

Deep Learning Practice - NLP

Recap of NN models, Roadmap, Datasets

Mitesh M. Khapra



AI4Bharat, Department of Computer
Science and Engineering, IIT Madras

Module 1 : Natural Language Processing

Mitesh M. Khapra

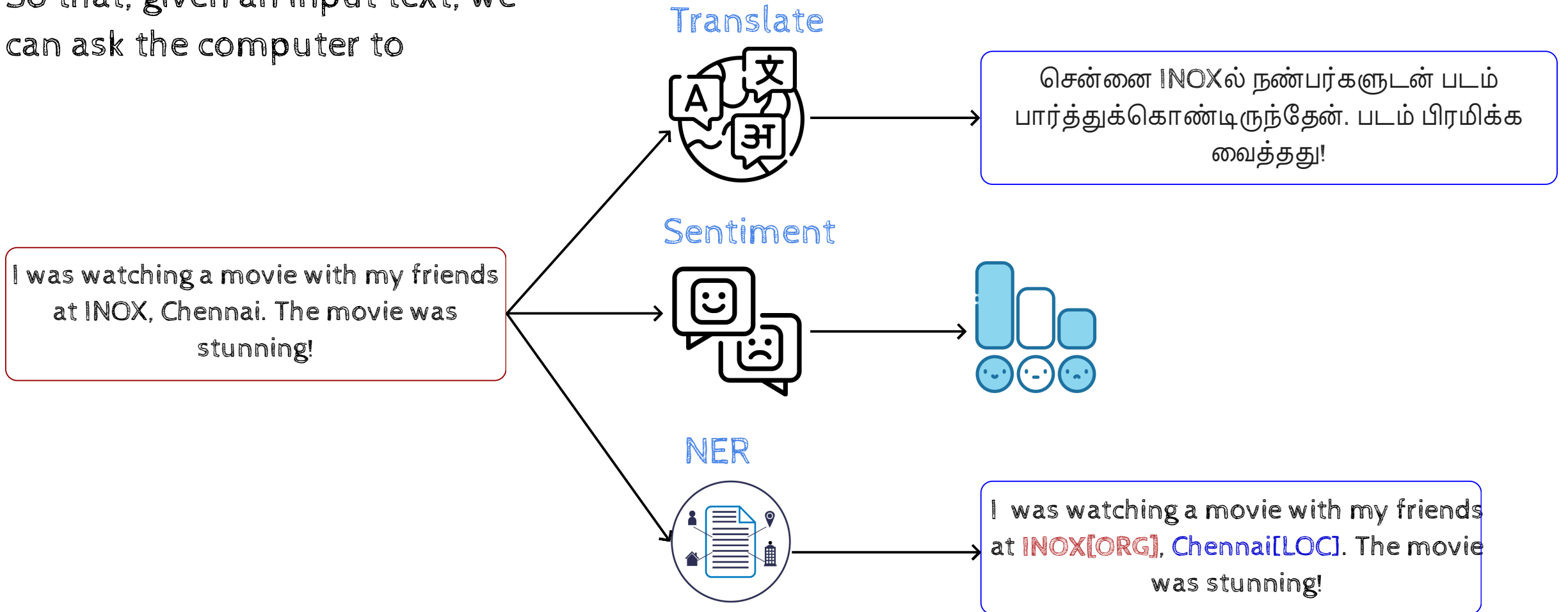


AI4Bharat, Department of Data Science
and Artificial Intelligence, IIT Madras

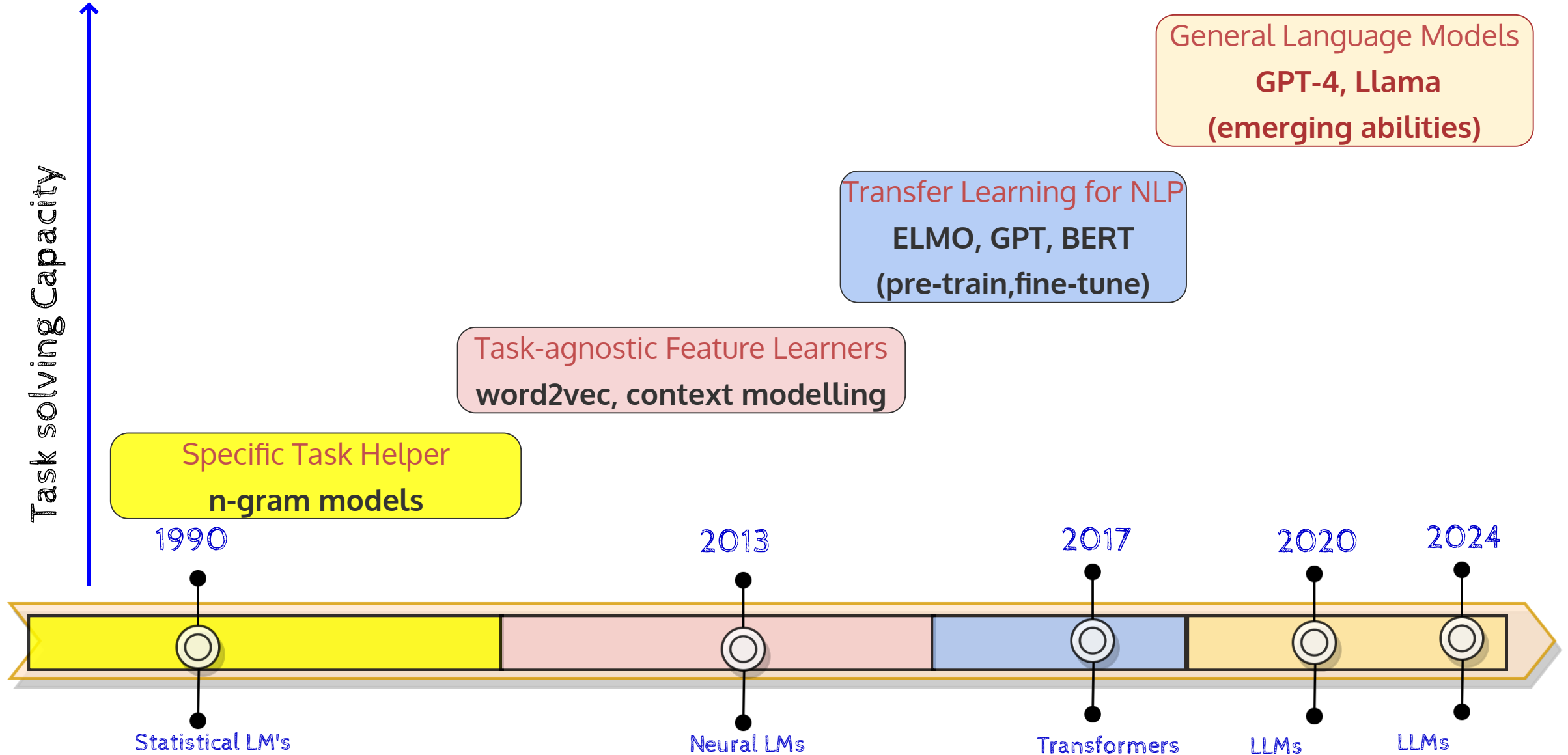
What is it?

enabling computers to understand, interpret, and generate human language

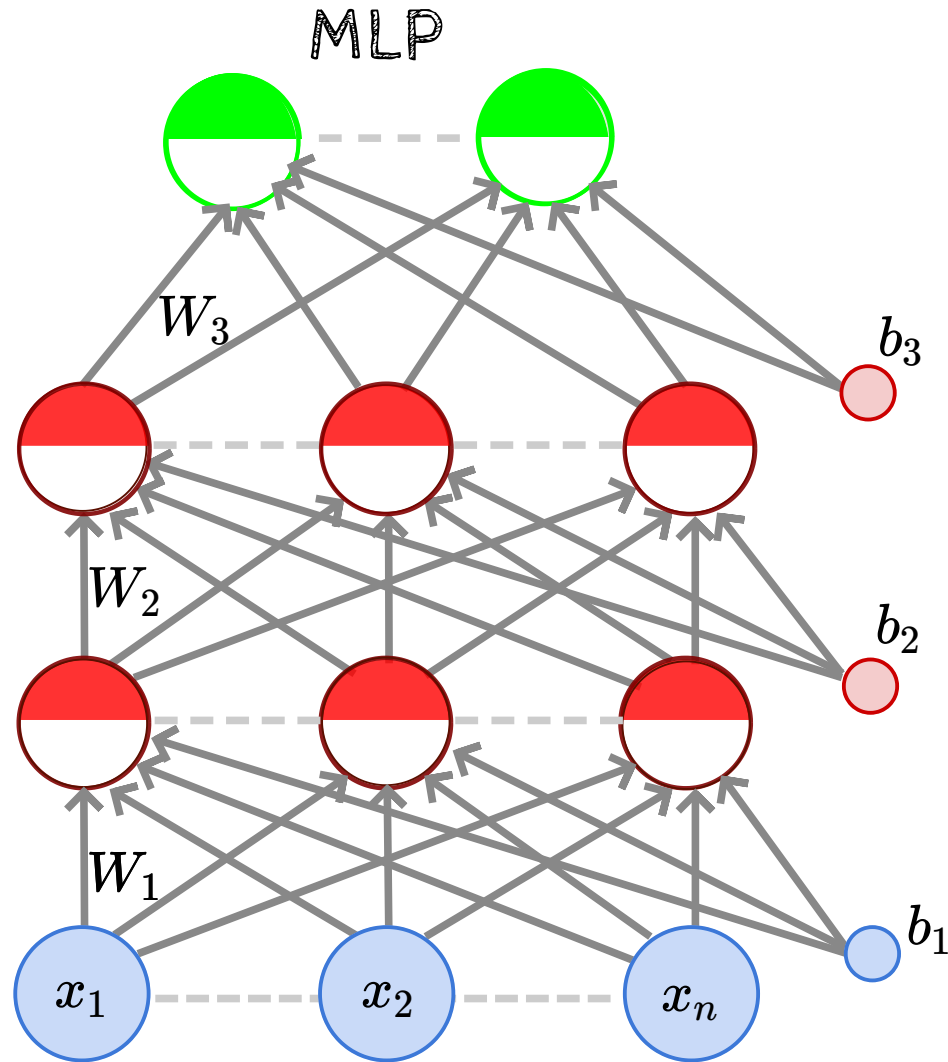
So that, given an input text, we can ask the computer to



What are the methods?



