

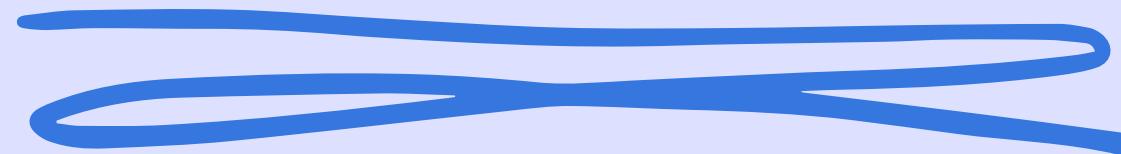


**ANL488**

**Business Analytics**

**Applied Project**

# Image to Audio Generator



## Via the Transformer Model Network



**Presented by:**

*Koh Wei Jie Edmund*

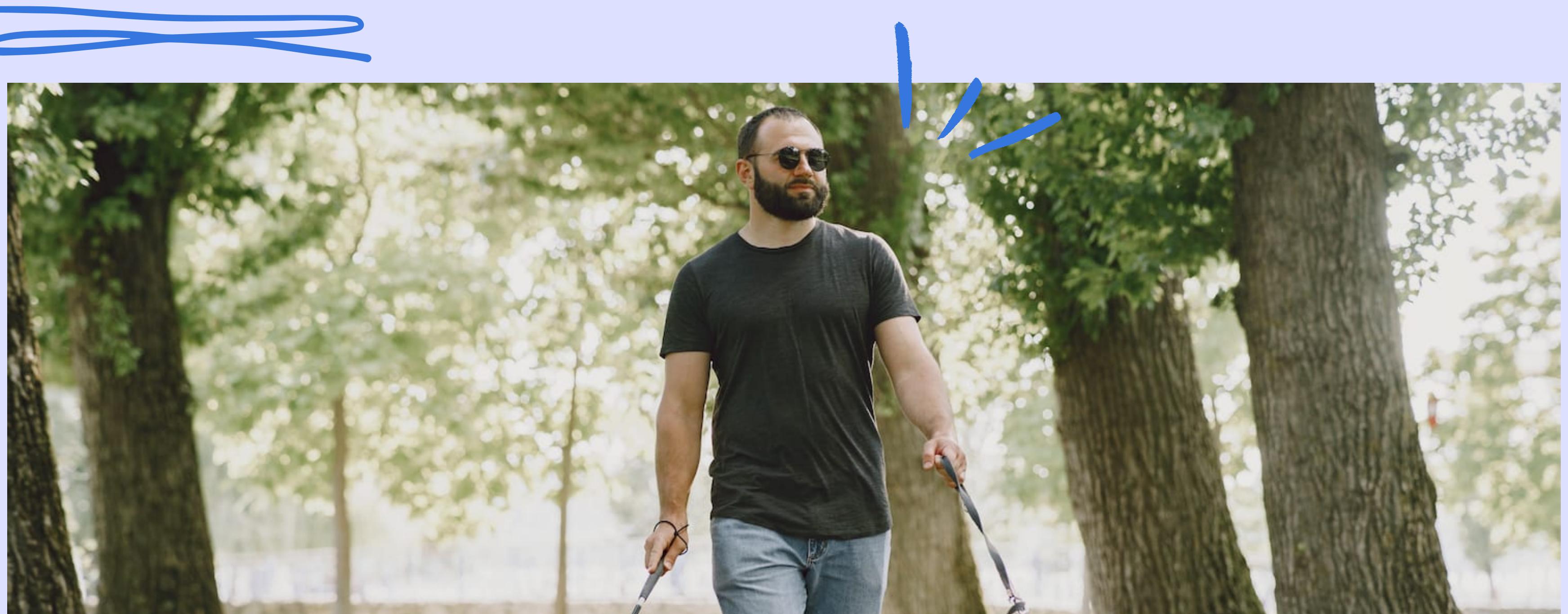


# **Agenda**

1. Introduction
2. Objective
3. Literature Review
4. Model Exploration
5. Model Evaluation
6. Proposed Modelling
7. Tool Optimisation
8. Recommendation

