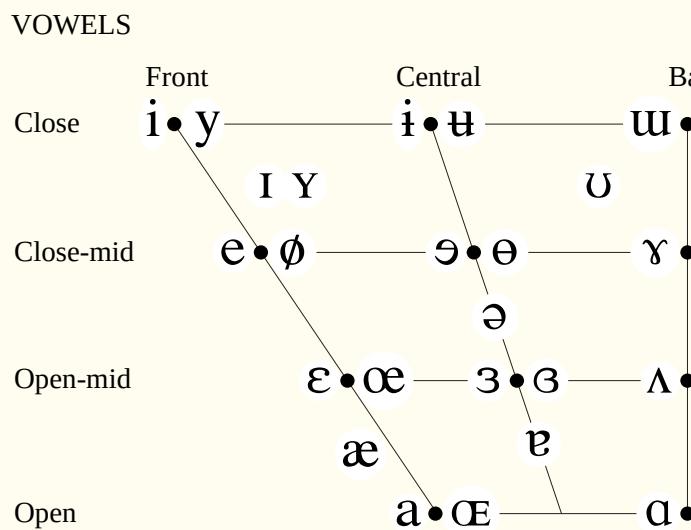
	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		t̪ d̪	c j	k g	q G		?
Nasal	m	m̪		n		n̪	n̥	n̪	N		
Trill	B			r					R		
Tap or Flap		v̆		t̆		t̆					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	s̪ z̪	ç j	x y	χ ʁ	h ʕ	h
Lateral fricative				ɬ ɭ							
Approximant		v̞		ɹ		ɻ	j	w̞			
Lateral approximant				ɬ		ɬ	ɻ	ɬ			

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

CONSONANTS (NON-PULMONIC)

Clicks	Voiced implosives	Ejectives
ʘ Bilabial	b Bilabial	' Examples:
Dental	d Dental/alveolar	p' Bilabial
! (Post)alveolar	f Palatal	t' Dental/alveolar
ǂ Palatoalveolar	g Velar	k' Velar
Alveolar lateral	g' Uvular	s' Alveolar fricative

OTHER SYMBOLS



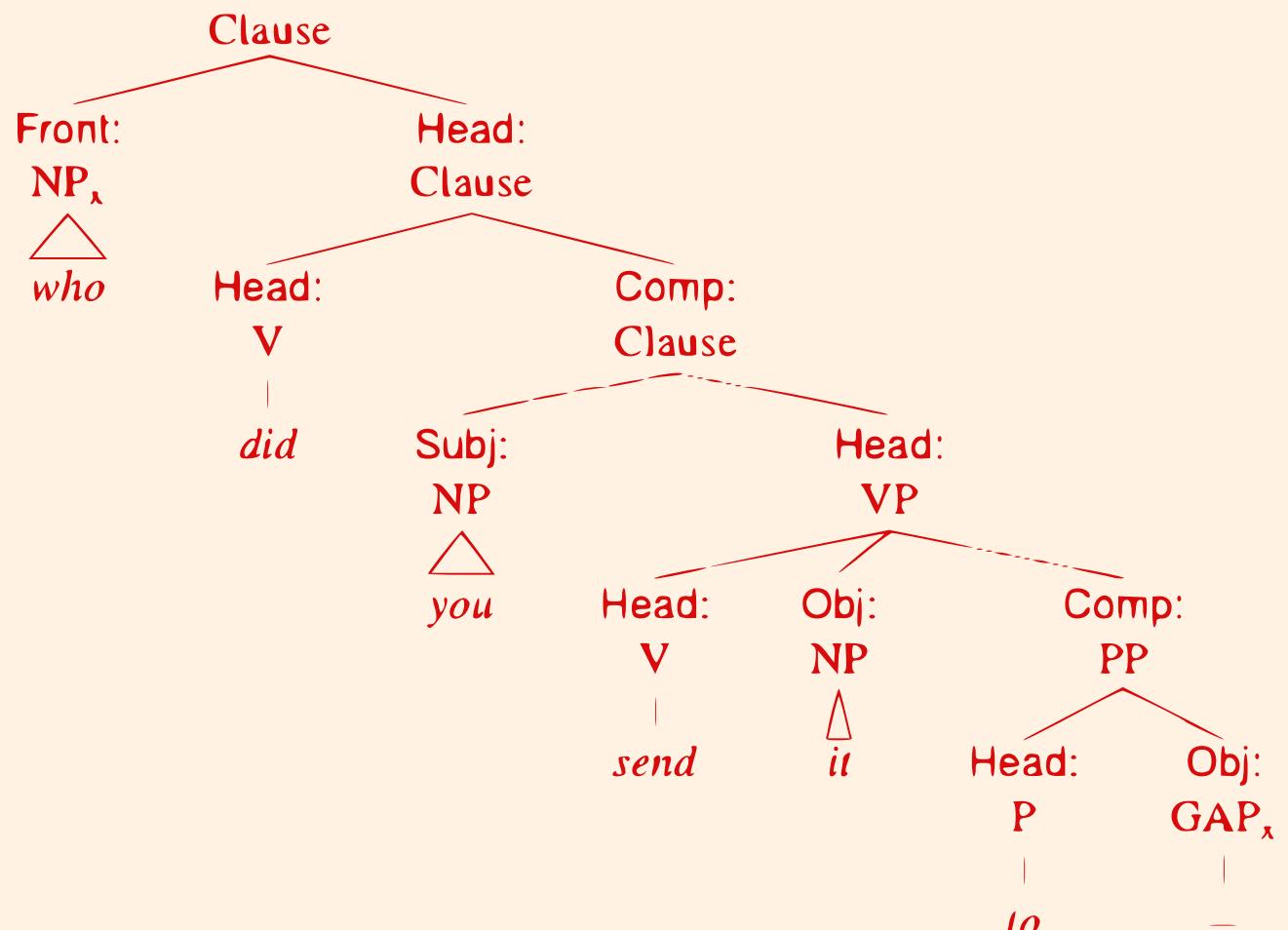
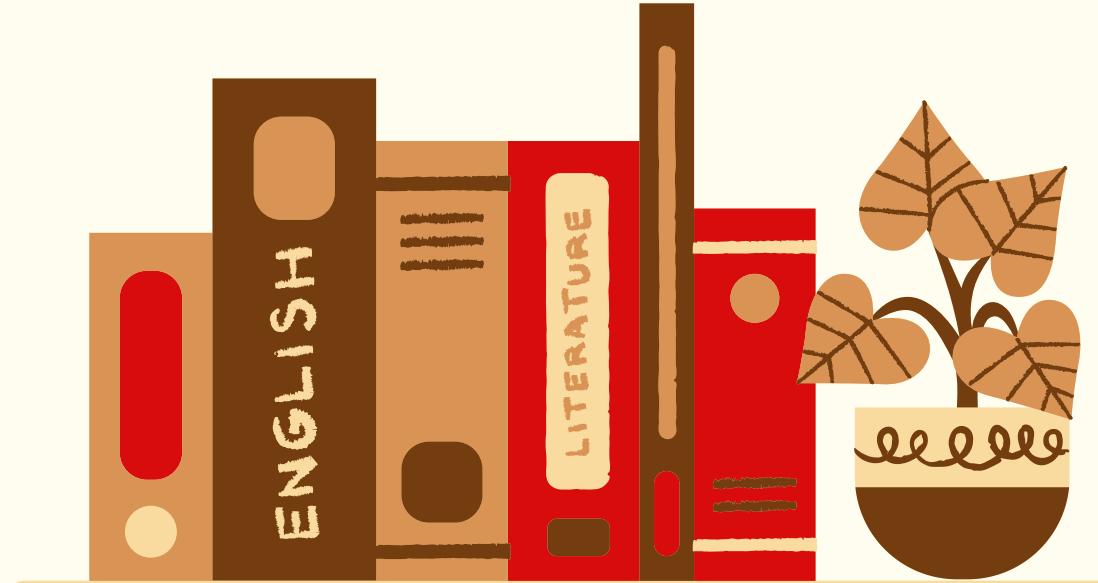
OTHER SYMBOLS

M	Voiceless labial-velar fricative	C Z	Alveolo-palatal fricatives
W	Voiced labial-velar approximant	J	Voiced alveolar lateral flap
ɥ	Voiced labial-palatal approximant	ħ	Simultaneous ʃ and X
H	Voiceless epiglottal fricative		
ʕ	Voiced epiglottal fricative		Affricates and double articulations can be represented by two symbols joined by a tie bar if necessary.
ʢ	Epiglottal plosive		

DIACRITICS Some diacritics may be placed above a symbol with a descender, e.g. ñ

o	Voiceless	n̥ d̥	..	Breathy voiced	b̥ ə̥	▫	Dental	t̥ d̥
✓	Voiced	s̥ t̥	~	Creaky voiced	b̥ ə̥	▫	Apical	t̥ d̥
h	Aspirated	tʰ dʰ	~	Linguolabial	t̥̩ d̥̩	▫	Laminal	t̥̩ d̥̩
,	More rounded	ɔ̥	W	Labialized	tʷ dʷ	~	Nasalized	ɛ̥
,	Less rounded	ɔ̥	j	Palatalized	tj̥ dj̥	n̥	Nasal release	d̥n̥
+	Advanced	u̥	Y	Velarized	tʸ̥ dʸ̥	l̥	Lateral release	d̥l̥
-	Retracted	e̥	᷇	Pharyngealized	t̥᷇ d̥᷇	᷈	No audible release	d̥᷈
..	Centralized	ë̥	~	Velarized or pharyngealized	†			
×	Mid-centralized	ë̥	+	Raised	ɛ̥ (j̥ = voiced alveolar fricative)			
,	Syllabic	n̥	-	Lowered	ɛ̥ (β̥ = voiced bilabial approximant)			
,	Non-syllabic	ɛ̥	-	Advanced Tongue Root	ɛ̥			
~	Rhoticity	ə̥ ḁ	-	Retracted Tongue Root	ɛ̥			

Syntax



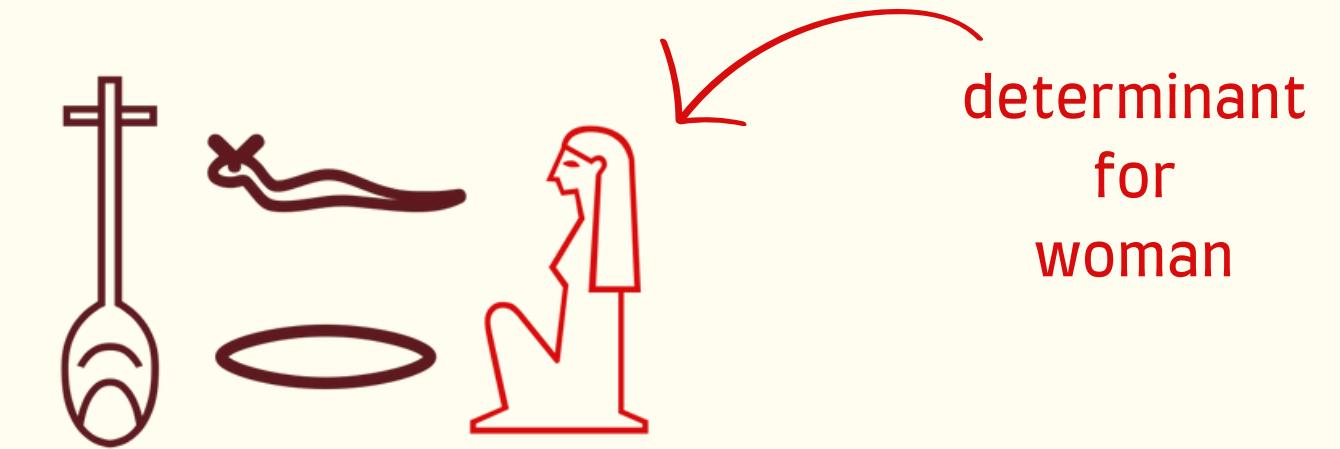
NOUN	ADJECTIVE	DETERMINER
VERB	ADVERB	CONJUNCTION
NUMERAL	PRONOUN	INTERJECTION
ADPOSITION	PARTICLE	...