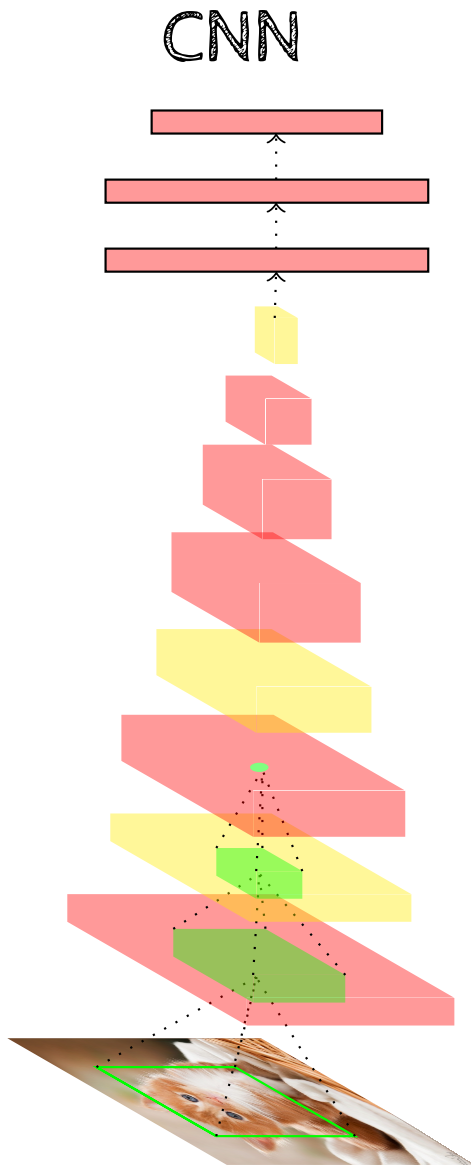
Neural Network Models



Let's quickly recap the four types of architectures that we learned in the deep learning theory course

We started with a simple Multi-Layer Perceptron (MLP) and the backpropagation algorithm for learning the parameters

It is a feedforward **fully connected** neural network

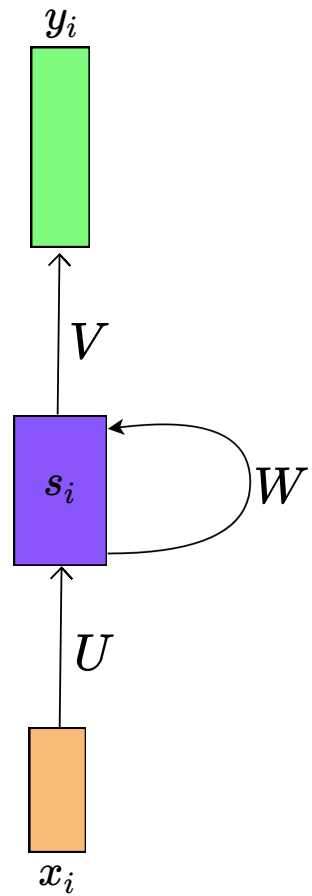


It is also a feedforward neural network commonly used in the domain of computer vision

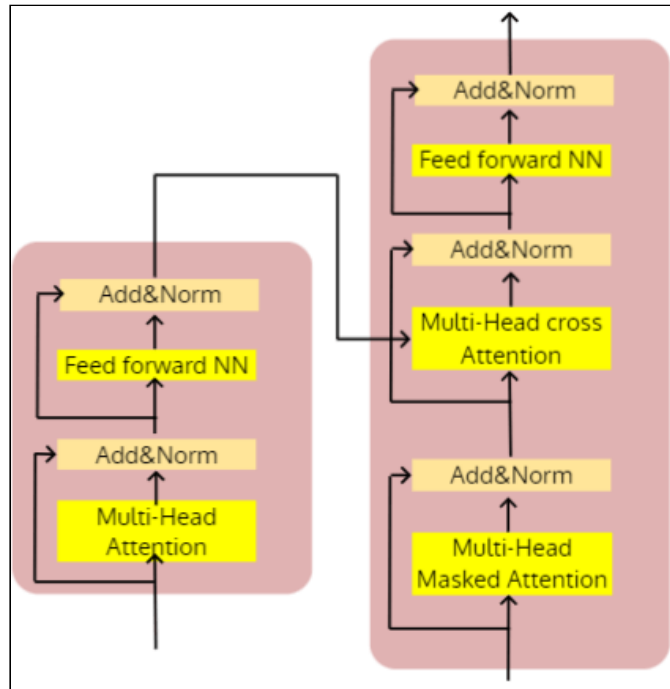
We introduced the weight sharing concept where the weights in the kernels are shared across the input.

This idea of weight sharing is prevalent in almost all modern neural networks

RNN



Transformer

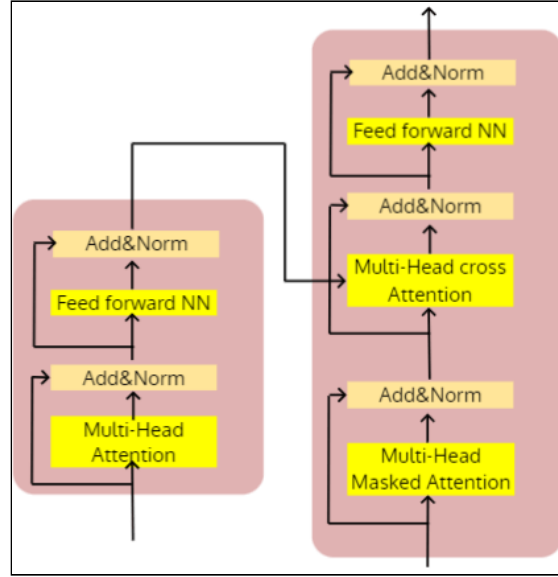
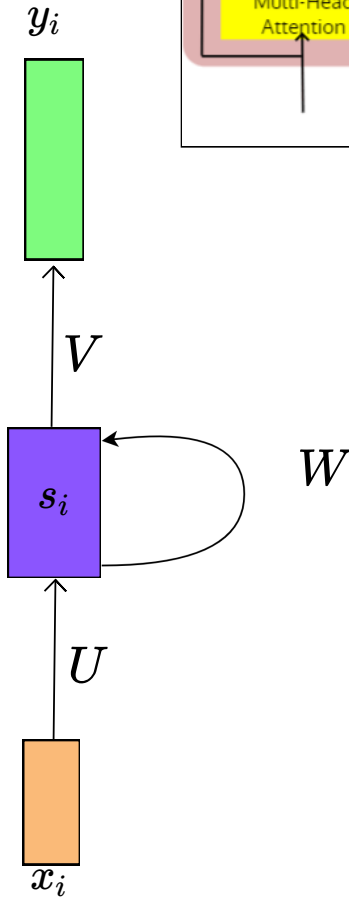
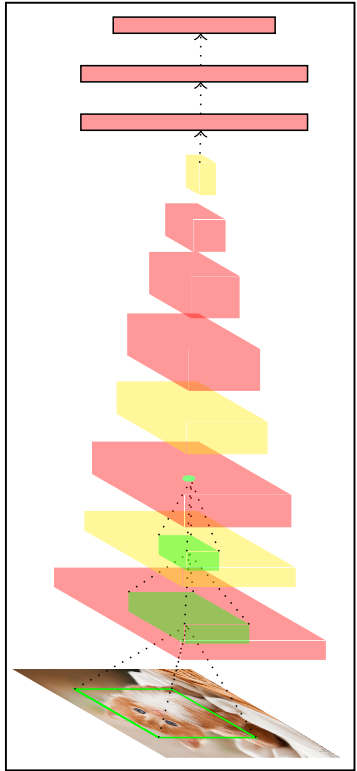
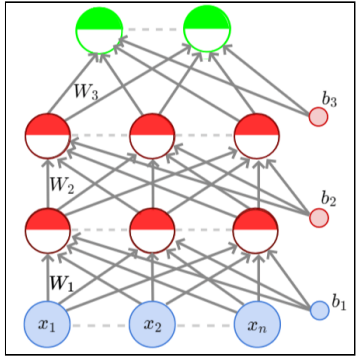


Recurrent Networks such as RNN, LSTM and transformers use the same weights for all elements/timesteps in the input sequence

Therefore, they are better suited for Natural Language Processing tasks, as weight sharing enables the models to generalize to sequences of any length

However, training of transformers can be parallelized as opposed to RNNs and LSTMs

This gives the transformer an edge over recurrent neural networks



Today, we can quickly build any of these architectures using modern deep learning frameworks such as Pytorch and TensorFlow