# *Motivation*

Image-to-Audio Generator Via the Transformer Model Network





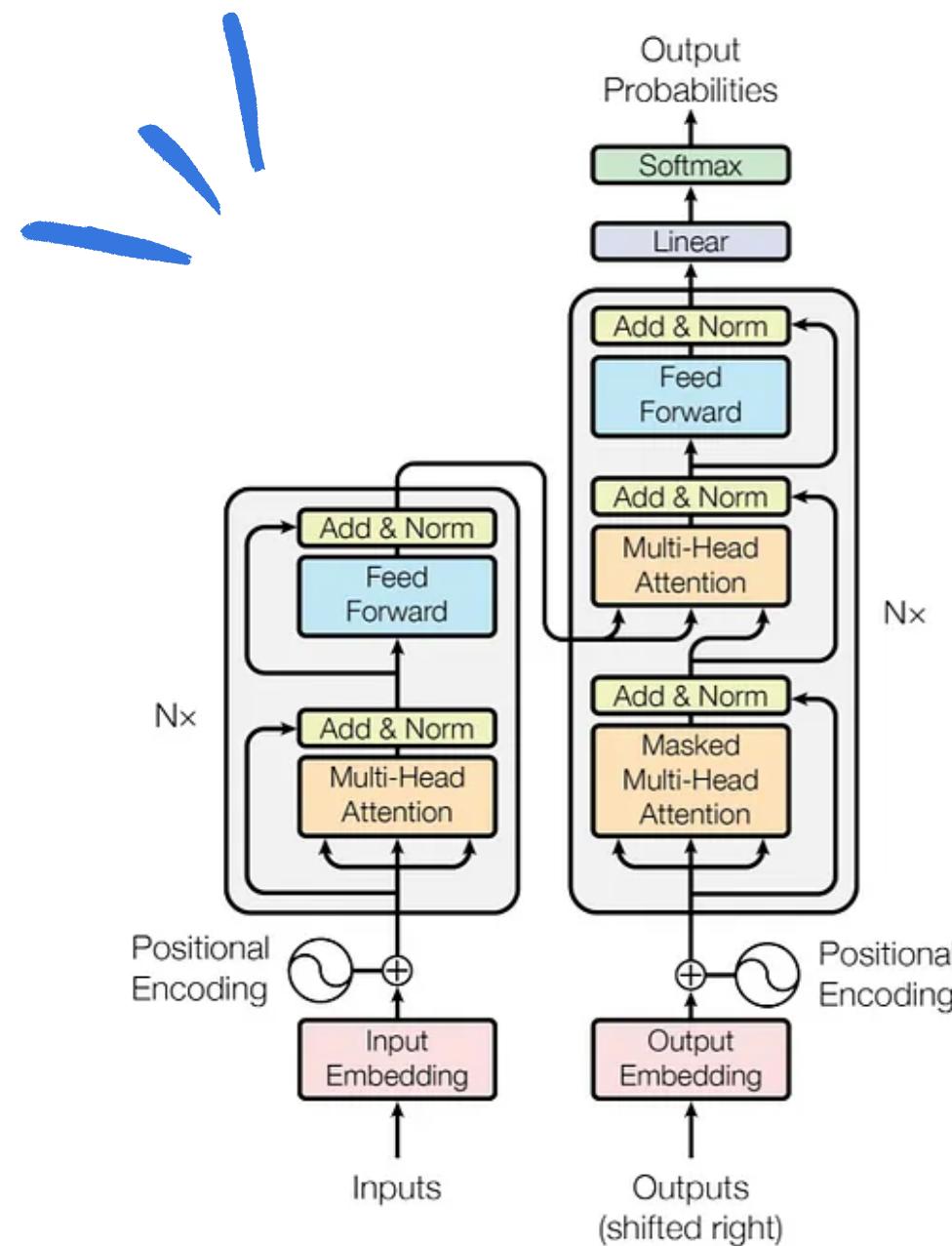
## ***Business Problem***

- The global cost of productivity losses due to vision impairment is projected to be US\$ 411 billion.

## ***Analytics Objective***

- Develop a state-of-the-art image captioning model that would accurately depict the details of an image

# *What is Transformer?*



\* “Attention Is All You Need,” published in 2017

\* No Labels, More Performance

\* State-of-the-art technique in the field of NLP

# *Literature Review (Related Work)*

**Dosovitskiy et al. (2020)**

ViT is found to be the best performer when the computational cost of pre-training the model is taken into account reaching state-of-the-art on the majority of recognition benchmarks at a lower pre-training cost.

**Vaswani et al. (2017)**

Transformer model used has achieved a BLEU score of 28.4 which outperform the best-performing model including the ensembles model by more than 2.0 BLEU while using only a fraction of the training cost of others.

**Ghandi et al. (2023)**

The future of image captioning would revolve around Vision-language pre-training (VLP) methods and Transformers as they are likely to be inseparable components of the models.

# Model Selection



Tasks 1 Libraries Datasets Languages Licenses Other

Filter Tasks by name  Reset Tasks

Multimodal

- Feature Extraction
- Text-to-Image
- Image-to-Text
- Text-to-Video
- Visual Question Answering
- Document Question Answering
- Graph Machine Learning

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Image-to-Image
- Unconditional Image Generation
- Video Classification
- Zero-Shot Image Classification

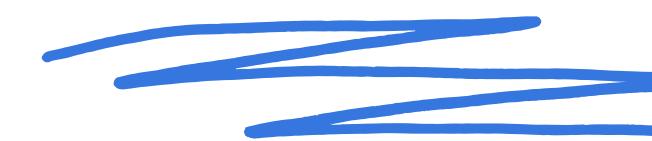
Natural Language Processing

- Text Classification
- Token Classification
- Table Question Answering
- Question Answering
- Zero-Shot Classification
- Translation
- Summarization
- Conversational
- Text Generation
- Text2Text Generation
- Fill-Mask
- Sentence Similarity

