



Republic of Tunisia
Ministry of Higher Education and Scientific Research
University of Tunis El Manar
National School of Engineers of Tunis



Deep Learning Project

From Classical to Cutting-Edge: Exploring Different Approaches for Questions Classification

Realised by :

Trigui Hatem
Hassouna Malek
Med Rabii Baccari

Supervised by :

Ms. Linda Marrakchi
M. Mustapha BENHAJ MINIAOUI

Class : 3ATEL.DASEC



01

General Overview & Preliminary operations

02

Use classical Machine Learning
Approaches

03

Use Deep Learning Approaches

04

Use Transformer-Based
Approach

05

Conclusion

General Overview & Preliminary operations

Dataset - TREC

Train Set
5452 Rows

File	Question
train_set5.txt	DESC:scanner How did serfdom develop in and then leave Russia ?
	ENTY:cine What films featured the character Popeye Doyle ?
	DESC:scanner How can I find a list of celebrities ' real names ?
	ENTY:animal what fowl grabs the spotlight after the Chinese Year of the Monkey ?



Label

Test Set
500 Rows

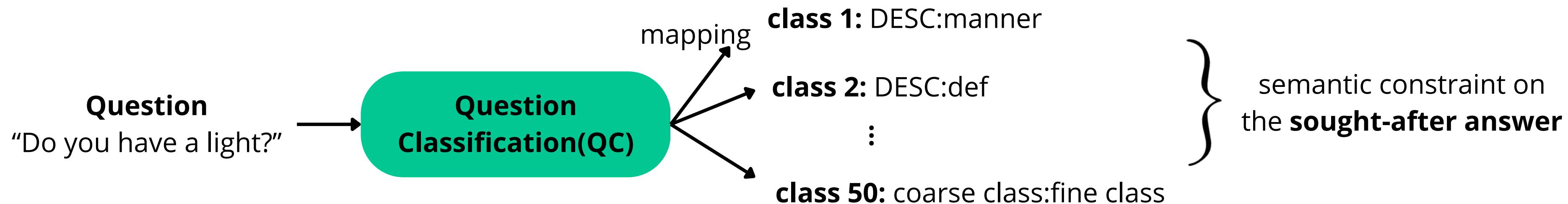
File	Question
test_set5.txt	NUM:dist How far is it from Denver to Aspen ?
	LOC:city What county is Modesto , California in ?
	HUM:desc Who was Galileo ?
	DESC:def what is an atom ?
	NUM:date when did hawaii become a state ?

6 coarse classes and 50 fine classes.

Class	#	Class	#
ABBREV	9	DESCRIPTION	7
abb	1	measure	3
exp	8	relation	6
ENTITY	64	HUMANS	65
animal	16	group	6
body	2	individual	23
color	16	title	1
creature	0	description	3
currency	0	LOCATION	41
dated	2	city	18
event	3	country	3
flood	4	mountain	3
instrument	1	other	26
language	2	state	7
letter	0	METRIC	11
other	12	code	6
plant	5	count	9
product	4	done	43
religion	0	distance	18
sport	1	military	3
substance	15	order	6
symbol	0	other	12
technology	1	period	8
term	7	present	5
vehicle	4	spend	6
word	0	temp	5
DESCRIPTIONS	13	size	6
definition	13	weight	4

The distribution of 500 TREC 10 (**test set**) questions over the question hierarchy. Coarse classes (in bold) are followed by their fine class refinements.

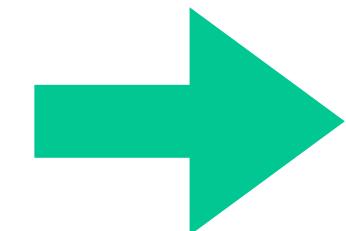
Purpose



The intension is that this we **do not** address questions like,
which calls for an action, but rather **only factual Wh-questions**.

Challenges : The Ambiguity Problem

1. What is bipolar disorder?
 - **definition**
 - **disease medicine**
2. What do bats eat?
 - **food**
 - **plant**
 - **animal**
3. What is the PH scale?
 - **numeric value**
 - **definition**



Multi-Label Classification

we allow our classifiers to assign multiple class labels for a single question.

Pre-processing

Question

How can I find a list of celebrities' real names?

1. **lower case**

how can I find a list of celebrities' real names?

Use **regex**

2. **Remove non-alphabetic characters**

how can I find a list of celebrities real names

Use **word_tokenize** from nltk.tokenize

3. **Tokenize**

[**'how'**, '**'can'**', '**'I'**,
'find'', '**'a'**', '**'list'**', '**'of'**',
'celebrities'', '**'real'**',
'names']

Use **stopwords** from nltk.corpus - {"**what**", "**where**", "**who**", "**when**", "**why**", "**how**"}

4. **remove stop words excluding wh questions**

[**'how'**, '**'find'**', '**'list'**',
'celebrities'', '**'real'**',
'names']

Use **WordNetLemmatizer** from nltk.stem

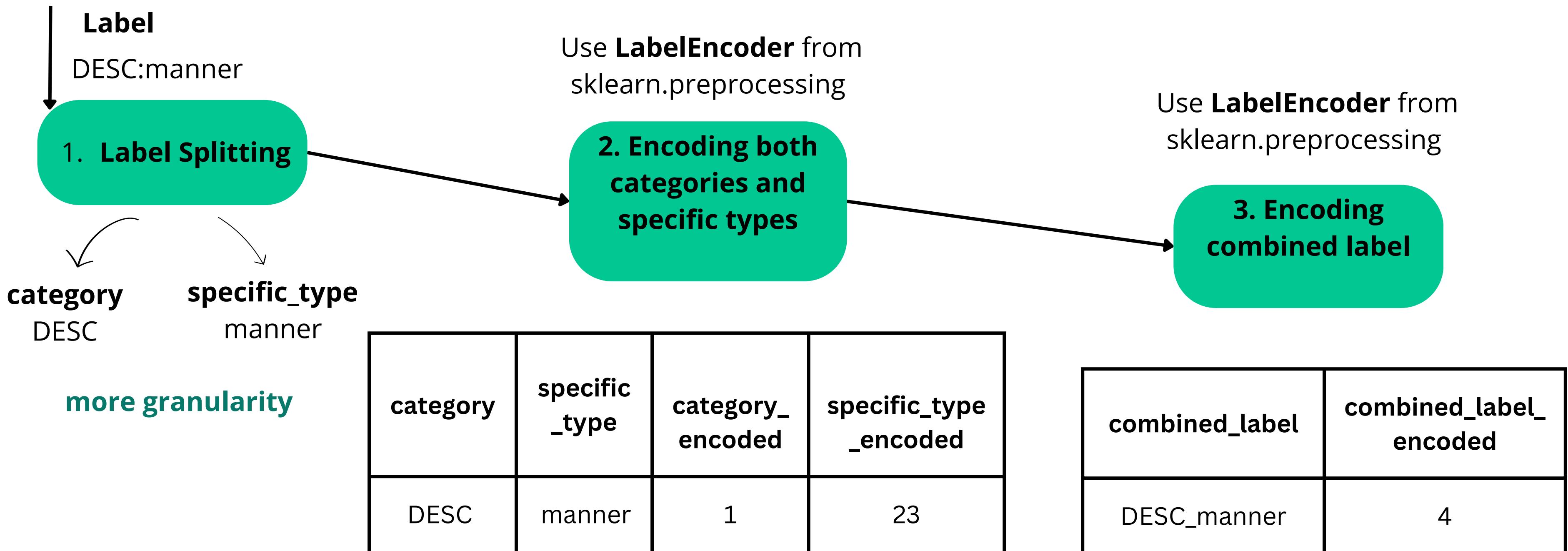
5. **Lemmatization**

[**'how'**, '**'find'**', '**'list'**',
'celebrity', '**'real'**',
'name']

Join

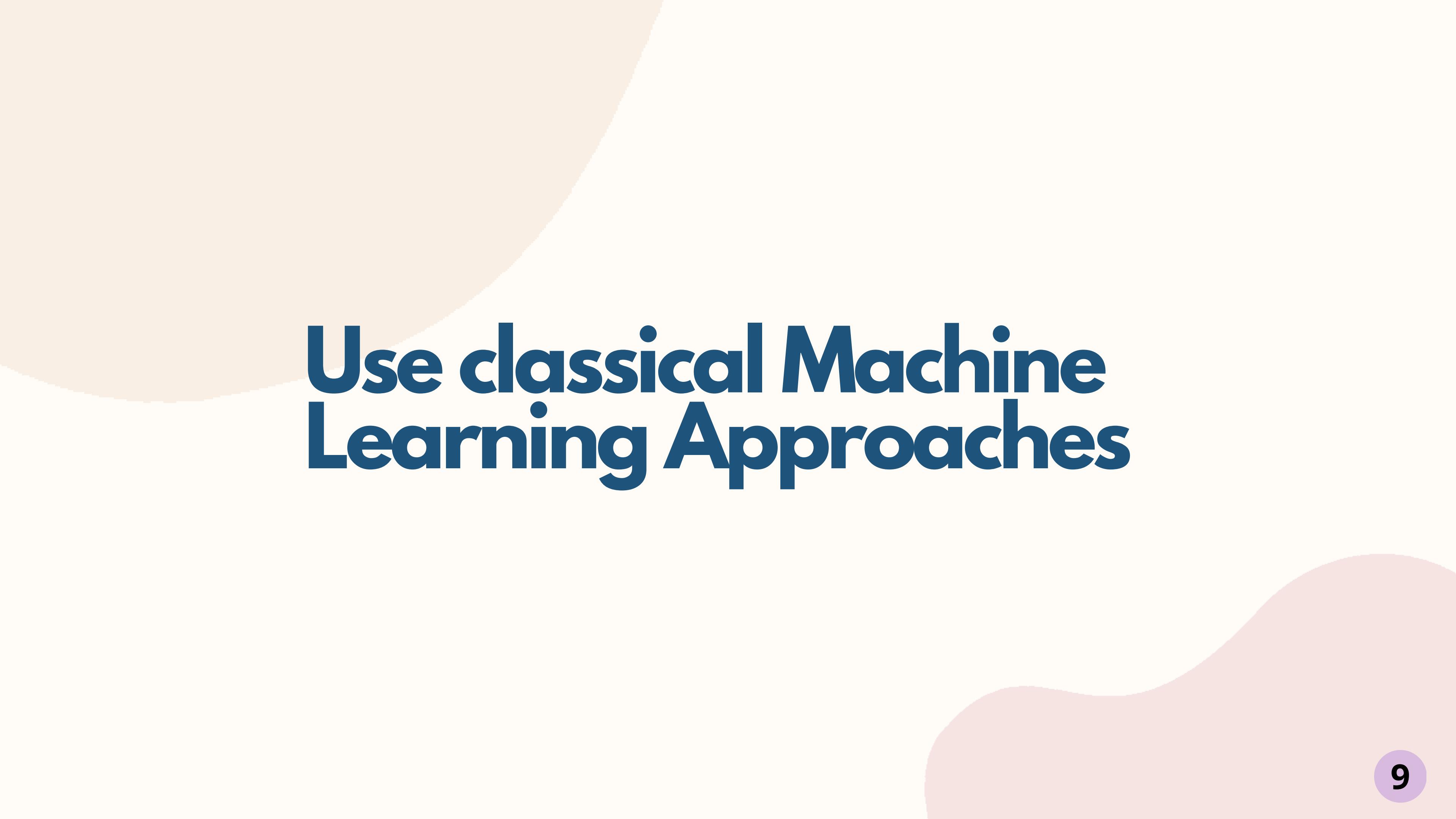
how find list celebrity real name

Embedding



Train Set Overview (after preprocessing + encoding)

