

Carleton University

COMP 4905 – Honours Project

TactileNet: Exploring Model Optimization and Dataset Expansion

Youssif Ashmawy

Supervisor: **Dr. Majid Komeili**

Date: August 04, 2025

Abstract

In this work, we expanded the TactileNet dataset (1,200 tactile images) to incorporate feedback text with the help of 132 benchmark-evaluated examples by human experts. We added structured annotations and incorporated them into a reward-guided training pipeline using Stable Diffusion fine-tuned with LoRA adapters trained on our dataset https://huggingface.co/datasets/youssifashmawy/training_dataset. Multiple loss functions and reward weightings were tested, along with extended training at reduced learning rates. Experiments showed that CLIP-heavy weighting improved contour clarity for some classes, while a balanced weighting better preserved structural fidelity in others. Longer training runs yielded small stability gains but also signs of overfitting, which degraded tactile readability. Despite these advances, dataset scarcity and the reliance on a frozen CLIP model as a reward signal limited performance, as CLIP did not evaluate tactile quality or guideline compliance. These results provided a foundation for future work on expanding datasets, applying preference-based reinforcement learning, and leveraging large language models to generate feedback-guided image edits toward fully guideline-compliant tactile graphics.

Acknowledgments

I want to take this opportunity to express my thanks to all those who supported me throughout my project journey. I would like to express my deepest gratitude to my supervisor, Dr. Majid Komeili, for his invaluable guidance, feedback, and support throughout the course of this project. I am also sincerely thankful to Adnan Khan, whose mentorship, technical insights, and constant encouragement were essential in navigating the challenges of this research.

Contents

1	9
1.1 Motivation	9
1.2 Background	9
1.3 Objectives	10
1.4 Gap Analysis	10
1.5 Report Structure	12
2 Literature Review	14
2.1 Tactile Graphics and Accessibility Standards	14
2.2 Evaluation Protocol	15
2.3 AI in Tactile Graphics	16
2.3.1 AI Applications in Tactile Graphics	17
2.3.2 Generative Models for Image Simplification	17
2.3.3 Diffusion Models for Tactile and Accessible Image Generation	19
2.3.4 Large Language and Multi-Modal Models for Accessibility	21
2.3.5 RL-based Frameworks	21
2.3.6 Summary	22
3 Methodology	26
3.1 Pipeline Overview	26
3.2 Dataset Expansion	27
3.2.1 Building and Structuring the Dataset	28
3.2.2 Challenges and Refinements	29
3.2.3 Summary	30
3.3 Initial Training Pipeline Implementation	30
3.3.1 Dataset Loading and Preprocessing	31
3.3.2 Reward Model Integration	32
3.3.3 Loss Function Variants	33
3.3.4 Forward and Backward Diffusion	35
3.3.5 LoRA Fine-Tuning Setup	36
3.3.6 Summary	36
3.4 Initial Testing Pipeline Implementation	37

3.4.1	Model Setup and Configuration	37
3.4.2	Test Dataset Loading	37
3.4.3	Image Generation and Evaluation	38
3.4.4	Results Compilation	38
3.4.5	Summary	39
3.5	Our Final Tuned InstructPix2pix Baseline	39
3.6	Our Final Adapter	40
3.6.1	Inputs	41
3.6.2	Outputs	41
3.6.3	Model Components and LoRA Setup	42
3.6.4	Training Procedure	43
3.6.5	Testing Procedure	43
4	Experimental Settings	46
4.1	Compute and Software Environment	46
4.2	Training Experiments	47
4.3	Testing Experiments	48
4.4	Hyperparameters	49
4.5	Model Selection	49
5	Results and Evaluations	50
5.1	Overview of Evaluation	50
5.2	Quantitative Results	50
5.2.1	Effect of Number of Epochs	51
5.2.2	Effect of Number of Inference Steps	52
5.2.3	Effect of Guidance Scale	53
5.2.4	Different Baseline Models	54
5.3	Qualitative Observations	55
5.3.1	Effect of Training Steps for Instruct-Tuned Img2Img Baseline (10k vs 15k)	56
5.3.2	Effect of Adding Final Adapter to Instruct-Tuned Img2Img Baseline	57
5.3.3	Effect of Reward Weighting on Adapter Performance	59
5.3.4	Extended Training Ablation	60
6	Summary and Conclusion	62
6.1	Limitations	62
6.2	Future Work	63
	Bibliography	64

A Source Code

67

List of Figures

1.1	Braille PAD Tactile Tablet from [VisonAid, 2025].	9
1.2	TactileNet Flow Diagram from [Khan et al., 2025].	10
1.3	Original TactileNet vs Our Project.	11
2.1	Tactile graphic of a dog.	16
2.2	Tactile graphic of a horse.	16
3.1	Initial Project's Flow Diagram.	27
3.2	Status and its corresponding rating from the CSV file.	28
3.3	Example of a dataset row from the CSV file.	28
3.4	Training Phase from the Project's Flow Diagram.	31
3.5	Testing Phase from the Project's Flow Diagram.	37
3.6	Project's Fine-Tuned Baseline Flow Diagram.	40
3.7	New Adapter Creation.	41
4.1	10-epochs models generated by " <code>train.py</code> ".	47
5.1	Qualitative comparison under different CLIP–LPIPS reward weightings.	60
5.2	Effect of prolonged adapter training (30 vs. 50 vs. 100 epochs) on a fixed airplane example.	61

List of Tables

2.1	Topics of References cited across subsections of Section 2.3: AI in Tactile Graphics.	23
2.2	Summary of References Cited in Section: AI in Tactile Graphics.	24
3.1	Bridging RL Concepts with Our Pipeline	27
3.2	An example of an accepted data entry.	29
3.3	An example of a rejected data entry due to a missing feature.	29
3.4	An example of a rejected data entry due to the presence of 3D aspects.	30
5.1	Average CLIP reward scores for all loss functions across 5, 8, and 10 epochs.	51
5.2	Average CLIP reward scores across different inference step counts for models trained for 10 epochs.	52
5.3	Average CLIP reward scores for models trained for 10 epochs under different guidance scale values.	53
5.4	CLIP scores across different loss functions using two baseline models.	54
5.5	Qualitative comparison between the original image, Stable Diffusion v1.5, and Instruction-Tuned SD model.	56
5.6	Qualitative comparison between the original image, 15k steps, and 10k steps of Instruction-Tuned SD model.	57
5.7	Qualitative comparison between the original image, baseline only, and after adding the adapter based on modified prompts.	58
5.8	Qualitative comparison between the original image, baseline only, and after adding the adapter based on having different reference image.	59

Listings

3.1	Loading and preprocessing image pairs	31
3.2	Reward model integration	32
3.3	MSE Loss	33
3.4	Clip-Weighted MSE	33
3.5	Reward-Weighted MSE	34
3.6	Inverse Reward Loss	34
3.7	Margin Reward Loss	35
3.8	LoRA Fine-Tuning Setup	36
3.9	Loading the test dataset	37
3.10	Image generation during testing	38
3.11	LoRA on UNet and Text Encoder; VAE frozen and adapters saved	42
3.12	Evaluation configuration toggles	44
A.1	train.py	67
A.2	test.py	75
A.3	analyzer.py	78
A.4	"train2.py" for final adapter creation	80
A.5	"test2.py" for final adapter and instruct-tuned img2img baseline testing .	93

Chapter 1

1.1 Motivation

Access to graphical information remains a significant barrier for individuals who are blind and visual impaired (BLV). Tactile graphics are crucial for conveying visual concepts such as diagrams, maps, and illustrations, but their production is time-consuming, costly, and dependent on specialist expertise. Devices like the Dot Pad (Figure 1.1), the first smart tactile graphics display, now allow users to instantly access tactile graphics from digital sources by converting any input on the connected device into touchable graphics [VisonAid, 2025]. While this technology opens up new possibilities, scaling the generation of tactile graphics still requires not only advanced generative models but also a robust, well-structured dataset that reflects tactile readability standards.



Figure 1.1: Braille PAD Tactile Tablet from [VisonAid, 2025].

1.2 Background

TactileNet [Khan et al., 2025] is a deep learning framework designed to generate tactile graphics from text and image prompts using Stable Diffusion (SD) fine-tuned with LoRA and DreamBooth. Earlier versions of TactileNet produced outputs that visually aligned with reference images but often failed to meet tactile readability requirements, including excessive detail and missing essential tactile features. A shift from multiple

category-specific LoRA adapters to a single unified adapter further reduced tactile quality, revealing that dataset limitations rather than model architecture were the primary bottleneck.

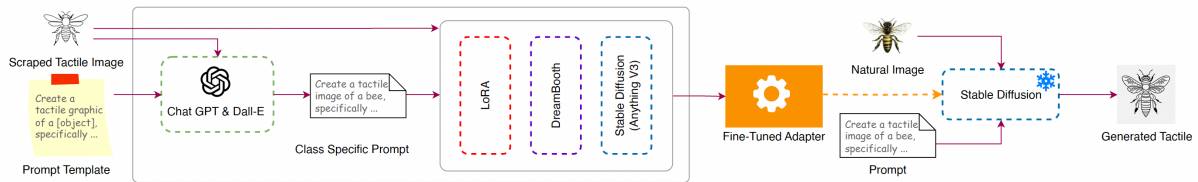


Figure 1.2: TactileNet Flow Diagram from [Khan et al., 2025].

1.3 Objectives

This project focuses on addressing the limitations of scalability and generalization by expanding and structuring the dataset and using it in the reinforcement learning (RL) inspired SD frameworks. Given the small sample of 132 expert evaluated images, the dataset was created containing over 1,200 image pairs, and each pair was annotated with targeted feedback from tactile professionals and classified into four evaluation categories: *Accept as Is*, *Minor Edit*, *Major Edit*, and *Reject*, according to some questions and guidelines that need to be met as seen in the sample size. This richer dataset now provides the foundation to experiment with different training strategies, including alternative loss functions, augmentation methods, and model refinements, with the objective of enhancing TactileNet’s ability to generate tactile graphics that are both accessible and practical for visually impaired users.

1.4 Gap Analysis

The initial TactileNet framework demonstrated the feasibility of generating tactile graphics using Stable Diffusion (SD) fine-tuned with LoRA and DreamBooth. However, several limitations identified in that study created clear gaps that needed to be addressed for the system to progress beyond proof-of-concept.

First, the dataset was both small, containing only 1,029 tactile images across 66 classes. This scarcity restricted the model’s ability to generalize and resulted in inconsistent quality across object types, particularly for organic shapes and classes with limited representation.

Second, the original approach relied on 66 separate LoRA adapters, one for each class, which achieved higher fidelity for individual categories but was not scalable. Attempts

to shift toward a single, unified adapter caused a noticeable drop in tactile readability, revealing that the dataset and feedback mechanisms were not sufficient to support a generalized model.

Third, while the evaluation process in the original study included expert feedback and a four-level quality rating that were mentioned earlier, these assessments were used only for reporting and not systematically incorporated into the dataset for model training or refinement. As a result, the model lacked targeted, structured feedback that could guide future iterations toward producing guideline-compliant tactile graphics.

Finally, the previous study proposed several avenues for improvement including expanding the dataset, experimenting with alternative loss functions, and refining data augmentation techniques, but these directions were left for future work and were not fully explored within the scope of that study.

This project directly addresses these gaps by:

1. **Expanding the feedback dataset tenfold** to over 1,200 image pairs and balancing weakly represented categories.
2. **Structuring each image pair with targeted feedback and evaluation categories** aligned with BANA guidelines.
3. **Experimenting with training strategies** including loss functions, augmentation methods, and unified adapter setups to explore potential improvements in tactile readability and scalability of the TactileNet framework.

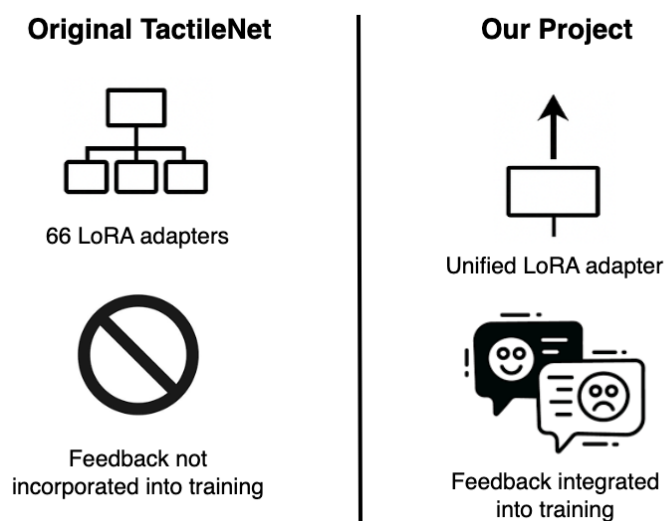


Figure 1.3: Original TactileNet vs Our Project.

1.5 Report Structure

The present report is organised into six chapters, including this Introduction chapter. An overview of each chapter is provided below:

Chapter 2 – Literature Review

This chapter reviews the key concepts, standards, and research relevant to generating tactile graphics for visually impaired users. It covers tactile graphics design guidelines and accessibility standards, the *TactileNet* evaluation protocol for assessing tactile images beyond conventional metrics, and how artificial intelligence is being applied through generative models, diffusion methods, reinforcement learning (RL)-based frameworks, and multimodal systems to produce clearer, touch-readable graphics.

Chapter 3 – Methodology

This chapter explains the overall system design, detailing dataset preparation, model setup, training strategies, and evaluation methods. It also introduces the reward model and custom loss functions integrated into the pipeline, as well as the final LoRA-based adapter configuration and the fine-tuned Stable Diffusion baseline.

Chapter 4 – Experimental Settings

This chapter describes the training and testing configurations, including hyperparameters, epoch counts, and loss function variations. The rationale behind these choices and how they relate to testing underfitting, overfitting, and model performance is also discussed. Different baseline models were also evaluated to observe their influence on overall generation quality and alignment with tactile feedback. Experiments further include the evaluation of the final LoRA-based adapter and the fine-tuned Stable Diffusion baseline trained on the expanded feedback dataset, alongside variations in reward weighting, extended training durations, and inference settings to assess their impact on tactile readability and adherence to BANA guidelines.

Chapter 5 – Results and Evaluation

This chapter presents quantitative and qualitative results from the experiments, comparing different models, loss functions, and training settings. It also evaluates the visual fidelity, CLIP reward scores, and the overall impact of each approach. Results include a detailed comparison between the final LoRA-based adapter and the fine-tuned Stable Diffusion baseline, as well as ablation studies on CLIP-LPIPS reward weighting, extended training durations, and inference parameter variations. Both numerical metrics and qualitative assessments are used to analyze tactile readability, structural fidelity, and compliance with BANA guidelines.