In this course, we will use Pytorch or any framework built on top of that

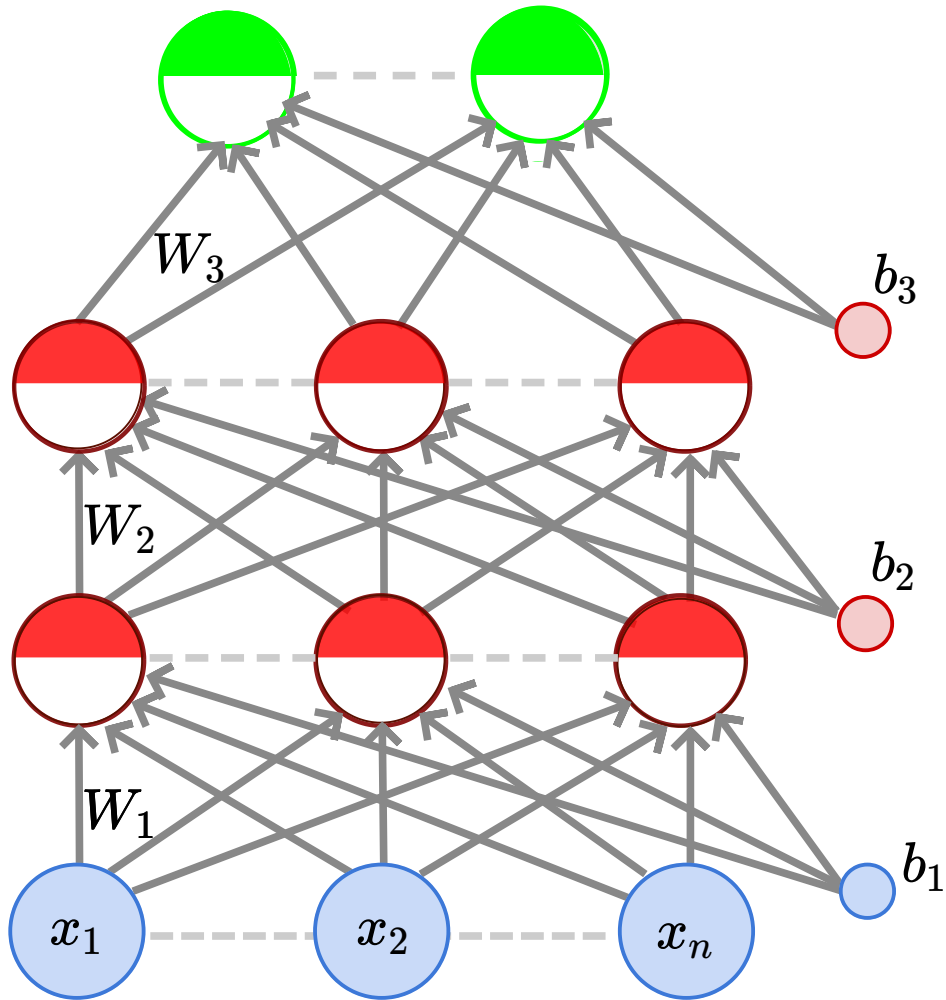
All architectures contain a sequence of learnable layers (linear, convolution, rnn, embedding ...)

Every learnable layer can be derived from the `torch.nn.Module` in PyTorch

Therefore, any architecture can be composed with the `torch.nn.Module`

let's build all four architectures with a few lines of code

MLP

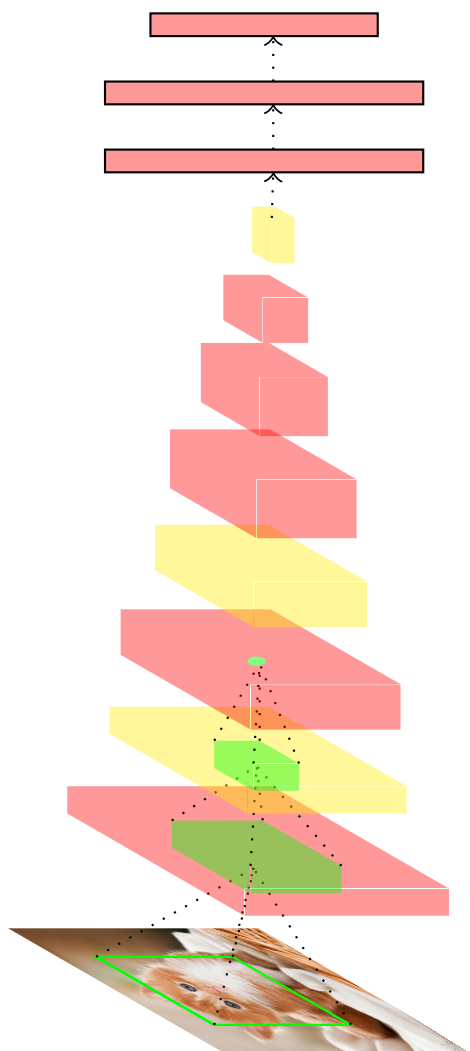


Implementation in Pytorch

The model is composed of linear layers and non-linear activations (assuming Relu). Therefore, the code contains only `nn.linear` and `torch.relu`

```
1 class MLP(nn.Module):
2
3     def __init__(self, x_dim, h1_dim, h2_dim, out_dim):
4         super().__init__()
5         self.w1 = nn.Linear(x_dim, h1_dim, bias=True)
6         self.w2 = nn.Linear(h1_dim, h2_dim, bias=True)
7         self.w3 = nn.Linear(h2_dim, out_dim, bias=True)
8
9     def forward(self, x):
10         out = torch.relu(self.w1(x))
11         out = torch.relu(self.w2(out))
12         out = torch.relu(self.w3(out))
13         return out
```

CNN

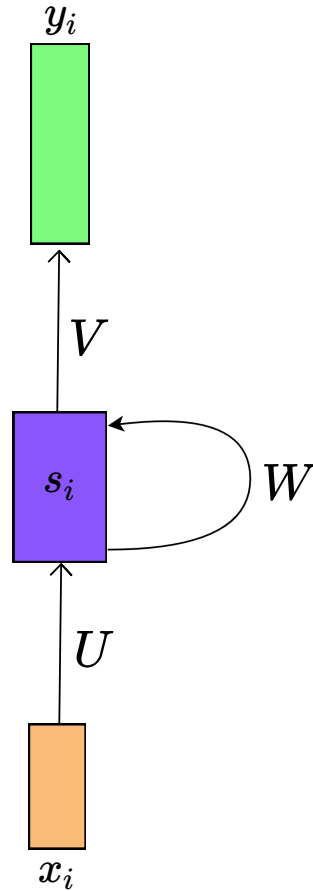


Implementation in Pytorch

The model is composed of `nn.linear`, `nn.Conv2d`, `MaxPool2D` and `torch.relu`

```
1 class CNN(nn.Module):
2     def __init__(self):
3         super(CNN, self).__init__()
4         self.conv1 = nn.Conv2d(3, 6, 5)
5         self.pool = nn.MaxPool2d(2, 2)
6         self.conv2 = nn.Conv2d(6, 16, 5)
7         self.fc1 = nn.Linear(16 * 5 * 5, 10)
8         self.fc2 = nn.Linear(10, 4)
9         self.fc3 = nn.Linear(4, 2)
10
11     def forward(self, x):
12
13         x = self.pool(F.relu(self.conv1(x)))
14         x = self.pool(F.relu(self.conv2(x)))
15         x = x.view(-1, 16 * 5 * 5)
16         x = F.relu(self.fc1(x))
17         x = F.relu(self.fc2(x))
18         x = self.fc3(x)
19         return x
```

RNN



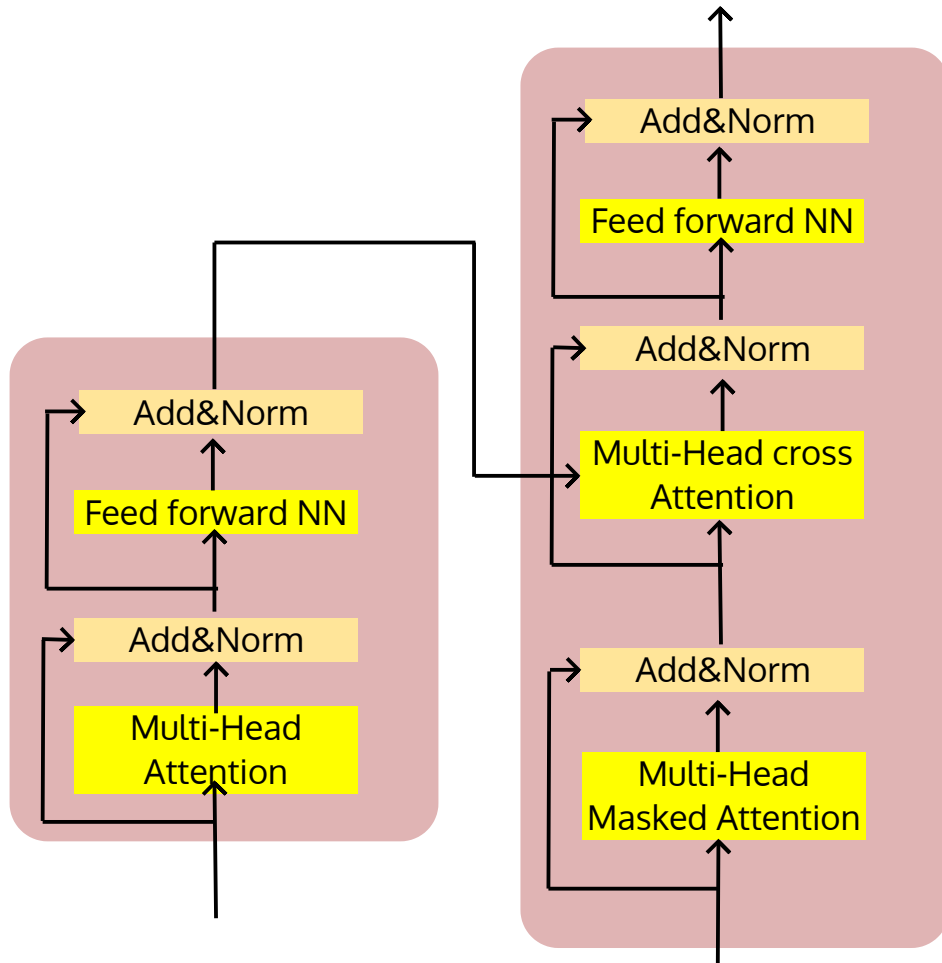
Implementation in Pytorch

The model is composed of `nn.linear`, `nn.Embedding`, `nn.RNN`

`nn.RNN` is composed of `nn.linear`, `torch.relu` and other required modules

```
1 class RNN(torch.nn.Module):
2     def __init__(self, vocab_size, embed_dim, hidden_dim, num_class):
3         super().__init__()
4         self.embedding = nn.Embedding(vocab_size,
5                                       embed_dim,
6                                       padding_idx=0)
7         self.rnn = nn.RNN(embed_dim,
8                           hidden_dim,
9                           batch_first=True)
10        self.fc = nn.Linear(hidden_dim, num_class)
11
12    def forward(self, x, length)
13        x = self.embedding(x)
14        x = pack_padded_sequence(x,
15                                lengths=length,
16                                enforce_sorted=False,
17                                batch_first=True)
18        x = self.rnn(x)
19        x = self.fc(x[-1])
20        return x
```

