

Image Captioning

Neural Networks and Deep Learning,
Artificial Intelligence, UAB, 2023

Eduard Florin Hogeia

Júlia Garcia Torné

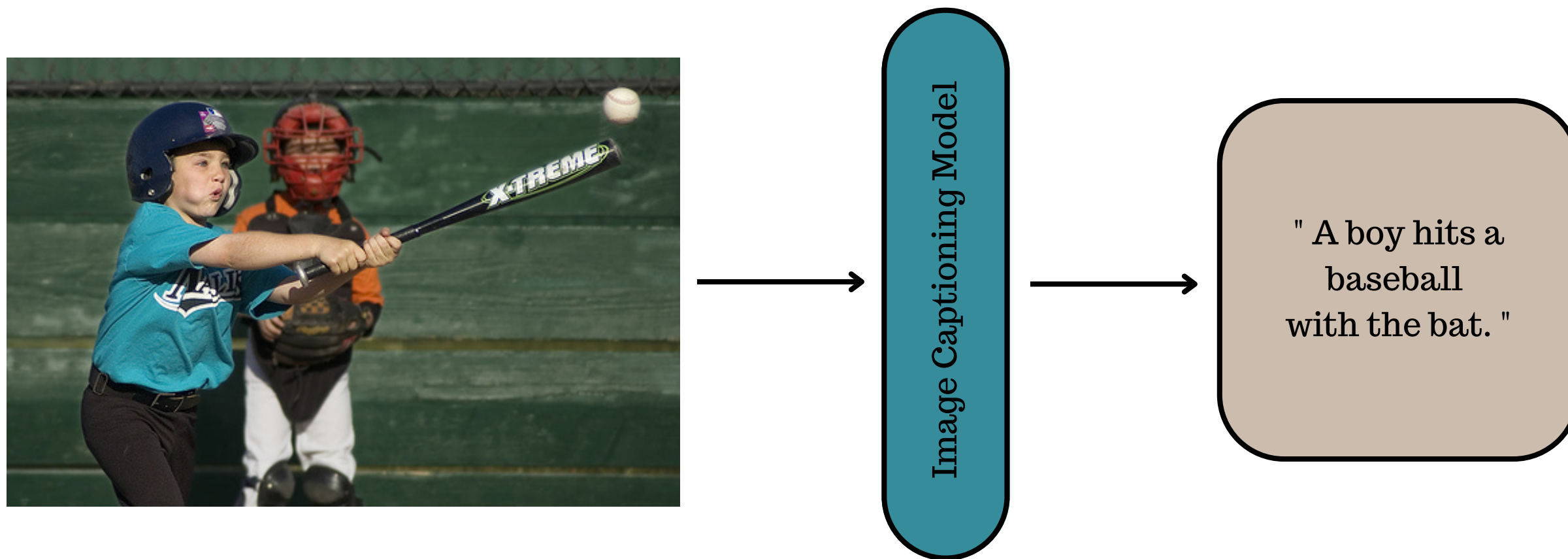
Josep Maria Rocafort Ferrer

CONTENT

- 1 Project Overview
- 2 Dataset Description
- 3 Model Architecture
- 4 Project Flowchart
- 5 Model architecture
- 6 Data pre-processing
- 7 Training process
- 6 Results and performance
- 7 Conclusions

PROJECT OVERVIEW

- Main objective: Develop a system that can generate natural language descriptions for input images.



DATASET DESCRIPTION

Flickr 8k Dataset



Multiple human-annotated
captions for each image

" A man in an orange hat starring
at something. "

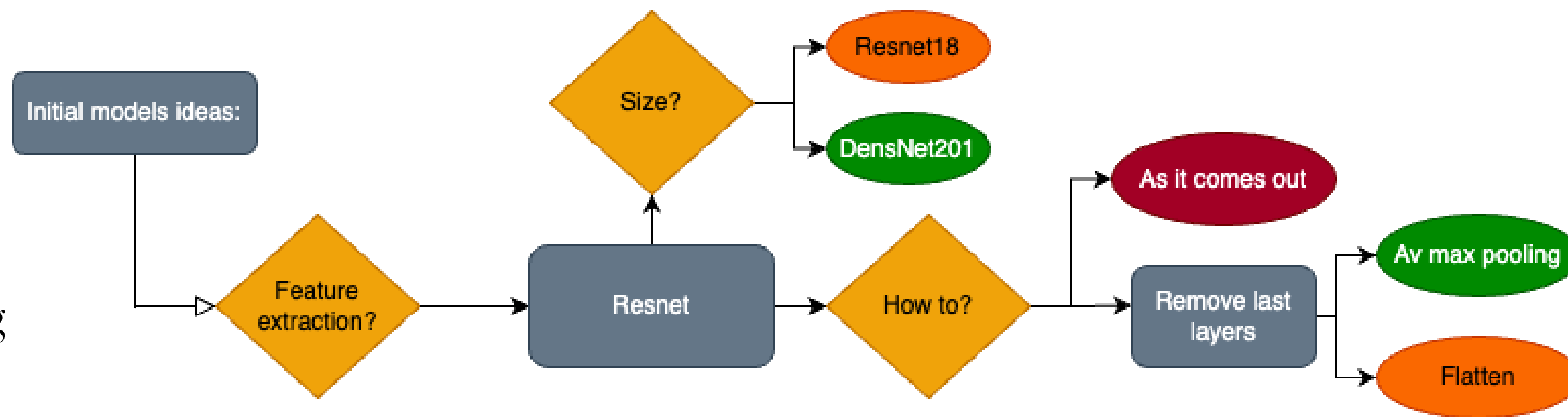
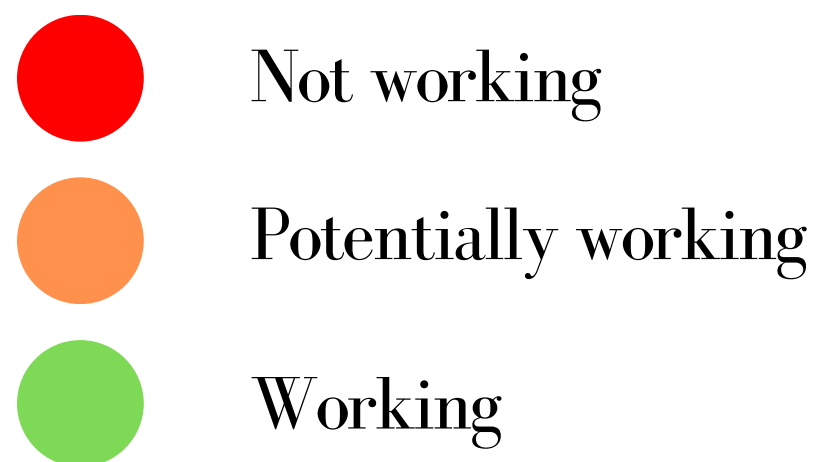
" A person standing on a frozen lake. "

