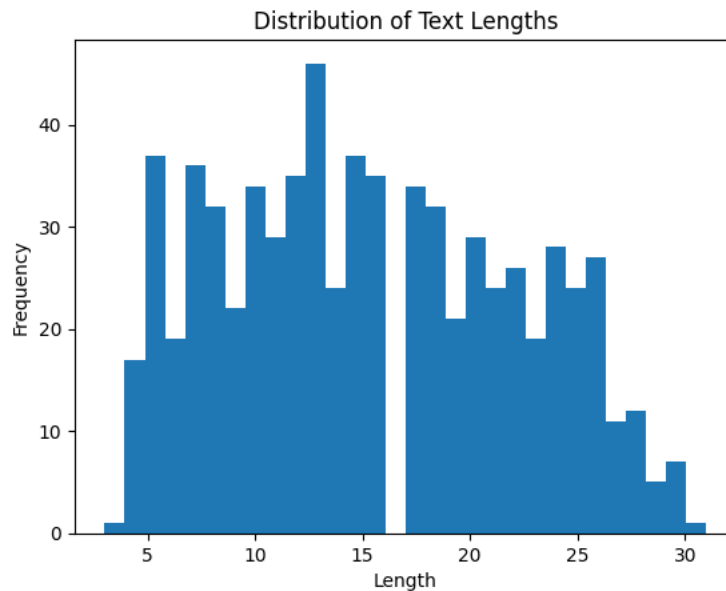


To address the overfitting challenge caused by this emotional imbalance, I chose to approach the problem by building a tailored SeqGAN system for emotion-specific text generation.

Tex Length Distribution in the Training Data

This histogram shows the distribution of text lengths used during training. Most samples fall between 5 and 25 tokens. Based on this distribution, a maximum sequence length of **26 tokens** was selected to capture the majority of expressions while maintaining computational efficiency.



Analyzing the Learning Process of a SeqGAN Model for Emotional Text Generation

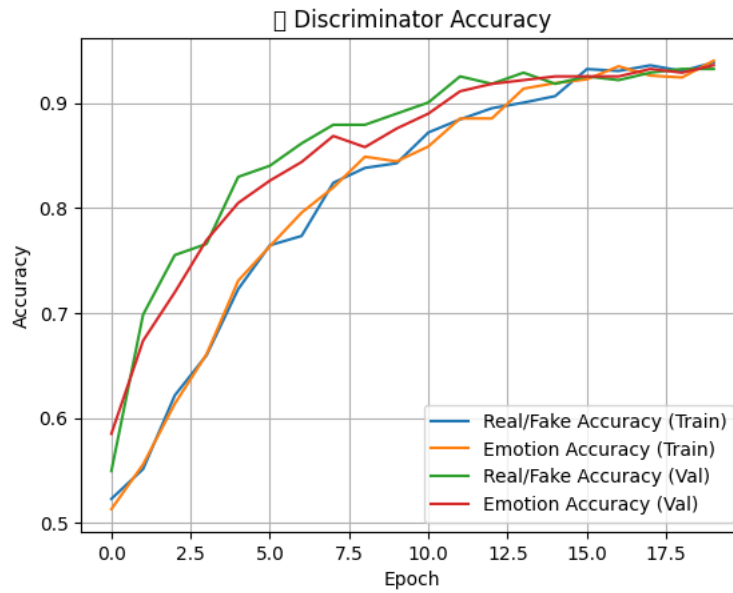
Generator Pretraining – Loss Reduction Over Epochs

```
code cell output actions Generator...
Epoch 1/10
22/22 ----- 5s 25ms/step - loss: 7.5226
Epoch 2/10
22/22 ----- 0s 20ms/step - loss: 6.0189
Epoch 3/10
22/22 ----- 0s 20ms/step - loss: 5.3159
Epoch 4/10
22/22 ----- 0s 20ms/step - loss: 4.8760
Epoch 5/10
22/22 ----- 0s 20ms/step - loss: 4.6094
Epoch 6/10
22/22 ----- 0s 20ms/step - loss: 4.5424
Epoch 7/10
22/22 ----- 0s 20ms/step - loss: 4.2858
Epoch 8/10
22/22 ----- 0s 20ms/step - loss: 4.1814
Epoch 9/10
22/22 ----- 0s 20ms/step - loss: 4.0327
Epoch 10/10
22/22 ----- 0s 20ms/step - loss: 3.9978
Restoring model weights from the end of the best epoch: 10.
Generated sentence: <start> there pool. of gross am wooshing roommate youve corrupt had once but commuting
```

The generator was pretrained using MLE for 10 epochs, with the loss decreasing from **7.52** to **3.99**. This stage sets a solid starting point before switching to reinforcement learning. The sample generated at the end is still unrefined, as expected at this phase.

