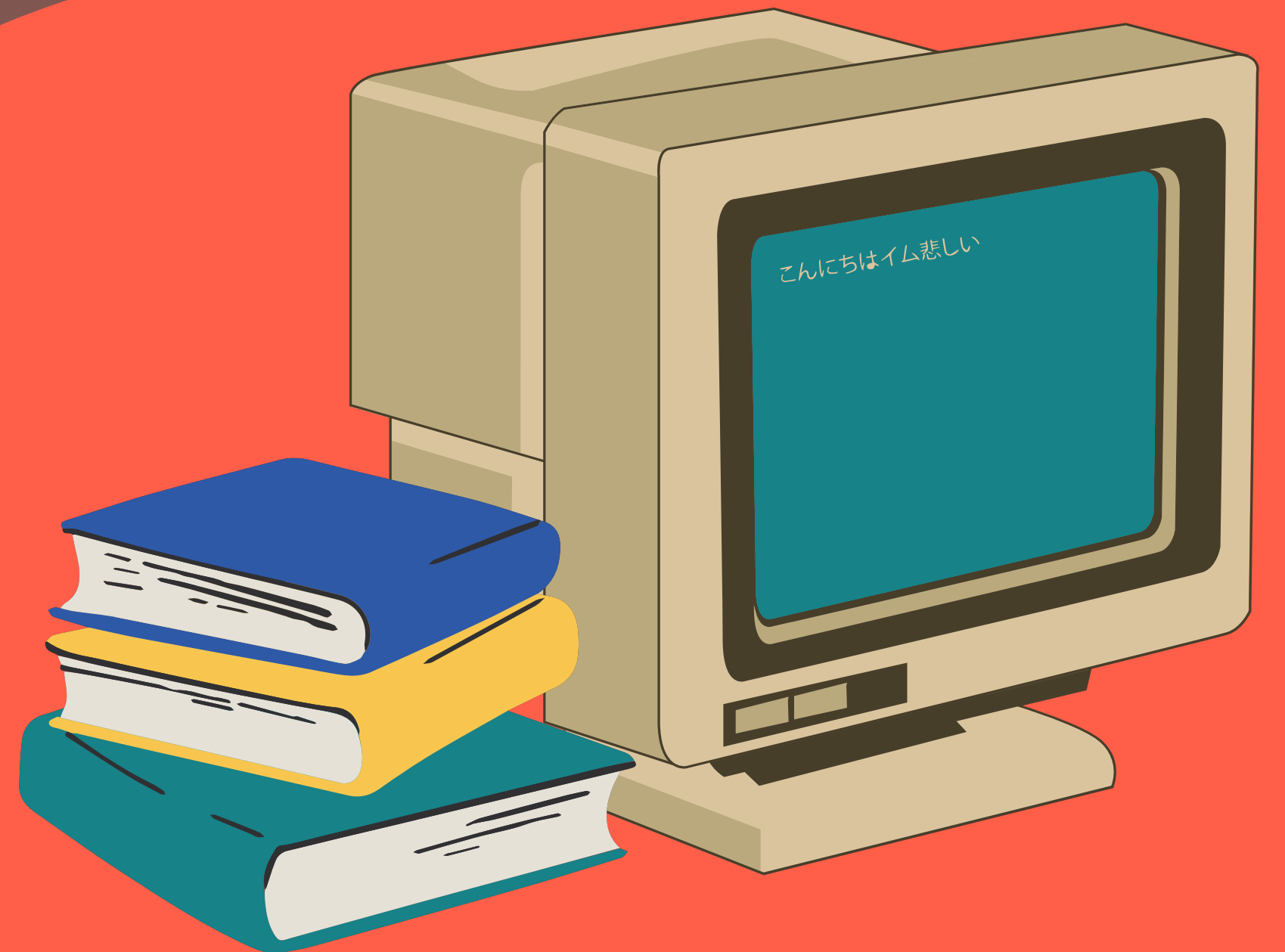


Natural Language Processing

Second School on Data Science and
Machine Learning, São Paulo, Brazil

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.

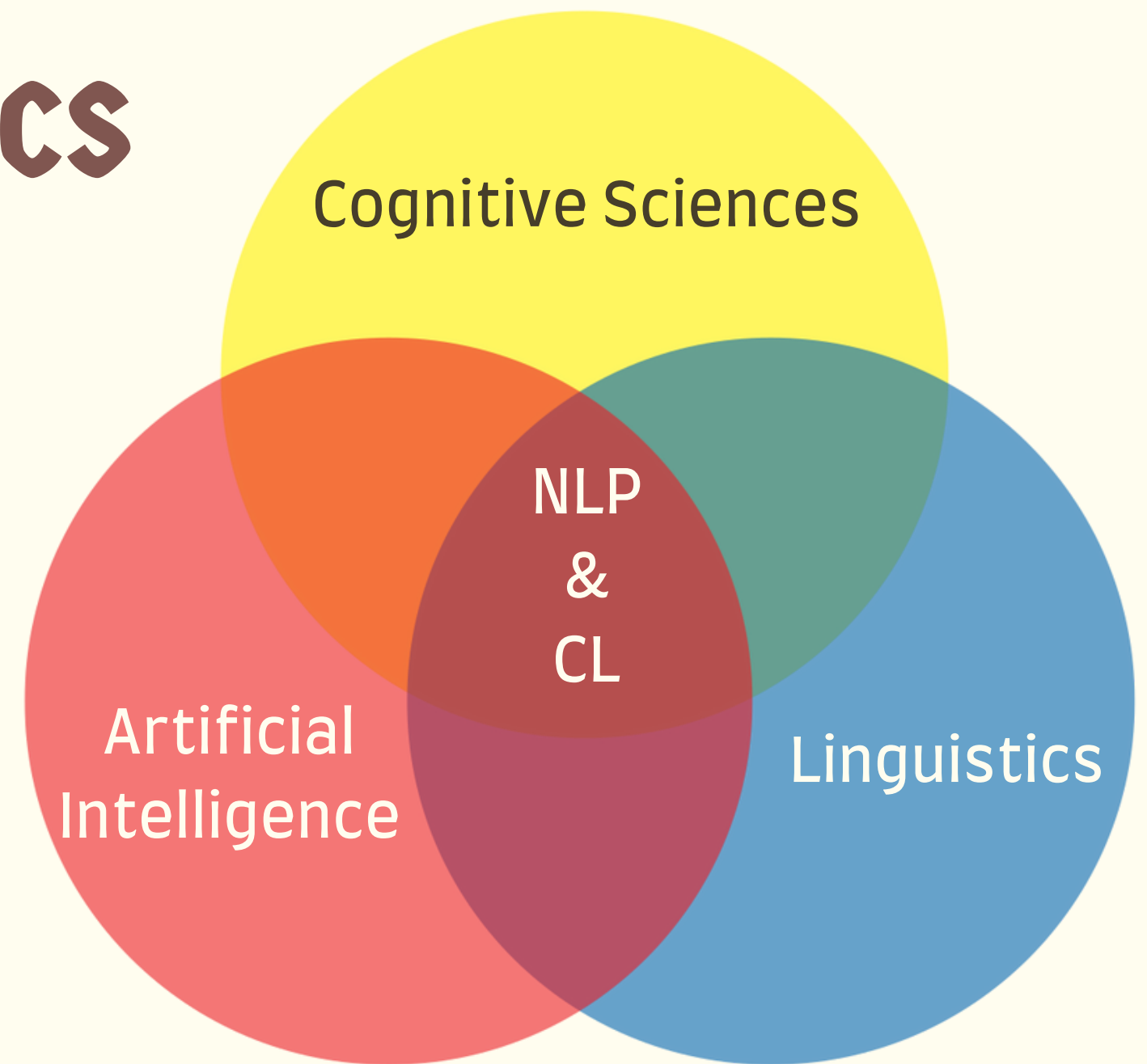
Which models are state-of-the-art in NLP and Computational Linguistics?


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 from the attempt to create computer programs to translate texts from Russian to English.





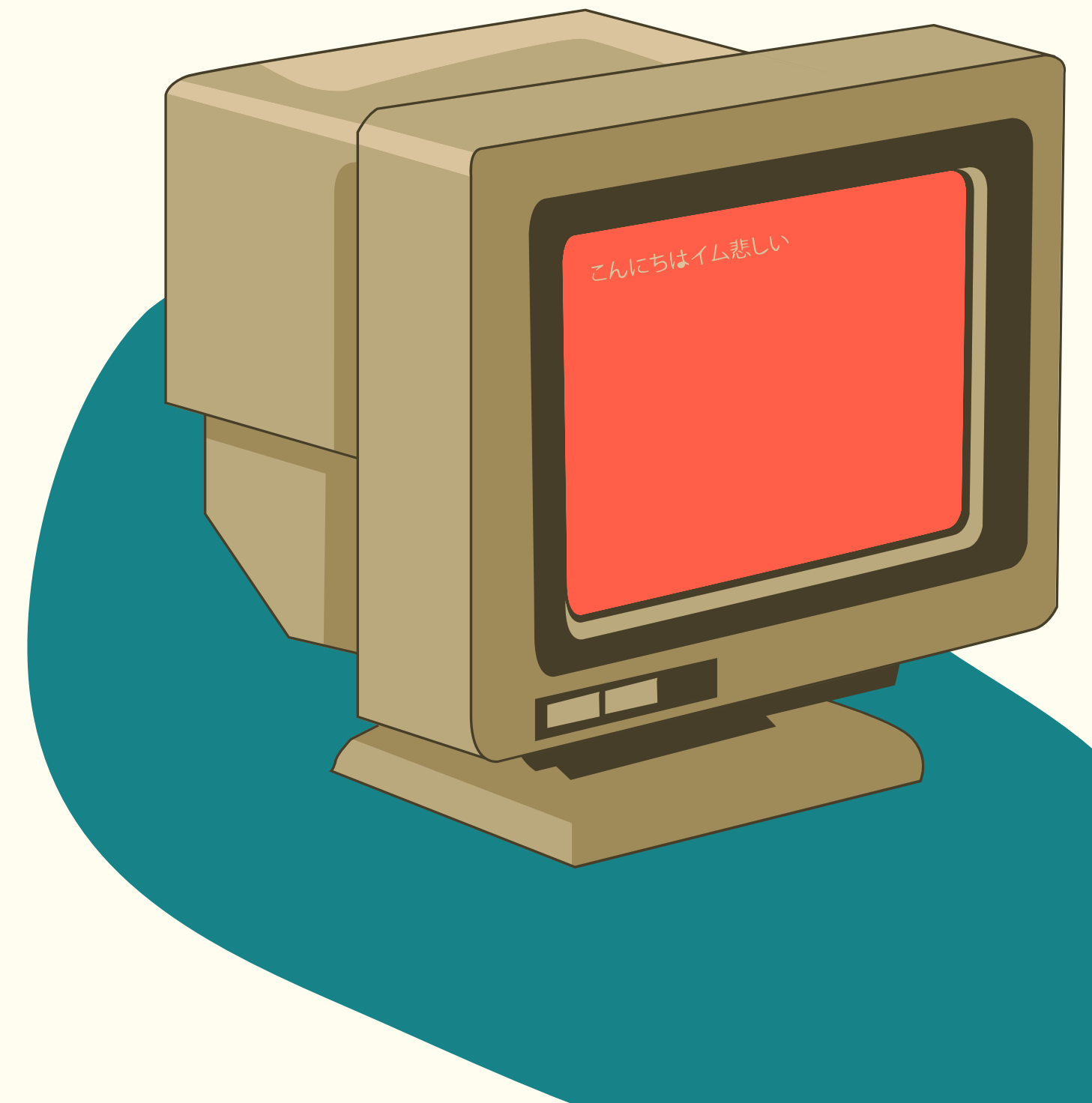
What is the difference between Natural Language Processing and Computational Linguistics

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

Tasks in Natural Language Processing

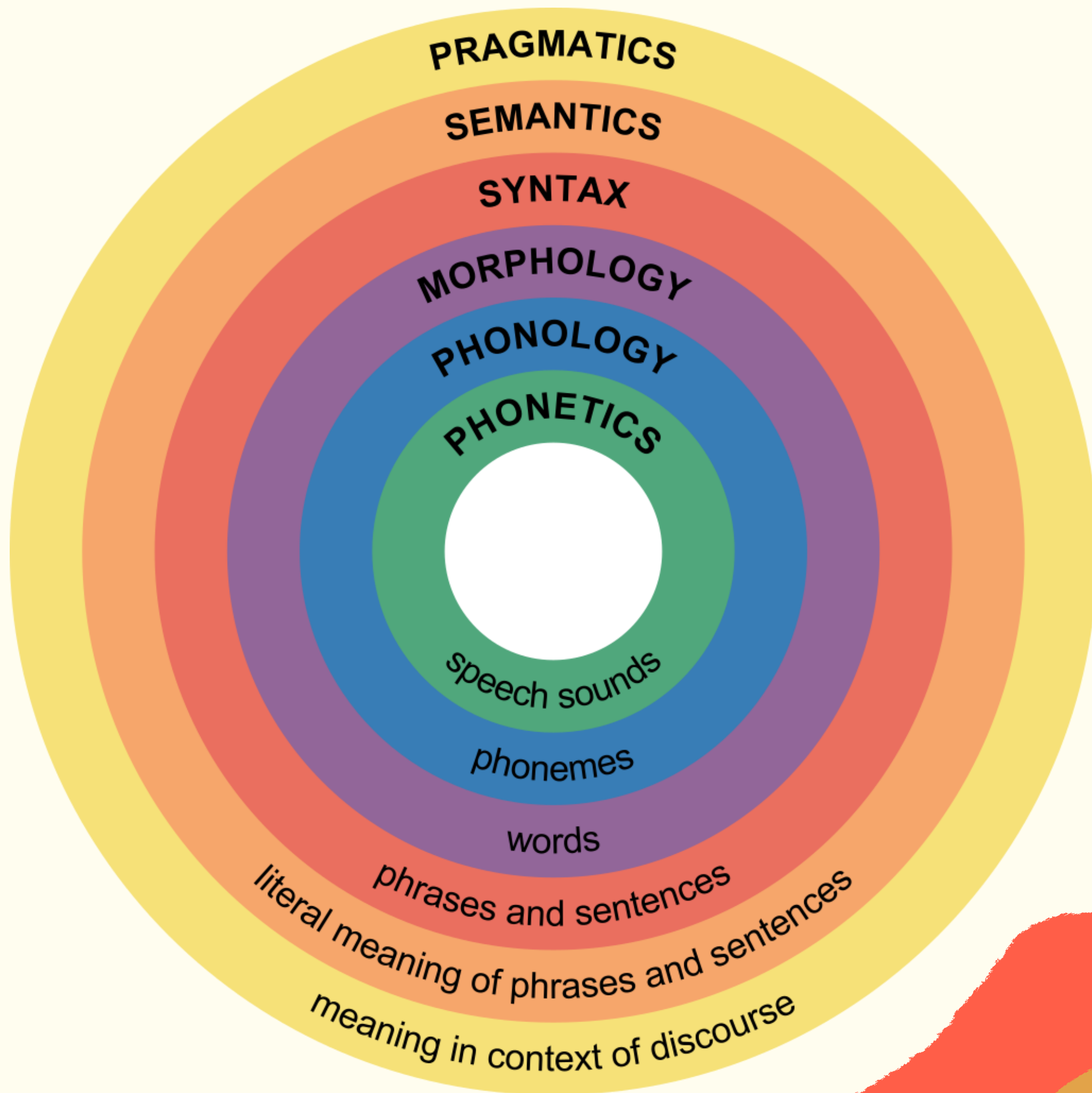
- **Translation**
- **Classification:** Sentiment Analysis, Spam Detection, Topic Classification
- **Regression:** Autograding
- **Clustering:** Topic Modeling, Authorship Attribution, Similarity-based Recommendations
- **Tagging:** Named Entity Recognition (NER), Part-of-Speech Tagging
- **Generation:**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