Morphology

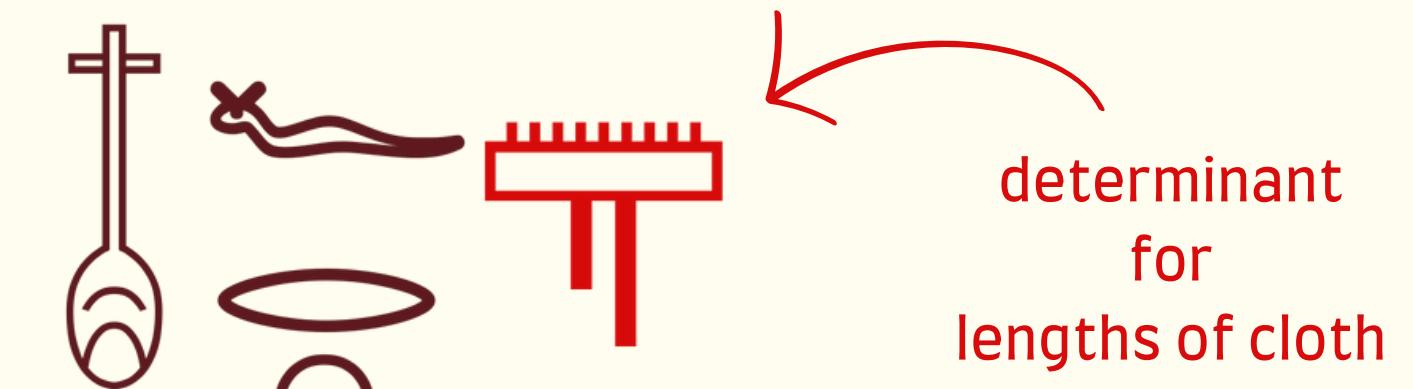
- Derivation
 - break → breakable → unbreakable
(English)
- Inflection
 - Portuguese:
 - Eu falo, tu falas, nós falamos
 - Spanish:
 - Yo hablo, tu hablas, nosotros
hablamos.



- Egyptian:



young woman of marriageable
age



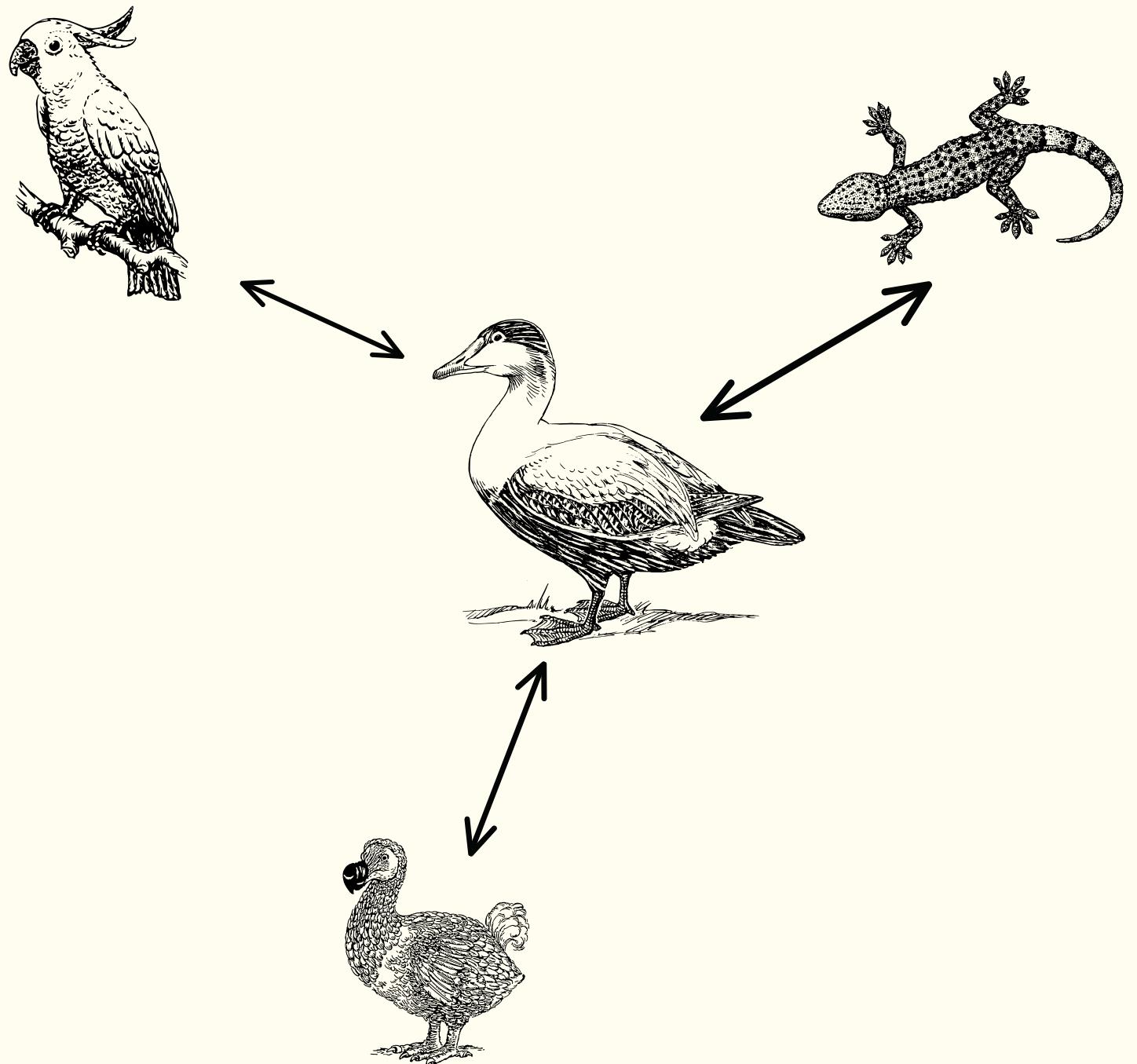
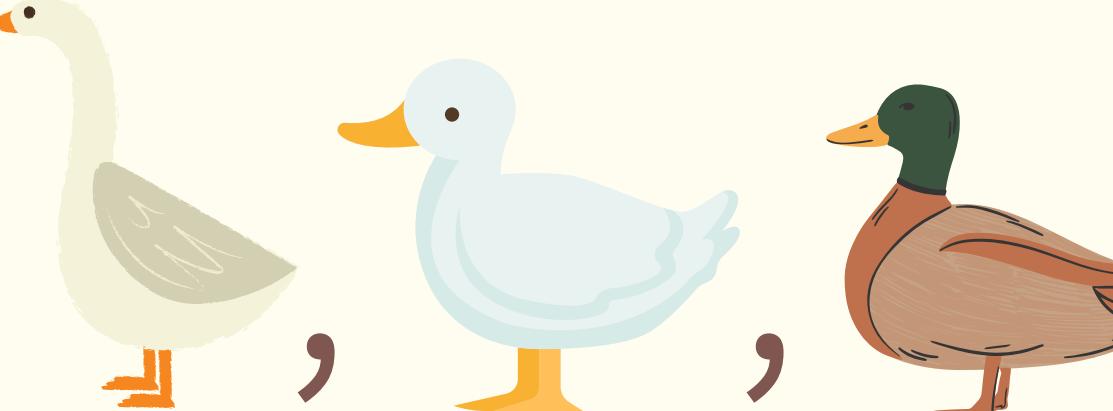
clothing

Semantics

- What is the meaning of “meaning”?
- What is the meaning of “duck”?

Is it the set of all possible ducks?