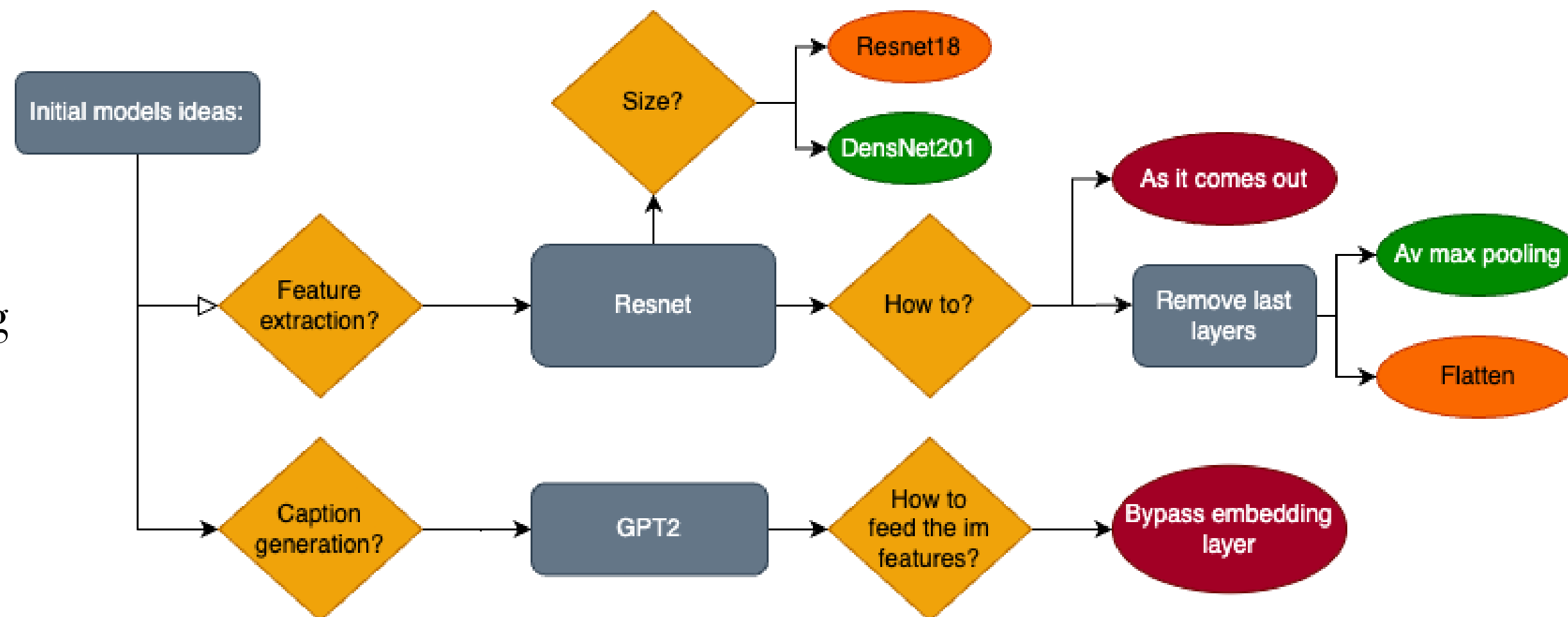
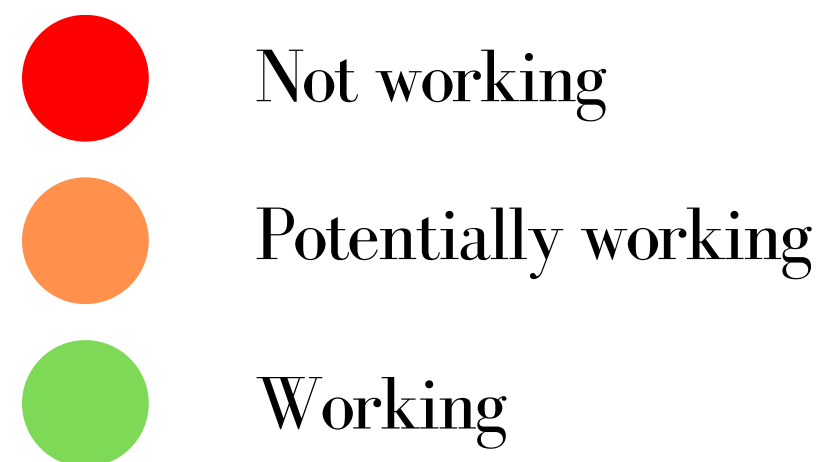
" A young boy runs across the street. "