Models 266 Filter by name  new Full-text search Sort: Most Downloads

<a href="#">nlpconnect/vit-gpt2-image-captioning</a> Image-to-Text • Updated Feb 27 • 1.92M • 548	<a href="#">Salesforce/blip-image-captioning-large</a> Image-to-Text • Updated Aug 1 • 1.17M • 359
<a href="#">Salesforce/blip-image-captioning-base</a> Image-to-Text • Updated Aug 1 • 506k • 224	<a href="#">microsoft/trocr-base-handwritten</a> Image-to-Text • Updated Jan 26 • 333k • 93
<a href="#">microsoft/git-large-coco</a> Image-to-Text • Updated Jun 27 • 146k • 65	<a href="#">Salesforce/blip2-opt-2.7b</a> Image-to-Text • Updated 14 days ago • 136k • 138
<a href="#">microsoft/trocr-large-printed</a> Image-to-Text • Updated Jan 25 • 123k • 31	<a href="#">Salesforce/blip2-flan-t5-xl</a> Image-to-Text • Updated 12 days ago • 118k • 29
<a href="#">microsoft/trocr-large-stage1</a> Image-to-Text • Updated Apr 1 • 116k • 9	<a href="#">fxmarty/pix2struct-tiny-random</a> Image-to-Text • Updated Jun 1 • 91.3k • 1
<a href="#">laion/mscoco_finetuned_CoCa-ViT-L-14-laion2B-s13B-b...</a> Image-to-Text • Updated Jun 11 • 84.9k • 14	<a href="#">Salesforce/instructblip-vicuna-7b</a> Image-to-Text • Updated Jul 17 • 76.6k • 36
<a href="#">microsoft/trocr-small-handwritten</a> Image-to-Text • Updated Jan 25 • 49.9k • 15	<a href="#">kha-white/manga-ocr-base</a> Image-to-Text • Updated Jun 22, 2022 • 46.9k • 47
<a href="#">bipin/image-caption-generator</a> Image-to-Text • Updated Jul 5 • 41.3k • 8	<a href="#">microsoft/git-base</a> Image-to-Text • Updated Apr 24 • 38.7k • 23

# How to load a Transformer



## </> How to use from the /transformers library

X

```
# Use a pipeline as a high-level helper
from transformers import pipeline

pipe = pipeline("image-to-text", model="Salesforce/blip-image-captioning-large")
```

Copy

```
# Load model directly
from transformers import AutoProcessor, AutoModelForSeq2SeqLM

processor = AutoProcessor.from_pretrained("Salesforce/blip-image-captioning-large")
model = AutoModelForSeq2SeqLM.from_pretrained("Salesforce/blip-image-captioning-large")
```

Copy

# *Model Selection*

```
vit_gpt2 = pipeline("image-to-text", model="nlpconnect/vit-gpt2-image-captioning")
```

```
blip_base = pipeline("image-to-text", model="Salesforce/blip-image-captioning-base")
```

```
blip_2 = pipeline("image-to-text", model="Salesforce/blip2-opt-2.7b")
```

# *Comparison of Models*

```
vit_gpt2(URL1, max_new_tokens=100)[0]['generated_text']
```

'a refrigerator freezer sitting in a kitchen next to a counter '



```
blip_base(URL1, max_new_tokens=100)[0]['generated_text']
```

'a white refrigerator sitting in a kitchen next to a sink'

```
blip_2(URL1, max_new_tokens=100)[0]['generated_text']
```

'a refrigerator sitting in a kitchen with cabinets\n'

# *Comparison of Models*



```
vit_gpt2(URL2, max_new_tokens=100)[0]['generated_text']
```

```
'a bench on a mountain with a view of the mountains'
```

```
blip_base(URL2, max_new_tokens=100)[0]['generated_text']
```

```
'a bench on a rocky mountain with a view of a mountain'
```

```
blip_2(URL2, max_new_tokens=100)[0]['generated_text']
```

```
'a bench sitting on a hill\n'
```

# *Comparison of Models*



```
vit_gpt2(URL3, max_new_tokens=100)[0]['generated_text']
```

'a large jetliner flying through a cloudy sky '