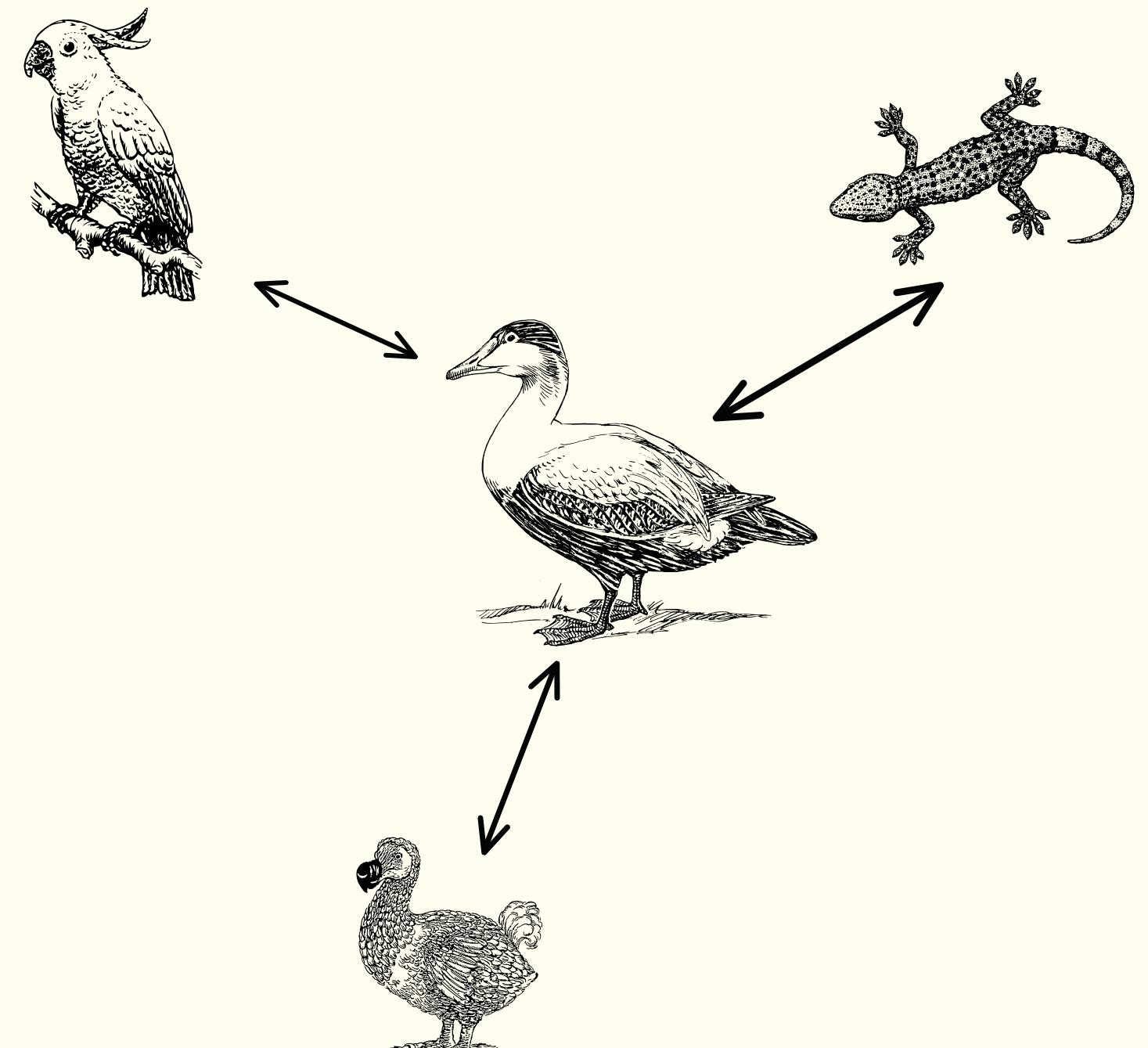
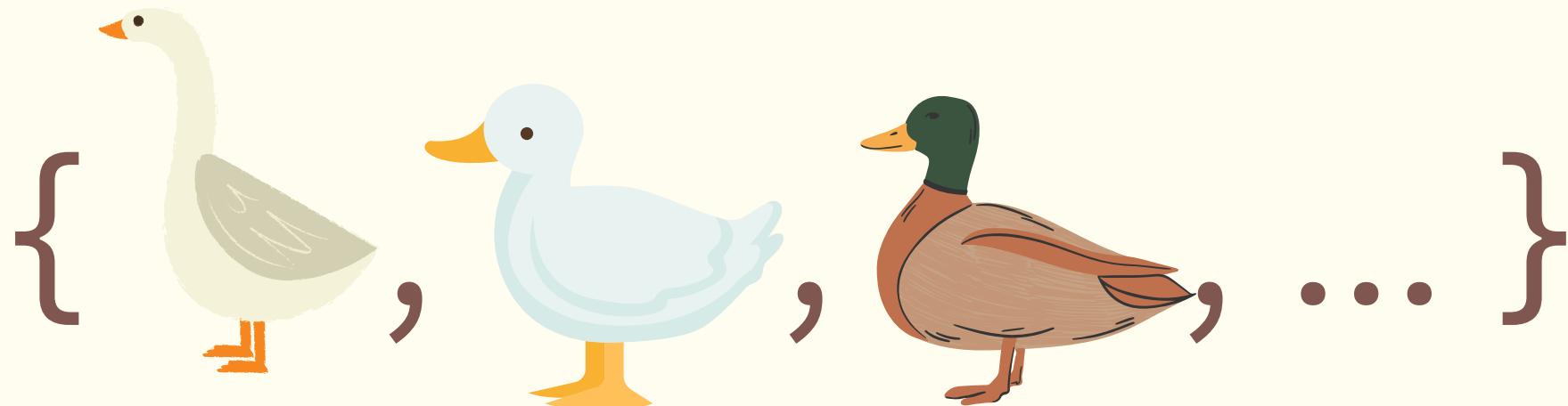
{
 , , ... }



Is it the distance to a prototypical duck?

Semantics

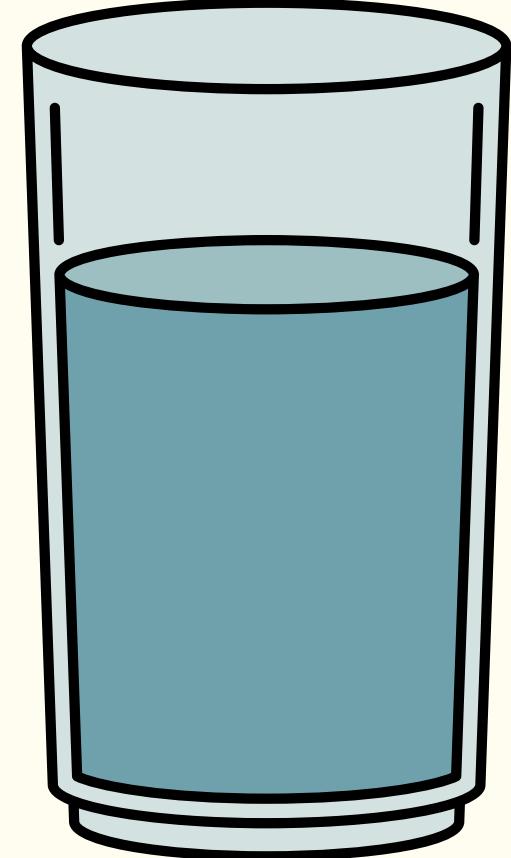
- What is the meaning of “duck”?
 - Obviously, a duck is a hypothesized, fuzzy prototype, updatable and represented as a point in an n-dimensional space



Is it the distance to a prototypical duck?

Pragmatics

- How meaning changes with context?
 - Irony, Implicature
- What is the relationship between meaning and context?
 - Distributional Hypothesis



Can you pass me
the water?

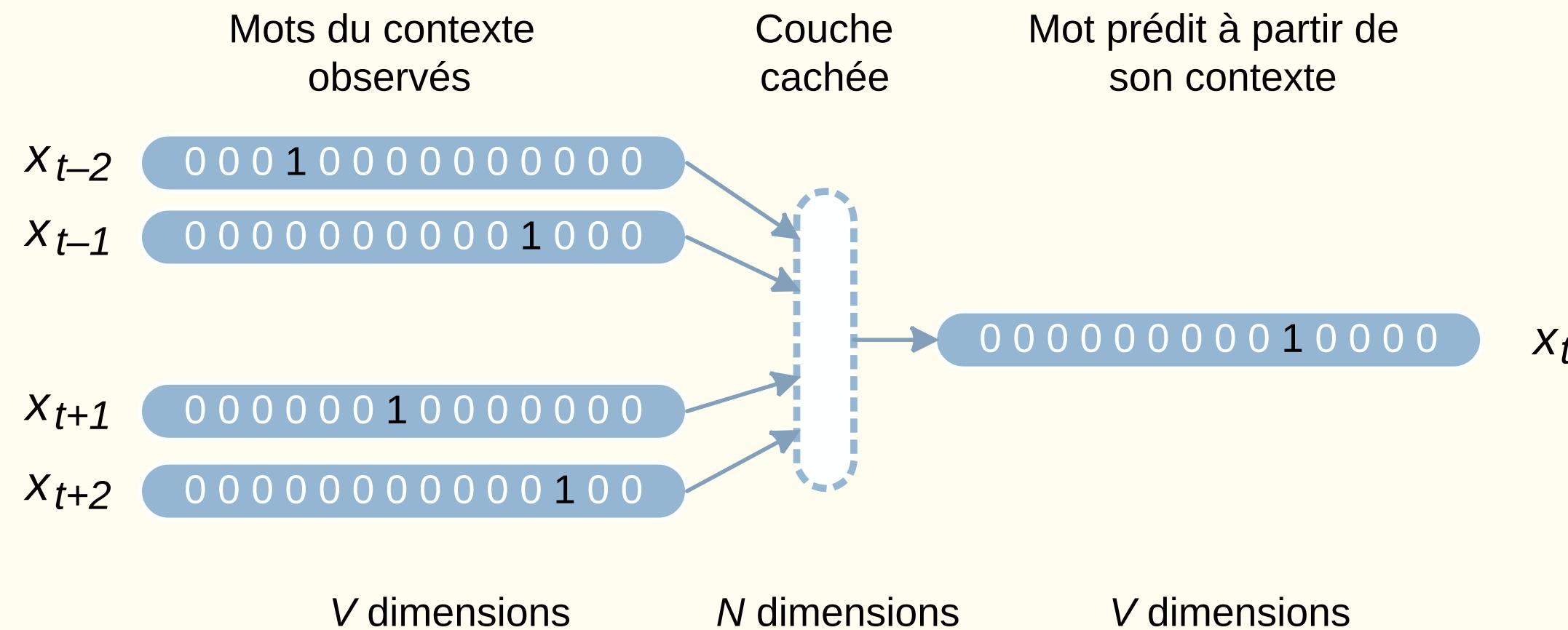
I am fully capable of passing the
glass of water. Thanks for
asking!

3. How do we perform Natural Language Processing and Computational Linguistics?



The Problem of Representation

- How to represent data in human language in a format that allows us to perform mathematical operations?
- Vector Semantics → Embeddings



Several Paradigms

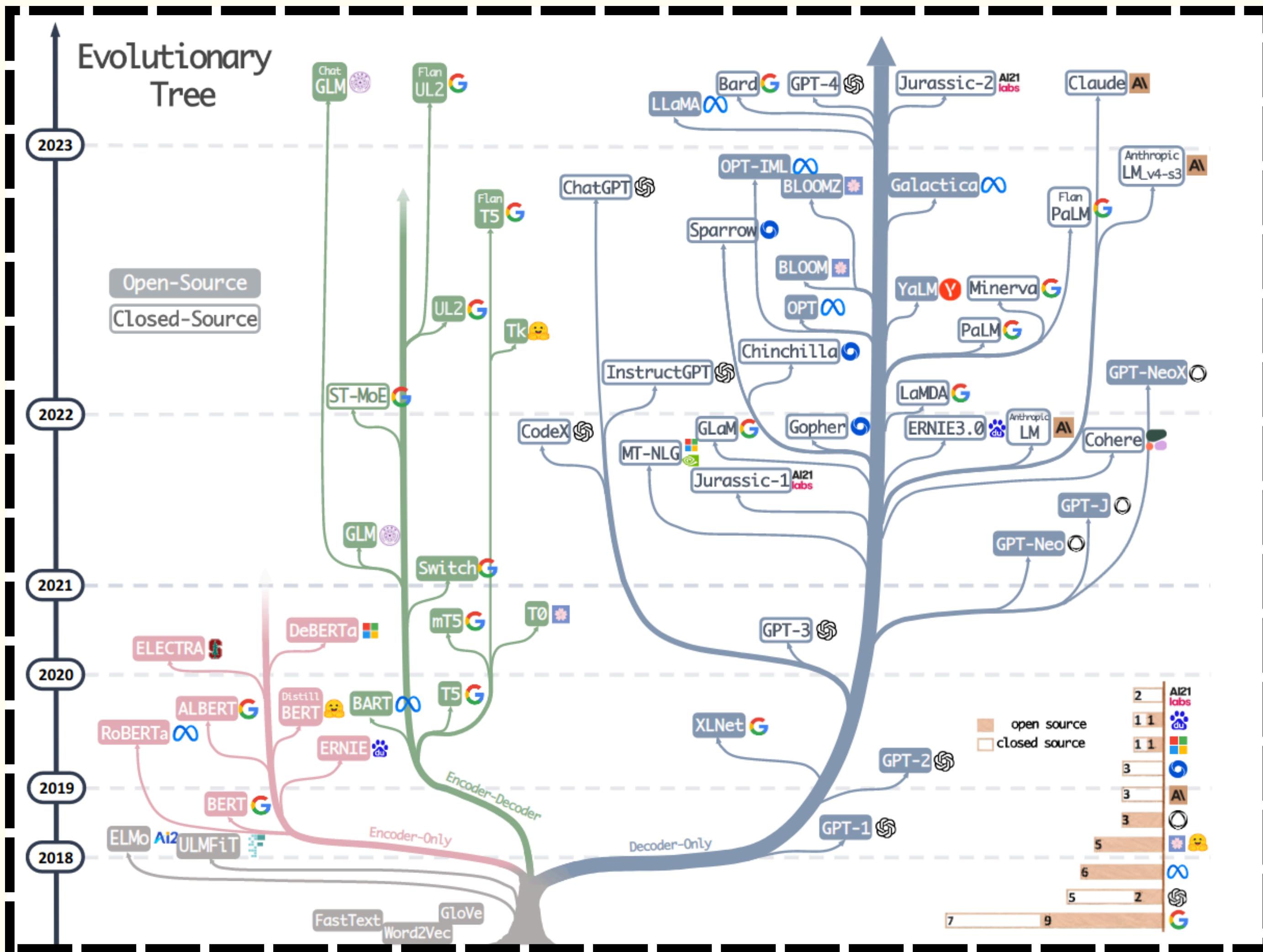
- Knowledge Representation
 - Ontologies, Logic Programming, Automated Theorem Proving
- Probabilistic Models
 - Bayesian Networks, Markov Chains, ...
- Neural Networks
 - FNN, CNN, RNN, Attention
- Neuro-symbolic