Transformer



Implementation in Pytorch

The model is composed of `nn.Embedding`, `nn.Transformer`, `nn.Linear`

`nn.Transformer` is composed of `nn.MultiHeadAttention` and other required modules

```
1 class TRANSFORMER(torch.nn.Module):
2     def __init__(self, vocab_size, embed_dim, hidden_dim, num_class):
3         super().__init__()
4         self.embedding = nn.Embedding(vocab_size,
5                                       embed_dim,
6                                       padding_idx=0)
7         self.transformer = nn.Transformer(dmodel,
8                                           nhead,
9                                           num_encoder_layers,
10                                          num_decoder_layers,
11                                          dim_feedforward)
12         self.fc = Linear(dmodel, vocab_size)
13
14     def forward(self, x, length)
15         x = self.embedding(x)
16         x = self.transformer(x)
17         x = self.fc(x[-1])
18         return x
```

```

1 class CNN(nn.Module):
2     def __init__(self):
3         super(CNN, self).__init__()
4         self.conv1 = nn.Conv2d(1, 16, 5, 1)
5         self.pool = nn.MaxPool2d(2, 2)
6         self.conv2 = nn.Conv2d(16, 32, 5, 1)
7         self.fc1 = nn.Linear(32 * 4 * 4, 1000)
8         self.fc2 = nn.Linear(1000, 1000)
9         self.fc3 = nn.Linear(1000, 10)
10
11     def forward(self, x):
12         x = self.pool(self.conv1(x))
13         x = self.pool(self.conv2(x))
14         x = x.view(-1, 32 * 4 * 4)
15         x = F.relu(self.fc1(x))
16         x = F.relu(self.fc2(x))
17         x = self.fc3(x)
18         return x

```

```

1 class RNN(torch.nn.Module):
2     def __init__(self, vocab_size, embed_dim, hidden_dim, num_class):
3         super(RNN, self).__init__()
4         self.embedding = nn.Embedding(vocab_size, embed_dim)
5
6         self.rnn = nn.LSTM(embed_dim, hidden_dim, batch_first=True)
7
8         self.fc = nn.Linear(hidden_dim, num_class)
9
10    def forward(self, x):
11        x = self.embedding(x)
12        x, _ = self.rnn(x)
13        x = self.fc(x[-1])
14        return x

```

```

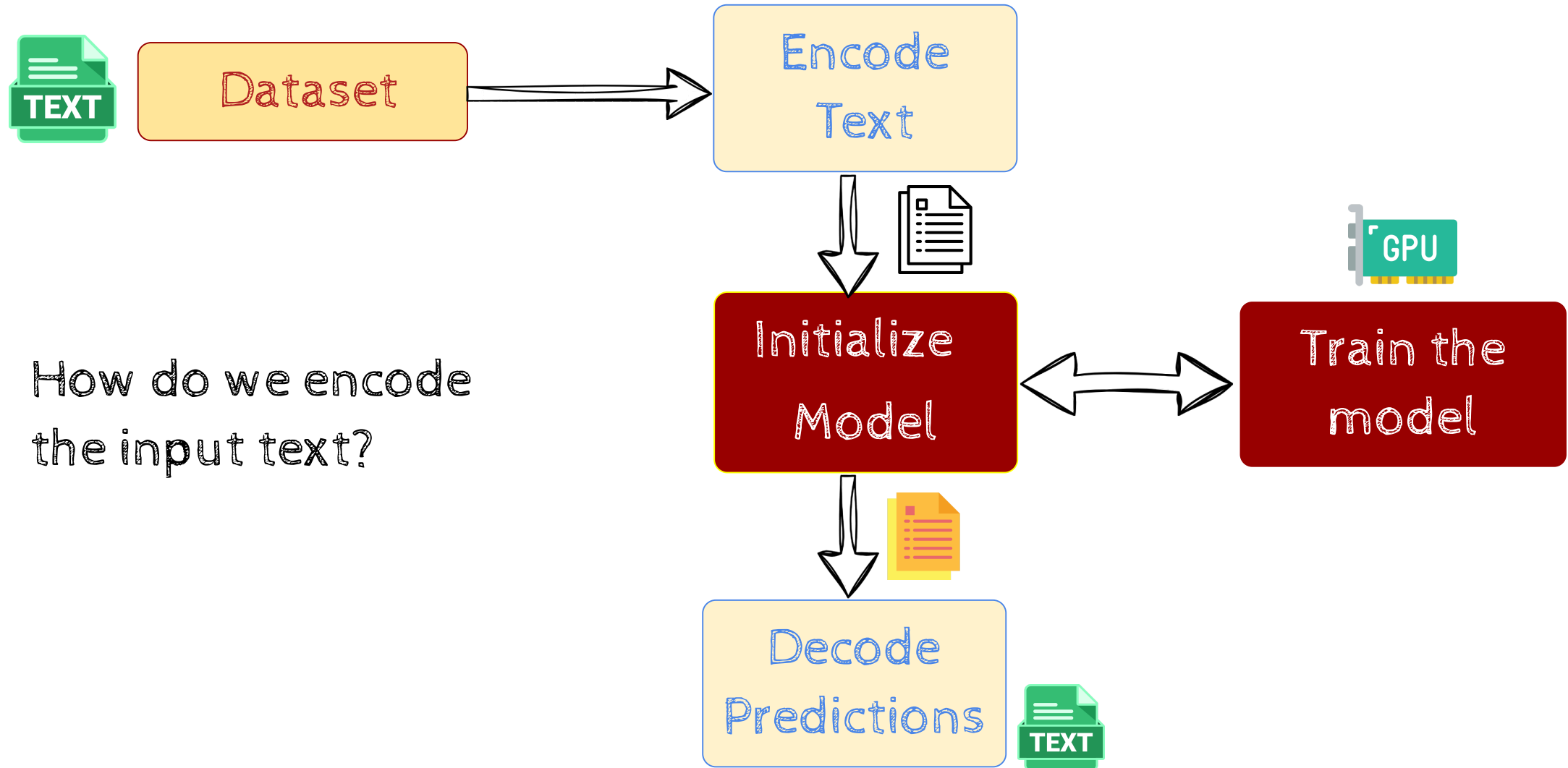
1 class TRANSFORMER(torch.nn.Module):
2     def __init__(self, vocab_size, embed_dim, hidden_dim, num_class):
3         super(TRANSFORMER, self).__init__()
4         self.embedding = nn.Embedding(vocab_size, embed_dim,
5                                       padding_idx=0)
6
7         self.transformer = nn.Transformer(dmodel=hidden_dim,
8                                           nhead=4,
9                                           num_encoder_layers=6,
10                                          num_decoder_layers=6,
11                                          dim_feedforward=2048)
12
13        self.fc = Linear(hidden_dim, vocab_size)
14
15    def forward(self, x, length):
16        x = self.embedding(x)
17        x = self.transformer(x)
18        x = self.fc(x[1])
19        return x

```

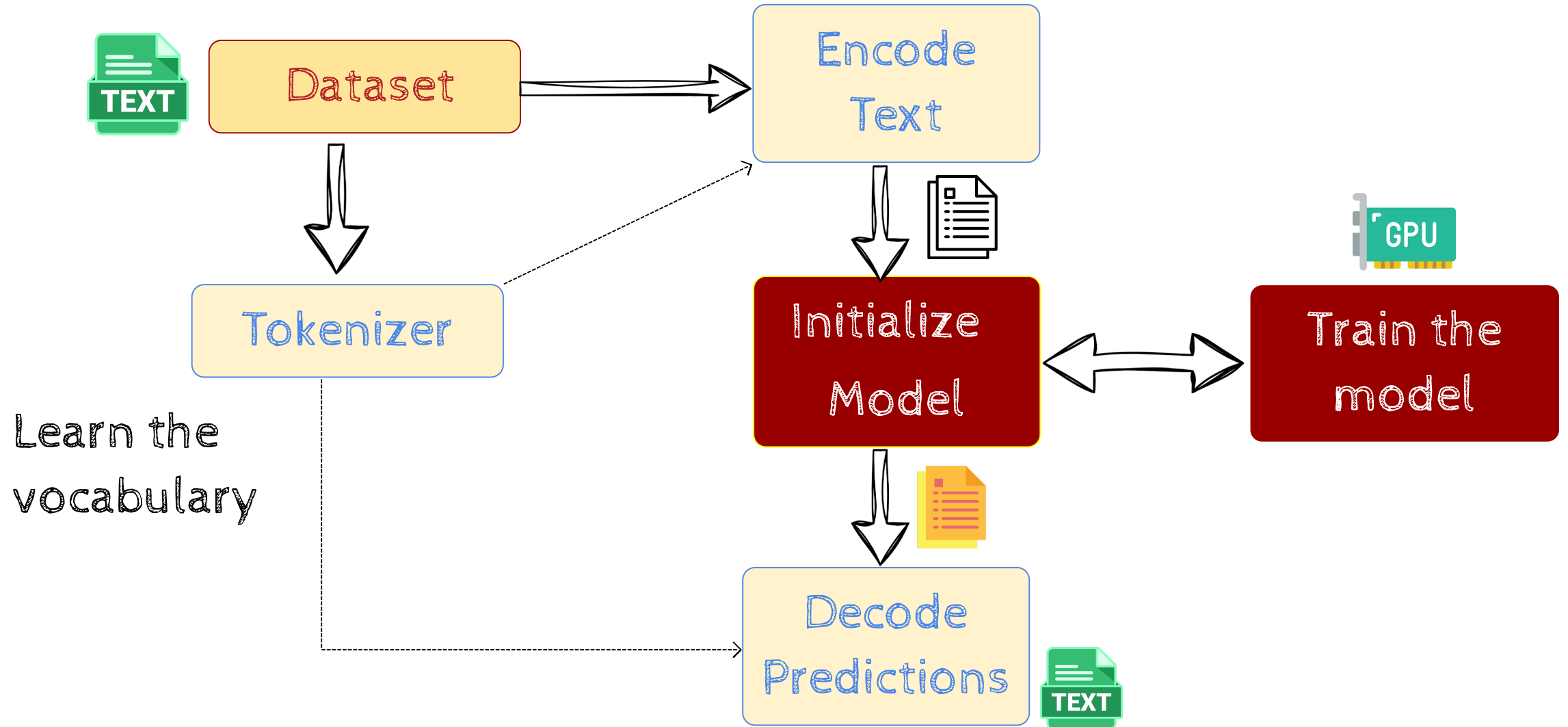
We can create any architecture with a few lines of code. However, this is just a part of the entire training setup.

