

# GEN-AI

## UNIT 1

### Lab 1

Name: Anshul Banda

SRN: PES2UG23CS081

Section: B

What is Hugging Face ?

Hugging Face is a forum where multiple models of AI are uploaded and shared for use among engineers.

The models uploaded to Hugging Face are pre-trained. This is done in order to save millions of dollars on training models from scratch and allow users to download models like GPT-2, BERT or RoBERTa directly from Hugging Face for use.

Transformers allow us to access these models by providing us with APIs to download and use these models. It supports framework interoperability, meaning you can often move between PyTorch, TensorFlow, and JAX

Pipelines allow to use these models by abstracting the complex math into three steps:

1. **Preprocessing:** Converts your raw text into numbers (Tokens & IDs) that the model can understand.
2. **Model Inference:** The model processes the numbers and outputs predictions (logits).
3. **Post-processing:** The raw predictions are converted back into human-readable text (labels, answers, summaries).

## HANDSON-1\_Unit1.ipynb:

The first step we perform in the .ipynb file is to download all the required transformers for the lab.

Next we upload the required .txt file for the lab.

First we perform a test between dumb and smart AI models. We compare **distilgpt2** and **gpt2** models.

Now we have to set a seed value for the AI models. A seed values sets the randomness of the AI. Same seeds means same random choices which means same outputs, different seeds means different random choices which means different outputs.

For a seed of value 42 and a prompt of “Generative AI is a revolutionary technology that”, the AI models produce the current output:

## Fast Model (distilgpt2):

```
/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:  
The secret 'HF TOKEN' does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your Google Colab and restart your session.  
You will be able to reuse this secret in all of your notebooks. 
```

Please note that authentication is recommended but still optional to access public models or datasets.

```
warnings.warn(  
    "warnings.warn(  
        config.json: 100% |██████████| 762/762 [00:00<00:00, 89.4kB/s]  
        model safetensors: 100% |██████████| 353M/353M [00:02<00:00, 231MB/s]  
        generation_config.json: 100% |██████████| 124/124 [00:00<00:00, 14.6kB/s]  
        tokenizer_config.json: 100% |██████████| 26.0/26.0 [00:00<00:00, 2.94kB/s]  
        vocab.json: 100% |██████████| 1.04M/1.04M [00:00<00:00, 1.22MB/s]  
        merges.txt: 100% |██████████| 456k/456k [00:00<00:00, 1.07MB/s]  
        tokenizer.json: 100% |██████████| 1.36M/1.36M [00:00<00:00, 9.83MB/s]  
    )  
Device set to use cuda:0  
Truncation was not explicitly activated but 'max_length' is provided a specific value, please use 'truncation=True' to explicitly truncate examples to max length. Defaulting to 'longest_first' truncation.  
Setting 'pad_token_id' to 'eos_token_id':50256 for open-end generation.  
Both 'max_new_tokens' (=256) and 'max_length' (=50) seem to have been set, 'max_new_tokens' will take precedence. Please refer to the documentation for more information. (https://huggingface.co/docs)  
Generative AI is a revolutionary technology that can take on the task of finding, learning, and learning in a given environment.
```

## Standard Model (gpt2):

```
config.json: 100%          665/665 [00:00:00.00, 69.2kB/s]
model.safetensors: 100%    548M/548M [00:03:00.00, 280kB/s] colab.research.google.com - To exit full screen, press Esc
generation_config.json: 100% 124/124 [00:00:00.00, 12.4kB/s]
tokenizer_config.json: 100% 26.02B.0 [00:00:00.00, 2.80kB/s]
vocab.json: 100%           1.04M/1.04M [00:00:00.00, 1.53MB/s]
merges.txt: 100%           456K/456K [00:00:00.00, 755kB/s]
tokenizer.onnx: 100%       1.36M/1.36M [00:00:00.00, 1.99MB/s]

Device set to use cuda:0
Truncation was not explicitly activated but 'max_length' is provided a specific value, please use 'truncation=True' to explicitly truncate examples to max length. Defaulting to 'longest_first' truncation.
Setting pad_token_id to eos_token_id=50256 for open-end generation.
Both `max_new_tokens` (=256) and `max_length` (=50) seem to have been set, `max_new_tokens` will take precedence. Please refer to the documentation for more information. (https://huggingface.co/docs)
Generative AI is a revolutionary technology that enables a wide range of intelligent systems to work independently from one another. It introduces a new way of thinking about AI and provides a new way of interacting with it.

In this article, we will discuss the main features of the new AI platform, and how it can be used to help us create a world that will improve our lives for the better.

1. How Can I Use It?

The concept of AI is not new. It has been used by many people to measure their mental health and health-related behaviors, and as a tool for medical research, it has been used by many of us to track our fitness levels and overall well-being.

It is based on the premise that AI is a way for humans to move towards a more efficient way of thinking, and therefore, a better way of living.

In this article, we will explain what AI can do.

What does it do

In this article, we will explain how all of our cognitive and emotional systems interact with the AI platform. The main features of AI are:

A new way of thinking about AI
A new paradigm for the development of intelligent AI
A new way of thinking about mental health and health-related behaviors
```

Now changing the seed value to 81:

## Fast Model (distilgpt2):

Device set to use cuda8  
Truncation was not explicitly activated but `max\_length` is provided a specific value, please use `truncation=True` to explicitly truncate examples to max length. Defaulting to `longest\_first` truncation.  
Setting `pad\_token\_id` to `eos\_token\_id`=50256 for open-end generation.  
Both `max\_new\_tokens` (-256) and `max\_length` (-50) seem to have been set. `max\_new\_tokens` will take precedence. Please refer to the documentation for more information. (<https://huggingface.co/docs>)  
Generative AI is a revolutionary technology that allows people to build complex, intelligent and innovative AI projects. It is the most powerful AI technology in the world and it will be the largest

## Standard Model (gpt2):

From the two models and the changes in the seed value we can see that the fast model provides us a short but straightforward response to given prompt. However, the standard model gives a longer response to the given prompt, but the standard model also has a tendency to just start drifting into generating nonsensical responses.