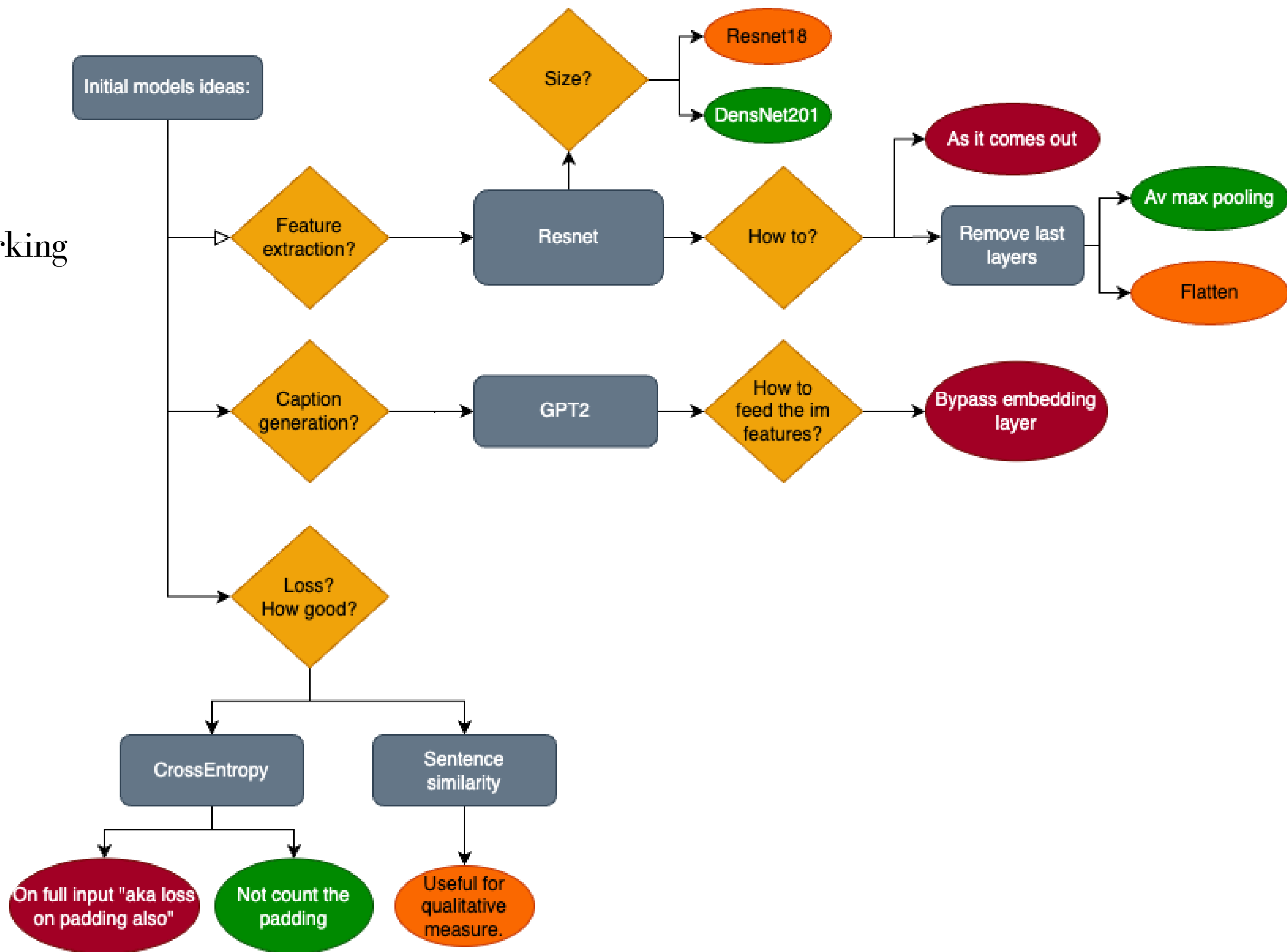
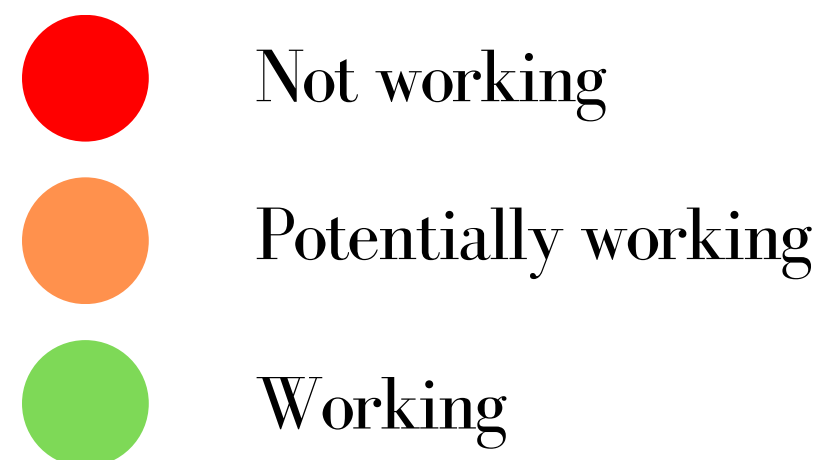
8,000 Images

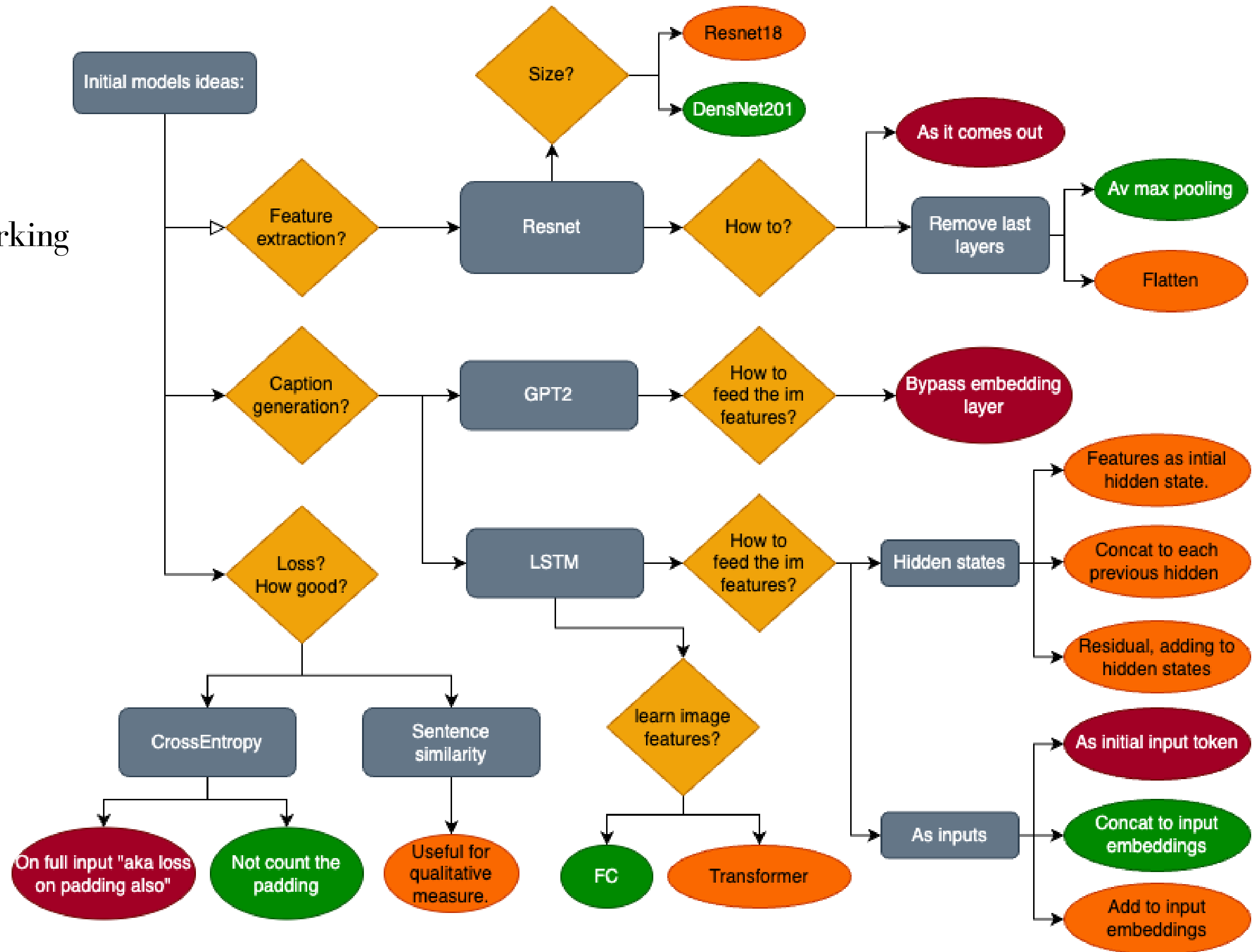
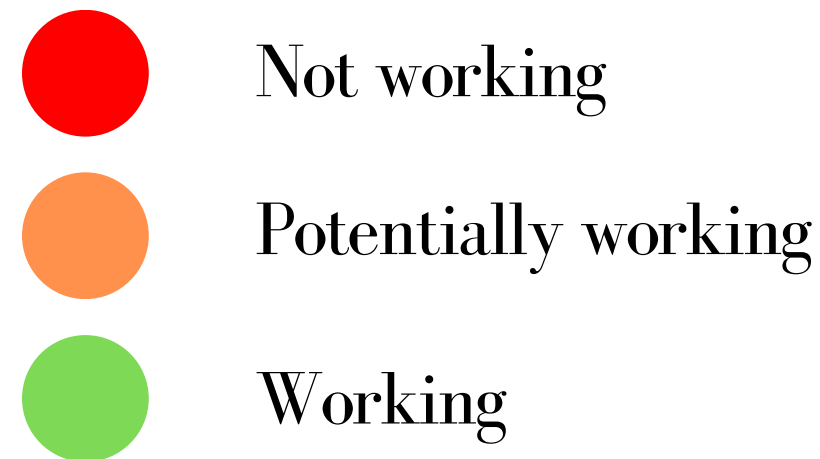
- 80 % for Training
- 20% for Validation

