
Generating coherent comic with rich story using ChatGPT and Stable Diffusion

Ze Jin

University of Toronto
ze.jin@mail.utoronto.ca

Zorina Song

University of Toronto
zo.song@mail.utoronto.ca

Abstract

Past work demonstrated that using neural networks, we can extend unfinished music pieces while maintaining the music style of the musician. With recent advancements in large language models and diffusion models, we are now capable of generating comics with an interesting storyline while maintaining the art style of the artist. In this paper, we used ChatGPT to generate storylines and dialogue and then generated the comic using stable diffusion. We introduced a novel way to evaluate AI-generated stories, and we achieved SOTA performance on character fidelity and art style by fine-tuning stable diffusion using LoRA, ControlNet, etc.

1 Introduction

Previously, Huawei created a model that extended Franz Schubert’s Symphony No. 8, a piece famously left unfinished.[Huawei, 2019] This is a good idea that can be applied to not just music, but also comics. There are lots of incentive to be able to extend a comic, it can be done when the author left the comic unfinished, it can be done in case people doesn’t like the ending of a particular comic, etc. In this paper, we picked the popular Japanese Manga One Piece as the target to extend and evaluate performance on, but the method we used can be applied to any comic.

2 Related Work

2.1 GPT

OpenAI recently released its latest GPT-4, it outperformed the previous benchmark GPT-3.5 in a variety of exams according to OpenAI.[OpenAI, 2023] But there is no existing work testing GPT-4 ability to generate stories, so we will be comparing GPT-3.5 against GPT-4.

2.2 Diffusion Models

With story generation handled by GPT, the next step is to pick a model that draws the comic. Midjourney[midjourney] is the current SOTA choice for generating comics, the results are very high fidelity[Katz, 2022], however, Midjourney is closed-sourced and accessible only via cloud services. Therefore, we will be working on Stable Diffusion[CompVis], where we can fine-tune the model.

2.3 Visual ChatGPT and HuggingGPT

The plan of this paper is to combine ChatGPT and Stable Diffusion, where we used the output of ChatGPT and feed it into Stable Diffusion. Visual ChatGPT[Wu et al., 2023] and HuggingGPT[Shen et al., 2023] already did that, but because they are using stable diffusion base models, they can’t generate Characters from our target Manga One Piece.

2.4 Dreambooth, Textual Inversion, LoRA - Low-Rank Adaptation

If stable diffusion base model is not capable of generating characters from One Piece, then we need fine-tuning. LoRA[Hu et al., 2021] is our choice for fine-tuning stable diffusion because LoRA offers a good trade-off between file size and training power. Dreambooth[Ruiz et al., 2023] is powerful but results in large model files (2-7 GBs). Textual inversions[Gal et al., 2022] are tiny (about 100 KBs), but you can't do as much. The LoRA we will be using in this paper is Wano Saga made by [Lykon].

3 Methods

3.1 Generating comic story with ChatGPT

We used ChatGPT to generate one page of the comic at a time, where each page of the comic contains 6 panels, and each panel contains a scene description and dialogue between one or two characters. The prompt we used is: "Write a comic book page with six comic book panels with descriptions and dialogs. Characters can only be Nami, Zoro, Monkey D. Luffy, and only one or two characters per panel." To reduce the amount of work, we told ChatGPT to limit the story to contain only three characters, Luffy, Zoro, and Nami. Then, the prompt "continue for another page" is used to generate another page.

3.1.1 Generating comic story with ChatGPT with extra biography prompt

Because we are not sure how much ChatGPT knows about One Piece, we wanted to provide knowledge of One Piece to ChatGPT. Therefore, we also tried entering the biography of Luffy, Zoro, and Nami that we got from One Piece Wiki[fandom] as three different prompts to ChatGPT prior to the actual comic generation prompt.

3.2 Calculate Story Score

Since there is no existing metric that evaluates AI-generated stories, we will create one. Introducing story score, the story score of story x against target story t and popular Manga set M :

$$story(x; t, M) = \gamma sim(x, t) + (1 - \gamma) plot(x; M) \quad (1)$$

where $sim(x, t)$ is the similarity score of story x against the plot of the target Manga One Piece, and $plot(x)$ is the plot score of story x against the story of a collection of popular Manga M ,

$$plot(x; M) = \frac{1}{|M|} \sum_{m \in M} sim(x, m) \quad (2)$$

with originality adjusting factor $0 \leq \gamma \leq 1$.