© 2015 IPA

	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ		
Trill	ʙ			r					ʀ		
Tap or Flap		ⱱ		ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral fricative				ɬ ɮ							
Approximant		ʋ		ɹ		ɻ	j	ɰ			
Lateral approximant				l		ɭ	ʎ	ʟ			

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
◌◌ Bilabial	ɓ Bilabial	ʼ Examples:
Dental	ɗ Dental/alveolar	pʼ Bilabial
! (Post)alveolar	ɟ Palatal	tʼ Dental/alveolar
≠ Palatoalveolar	ɡ Velar	kʼ Velar
Alveolar lateral	ɠ Uvular	sʼ Alveolar fricative

OTHER SYMBOLS

ʍ Voiceless labial-velar fricative ʎ Alveolo-palatal fricative

OTHER SYMBOLS

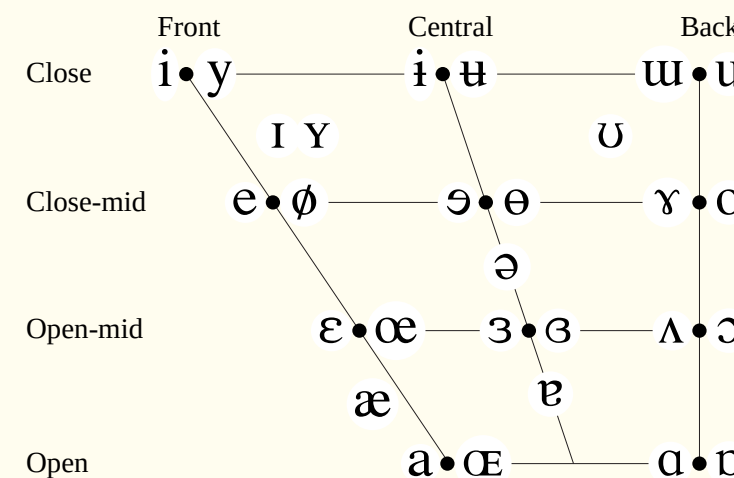
ʍ Voiceless labial-velar fricative ʎ Alveolo-palatal fricatives
 W Voiced labial-velar approximant ɺ Voiced alveolar lateral flap
 ɥ Voiced labial-palatal approximant ɥ Simultaneous ʃ and x
 H Voiceless epiglottal fricative
 ʕ Voiced epiglottal fricative Affricates and double articulations can be represented by two symbols joined by a tie bar if necessary.
 ʔ Epiglottal plosive

ts k̟p

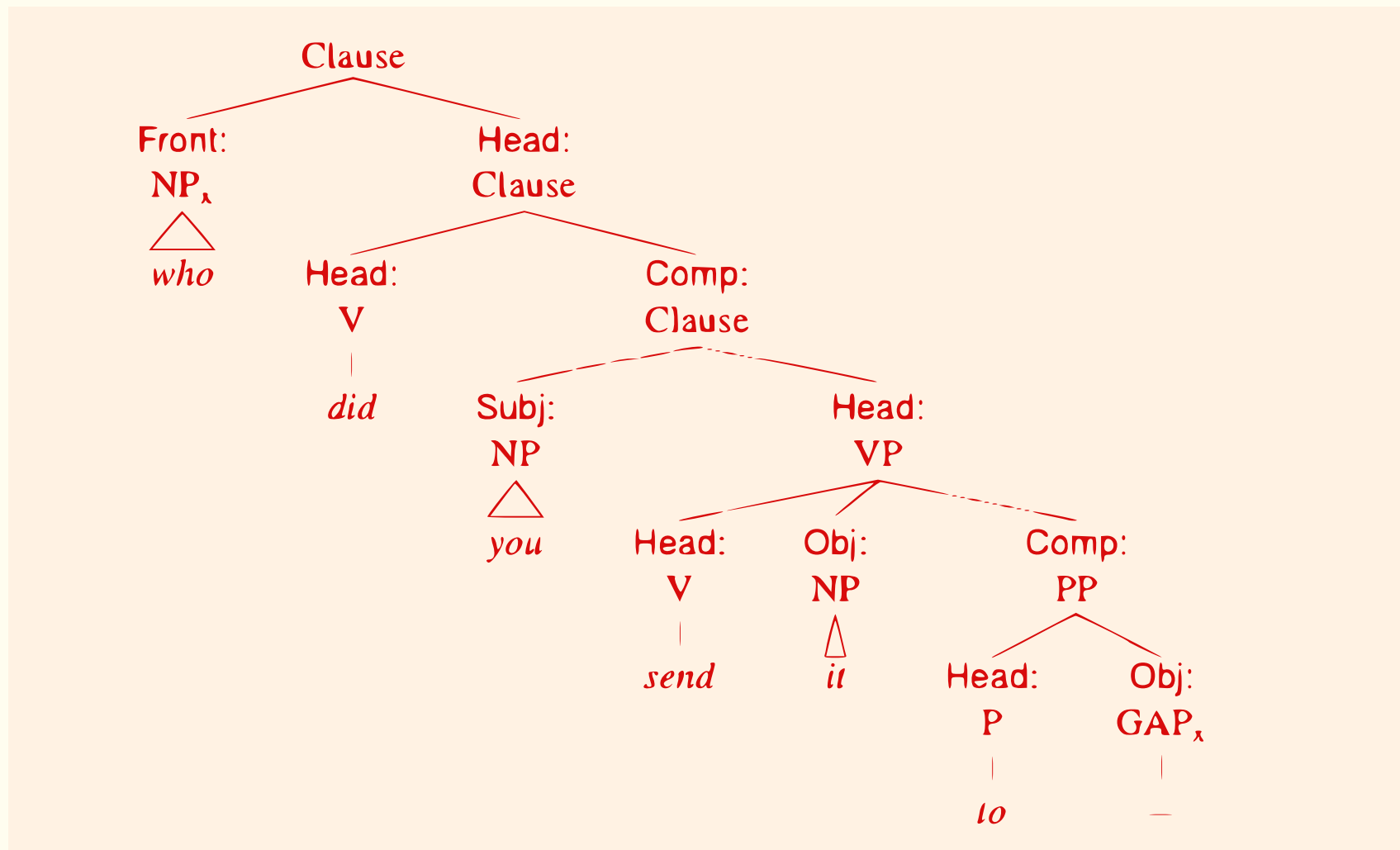
DIACRITICS Some diacritics may be placed above a symbol with a descender, e.g. ɲ̥

◌̥ Voiceless	◌̤ Voiced	◌̰ Breathy voiced	◌̠ Dental	◌̡ Apical
◌̧ Voiced	◌̨ Aspirated	◌̩ Linguolabial	◌̪ Laminal	◌̫ Nasalized
◌̜ More rounded	◌̝ Less rounded	◌̞ Labialized	◌̟ Palatalized	◌̠ Nasal release
◌̡ Advanced	◌̢ Retracted	◌̣ Velarized	◌̤ Pharyngealized	◌̥ Lateral release
◌̦ Centralized	◌̧ Mid-centralized	◌̨ Raised	◌̩ Lowered	◌̪ Advanced Tongue Root
◌̫ Syllabic	◌̬ Non-syllabic	◌̭ Velarized or pharyngealized	◌̮	◌̯ Retracted Tongue Root
◌̰ Rhoticity	◌̱	◌̲	◌̳	◌̴

VOWELS



Syntax

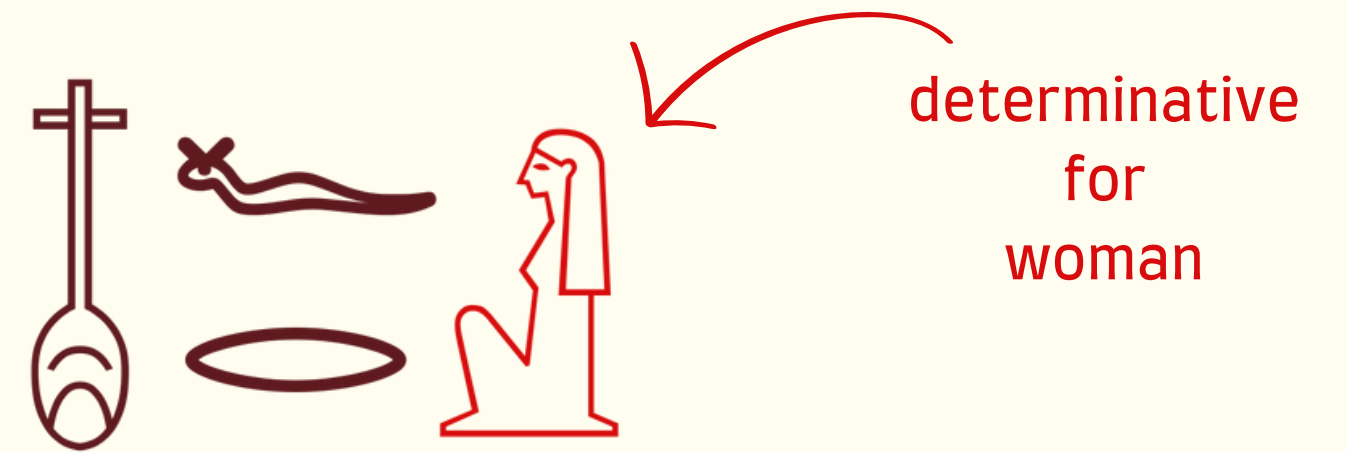


NOUN	ADJECTIVE	DETERMINER
VERB	ADVERB	CONJUNCTION
NUMERAL	PRONOUN	INTERJECTION
ADPOSITION	PARTICLE	...

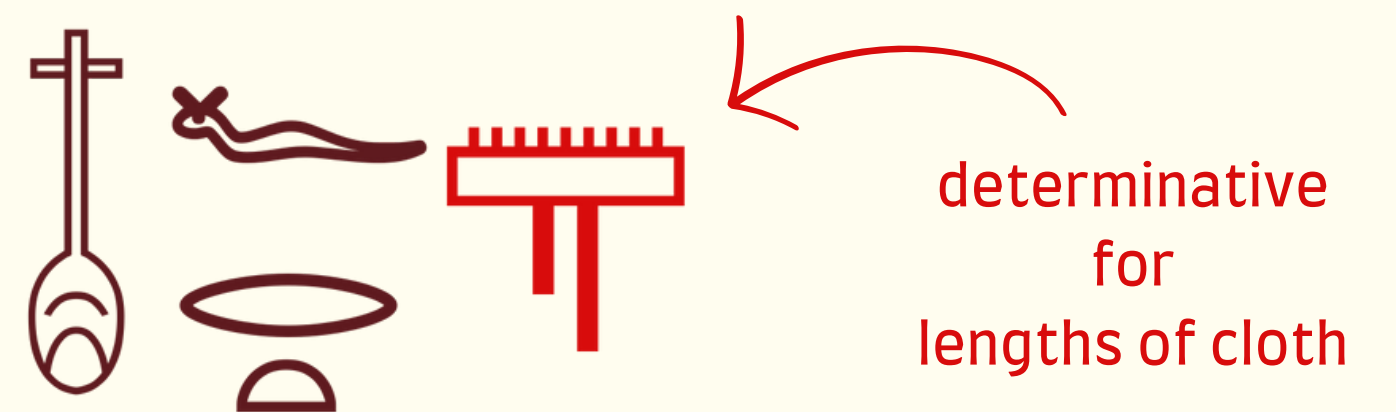
Morphology

- Derivation
 - break → breakable → unbreakable (English)
- Inflection
 - Portuguese:
 - Eu fal**o**, tu fal**as**, nós fal**amos**
 - Spanish
 - Yo habl**o**, tu habl**as**, nosotros habl**amos**.