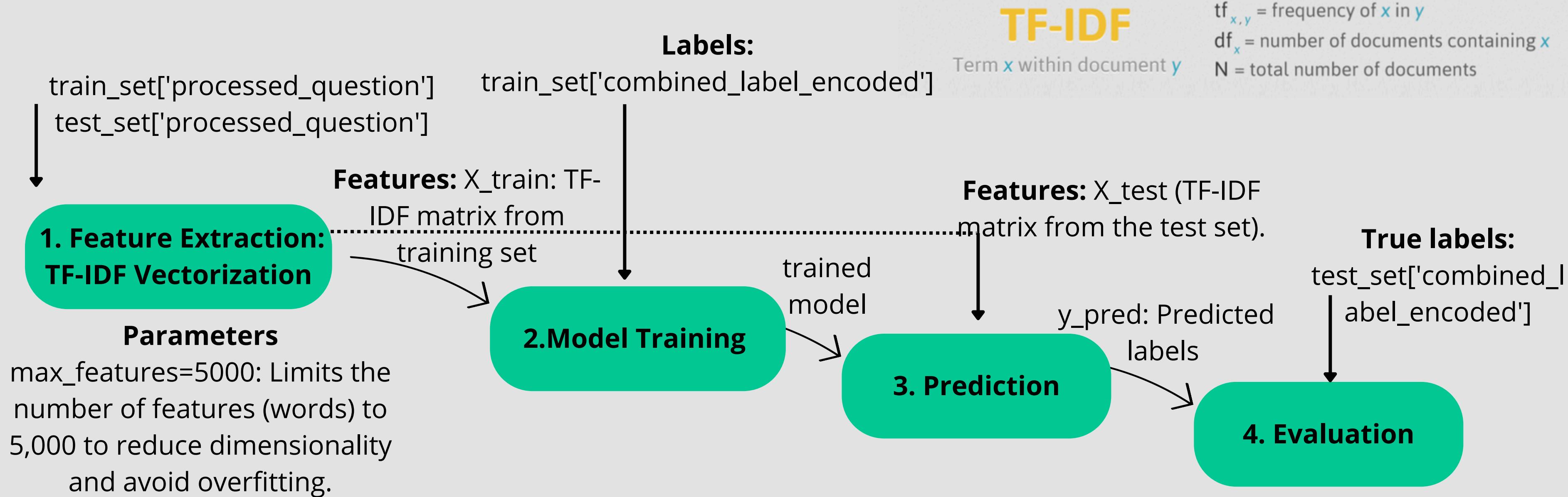
	question	label	processed_question	category	specific_type	category_encoded	specific_type_encoded	combined_label	combined_label_encoded
0	How did serfdom develop in and then leave Russ...	DESC:manner	how serfdom develop leave russia	DESC	manner	1	23	DESC_manner	4
1	What films featured the character Popeye Doyle ?	ENTY:cremat	what film featured character popeye doyle	ENTY	cremat	2	8	ENTY_creatmat	9
2	What fowl grabs the spotlight after the Chines...	ENTY:animal	what fowl grab spotlight chinese year monkey	ENTY	animal	2	1	ENTY_animal	6



Use classical Machine Learning Approaches

I. First Word Representation Approach (TF-IDF)

General Process



I. First Word Representation Approach (TF-IDF)

Classification Report

	precision	recall	f1-score	support
0	0.00	0.00	0.00	1
1	1.00	0.62	0.77	8
2	0.74	0.98	0.85	123
49	0.00	0.00	0.00	4
Accuracy			0.71	500
macro avg	0.53	0.45	0.45	500
weighted avg	0.70	0.71	0.66	500

50 classes

Since our data is **imbalanced** so we'll consider weighted avg values weighing each class's contribution by its support (i.e., the number of true instances in that class)

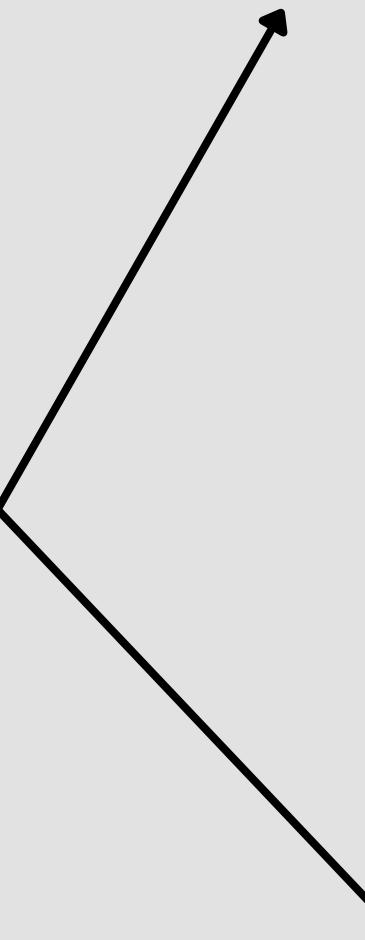
Logistic Regression model

I. First Word Representation Approach (TF-IDF)

Imbalanced Data ?

category	specific_type	count
ENTY 1250	ind	962
HUM 1223	other	733
DESC 1162	def	421
NUM 896	count	363
LOC 835	desc	321
ABBR 86	manner	276
	date	218
	...	
	letter	9
	code	9
	speed	9
	temp	8
	ord	6
	religion	4
	currency	4

1. Class Imbalance:
Minority classes may lead to poor model performance for those classes.



2. Granularity:
Some specific types, such as techmeth and volsize, are too granular, which may increase the complexity of classification.

I. First Word Representation Approach (TF-IDF)

First Solution (Oversampling)

- X_train: The feature matrix of the training set (inputs for the model).
- y_train: train_set['combined_label_encoded']

1. Initializing the SMOTE Object

k_neighbors=3: NNN to use when generating synthetic samples for the minority class.

smote object

2. Applying SMOTE to the Training Data

For each minority class sample, 3 nearest neighbors will be found, and synthetic samples will be created based on these neighbors.

X_train_balanced
y_train_balanced

equal number of samples for each class (or at least closer to equal)

I. First Word Representation Approach (TF-IDF)

First Solution (Oversampling)

Classification Report

	precision	recall	f1-score	support
0	1.00	1.00	1.00	1
1	0.67	0.75	0.71	8
2	0.88	0.76	0.82	123
49	1.00	1.00	1.00	4
Accuracy			0.76	500
macro avg	0.68	0.72	0.67	500
weighted avg	0.82	0.76	0.77	500

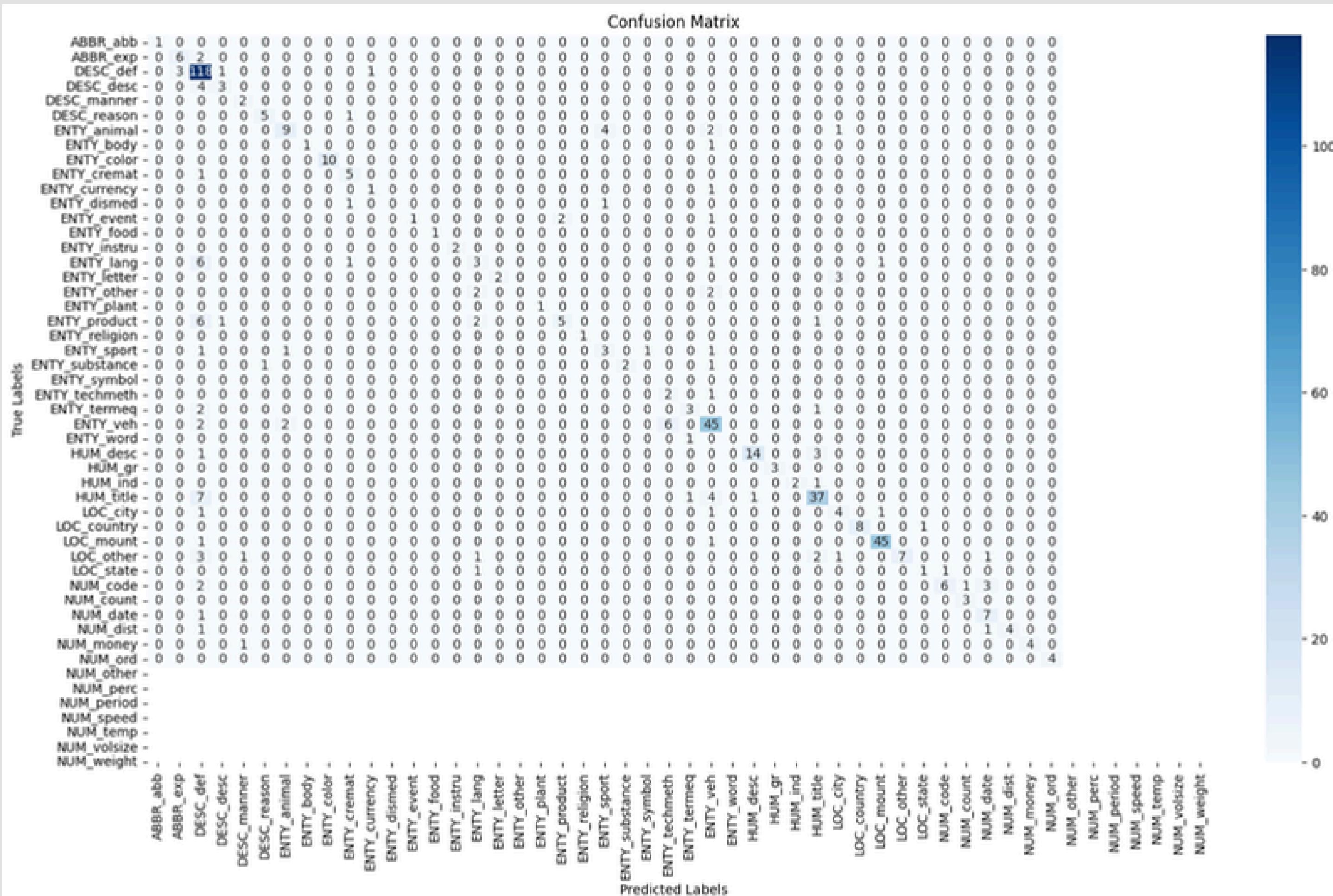
Logistic Regression model

Classification Report

	precision	recall	f1-score	support
0	1.00	1.00	1.00	1
1	0.67	0.75	0.71	8
2	0.74	0.96	0.84	123
49	1.00	1.00	1.00	4
Accuracy			0.76	500
macro avg	0.72	0.64	0.66	500
weighted avg	0.77	0.76	0.75	500

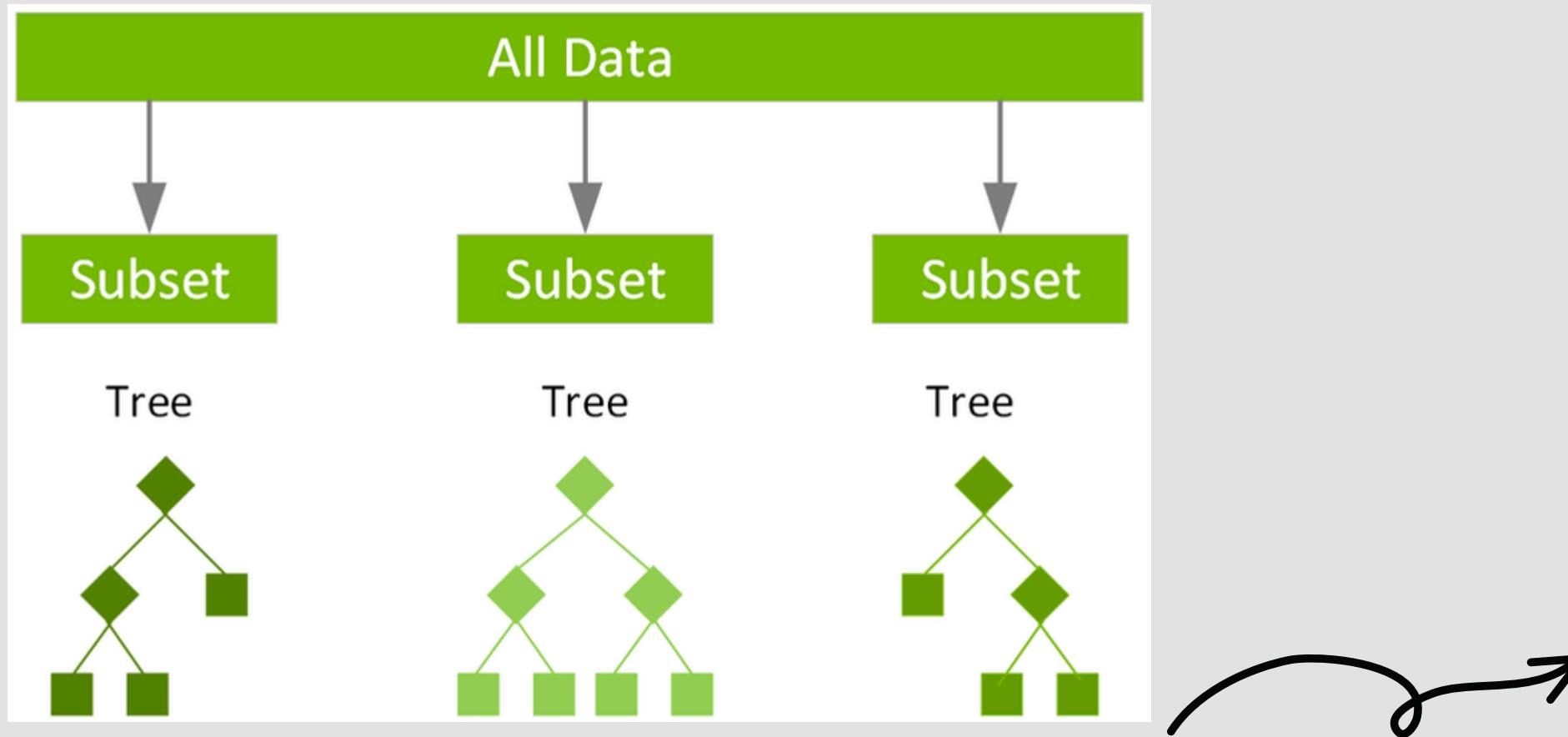
Random Forest model

I. First Word Representation Approach (TF-IDF)



I. First Word Representation Approach (TF-IDF)

Second Solution (Use Imbalanced data supportive model “XGboost”)



- XGBoost assigns higher weights to the minority class using the **scale_pos_weight** parameter.
- Decision trees split the data to maximize information gain, and XGBoost focuses on reducing errors iteratively.

