

Initial attempt at an image captioning model

Summary:

Basic model architecture involves a ViT (vision transformer) model to transform raw image data into meaningful features and GPT2 to generate captions based on these features. Base code to train the model was sourced from [1].

Initially the model was trained on 2000 images (took 3.5 hours) from the coco dataset [2] and then a fully trained model (also trained on the same dataset [2]) was sourced [3] to see how good a fully trained model is at image captioning and how well a partially trained model can estimate what an image includes. Then the images that the client's model did well with vs struggled on was also included to understand what the client's model is capable of. I've come up with 2-3 bullet points at the end for us to try on our next attempt.

Datasets were downloaded and models trained locally on a Yoga Slim ProX 16GB RAM.

Testing model with coco data set images:



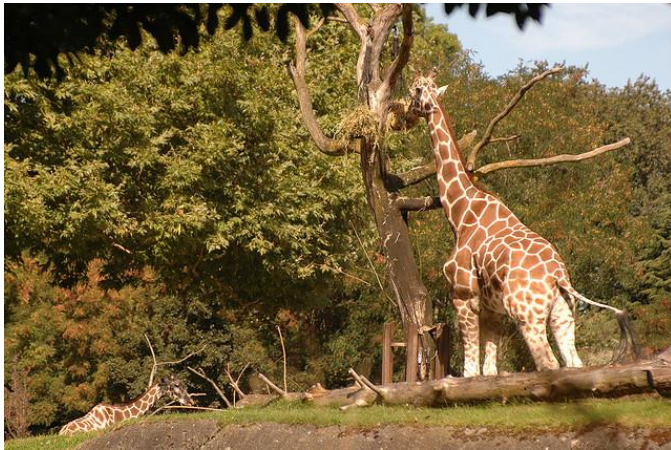
Model trained on complete dataset: a soccer game with a player jumping to catch the ball

Partially trained model: A man is holding a baseball bat in the air



Fully trained model: a lunch box filled with different types of food

Partially trained model: A man is holding a sandwich on a table.



Fully trained model: a giraffe standing next to a tree in a field

Partially trained model: A giraffe is grazing on a grassy field



Fully trained model: a bathroom with a tub, sink, and mirror

Partially trained model: A bathroom with a shower and sink.



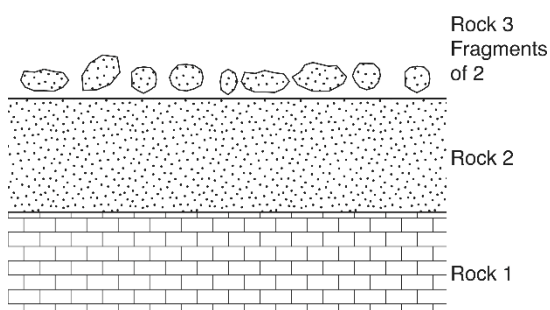
Fully trained model: a train on a track near a fence

Partially trained model: A train pulling into a station

Evaluation:

- Partially trained model was able to accurately describe some of the test images
- When it failed, it understood which category the images belonged to i.e. saying sandwich when the image had other foods displayed
- Fully trained model got all the images correct since these images are very similar to what the model was trained on
- Partially trained model was used to see how well the model can predict images when the training set is very limited for real world applications – more testing needs to be done in this arena, but our time is better spent in testing for harder cases.

Testing model with images the client's model could accurately classify:



Fully trained model: a series of photographs showing a series of small, white objects

Partially trained model: A man standing in front of a large mirror

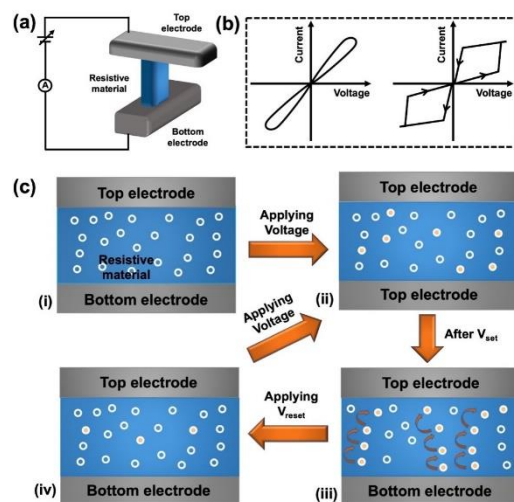
Client model: a diagram of a rock wall with a line of rocks on it