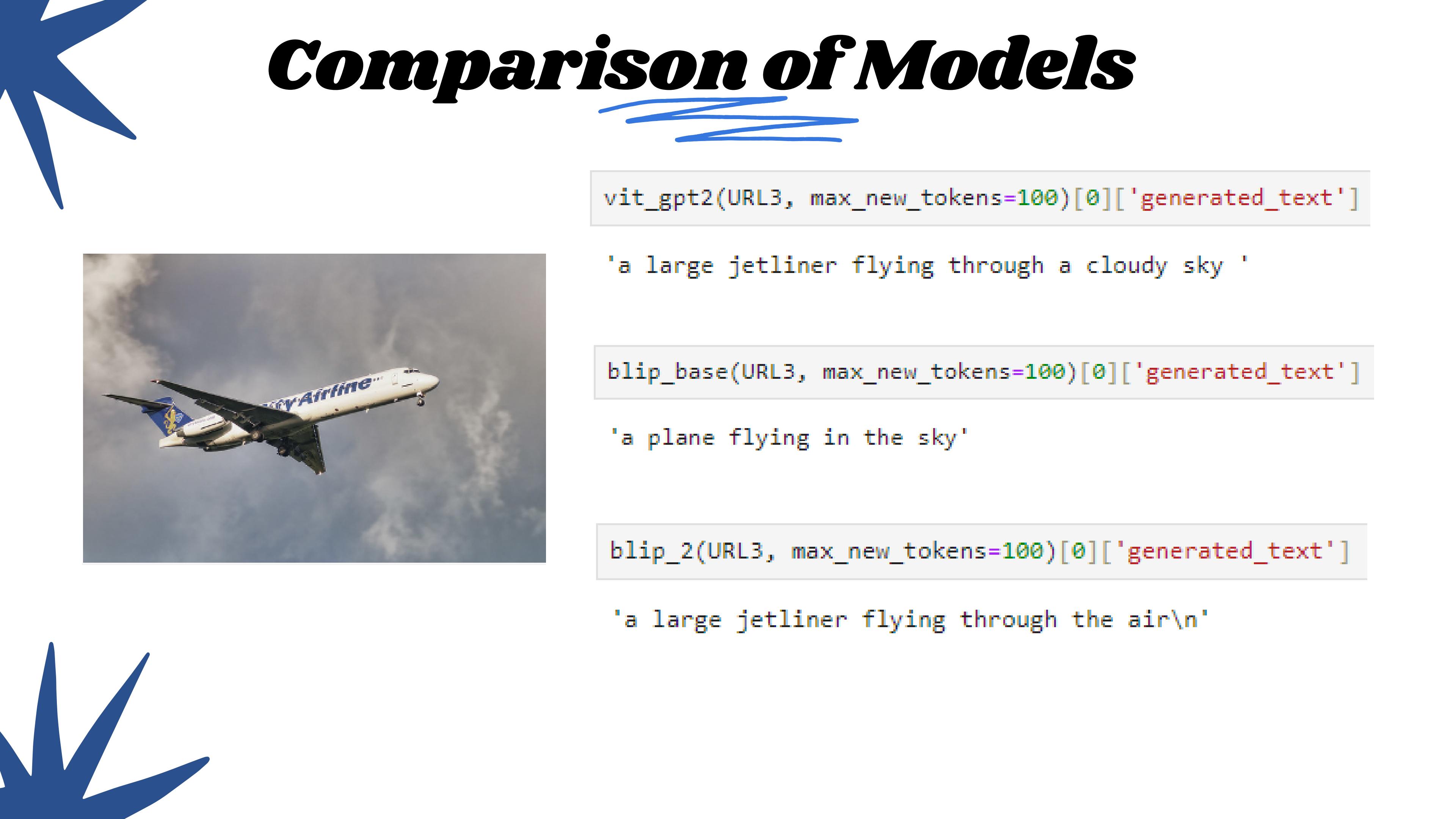


```
blip_base(URL3, max_new_tokens=100)[0]['generated_text']
```

'a plane flying in the sky'

```
blip_2(URL3, max_new_tokens=100)[0]['generated_text']
```

'a large jetliner flying through the air\n'



# *Model Evaluation (BLEU)*

```
# Calculating BLEU Score using the frige example
reference = [
    ['this', 'is', 'the', 'back', 'of', 'a', 'freezer', 'in', 'the', 'middle', 'of', 'a', 'kitchen', 'floor'],
    ['the', 'appliance', 'in', 'the', 'kitchen', 'may', 'be', 'broken', '.'],
    ['the', 'back', 'of', 'a', 'refrigerator', 'in', 'a', 'kitchen'],
    ['a', 'refrigerator', 'sitting', 'in', 'the', 'middle', 'of', 'a', 'kitchen', 'floor'],
    ['the', 'back', 'of', 'a', 'refrigerator', 'standing', 'in', 'the', 'middle', 'of', 'a', 'kitchen', '.'],
]

vit_candidate = ['a', 'refrigerator', 'freezer', 'sitting', 'in', 'a', 'kitchen', 'next', 'to', 'a', 'counter']

vit_score = sentence_bleu(reference, vit_candidate)
print(vit_score)
```