Chapter 6 – Summary and Conclusion

This chapter summarizes the key findings, discusses limitations, and draws final conclusions. It also outlines potential directions for future work.

Chapter 2

Literature Review

This chapter reviews the foundations and recent developments relevant to tactile graphics for people with visual impairments. It begins by outlining tactile graphics and accessibility standards that guide the creation of touch-readable materials. It then summarizes the evaluation protocol proposed in the *TactileNet* [Khan et al., 2025], which defines how the quality of tactile images should be assessed beyond traditional metrics. Finally, the chapter examines how artificial intelligence (AI) is being applied to tactile graphics, covering generative models for image simplification, diffusion-based methods for controllable image editing, and large language and multi-modal models that enable richer, more accessible descriptions. Together, these sections provide context for how tactile graphics are designed, evaluated, and increasingly automated using AI.

2.1 Tactile Graphics and Accessibility Standards

Tactile graphics are essential for providing access to visual information such as diagrams, charts, and illustrations for people with visual impairments. Easy access to such information improves literacy, education, employment, and independence for individuals who are blind or visually impaired [Edman, 1992]. Rosenblum and Herzberg (2015) found that students who used tactile graphics in their mathematics and science classes defined a "good" tactile graphic primarily by the clarity of the information presented.

Traditionally, creating tactile graphics is a manual and highly specialised process. Producing even a simple tactile diagram can take 10–15 minutes, while complex images may require hours of work by accessibility specialists. Tools like CorelDRAW, Adobe Illustrator, and embossers such as ViewPlus are commonly used, but the process remains time-consuming and costly [Khan et al., 2025].

There are also established resources and libraries for tactile graphics, including those from the American Printing House for the Blind (APH) [APH, 2024] and Perkins [Perkins, 2025]. The Braille Authority of North America (BANA) [BANA, 2025] provides formal guidelines to standardise Braille and tactile graphics. These guidelines cover requirements such as line thickness, symbol spacing, and overall clarity to ensure that tactile graphics are both usable and understandable. Teachers and instructors working with visually impaired students are expected to follow these standards when developing tactile materials.

This project draws directly on these guidelines to ensure that generated images and dataset annotations reflect best practices for tactile readability.

2.2 Evaluation Protocol

In the *TactileNet* paper, Khan et al. designed a custom evaluation protocol to assess the quality and usability of generated tactile graphics, since conventional metrics like FID and SSIM do not capture tactile-specific requirements such as line clarity, texture, and adherence to accessibility standards. Expert reviewers evaluated 132 tactile graphics (66 generated and 66 sourced) against 66 reference natural images using a side-by-side interface, without knowing whether the tactile graphic was generated or sourced.

Evaluators answered four key questions:

1. **Natural Features & Posture Alignment:** Does the tactile image accurately match the pose and key features of the reference image?
2. **Adherence to Tactile Graphics Guidelines:** Does the graphic follow BANA and similar tactile design standards?
3. **Quality Rating:** Categorized as *Accept as Is*, *Minor Edits*, *Major Edits*, or *Reject*.
4. **Optional Feedback:** Free-text comments suggesting specific improvements (e.g., line clarity, texture, or structure).

Natural images were carefully selected to pair with gold-standard tactile graphics from established libraries such as APH and Perkins. For each class, evaluators compared:

- **Sample 1:** A generated tactile graphic (created from original and paraphrased prompts).
- **Sample 2:** A sourced tactile graphic from trusted libraries.

This combined approach leveraged expert judgment along with limited SSIM checks to ensure that generated tactile graphics were not only structurally sound but also practically usable by individuals with visual impairments.

There are two examples below of tactile images that follow BANA and similar tactile design guidelines. In Figure 2.1, the dog illustration demonstrates several important concepts. The image is drawn from a side view, eliminating all 3D aspects. The ear, eye, and nose are represented with solid black shapes that are clearly separated from the body for distinction. The outer contour of the body is drawn with a thicker, darker line than the inner details. A texture pattern appears on the dog's body, but it is omitted from the far legs to make them distinguishable from the nearer legs. The tail is positioned slightly apart from the body to prevent confusion for tactile readers.

In Figure 2.2, the horse image also applies these principles. The image is also drawn from a side view, eliminating all 3D aspects. Light grey is used to indicate the far legs, while black defines the lower parts of the legs. A darker grey is used for the mane and tail to convey a texture different from the body. Finally, black is used again to mark the nostril, mouth, eyes, and ears for clear identification.

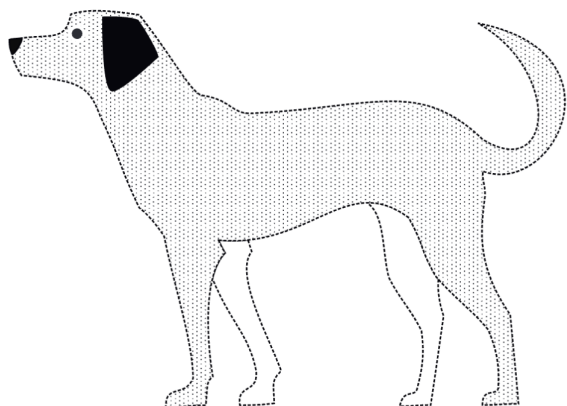


Figure 2.1: Tactile graphic of a dog.

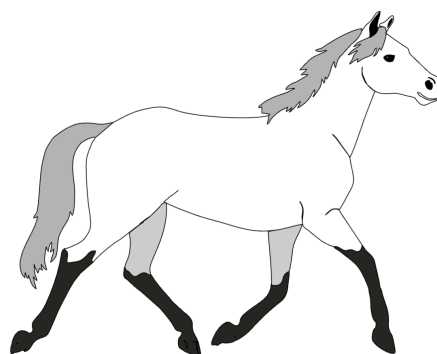


Figure 2.2: Tactile graphic of a horse.

2.3 AI in Tactile Graphics

Recent research has increasingly explored how generative AI can support the creation of tactile graphics for visually impaired users. This includes systems that generate tactile-ready images and interactive tools, methods for simplifying complex visuals using models like GANs, and newer diffusion-based approaches that allow precise, instruction-driven edits. Together, these advances show how AI is shifting tactile graphics from a manual

process to one that is more scalable and accessible.

2.3.1 AI Applications in Tactile Graphics

Dzhurynskyi et al. (2024) introduced a model that combines a Bidirectional and Auto-Regressive Transformer (BART) with a Vector Quantized Variational Auto-Encoder (VQ-VAE) to generate tactile graphics from text prompts. Their system was evaluated using metrics such as cross-entropy, perplexity, and CLIP Score, and testing with educational and rehabilitation institutions showed it could drastically reduce production time and the need for specialized expertise. This work highlights how advanced generative architectures can automate what has traditionally been a highly manual process.

Other studies have focused on converting existing visual materials into tactile forms. For instance, Pakenaite et al. (2024) presented Pic2Tac, a system that translates semantic content from photographs into tactile symbols, demonstrating strong results in user studies with both sighted and visually impaired participants. Similarly, Tanaka et al. (2023) explored how Generative AI and 3D printing can transform ordinary picture books into tactile picture books, opening up new opportunities for inclusive learning experiences.

Beyond purely visual-to-tactile conversion, researchers have also investigated ways to enhance tactile graphics with real-time interactivity. The TAURIS project [Zeinullin, 2025], a PhD study involving 20 participants, developed a fingertip-tracking system for mobile devices that provides instant audio descriptions as users explore tactile graphics. This system not only improved the speed and accuracy of interaction but also demonstrated a notable advantage in memory retention compared to screen readers.

Work has also begun to explore the use of large language models for multimodal educational tools. Wu et al. (2025) developed TaleVision, which uses LLMs to generate touchable images and sound effects for visually impaired children. In testing with 19 participants, nearly all reported high satisfaction and positive experiences, suggesting strong potential for combining tactile and auditory outputs to create richer learning tools.

Collectively, these projects illustrate the breadth of AI applications in tactile graphics: from automating graphic generation and enhancing printed materials, to introducing interactivity and multimodal content, marking an important shift toward scalable, user-centered solutions for accessibility.

2.3.2 Generative Models for Image Simplification

Generative Adversarial Networks (GANs) have become a cornerstone for image simplification tasks relevant to tactile graphics, such as converting photographs into line drawings

or stripping away visual clutter. Early GAN frameworks relied on end-to-end learning, generating images directly from uniform noise distributions. While effective, these early models lacked interpretability and offered little control over what aspects of the image were simplified.

Decomposing Structure and Style

More recent work has sought to decompose the generation process for greater structure and flexibility. For example, [Wang and Gupta, 2016] introduced the Style and Structure GAN, which splits image generation into two coordinated components: Structure-GAN, responsible for producing a surface normal map, and Style-GAN, which applies textures onto that structure. This separation of structure from style allows models to preserve essential shapes while modifying or simplifying textures, a capability highly relevant to tactile graphics, where clear boundaries must be retained even as unnecessary details are abstracted away.

Image-Driven Simplification Approaches

Other researchers have explored alternatives to purely geometry-driven simplification methods, which often strip away perceptually important information.

[Lindstrom and Turk, 2000] proposed an image-driven simplification approach that compares rendered images of the original and simplified models to guide the simplification process. This framework maintains high-fidelity silhouettes, removes visually insignificant regions, and strikes a balance between geometric and appearance-based trade-offs. Although developed for 3D models, the principle, prioritizing perceptually important structures while simplifying less critical elements, aligns closely with the goals of tactile graphics, where maintaining clear, touchable shapes matters more than reproducing every visual nuance.

Stylization for Tactile-Friendly Outputs

GAN-based models have also been adapted for stylization tasks that naturally support tactile interpretation. CartoonGAN [Chen et al., 2018], for instance, converts real-world photos into cartoon-style images with bold edges and simplified textures, producing outputs that more closely resemble tactile-ready line art.

Pix2Pix and Its Improvements

Another foundational model, Pix2Pix, has been widely used for image-to-image translation tasks such as sketch-to-image and outline-to-photo conversion. However, standard Pix2Pix often loses critical structural information during encoding and decoding, and its training lacks constraints to maintain fine details. Zhao et al. (2023) addressed these limitations by integrating a U-Net generator with denser skip connections to reduce information loss and adding a differential image discriminator to improve supervision. Their improved Pix2Pix variant demonstrated significantly better image quality and structural preservation: an important advance for tactile graphics, where even the smallest lost edge or contour can compromise the tactile readability of the output.

2.3.3 Diffusion Models for Tactile and Accessible Image Generation

Diffusion models have rapidly become the dominant paradigm for high-quality image generation, offering state-of-the-art performance in controllable, prompt-driven synthesis. Unlike GANs, which map noise to images in a single step, diffusion models iteratively denoise random noise, providing more stable training and finer control over the generation process. This controllability has made diffusion approaches especially attractive for accessibility-focused research, where precision, structure, and repeatability are essential.

From Visual Fidelity to Tactile Content

Early research has started extending diffusion models beyond purely visual fidelity to tackle tactile and multimodal tasks. Dou et al. (2024) introduced *Tactile-Augmented Radiance Fields* (TARF), a system that unifies vision and touch into a shared 3D space. Their conditional diffusion model takes an RGB image and depth map from a neural radiance field and produces corresponding tactile "images." By training on the largest dataset of spatially aligned visual–tactile pairs to date, TARF demonstrated that diffusion models can accurately translate visual cues into tactile features, an important step toward AI-driven tactile data generation.

Instruction-Driven Editing

A second major development has been the use of diffusion models for instruction-driven image editing, where text commands guide specific changes to images. Paul (2023) presented the open-source project *Instruction-Tuning Stable Diffusion*, inspired by FLAN

[Wei et al., 2021] and InstructPix2Pix [Brooks et al., 2023], which fine-tunes Stable Diffusion (SD) to follow explicit instructions paired with input images (e.g., "apply a cartoon filter"). This ability to modify images through natural language commands is directly relevant to tactile graphics: simplification, outlining, or re-texturing for tactile use can be framed as targeted "instructions."

Brooks et al. (2023) expanded this idea with *InstructPix2Pix*, a conditional diffusion model trained on a large dataset of synthetic image–instruction pairs created using GPT-3 and Stable Diffusion. InstructPix2Pix performs edits in a single forward pass, without requiring fine-tuning or inversion for each image, making it extremely fast and flexible. This framework is highly applicable to tactile graphics generation, where edits like removing shading, emphasizing contours, or simplifying textures must be executed quickly and consistently.

Shagidanov et al. (2024) introduced *Grounded-Instruct-Pix2Pix*, building on Brooks et al.'s work by adding automatic "target grounding." This ensures that edits only affect the intended area (e.g., just the "cat" in a complex background). Such localized precision is crucial for tactile graphics, as it allows selective simplification or re-texturing of certain objects while leaving the surrounding scene intact.

Spatial Conditioning with ControlNet

While instruction-driven editing handles "what" to change, spatial conditioning handles "where" to change it. Zhang et al. (2023) addressed this with *ControlNet*, a neural architecture that adds spatial controls to large, pretrained diffusion models like Stable Diffusion. ControlNet preserves the base model's pretrained knowledge while adding "zero convolution" modules (zero-initialized layers that progressively learn how to respond to control inputs like edges, segmentation maps, or human poses). This means users can tightly guide the generation process with structural cues, an ability particularly promising for tactile graphics. For example, a ControlNet-based system could generate tactile graphics that precisely follow edge maps or segmentation boundaries, ensuring tactile clarity and alignment with visual structure.

Implications for Accessibility

Together, these developments position diffusion models as a powerful toolkit for accessibility research. Their ability to incorporate multimodal inputs, follow human-readable instructions, and respect structural guidance means they can not only generate high-quality images but do so in ways that meet the unique demands of tactile graphics. This opens the door to scalable, user-centered systems capable of producing tactile images that are

both accurate and touch-friendly.

2.3.4 Large Language and Multi-Modal Models for Accessibility

Large Language Models (LLMs) are increasingly being explored as tools for bridging text and vision, offering new ways to structure, describe, and prioritize visual information for accessibility. Rather than simply labeling images, these models can provide richer semantic context and enable new forms of interaction that are particularly relevant to tactile graphics.

Bertolotto et al. (2021) demonstrated how LLMs can enrich vision-based systems by generating semantic descriptions for class labels. Instead of treating classes as isolated categories, their method used an LLM to create descriptive text for each label, which was then embedded and used as the training target for image classification. This approach encouraged the model to learn a hierarchically structured feature space, leading to more semantically meaningful errors (e.g., confusing a "dog" for a "wolf" rather than for a "car"). For tactile graphics, such semantic structuring could help generative models decide which visual concepts to emphasize and how to frame prompts, ensuring that outputs preserve the most relevant and interpretable features for touch.