- Egyptian:



young woman of marriageable age



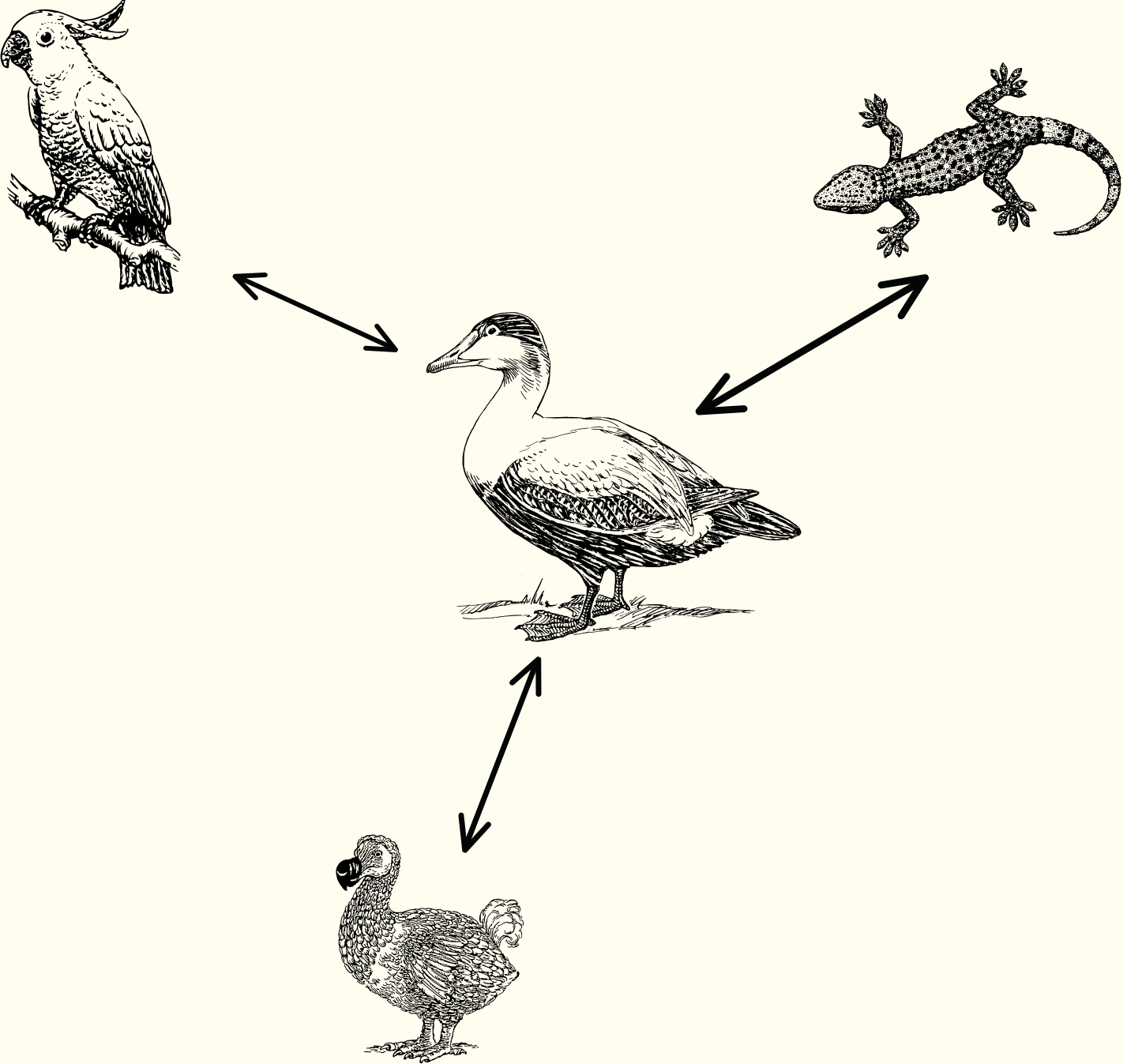
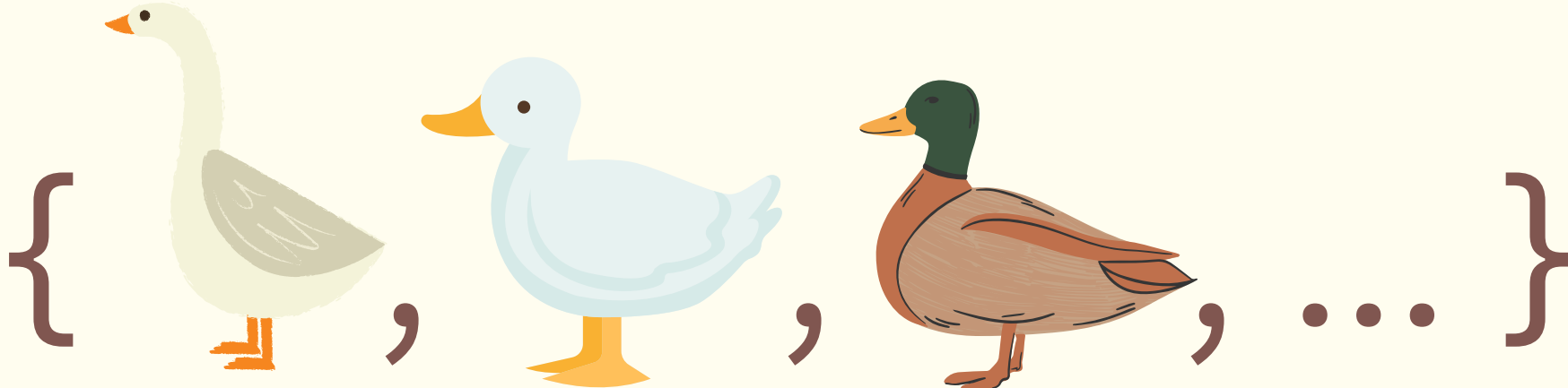
clothing



Semantics

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

Is it the set of all possible ducks?



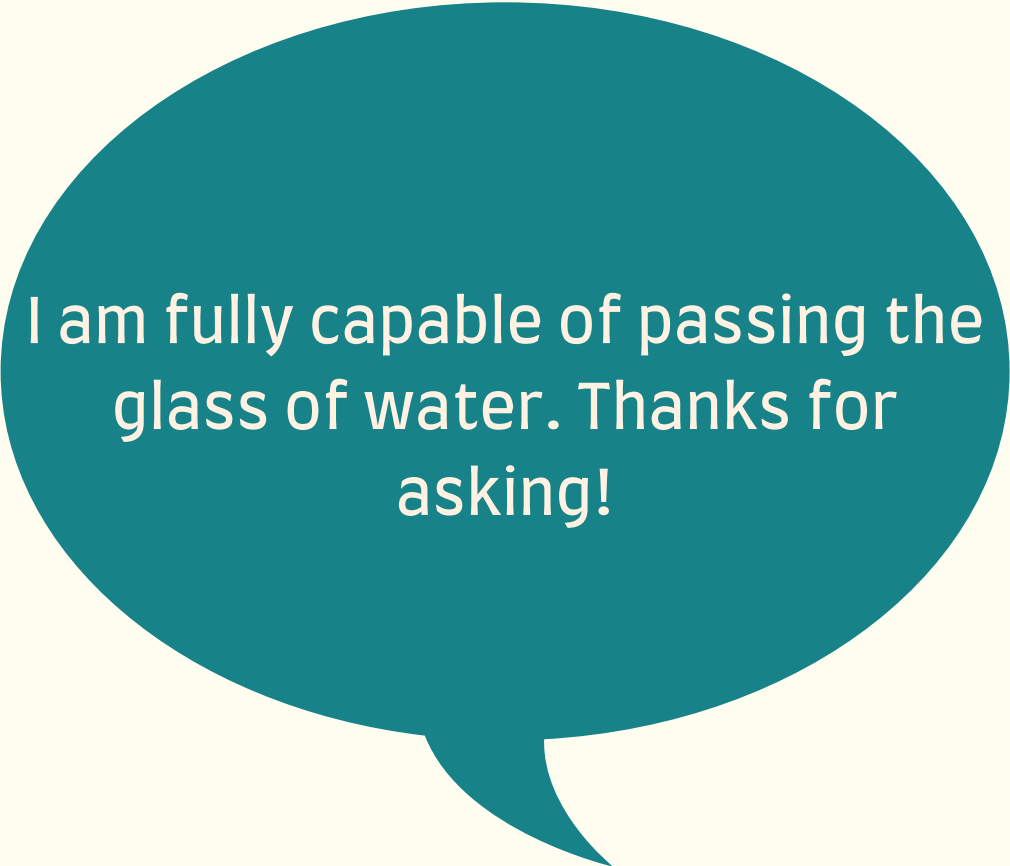
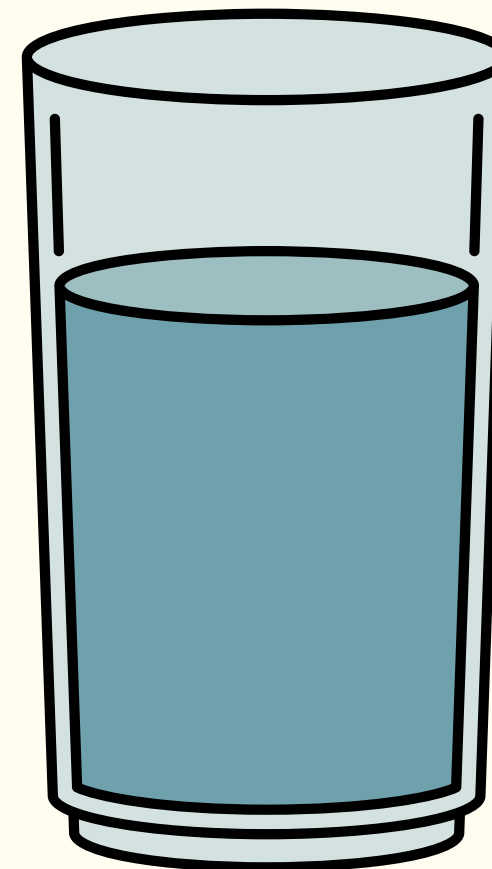
Is it the distance to a prototypical duck?

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!