Classification Report

	precision	recall	f1-score	support
0	0.00	0.00	0.00	1
1	0.75	0.75	0.75	8
2	0.75	0.99	0.86	123
49	1.00	0.75	0.86	4
Accuracy			0.76	500
macro avg	0.66	0.59	0.59	500
weighted avg	0.78	0.77	0.75	500

XGBoost model

I. First Word Representation Approach (TF-IDF)

XGBoost's hyperparameters tuning

1. Define Parameter Grid

- **max_depth**: [3, 6, 10]
- **learning_rate**: [0.01, 0.1, 0.3]
- **n_estimators**: [100, 200, 500]

2. Perform Grid Search (cv=3)

3. Select the Best Model

Best Parameters: {'learning_rate': 0.1, 'max_depth': 6, 'n_estimators': 200}

	precision	recall	f1-score	support
0	0.00	0.00	0.00	1
1	0.75	0.75	0.75	8
2	0.73	0.99	0.86	123
49	1.00	0.75	0.86	4
Accuracy			0.76	500
macro avg	0.66	0.57	0.86	500
weighted avg	0.78	0.76	0.75	500

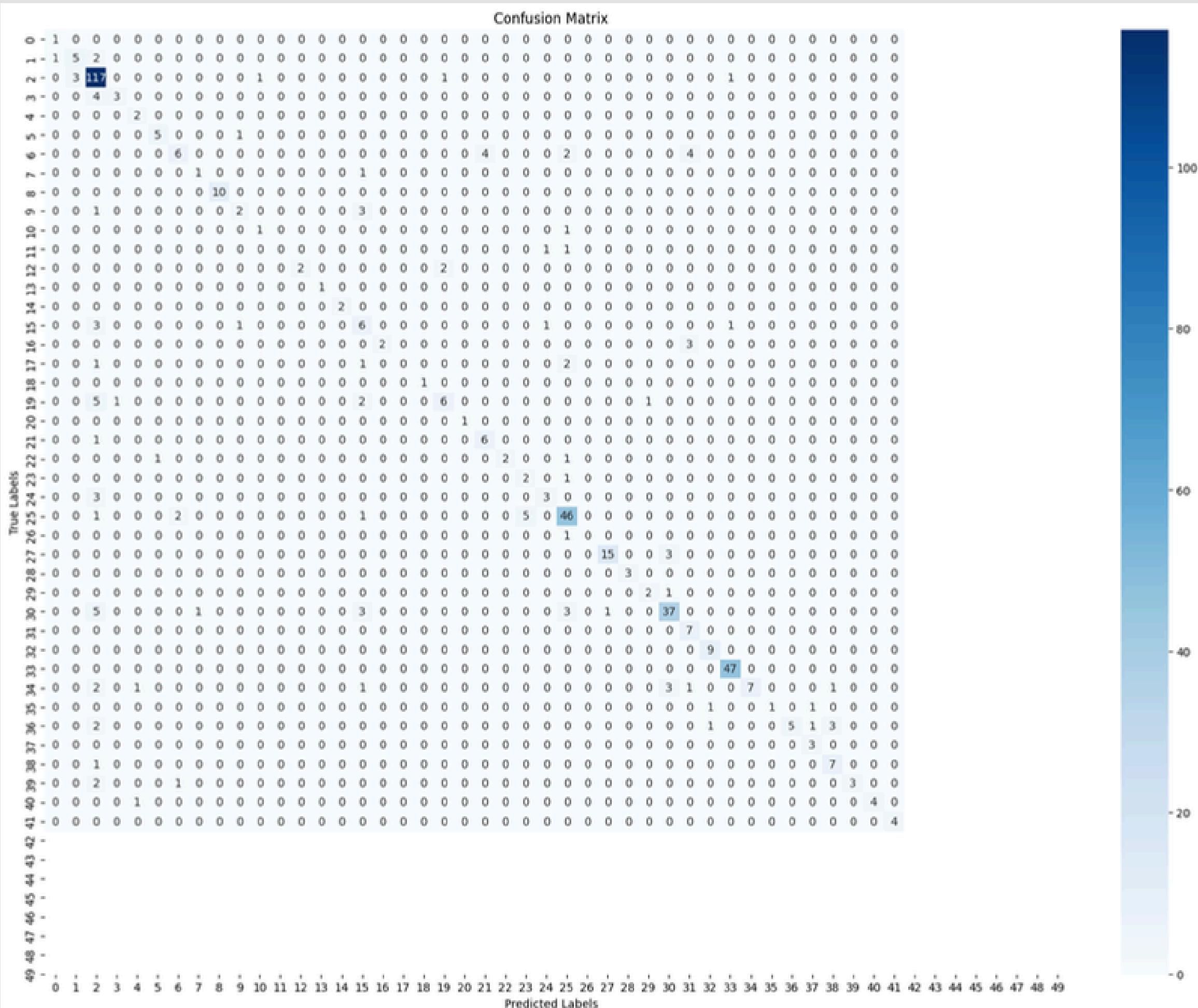
SMOTE + XGBoost

XGBoost works well when **combined** with SMOTE to balance the dataset before training to improve the **representation of the minority classes** while still benefiting from XGBoost's **efficient learning process**.

	precision	recall	f1-score	support
0	0.50	1.00	0.67	1
1	0.62	0.62	0.62	8
2	0.78	0.95	0.86	123
49	1.00	1.00	1.00	4
Accuracy			0.77	500
macro avg	0.72	0.67	0.67	500
weighted avg	0.79	0.77	0.76	500

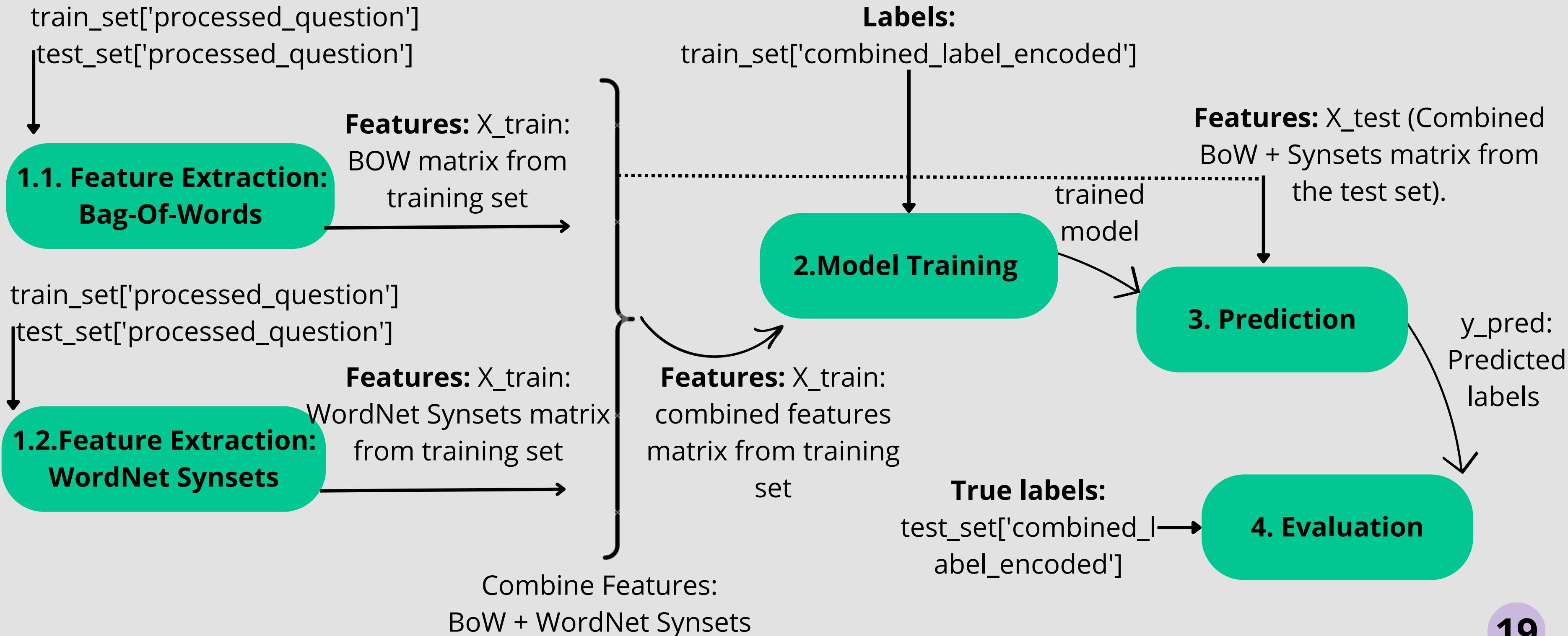
I . First Word Representation Approach (TF-IDF)

SMOTE + XGBoost



II. Second Word Representation Approach (BOW + WordNet Synsets):

- BoW focuses on word frequency but misses context.
- WordNet captures semantic context but might generalize too much.



I. First Word Representation Approach (TF-IDF)

Classification Report

	precision	recall	f1-score	support
0	0.00	0.00	0.00	1
1	0.86	0.75	0.80	8
2	0.77	1.00	0.87	123
49	1.00	0.75	0.86	4
Accuracy			0.81	500
macro avg	0.75	0.65	0.66	500
weighted avg	0.83	0.81	0.80	500

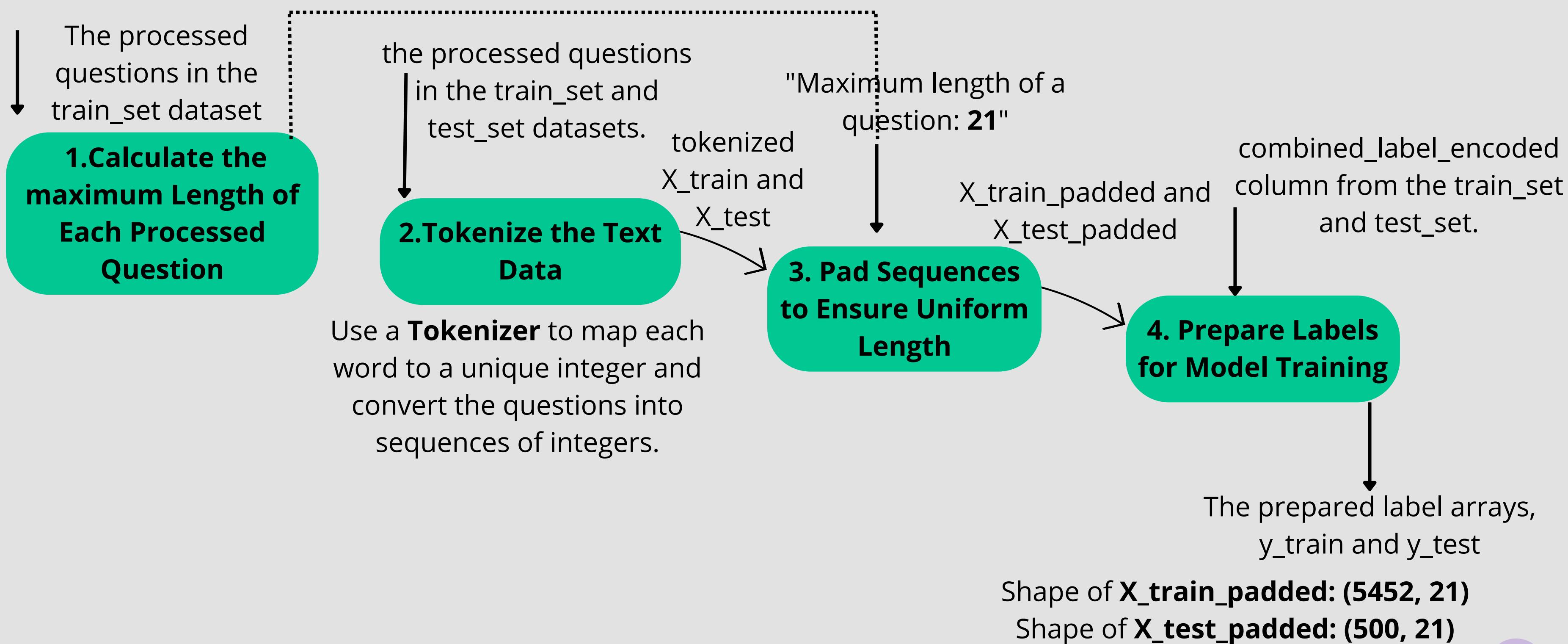
Support Vector Machine: SVM model

Recap

Approach	Embedding Method (Feature Extraction)	Classifier/model	precision	recall	f1-score	Accuracy
Classical Machine Learning	TF-IDF	Logistic Regression	0.70	0.71	0.66	0.71
		Logistic Regression + SMOTE	0.82	0.76	0.77	0.76
		Random Forest + SMOTE	0.77	0.76	0.75	0.76
		XGBoost	0.78	0.77	0.75	0.76
		XGBoost + hyperparameters tuning	0.78	0.76	0.75	0.76
		XGBoost + SMOTE	0.79	0.77	0.76	0.77
	BOW + WordNet Synsets	SVM	0.83	0.81	0.80	0.81