Building on this theme of richer interpretation, Ricci et al. (2024) expanded the scope of image captioning by introducing a dialogue-based approach using ChatGPT. Their system moved beyond the typical single-sentence caption, instead generating a multi-turn question-and-answer dialogue about an image and summarizing the exchange into enriched, context-aware captions. For tactile graphics, such a dialogue-driven captioning framework could generate layered explanations for complex images: first summarizing the overall scene, then adding detail about key objects or relationships, helping visually impaired users build a more comprehensive mental model of the content.

Collectively, these works highlight how LLMs and multi-modal systems can go beyond basic captioning or labeling, enabling structured, meaningful descriptions that could feed directly into tactile graphics generation pipelines, supporting not only clearer image simplification but also richer, more accessible explanations of visual content.

2.3.5 RL-based Frameworks

While GANs, diffusion models, and large language models provide powerful generative and editing capabilities, reinforcement learning (RL) offers a complementary paradigm: directly optimizing policies to maximize accessibility-focused objectives. A central challenge in applying RL to real-world problems is specifying goals that are both precise and

practical to obtain.

Christiano et al. (2017) introduced deep reinforcement learning from human preferences, in which agents learn a reward model from pairwise comparisons of short trajectory segments provided by non-expert annotators. This separates learning the reward function from learning the policy, enabling the agent to acquire complex behaviors with minimal human oversight. Ziegler et al. (2019) extended this paradigm to language models, demonstrating that fine-tuning from human preferences (RLHF) can align generative models with nuanced stylistic and semantic goals.

Rocamonde et al. (2024) further showed that pretrained vision language models (VLMs) can act as zero-shot reward models, producing a reward signal by measuring similarity between generated outputs and natural language task descriptions. Such VLM-based rewards could complement human ratings by automatically assessing whether a tactile output meets accessibility-focused prompts, for example “a simplified, high-contrast outline with Braille labels”, and provide scalable and consistent feedback.

Together, these works suggest that RL-inspired optimization, whether through human preferences [Christiano et al., 2017; Ziegler et al., 2019] or multimodal reward models [Rocamonde et al., 2024], provides a promising route to develop tactile generation systems that are adaptive, user-centered, and scalable.

2.3.6 Summary

This chapter reviewed how AI is transforming the creation and use of tactile graphics for visually impaired users (Table 2.1). First, it explored applications of AI in tactile graphics, including systems that generate tactile-ready images, convert visual materials like books and photos into tactile forms, and add interactive features such as fingertip-tracked audio descriptions. Next, it examined generative models for image simplification, showing how GAN-based approaches like StyleGAN, Pix2Pix, and CartoonGAN simplify visuals while preserving essential structure for touch interpretation. It then discussed diffusion models, which offer instruction-driven editing, spatial control, and multimodal capabilities through frameworks like InstructPix2Pix, ControlNet, and TARF, opening new possibilities for precise, scalable tactile content creation. Finally, it highlighted the role of large language and multimodal models, demonstrating how LLMs like ChatGPT can enrich image descriptions and guide generative pipelines toward semantically meaningful, touch-friendly outputs. Building on this, reinforcement learning (RL)-based frameworks offer a complementary paradigm by directly optimizing models with accessibility-focused reward signals, whether derived from human preferences or multimodal evaluators.

Together, these advances mark a shift from manual, specialist-driven production of tactile

graphics to AI-powered, scalable systems capable of generating richer, clearer, and more accessible tactile experiences (Table 2.2).

Table 2.1: Topics of References cited across subsections of Section 2.3: AI in Tactile Graphics.

Topic	# of References
AI Applications in Tactile Graphics	5
Generative Models for Image Simplification	4
Diffusion Models for Tactile and Accessible Image Generation	5
Large Language and Multi-Modal Models for Accessibility	2
RL-based Frameworks	3
Total	19

Table 2.2: Summary of References Cited in Section: AI in Tactile Graphics.

Citation	Description
Dzhurynskyi et al. (2024)	Introduced a BART + VQ-VAE model for generating tactile graphics from text prompts.
Pakenaite et al. (2024)	Developed Pic2Tac, converting photos into tactile symbols using semantic information.
Tanaka et al. (2023)	Explored converting printed picture books into tactile books using generative AI and 3D printing.
Zeinullin (2025)	Developed TAURIS, a fingertip-tracking system for tactile graphics with audio feedback.
Wu et al. (2025)	Built TaleVision, using LLMs to generate touchable and audible content for visually impaired children.
Wang et al. (2016)	Proposed Style and Structure GAN to separate texture and shape during image generation.
Lindstrom (2000)	Introduced image-driven simplification by comparing rendered images to preserve perceptually important shapes.
Chen et al. (2018)	Created CartoonGAN for photo-to-cartoon conversion with bold outlines and flat textures.
Zhao et al. (2023)	Improved Pix2Pix using U-Net and differential image discriminator to enhance structure preservation.
Dou et al. (2024)	Proposed TARF: a conditional diffusion model to generate tactile features from RGB + depth inputs.
Paul (2023)	Released Instruction-Tuned Stable Diffusion for applying text-driven modifications to images.
Brooks et al. (2023)	Created InstructPix2Pix, enabling fast instruction-based edits without per-image fine-tuning.
Shagidanov et al. (2024)	Introduced Grounded-Instruct-Pix2Pix for localized, instruction-specific image editing.
Zhang et al. (2023)	Developed ControlNet, enabling spatial conditioning for diffusion models using edge/pose maps.
Bertolotto et al. (2021)	Used LLMs to create semantic-rich label embeddings, improving classification hierarchy.
Ricci et al. (2024)	Proposed dialogue-driven image captioning with ChatGPT for richer semantic outputs.
Christiano et al. (2017)	Introduced deep reinforcement learning from human preferences, enabling agents to learn reward models from pairwise human comparisons.

Summary of References Cited in Section: AI in Tactile Graphics (continued).

Citation	Description
Ziegler et al. (2019)	Applied preference-based reinforcement learning to fine-tune language models, aligning outputs with nuanced stylistic and semantic goals.
Rocamonde et al. (2024)	Demonstrated that vision-language models can serve as zero-shot reward models for reinforcement learning, enabling scalable multimodal reward estimation.

Chapter 3

Methodology

3.1 Pipeline Overview

Referring to Figure 3.1, the workflow of this project can be divided into three main phases. The first phase is dataset expansion. Each image pair in the dataset contains four elements: the tactile image, the corresponding natural image, written feedback highlighting what is wrong with the tactile image based on evaluator questions from the *TactileNet* paper [Khan et al., 2025] (whether the tactile image accurately matches the pose and key features of the reference image, and whether it follows tactile graphics guidelines such as those from BANA), and finally a rating of the tactile image using the natural image as the reference.

The second phase is the training phase. We used a Python script to fine tune a Stable Diffusion model with LoRA adapters to generate tactile images using a reward guided approach inspired by reinforcement learning (RL) (Table 3.1). The script loads paired tactile and natural images, uses a CLIP reward model to score how well the generated images align with the provided instructions, and applies several custom loss functions within a tailored training loop. Nonessential parts of the model are frozen to improve efficiency, training progress is logged with Weights & Biases, and the resulting LoRA adapters are saved to create a streamlined pipeline for instruction guided and accessibility focused image generation.

The third phase is the testing phase. A separate Python script evaluates the LoRA fine tuned Stable Diffusion models by generating and scoring images from a test dataset. It processes a set of prompts and input images, applies each model, and uses the CLIP reward model to calculate semantic alignment scores between the generated images and their prompts. For each model, the generated images are saved and the evaluation results are compiled into a CSV file. This phase forms the evaluation and ablation stage of the

pipeline, providing a clear comparison of how different training strategies affect image quality and alignment with prompts.

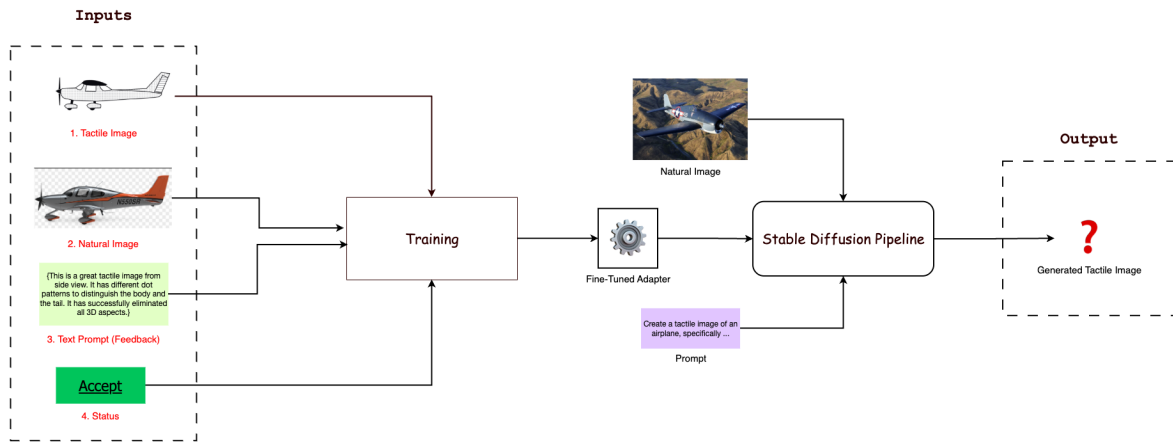


Figure 3.1: Initial Project's Flow Diagram.

Table 3.1: Bridging RL Concepts with Our Pipeline

RL Concept	Our Setup
Agent	Stable Diffusion model (with or without LoRA adapter)
Environment	Training pipeline with feedback-based dataset
Action	The generated tactile image
State	The natural image + instruction (prompt)
Reward	CLIP similarity score between generated image and instruction

3.2 Dataset Expansion

The *TactileNet* paper highlighted the need for a larger dataset, which led me to expand the collection by evaluating additional image pairs. Understanding how to evaluate these pairs was challenging, as the criteria became clearer only through the process itself, with each new example revealing gaps in my understanding and prompting refinements to the evaluation framework.

Two main sources guided this process: (1) earlier expert feedback and (2) the BANA (Braille Authority of North America) guidelines for tactile graphics. The expert feedback, while useful, was simplified and often vague, which made it less effective for training purposes. To address this, we reworked the feedback into clearer, more directed instructions and linked each comment to specific BANA standards when relevant.

3.2.1 Building and Structuring the Dataset

Originally, the expert feedback was stored in a Word document. We reformatted it into an Excel sheet to make it compatible with the baseline code and easier for future contributors to expand. This sheet was then converted into CSV and JSONL formats so the script could read the data directly. These files, along with the image folders, were uploaded to Google Drive and linked in the project's GitHub README for open access.

The dataset originally had four outcome labels: *Accept as is*, *Minor Edit*, *Major Edit*, and *Reject*. For simplicity and better training consistency, these were consolidated into two categories:

- **Accept:** Accept as is & Minor Edit
- **Reject:** Major Edit & Reject

Status	rating
Reject	reject
Major Edit	reject
Accept as is	accept
Minor Edit	accept

Figure 3.2: Status and its corresponding rating from the CSV file.

Each entry in the dataset now includes:

1. The natural image
2. The tactile image
3. A clear directed feedback statement explaining what is wrong and how to fix it
4. A rating (Accept/Reject)

Category	T Filename	N Filename	feedback_text	rating
Banana	1.png	1.jpg	Rotate image 45 degrees anticlockwise to match the posture of the original image.	accept

Figure 3.3: Example of a dataset row from the CSV file.

3.2.2 Challenges and Refinements

The most significant challenge was the vague nature of the expert feedback. It lacked the directness needed for a model to learn corrective actions. We rewrote much of this feedback to explain not only *what* was wrong but also *why* it was wrong and *how* to fix it.

Common problems found in tactile graphics included:

- Missing or misleading textures
- Inaccurate posture or object orientation
- Unnecessary shading that confused tactile
- Poor representation of key features compared to the natural image

Example of an Accepted Case:

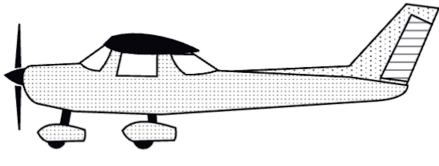

 <p style="text-align: center;">Airplane Tactile Image</p>	 <p style="text-align: center;">Airplane Natural Image</p>
<p>Feedback: This is a great tactile image from side view. It has different dot patterns to distinguish the body and the tail. It has successfully eliminated all 3D aspects.</p>	
<p style="text-align: center;">Rating: accept</p>	

Table 3.2: An example of an accepted data entry.

Example of Reject Cases:



 <p style="text-align: center;">Airplane Tactile Image</p>	 <p style="text-align: center;">Airplane Natural Image</p>
<p>Feedback: Add the missing landing gear.</p>	
<p style="text-align: center;">Rating: reject</p>	

Table 3.3: An example of a rejected data entry due to a missing feature.



 <p style="text-align: center;">Airplane Tactile Image</p>	 <p style="text-align: center;">Airplane Natural Image</p>
<p>Feedback: Remove 3D aspects by showing a proper front view of the natural image. Adjust proportions to match the reference image. It doesn't resemble the reference airplane in key structural ways (engines, wing shape, windows, and emergency doors).</p>	
<p style="text-align: center;">Rating: reject</p>	

Table 3.4: An example of a rejected data entry due to the presence of 3D aspects.

3.2.3 Summary

This section outlined the process of expanding and restructuring the TactileNet dataset to support more effective training and evaluation. The dataset was reformatted from scattered Word documents into structured CSV and JSONL files, with each entry containing a natural image, a tactile image, clear feedback, and an accept/reject rating. Feedback from experts was rewritten to provide actionable, targeted instructions. Labels were consolidated from four to two categories (Accept and Reject) for consistency. Common issues identified in the dataset included missing textures, 3D elements, and inaccurate object proportions. Accepted and rejected examples illustrated how these refinements were tested in an attempt to improve clarity and standardization, creating a more reliable and scalable dataset for model development.

3.3 Initial Training Pipeline Implementation

The training pipeline was implemented in Python and designed to fine-tune a Stable Diffusion model for generating tactile graphics using LoRA adapters and reward-guided optimization. The script "`train.py`" is structured into distinct stages, each handling a key component of the training workflow.