Discriminator Pretraining – Accuracy Over Epochs

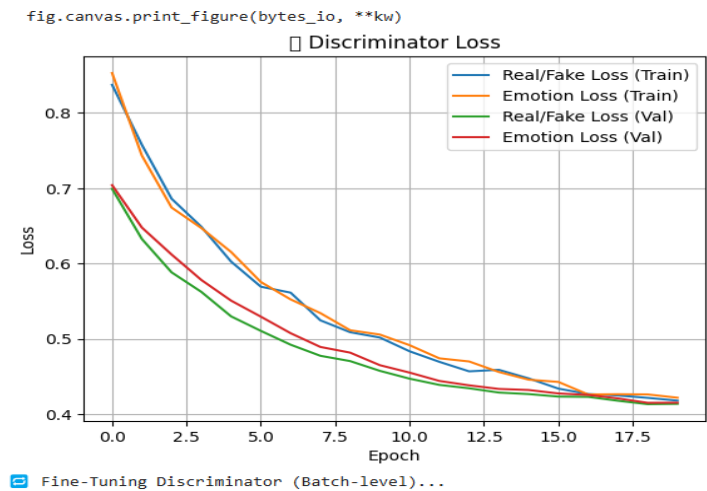
```
/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning  
fig.canvas.print_figure(bytes_io, **kw)
```



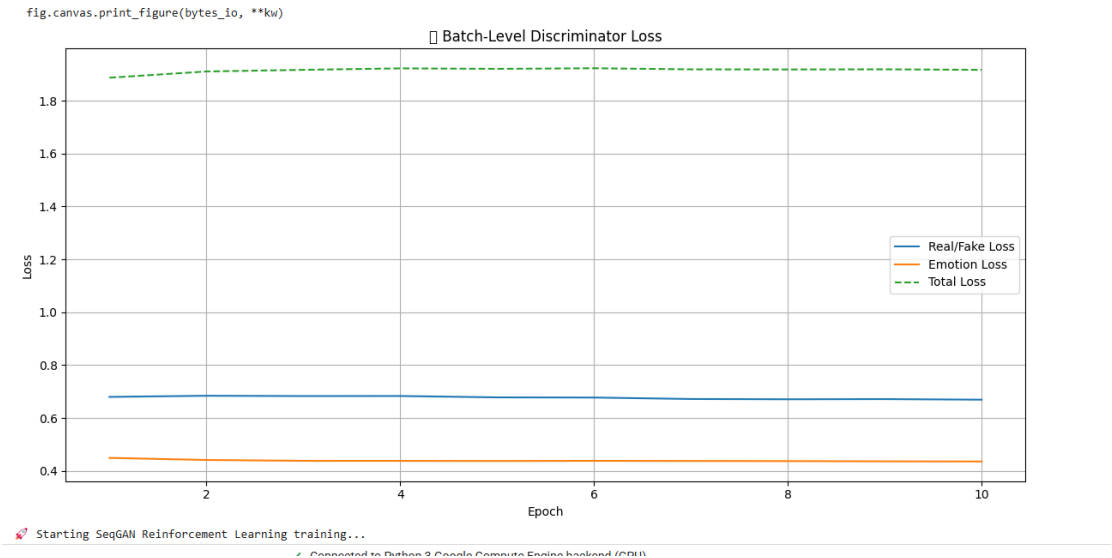
During pretraining, the discriminator learned to classify both **Real/Fake** and **Emotion** labels. Accuracy improved steadily across 20 epochs, reaching over **92%** on both training and validation sets for both tasks. This dual-head setup supports richer reward signals during adversarial training.

Discriminator Loss – Training vs Validation

The chart illustrates a consistent decrease in loss values for both **Real/Fake** and **Emotion** predictions over time. The convergence of training and validation lines suggests stable learning without signs of overfitting.