The training setup contains many other components. Let's see

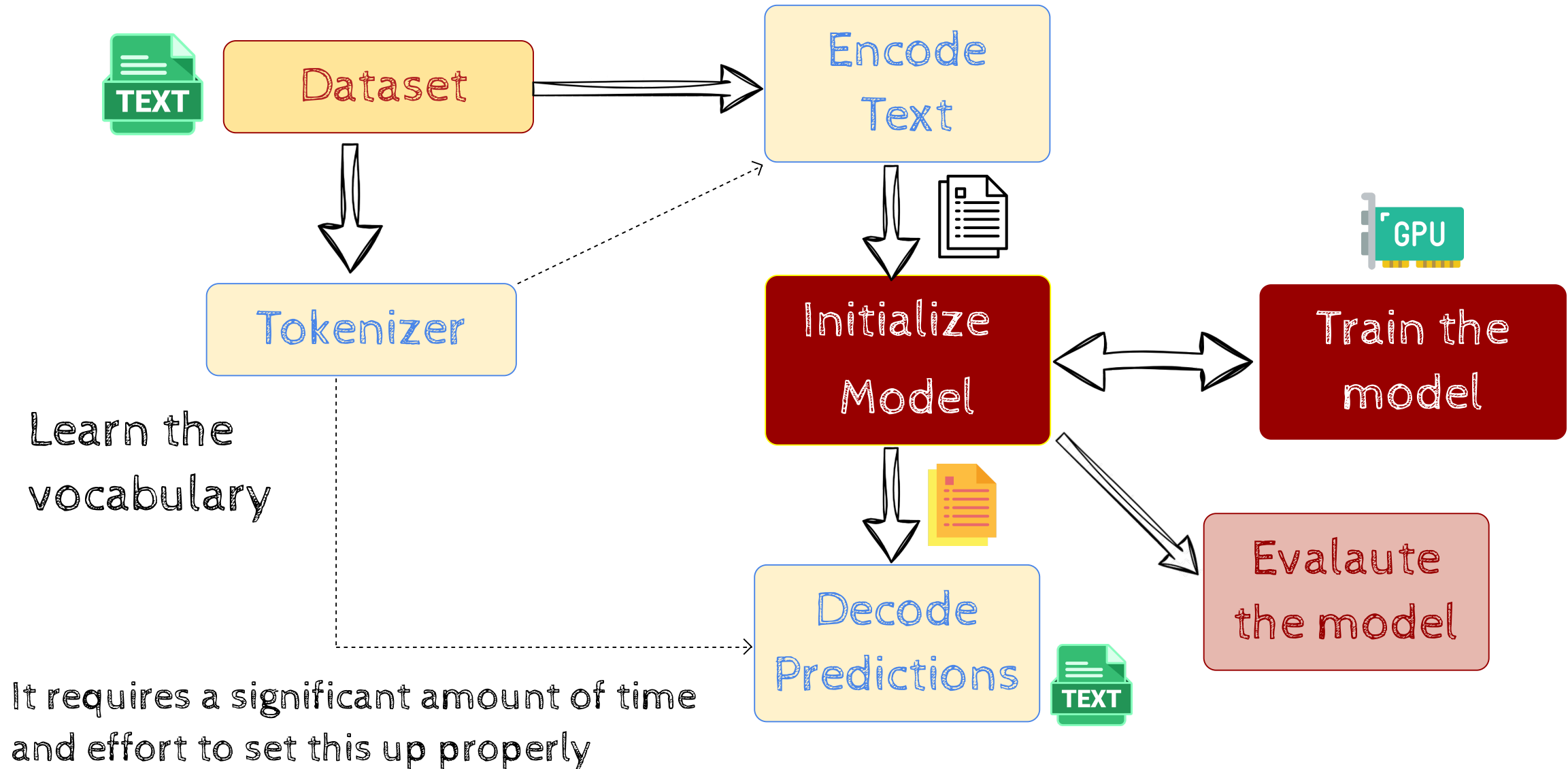
Experimental Setup



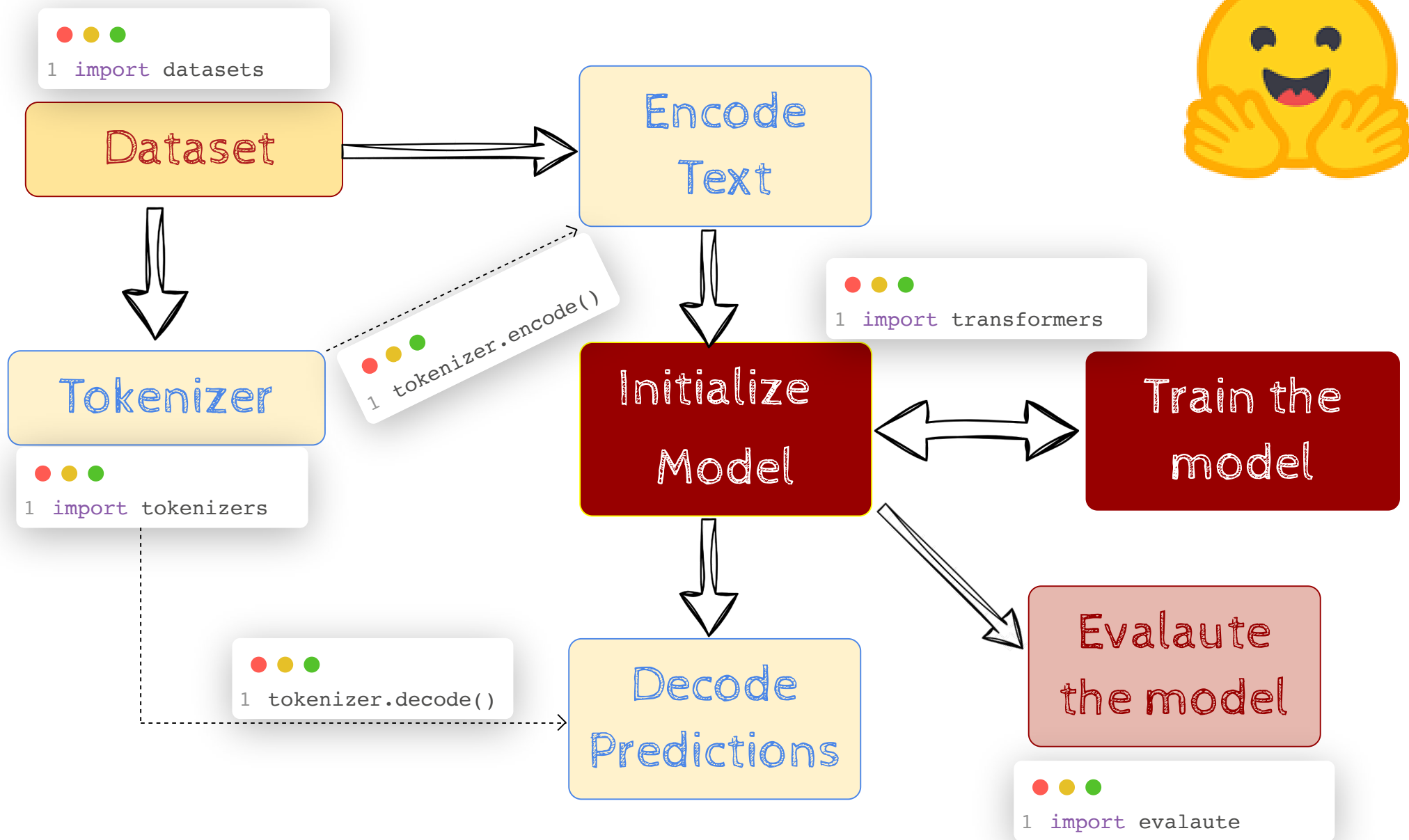
Experimental Setup



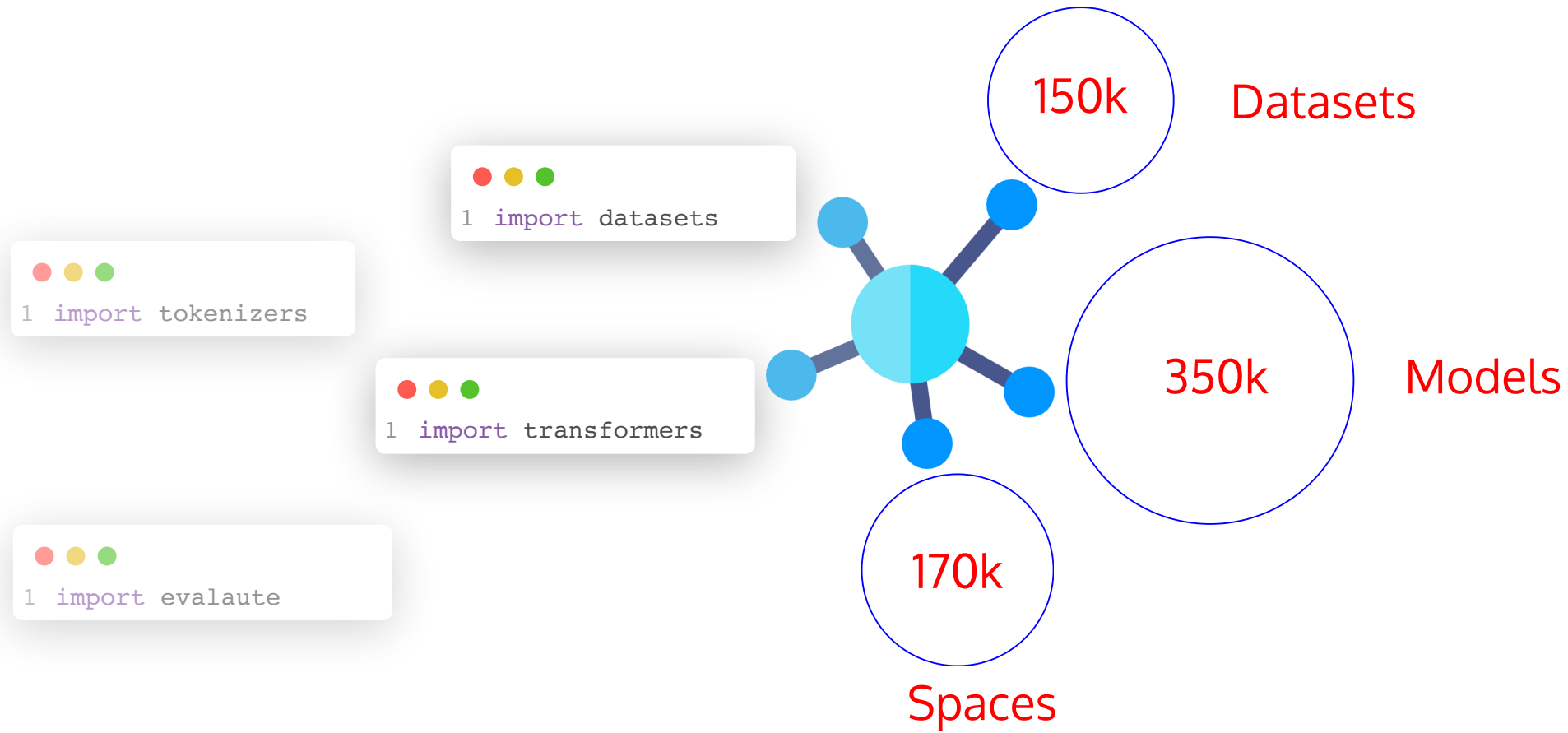
Experimental Setup



Hugging Face Does the Heavy Lifting




Hugging Face Hub



We have access to 150k datasets and 350k models via hugging face hub

Rich Ecosystem



```
1 import datasets
```




```
1 from accelerate ..
```



```
1 import tokenizers
```



```
1 from optimum ..
```




```
1 import transformers
```



```
1 from peft ..
```



```
1 import evalaute
```



```
1 import bitsandbytes
```