PROJECT FLOWCHART









FINAL MODEL ARCHITECTURE

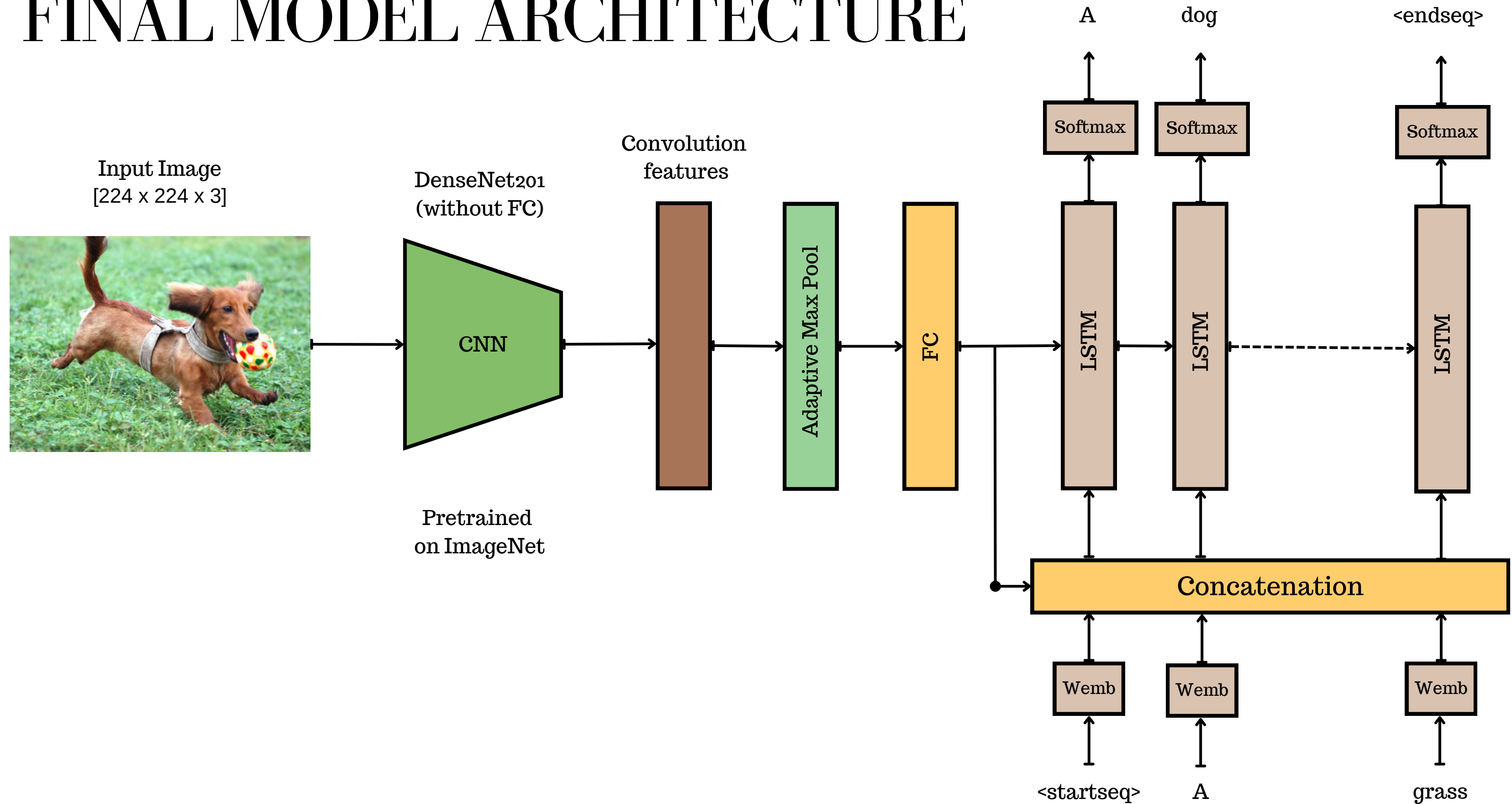
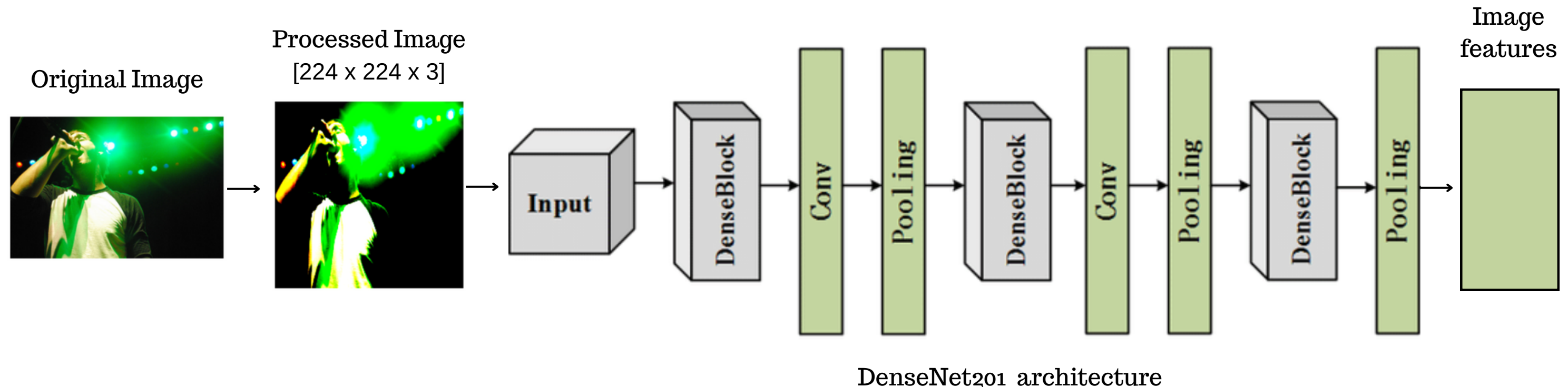