Fine-Tuning Phase: Discriminator Adaptation During GAN Training



Losses Log:

epoch	d_real_fake_loss	d_emotion_loss	real_fake_acc	emotion_acc	total_loss
1	0.67970693	0.44896722	0.90735066	0.27565083	1.8868984
2	0.68399763	0.44116122	0.91716737	0.19656652	1.9105315
3	0.68285525	0.43792963	0.91741204	0.16547406	1.9166743
4	0.6830655	0.43776506	0.9166286	0.14892645	1.9223051
5	0.6779311	0.43735477	0.91688263	0.14050351	1.9204444
6	0.6772174	0.43789282	0.9159664	0.13336445	1.9226999
7	0.6717166	0.43742612	0.91744965	0.12778524	1.9185381
8	0.6707296	0.43681073	0.91810244	0.12379042	1.9181085
9	0.6713378	0.4359572	0.91998315	0.1205517	1.9186668
10	0.66918284	0.43537518	0.9208325	0.11870367	1.9166579

<ipython-input-127-9a80086d4783>:423: UserWarning: Glyph 128201 (\N{CHART WITH
plt.tight_layout()
/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWar

Once the adversarial training phase begins, the discriminator continues learning in a more challenging setting—facing increasingly realistic fake samples. This graph presents the batch-level loss progression during fine-tuning, showing the balance between real/fake and emotional classification tasks.

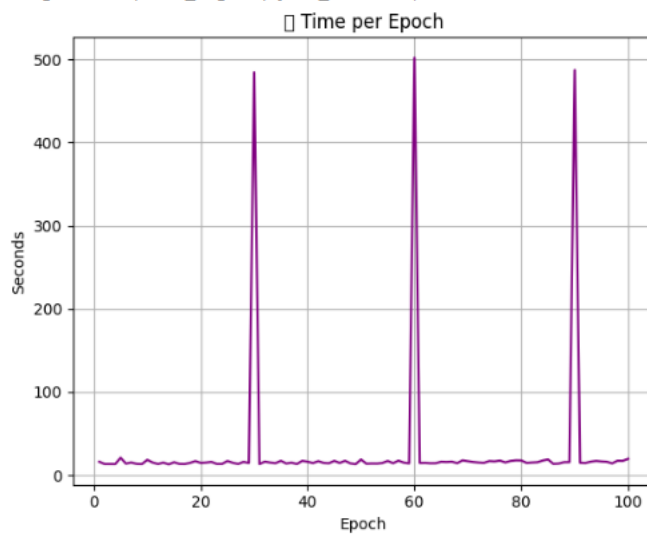
Loss values stabilize over epochs, with the real/fake loss slightly decreasing and the emotion loss remaining consistent. The table confirms strong classification performance (real/fake accuracy > 0.91), while emotional classification remains more difficult, with lower accuracy (~0.12).

This indicates that while the discriminator is highly capable of detecting authenticity, classifying subtle emotional cues remains more complex.

Final Training Insights: Stability, Diversity & Quality

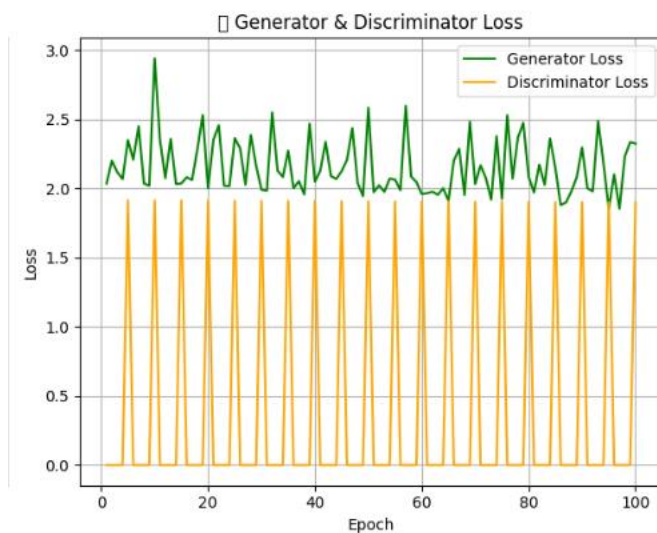
Time per Epoch

```
fig.canvas.print_figure(bytes_io, **kw)
```



[<ipython-input-127-9a80086d4783>:1063: UserWarning: Glvoh 128680](#) This plot shows the training time per epoch. The noticeable spikes suggest periodic operations such as evaluation, saving checkpoints, or heavy batch processing.