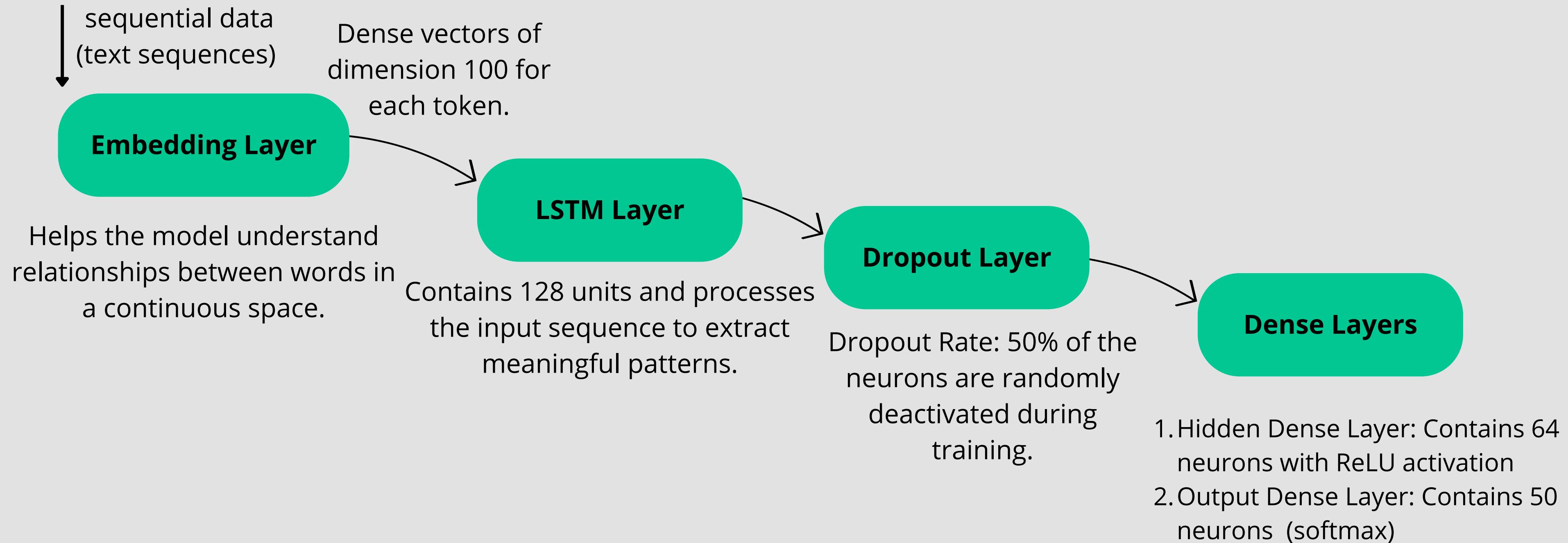
Use Deep Learning Approaches

Pre-processing



I. First Approach (Embedding: Use Embedding Layer)

Initial model (baseline)



I. First Approach (Embedding: Use Embedding Layer)

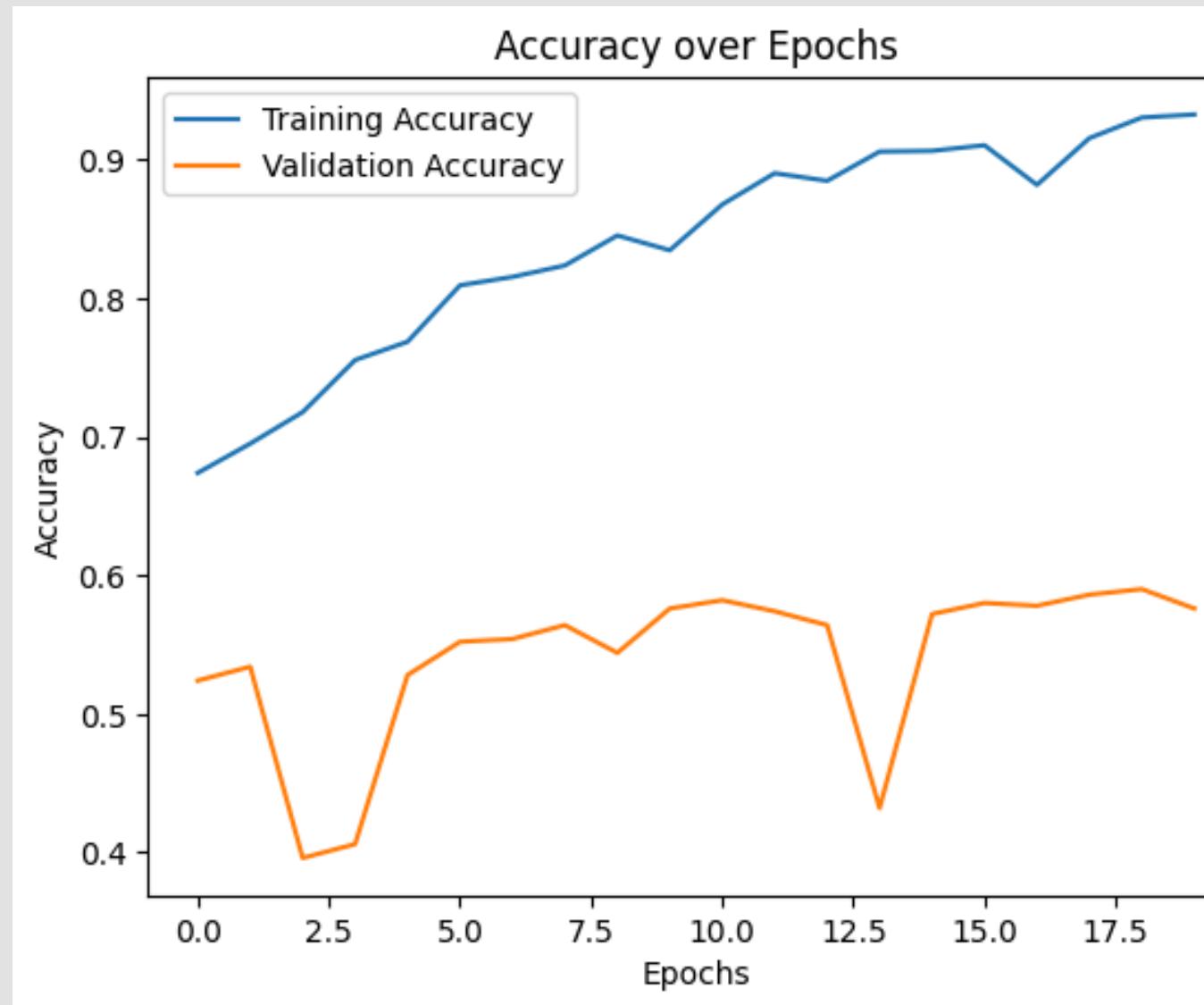
Initial model (baseline)

Test Accuracy: **57.60%**

Question	Predicted label
How far is it from Denver to Aspen ?	ENTY_body
What county is Modesto , California in ?	LOC_other
Who was Galileo ?	HUM_desc
What is an atom ?	DESC_desc
When did Hawaii become a state ?	NUM_date
How tall is the Sears Building ?	DESC_manner
George Bush purchased a small interest in which baseball team ?	DESC_def
What is Australia 's national flower ?	NUM_speed
Why does the moon turn orange ?	DESC_reason
What is autism ?	DESC_def

I. First Approach (Embedding: Use Embedding Layer)

Initial model (baseline)



high training accuracy



Overfitting

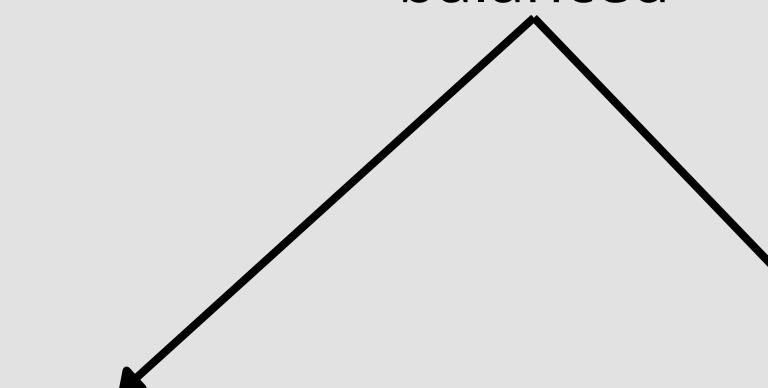
I. First Approach (Embedding: Use Embedding Layer)

Step 1: Class Weights for Imbalanced Data

Used class weights to address class imbalance with `class_weight` from `sklearn.utils` to make it balanced

Still Overfitting

Test Accuracy: 48.80%



Question	Predicted label
How far is it from Denver to Aspen ?	NUM_dist
What county is Modesto , California in ?	LOC_other
Who was Galileo ?	HUM_desc
What is an atom ?	ABBR_exp
When did Hawaii become a state ?	NUM_date
How tall is the Sears Building ?	NUM_dist
George Bush purchased a small interest in which baseball team ?	NUM_date
What is Australia 's national flower ?	ENTY_plant
Why does the moon turn orange ?	DESC_reason
What is autism ?	DESC_desc

was wrong
but wrongly
classified it
with another
class

I. First Approach (Use Embedding Layer)

Step	Key changes	Results
Adding L2 Regularization to LSTM	LSTM Layer: Added L2 regularization (kernel_regularizer=l2(0.01))	No significant improvement, overfitting still present.
Increasing Dropout	Dropout Rate: Increased to 0.6 (Dropout(0.6)).	Reduced overfitting slightly, but still not enough to generalize well to validation data.
Using Adam Optimizer with Smaller Learning Rate:	Optimizer: Adam optimizer with a reduced learning rate (learning_rate=0.0005).	Slower learning curve, but overfitting persisted.
Adding L1 and L2 Regularization Together	LSTM Layer: Used both L1 and L2 regularization (kernel_regularizer=l1_l2(l1=0.01, l2=0.01)).	Reduced some overfitting, but it remained evident.
Tested Various Batch Sizes	Experimented with different batch sizes, including 64 and 32.	Slight differences in performance, but no major effect on overfitting.
Reduced LSTM Units (Smaller Model)	Reduced the number of LSTM units to 64 (down from 128).	Slight improvement in generalization, but overfitting persisted.
Learning Rate Schedulers	Learning Rate Scheduler: Attempted to adjust the learning rate dynamically using a scheduler.	Slower convergence, but did not solve overfitting.