IMAGE FEATURE EXTRACTION

1 Pre-process images →

```
preprocess = transforms.Compose([
    transforms.Resize((224, 224)), # Resize the image: ResNet model - > (224,224,3)
    transforms.ToTensor(), # Img to Python Tensor
    transforms.Normalize(mean=mean, std=std), # image = (image - mean) / std
])
```

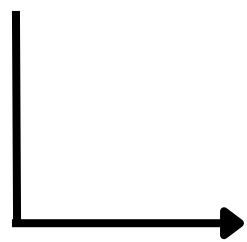
2 Feature Extraction with a pretrained CNN: DenseNet201



CAPTIONS PRE-PROCESSING

- Convert all text to lowercase.
- Remove non-letter characters: Delete punctuation marks, special characters, or any other non-alphabetical characters.
- Replace whitespace sequences.
- Remove single-letter words.
- Add start and end tokens: "startseq" and "endseq" tokens to the beginning and end of each caption.

" A child in a pink dress is climbing up a set of stairs in an entry way. "



" startseq child in pink dress is climbing up set of stairs in an entry way endseq "

- Tokenization process:

The tokenizer will tokenize the caption by splitting it into words and assign a unique index to each word based on its position in the tokenizer's word index.

- The captions are transformed into tensors of length `max_length`.

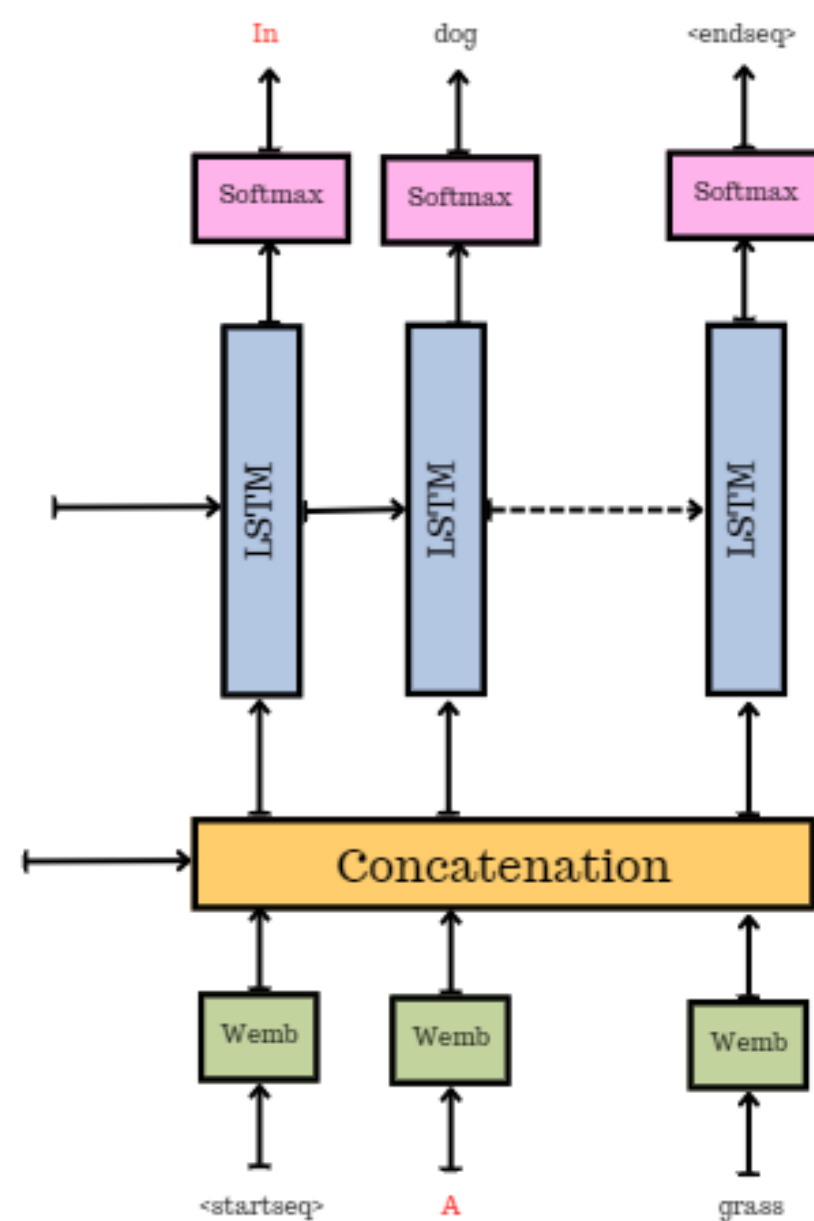


" startseq child in pink dress is climbing up set of stairs in an entry way endseq "