## NLP Fundamentals: Under the Hood

In order for models to process text, we first need to tokenize them. The pipeline does this for us.

Split a sample sentence "Transformers revolutionized NLP." into tokens.

```
tokens = tokenizer.tokenize(sample_sentence)
print(f"Tokens: {tokens}")

Tokens: ['Transform', 'ers', 'revolution', 'ized', 'N', 'LP', '.']
```

And then we assign them IDs:

```
▶ token_ids = tokenizer.convert_tokens_to_ids(tokens)
print(f"Token IDs: {token_ids}")

... Token IDs: [41762, 364, 5854, 1143, 399, 19930, 13]
```

## PART OF SPEECH(POS):

In order for models to understand the text, it needs to understand context of the text, the words in the text, and how the words in the text are used.

We can see how it recognises verbs, nouns and such:

```
▶ pos_tags = nltk.pos_tag(nltk.word_tokenize(sample_sentence))
print(f"POS Tags: {pos_tags}")

... POS Tags: [('Transformers', 'NNS'), ('revolutionized', 'VBD'), ('NLP', 'NNP'), ('.', '.')]
```

## Types of part of speech:

### # NOUNS

NN = Noun, singular or mass

NNS = Noun, plural

NNP = Proper noun, singular

NNPS = Proper noun, plural

### # PRONOUNS

PRP = Personal pronoun

PRP\$ = Possessive pronoun

### # VERBS

VB = Verb, base form

VBD = Verb, past tense

VBG = Verb, gerund / present participle

VBN = Verb, past participle

VBP = Verb, non-3rd person present

VBZ = Verb, 3rd person singular present

### # ADJECTIVES

JJ = Adjective

JJR = Adjective, comparative

JJS = Adjective, superlative

### # ADVERBS

RB = Adverb

RBR = Adverb, comparative

RBS = Adverb, superlative

WRB = Wh-adverb

### # DETERMINERS / ARTICLES

DT = Determiner

PDT = Predeterminer

WDT = Wh-determiner

# CONJUNCTIONS / PREPOSITIONS

CC = Coordinating conjunction

IN = Preposition / subordinating conjunction

# NUMBERS

CD = Cardinal number

# PARTICLES / SMALL WORDS

RP = Particle

TO = "to"

# MODALS

MD = Modal verb

# WH-WORDS

WP = Wh-pronoun

WP\$ = Possessive wh-pronoun

# INTERJECTIONS

UH = Interjection

# SYMBOLS / FOREIGN

SYM = Symbol

FW = Foreign word

LS = List item marker

# PUNCTUATION

. = Sentence-ending punctuation

, = Comma

: = Colon / semicolon

## Summarization: Efficiency vs. Quality

Now we test summarization of different types of models.

We check how well the AI can summarize the output based on this text:

"

The introduction of the Transformer architecture in the 2017 paper "Attention is all you need" was a watershed moment in AI. It provided a more effective and scalable way to handle sequential data like text, replacing older, less efficient methods like recurrence (RNNs) and convolutions.

The fundamental innovation of the Transformer is the attention mechanism. This component allows the model to weigh the importance of different words (tokens) in the input sequence when making a prediction. In essence, for each word it processes, the model can "pay attention" to all other words in the input, helping it understand context, resolve ambiguity, and handle long-range dependencies. This is crucial for tasks like translation, summarization, and question answering.

The Transformer architecture consists of an encoder stack (to process the input) and a decoder stack (to generate the output), both of which heavily utilize multi-head attention and feed-forward networks.

"

Fast Summarizer (**distilbart-cnn-12-6**):

Output:

The introduction of the Transformer architecture in the 2017 paper "Attention is all you need" was a watershed moment in AI . It provided a more effective and scalable way to handle sequential data like text, replacing older, less efficient methods like recurrence (RNNs) and conv

## Quality Summarizer (**bart-large-cnn**):

Output:

The introduction of the Transformer architecture in the 2017 paper "Attention is all you need" was a watershed moment in AI. It provided a more effective and scalable way to handle sequential data like text.

## Question Answering:

In this part, we will check if the model can answer our questions.

Based on the questions:

"What is the fundamental innovation of the Transformer?",

"What are the risks of using Generative AI?"

The output provided by AI model distilbert-base-cased-distilled-squad is:

```
▶  questions = [
    "What is the fundamental innovation of the Transformer?",
    "What are the risks of using Generative AI?"
]

for q in questions:
    res = qa_pipeline(question=q, context=text[:5000])
    print(f"\nQ: {q}")
    print(f"A: {res['answer']}")

...
Q: What is the fundamental innovation of the Transformer?
A: to identify hidden patterns, structures, and relationships within the data

Q: What are the risks of using Generative AI?
A: data privacy, intellectual property, and academic integrity
```

## Masked Language Modeling (The 'Fill-in-the-Blank' Game):

In this part, we will see if the AI model can fill in some blanks

```
masked_sentence = "The goal of Generative AI is to create new [MASK]."  
preds = mask_filler(masked_sentence)  
  
for p in preds:  
    print(f"{p['token_str']}: {p['score']:.2f}")  
  
...  
applications: 0.06  
ideas: 0.05  
problems: 0.05  
systems: 0.04  
information: 0.03
```

From this we observe the probabilities of words the AI thinks are correct for this task.