•
•
•

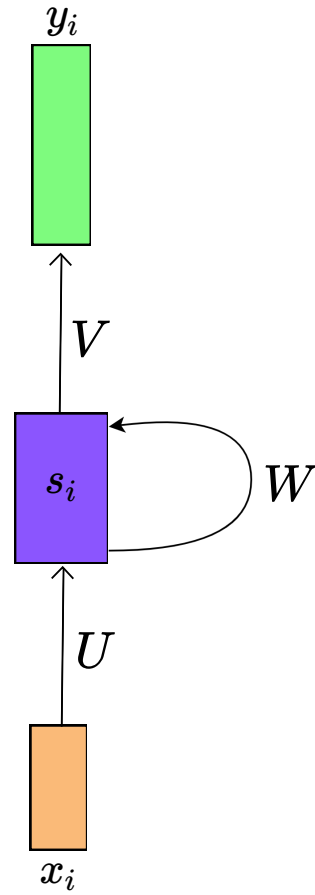
Module 2 : Roadmap

Mitesh M. Khapra



AI4Bharat, Department of Data Science
and Artificial Intelligence, IIT Madras

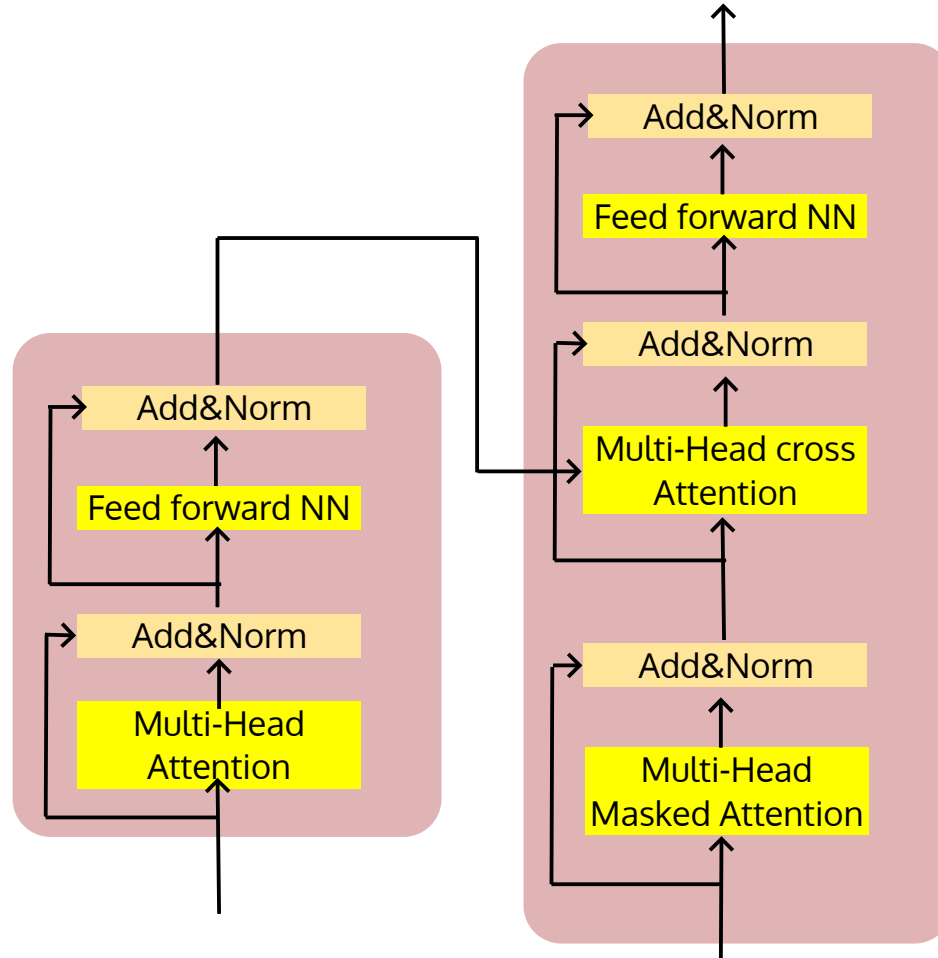
RNN based models for NLP



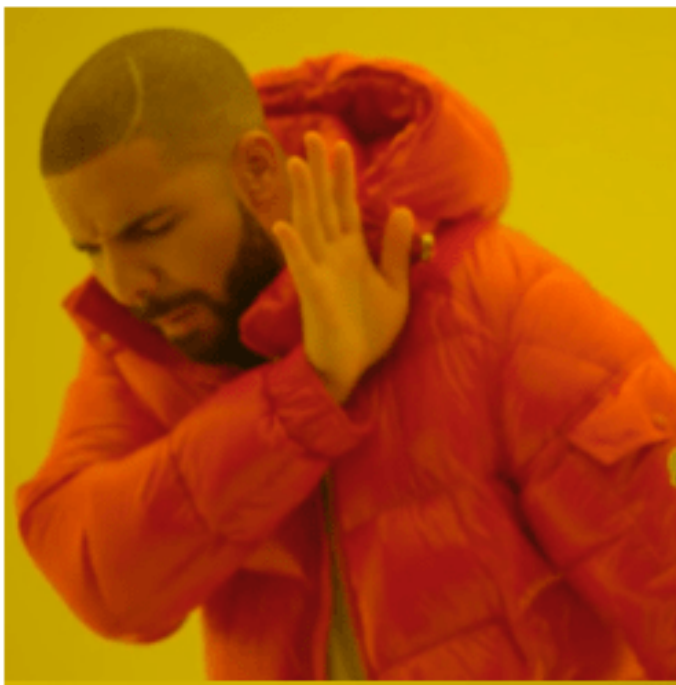
Paradigm shift



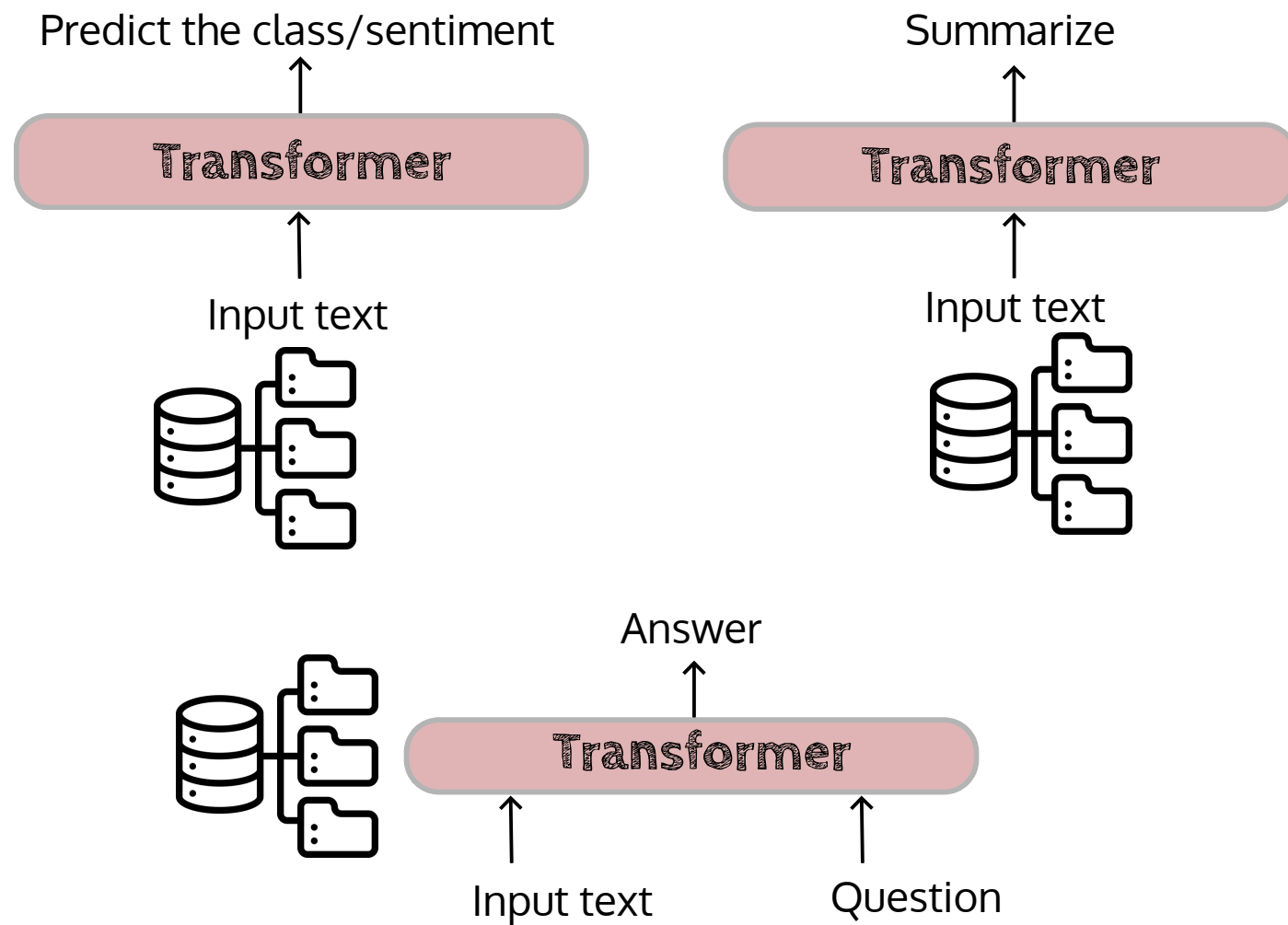
Transformer based models for NLP



Traditional NLP



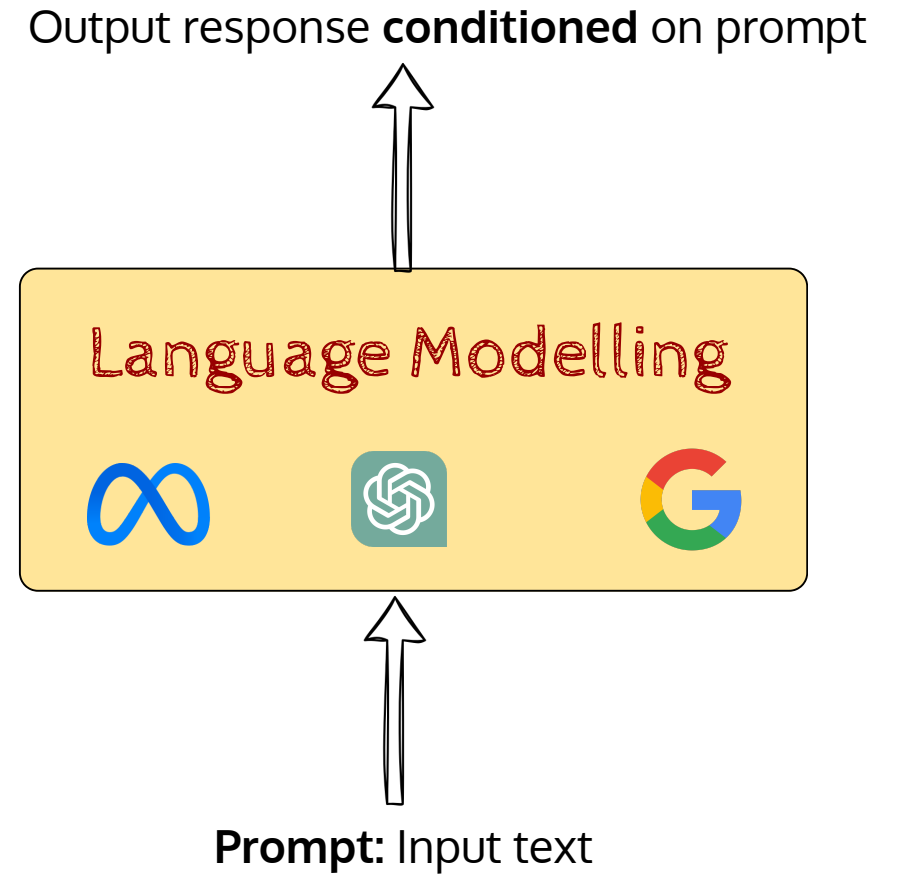
Task-specific Supervised Training



Modern NLP



Pre-training, Fine Tuning, Instruction tuning,



Prompt: Predict sentiment, summarize, fill in the blank, generate story

Syllabus Outline

```
1 import datasets
```

Dataset

Encode
Text

We will start with the datasets module and see how HF simplifies loading different datasets easily. We will also learn to use some commonly used functions from this module