The reason behind the story score is, firstly, we need the newly generated story to be a continuation of the target story, therefore we want it to have some similarity with the original story. Then, we also want the generated story to be good, therefore we compare it against a set of popular Manga M , we calculate the plot score as average similarity score on the set of popular Manga M .

In this paper, we calculated a simplified version of the story score to reduce the amount of computation, where we calculate the story score using the summary of ChatGPT generated story against the summary of the target story, and the summary of popular Manga set M' .

$$story(x; t, M') = \gamma sim(smry(x), smry(t)) + (1 - \gamma) plot(smry(x); M') \quad (3)$$

We summarized the ChatGPT generated story using ChatGPT, because according to this [Chen] and other articles, GPT is currently the best option for text summarization. Then we calculated the similarity score using pre-trained BERT model[Devlin et al., 2019], with sentence-encoder[pyp]. We get the summary's embedding using BERT model, then calculate the cosine similarity of the embeddings, our code is on Github[Song and Jin]. The popular Manga set we compared against is

$M' = \{Naruto, Bleach, HunterXHunter, AttackOnTitan, FullmetalAlchemist\}$

And we used $\gamma = 0.5$.

3.3 Generating Character with Stable diffusion

The scene description generated using ChatGPT contains the character, background, and action. Our hypothesis is that, before we generate the entire scene, we should fine-tune the model so that it can generate the characters first, once our model is capable of generating characters, then we will ask it to draw the entire scene which contains the characters along with the background and actions.

3.3.1 Generating Character - Checkpoints/base models

To use Stable Diffusion, we can either start with a base model like v2.1, or checkpoints like waifu-diffusion[hakurei] and Midnight Mixer Melt[DrBob2142] that already had fine-tuning done on anime images.

3.3.2 Generating Character - Fine-Tuning

Then, we came up with a sequence of steps to fine-tune stable diffusion models. First, apply LoRA made by [Lykon] that was specifically fine-tuned on One Piece characters. Next, adjust the parameters, Sampling Steps from 20 to 30, CFG Scale from 7 to 9, and sampler from Eucla A to DPM++2M Karras. Next, change Clip-Skip[El-Nouby et al., 2019] to 2 because our base model was trained to perform better on Clip-Skip. Next, use waifu-diffusion-v1-4's VAE [Kingma and Welling, 2022] made by [hakurei] that was trained on anime images to produce better color. Then, add extra prompts with keywords "wanostyle", "solo", "stand straight", "long green clothes", "white clear background". Lastly, the character already looks good by this point, but the pose of characters still varies based on the random seed, to have the character in the exact pose we want, use ControlNet[Zhang and Agrawala, 2023] with Canny model[llyasviel]. Where Canny's pre-processor extracts the pose in any given image, and Canny's model makes stable diffusion generate character based on the given pose.

3.4 Calculating SSIM and FID

We used SSIM[Wang et al., 2004] and FID[Heusel et al., 2018] to evaluate the result after fine-tuning stable diffusion against the target. To calculate SSIM, we feed the result and target batch into Pytorch's conv2d with a filter sampled from a Gaussian distribution, then used the output's mean and variance to compute SSIM. To calculate FID, we feed the result and target batch into Keras' InceptionV3, then we used the output's mean and variance to compute FID, our code is on Github[Song and Jin].

4 Experiment

4.1 Generating comic story

We generated comic stories using GPT-3.5 and GPT-4 through ChatGPT, with and without extra character biography prompts, then calculated story scores for each story using the method in 3.2.

4.2 Generating Character

We picked the character Roronoa Zoro from One Piece as the target, we used 5 original images of Roronoa Zoro as the target batch. We posted the target batch we used in the appendix.

4.2.1 Generating Character - Checkpoints/base models

The first step is to evaluate checkpoints/base models. We evaluated Midjourney V4, Midjourney V5, stable diffusion v2.1, and Midnight Mixer Melt, on their ability to generate Roronoa Zoro from One Piece. We used the default parameters, with the bare minimum prompt "Roronoa Zoro, One Piece, Katana, Full body". For each model, we generated a batch of 5 images of Roronoa Zoro using 5 random seeds. And then we used the batch of 5 images generated by each model to calculate SSIM and FID scores against the batch of 5 target images.