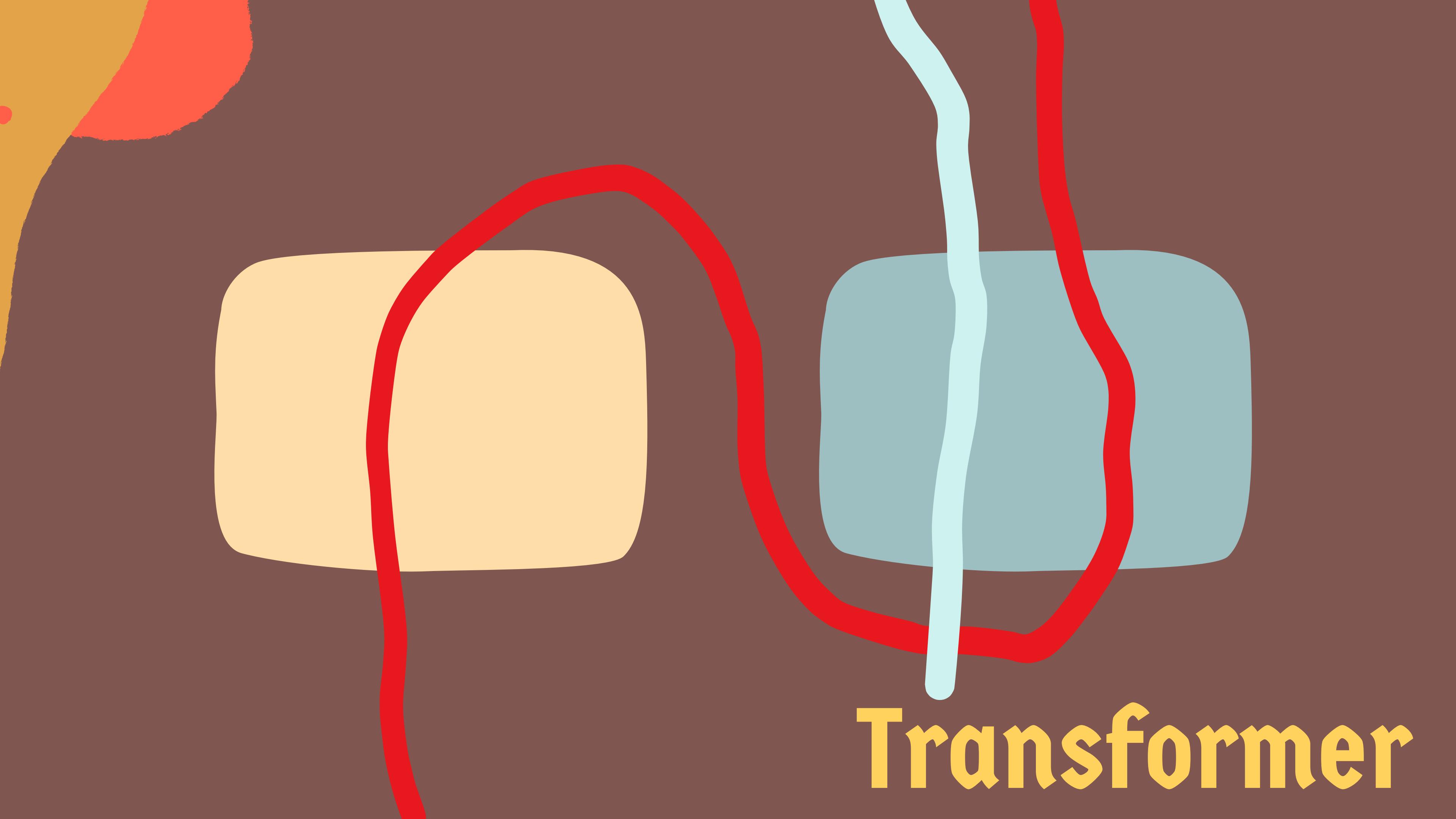


4. Which models are state-of-the-art in Natural Language Processing and Computational Linguistics?



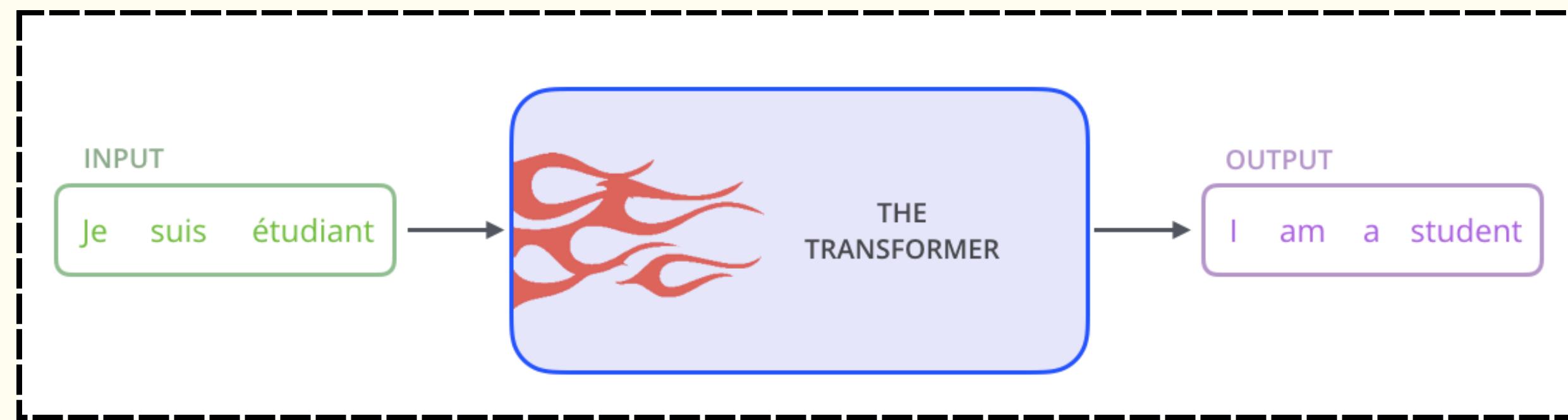
Evolutionary Tree

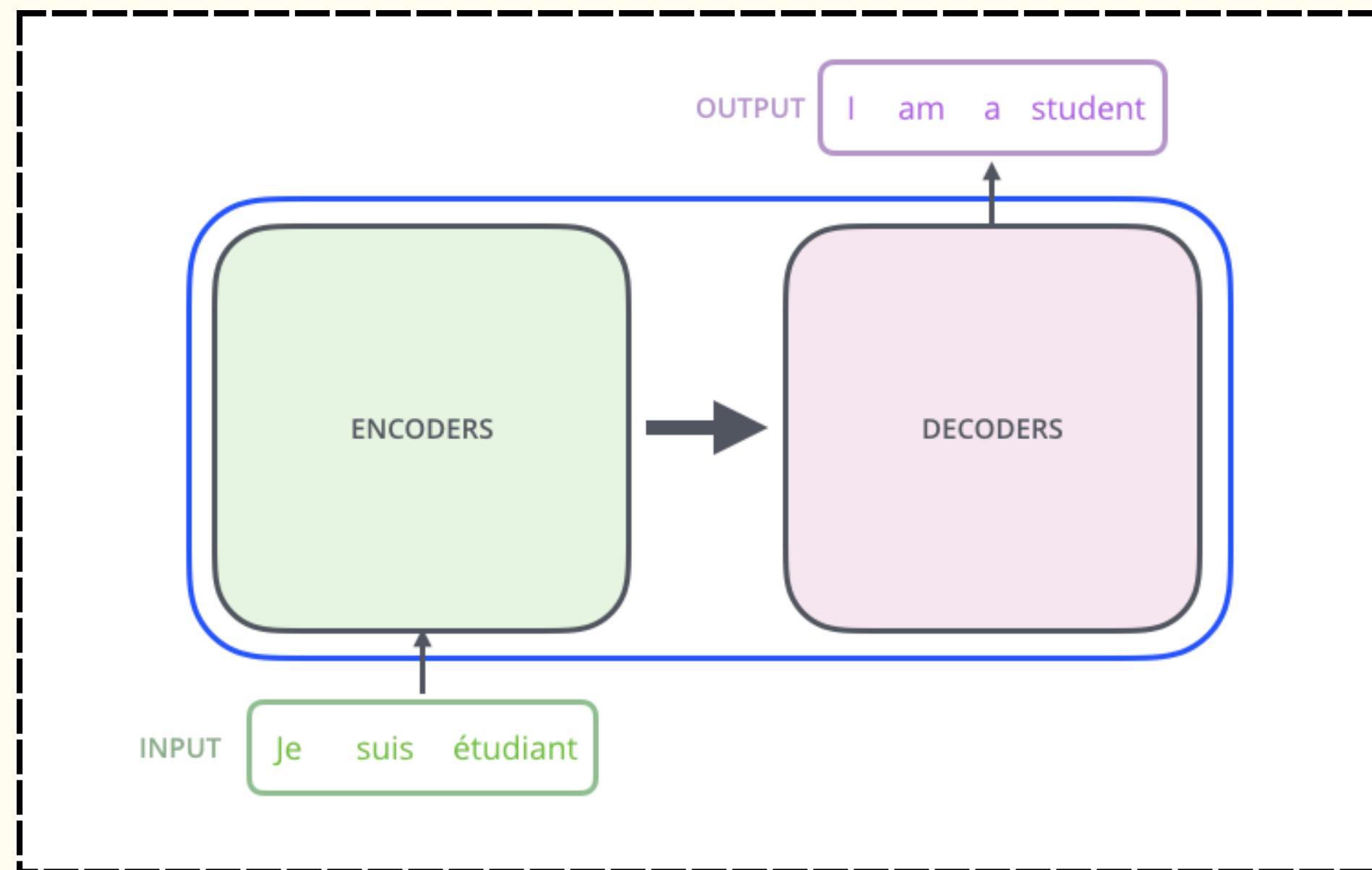


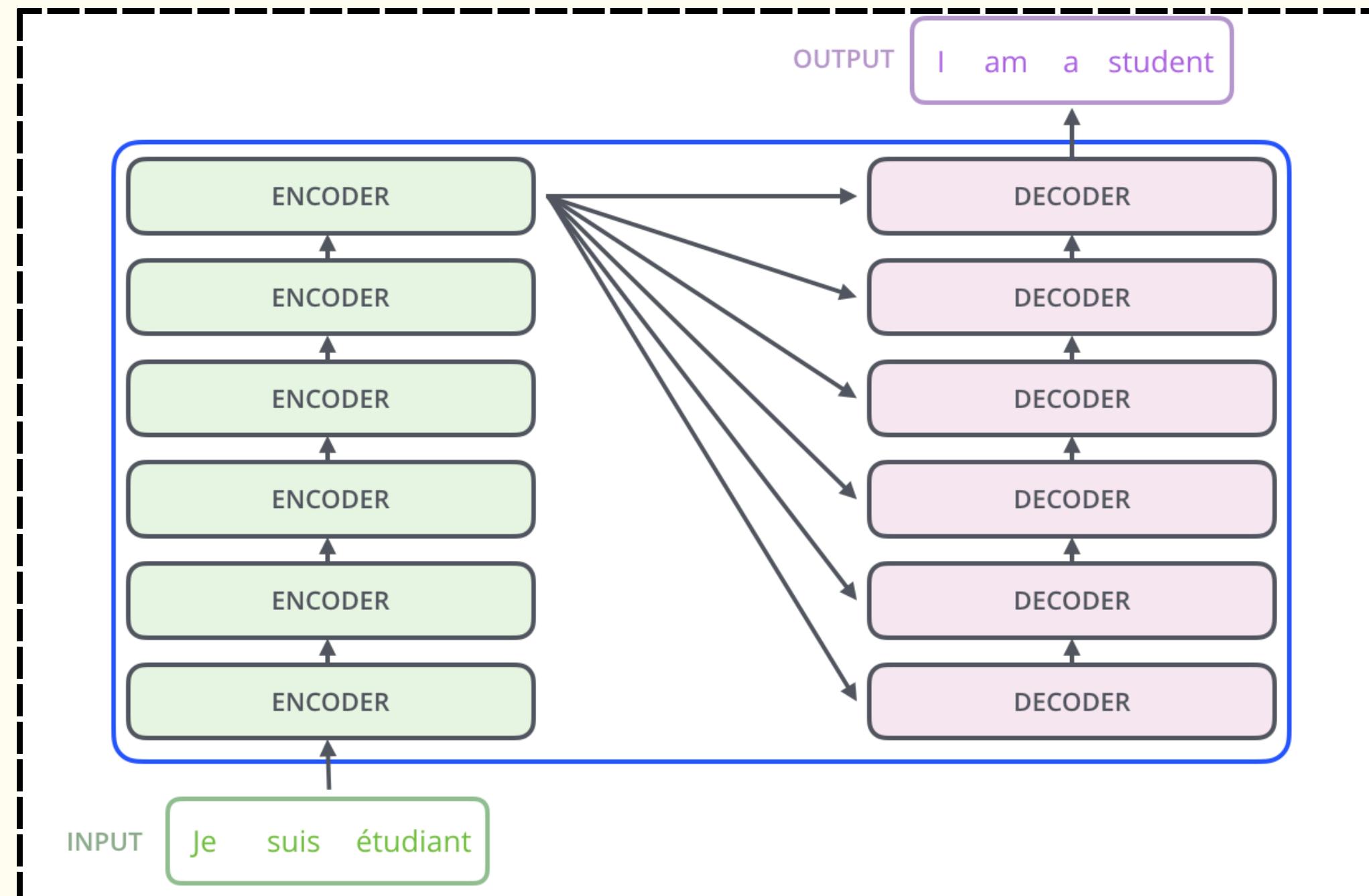


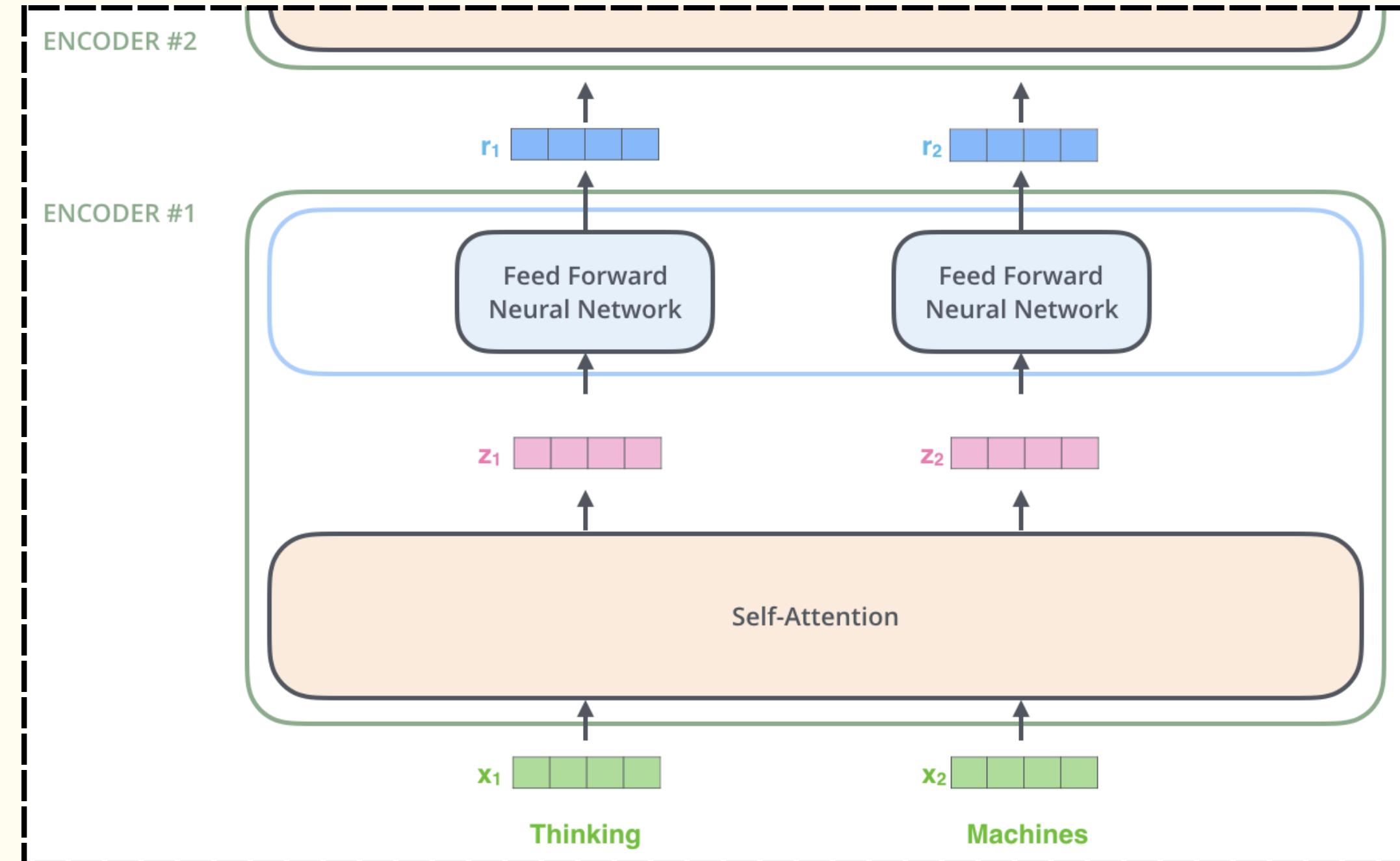
A stylized illustration of a transformer core. It features two main vertical columns of segments. The left column has two light yellow segments at the top, followed by a single large yellow segment, and then two more light yellow segments at the bottom. The right column has three light blue segments at the top, followed by a single light blue segment in the middle, and then two more light blue segments at the bottom. Red