Fully trained model: a rocky cliff with a large rock wall

Partially trained model: A man riding a skateboard on a snowy mountain

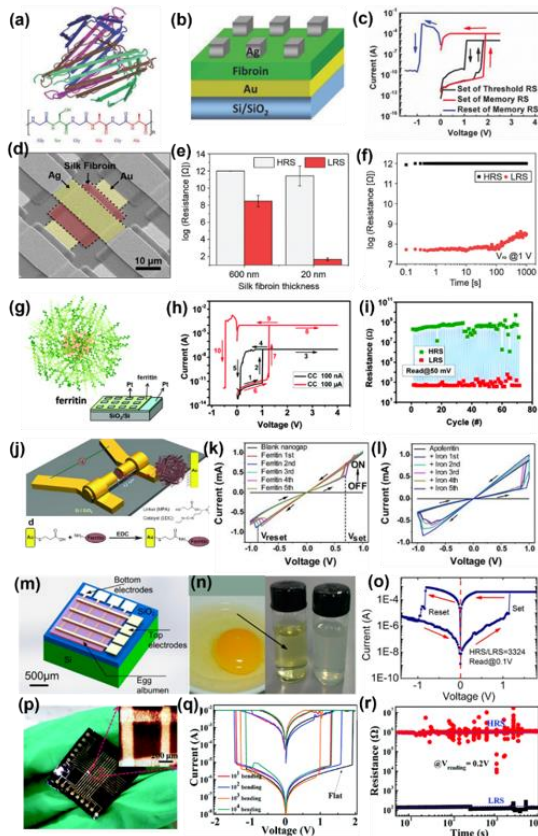
Client model: there is a large rock formation with a hole in it



Fully trained model: a series of photos showing different types of computers

Partially trained model: A man standing in front of a large mirror

Client model: a diagram of a top electrode and bottom electrode are shown



Fully trained model: a collage of photos of various types of objects

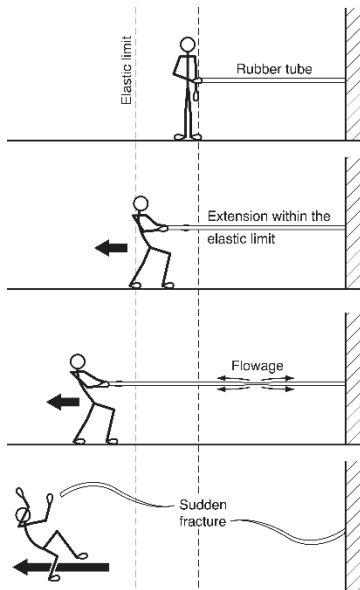
Partially trained model: A large wooden table with a clock and a clock tower

Client model: a bunch of pictures of different types of materials and processes

Evaluation:

- Our fully trained model failed to recognise what the white and black diagram represented, only mentioned colours and shapes whereas the clients model succeeded – needs to be trained on more diagrams.
- Both our fully trained model and client's model can accurately describe landscapes which is not true of our partially trained model however it recognised that the general category of the image was landscapes
- Fully trained model can successfully identify when it is given a collage of images, but it struggles to identify what scientific technique is being shown in this
- For fully trained model to match the performance of the client's model, it needs to be trained on textbook diagrams – it can even exceed the clients model in terms of performance if we included OCR models.

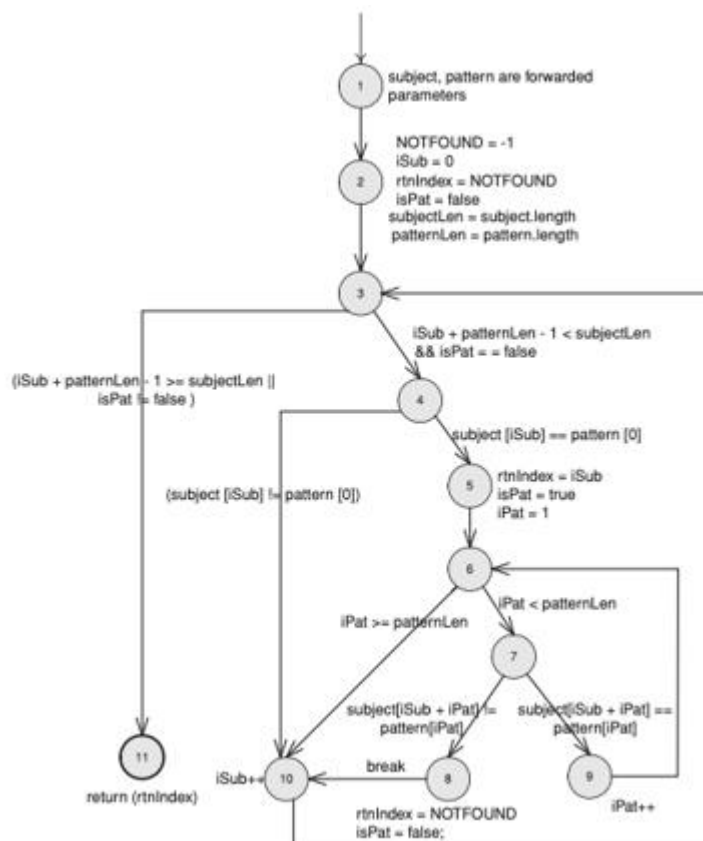
Testing model with tricky client images:



Fully trained model: a collage of photos showing a woman cutting a pair of scissors

Partially trained model: A man standing in front of a large building with a clock

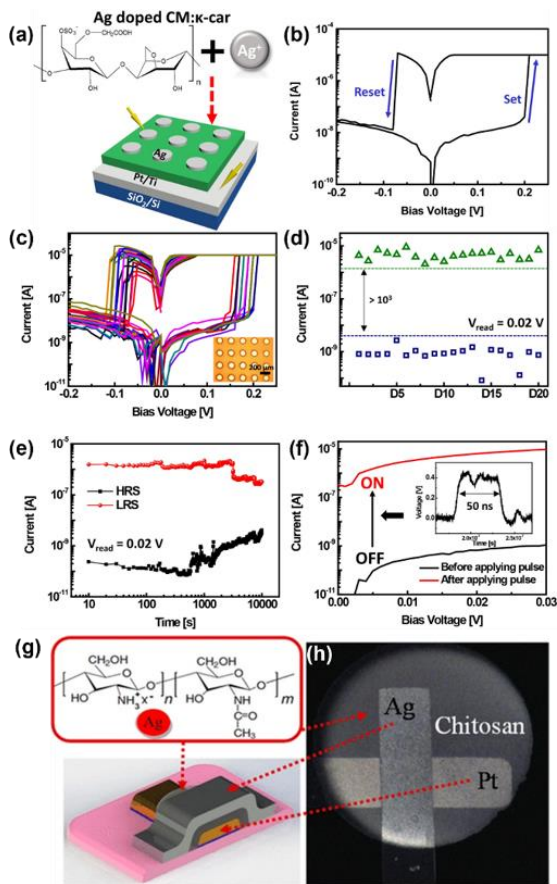
Client model: a diagram of a man doing a trick on a skateboard



Fully trained model: a series of photographs showing a computer screen

Partially trained model: A man standing in front of a large mirror

Client model: a diagram of a tree with a number of different types of trees



Fully trained model: a collage of various pictures of various electronic devices

Partially trained model: A man in a suit and tie holding a cell phone

Client model: a close up of a cell phone with a cell phone on it

Evaluation:

- All three models can identify images of a man however both make mistakes that humans would make if the picture was blurry – fully trained model assumed something being cut because of the torn rope and clients model identified the man being on a skateboard because he was moving with a black arrow underneath him – best possible solution to this would be ‘increase resolution’ of the images by adding context from words mentioned in the images [possibly through OCR]
- Both the fully and partially trained models failed when it came to the tree diagram whereas clients model got part of the answer correct by mentioning trees – need to diversify training data
- In final image, interestingly the clients model matched the answer of the partially trained model, telling us the potentially the client’s training data wasn’t as diversified as it should have been. All three completely failed but the layout of the metal strips matched a gaming style keypad which explains why all three models classed it as electronics.

Overall evaluation:

- Need to diversify training data to include images of a diagram
- We want data with more descriptive captions that aim to describe everything in the image, so we'll use [4] dataset instead
- Need to incorporate OCR to include context based off the text in the images
- Explore any models which are more 'powerful' than GPT2 when it comes to learning from images

References

- [1] Kumar, A. (2022). The Illustrated Image Captioning using transformers. [online] Ankur | NLP Enthusiast. Available at: <https://ankur3107.github.io/blogs/the-illustrated-image-captioning-using-transformers/>
- [2] huggingface.co. (2023). ydshieh/coco_dataset_script · Datasets at Hugging Face. [online] Available at: https://huggingface.co/datasets/ydshieh/coco_dataset_script [Accessed 22 Jun. 2024].
- [3] huggingface.co. (n.d.). nlpconnect/vit-gpt2-image-captioning · Hugging Face. [online] Available at: <https://huggingface.co/nlpconnect/vit-gpt2-image-captioning>.
- [4] google.github.io. (n.d.). DOCCI. [online] Available at: <https://google.github.io/docci/> [Accessed 22 Jun. 2024].