4.2.2 Generating Character - Fine-Tuning

Then we picked Midnight Mixer Melt to fine-tune and evaluated the effect of each fine-tune step. In order to fairly evaluate the effect of each fine-tuning step, we fixed 5 random seeds to be reused in every step. We performed the sequence of fine-tuning steps in 3.3.2, and used the batch of 5 outputs to calculate the SSIM and FID against the target batch after each fine-tuning step.

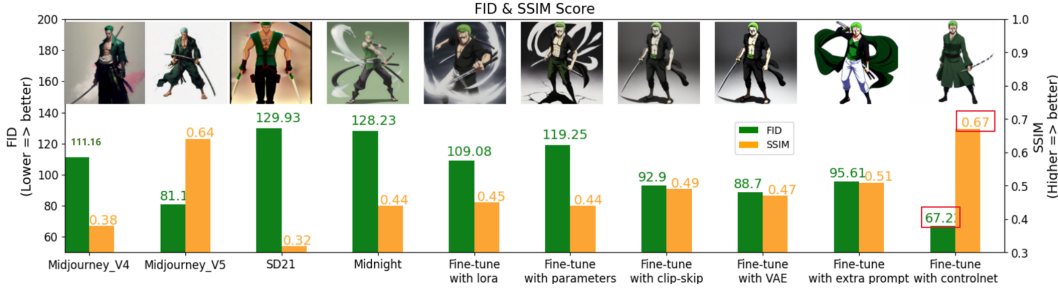
5 Results

5.1 Generating Story - Story Score

Model	Similarity Score	Plot Score	Story Score
GPT-4 with biography prompt	0.57	0.67	0.62
GPT-4 without biography prompt	0.49	0.66	0.57
GPT-3.5 with biography prompt	0.52	0.58	0.55
GPT-3.5 without biography prompt	0.54	0.67	0.6

Giving extra prompts and extra knowledge about One Piece does not necessarily generate a better story both by looking at the story score and reading the story by humans. Also, GPT-3.5 and GPT-4 performed similarly both by looking at the story score and reading the story by humans. But overall, GPT-4 with extra biography prompt performed slightly better. We posted sample stories that we generated in the appendix.

5.2 Generating Character - SSIM and FID Score



Midjourney V5[midjourney] initially performed the best out of all the checkpoints/base models. But in the fine-tuning steps, our result gradually improves for each step both by looking at the score, and in human eyes. Then, after all the fine-tuning steps, our model did outperform the benchmark Midjourney V5. We posted the images generated before and after fine-tuning in the appendix.

5.3 Comic Results

We don't have space to evaluate the result of generating full comics for this paper, but it does look like our hypothesis from 3.3 holds true, we posted sample comics that we generated in the appendix.

6 Conclusion

In conclusion, we found a good way to generate comic scenes with ChatGPT and evaluated the story using the story score we introduced, then we found an effective sequence of steps to fine-tune stable diffusion to generate the character in the target Manga and outperformed Midjourney V5.

6.1 Limitation and Future work

For comic characters, our experiment could be better if we can test more characters and larger random batches. Also, SSIM and FID score is not the perfect way to evaluate the AI-generated comic, we don't have better metrics. For the story, our story score metric needs more testing and proof. Due to model and compute limitations, we calculated story scores using stories summarized by ChatGPT.