# *Model Evaluation (BLEU)*

```
vit_candidate = ['a', 'refrigerator', 'freezer', 'sitting', 'in', 'a', 'kitchen', 'next', 'to', 'a', 'counter']
```

```
vit_score = sentence_bleu(reference, vit_candidate)
print(vit_score)
```

```
5.008605395783359e-78
```

```
print('Cumulative 1-gram: %f' % sentence_bleu(reference, vit_candidate, weights=(1, 0, 0, 0)))
print('Cumulative 2-gram: %f' % sentence_bleu(reference, vit_candidate, weights=(0.5, 0.5, 0, 0)))
print('Cumulative 3-gram: %f' % sentence_bleu(reference, vit_candidate, weights=(0.33, 0.33, 0.33, 0)))
print('Cumulative 4-gram: %f' % sentence_bleu(reference, vit_candidate, weights=(0.25, 0.25, 0.25, 0.25)))
```

```
Cumulative 1-gram: 0.636364
```

```
Cumulative 2-gram: 0.504525
```

```
Cumulative 3-gram: 0.308321
```

```
Cumulative 4-gram: 0.000000
```

# Model Evaluation (BLEU)

```
blip_base_candidate = ['a', 'white', 'refrigerator', 'sitting', 'in', 'a', 'kitchen', 'next', 'to', 'a', 'sink']
```

```
blip_base_score = sentence_bleu(reference, blip_base_candidate)
print(blip_base_score)
```

```
5.731095145962094e-78
```

```
print('Cumulative 1-gram: %f' % sentence_bleu(reference, blip_base_candidate, weights=(1, 0, 0, 0)))
print('Cumulative 2-gram: %f' % sentence_bleu(reference, blip_base_candidate, weights=(0.5, 0.5, 0, 0)))
print('Cumulative 3-gram: %f' % sentence_bleu(reference, blip_base_candidate, weights=(0.33, 0.33, 0.33, 0)))
print('Cumulative 4-gram: %f' % sentence_bleu(reference, blip_base_candidate, weights=(0.25, 0.25, 0.25, 0.25)))
```

```
Cumulative 1-gram: 0.545455
```

```
Cumulative 2-gram: 0.467099
```

```
Cumulative 3-gram: 0.368341
```

```
Cumulative 4-gram: 0.000000
```

# *Model Evaluation (BLEU)*

```
blip_2_candidate = ['a', 'refrigerator', 'sitting', 'in', 'a', 'kitchen', 'with', 'cabinets\n']
```

```
blip_2_score = sentence_bleu(reference, blip_2_candidate)
print(blip_2_score)
```

```
0.4810977290978808
```

```
print('Cumulative 1-gram: %f' % sentence_bleu(reference, blip_2_candidate, weights=(1, 0, 0, 0)))
print('Cumulative 2-gram: %f' % sentence_bleu(reference, blip_2_candidate, weights=(0.5, 0.5, 0, 0)))
print('Cumulative 3-gram: %f' % sentence_bleu(reference, blip_2_candidate, weights=(0.33, 0.33, 0.33, 0)))
print('Cumulative 4-gram: %f' % sentence_bleu(reference, blip_2_candidate, weights=(0.25, 0.25, 0.25, 0.25)))
```

```
Cumulative 1-gram: 0.750000
```

```
Cumulative 2-gram: 0.731925
```

```
Cumulative 3-gram: 0.647453
```

```
Cumulative 4-gram: 0.481098
```