II. Second Approach (Use Pretrained GloVe Embeddings)

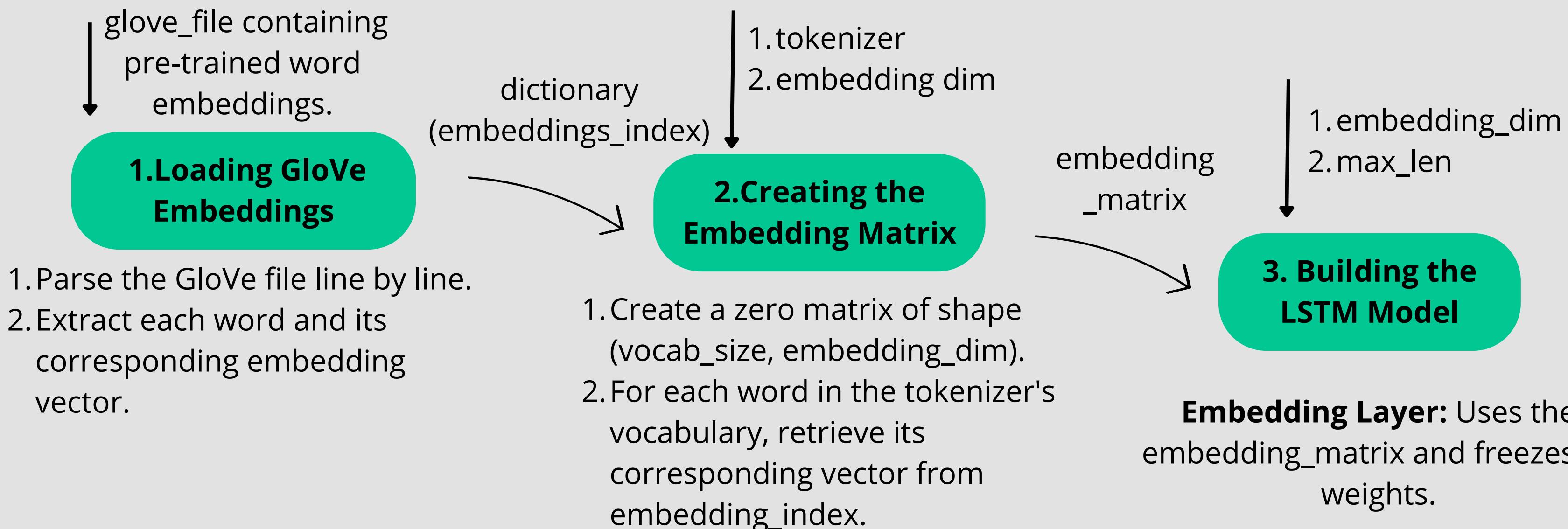
GloVe: Global Vectors for Word Representation : an **unsupervised learning** algorithm for obtaining vector representations for words. (deeper contextual embeddings)

Wikipedia 2014 + Gigaword 5 (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors)

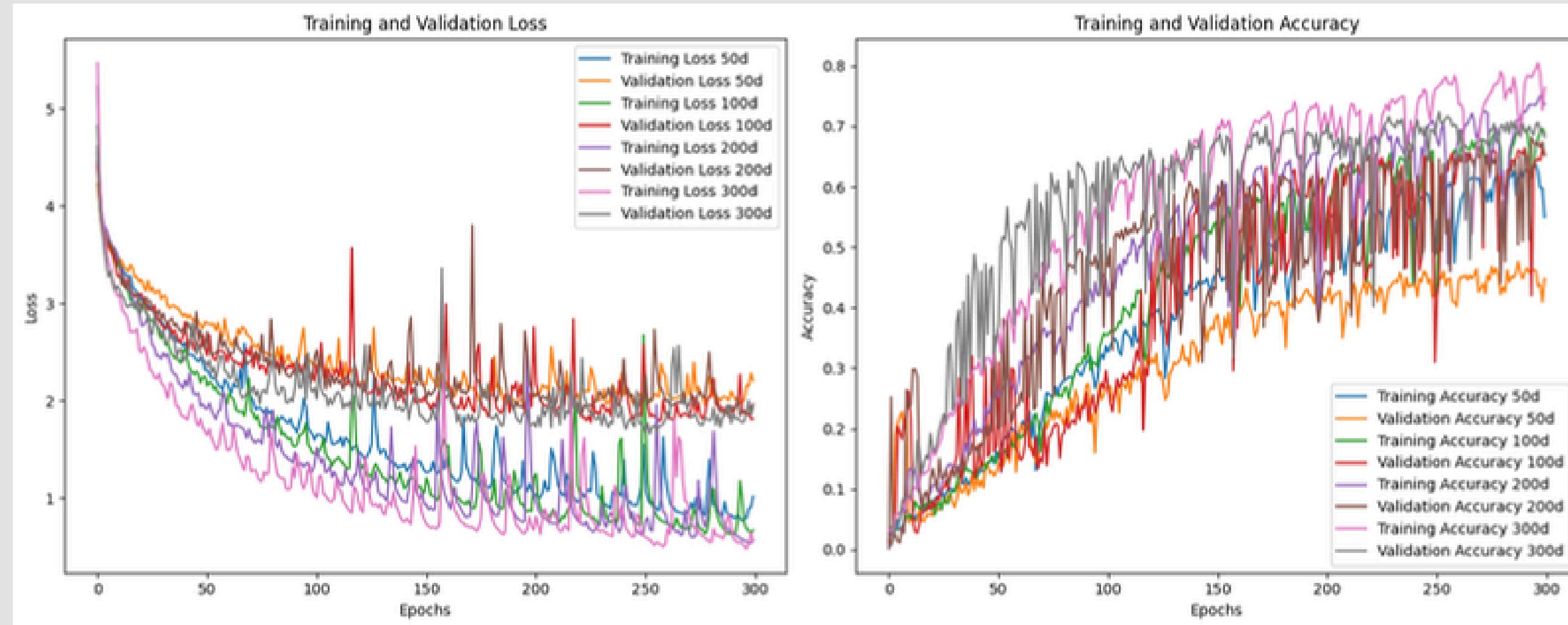
<word> <value_1> <value_2> ... <value_N>

<word>: The word or token.

<value_1>, <value_2>, ..., <value_N>: The numerical values representing the word's embedding in an N-dimensional space.



II. Second Approach (Use Pretrained GloVe Embeddings)

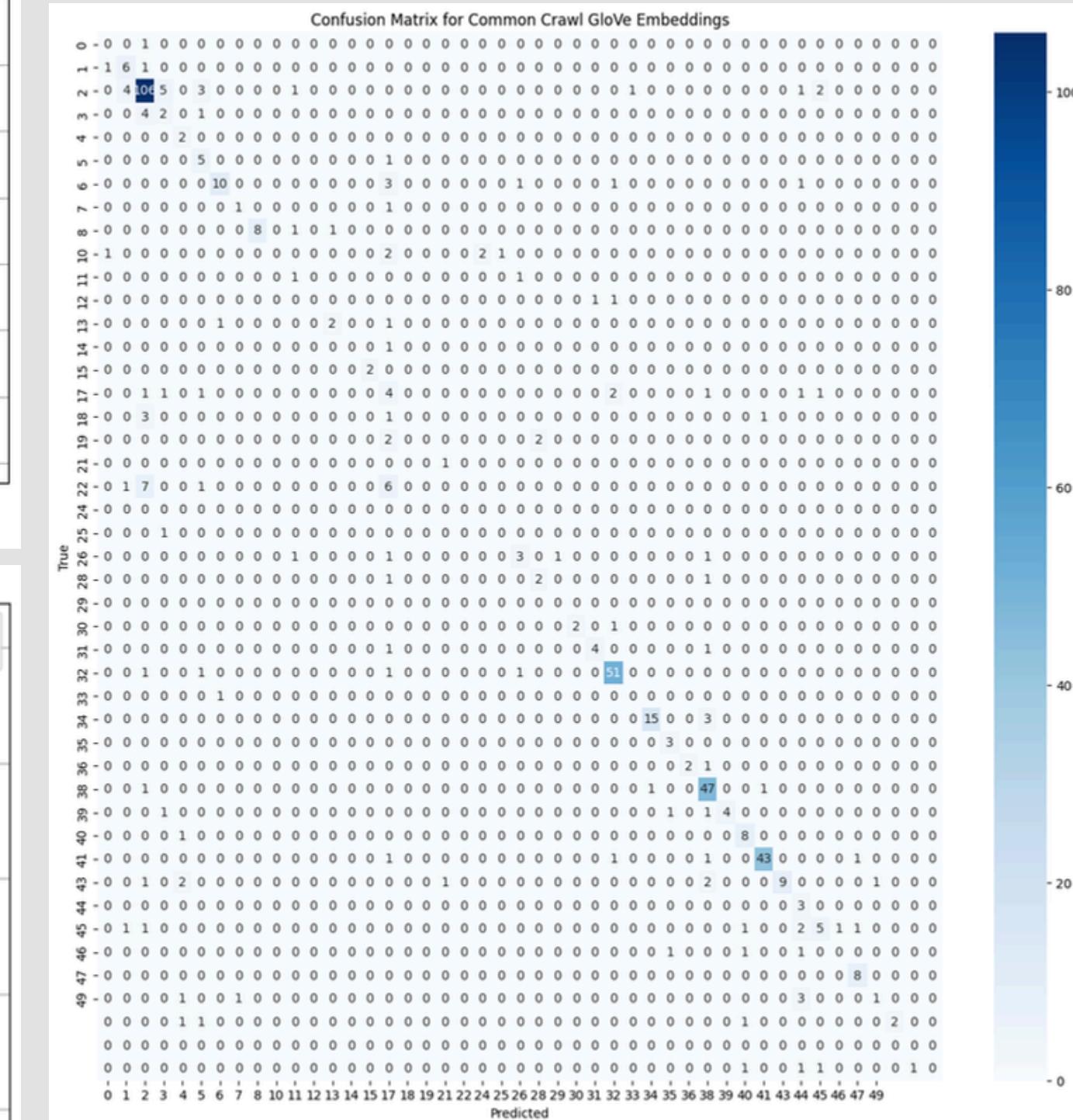
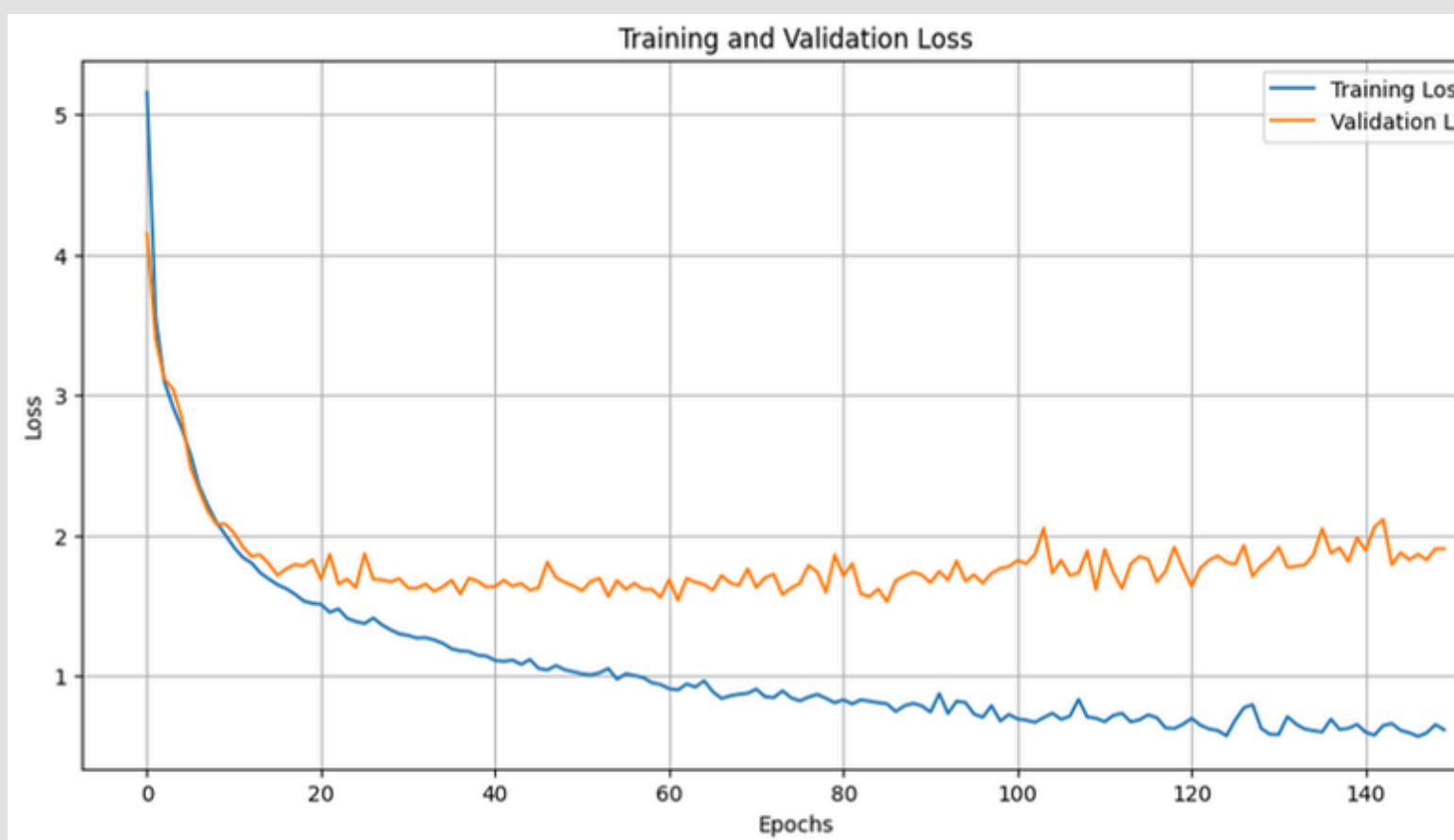
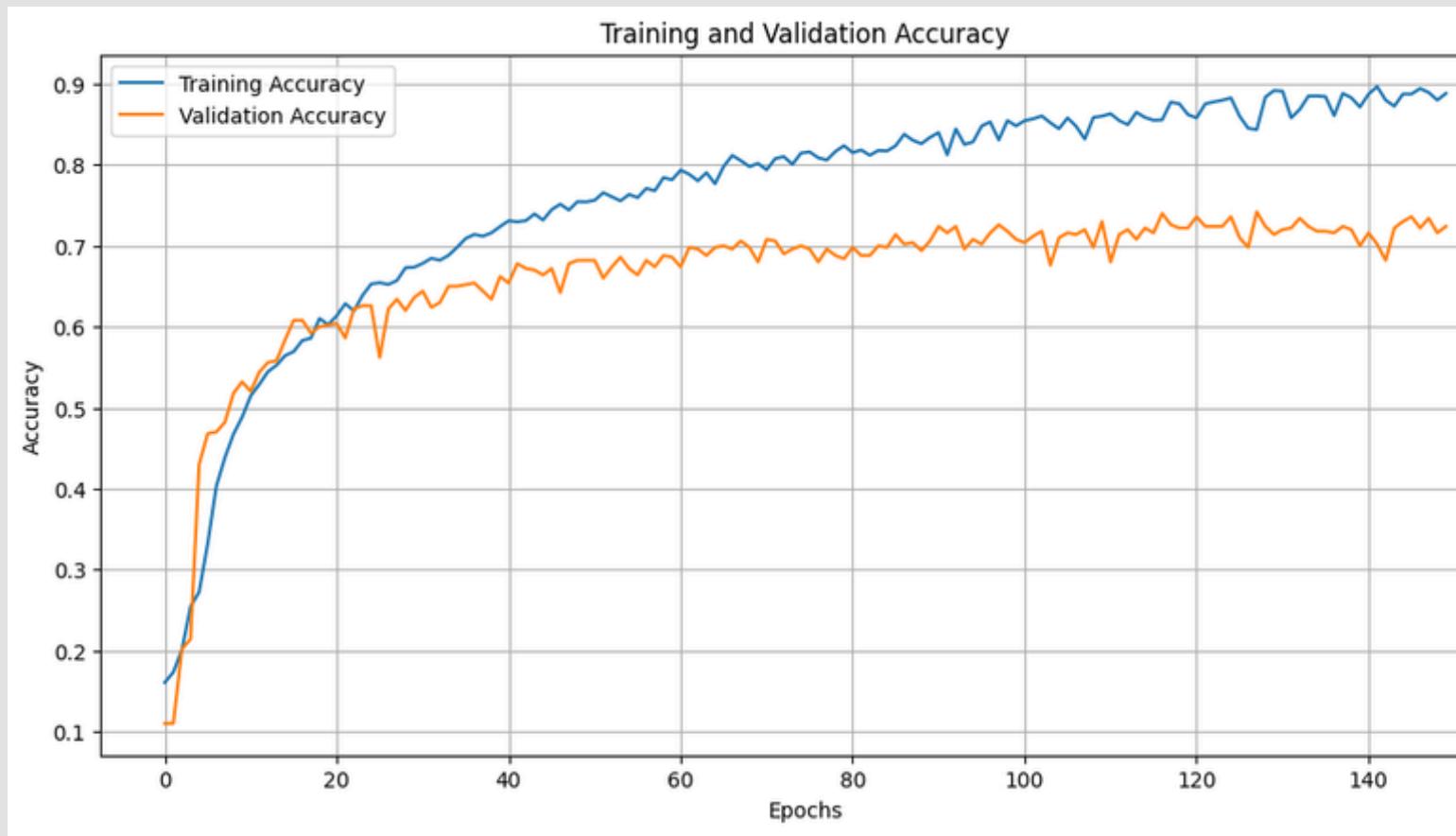