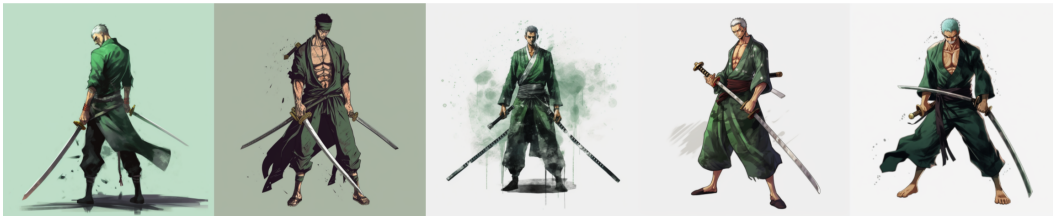
7 Appendix

7.1 Sample Character - Roronoa Zoro from One Piece

7.1.1 target batch



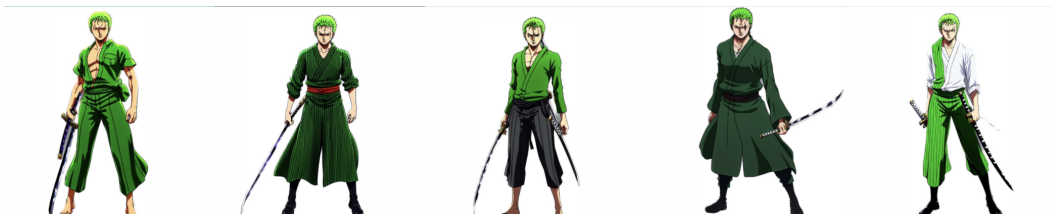
7.1.2 Midjourney V5 benchmark




7.1.3 before fine-tuning



7.1.4 after fine-tuning



7.2 Sample Story Generated using ChatGPT

>  write a comic book page with six comic book panels with descriptions and dialogs. Characters can only be Nami, Zoro, Monkey D. Luffy, and only one or two characters per panel.



Panel 1:

Description: Nami is at the ship's helm, navigating the crew through a storm. The waves are rough, and the sky is filled with lightning.

Dialog: Nami (shouting): "Hang on, everyone! This storm is getting worse!"

Panel 2:

Description: Zoro is securing the sails to prevent damage to the ship. The wind is howling, and rain is pouring down.

Dialog: Zoro (grunting): "I've got the sails! Just focus on steering, Nami!"

Panel 3:

Description: Luffy is standing on the deck, staring out at the raging sea with determination in his eyes.

Dialog: Luffy (yelling): "We can't let a storm like this stop us! We're going to be the Pirate King's crew!"

Panel 4:

Description: Nami and Zoro exchange a glance, acknowledging Luffy's determination.

Dialog: Nami (smirking): "He's right! We've been through worse!"

Zoro (grinning): "Yeah, and we'll get through this too!"

Panel 5:

Description: Luffy starts to inflate his arm, preparing for a powerful attack on a massive wave that threatens to capsize the ship.

Dialog: Luffy (shouting): "Gomu Gomu no..."

Zoro (calling out): "Luffy, what are you doing?!"

Panel 6:

Description: Luffy's attack sends the wave crashing away from the ship, as Nami steers them clear of danger. The storm begins to subside.

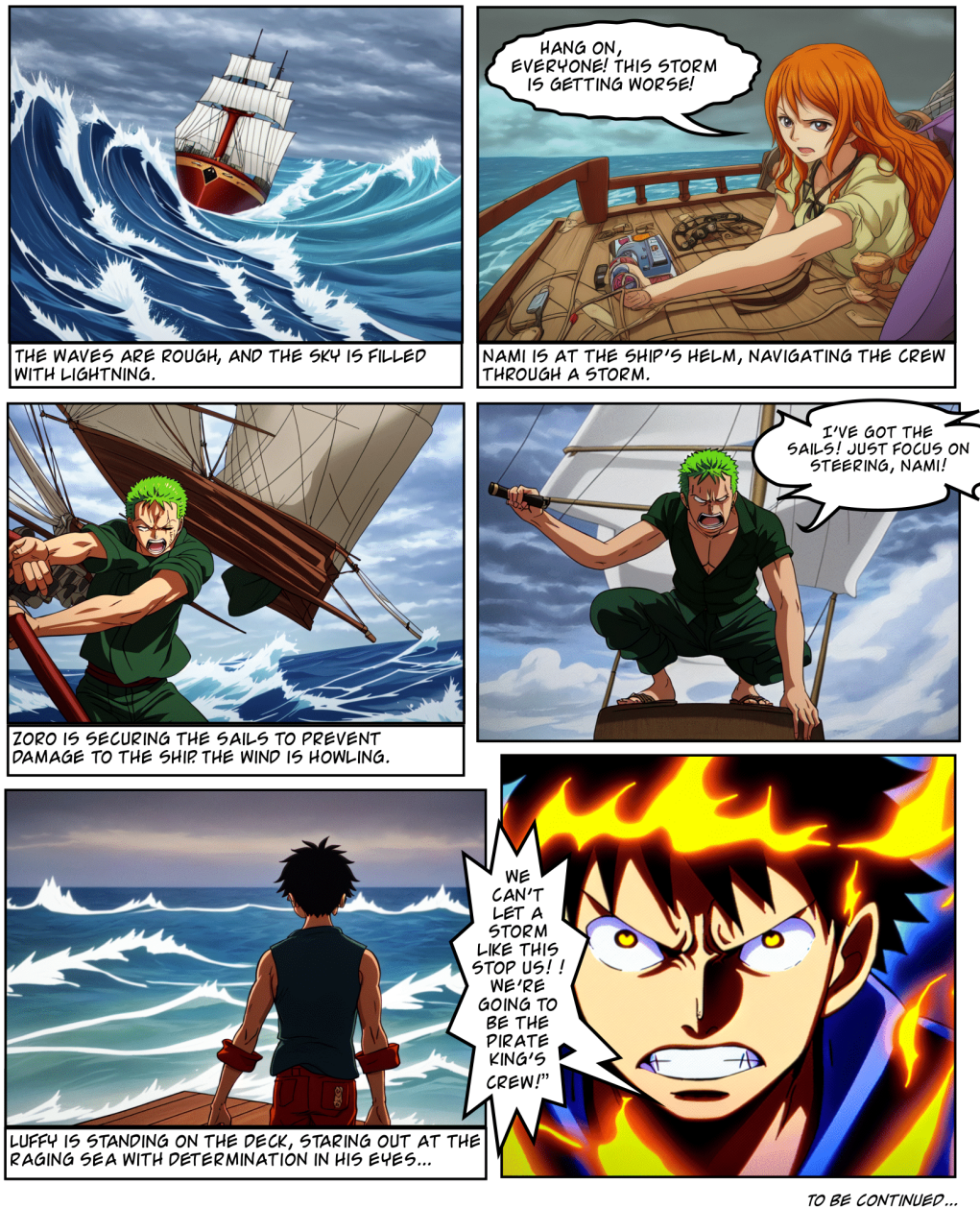
Dialog: Luffy (grinning): "Giant Pistol! See, I told you we'd make it!"