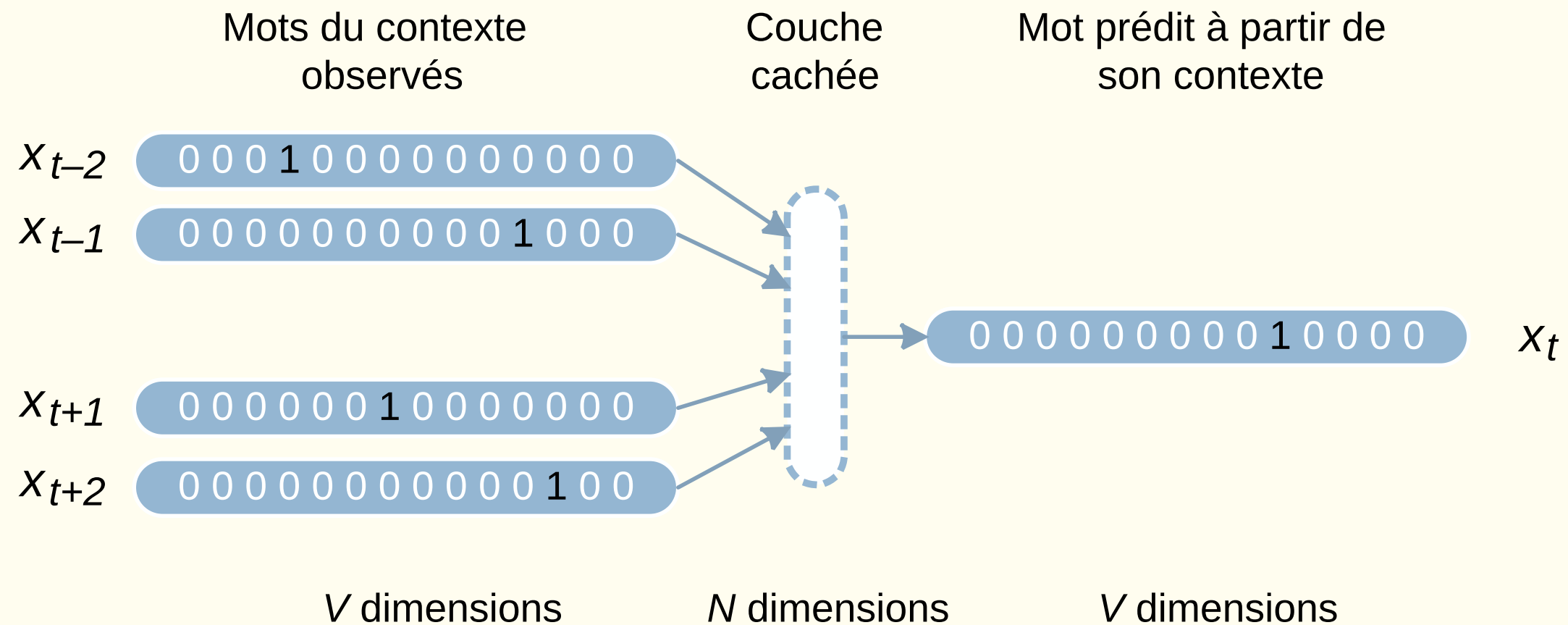
3. How

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

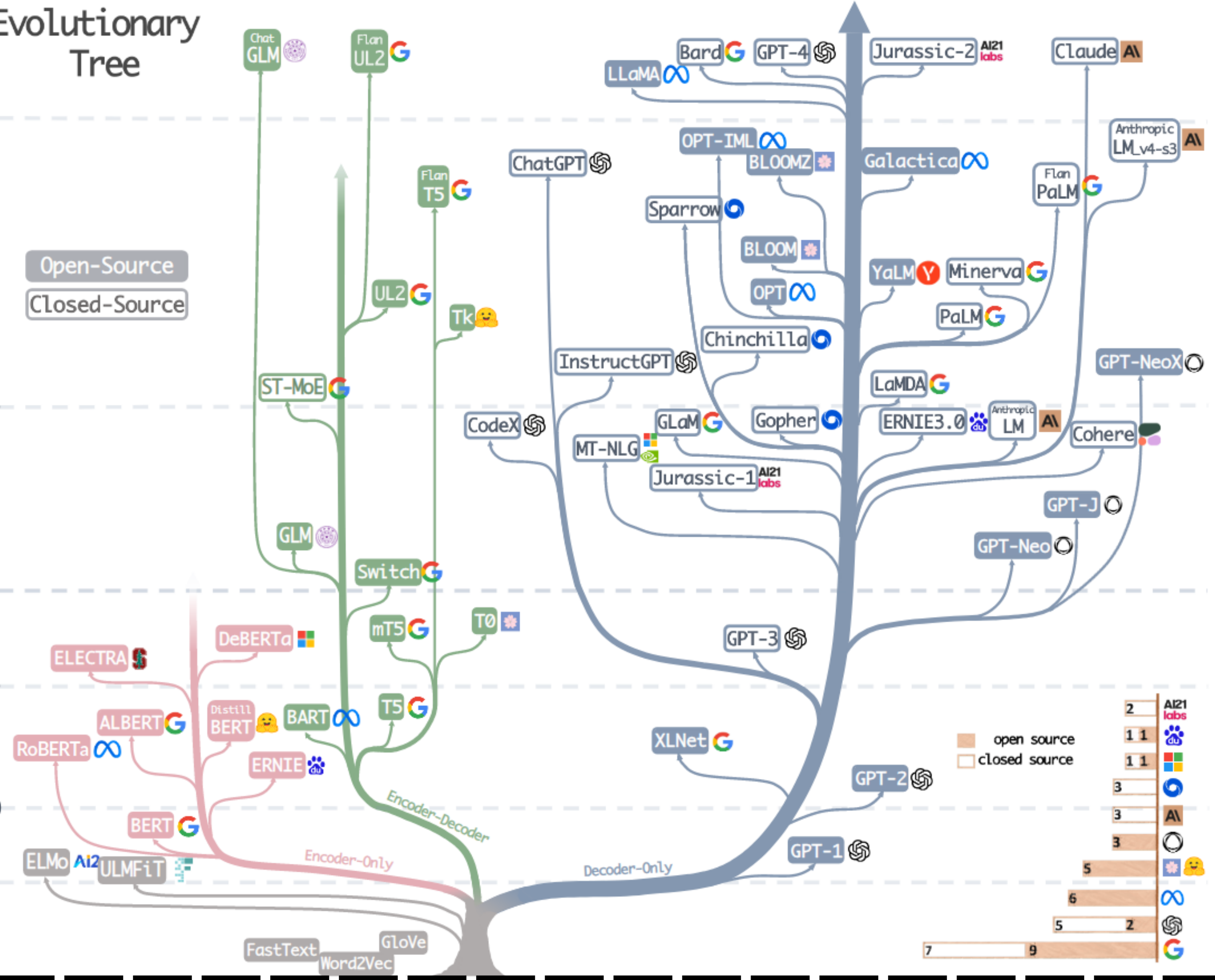
models are state-of-the-art in
Natural Language Processing
and Computational
Linguistics?



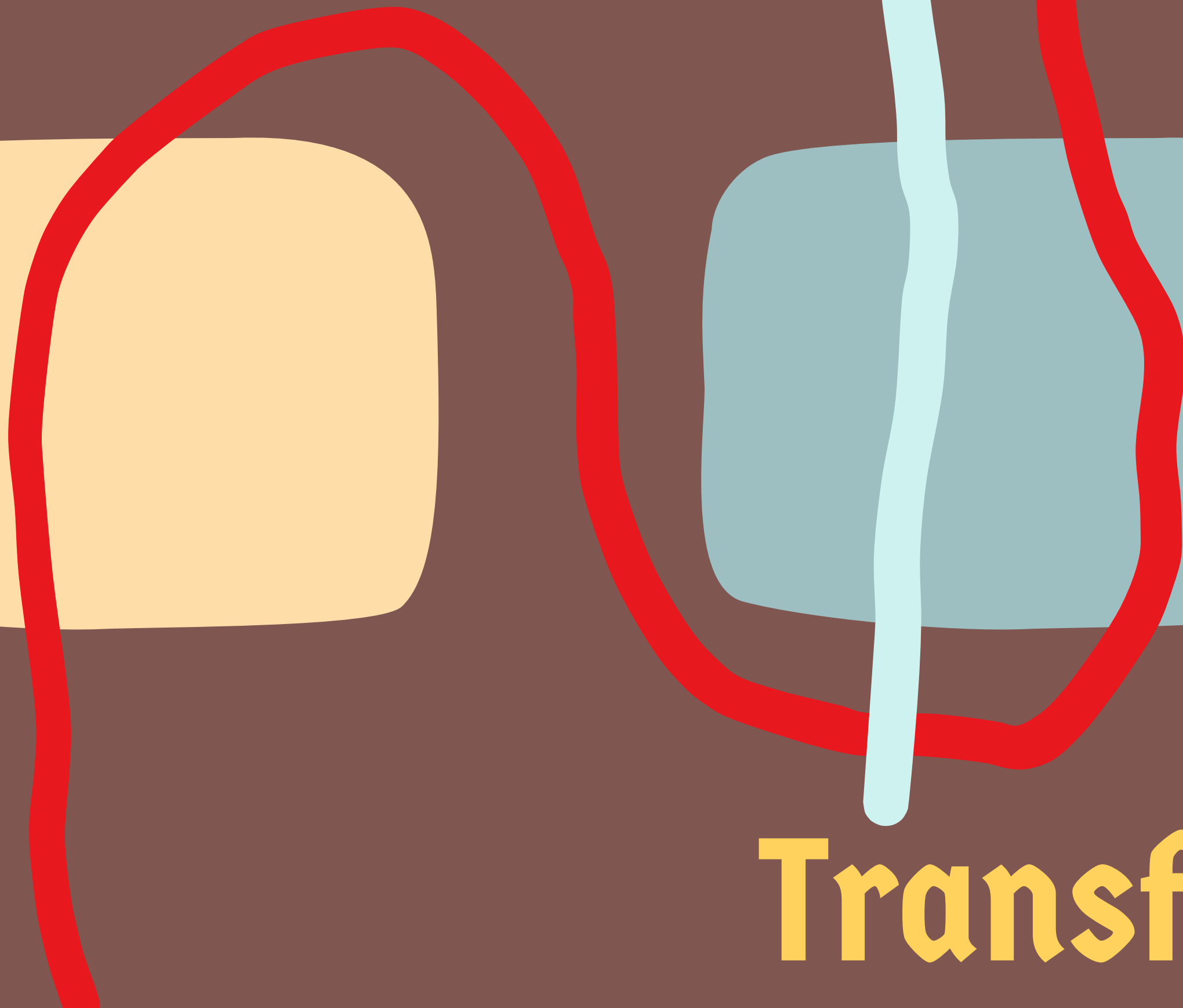
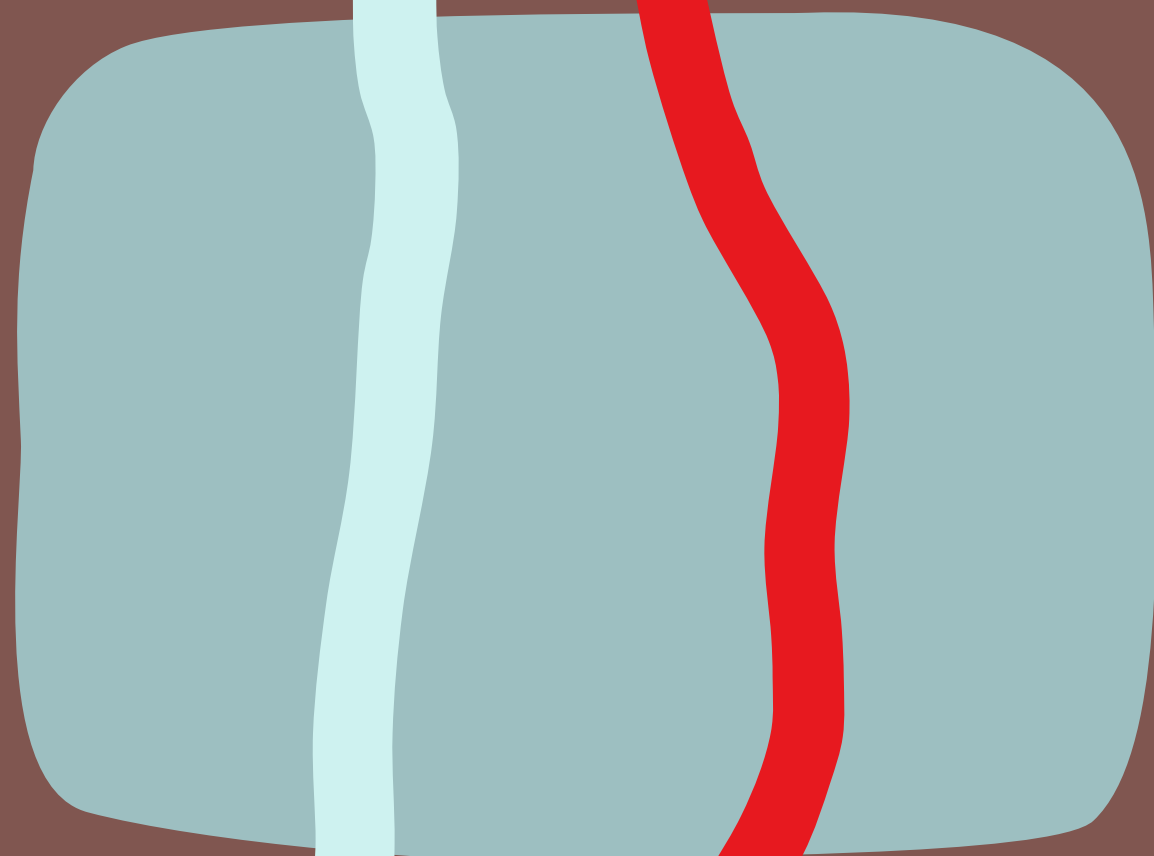
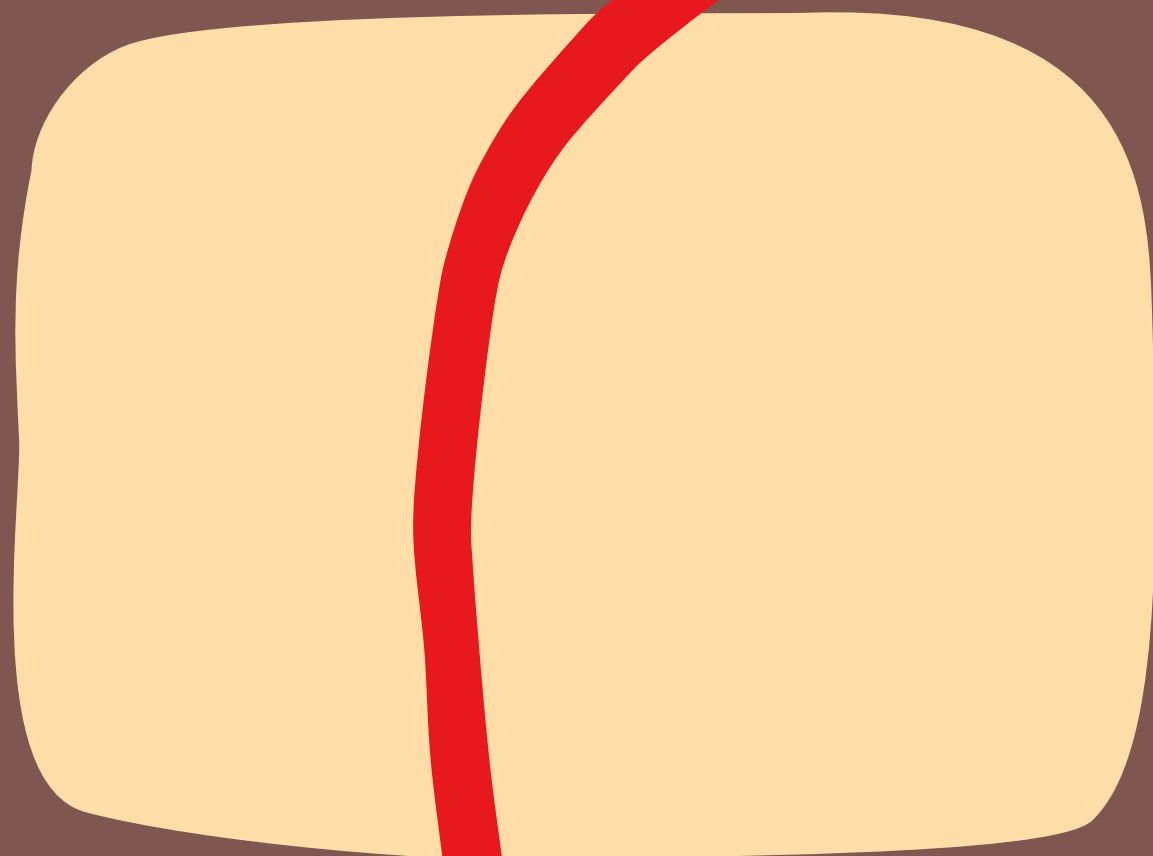
Evolutionary Tree