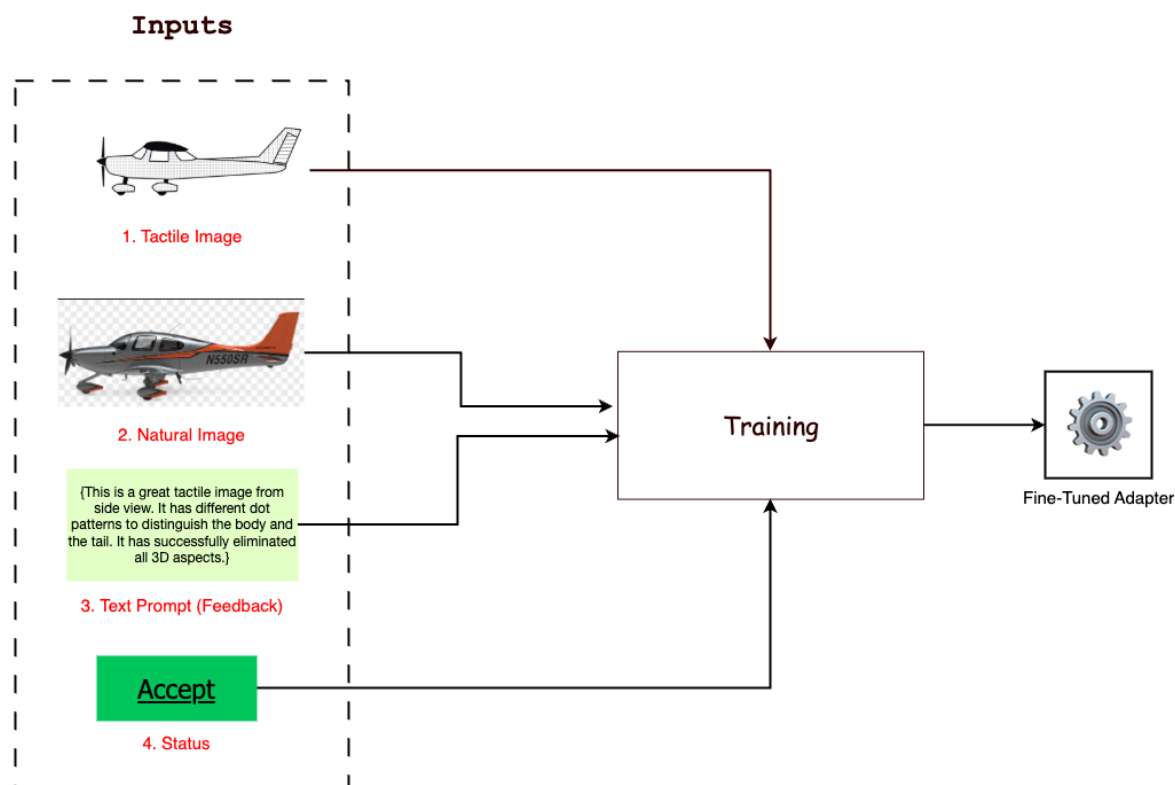


Figure 3.4: Training Phase from the Project's Flow Diagram.

3.3.1 Dataset Loading and Preprocessing

The pipeline begins by loading tactile–natural image pairs from a JSONL file. Each entry includes the tactile image, natural reference image, feedback text, and an acceptance rating. Images are resized to 512×512 pixels, converted to RGB, and preprocessed with normalization for the Variational Autoencoder (VAE). Augmentations such as horizontal flips and color jittering are applied to expand data variability.

```

1 def load_image_pairs(jsonl_path: Path) -> List[Dict]:
2     pairs = []
3     with open(jsonl_path) as f:
4         for line in f:
5             try:
6                 data = json.loads(line)
7                 rating = 1 if data["rating"].lower() == "accept"
8                     else 0
9                 input_img_path = Path(data["input_image"])
10                gen_img_path = Path(data["generated_image"])
11
12                if not input_img_path.exists() or not gen_img_path
                    .exists():
                    continue

```

```
13
14         pairs.append({
15             "input_image": Image.open(input_img_path).
16                 convert("RGB").resize(IMAGE_SIZE),
17             "generated_image": Image.open(gen_img_path).
18                 convert("RGB").resize(IMAGE_SIZE),
19             "instruction": data["feedback_text"] if data["
20                 feedback_text"].lower() != "nan" else "",
21             "rating": rating
22         })
23     return pairs
```

Listing 3.1: Loading and preprocessing image pairs

3.3.2 Reward Model Integration

To guide training, a CLIP-based reward model was implemented. The reward model calculates semantic similarity scores between generated images and corresponding text instructions, providing a quantitative "reward" signal. These scores are used later in the custom loss functions to prioritize alignment between model outputs and tactile feedback instructions.

```
1 class CLIPRewardModel:
2     def __init__(self):
3         self.model = CLIPModel.from_pretrained("openai/clip-vit-
4             base-patch32").to(DEVICE)
5         self.processor = CLIPProcessor.from_pretrained("openai/
6             clip-vit-base-patch32")
7
8     def calculate_reward(self, images: List[Image.Image], texts:
9         List[str]) -> torch.Tensor:
10         inputs = self.processor(text=texts, images=images,
11             return_tensors="pt", padding=True, truncation=True).to(
12             DEVICE)
13         with torch.no_grad():
14             outputs = self.model(**inputs)
15             sims = outputs.logits_per_image
16             diagonal = sims.diagonal()
17             return diagonal
```

Listing 3.2: Reward model integration

3.3.3 Loss Function Variants

Several custom loss functions were implemented to experiment with different training dynamics. Each formulation is described with its purpose and variable definitions.

- **MSE Loss** – standard pixel-level loss.

$$L_{MSE} = \|\hat{\epsilon} - \epsilon\|^2$$

Where:

- $\hat{\epsilon}$: noise predicted by the UNet
- ϵ : ground-truth noise sampled during training

```

1     def mse_loss_fn(pred, target, **kwargs):
2         return F.mse_loss(pred, target)

```

Listing 3.3: MSE Loss

- **Clip-Weighted MSE** – adjusts the MSE based on the CLIP similarity score.

$$L_i = (1 - \text{CLIP}(x_i, t_i)) \cdot L_{MSE}$$

Where:

- x_i : generated tactile image
- t_i : instruction or feedback text
- $\text{CLIP}(x_i, t_i)$: similarity score from CLIP model (ranges from 0 to 1)

```

1     def clip_weighted_mse_loss(pred, target, ratings, **kwargs
2         ):
3         accepted_mask = torch.tensor(ratings, device=DEVICE,
4             dtype=torch.bool)
5         loss_weights = torch.where(
6             accepted_mask,
7             torch.tensor(0.5, device=DEVICE),
8             torch.tensor(2.0, device=DEVICE)
9         )
10        base_loss = F.mse_loss(pred, target)
11        return base_loss * loss_weights.mean()

```

Listing 3.4: Clip-Weighted MSE

- **Reward-Weighted MSE** – scales the loss by rating-based weights.

$$L = w_i \cdot L_{MSE} \quad \text{where} \quad w_i = \begin{cases} 2.0, & \text{if rejected} \\ 0.5, & \text{if accepted} \end{cases}$$

Where:

- w_i : weight applied based on feedback (higher for rejected images)
- L_{MSE} : standard Mean Squared Error loss

```

1     def reward_weighted_mse_loss(pred, target, rewards=None,
2         **kwargs):
3         rewards = rewards.to(device=pred.device, dtype=pred.
4             dtype)
5         base_loss = F.mse_loss(pred, target, reduction='none')
6         weight = rewards.reshape(-1, 1, 1, 1)
7         return (base_loss * weight).mean()

```

Listing 3.5: Reward-Weighted MSE

- **Inverse Reward Loss** – penalizes predictions inversely proportional to the reward values.

$$L_i = \frac{1}{\text{CLIP}(x_i, t_i) + \epsilon} \cdot L_{MSE}$$

Where:

- $\text{CLIP}(x_i, t_i)$: CLIP similarity score
- ϵ : small constant to prevent division by zero
- L_{MSE} : Mean Squared Error loss

```

1     def inverse_reward_loss(pred, target, rewards=None, **
2         kwargs):
3         rewards = rewards.to(device=pred.device, dtype=pred.
4             dtype)
5         base_loss = F.mse_loss(pred, target, reduction='none')
6         inverse_weight = (1.0 - rewards).unsqueeze(1).
7             unsqueeze(2).unsqueeze(3)
8         return (base_loss * inverse_weight).mean()

```

Listing 3.6: Inverse Reward Loss

- **Margin Reward Loss** – applies a margin-based weighting strategy.

$$L = \text{MSE}(x_i, x'_i) \cdot \max(0, m - r_i)$$

Where:

- x_i : ground truth noise
- x'_i : model's predicted noise
- m : reward margin threshold (e.g. 0.2)
- r_i : CLIP reward score for image-text alignment

```

1     def margin_reward_loss(pred, target, rewards=None, **
      kwargs):
2         rewards = rewards.to(device=pred.device, dtype=pred.
          dtype)
3         base_loss = F.mse_loss(pred, target, reduction='none')
4         margin = 0.2
5         clipped = torch.clamp(margin - rewards.unsqueeze(1).
          unsqueeze(2).unsqueeze(3), min=0)
6         return (base_loss * clipped).mean()

```

Listing 3.7: Margin Reward Loss

These loss functions allow exploration of how reward-guided training can improve tactile image generation.

3.3.4 Forward and Backward Diffusion

This stage describes how the diffusion process is applied during training to teach the model to denoise images and update its parameters. It involves three main components:

1. **Forward Diffusion:** Each input image is encoded into a latent representation using the VAE. Random noise is then added at varying timesteps, simulating the forward diffusion process that gradually destroys image information.
2. **Reverse Diffusion Prediction:** The UNet receives the noisy latent, timestep, and text instruction embeddings, and predicts the noise that should be removed, effectively learning to reverse the diffusion process.
3. **Backward Pass (Training):** The predicted noise is compared to the true noise using the selected loss function. Gradients are computed, and LoRA parameters are updated via backpropagation with the AdamW optimizer and a cosine scheduler.

This combined approach ensures that the model learns how to denoise (reverse diffusion) while integrating guidance from instructions, ratings, and rewards during training. Training progress is logged with Weights & Biases (W&B) for monitoring.

3.3.5 LoRA Fine-Tuning Setup

Stable Diffusion’s UNet is prepared for LoRA fine-tuning by freezing the VAE and text encoder. A LoRA configuration (rank 8, dropout 0.1) targets key attention layers (`to_q`, `to_k`, `to_v`) across 16 layers, allowing the model to learn tactile-specific adjustments without overfitting or overwriting the base model’s general capabilities.

```
1   pipe.unet = prepare_unet_for_lora(pipe.unet)
2   lora_config = LoraConfig(
3       r=8,
4       lora_alpha=32,
5       target_modules=["to_q", "to_k", "to_v"],
6       layers_to_transform=list(range(16)),
7       lora_dropout=0.1,
8       bias="none",
9       fan_in_fan_out=False
10  )
11  pipe.unet = get_peft_model(pipe.unet, lora_config)
12  pipe.unet = pipe.unet.to(device=DEVICE, dtype=torch.float32)
13  pipe.vae.requires_grad_(False)
14  pipe.text_encoder.requires_grad_(False)
```

Listing 3.8: LoRA Fine-Tuning Setup

After completing all epochs, the trained LoRA adapters are saved to disk, along with metadata including run name, loss type, number of epochs, training duration, and final loss. This ensures reproducibility and supports structured ablation studies on different training configurations.

3.3.6 Summary

This section outlines the full training pipeline for fine-tuning Stable Diffusion to generate tactile graphics. It covers dataset loading and preprocessing, where tactile–natural image pairs are structured and augmented, and the integration of a CLIP-based reward model to provide semantic guidance. Multiple custom loss functions are implemented to test different training strategies, and the diffusion process is detailed, explaining how the model learns to denoise (forward and backward diffusion). Finally, a LoRA fine-tuning setup is described, targeting attention layers while freezing the VAE and text encoder.

3.4 Initial Testing Pipeline Implementation

The testing pipeline, implemented in the "test.py" script, evaluates multiple LoRA-fine-tuned Stable Diffusion models on a held-out test dataset. It systematically generates images, scores them using a CLIP-based reward model, and compiles the results for comparison across different loss functions and training strategies.

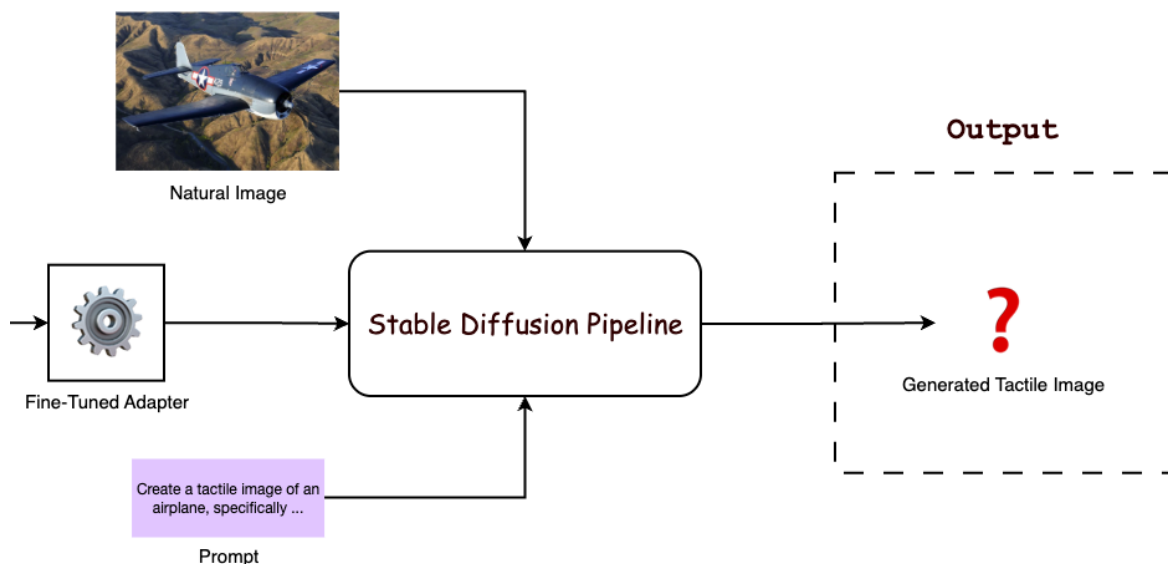


Figure 3.5: Testing Phase from the Project's Flow Diagram.

3.4.1 Model Setup and Configuration

The pipeline supports evaluating several models trained with different loss functions, including MSE, Clip-Weighted MSE, Reward-Weighted MSE, Inverse Reward Loss, and Margin Reward Loss. Each model is linked to its corresponding LoRA adapter path, while a baseline model is also included for comparison. Stable Diffusion v1-5 serves as the base model, and LoRA adapters are loaded dynamically after being created by the "train.py" script.