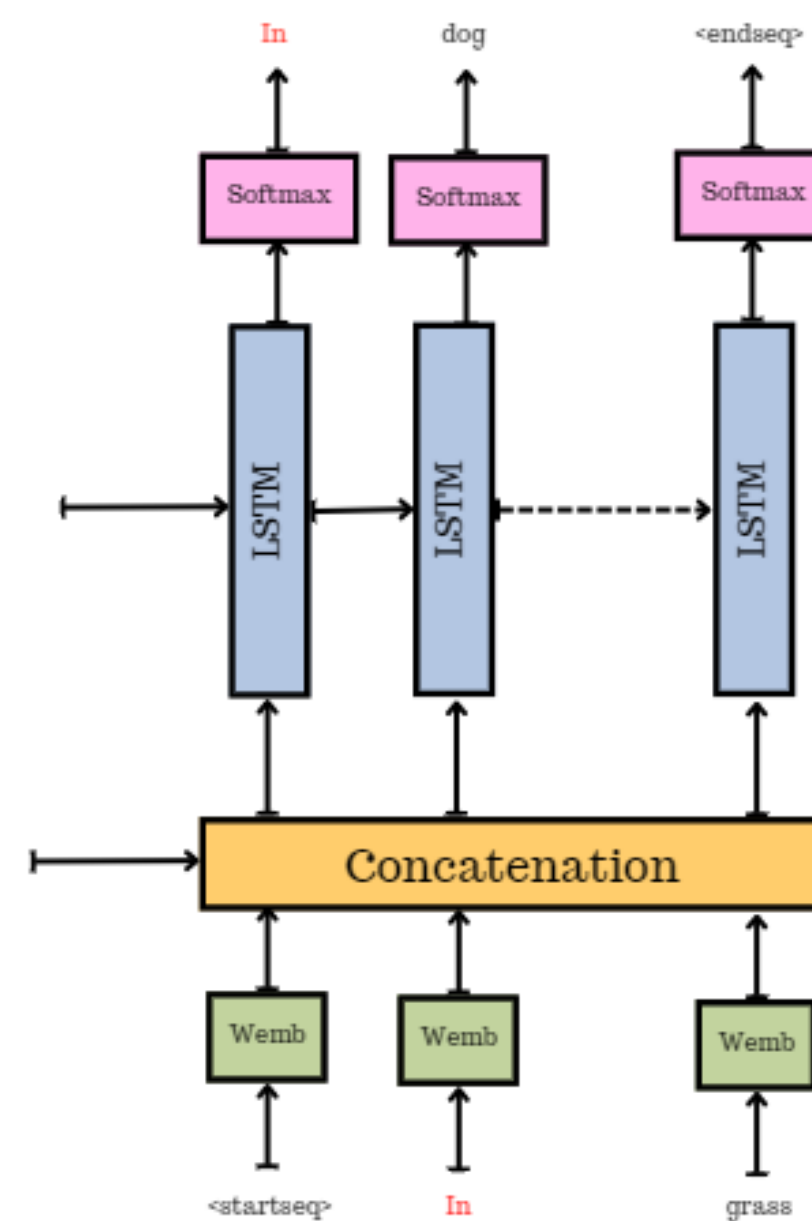
tensor([2, 43, 4, 90, 171, 7, 119, 51, 391, 12, 392, 4, 28, 5123, 668, 3, 0,0, ...])

TRAINING PROCESS: Two approaches

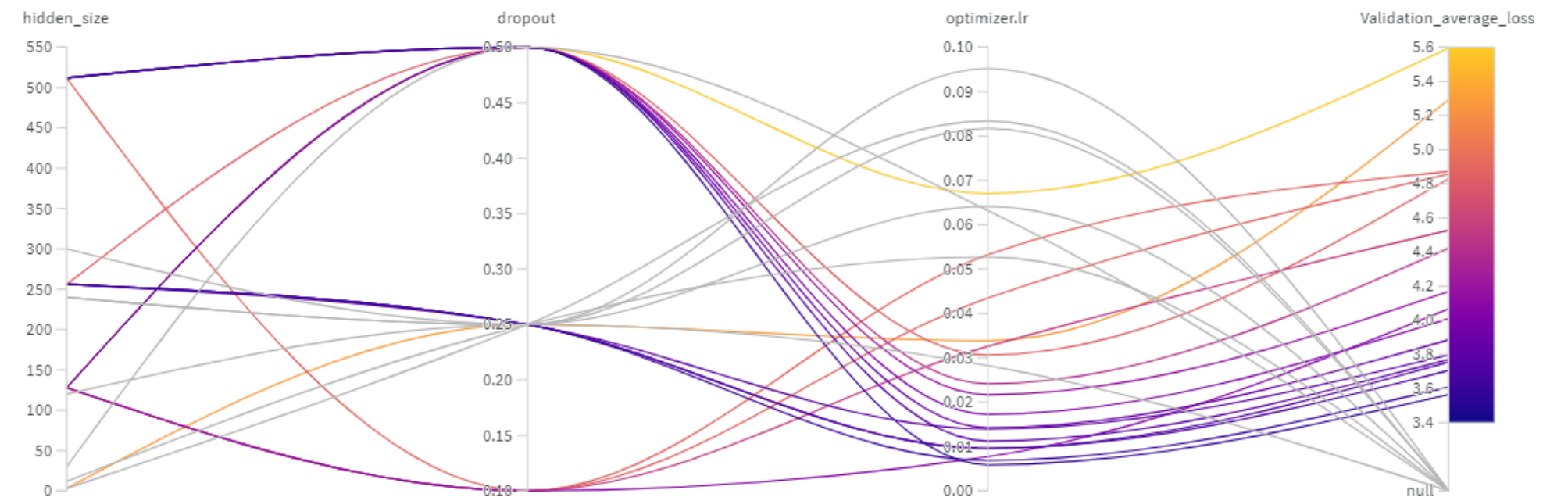
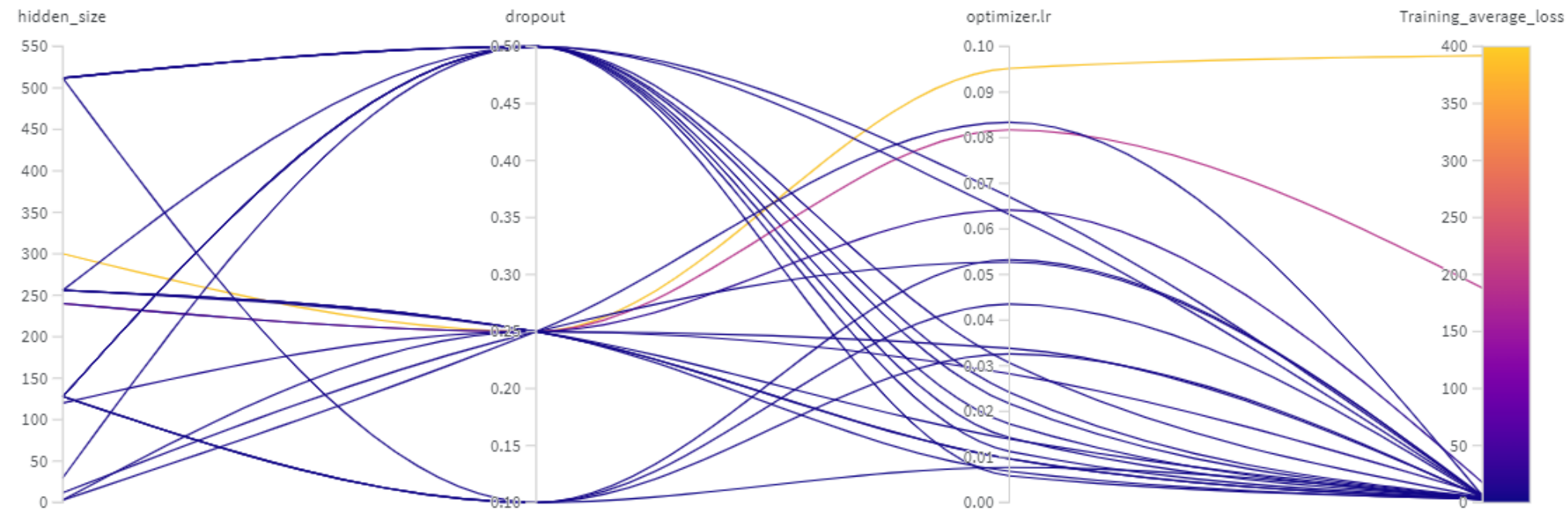
Ground Truth tokens as Input