2023
2022
2021
2020
2019
2018

Open-Source
Closed-Source

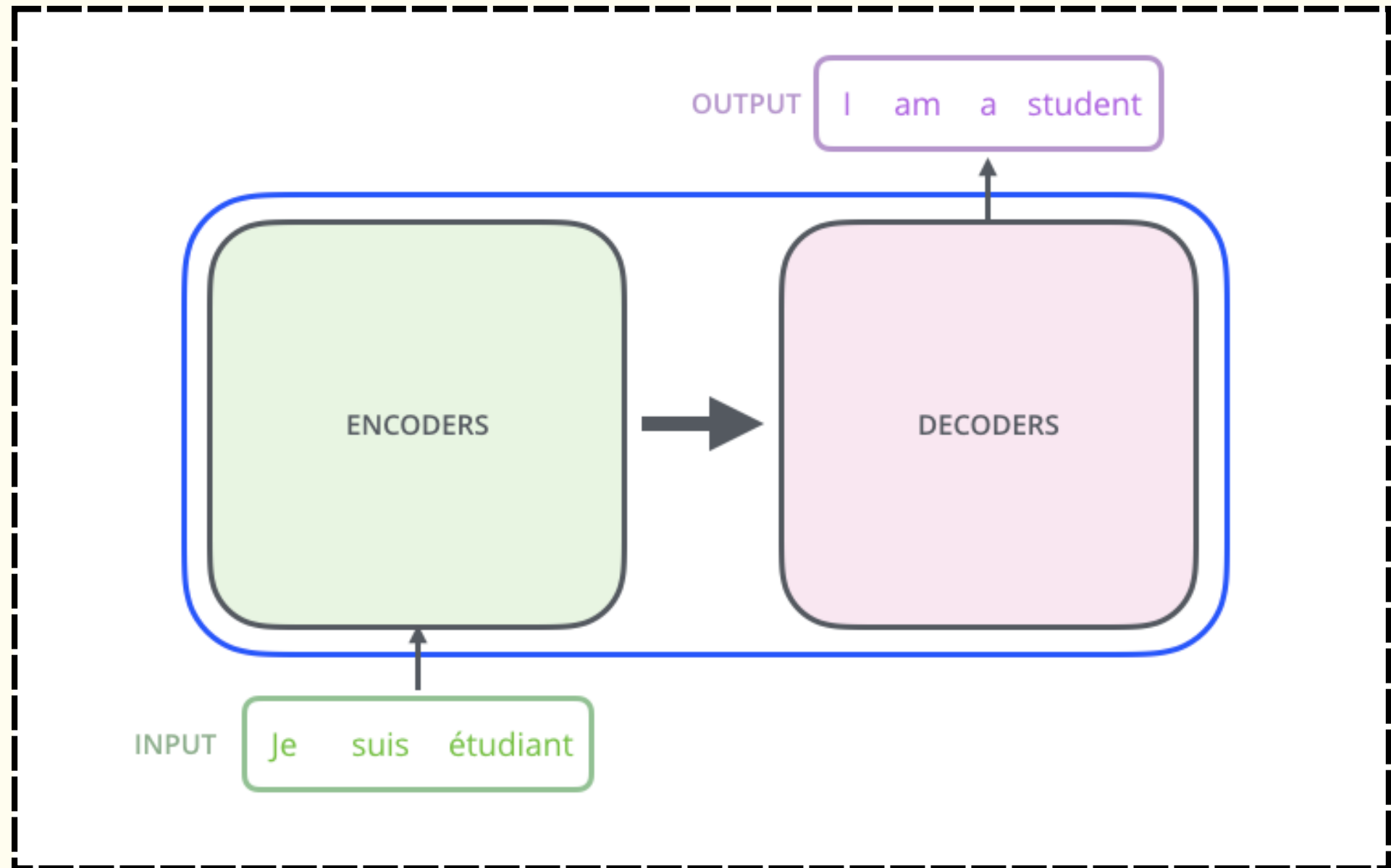


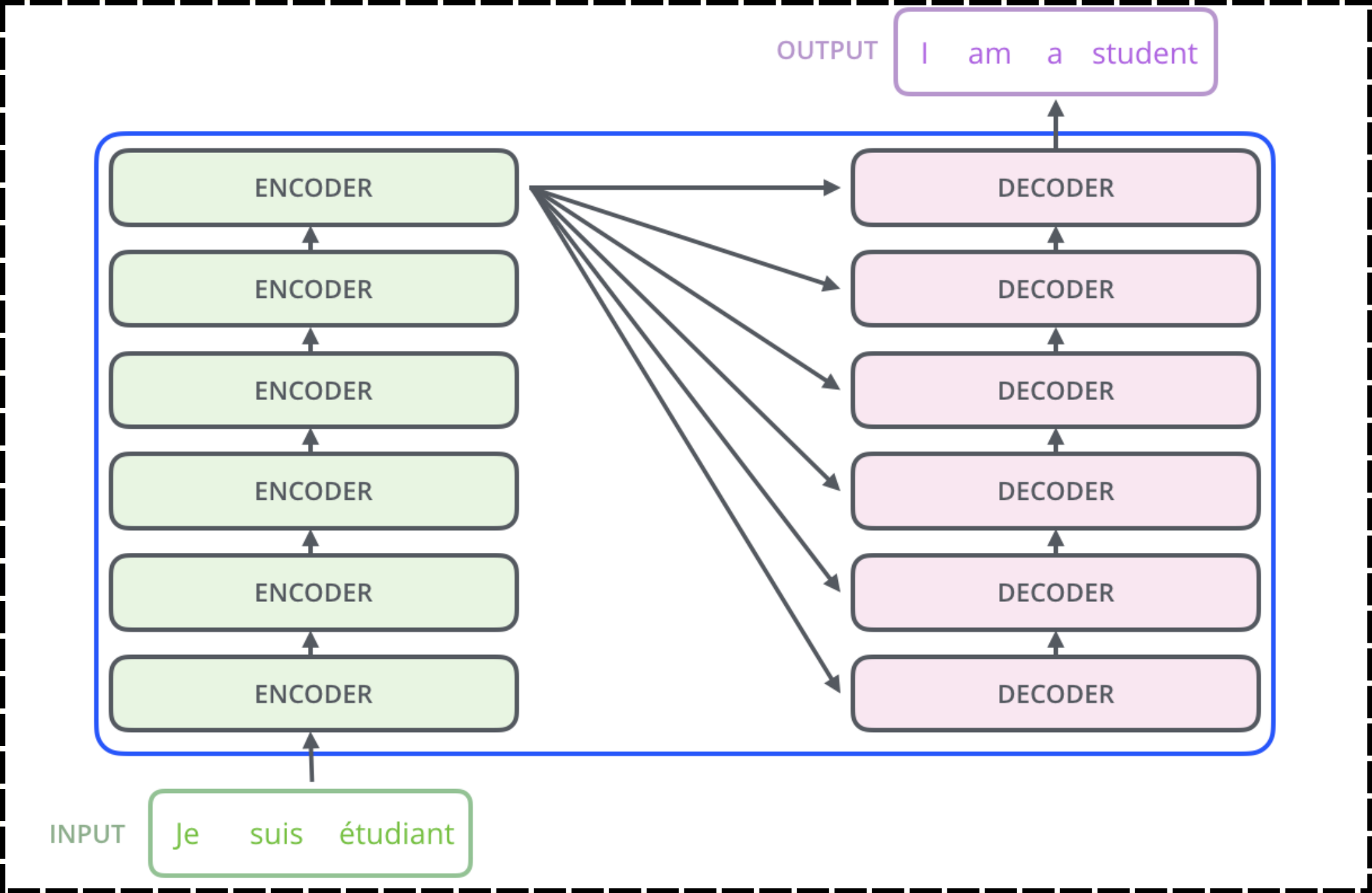
2	AI21 labs
1 1	Anthropic
1 1	Google
3	OpenAI
3	AI
3	OpenAI
5	Google
6	OpenAI
5	2
7	9