wavy lines connect the top and bottom segments of each column, forming a continuous loop. The background is a solid dark brown color.

Transformer









$$X \times W^Q = Q$$
$$X \times W^K = K$$
$$X \times W^V = V$$

$$\text{softmax}\left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}}\right) = \mathbf{z}$$

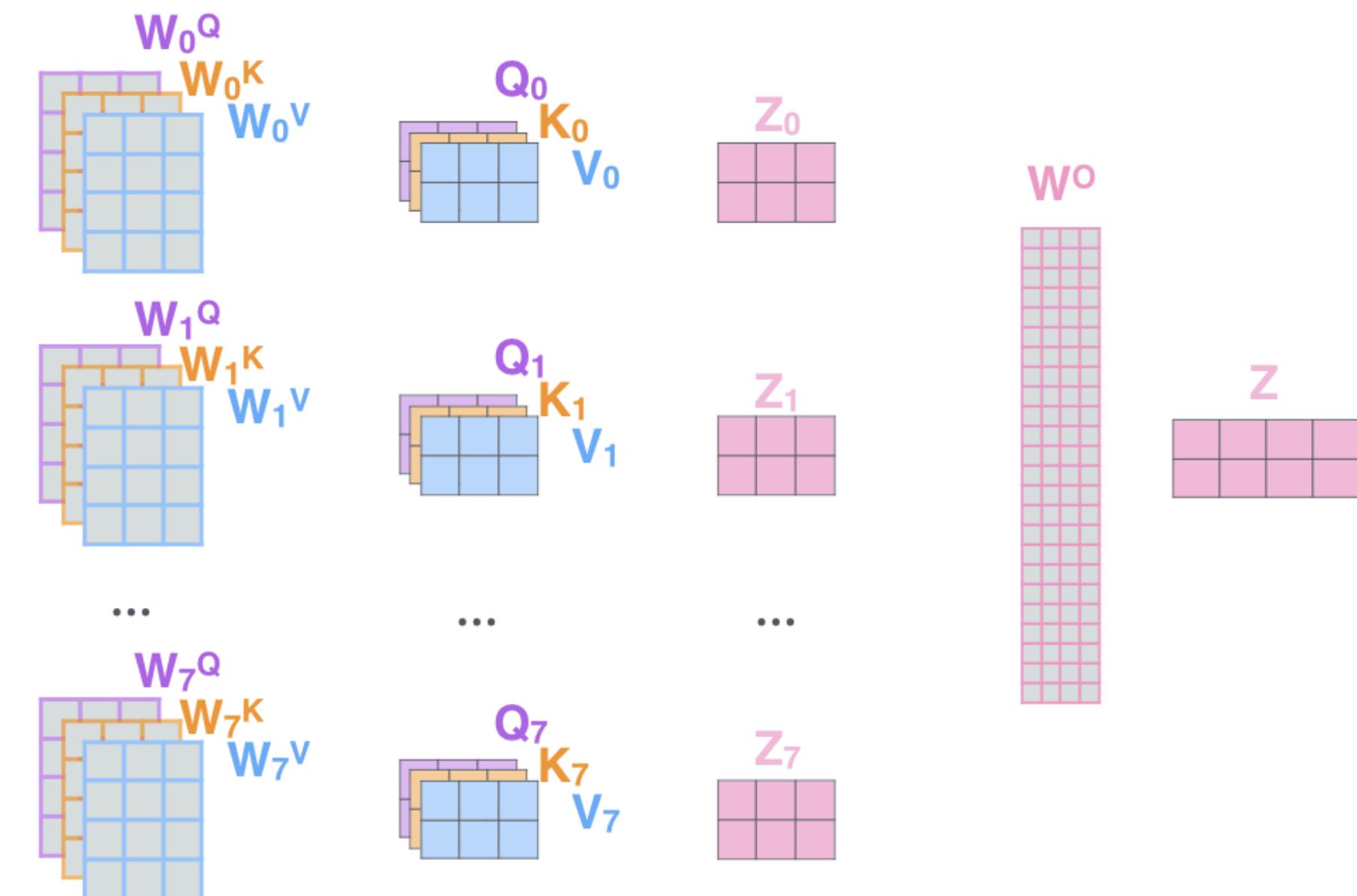
The diagram illustrates the computation of the attention matrix \mathbf{z} . It shows three input tensors: \mathbf{Q} (purple 3x3 matrix), \mathbf{K}^T (orange 3x3 matrix), and \mathbf{V} (blue 3x3 matrix). The multiplication of \mathbf{Q} and \mathbf{K}^T is scaled by $\sqrt{d_k}$ before being passed through a softmax function to produce the output matrix \mathbf{z} (pink 3x3 matrix).