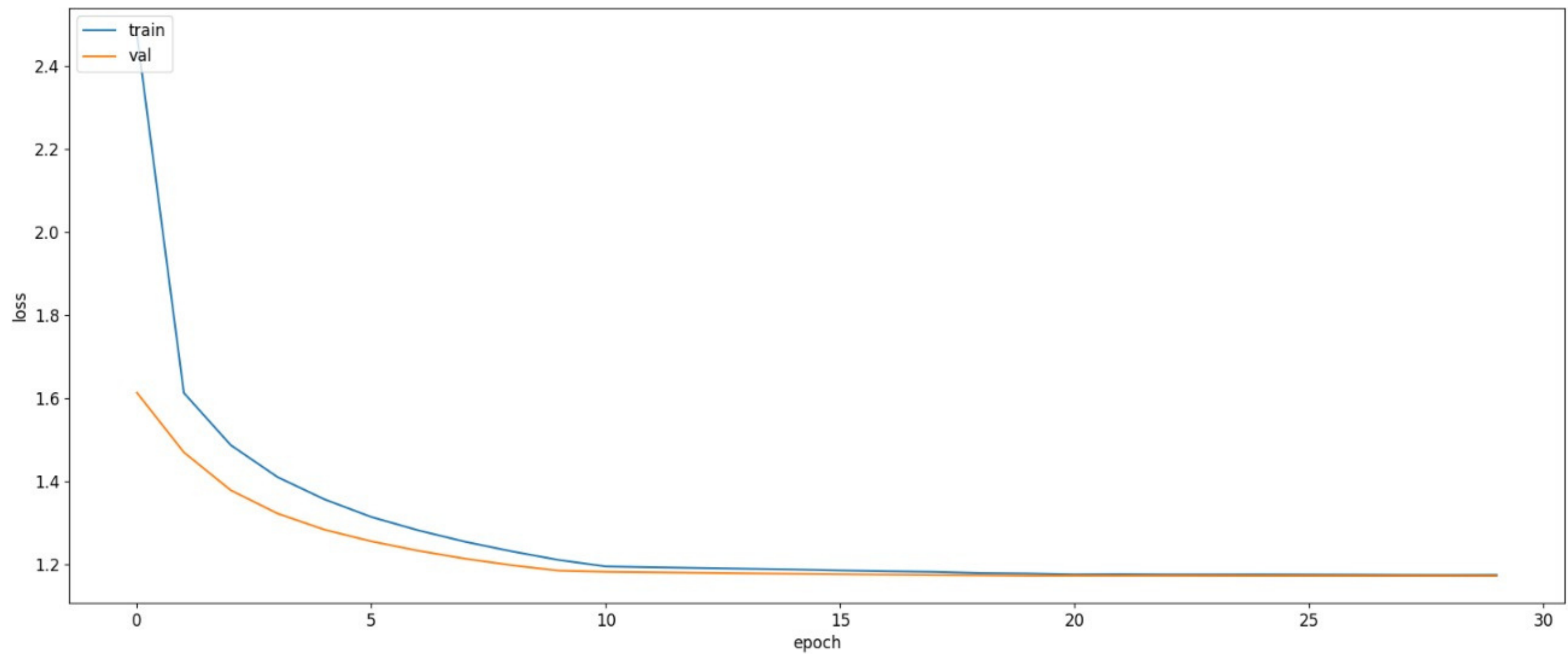
Learning from Model Predictions



WEIGHTS AND BIASES: Optimizing Hyperparameters



Training and Validation Loss

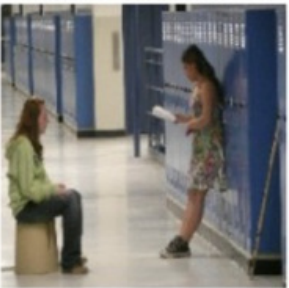


RESULTS AND PERFORMANCE

on to the young
down and holding is
on runs in



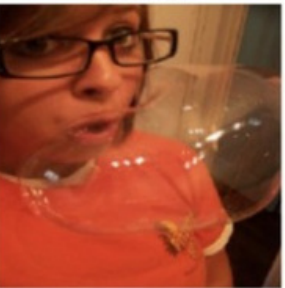
girl the dogs down
dog black down and
ball is on runs in