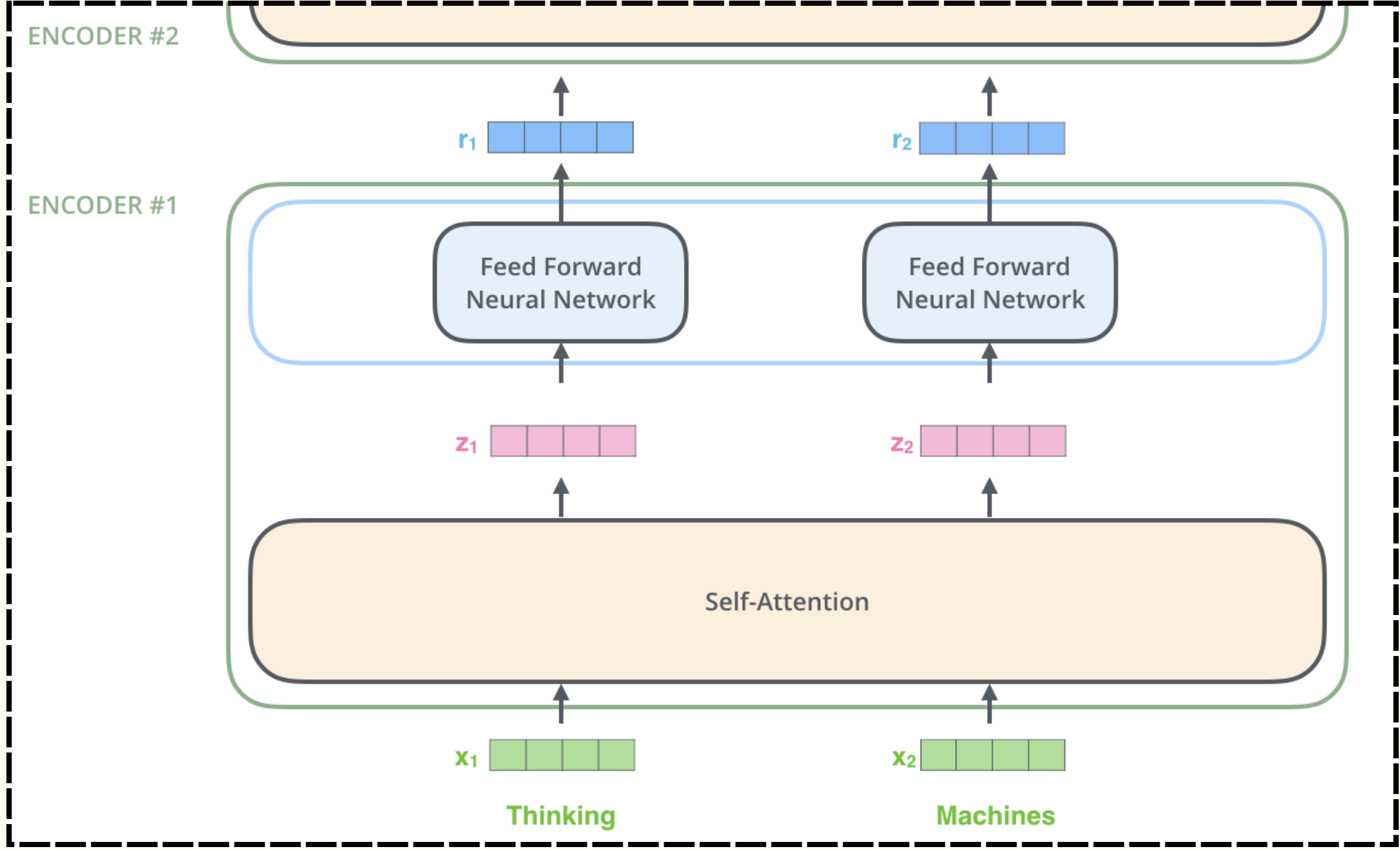


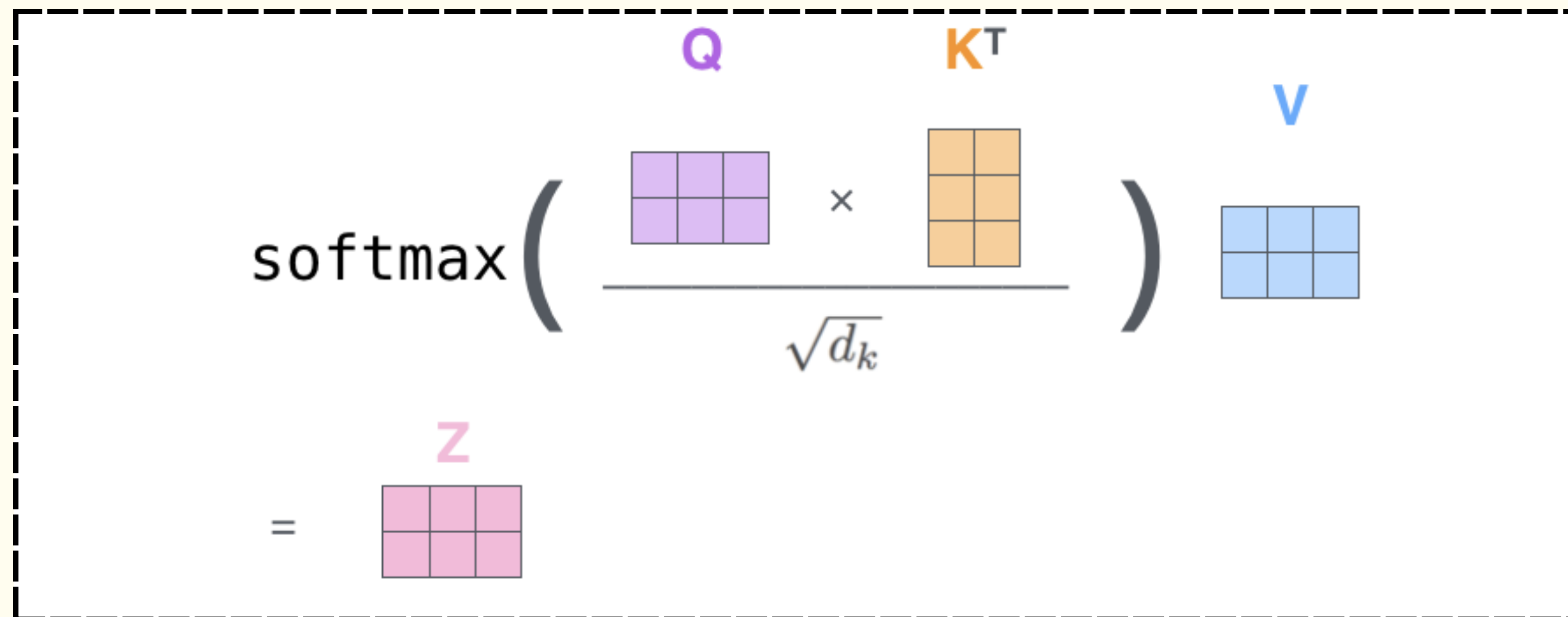
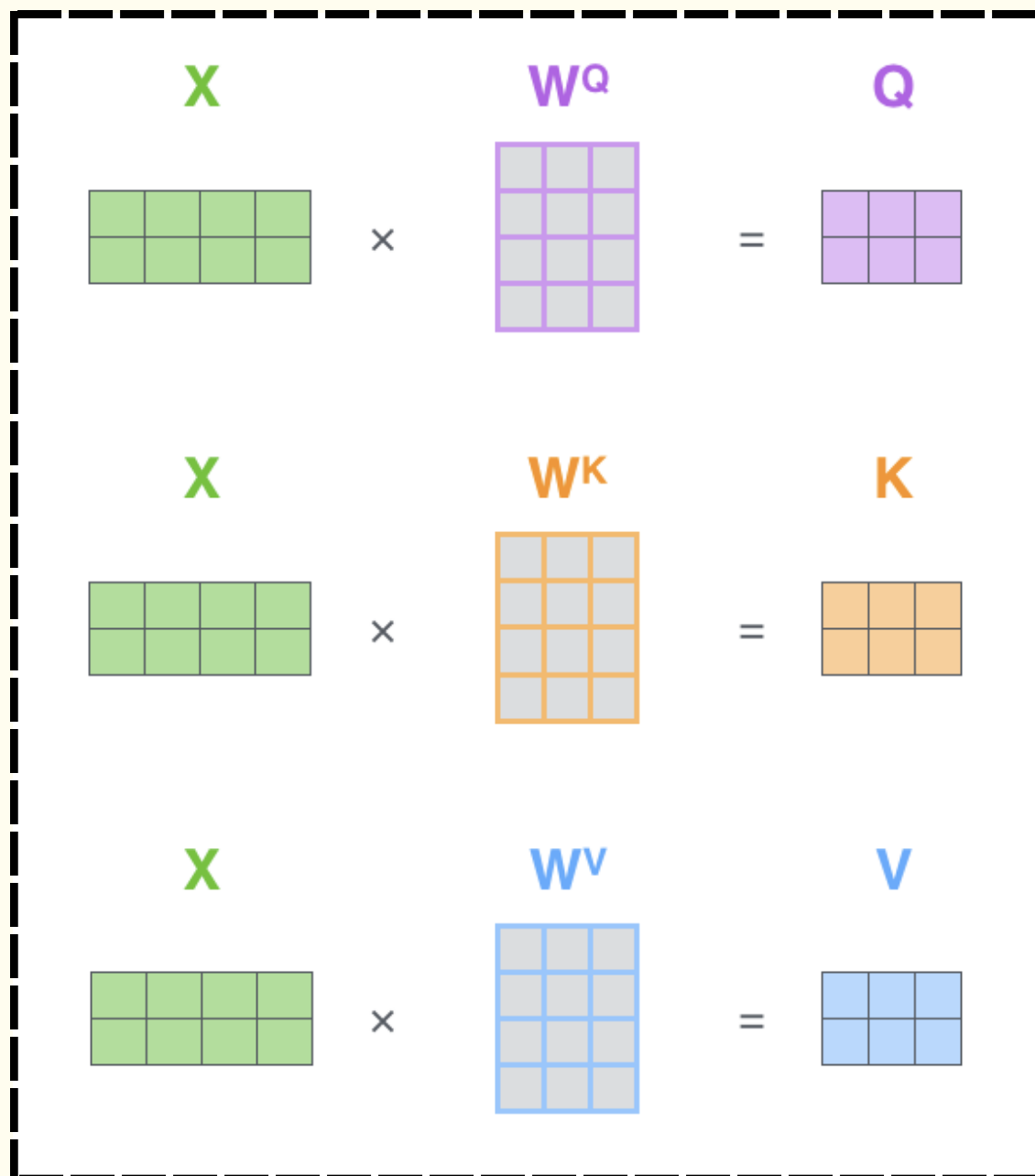
Transformer