# *Human Evaluation*



**VIT\_GPT2**

a refrigerator freezer sitting in a kitchen next to a counter



**BLIP\_BASE**

a white refrigerator sitting in a kitchen next to a sink



**BLIP\_2**

a refrigerator sitting in a kitchen with cabinets\n



a bench on a mountain with a view of the mountains



a bench on a rocky mountain with a view of a mountain



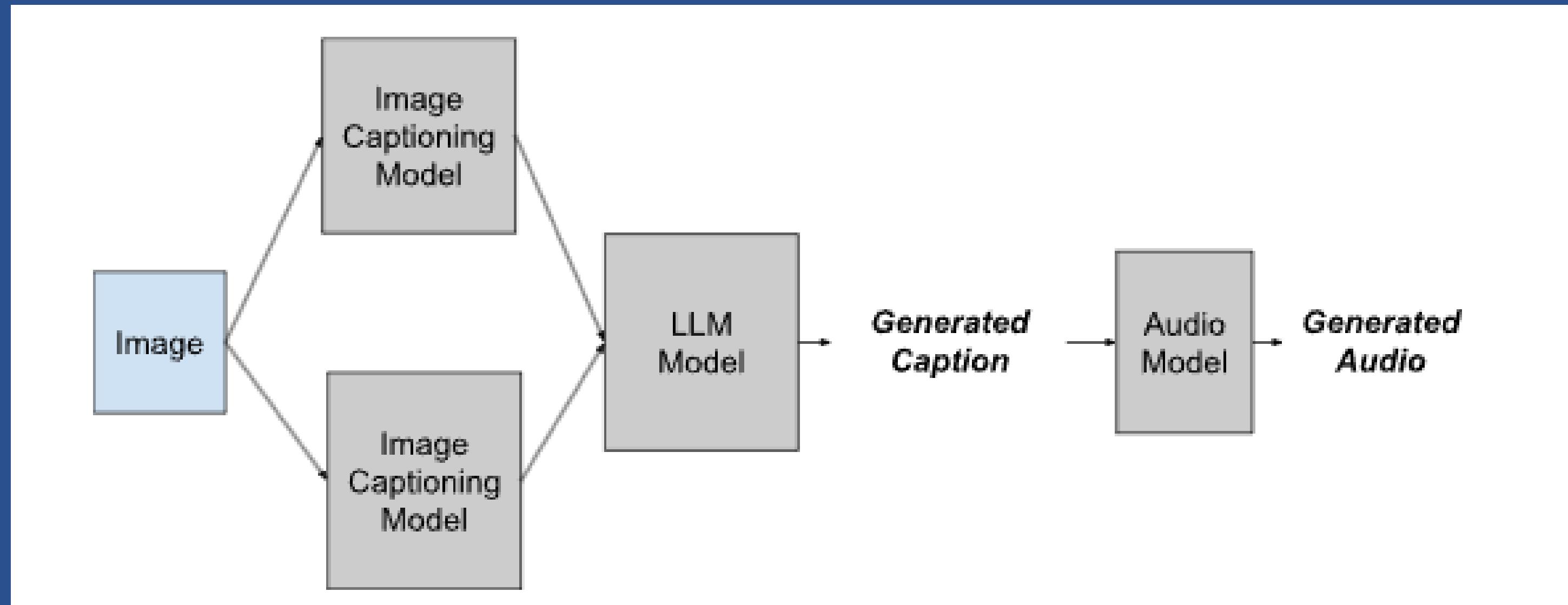
a bench sitting on a hill\n

a large jetliner flying through a cloudy sky

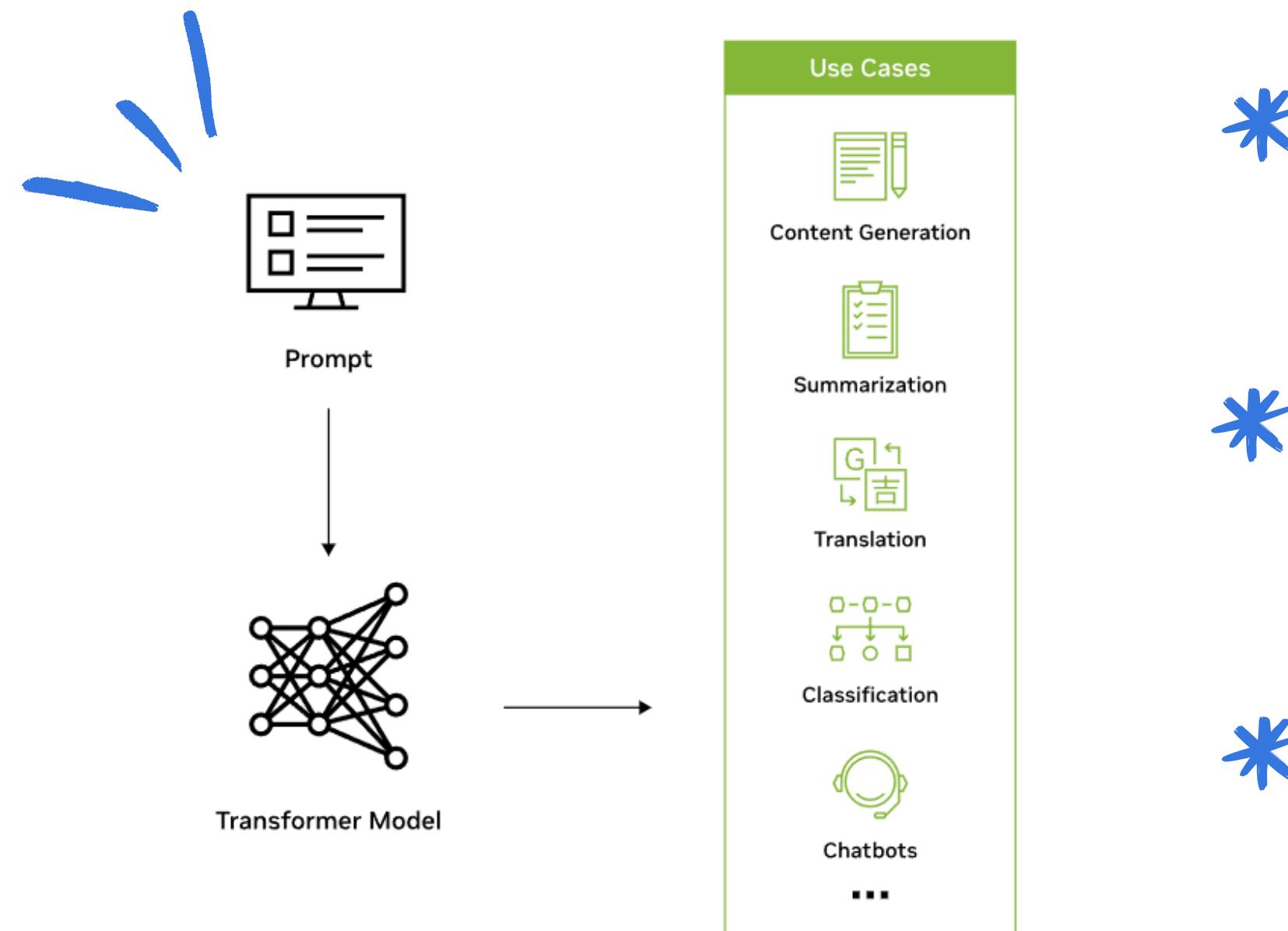
a plane flying in the sky

a large jetliner flying through the air\n

# *Proposed Pipeline*



# *Why use Text-Davinci-003?*

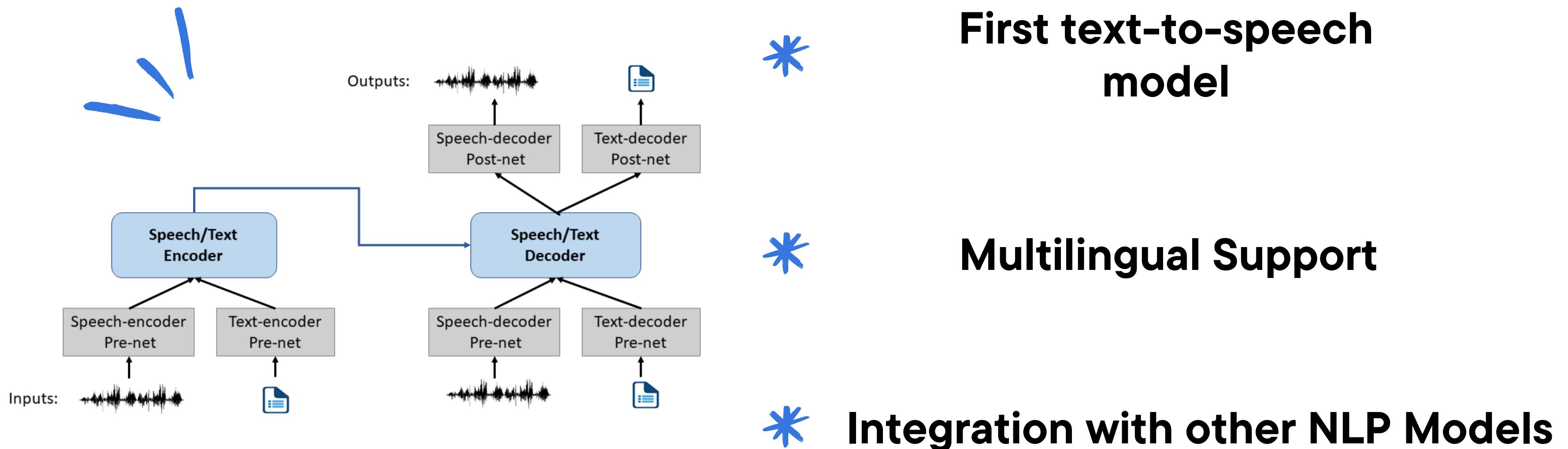


**Produces higher quality writing**

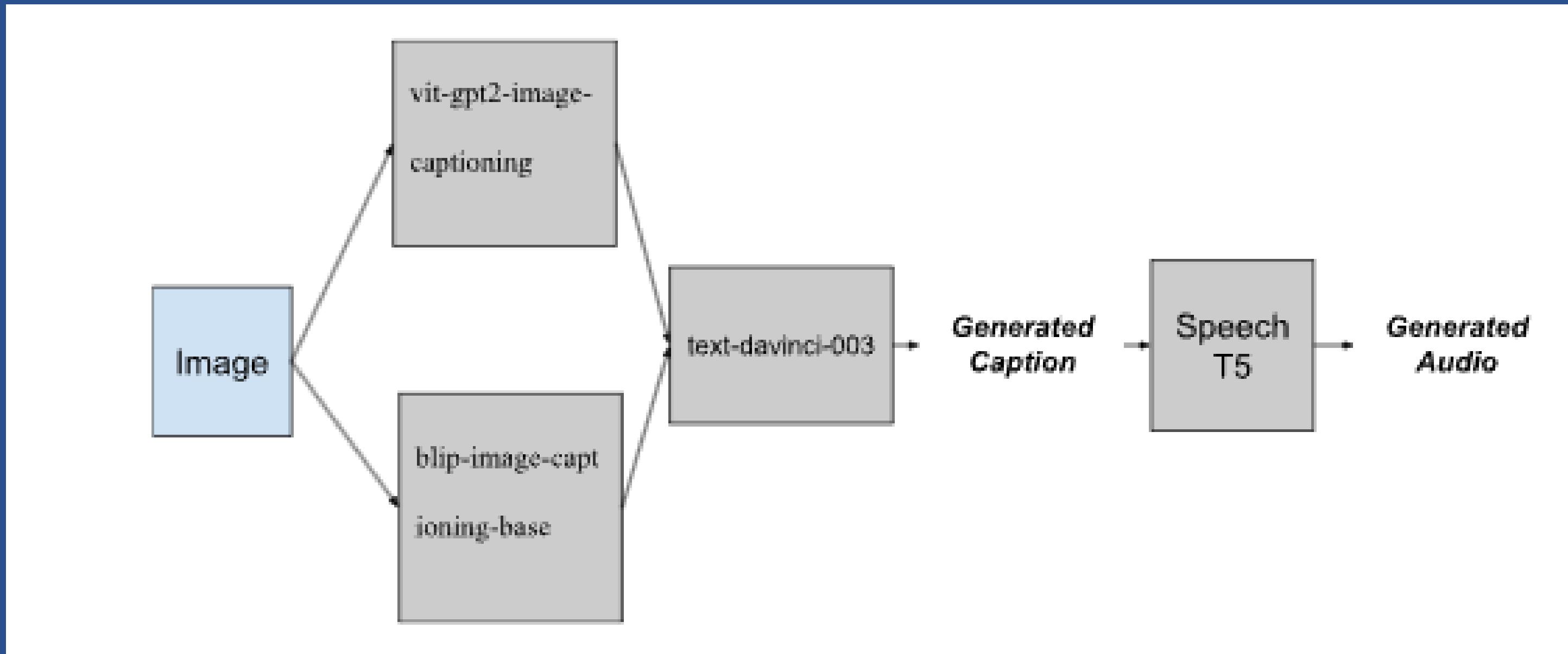
**Handle more complex instructions**

**Better at longer form content generation**

# *Why use Speech-T5?*



# *Finalised Pipeline*



# ***How do we go about Fine-tuning?***

---



- \* Pre-Training
- \* Prompt-Engineering

# Promp Engineering

## PROMPT USED



**"Combine both text:" +  
vit\_gpt2\_output + "and "  
blip\_base\_output**

A white refrigerator freezer is  
sitting in a kitchen next to a  
counter and a sink.

lake  
AI: A bench on a rocky mountain  
with a view of both the mountains  
and a mountain lake.

A large jetliner flying through a  