Nami (relieved): "Thanks, Luffy. You saved us!"

Zoro (smirking): "You're crazy, Captain... but it worked!"

 Regenerate response

7.3 Sample Comic Generated using fine-tuned Stable Diffusion



References

Sentence-transformers. URL <https://pypi.org/project/sentence-transformers/>.

M. Chen. How to validate openai gpt model performance with text summarization ... URL <https://towardsdatascience.com/how-to-validate-openai-gpt-model-performance-with-text-summarization-298978fea764>.

CompVis. Compvis/stable-diffusion: A latent text-to-image diffusion model. URL <https://github.com/CompVis/stable-diffusion>.

- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- DrBob2142. Drbob2142/midnight_mixes · hugging face. URL https://huggingface.co/DrBob2142/Midnight_Mixes.
- A. El-Nouby, S. Zhai, G. W. Taylor, and J. M. Susskind. Skip-clip: Self-supervised spatiotemporal representation learning by future clip order ranking, 2019.
- fandom. One piece wiki. URL https://onepiece.fandom.com/wiki/One_Piece_Wiki.
- R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022.
- hakurei. Hakurei/waifu-diffusion-v1-4 · hugging face. URL <https://huggingface.co/hakurei/waifu-diffusion-v1-4>.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018.
- E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models, 2021.
- Huawei. Huawei presents ‘unfinished symphony’. 2019. URL <https://consumer.huawei.com/uk/campaign/unfinishe Symphony/>.
- L. Katz. Ai drew this gorgeous comics series. you’d never know it, Dec 2022. URL <https://www.cnet.com/culture/ai-drew-this-gorgeous-comics-series-youd-never-know-it/>.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes, 2022.
- llyasviel. Llyasviel/sd-controlnet-canny · hugging face. URL <https://huggingface.co/llyasviel/sd-controlnet-canny>.
- Lykon. One piece (wano saga) style lora | stable diffusion lora | civitai. URL <https://civitai.com/models/4219/one-piece-wano-saga-style-lora>.
- midjourney. URL <https://www.midjourney.com/>.
- OpenAI. Gpt-4 technical report, 2023.
- N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation, 2023.
- Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface, 2023.
- Z. Song and Z. Jin. Zorinaaaaa/csc2516-project. URL <https://github.com/Zorinaaaaa/csc2516-project>.
- Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. doi: 10.1109/TIP.2003.819861.
- C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models, 2023.
- L. Zhang and M. Agrawala. Adding conditional control to text-to-image diffusion models, 2023.