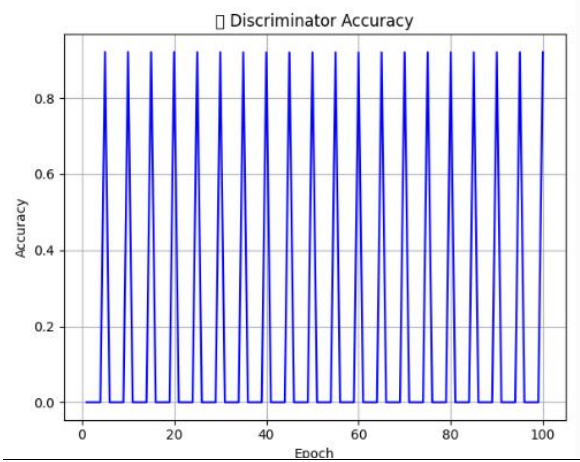
Generator & Discriminator Loss



Generator loss remains relatively stable, while the discriminator shows sharp periodic resets — likely due to adversarial training dynamics.

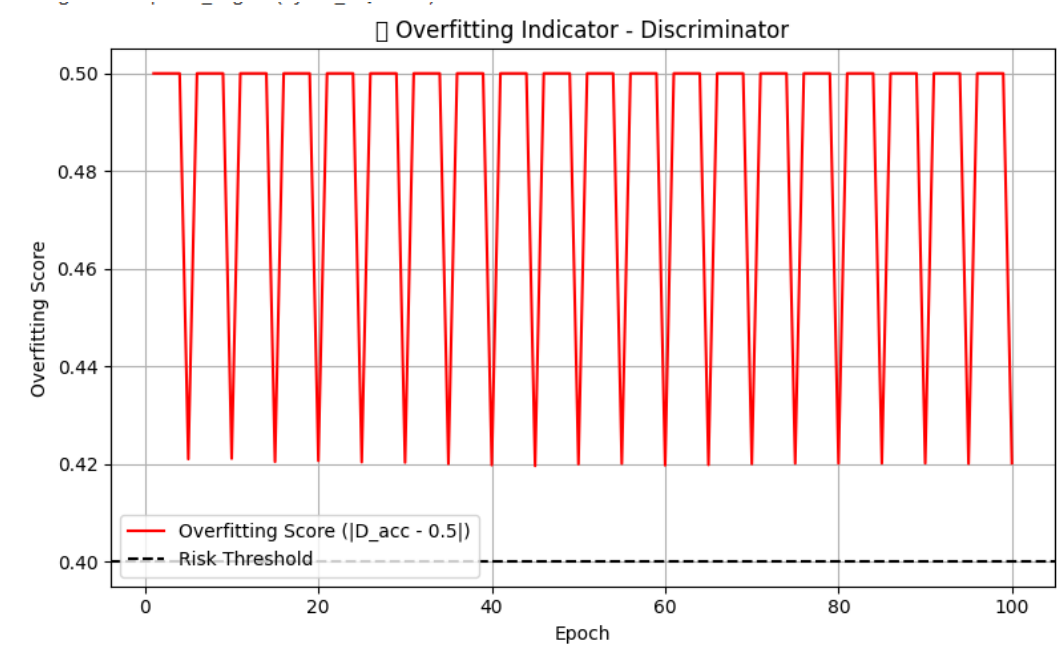
Discriminator Accuracy

font(s) DejaVu Sans.