On **test Set**

	precision	recall	f1-score	Accuracy
50d GloVe	0.6785	0.4480	0.4621	0.4480
100d GloVe	0.7198	0.6520	0.6732	0.6520
200d GloVe	0.7455	0.6520	0.6792	0.6520
300d GloVe	0.7564	0.6820	0.6993	0.6820

II. Second Approach (Use Pretrained GloVe Embeddings)

Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors)



Evaluation Metrics:
Accuracy: 0.7180
F1 Score: 0.6973
Precision: 0.7011
Recall: 0.7180

Use Transformer-Based Approach

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

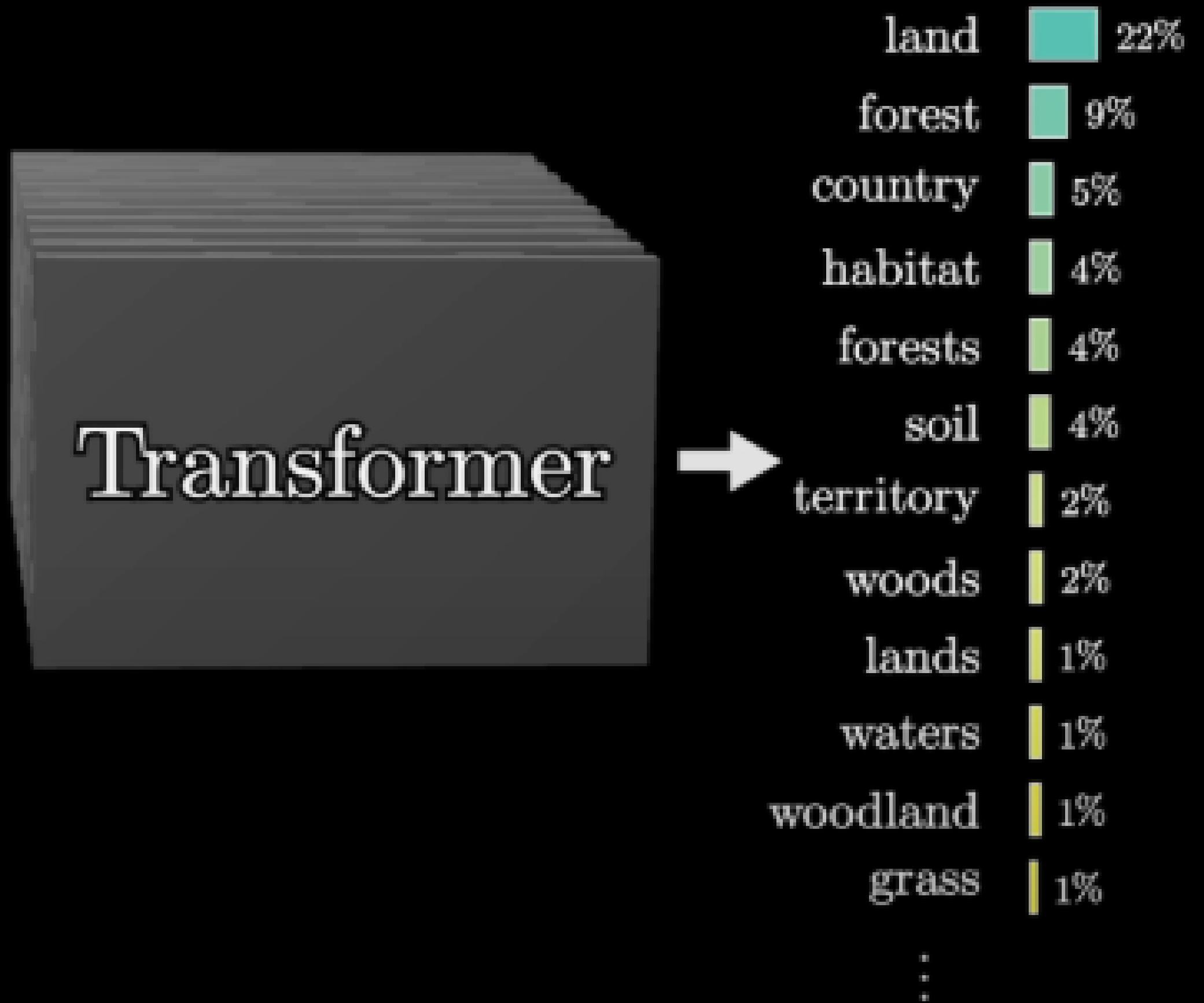
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to

Input Sequence:

Behold, this wild creature
foraging on his ____.



But First We tokenize :

To date, the cleverest thinker of all time was ...

The tokenization of this input would be:

To| date|,| the| cle|ve|rest| thinker| of| all| time| was ...

Tokens assigned to Vectors

To	date		the	cle	ve	rest	thinker	of	all	time	was
5.4	7.8	9.7	2.6	3.6	5.6	1.6	9.7	3.2	6.7	4.4	
7.1	5.2	7.9	7.7	4.3	4.3	1.1	4.6	6.6	2.7	8.4	
6.0	5.6	4.6	4.5	6.9	9.8	6.5	9.7	1.3	7.3	6.9	
5.4	9.2	7.7	5.6	0.6	1.0	1.4	6.0	7.1	9.5	2.9	
4.2	0.7	1.2	0.2	6.6	2.1	1.9	7.3	...	2.9	2.5	8.1
6.4	0.9	6.3	6.1	6.6	1.6	3.7	0.4	1.8	5.7	3.9	
4.3	0.2	1.4	6.1	2.1	6.5	8.1	2.8	5.8	5.9	8.7	
8.8	8.2	9.4	6.1	1.3	2.5	1.0	1.2	0.2	5.7	5.8	
:	:	:	:	:	:	:	:	:	:	:	
3.8	8.6	4.1	6.8	3.6	2.4	1.0	1.2	0.0	9.4	6.9	

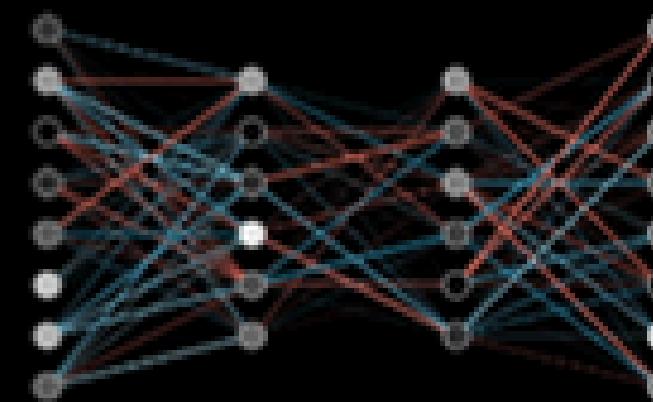
bound
jump
skip

Words with similar meanings have close vectors

To date | the do we remember? | the blue

Attention

Attention Block



A machine learning model ...