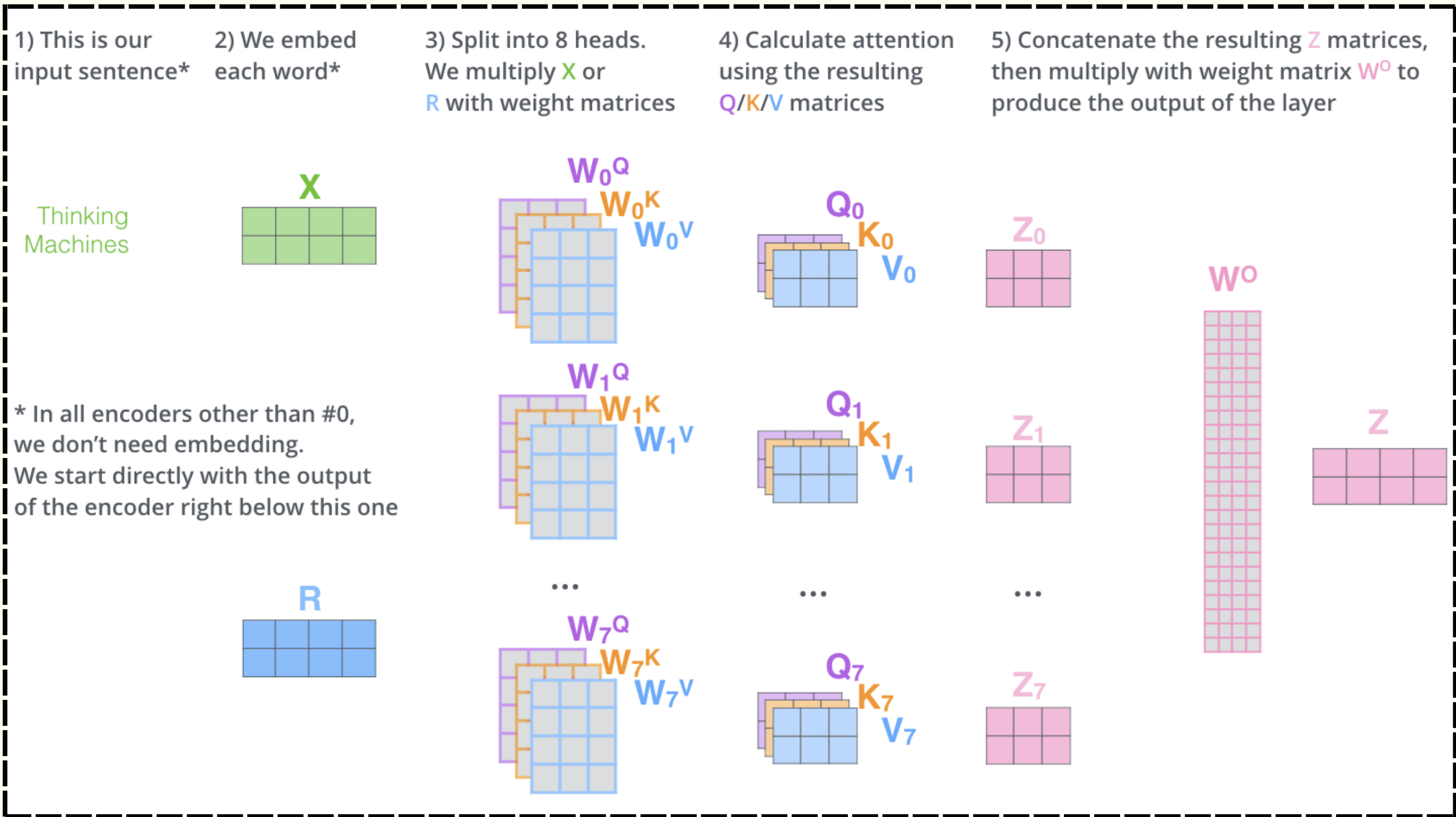


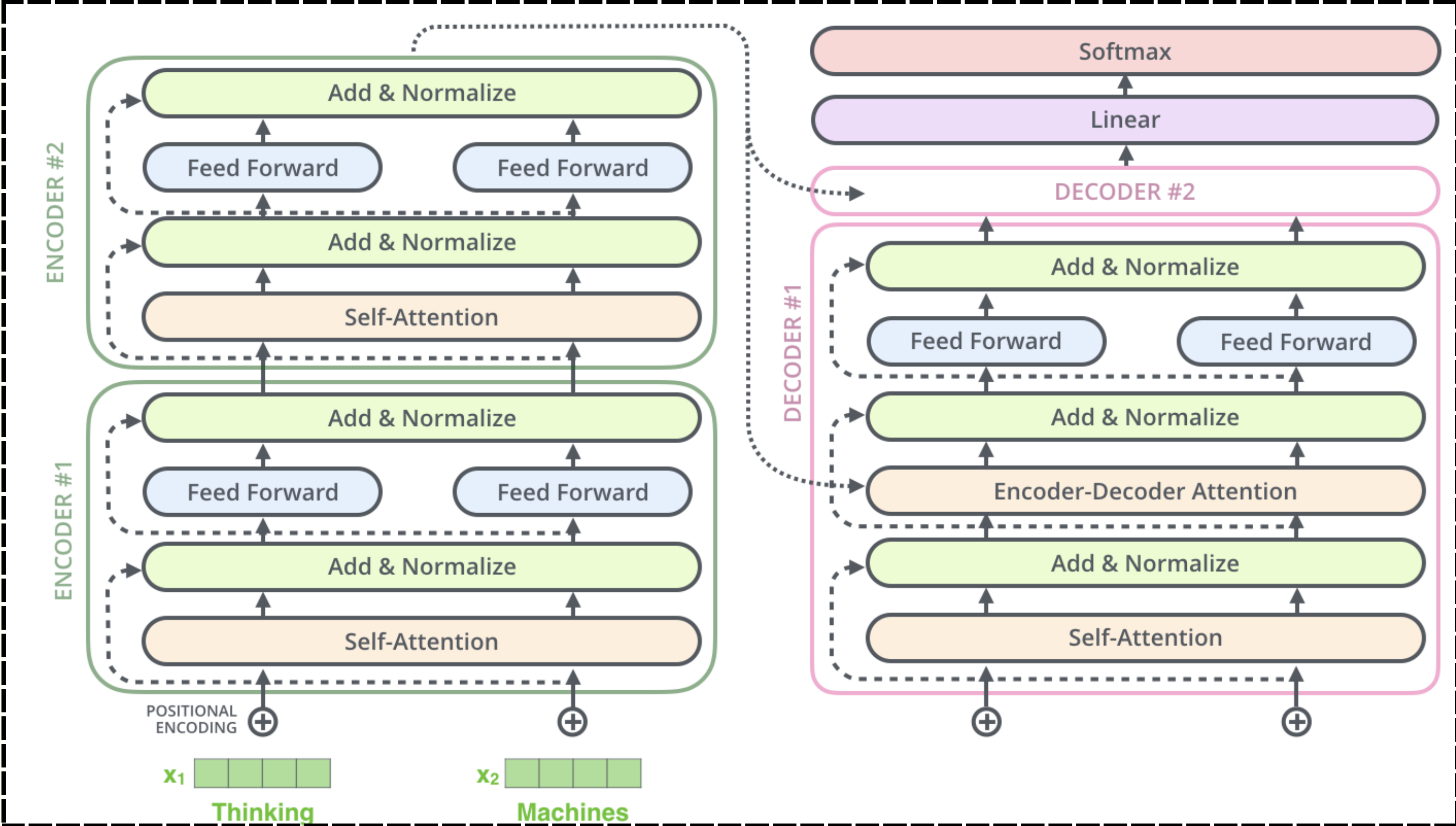


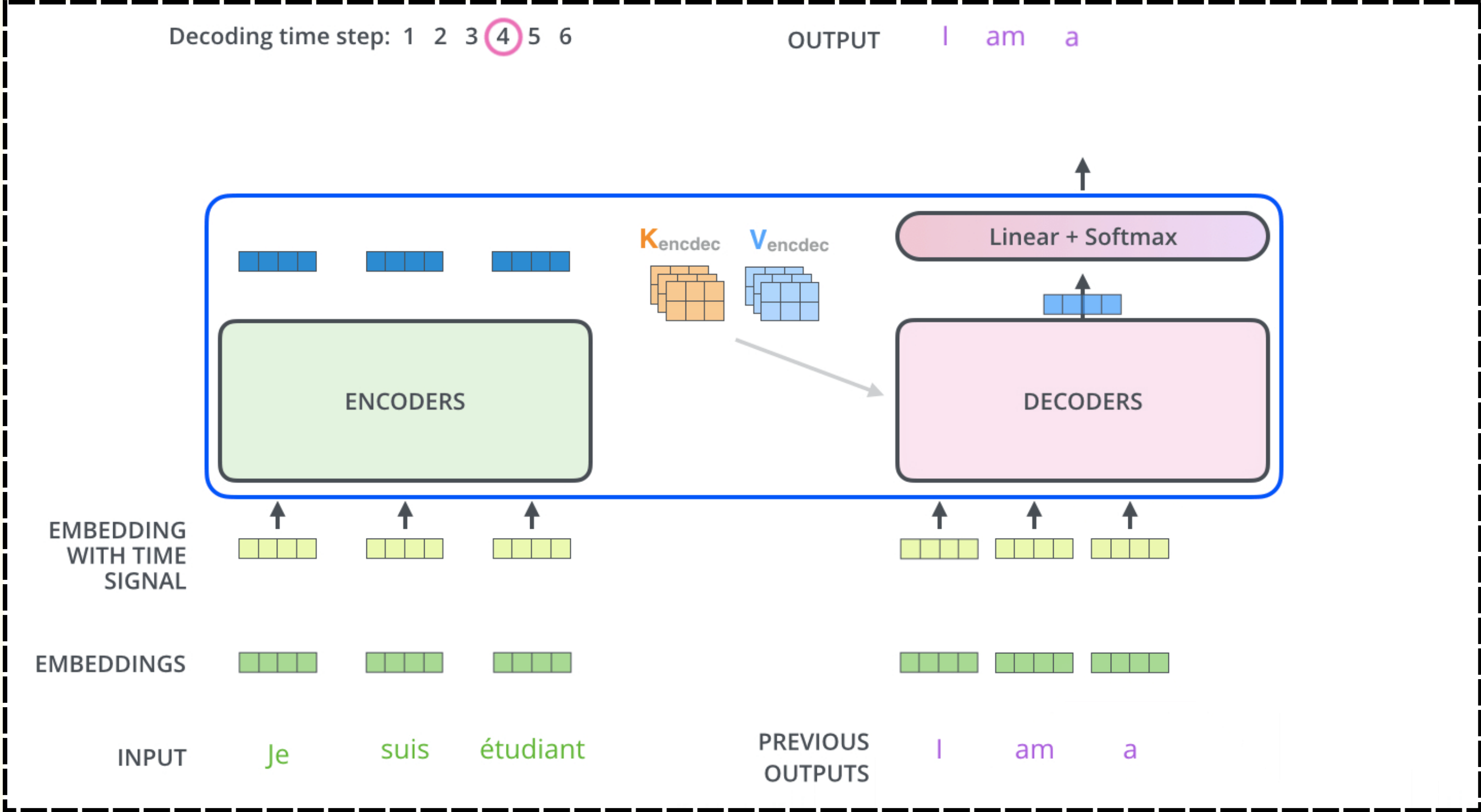


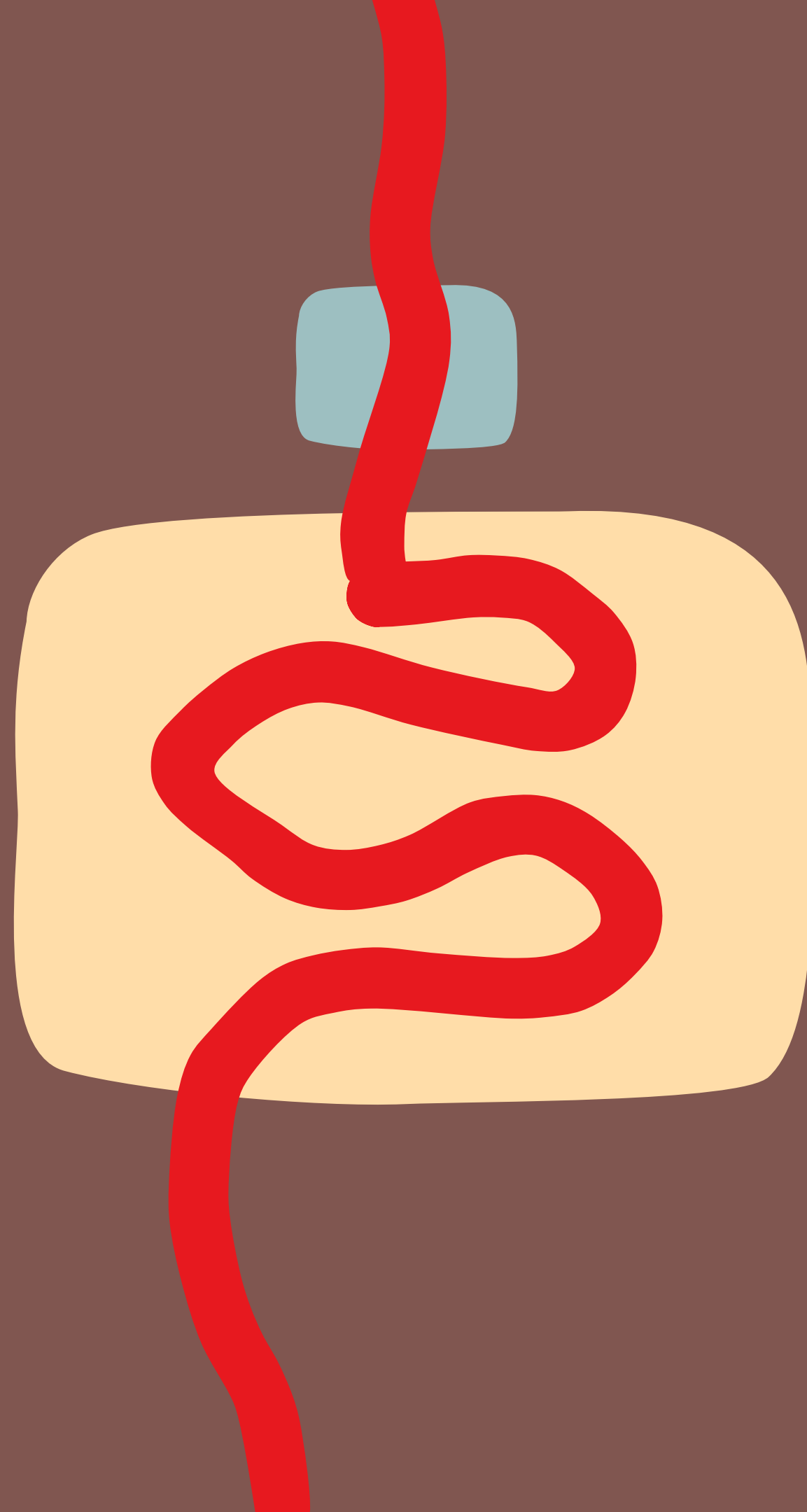




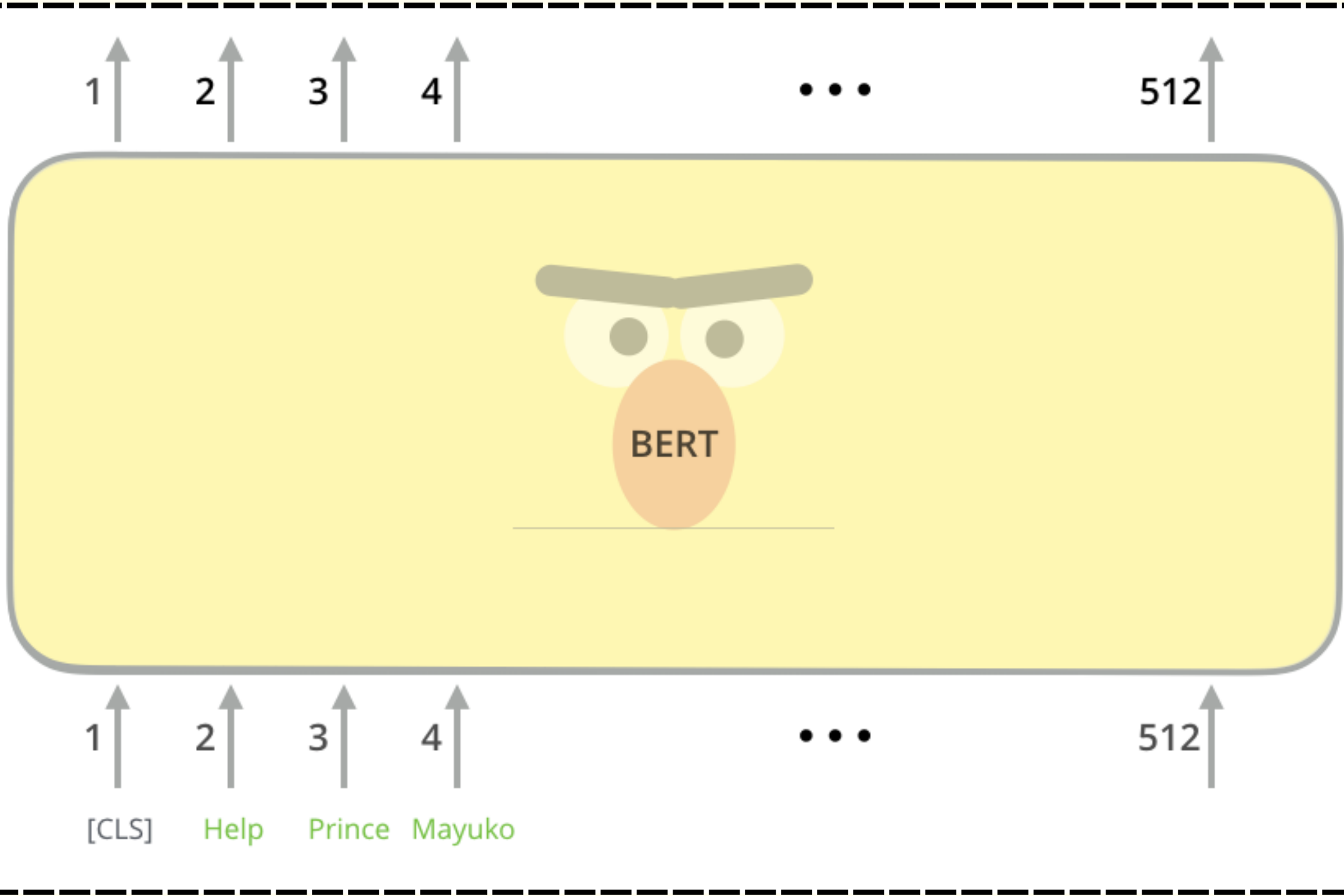


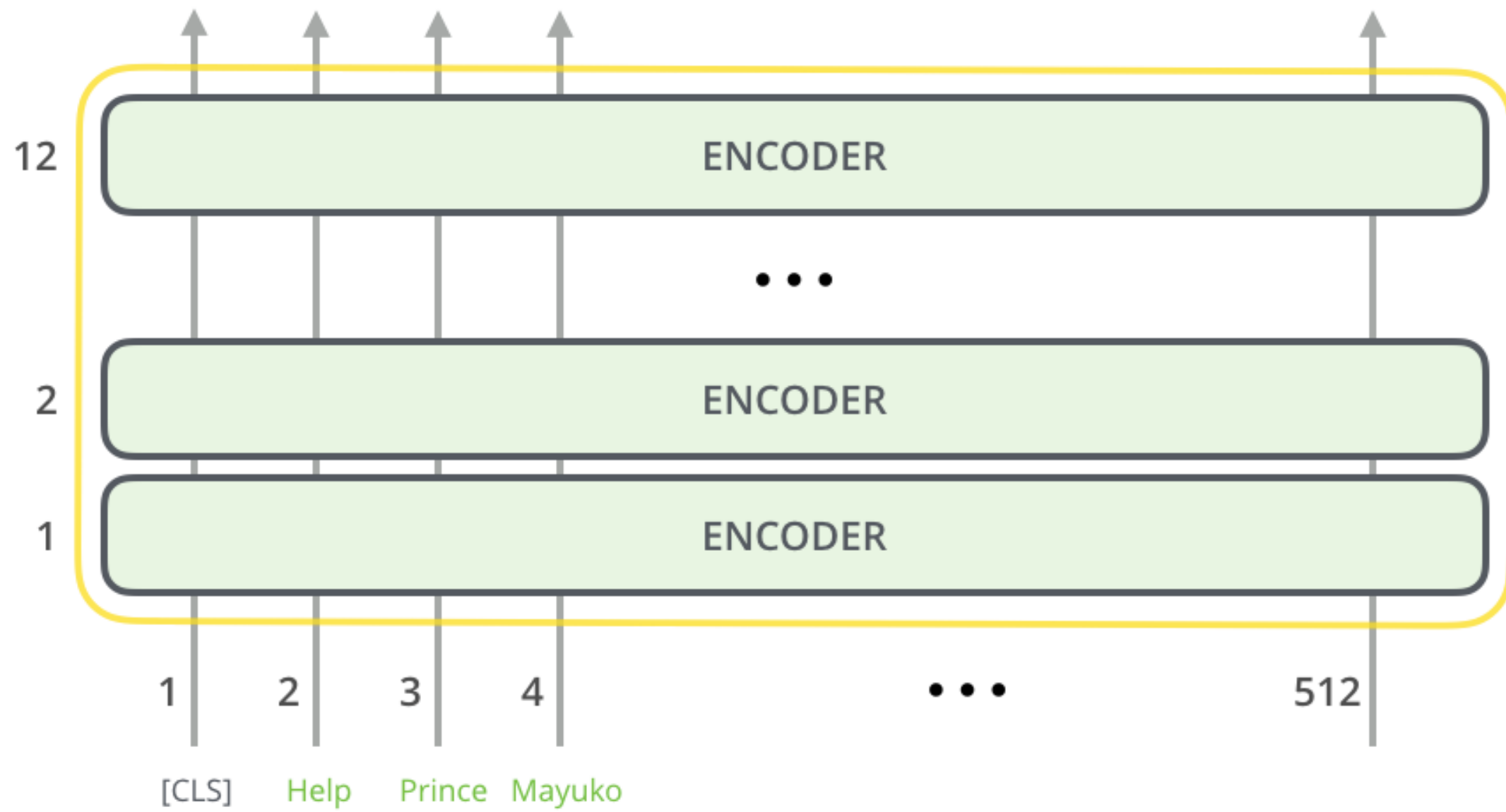




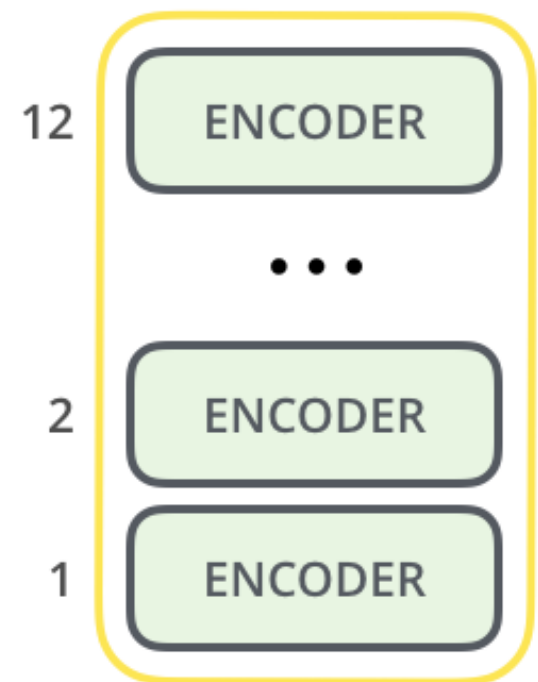


BERT

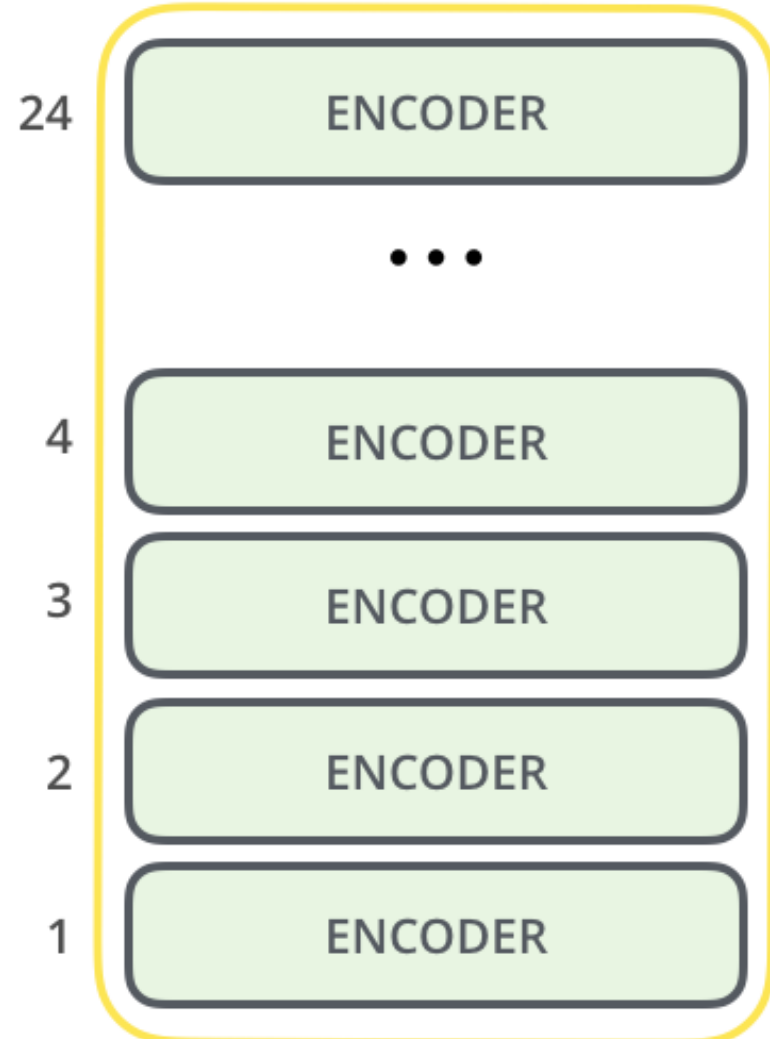




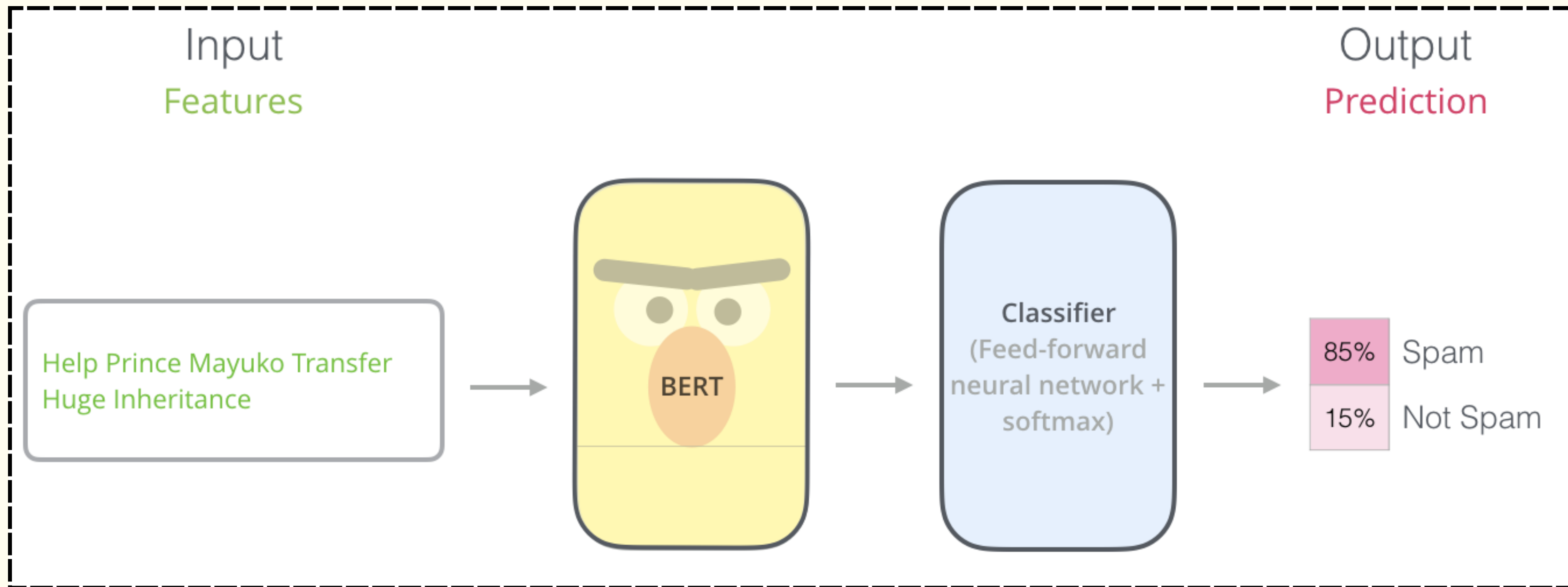
BERT



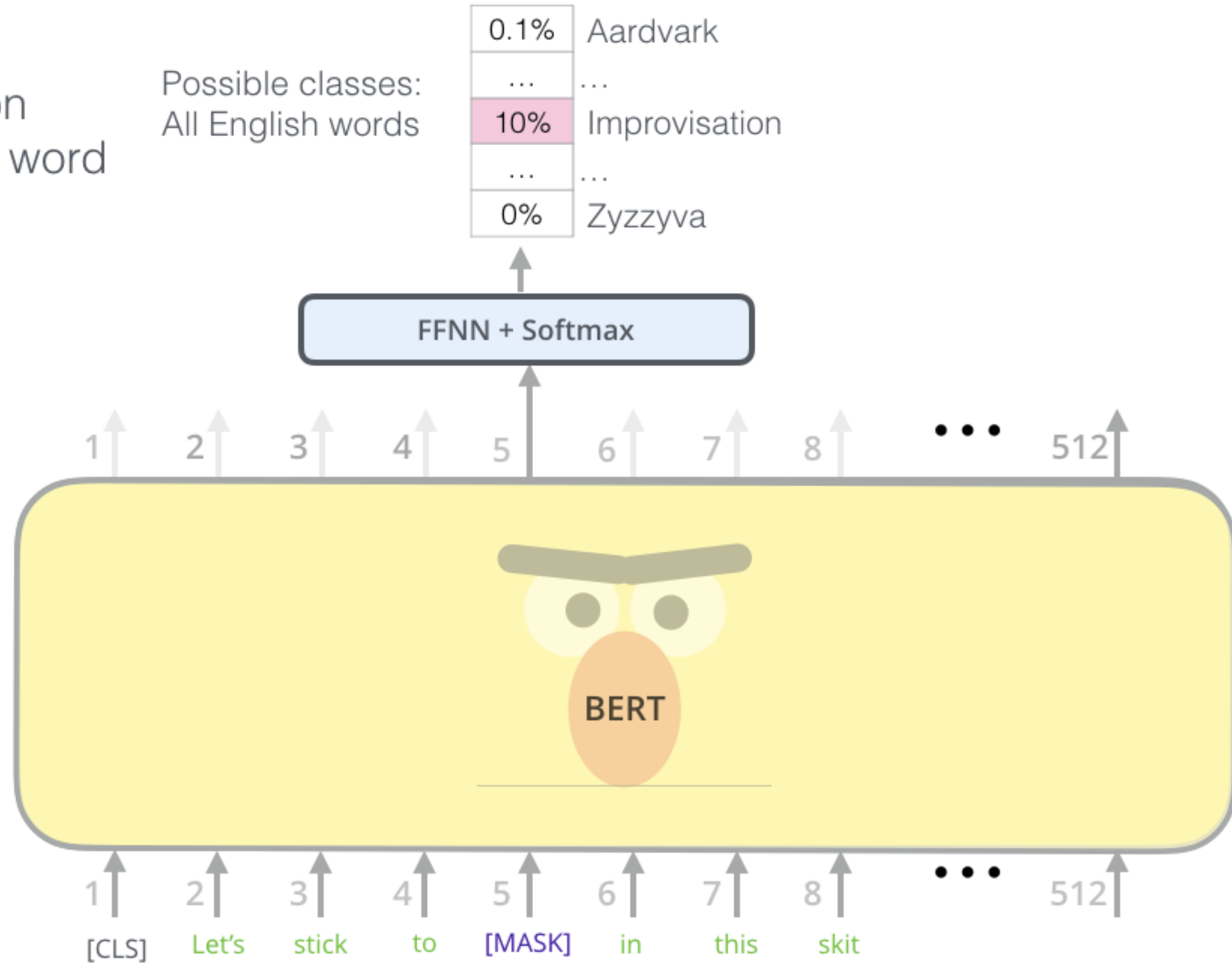
BERT_{BASE}



BERT_{LARGE}



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



Randomly mask 15% of tokens

Input

[CLS] Let's stick to improvisation in this skit

The end