cloudy sky and a plane flying in  
the sky.

**"Combine and refine the text:"  
+ vit\_gpt2\_output + "and "  
blip\_base\_output**

A white refrigerator freezer is  
situated in a kitchen next to both  
a counter and a sink.

A bench on a rocky mountain  
with a view of both the mountains  
and the surrounding landscape.

A large jetliner and a plane flying  
through a cloudy sky.

**"Combine both:" +  
vit\_gpt2\_output\_example1 +  
"and "  
blip\_base\_output\_example1 + "  
Please refine the text."**

A white refrigerator freezer sitting  
in a kitchen next to both a  
counter and a sink.

A bench on a rocky mountain  
with a view of both the mountains  
and the surrounding landscape.

A large jetliner and a plane flying  
through a cloudy sky.

# Promp Engineering

PROMPT USED			
<p>"You are an AI assistant that combines information. The following are the outputs:" + vit_gpt2_output_example1 + "and " + blip_base_output_example1 + "Please refine the text."</p>	<p>A white refrigerator freezer sitting in a kitchen next to both a counter and a sink.</p>	<p>A bench on a rocky mountain with a view of both the mountains and the surrounding landscape.</p>	<p>I can see a large jetliner and a plane flying through the cloudy sky.</p>
<p>"Act as a human, provide output that a human would say to a visually impaired from the combination of both text:" + vit_gpt2_output + "and " + blip_base_output</p>	<p>I can see a white refrigerator freezer sitting in a kitchen next to both a counter and a sink.</p>	<p>I can see a bench on a rocky mountain with a spectacular view of both the mountains and the landscape around it.</p>	<p>I can see a large jetliner and a plane soaring through the cloudy sky, their wings cutting through the clouds as they gracefully fly through the sky.</p>
<p>"Act as a human, be as elaborative as possible, combined and refined the text output that a human would say to a visually impaired:" + vit_gpt2_output_example1 + "and " + blip_base_output_example1</p>	<p>I can see a white refrigerator freezer sitting in a kitchen, it's right next to a counter on one side and a sink on the other.</p>	<p>I can see a bench on a rocky mountain with a beautiful view of the distant mountains and the rolling hills of the surrounding landscape.</p>	<p>I can see a large jetliner and a plane gracefully soaring through the cloudy sky, the jetliner's wings cutting through the clouds as they fly through the sky.</p>

# **Demo**

**Loading ...**

## *Future Work*

- \* **Further Prompt Engineering**
- \* **Pre-training with a customized dataset**
- \* **Creation of a new pipeline**

***Thank you for your time!***

---

Any Questions?