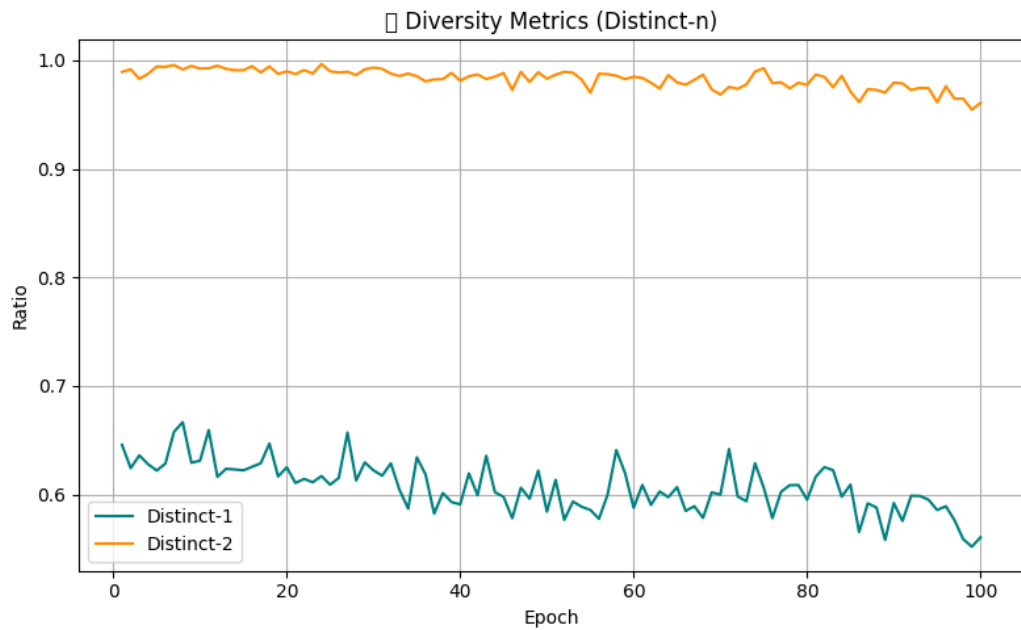
Accuracy spikes show the discriminator quickly overfits between resets. This instability is common in GAN setups and needs balancing strategies.

Overfitting Indicator – Discriminator



The overfitting score remains high across epochs, hovering near the threshold. This signals consistent overfitting in the discriminator component.

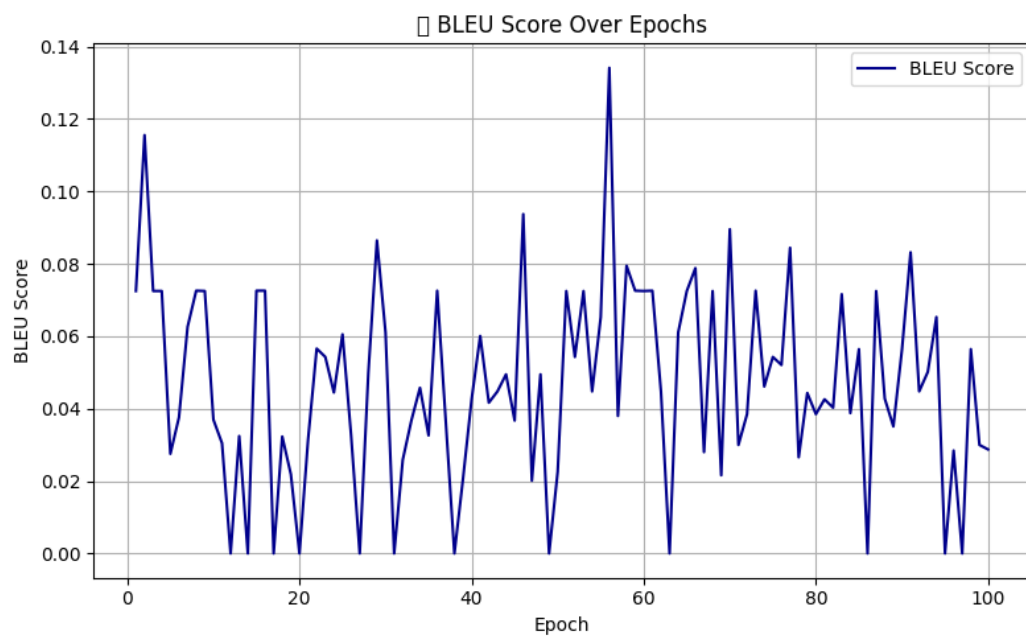
Diversity Metrics (Distinct-n)



Distinct-1 and Distinct-2 scores gradually decline, indicating a drop in text variability. This trend could point to repetitive outputs or mode collapse.