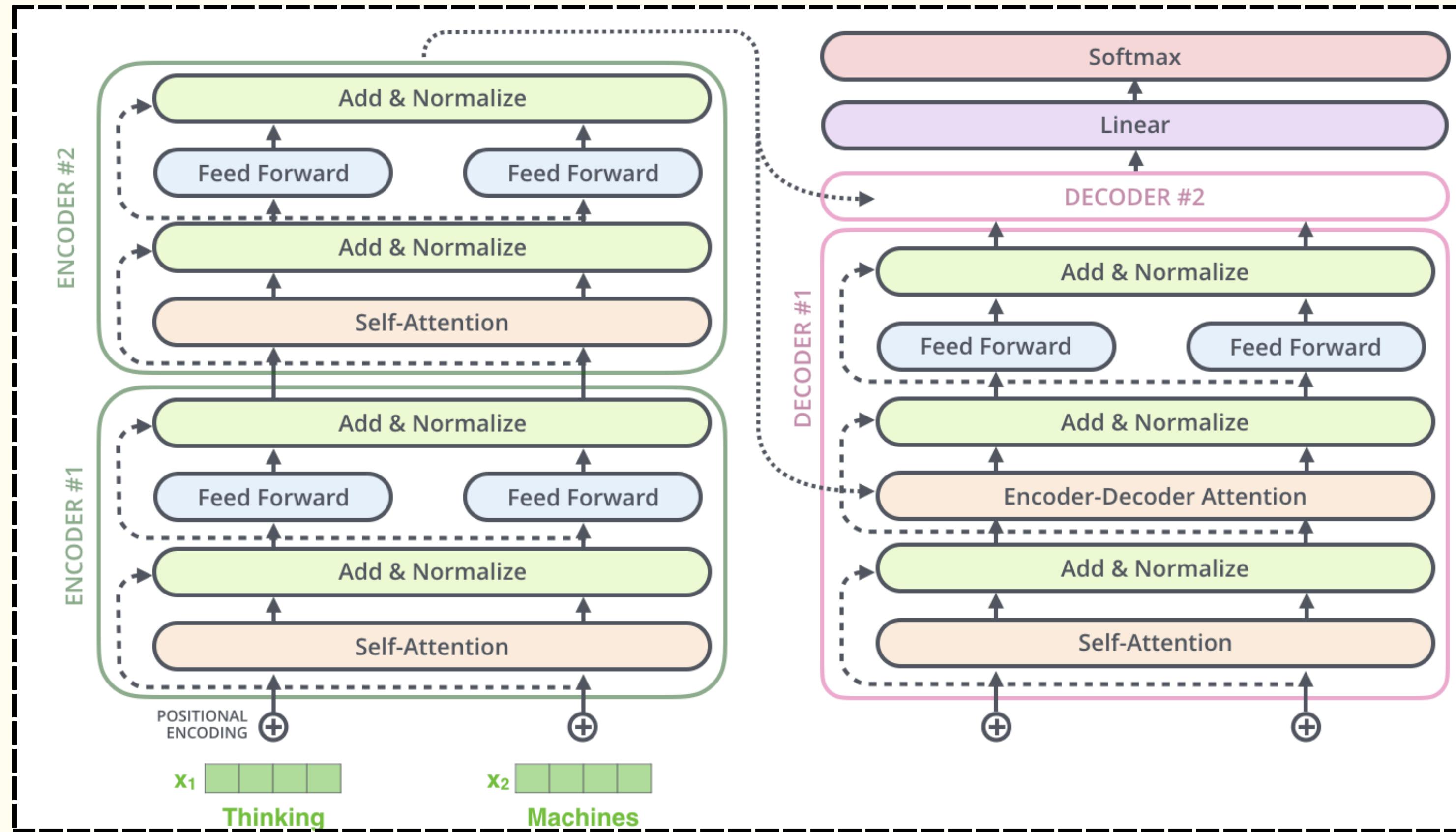
- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