3.4.2 Test Dataset Loading

The test dataset is stored in a JSONL file including 60 entries, where each entry contains the natural image, feedback prompt, and a category. These samples are loaded and iterated over for evaluation.

```

1 def load_test_set(path):
2     samples = []
3     with open(path) as f:
  
```

```
4     for line in f:
5         data = json.loads(line)
6         samples.append({
7             "input_image": data["input_image"],
8             "prompt": data["feedback_text"],
9             "category": data.get("Category", "unknown")
10        })
11    return samples
```

Listing 3.9: Loading the test dataset

3.4.3 Image Generation and Evaluation

A CLIP model is used to provide semantic alignment scores between the generated images and the feedback prompts. These scores act as a quantitative measure of how well the model output aligns with textual instructions. For each model, the pipeline uses the Stable Diffusion Img2ImgPipeline to generate images based on the test inputs. The LoRA adapter is loaded, and images are generated using the prompt and the input image as guidance. Each output is scored with the CLIP reward model.

```
1 def generate_image(pipe, image_path, prompt):
2     image = Image.open(image_path).convert("RGB").resize(
3         IMAGE_SIZE)
4     generator = torch.Generator(DEVICE).manual_seed(SEED)
5     with torch.no_grad():
6         result = pipe(prompt=prompt, image=image, strength=0.8,
7             guidance_scale=7.5)
8     return result.images[0].convert("RGB")
```

Listing 3.10: Image generation during testing

3.4.4 Results Compilation

Generated images are saved into model-specific directories, and all evaluation metrics (prompt, category, reward score, and image path) are written into CSV files for later comparison. This structure enables ablation studies by making it easy to assess how different loss functions and adapters impact image quality and semantic alignment.

This testing pipeline provides a structured evaluation framework for all LoRA-tuned models. It generates visual outputs, assigns reward-based scores, and organizes results for direct comparison. The process is essential for validating which training strategies

(loss functions, adapters) best support generating tactile graphics that align with feedback instructions and accessibility standards.

3.4.5 Summary

The testing pipeline systematically benchmarks all LoRA-tuned models by generating images, scoring them with a CLIP reward model, and compiling results for comparison. This structured evaluation ensures that the most effective loss functions and training strategies are clearly identified.

3.5 Our Final Tuned InstructPix2pix Baseline

Our baseline uses the InstructPix2Pix pipeline to translate a *natural image* and a *text prompt* into a *tactile-style image* (Fig. 3.6). The model conditions denoising on both the prompt and the input image, following the original InstructPix2Pix formulation.

- **Inputs:** a natural (reference) image, tactile image and an instruction prompt (e.g., “create a tactile image of a bird”).
- **Output:** a generated tactile image aligned with the prompt and the input visual content.

Architecture and training configuration

- **Backbone:** Stable Diffusion InstructPix2Pix with an 8-channel UNet input (4 for noisy latents + 4 for the conditioned image embedding).
- **Frozen components:** VAE and CLIP text encoder (`requires_grad_(False)`).
- **Fine-tuned component:** UNet only (all UNet parameters optimized).
- **Objective:** mean-squared error on the predicted noise (standard diffusion loss).
- **Conditioning:** text tokens from the prompt (via the frozen CLIP text encoder) and latent features of the input image (via the frozen VAE). Optional conditioning dropout is supported to enable classifier-free guidance at inference.

Implementation and Usage

Training uses triplets $\langle \text{input image, target image, edit prompt} \rangle$ provided via a `metadata.jsonl` file. Images are resized and normalized, and prompts are tokenized with the CLIP tokenizer.

At inference, the pipeline takes a new natural image and prompt, generating a tactile image through a single `img2img` pass with fixed guidance and step settings for consistent comparison across experiments.

This baseline isolates the capabilities of a strong instruction-guided `img2img` model *without* reward modeling or LoRA adapters, serving as a controlled reference point for later feedback-guided fine-tuning and adapter-based variants.

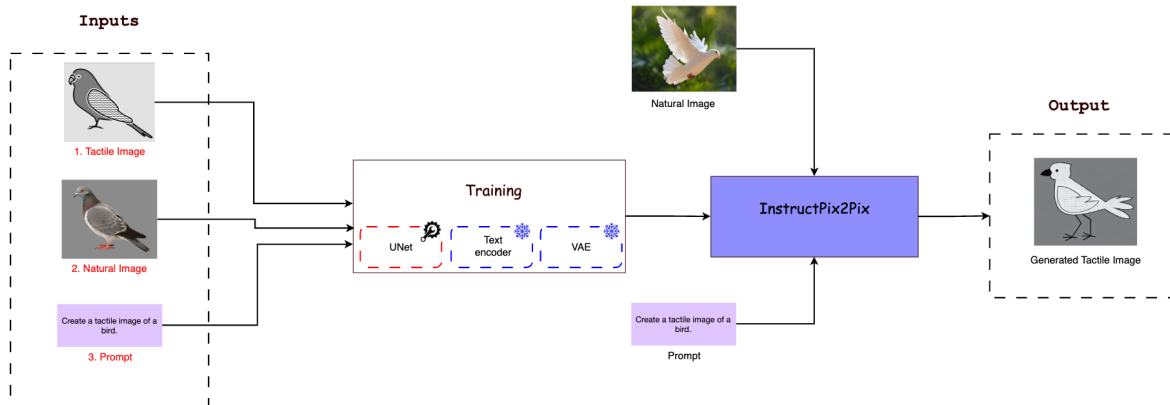


Figure 3.6: Project’s Fine-Tuned Baseline Flow Diagram.

3.6 Our Final Adapter

We create a new adapter by fine-tuning Stable Diffusion with LoRA applied to both the UNet and the CLIP text encoder, while keeping the VAE frozen (Fig. 3.7). Compared to the earlier `train.py` pipeline, most dataset handling and training loop structure remain unchanged, and the hyperparameters that previously yielded the best results are reused. The main architectural change is that the text encoder is now unfrozen and modified to accept LoRA layers, allowing it to be trained jointly with the UNet. Training is guided by a modular reward suite—any of the four available reward models (CLIP, LPIPS, SSIM, Pixel) can be used individually or in combination, with weights controlling their influence. Data augmentation is adapted for each modality: the tactile image receives different transformations than the natural image, with the latter including saturation and Gaussian blur to enhance robustness.

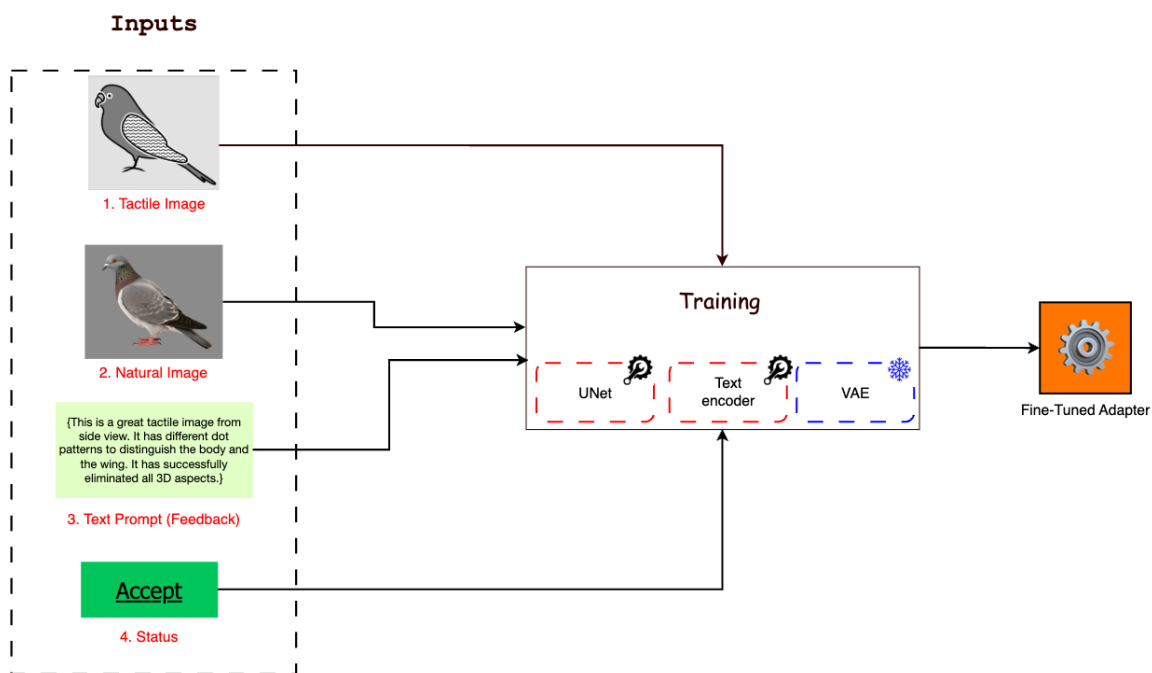


Figure 3.7: New Adapter Creation.

3.6.1 Inputs

- **Base model:** pre-trained Stable Diffusion (Img2Img) checkpoint.
- **Dataset:** The same JSONL file `training_dataset.jsonl` that was used in the initial training phase with records containing:
 - `input_image`: path to source (natural) image.
 - `generated_image`: path to target (tactile) image.
 - `feedback_text`: instruction/feedback string.
 - `rating`: `accept` / `reject`.
- **Training hyperparameters:** Best performing epochs, batch size, learning rate, scheduler from initial tests.
- **Reward configuration:** `REWARD_TYPE` (e.g., `clip+lpips`) and `REWARD_WEIGHTS`.

3.6.2 Outputs

- **LoRA adapters:**
 - `UNET/`: UNet LoRA weights.
 - `TEXT_ENCODER/`: CLIP text encoder LoRA weights.

- **Training log:** `results.jsonl` with *run name*, *loss type*, *reward type*, *epochs*, *final loss*, and *timestamps*.

3.6.3 Model Components and LoRA Setup

- **UNet:** prepared for LoRA (`prepare_model_for_kbit_training`), adapters attached to attention projections `to_q`, `to_k`, `to_v` across specified layers.
- **Text encoder (unfrozen):** LoRA adapters attached to `q_proj`, `k_proj`, `v_proj`, `out_proj`. The text encoder is trainable in this method.
- **VAE:** frozen (`requires_grad_(False)`).

```

1 pipe.unet = prepare_unet_for_lora(pipe.unet)
2
3 lora_config_unet = LoraConfig(
4     r=8, lora_alpha=32,
5     target_modules=["to_q", "to_k", "to_v"],
6     layers_to_transform=list(range(16)),
7     lora_dropout=0.1, bias="none", fan_in_fan_out=False
8 )
9 pipe.unet = get_peft_model(pipe.unet, lora_config_unet)
10
11 lora_config_te = LoraConfig(
12     r=8, lora_alpha=32,
13     target_modules=["q_proj", "k_proj", "v_proj", "out_proj"],
14     lora_dropout=0.1, bias="none", fan_in_fan_out=False
15 )
16 pipe.text_encoder = get_peft_model(pipe.text_encoder,
17     lora_config_te)
18 pipe.vae.requires_grad_(False)          # freeze VAE
19 pipe.unet.train(); pipe.text_encoder.train()
20
21 # training loop (rewards, latents, denoising, loss, optimizer)
22
23 adapter_root = output_dir / "lora_adapters"
24 (pipe.unet).save_pretrained(str(adapter_root / "unet"))
25 (pipe.text_encoder).save_pretrained(str(adapter_root / "
    text_encoder"))

```

Listing 3.11: LoRA on UNet and Text Encoder; VAE frozen and adapters saved

3.6.4 Training Procedure

The training process begins with loading the paired natural and tactile images, applying modality-specific augmentations. Light transformations are applied to the tactile images, while additional saturation and Gaussian blur are applied to the natural images. Each sample is then assigned a reward score $r_i \in [0, 1]$ from the configured reward suite (CLIP similarity, $1 - \text{LPIPS}$, SSIM, or pixel-based metrics). Both reference and target images are VAE-encoded into latents, with random noise added to the target latents according to a sampled diffusion timestep. Instructions are tokenized and processed by the trainable text encoder to produce conditioning embeddings, which guide the UNet in predicting the noise for denoising. The selected loss variant (e.g. reward-weighted MSE) uses r_i to scale the objective. Optimization updates only the LoRA parameters in the UNet and text encoder, while logging training metrics across epochs. Once training completes, the resulting LoRA adapters are saved for downstream evaluation and inference.

3.6.5 Testing Procedure

The testing phase aims to rigorously evaluate the performance, robustness, and generalization capabilities of the proposed feedback-guided LoRA adapter in comparison to the baseline Stable Diffusion model. Testing is conducted on a fixed, curated dataset, allowing reproducible comparison across different configurations. Two complementary components structure the evaluation: a *Custom Test Protocol*, designed specifically for tactile image generation tasks, and a *New Testing Setup* that enables systematic baseline–adapter comparisons under controlled conditions. Together, these components provide both targeted, qualitative insights and broad, quantitative performance metrics.

Custom Test Protocol

A tailored evaluation protocol was developed to rigorously assess the performance and generalization ability of the newly trained feedback-guided adapter for Stable Diffusion. Unlike generic benchmarking, this protocol was explicitly designed for the tactile image generation setting and integrates both seen and unseen examples, as well as controlled variations in prompts and reference images. The dataset contains 29 carefully curated entries spanning 14 classes, each including a natural image, its corresponding tactile image, a generation prompt, and a human rating. The protocol incorporates:

- **Held-out evaluation** for unseen examples with prompt adaptation to a baseline template.

- **Seen evaluation** using identical inputs from training to measure retention of learned mappings.
- **Targeted interpretation subset** selection for deeper qualitative inspection.
- **Prompt variation testing** to assess robustness to linguistic changes.
- **Reference variation testing** to evaluate sensitivity to visual changes.

This structured evaluation approach provides a reproducible framework for analyzing both fidelity and adaptability in feedback-driven tactile image generation systems.

New Testing Setup

The evaluation script is designed to systematically compare the baseline Stable Diffusion model with the newly trained LoRA adapter. It supports three test configurations, controlled by boolean flags:

- **baseline_only** – runs without any LoRA adapter.
- **adapter_only** – applies the LoRA adapter without merging weights.
- **both_together** – loads the LoRA adapter and merges its weights into both the UNet and the CLIP text encoder before testing.

This toggle system allows direct comparison of performance impacts from the adapter and the merge process. The merging step integrates LoRA-trained parameters into both the UNet and text encoder for potentially faster inference and unified weights during evaluation.

The setup of these configurations is defined as follows:

```
1 RUNS = {
2     "baseline_only": {"use_lora": False, "merge": False, "
3     lora_path": None},
4     "adapter_only": {"use_lora": True, "merge": False, "lora_path
5     ": LORA_PATH},
6     "both_together": {"use_lora": True, "merge": True, "lora_path
7     ": LORA_PATH},
8 }
```

Listing 3.12: Evaluation configuration toggles

For each configuration, the pipeline loads the appropriate weights, generates outputs for the fixed test set, and computes semantic alignment scores using a CLIP-based reward model. All images and scores are stored for later quantitative and qualitative comparison.