BLEU Score Over Epochs

```
fig.canvas.print_figure(bytes_io, **kw)
```



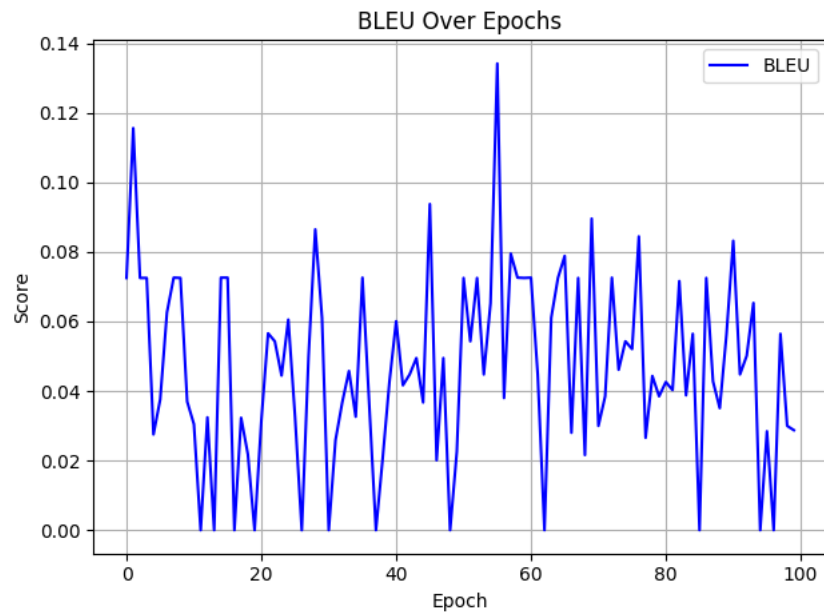
Training metrics exported to 'training_state_rev'

BLEU scores fluctuate across epochs without clear upward trend, suggesting that the generated texts vary in structure but lack consistent linguistic alignment with the training data.

Reward Components Overview

This section presents the different metrics used to compose the final reinforcement reward function.

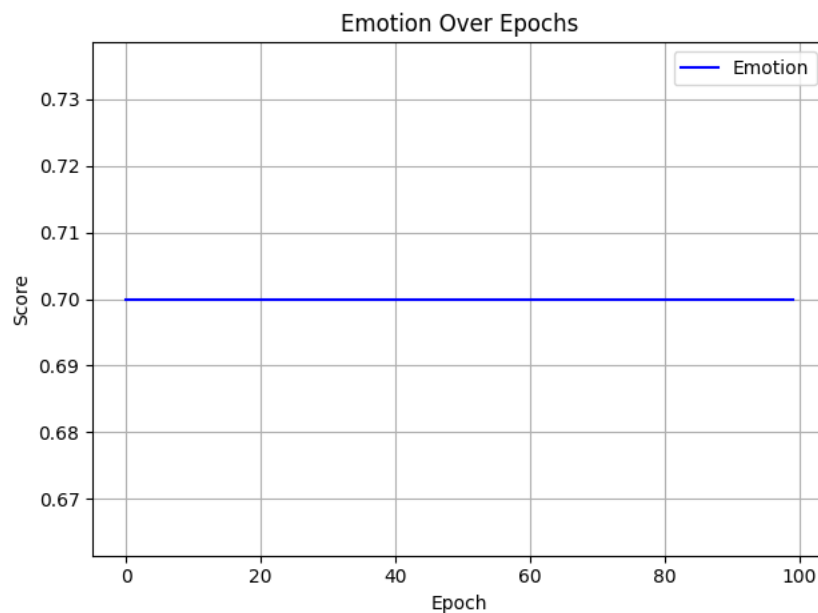
1. BLEU Score



Measures the n-gram overlap between generated and reference texts.

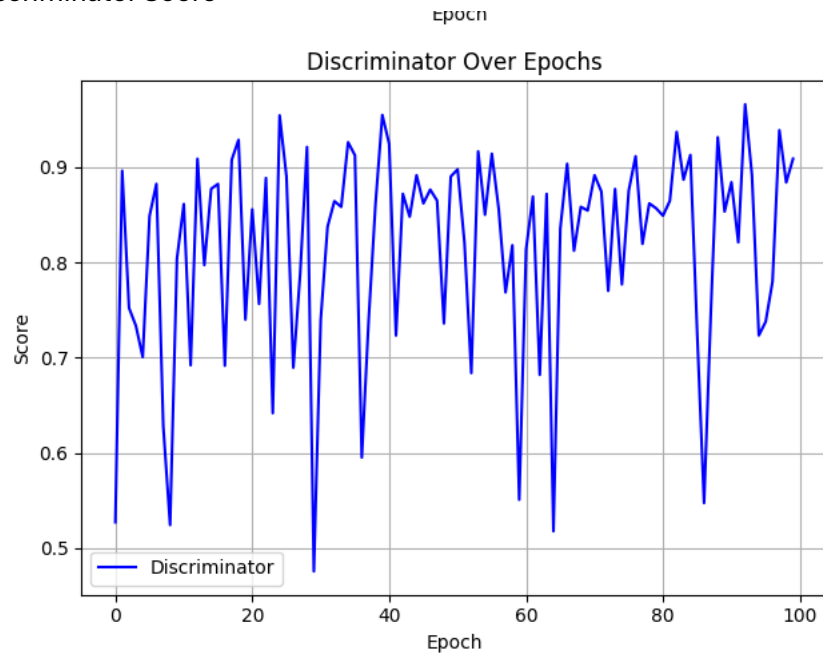
Useful for evaluating grammatical correctness and lexical similarity.