The hf hub contains 150K+ datasets

```
1 import tokenizers
```

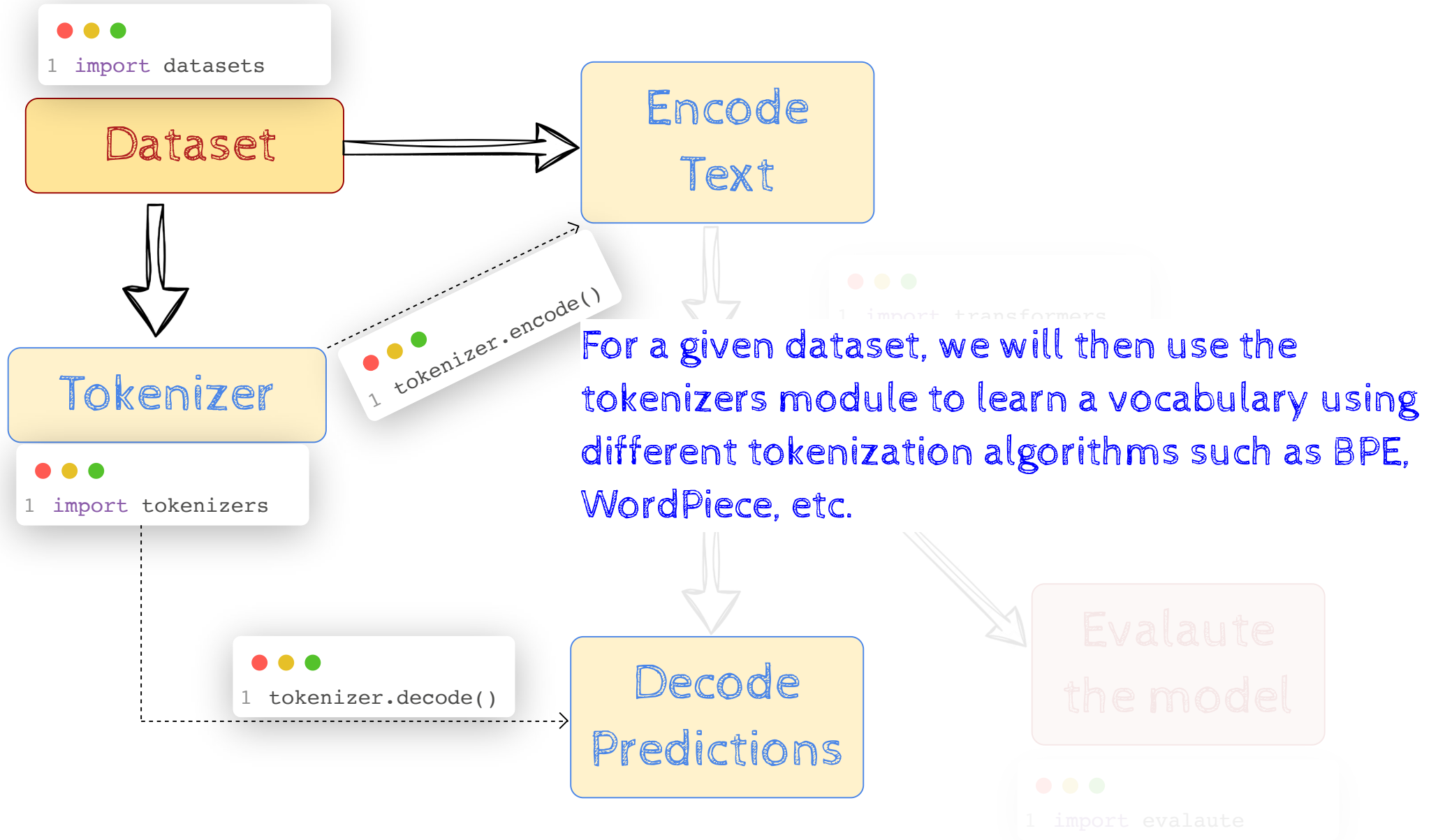
```
1 tokenizer.decode()
```

Decode
Predictions

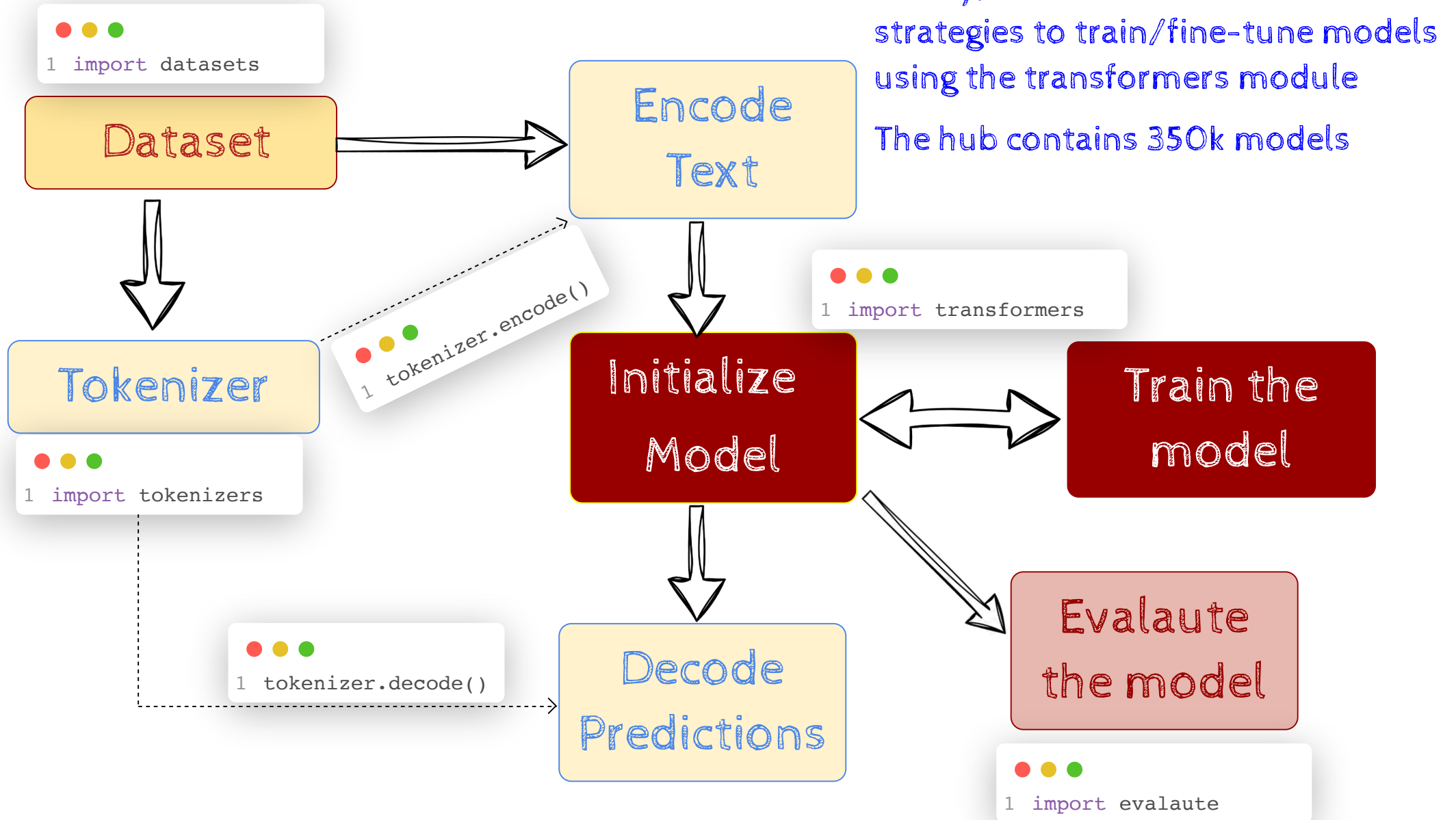
Evaluate
the model

```
1 import evaluate
```

Syllabus Outline



Syllabus Outline



In each week, we cover the required concepts to understand various terminologies that are helpful while developing a solution for NLP problems

We expect you to revisit the transformer architecture in detail as covered in the deep learning theory course

This week, we will delve deeper into the dataset part of the development cycle.

Module 3 : Datasets

Mitesh M. Khapra



AI4Bharat, Department of Data Science
and Artificial Intelligence, IIT Madras

Task Specific Datasets

Learning starts with a dataset

We have a wide range of NLP tasks such as

- Sentiment classification
- Machine Translation
- Named Entity Recognition
- Question Answering
- Textual entailment
- Summarization
- Generation

For each of these tasks, we may have hundreds of datasets with thousands of samples