It loads the fixed test dataset that was mentioned earlier and generates tactile image outputs for each configuration using the same inference settings. A CLIP-based reward model then scores each output for semantic alignment with its prompt. Generated images and their scores are stored per configuration for later quantitative and qualitative analysis.

Chapter 4

Experimental Settings

This chapter outlines the experimental setup used to evaluate and refine the LoRA fine-tuned Stable Diffusion models for tactile graphic generation. The experiments were designed to explore how different training and testing parameters impact performance, with a focus on model quality, semantic alignment, and tactile readability.

4.1 Compute and Software Environment

All experiments were conducted on a local machine equipped with an **NVIDIA RTX 4090 GPU** with **24GB VRAM** and **CUDA version 12.2**. The system used **Python 3.10.16** and **PyTorch 2.7.0**.

The following major libraries were used:

- **Transformers** (v4.52.3): for CLIP and related model access
- **Diffusers** (v0.33.1): for Stable Diffusion and image generation workflows
- **Accelerate, PEFT, Bitsandbytes**: for efficient model loading and fine-tuning

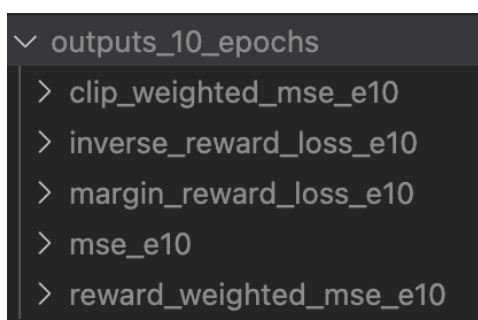
All libraries were managed via a **Conda** environment.

4.2 Training Experiments

Initial Adapter Training

The `train.py` script was executed a total of 15 times to systematically train models across different settings. Five distinct loss functions were tested:

- Mse Loss
- Clip Weighted Mse
- Reward Weighted Mse
- Inverse Reward Loss
- Margin Reward Loss



```
▼ outputs_10_epochs
  > clip_weighted_mse_e10
  > inverse_reward_loss_e10
  > margin_reward_loss_e10
  > mse_e10
  > reward_weighted_mse_e10
```

Figure 4.1: 10-epochs models generated by "`train.py`".

Each loss function was trained for three different epoch counts (5, 8, and 10) using a fixed batch size of 4. This setup enabled evaluation of underfitting versus overfitting, with metrics including:

- Final training loss
- CLIP reward score
- Visual fidelity of generated images

New `img2img` Baseline Training

A new `InstructPix2Pix` `img2img` baseline was trained from the instruction-tuned Stable Diffusion model using paired natural-tactile image data. Two main variants were explored:

- **10k steps** training

- **15k steps** training

This allowed us to study whether extended training improves fidelity and semantic alignment without overfitting.

Final Adapter Training

The `train2.py` script was used to fine-tune the LoRA adapter jointly on the UNet and CLIP text encoder (VAE frozen) with reward-weighted objectives. Different reward weight distributions between **CLIP** and **LPIPS** were tested to determine their effect on output quality.

4.3 Testing Experiments

Initial Adapter Testing

Once training was completed, the `test.py` script was used to evaluate all LoRA fine-tuned models. Testing was performed three times, once for each trained epoch count (5, 8, and 10), to measure how the number of epochs affected model performance.

The evaluation used a test dataset containing 60 entries: 30 images that had been part of the training dataset used by `train.py`, and 30 entirely unseen images that were not used during training. This split allowed assessment of both how well the models reproduced familiar examples and how well they generalized to new data.

Final Adapter & Baseline Testing

The `test2.py` script was used to evaluate:

1. **New Baseline Only:** `Img2Img` model without adapter.
2. **Adapter Only:** applying LoRA without merging weights.
3. **Merged Adapter + Baseline:** LoRA merged into both UNet and text encoder for faster inference.

The evaluation dataset contained 29 entries across 14 classes, each comprising a natural image, corresponding tactile image, generation prompt, and human rating. The protocol assessed performance on both seen and unseen examples, included a targeted high-quality subset for qualitative inspection, and introduced controlled variations in prompts and reference images to measure linguistic robustness and visual sensitivity.

4.4 Hyperparameters

To further explore the model behavior, additional testing runs varied key inference hyperparameters:

- **Guidance Scale:** Default value of 7.5 was adjusted to 5 and 10 (while keeping the other hyperparameters fixed) for the 10-epochs models.
- **Number of Inference Steps:** Default value of 30 was tested against 10 and 50 (while keeping the other hyperparameters fixed) for the 10-epochs models.

These controlled variations assessed how guidance strength and denoising depth influenced image quality, semantic alignment, and tactile readability.

Ablation Studies

Two ablation studies were conducted if using the best performing hyperparameters:

- **Reward Weight Impact:** varying CLIP vs LPIPS weights.
- **Extended Training:** 50 and 100 epochs at reduced learning rate to examine long-term training effects on stability, overfitting, and output quality.

Altogether, this experimental setup provides a comprehensive view of how training dynamics, loss design, and inference hyperparameters shape the effectiveness of Stable Diffusion for tactile graphic generation.

4.5 Model Selection

The baseline model used in initial experiments was **Stable Diffusion version 1.5**, a widely adopted general-purpose diffusion model. To better support instruction-based editing for tactile image generation, we replaced it with the **Instruction-Tuned Stable Diffusion model**. This version was fine-tuned on paired image–instruction datasets, enabling it to follow semantic editing prompts more reliably.

Chapter 5

Results and Evaluations

5.1 Overview of Evaluation

This chapter reports two evaluation phases. The first phase benchmarks five LoRA-fine-tuned models trained on the *original* Stable Diffusion v1.5 baseline using five loss functions (MSE, CLIP-Weighted MSE, Reward-Weighted MSE, Inverse Reward Loss, Margin Reward Loss). Models were trained for 5, 8, and 10 epochs and tested on a 60-entry set (30 seen, 30 unseen) to assess memorization vs. generalization. Performance was summarized with CLIP-based semantic alignment scores and qualitative inspection.

The second phase adopts a stronger *instruction-tuned* Stable Diffusion img2img baseline and evaluates both the baseline itself and the newly proposed LoRA adapter. A fixed test set of **29** entries spanning **14** classes is used; each entry includes a natural (reference) image, a corresponding tactile image, a generation prompt, and a human rating. This set also encodes controlled *prompt variations* (synonyms/feedback-guided refinements) and *reference-image variations* to probe robustness. Using consistent inference settings, we compare: (i) new baseline alone vs. new baseline + adapter, (ii) adapter variants trained with different reward weightings (CLIP/LPIPS and their combinations), and (iii) training durations (e.g., 10k vs. 15k img2img steps for the baseline, and extended runs such as 50/100 epochs for ablations). Results are reported with the same CLIP-based metric and targeted qualitative analysis.

5.2 Quantitative Results

To systematically evaluate model performance, quantitative experiments were carried out using the CLIP reward score as the primary metric. This score reflects how closely generated tactile graphics align with the provided textual prompts and serves as a consistent

benchmark across all tests. Results were aggregated for each LoRA-fine-tuned model under different training and testing conditions, including variations in epoch count, inference steps, and guidance scale. Together, these experiments provide a structured view of how different hyperparameters and loss functions influence output quality and alignment.

5.2.1 Effect of Number of Epochs

To evaluate the influence of the number of epochs on model performance, the "test.py" script was run on models trained with each loss function across three training durations: 5, 8, and 10 epochs. The resulting CLIP reward scores for each setting are summarized in Table 5.1.

Model	5 Epochs	8 Epochs	10 Epochs
Baseline	29.17	29.17	29.17
Margin Reward Loss	29.17	29.17	29.17
MSE	28.35	27.80	27.64
Clip-Weighted MSE	28.17	27.81	27.79
Reward-Weighted MSE	28.13	27.55	27.76
Inverse Reward Loss	22.64	20.93	21.57

Table 5.1: Average CLIP reward scores for all loss functions across 5, 8, and 10 epochs.

Observations and Trends

Several key trends emerge from comparing results across training durations:

- **Baseline and Margin Reward Loss remain stable.** The baseline model and the Margin Reward Loss model consistently achieved the same highest mean score (29.17) across all epochs, showing that fine-tuning with this loss did not harm semantic alignment.
- **Minimal improvement after 5 epochs.** For most loss functions, the 5-epoch models achieved slightly higher or similar scores compared to 8 and 10 epochs. For instance, MSE Loss started at **28.35** (5 epochs) but slightly declined to **27.64** (10 epochs).
- **Inverse Reward Loss consistently underperformed.** Across all epochs, Inverse Reward Loss produced the lowest scores, dropping from **22.64** (5 epochs) to **20.93** (8 epochs), showing that the approach negatively impacts alignment regardless of training duration.

- **Signs of overfitting beyond 5 epochs.** While 5-epoch results were generally strongest for MSE Loss, Clip-Weighted MSE, and Reward-Weighted MSE, slight declines at 8 and 10 epochs suggest that extending training did not lead to better generalization and may have caused mild overfitting.

Key Findings

- **5 epochs provided the best balance.** Shorter training preserved generalization and delivered the highest or near-highest CLIP scores for most loss functions.
- **Margin Reward Loss maintained strong performance across all epochs.** Unlike other losses, it neither improved nor degraded significantly, making it a stable choice for future work.
- **Longer training (8–10 epochs) did not yield major benefits.** Scores plateaued or declined slightly, particularly for MSE Loss and Reward-Weighted MSE.
- **Inverse Reward Loss should be avoided.** It was consistently the lowest-performing loss, regardless of training length.

5.2.2 Effect of Number of Inference Steps

To evaluate the influence of inference steps on model performance, the "test.py" script was run on the models trained for 10 epochs using three different inference step counts: 10, 30 (default), and 50. The CLIP reward scores for each setting are summarized in Table 5.2.

Model	10 Steps	30 Steps	50 Steps
Baseline	26.53	29.17	28.75
Margin Reward Loss	26.53	29.17	28.75
MSE	26.95	27.64	27.59
Clip-Weighted MSE	26.77	27.79	27.40
Reward-Weighted MSE	27.12	27.76	27.42
Inverse Reward Loss	21.80	21.57	21.54

Table 5.2: Average CLIP reward scores across different inference step counts for models trained for 10 epochs.

Observations

- **30 steps offered the best overall balance.** Models using the default 30-step setting consistently achieved the highest scores, particularly for the Baseline and Margin Reward Loss models (29.17).
- **Fewer steps (10) caused noticeable quality degradation.** All models saw a sharp drop when only 10 steps were used, with the Baseline and Margin Reward Loss scores falling from 29.17 (30 steps) to 26.53.
- **Increasing to 50 steps provided no major benefit.** Scores slightly declined or plateaued compared to the default 30-step setting, suggesting that additional steps did not translate into better alignment.
- **Inverse Reward Loss remained consistently low.** Regardless of the step count, Inverse Reward Loss produced the weakest results (around 21.5–21.8).

Key Takeaway

The default 30 inference steps strike the optimal balance between computational cost and image quality. Reducing the steps significantly degrades performance, while increasing beyond 30 yields no meaningful improvement.

5.2.3 Effect of Guidance Scale

The impact of guidance scale on model performance was also assessed using the 10-epochs models. The "test.py" script was run three times, each with a different guidance scale: 5, 7.5 (Default), and 10. The CLIP reward scores are shown in Table 5.3.

Model	Guidance 5	Guidance 7.5	Guidance 10
Baseline	29.10	29.17	29.13
Margin Reward Loss	29.10	29.17	29.13
MSE	27.08	27.64	27.86
Clip-Weighted MSE	26.99	27.79	27.73
Reward-Weighted MSE	27.05	27.76	27.69
Inverse Reward Loss	21.75	21.57	21.44

Table 5.3: Average CLIP reward scores for models trained for 10 epochs under different guidance scale values.

Observations

- **The baseline and margin loss models were largely unaffected.** Scores stayed almost flat across scales (29.1–29.2).
- **For most loss types, performance peaked around guidance 10.** MSE and Reward-Weighted MSE gained a small bump in score (27.86 and 27.69) at guidance 10.
- **Guidance 5 slightly reduced performance.** All models except Baseline and Margin Reward Loss saw marginal dips when guidance dropped to 5.
- **Inverse Reward Loss stayed weak across the board.** It consistently lagged, scoring around 21.4–21.8 regardless of scale.

Key Takeaway

Guidance scale adjustments had only subtle effects. The default 7.5 already provided balanced performance, while moving to 10 yielded minor gains for some loss functions (especially MSE), and reducing to 5 slightly degraded output alignment.

5.2.4 Different Baseline Models

Table 5.4 compares the performance of models trained with different loss functions using two baseline variants: Stable Diffusion v1.5 and the instruction-tuned SD model.

Table 5.4: CLIP scores across different loss functions using two baseline models.

Model	Stable Diffusion v1.5	Instruction-tuned SD
Baseline	29.17	30.27
Margin Reward Loss	29.17	30.27
CLIP Weighted MSE	27.79	30.29
Reward Weighted MSE	27.76	30.37
MSE	27.64	30.20
Inverse Reward Loss	21.57	30.87

Observations

- The instruction-tuned SD baseline consistently outperforms the Stable Diffusion v1.5 across all loss functions.

- The inverse reward loss shows the largest improvement when switching to the instruction-tuned SD model, jumping from 21.57 to 30.87.
- Loss functions that performed moderately on the original baseline (e.g., CLIP Weighted MSE, Reward Weighted MSE) achieve significantly better results on the instruction-tuned model.

Key Takeaway


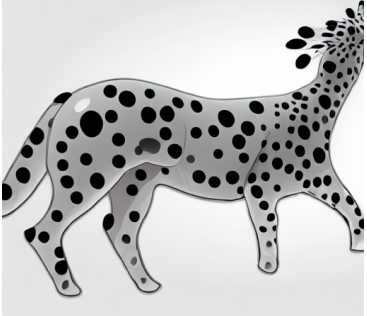
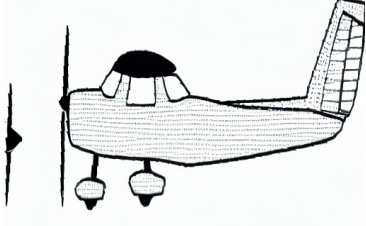

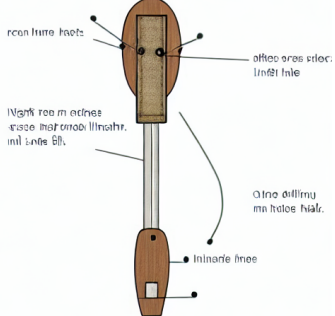

This indicates that a stronger base model allows the learning signal from reward-based losses to be better utilized.