2. Emotion Accuracy



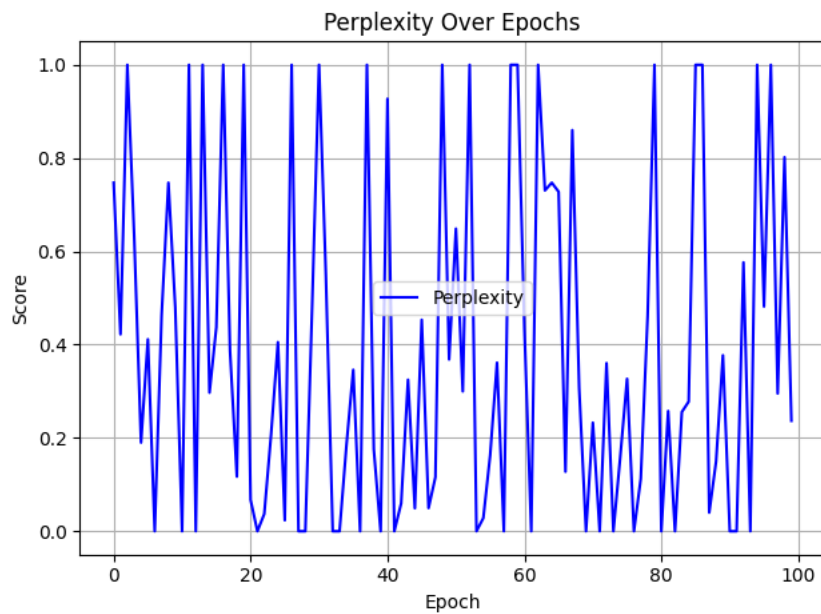
Represents how well the generated sentence aligns with the target emotion. Calculated using a pre-trained emotion classifier. Higher is better.

3. Discriminator Score



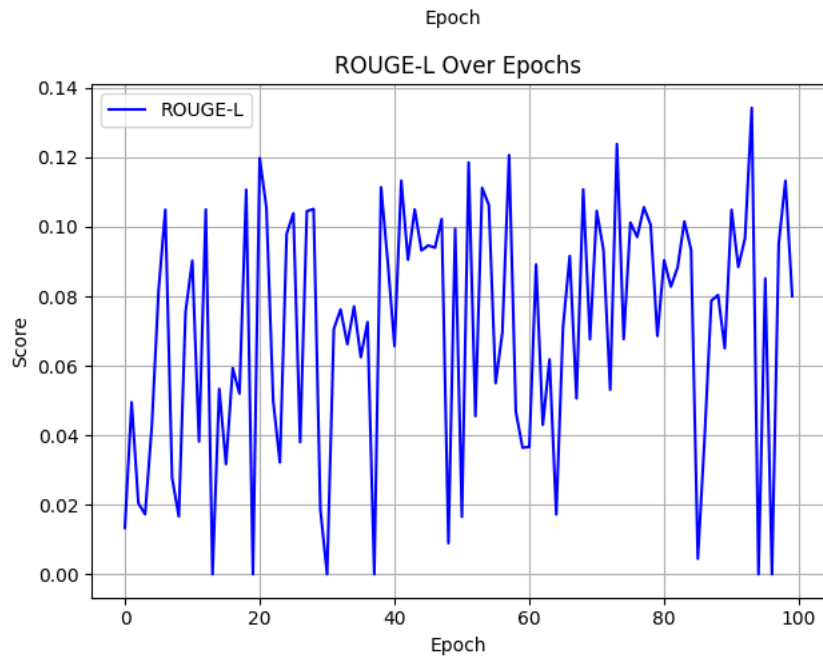
Probability from the discriminator that the sentence is *realistic* (not fake). High scores indicate that the output is coherent and plausible.