0.7	0.2	0.0	0.6
0.7	0.8	0.6	0.7
0.4	0.6	0.5	0.6
0.5	0.9	0.8	0.7
0.3	0.7	0.6	0.5
0.1	0.5	0.4	0.3
0.2	0.3	0.2	0.4

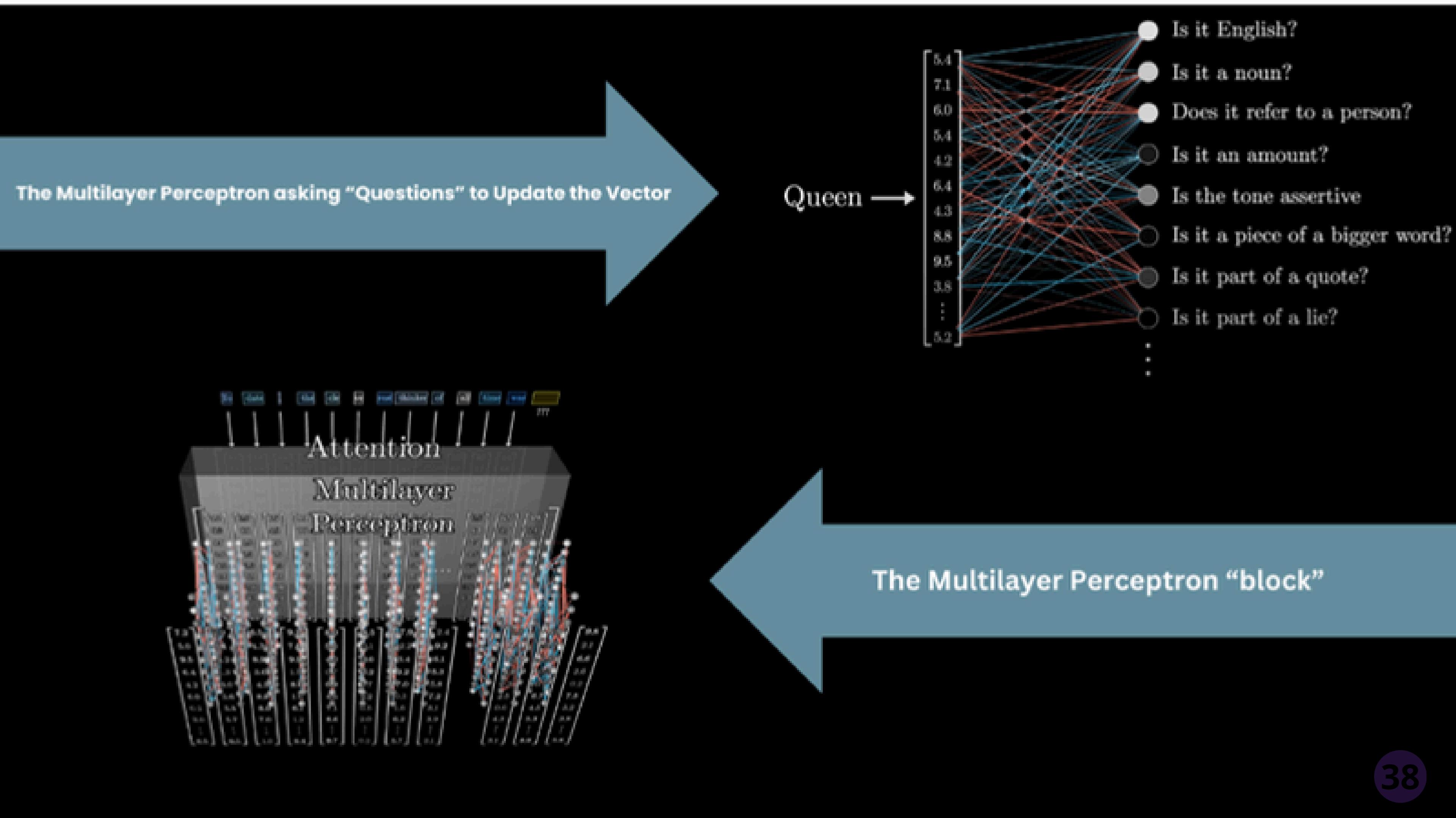


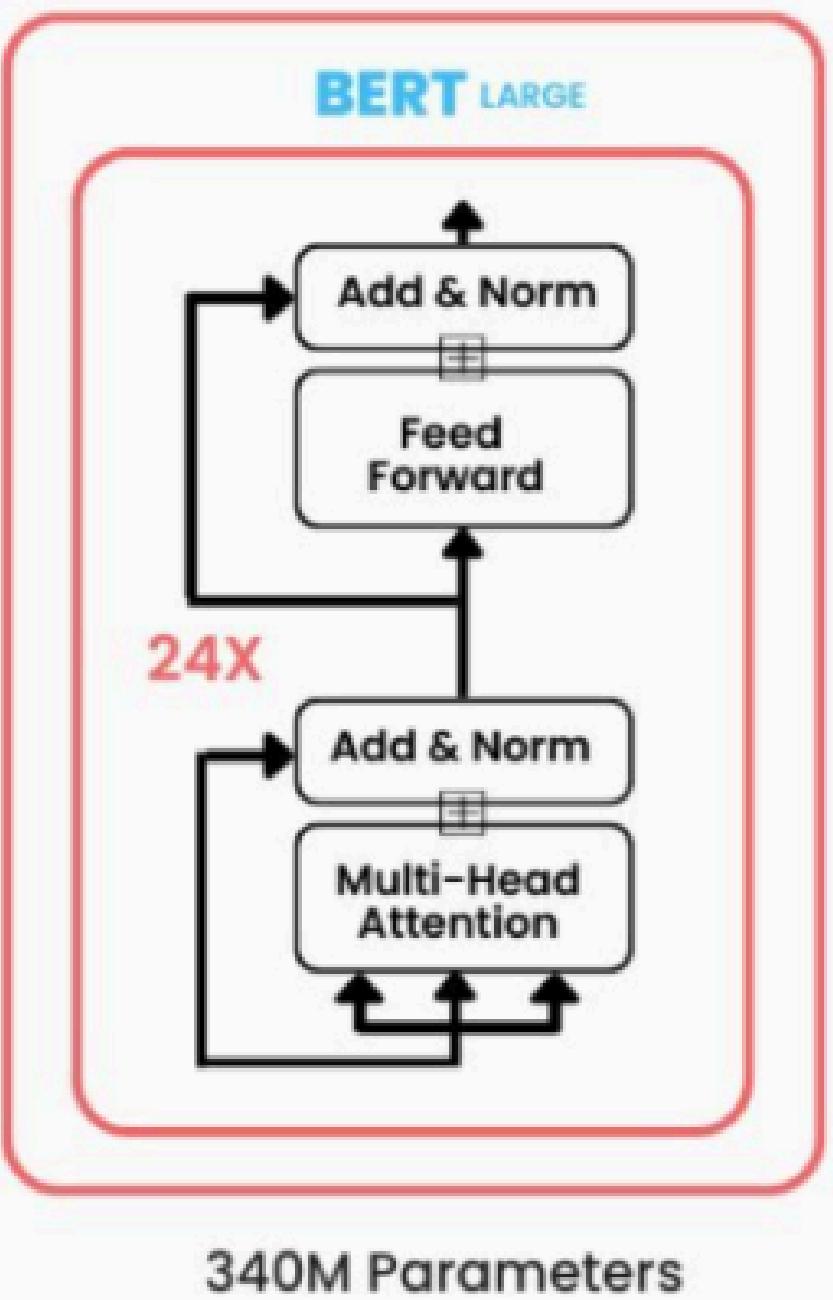
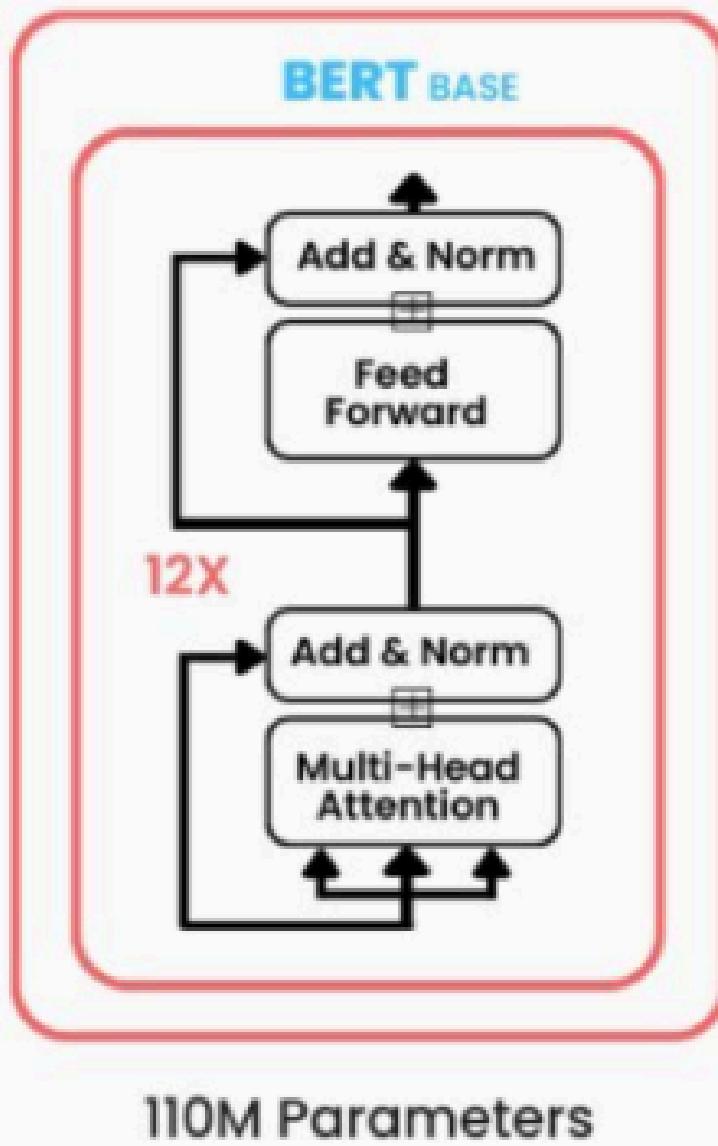
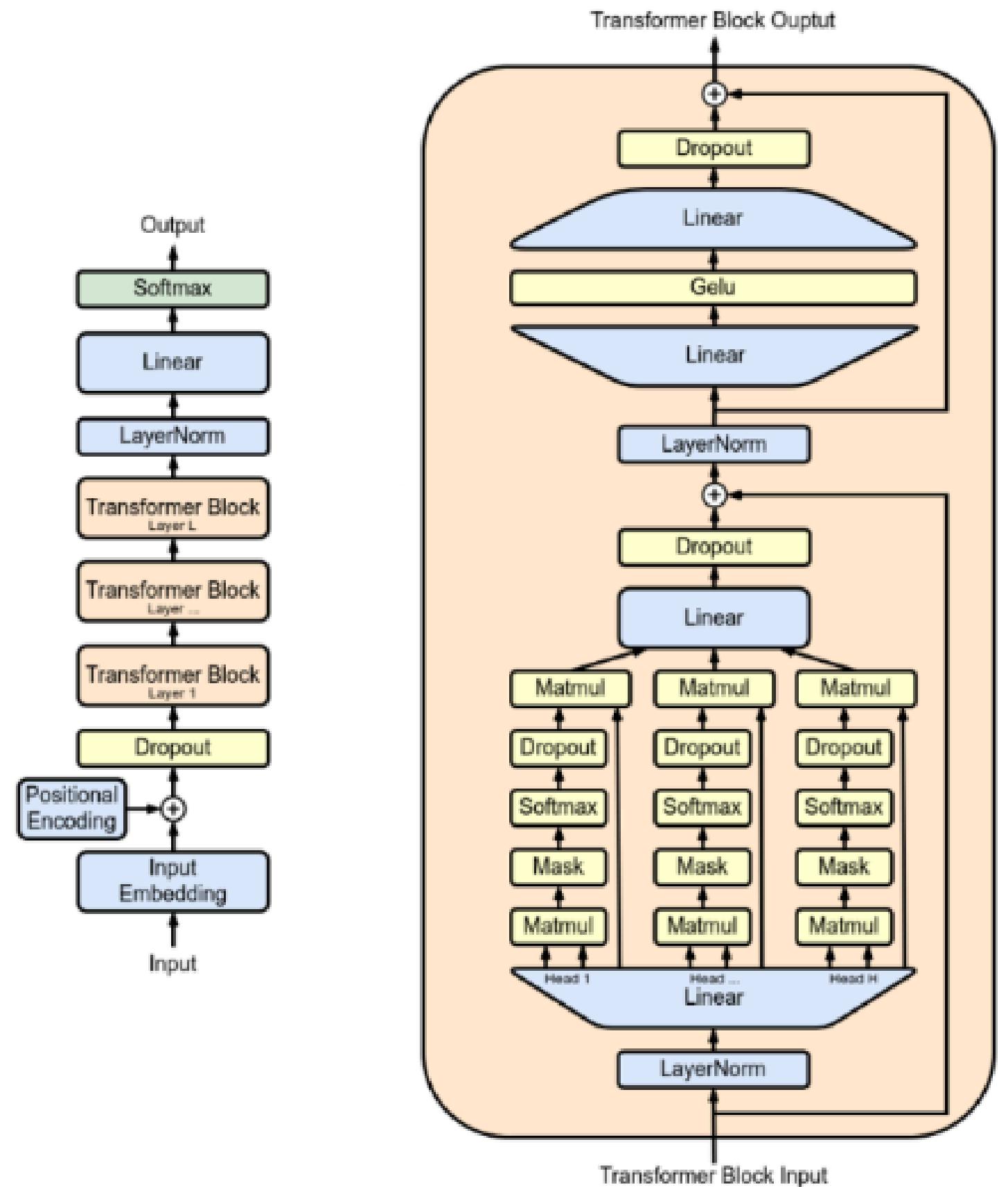
A fashion model ...