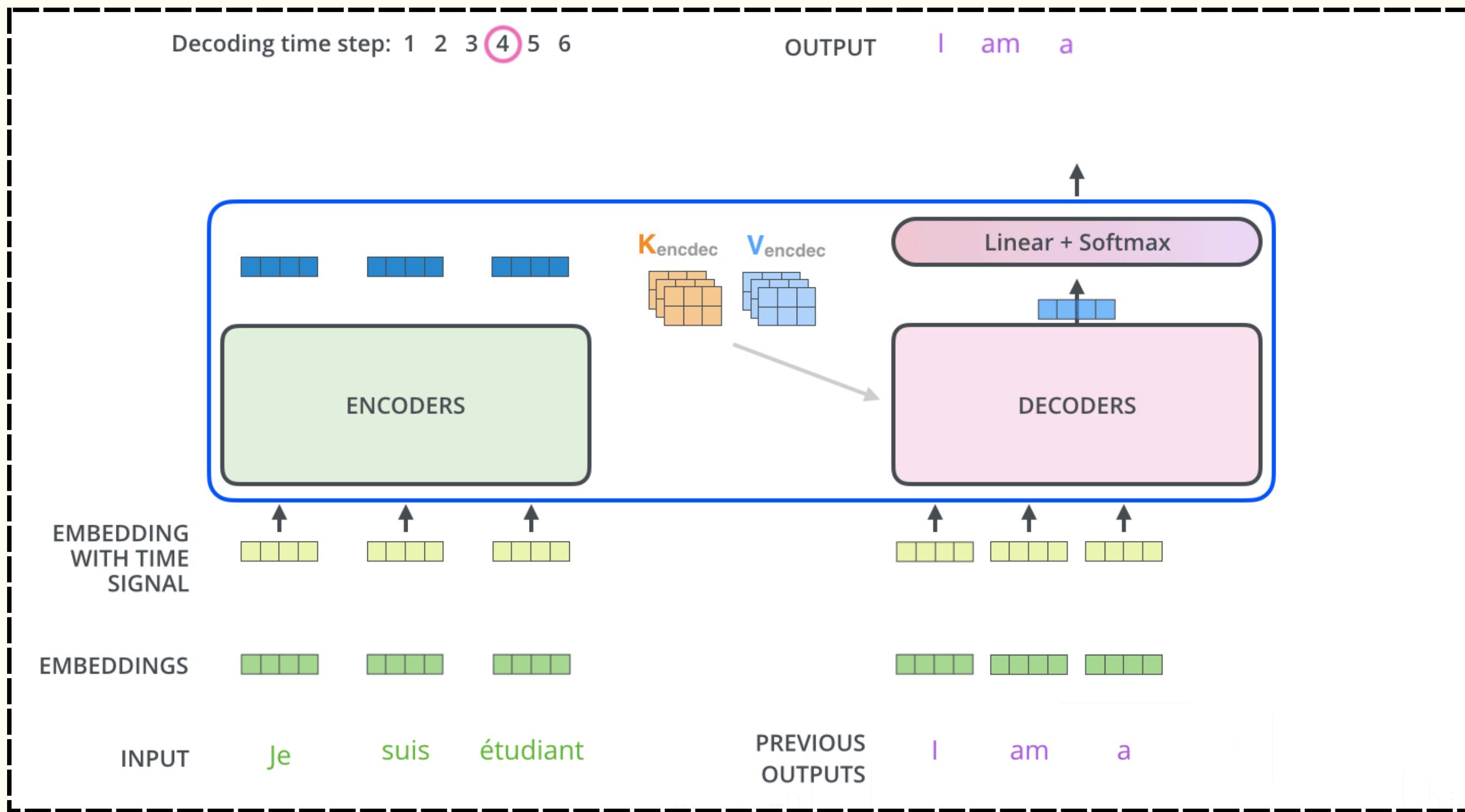
Thinking
Machines

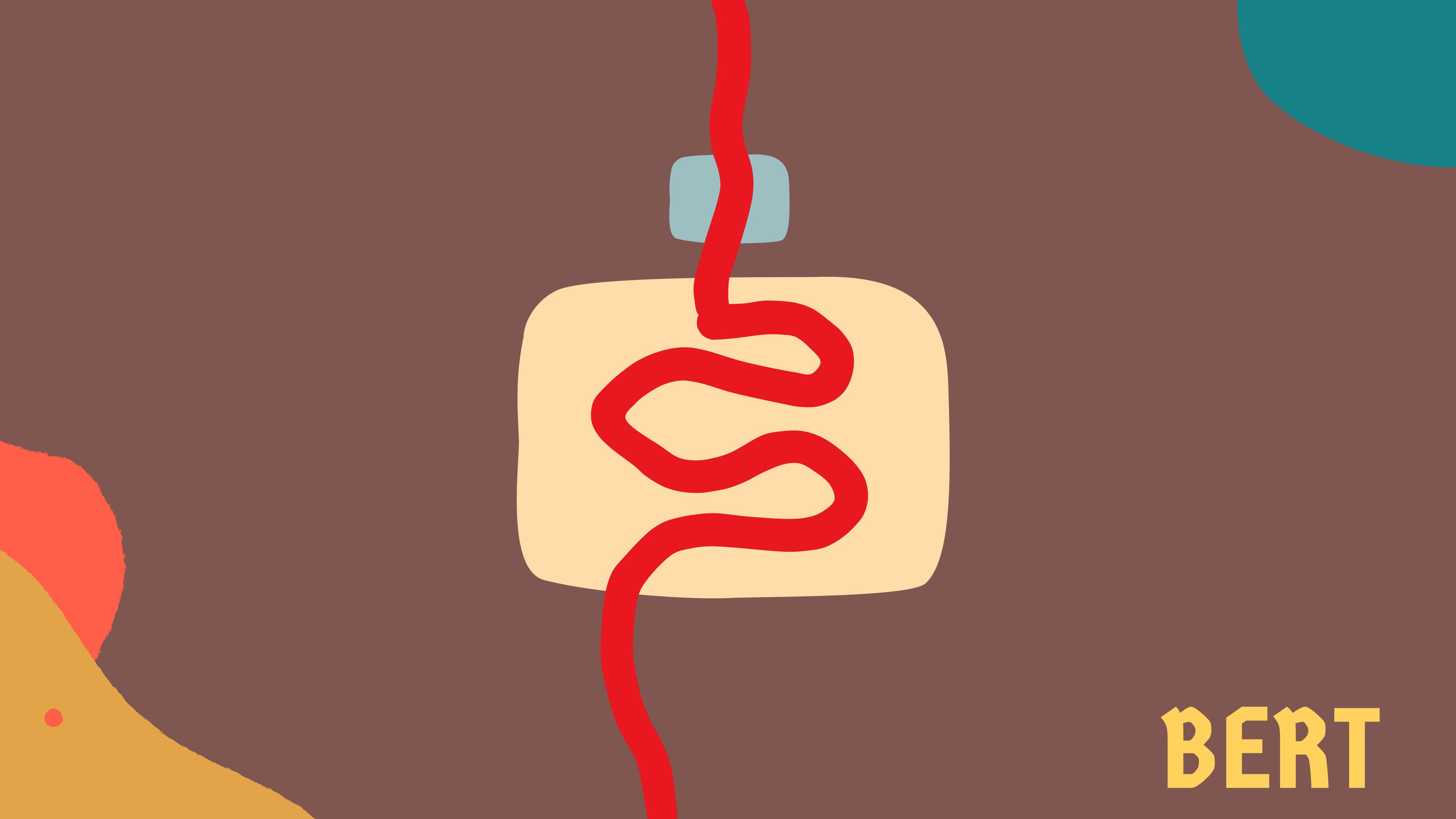


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

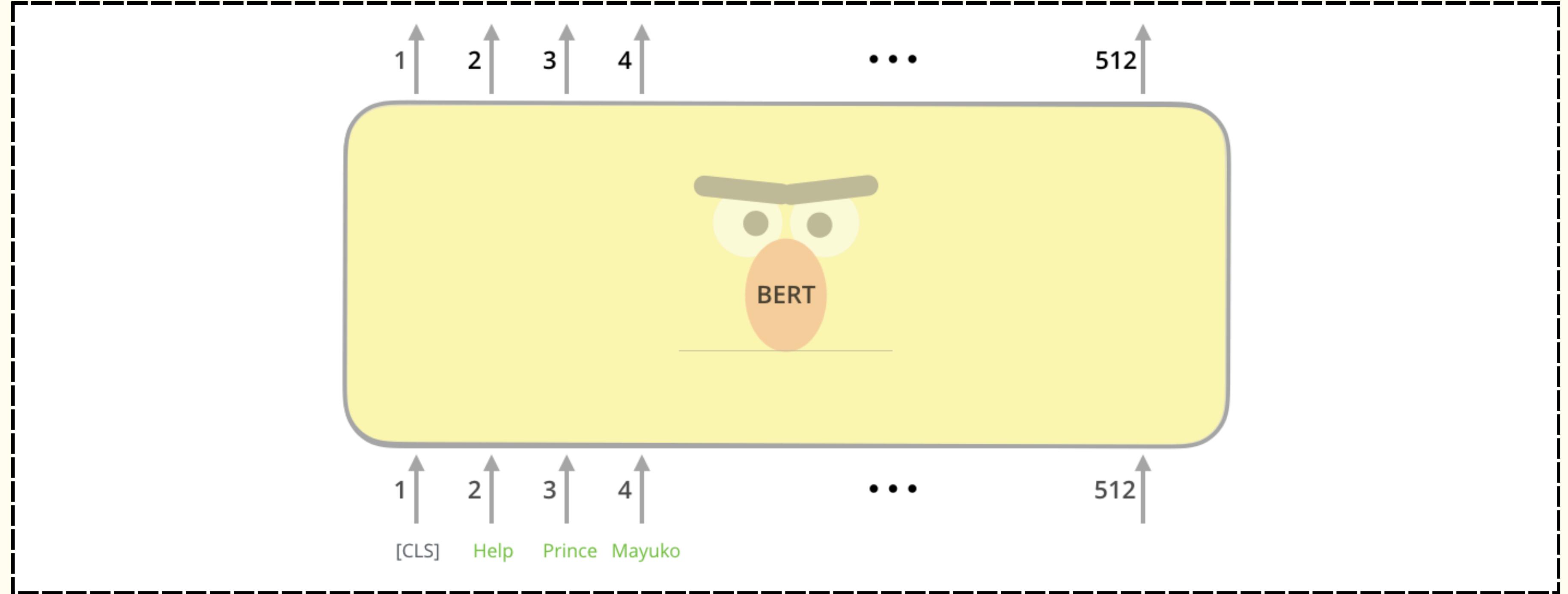


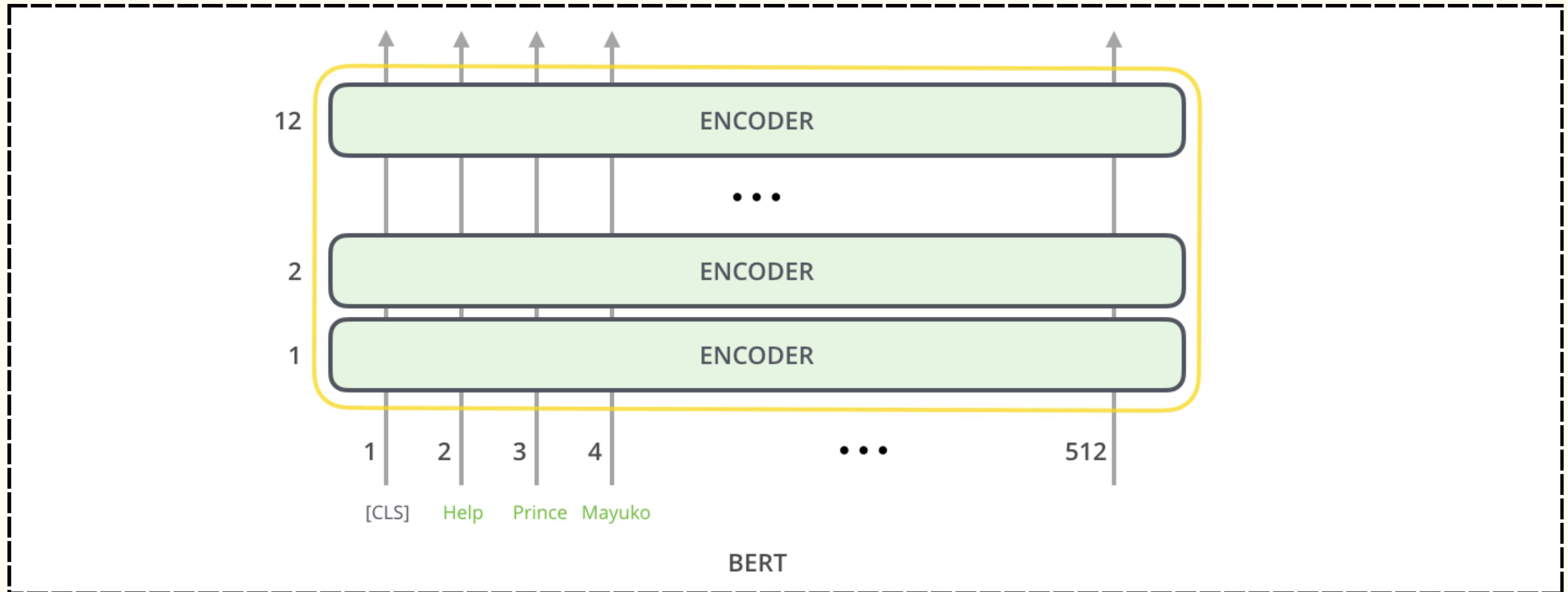


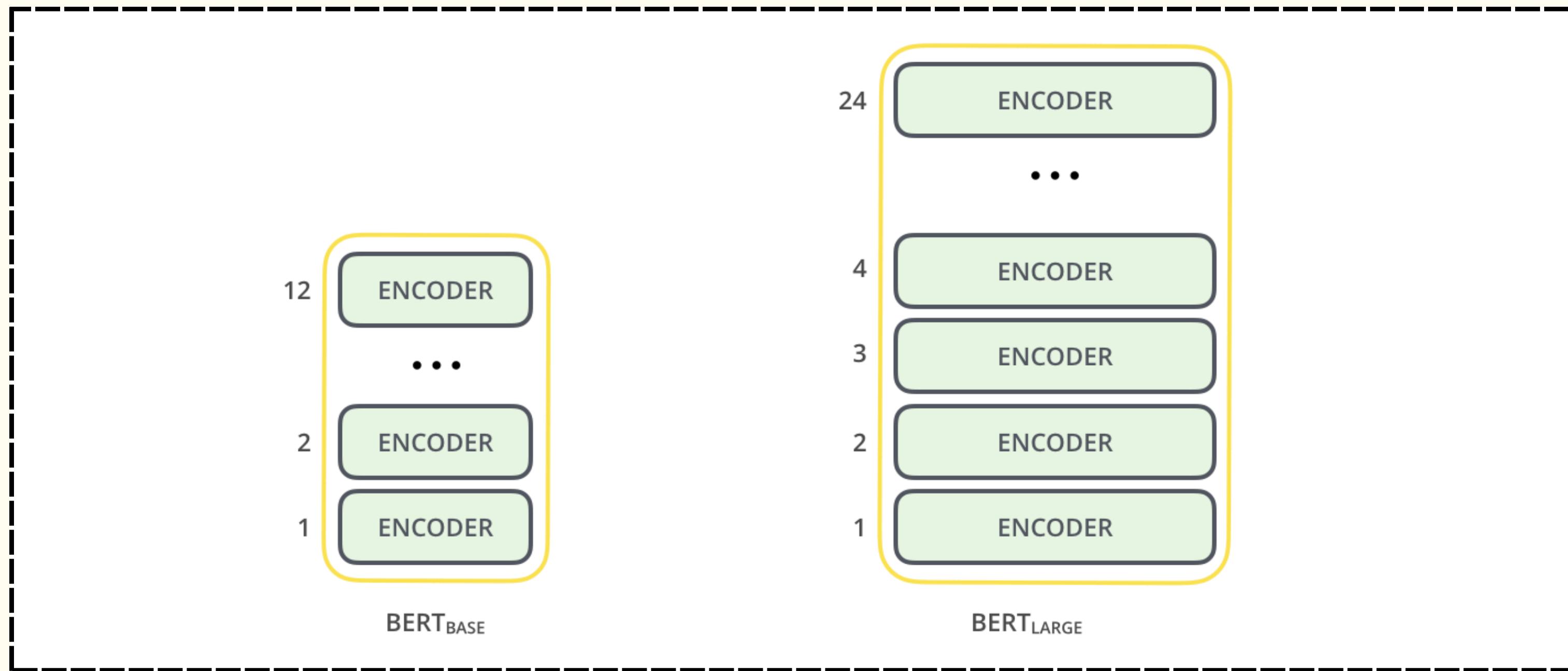


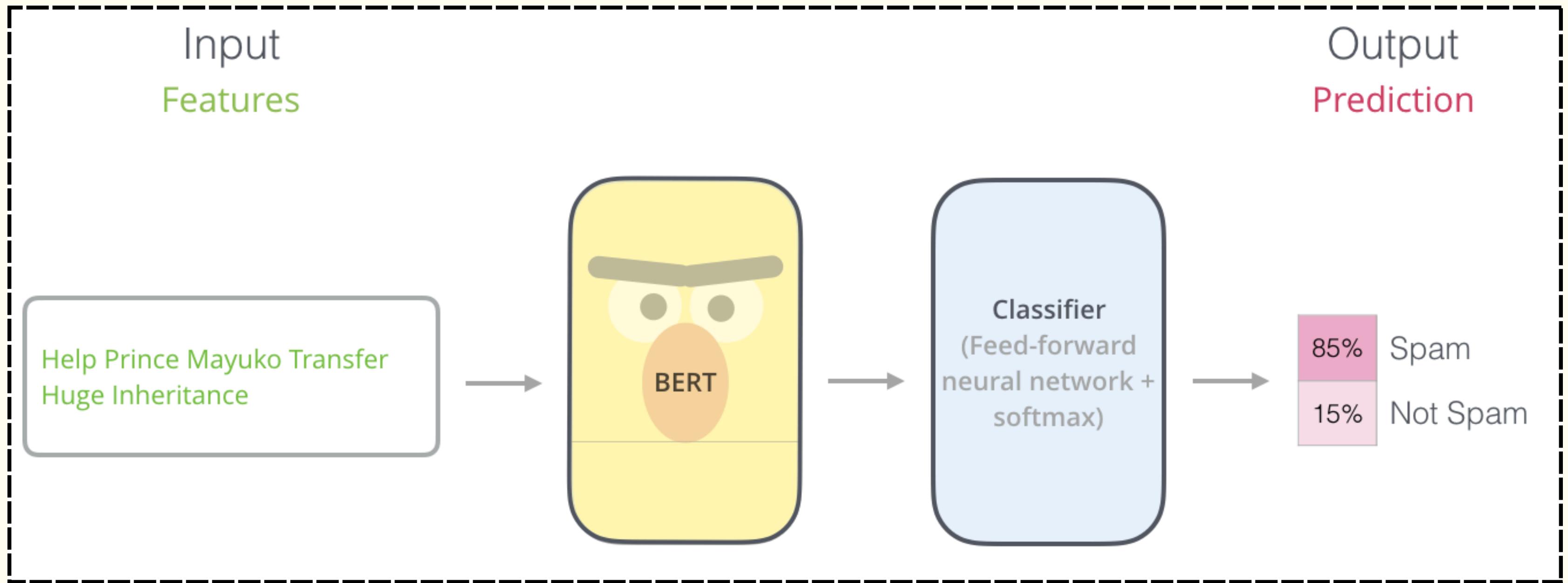


BERT









Use the output of the masked word's position to predict the masked word

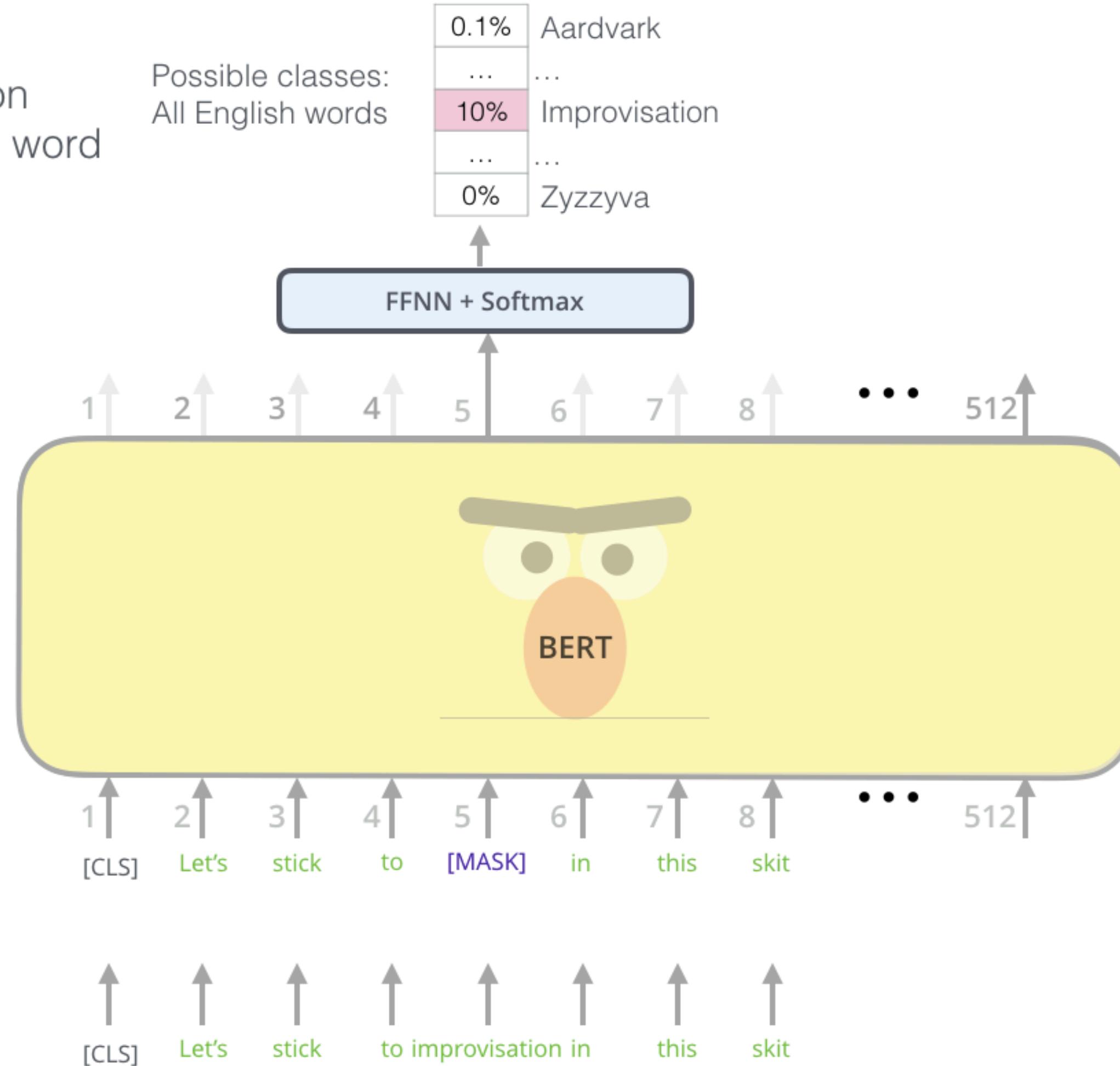
Possible classes:
All English words



FFNN + Softmax

Randomly mask
15% of tokens

Input



The end
El fin
O fim
M

Any Questions?

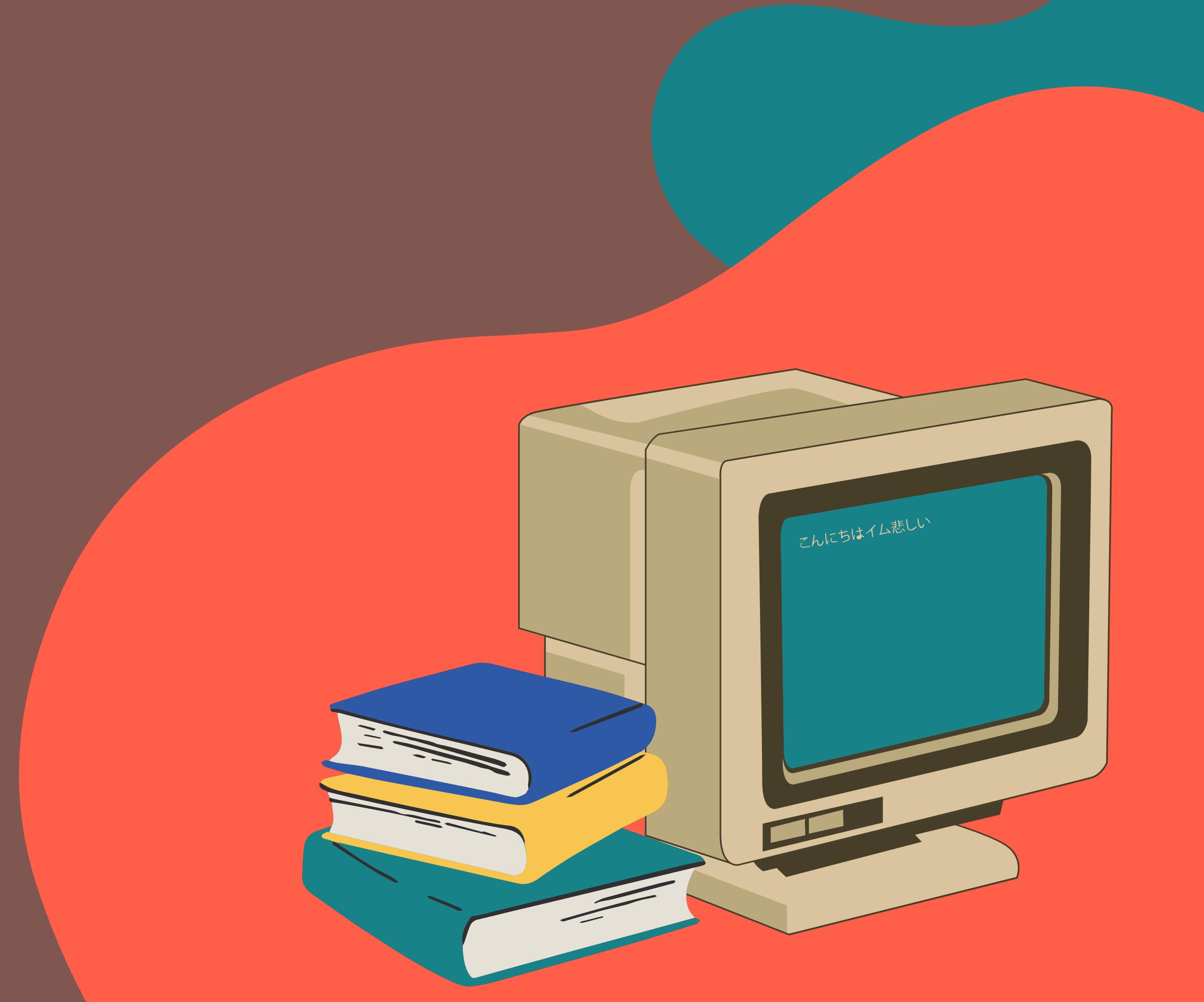


Image Sources

- Blank Venn Diagram Red, Blue and Yellow, by Amousey, under Public Domain
- Major Levels of Linguistic Structure, by James J. Thomas and Kristin A. Cook, and McSush, in Public Domain
- Word Embeddings CBOW, by Jeran Renz, under Creative Commons Attribution-Share Alike 4.0 International
- Transformers Evolutionary Tree, by Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, Xia Hu; extracted from <https://arxiv.org/abs/2304.13712>
- Transformer Illustrated by Jay Alammar, under Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License
 - Alammar, J (2018). The Illustrated Transformer [Blog post]. Retrieved from <https://jalammar.github.io/illustrated-transformer/>
- Illustrated BERT, by Jay Alammar, under Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License
 - Alammar, J (2018). The Illustrated BERT, ELMo and co. [Blog post]. Retrieved from <http://jalammar.github.io/illustrated-bert/>

These slides were made with Canva. Graphic Elements are licensed under the Canva Free License.