5.3 Qualitative Observations

The qualitative results in Table 5.5 demonstrate a clear advantage of the Instruction-Tuned SD model over the standard Stable Diffusion v1.5. While the baseline model often introduces artifacts or semantically unrelated content, such as animal patterns or annotated diagrams, the instruction-tuned model generates outputs that better preserve the semantics and structural characteristics of the original images. This makes it a more appropriate foundation for tactile representation tasks.

Building on these observations, we now present results obtained with the new instruction-tuned img2img baseline under different training regimes, as well as its combination with the proposed LoRA adapter. These experiments focus on qualitative analysis, given the improved semantic alignment and structural fidelity of the new baseline.





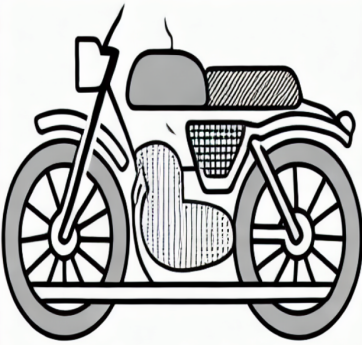

Table 5.5: Qualitative comparison between the original image, Stable Diffusion v1.5, and Instruction-Tuned SD model.

Natural Image	Stable Diffusion v1.5 Generated Image	Instruction-Tuned SD Generated Image
		
		

5.3.1 Effect of Training Steps for Instruct-Tuned Img2Img Baseline (10k vs 15k)

As shown in Table 5.6, the 15k-step baseline provides clearer semantic fidelity and finer structural details than the 10k-step counterpart. For example, in the *camel* class, the 15k model better preserves species-specific traits and local contours, whereas the 10k model exhibits simplified outlines. In the *motorcycle* class, the 10k model fails to capture key texture cues and the reference posture, while the 15k model more faithfully reflects both the surface patterning and overall geometry. Overall, increasing training to 15k steps yields qualitatively closer matches to the source images before adapter integration.

Table 5.6: Qualitative comparison between the original image, 15k steps, and 10k steps of Instruction-Tuned SD model.

Natural Image	Trained for 15k steps	Trained for 10k steps
		
		

5.3.2 Effect of Adding Final Adapter to Instruct-Tuned Img2Img Baseline

We evaluate the qualitative effect of adding the LoRA adapter to the instruction-tuned img2img baseline. Across both prompt modifications and changes in the reference image, the adapter improves semantic faithfulness.






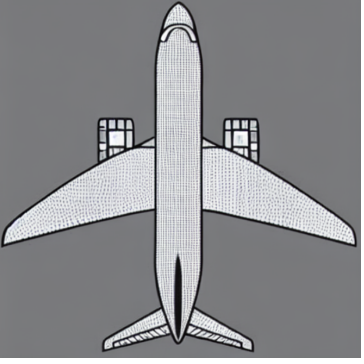
Case 1: Changing the prompt (Table 5.7)

- **Hut:** The adapter produces clearer separation between roof and walls, with thinner roof-texture strokes and a door rendered as an outline only, better matching the modified prompt.
- **Airplane:** The adapter sharpens the body/tail alignment, enforces single triangular wings, and renders engines as solid filled shapes, reducing spurious contours seen in the baseline.

Case 2: Changing the reference image (Table 5.8)

- **Hut:** When the reference image changes, the adapter preserves correct proportions and texture separation, producing more consistent roof patterns and door outlines despite variations in the source photo.
- **Watch:** The adapter corrects the distorted hand positions present in the baseline, while preserving the overall posture of the reference image and the casing shape. However, the generated image introduces some noise and extra details.

Table 5.7: Qualitative comparison between the original image, baseline only, and after adding the adapter based on modified prompts.

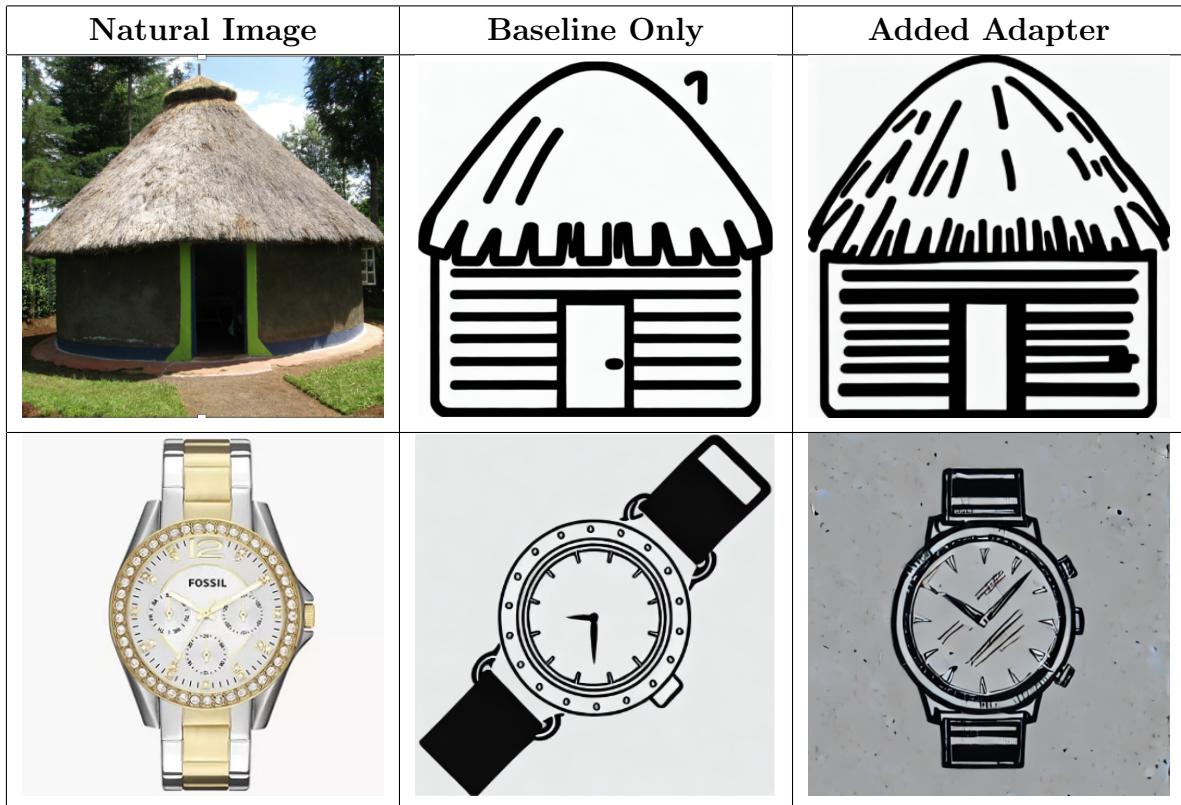
Natural Image	Baseline Only	Added Adapter
		
		

Modified Prompts:

Hut Class – Turn the natural image of a hut into a tactile-style drawing with a distinct base and top, using thinner lines for texture differentiation, and show the door only as an outline with no doorknob.

Airplane Class – Render edges and parts for tactile recognition of an airplane on a white background, with a distinct body and tail, triangular single main wings on each side, and solid black-filled engines.

Table 5.8: Qualitative comparison between the original image, baseline only, and after adding the adapter based on having different reference image.



5.3.3 Effect of Reward Weighting on Adapter Performance

We compared CLIP–LPIPS reward weighting schemes to determine which yields the best *tactile* output quality: crisp contours, correct part geometry, and high legibility. As shown in Fig. 5.1, the airplane class, the CLIP-heavy setting (0.8 CLIP) produces the *crispest and most legible* tactile rendering: clean silhouettes, reduced extraneous strokes, and better preservation of engine/body geometry. The balanced setting shows a slightly more accurate cockpit window but introduces extra line noise on the engines. Increasing LPIPS weight does not confer a clear advantage in our setting; it tends to thicken outlines, add spurious contours, and even drifts part shapes (e.g., engine mismatch). For the hut class, the balanced setting (0.6 CLIP / 0.4 LPIPS) yields the most faithful structure of the grass roof and respects the prompt requirement to render the door only as an outline with no doorknob. Higher LPIPS weight adds unnecessary thickness and visual clutter, while the CLIP-heavy setting simplifies details but slightly loses roof structure fidelity. We therefore adopt the 0.6 CLIP / 0.4 LPIPS weighting in our adapter.

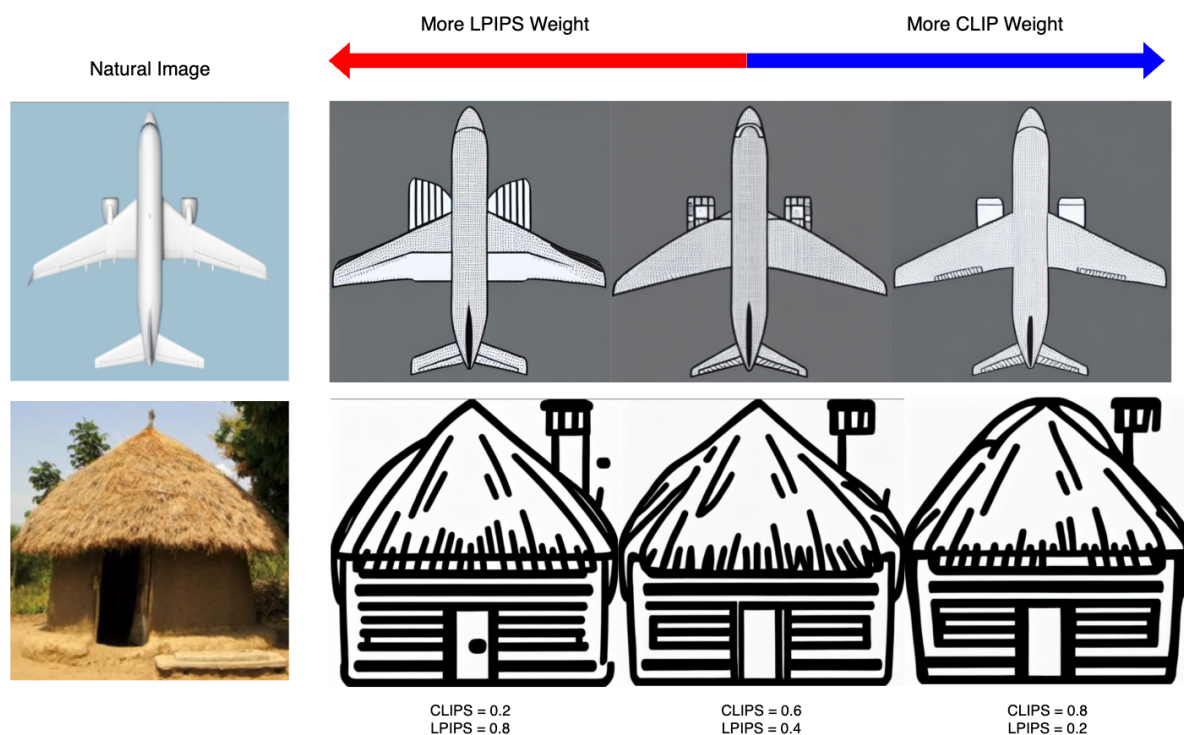


Figure 5.1: Qualitative comparison under different CLIP–LPIPS reward weightings.

5.3.4 Extended Training Ablation

We examined adapter training beyond the default schedule by continuing for 50 and 100 epochs at a reduced learning rate. As shown in Fig. 5.2, the 30-epoch model produces the cleanest, most stable outlines. Pushing to *50 epochs* slightly smooths interior dot textures but degrades global geometry (e.g., wing proportions and engine shapes become more blocky), and some contours soften. At *100 epochs* clear signs of overfitting appear: spurious structures emerge (e.g., a dorsal intake-like shape), outlines thicken, and part alignment drifts (wing and engine misplacement). Overall, extended training yielded no qualitative gains and instead increased artifacts and semantic drift.

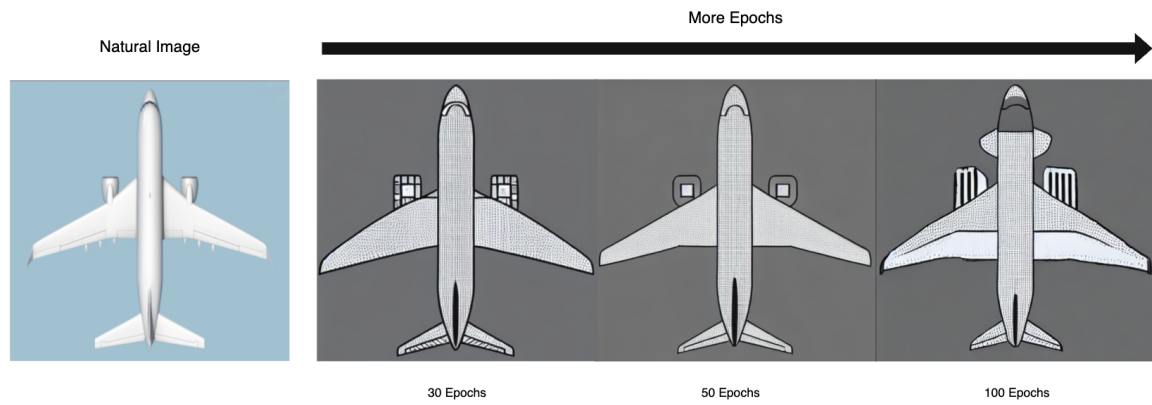


Figure 5.2: Effect of prolonged adapter training (30 vs. 50 vs. 100 epochs) on a fixed airplane example.

Chapter 6

Summary and Conclusion

This project advanced the TactileNet framework by addressing key limitations in dataset size, training strategy, and evaluation methodology. The dataset was expanded by a factor of ten and restructured with standardized feedback aligned to BANA guidelines, providing a more consistent and reliable foundation for model learning. A reward-guided training pipeline was implemented using Stable Diffusion fine-tuned with LoRA adapters, integrating multiple loss functions and forward/backward diffusion processes.

Extensive experiments compared training durations, loss types, and inference hyperparameters, revealing how these factors influence tactile image quality, CLIP reward scores, and structural readability. The instruction-tuned img2img Stable Diffusion baseline consistently outperformed the original Stable Diffusion v1.5 across most configurations, particularly when paired with reward-weighted objectives. Adapter integration improved semantic alignment and contour clarity in many cases. However, optimal results required careful balancing of reward weights (0.6 CLIP / 0.4 LPIPS) and restraint in training duration to avoid overfitting or semantic drift.

6.1 Limitations

A central limitation of this work is **data scarcity**. The available tactile graphics dataset remains small in both scale and class diversity, which restricts the model’s exposure to varied structural patterns and hinders its ability to generalize beyond the seen classes. This bottleneck cascades into multiple issues, such as overfitting to frequent patterns, poor handling of rare objects, and limited robustness to prompt or input variations. Another limitation lies in the heavy reliance on CLIP’s similarity score as a fixed, frozen proxy reward. This metric, computed as the cosine similarity between CLIP’s vision and text encoder embeddings, offers only a coarse measure of image–text alignment. It does

not evaluate tactile quality, structural clarity, or compliance with BANA guidelines, and it cannot adapt to the tactile graphics domain over time. Consequently, images that scored well under CLIP were sometimes unsuitable for tactile interpretation, revealing its role as a weak oracle rather than a reliable evaluator of tactile graphic quality.