4. Perplexity



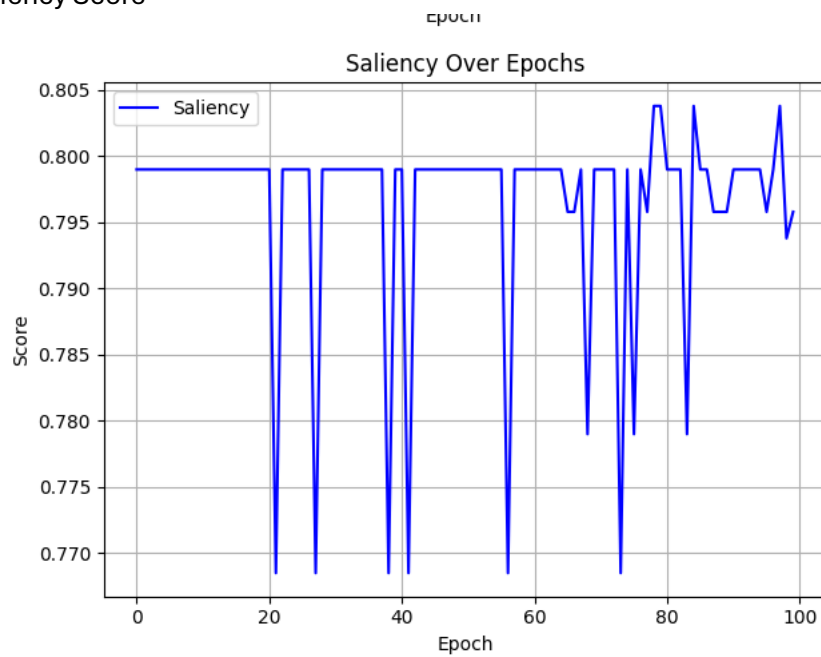
Reflects the fluency and confidence of a language model over the sentence. Lower perplexity suggests better syntactic structure.

5. Repetition Penalty



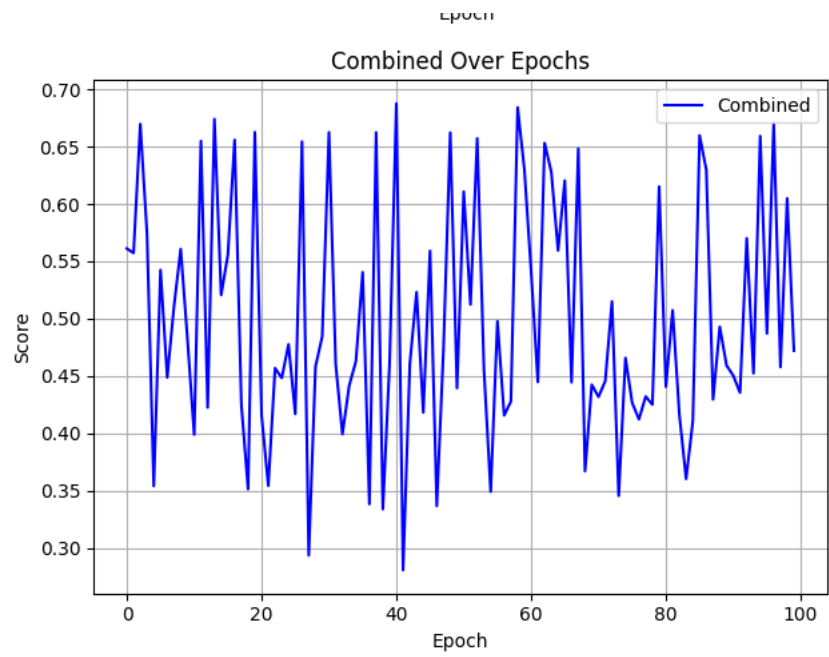
Penalizes overuse of repeated tokens or n-grams Encourages lexical variety and prevents looping.

6. Saliency Score



Measures how semantically important the generated words are Ensures the sentence contains informative and meaningful content.

Combined Reward



Weighted sum of all the above metrics.

Final reward signal that guides the generator during RL training.