For example, for sentiment classification we have

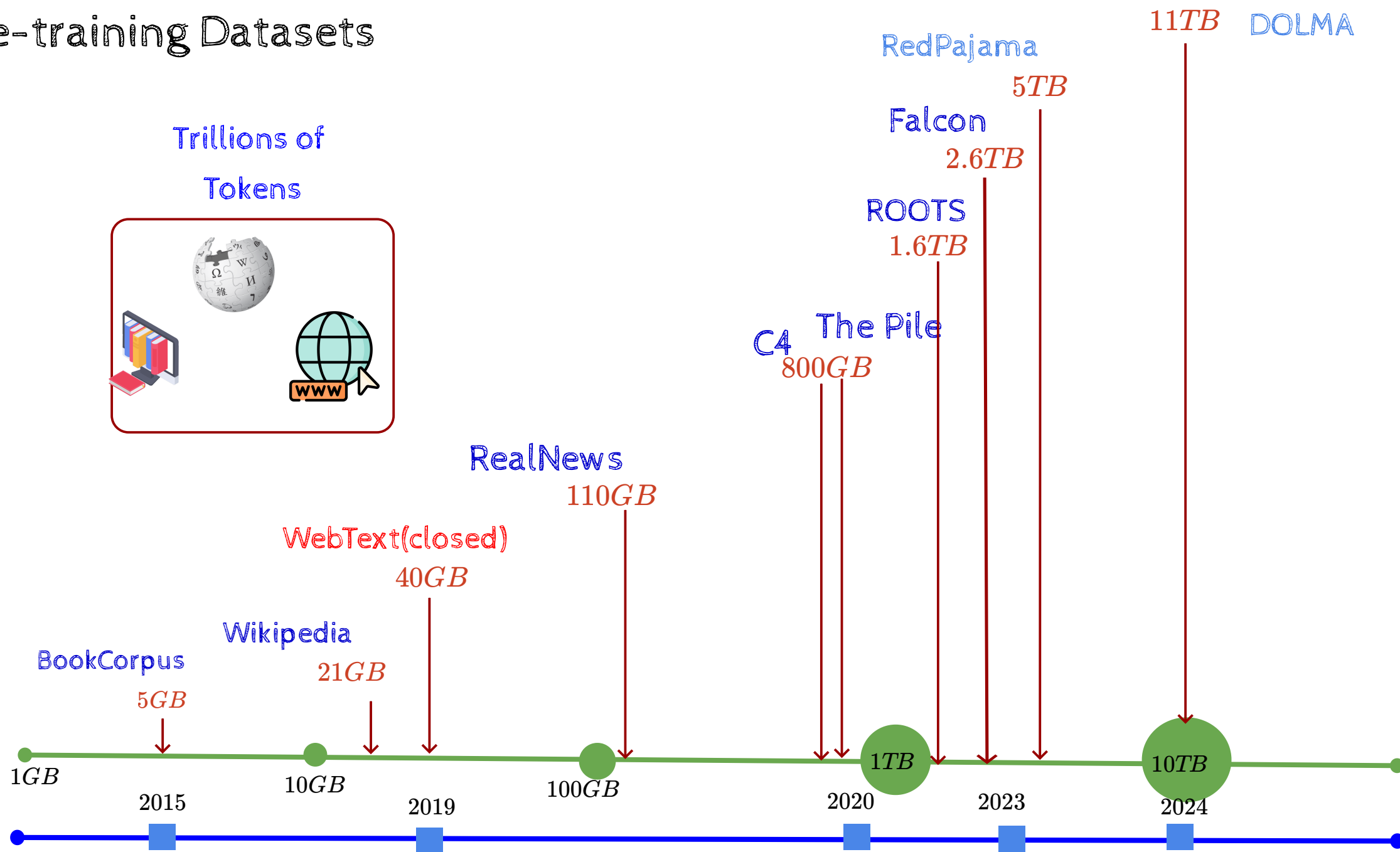
- Amazon review
- IMDB Movie review
- Twitter financial news sentiment
- poem sentiment

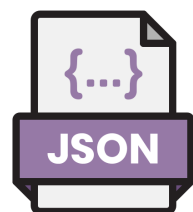
Similarly, one can find numerous datasets for other tasks

All these datasets typically have thousands to millions (if not billions) of words

In modern NLP, the pretraining datasets have billions and trillions of tokens

Pre-training Datasets





The datasets come with different formats and different sizes

We need to understand the format and may need to write a script to parse and load the text data

If the memory is limited, we need to implement mechanisms to stream the samples in the dataset .

However, it would be convenient if we have a single place to load any of these datasets with a consistent call signature

That's where Hugging Face's datasets module helps us

What is pre-training?

What are tokens?

Are tokens and the words in a sentence one and the same?

Well, we will discuss all these in subsequent lectures.

For now, let us just look at the datasets module from HF