0.7	0.3	0.6
0.5	0.5	0.5
0.3	0.3	0.3
0.5	0.5	0.5
0.2	0.2	0.2
0.1	0.1	0.1

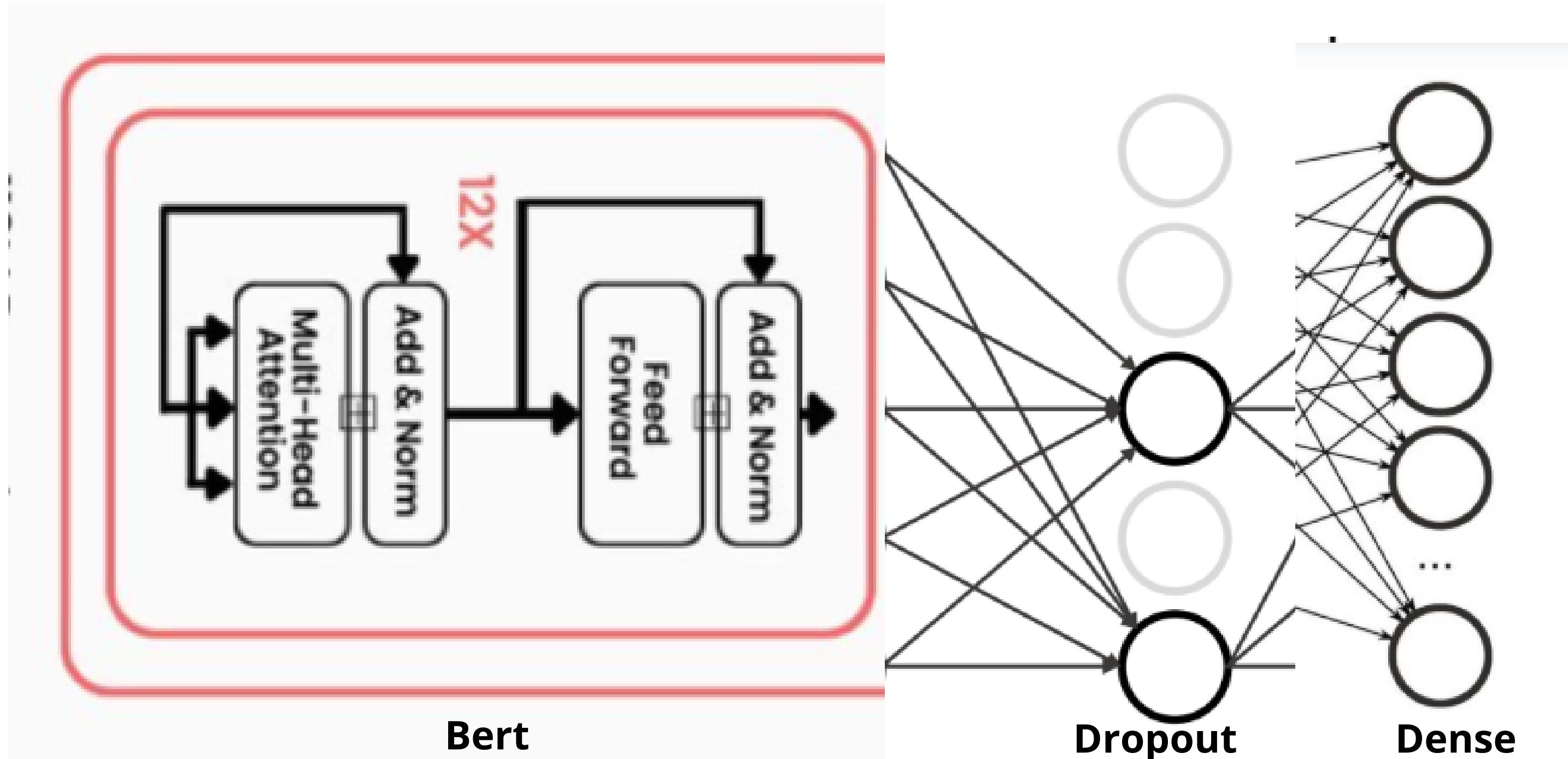
Example of Attention Block changing the vector of the word "model"







Architecture



Data

B	C	D	E
Questions	Category0	Category1	Category2
0 How did serfdom develop in and then leave Russia ?	DESCRIPTION	DESC	manner
1 What films featured the character Popeye Doyle ?	ENTITY	ENTY	cremat
2 How can I find a list of celebrities ' real names ?	DESCRIPTION	DESC	manner
3 What fowl grabs the spotlight after the Chinese Year of the M	ENTITY	ENTY	animal
4 What is the full form of .com ?	ABBREVIATION	ABBR	exp
5 What contemptible scoundrel stole the cork from my lunch ?	HUMAN	HUM	ind
6 What team did baseball 's St. Louis Browns become ?	HUMAN	HUM	gr
7 What is the oldest profession ?	HUMAN	HUM	title
8 What are liver enzymes ?	DESCRIPTION	DESC	def
9 Name the scar-faced bounty hunter of The Old West .	HUMAN	HUM	ind

Loss Function- Multilabel Task

MultiLabel : each sample can be in any number of the specified classes, including zero. it can be understood as a series of binary classifications:

Is sample 1 in class A – yes or no? Is sample 1 in class B – yes or no?
And so on.

Binary Cross Entropy Loss :

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

- N is the number of samples
- y_i is the true label for the i^{th} sample (0 or 1)
- p_i is the predicted probability that the i^{th} sample belongs to Positive class

Loss Function- Multilabel Task

BCElogits Loss :

$$L = -\frac{1}{N} \sum_{i=1}^N \left[y_i \log_e \left(\frac{1}{1 + e^{-z}} \right) + (1 - y_i) \log_e \left(1 - \frac{1}{1 + e^{-z}} \right) \right]$$

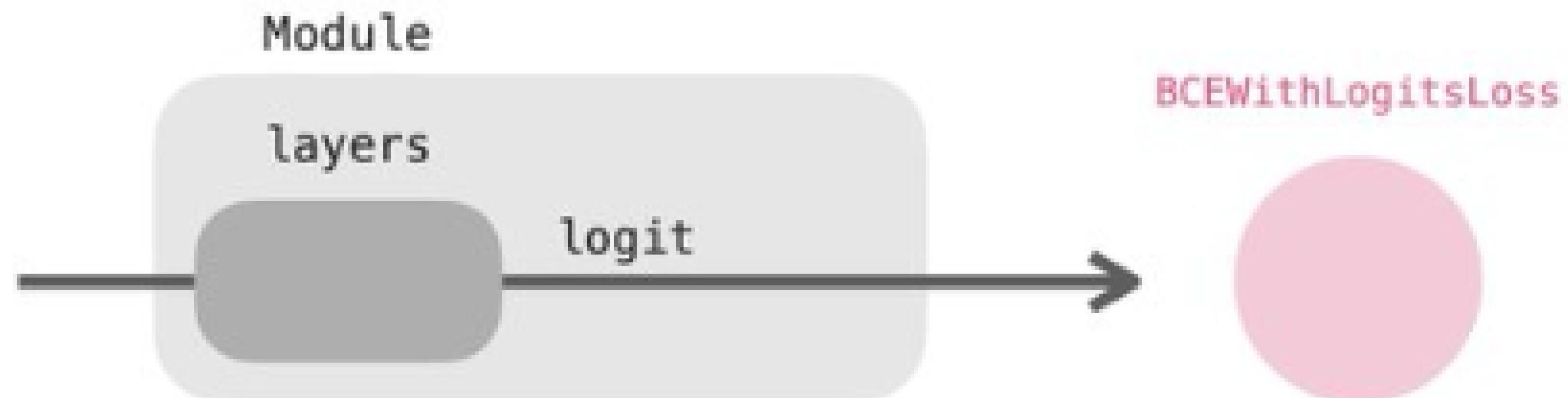
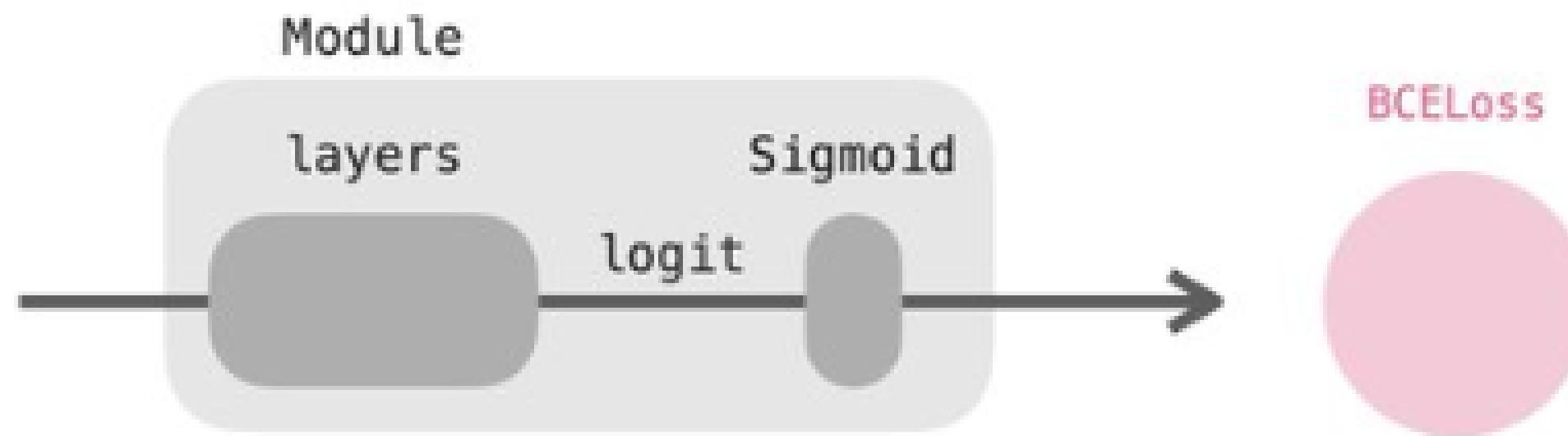
> Using Quotient Rule $\log_a \left(\frac{m}{n} \right) = \log_a m - \log_a n$

Each class label is effectively an independent binary classification problem (is this class present: yes/no?)

BCEWithLogitsLoss combines two operations efficiently:

- A sigmoid function to convert raw logits to probabilities (0-1 range)
- Binary cross entropy loss calculation

Loss Function- Multilabel Task



Conclusion

Perspectives

Areas of research :

- Since LLMs are too costly how would Small LMs behave in this task ? like SmoLLm , phi4
- LLMs being great Zero shot learners and few-shot learners ?
- How can we address the inexplicability nature of Transformers/LLMs?
- Can the Small Models substitute the large ones if trained on a very specific task thus bringing down energy consumption and cost?

Work Methodology



version control

A screenshot of a GitHub repository named 'QuestionsClassification'. The repository is public. It shows 4 branches and 0 tags. A search bar at the top says 'Find or create a branch...'. Below it, there are tabs for 'Branches' and 'Tags'. A red oval highlights the 'Feat/HatemTrigui' branch, which is circled by a red arrow pointing to the right. Other branches listed are 'main', 'Feat/RubensBaccari', and 'Feat/malekchassouna'. At the bottom, there's a link 'View all branches'.

A screenshot of the same GitHub repository 'QuestionsClassification'. The repository is public. It shows 4 branches and 0 tags. A green button at the top right says 'Compare & pull request'. Below it, a message says 'This branch is 9 commits ahead of main.' A list of commits is shown, all made by the user 'hatrigui'. The commits are:

- add paper manipulations (feat) · 7 minutes ago · 11 Commits
- data · 2 weeks ago
- C02-1150.pdf · yesterday
- Deep Learning Approaches.ipynb · 8 minutes ago
- Different Classical Machine Learning Approach... · 6 hours ago
- Paper manipulations.ipynb · 7 minutes ago
- QuestionsClassification.ipynb · 2 weeks ago
- README.md · 2 weeks ago

Github Repo Link: <https://github.com/hatrigui/QuestionsClassification>



**THANK YOU FOR
YOUR ATTENTION**