of the young down
and ball is on runs
in



of the young down
and holding is on
their in



girl the young down
and holding is on
their in



girl the dogs down
and ball is on runs
in



on of the on red
and over air on red
in



with through the on
snow in



on of the on red
and ball the on red
in



girl the dogs down
dog black down and
holding is on their
in



girl the young down
dog black down and
ball the on their in



girl the young down
and ball is on runs
in



boy with and
through playing on
snow in



of the young down
and ball is on runs
in



white running woman
while the on child
in



- BLEU score

RESULTS: Testing on a different dataset

- Out-of-context dataset
- Poorly results

startseq on of and over three on beach in



startseq of the young down and ball is on runs in



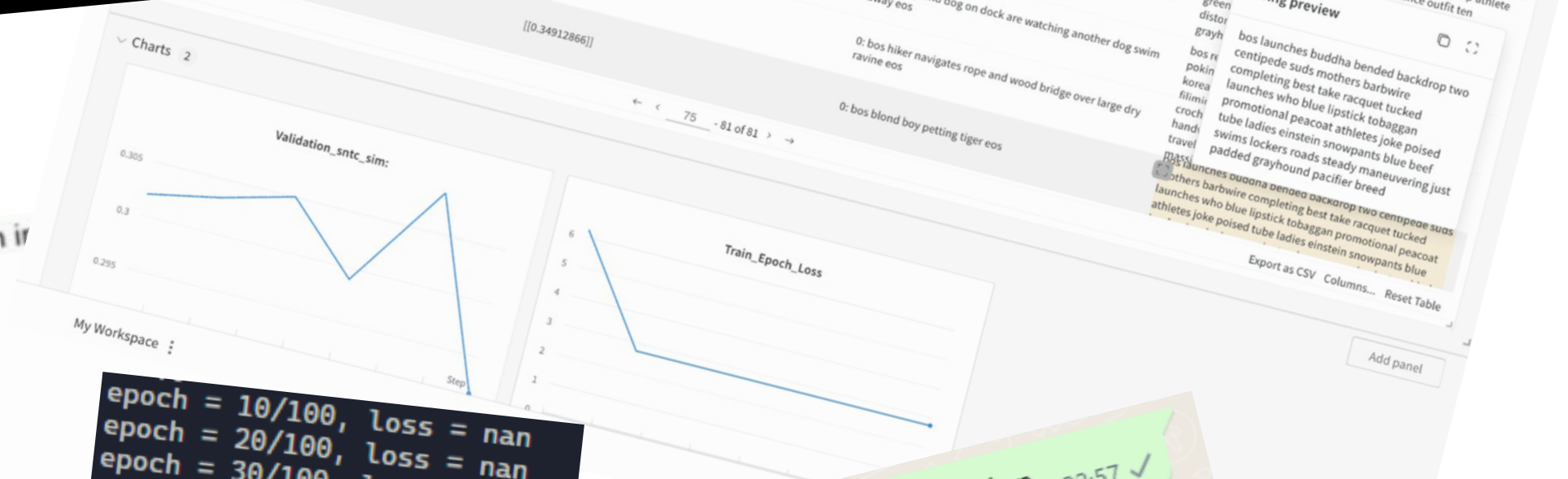
CHALLENGES FACED

[illegible]

It is not learning 🦴🦴🦴🦴

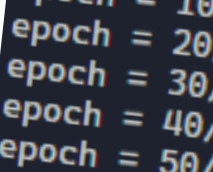
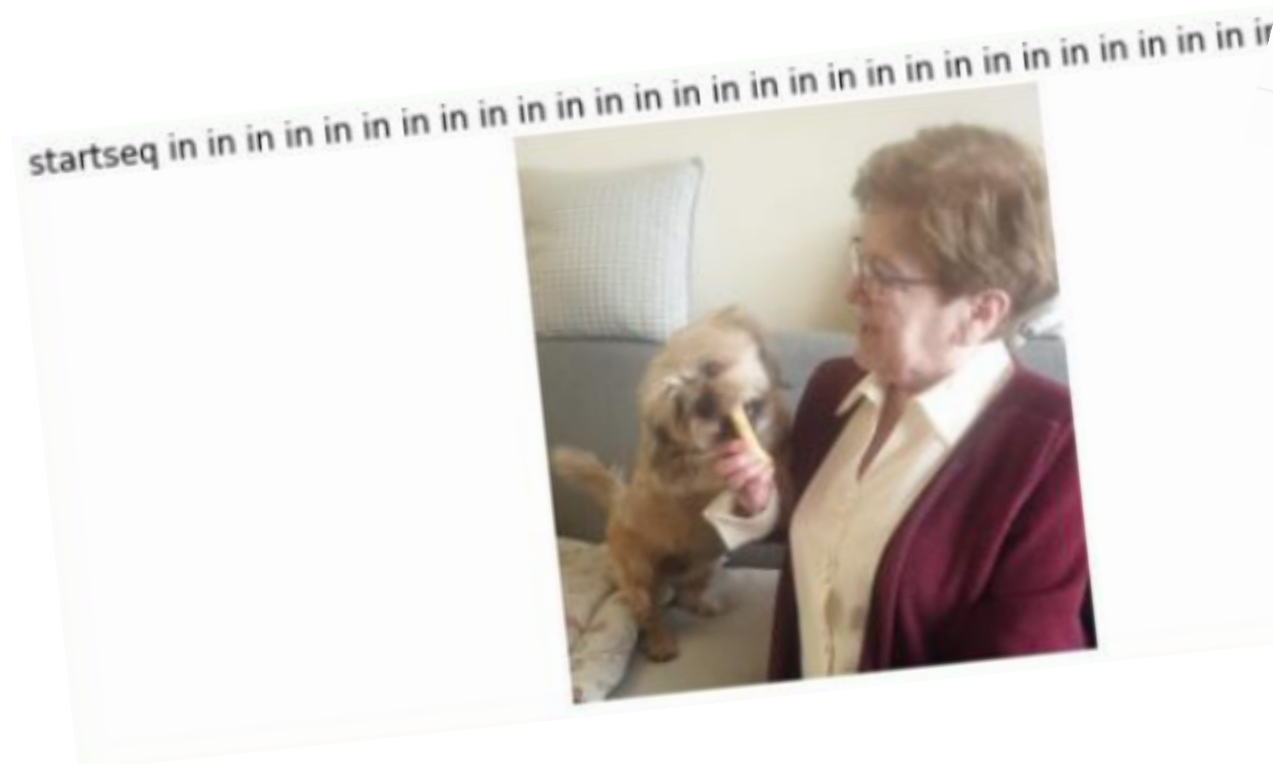
22:49 ✓

I'm close to punching my screen so I would avoid discord for a while :))



But it is a shitty prediction 22:57 ✓

I fucked it up 22:53 ✓



```
epoch = 10/100, loss = nan
epoch = 20/100, loss = nan
epoch = 30/100, loss = nan
epoch = 40/100, loss = nan
epoch = 50/100, loss = nan
epoch = 60/100, loss = nan
epoch = 70/100, loss = nan
epoch = 80/100, loss = nan
epoch = 90/100, loss = nan
epoch = 100/100, loss = nan
Final loss = nan
```

CONCLUSIONS

- Successfully developed a caption generation system.
- The system's performance did not meet our highest expectations.
- Acquired knowledge.
- Future improvements.