El fin

O fim



Any Questions?

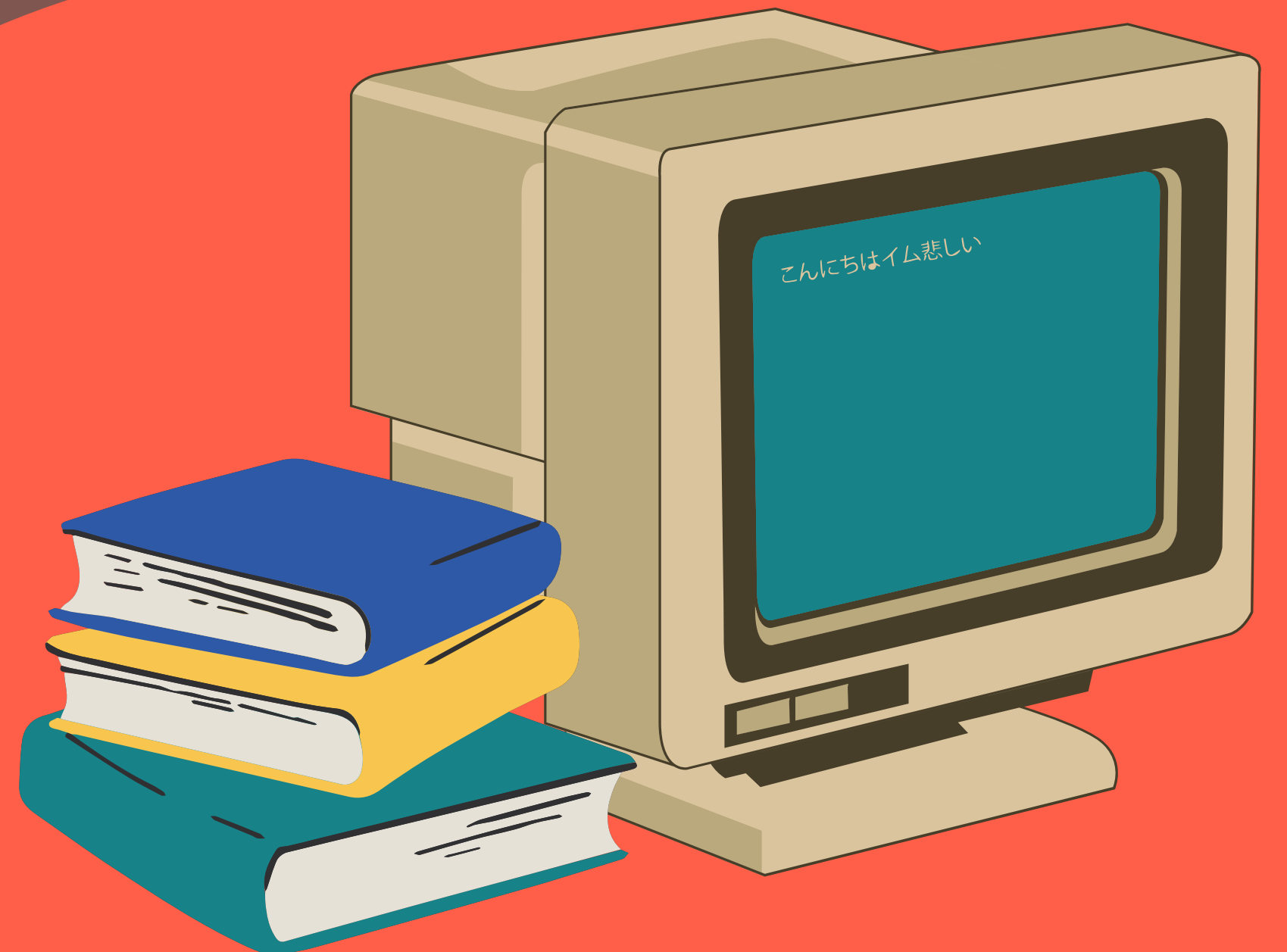


Image Sources

- 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 Alamar, under Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License
 - Alamar, J (2018). The Illustrated Transformer [Blog post]. Retrieved from <https://jalamar.github.io/illustrated-transformer/>
- Illustrated BERT, by Jay Alamar, under Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License
 - Alamar, J (2018). The Illustrated BERT, ELMo and co. [Blog post]. Retrieved from <http://jalamar.github.io/illustrated-bert/>

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