6.2 Future Work

Potential future work includes expanding the dataset to increase coverage and class diversity, thereby improving model robustness in the tactile graphics domain. Preference-based reinforcement learning optimization techniques could be explored, with this dataset serving as an initial step to create paired examples for fine-tuning. Furthermore, off-the-shelf models such as ChatGPT [ChatGPT, 2025] and Google Gemini [Gemini, 2025] could be employed to automatically generate edited versions of images based on the feedback text in the dataset, enabling scalable data augmentation and more effective alignment with tactile design requirements.

Bibliography

American Printing House for the Blind (APH), *The Tactile Graphics Image Library*, 2024.
<https://imagelibrary.aph.org/>

Braille Authority of North America (BANA), *Braille Authority of North America*, 2025.
<https://www.brailleauthority.org/>

Bertolotto, S., Panisson, A., & Perotti, A. (2021). LLM-Generated Class Descriptions for Semantically Meaningful Image Classification. In *Proceedings of the 1st International Conference on Explainable AI Neural Symbolic Methods (EXPLAINS)* (pp. 50–61). Porto, Portugal.

Brooks, T., Holynski, A., & Efros, A. A. (2023). InstructPix2Pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 18392–18402).

Chen, Y., Lai, Y. K., & Liu, Y. J. (2018). CartoonGAN: Generative adversarial networks for photo cartoonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 9465–9474).

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving, *Fine-tuning language models from human preferences*, arXiv preprint arXiv:1909.08593, 2019.

Dot Pad: The First Smart Tactile Graphics Display. Retrieved August 2, 2025, from <https://www.visionaid.co.uk/dot-pad>

Dzhurynskyi, Y., Mayik, V., & Mayik, L. (2024). Enhancing accessibility: Automated tactile graphics generation for individuals with visual impairments. *Computation*, 12(12), 251.

Edman, P. (1992). *Tactile graphics*. American Foundation for the Blind.

Google, *Gemini*, 2025. <https://gemini.google.com/app>

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le, *Finetuned language models are zero-shot learners*, arXiv preprint arXiv:2109.01652, 2021.

Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner, *Vision-Language Models are Zero-Shot Reward Models for Reinforcement Learning*, arXiv preprint arXiv:2310.12921, 2024.

Khan, A., Choubineh, A., Shaaban, M.A., Akkasi, A., & Komeili, M. (2025). *TactileNet: Bridging the Accessibility Gap with AI-Generated Tactile Graphics for Individuals with Vision Impairment*.

Lindstrom, P., & Turk, G. (2000). Image-driven simplification. *ACM Transactions on Graphics (TOG)*, 19(3), 204–241.

OpenAI, *ChatGPT*, 2025. <https://chat.openai.com/>

Pakenaite, K., Kamperou, E., Proulx, M. J., Sharma, A., & Hall, P. (2024, February). Pic2Tac: Creating accessible tactile images using semantic information from photographs. In *Proceedings of the Eighteenth International Conference on Tangible, Embedded, and Embodied Interaction* (pp. 1–12).

Paul F. Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei, *Deep reinforcement learning from human preferences*, *Advances in Neural Information Processing Systems*, 30, 2017.

Paul, S. (2023). Instruction-tuning Stable Diffusion with InstructPix2Pix. *Hugging Face Blog*. Retrieved from <https://huggingface.co/blog/instruction-tuning-sd>.

Perkins School for the Blind, *Tactile Graphics Library*, 2025. <https://www.perkins.org/resource/tactile-graphics-library/>

Ricci, R., Bazi, Y., & Melgani, F. (2024). Machine-to-Machine Visual Dialoguing with ChatGPT for Enriched Textual Image Description. *Remote Sensing*, 16(3), 441. <https://doi.org/10.3390/rs16030441>.

- Rosenblum, L. P., & Herzberg, T. S. (2015). Braille and tactile graphics: Youths with visual impairments share their experiences. *Journal of Visual Impairment & Blindness*, 109(3), 173–184.
- Shagidanov, A., Poghosyan, H., Gong, X., Wang, Z., Navasardyan, S., & Shi, H. (2024, April). Grounded-Instruct-Pix2Pix: Improving instruction-based image editing with automatic target grounding. In *ICASSP 2024 – IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 6585–6589). IEEE.
- Tanaka, K., Fushimi, T., Tsutsui, A., & Ochiai, Y. (2023). Text to haptics: Method and case studies of designing tactile graphics for inclusive tactile picture books by digital fabrication and generative AI. In *ACM SIGGRAPH 2023 Labs* (pp. 1–2).
- Wang, X., & Gupta, A. (2016). Generative image modeling using style and structure adversarial networks. In *European Conference on Computer Vision* (pp. 318–335). Cham: Springer International Publishing.
- Wu, H., Yang, H., Chang, F., Zhu, D., & Liu, Z. (2025). AI-generated tactile graphics for visually impaired children: A usability study of a multimodal educational product. *International Journal of Human-Computer Studies*, 103525.
- Zhang, L., Rao, A., & Agrawala, M. (2023). Adding Conditional Control to Text-to-Image Diffusion Models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (pp. 3836–3847).
- Zeinullin, M. (2025). *Improving exploration of tactile graphics by visually impaired people: Theoretical advances and a novel mobile application* (Doctoral dissertation).
- Zhao, X., Yu, H., & Bian, H. (2023). Image to Image Translation Based on Differential Image Pix2Pix Model. *Computers, Materials & Continua*, 77(1).

Appendix A

Source Code

"train.py"

```
1 import torch
2 from torchvision import transforms
3 import json
4 from pathlib import Path
5 from typing import List, Dict
6 from PIL import Image
7 from datasets import Dataset, Features, Image as ImageFeature,
   Value
8 from transformers import CLIPProcessor, CLIPModel, get_scheduler
9 from diffusers import StableDiffusionImg2ImgPipeline
10 from peft import LoraConfig, prepare_model_for_kbit_training,
   get_peft_model
11 from torch.nn import functional as F
12 from torch.utils.data import DataLoader
13 import wandb
14 from tqdm import tqdm
15 import random
16 from torchvision.transforms import Compose, ToTensor, Normalize
17 import time
18 from datetime import datetime
19 import os
20
21 # train.py
22 # Config
23 MODEL_PATH = "runwayml/stable-diffusion-v1-5"
24 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
```

```
25 BATCH_SIZE = 4
26 IMAGE_SIZE = (512, 512)
27 NUM_EPOCHS = 10
28
29 wandb.init(project="tactile-dpo-training")
30
31 # -----
32 # 1. Load Dataset
33 # -----
34 def load_image_pairs(jsonl_path: Path) -> List[Dict]:
35     pairs = []
36     with open(jsonl_path) as f:
37         for line in f:
38             try:
39                 data = json.loads(line)
40                 rating = 1 if data["rating"].lower() == "accept"
41                     else 0
42                 input_img_path = Path(data["input_image"])
43                 gen_img_path = Path(data["generated_image"])
44
45                 if not input_img_path.exists() or not gen_img_path
46                     .exists():
47                     continue
48
49                 pairs.append({
50                     "input_image": Image.open(input_img_path).
51                         convert("RGB").resize(IMAGE_SIZE),
52                     "generated_image": Image.open(gen_img_path).
53                         convert("RGB").resize(IMAGE_SIZE),
54                     "instruction": data["feedback_text"] if data["
55                         feedback_text"].lower() != "nan" else "",
56                     "rating": rating
57                 })
58             except Exception as e:
59                 print(f"Error: {e}")
60     return pairs
61
62 dataset_features = Features({
63     "input_image": ImageFeature(),
64     "generated_image": ImageFeature(),
65     "instruction": Value("string"),
```

```
61     "rating": Value("int32")
62 })
63
64 vae_preprocess = Compose([
65     ToTensor(),
66     Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
67 ])
68 # -----
69 # 2. Reward Model
70 # -----
71 class CLIPRewardModel:
72     def __init__(self):
73         self.model = CLIPModel.from_pretrained("openai/clip-vit-
74             base-patch32").to(DEVICE)
75         self.processor = CLIPProcessor.from_pretrained("openai/
76             clip-vit-base-patch32")
77
78     def calculate_reward(self, images: List[Image.Image], texts:
79         List[str]) -> torch.Tensor:
80         inputs = self.processor(text=texts, images=images,
81             return_tensors="pt", padding=True, truncation=True).to(
82             DEVICE)
83         with torch.no_grad():
84             outputs = self.model(**inputs)
85             sims = outputs.logits_per_image
86             diagonal = sims.diagonal()
87             return diagonal
88
89 reward_model = CLIPRewardModel()
90
91 augment = transforms.Compose([
92     transforms.RandomHorizontalFlip(),
93     transforms.ColorJitter(0.1, 0.1, 0.1),
94 ])
95
96 def collate_fn(batch: List[Dict]) -> Dict:
97     augmented = []
98     for item in batch:
99         augmented.append(item)
100         augmented.append({
101             "input_image": augment(item["input_image"]),
```

```
97         "generated_image": augment(item["generated_image"]),
98         "instruction": item["instruction"],
99         "rating": item["rating"]
100     })
101
102     rewards = reward_model.calculate_reward(
103         [x["generated_image"] for x in augmented],
104         [x["instruction"] for x in augmented]
105     )
106
107     return {
108         "input_images": [x["input_image"] for x in augmented],
109         "generated_images": [x["generated_image"] for x in
110             augmented],
111         "instructions": [x["instruction"] for x in augmented],
112         "ratings": [x["rating"] for x in augmented],
113         "rewards": rewards.squeeze().detach().cpu()
114     }
115
116 # -----
117 # 3. Custom Training Loop
118 # -----
119 def prepare_unet_for_lora(unet):
120     return prepare_model_for_kbit_training(unet)
121
122 # -----
123 # Loss Variants
124 # -----
125 def mse_loss_fn(pred, target, **kwargs):
126     return F.mse_loss(pred, target)
127
128 def clip_weighted_mse_loss(pred, target, ratings, **kwargs):
129     accepted_mask = torch.tensor(ratings, device=DEVICE, dtype=
130         torch.bool)
131     loss_weights = torch.where(
132         accepted_mask,
133         torch.tensor(0.5, device=DEVICE),
134         torch.tensor(2.0, device=DEVICE)
135     )
136     base_loss = F.mse_loss(pred, target)
137     return base_loss * loss_weights.mean()
```

```
136
137
138 def reward_weighted_mse_loss(pred, target, rewards=None, **kwargs)
    :
139     rewards = rewards.to(device=pred.device, dtype=pred.dtype) #
        <- FIX: move to same device
140     base_loss = F.mse_loss(pred, target, reduction='none') #
        shape: [B, C, H, W]
141     weight = rewards.reshape(-1, 1, 1, 1) # broadcast shape
142     return (base_loss * weight).mean()
143
144 def inverse_reward_loss(pred, target, rewards=None, **kwargs):
145     rewards = rewards.to(device=pred.device, dtype=pred.dtype)
146     base_loss = F.mse_loss(pred, target, reduction='none')
147     inverse_weight = (1.0 - rewards).unsqueeze(1).unsqueeze(2).
        unsqueeze(3)
148     return (base_loss * inverse_weight).mean()
149
150 def margin_reward_loss(pred, target, rewards=None, **kwargs):
151     rewards = rewards.to(device=pred.device, dtype=pred.dtype)
152     base_loss = F.mse_loss(pred, target, reduction='none')
153     margin = 0.2
154     clipped = torch.clamp(margin - rewards.unsqueeze(1).unsqueeze
        (2).unsqueeze(3), min=0)
155     return (base_loss * clipped).mean()
156
157 # Update map:
158 loss_fn_map = {
159     "mse": mse_loss_fn,
160     "clip_weighted_mse": clip_weighted_mse_loss,
161     "reward_weighted_mse": reward_weighted_mse_loss,
162     "inverse_reward_loss": inverse_reward_loss,
163     "margin_reward_loss": margin_reward_loss
164 }
165
166 def custom_train_loop(pipe, dataset, reward_model, run_name,
    loss_type="clip_weighted_mse", num_epochs=NUM_EPOCHS):
167     dataloader = DataLoader(dataset, batch_size=BATCH_SIZE,
        shuffle=True, collate_fn=collate_fn)
168     optimizer = torch.optim.AdamW(pipe.unet.parameters(), lr=1e-5)
```

```
169     scheduler = get_scheduler("cosine", optimizer,
170                               num_warmup_steps=10, num_training_steps=len(dataloader) *
171                               num_epochs)
172     pipe.unet.train()
173
174     start_time = time.time() # Start timing
175     for epoch in range(num_epochs):
176         for batch in tqdm(dataloader, desc=f"Epoch {epoch+1}/{
177                               num_epochs}"):
178             input_images = [vae_preprocess(img).to(device=DEVICE,
179                                                     dtype=torch.float32) for img in batch["input_images
180                                                         "]]
181             input_images = torch.stack(input_images) # shape: [B,
182                                                         3, 512, 512]
183             with torch.no_grad():
184                 latents = pipe.vae.encode(input_images).
185                     latent_dist.sample()
186                 latents = latents * 0.18215
187             latents = latents.to(dtype=torch.float32)
188             timesteps = torch.randint(
189                 0,
190                 pipe.scheduler.config.num_train_timesteps,
191                 (latents.shape[0],),
192                 device=DEVICE,
193                 dtype=torch.long
194             )
195             noise = torch.randn_like(latents).to(dtype=torch.
196                 float32, device=DEVICE)
197             latents = latents.to(device=DEVICE, dtype=torch.
198                 float32)
199             noisy_latents = pipe.scheduler.add_noise(latents,
200                 noise, timesteps)
201             text_inputs = pipe.tokenizer(batch["instructions"],
202                 padding="max_length", max_length=pipe.tokenizer.
203                 model_max_length, truncation=True, return_tensors="
204                 pt").to(DEVICE)
205             with torch.no_grad():
206                 encoder_hidden_states = pipe.text_encoder(**
207                     text_inputs).last_hidden_state
208             encoder_hidden_states = encoder_hidden_states.to(
209                 device=DEVICE, dtype=torch.float32)
```

```
195         noise_pred = pipe.unet(
196             noisy_latents,
197             timesteps,
198             encoder_hidden_states=encoder_hidden_states
199         ).sample
200         loss_fn = loss_fn_map[loss_type]
201         loss_fn = loss_fn_map[loss_type]
202
203         # Dynamically check if loss function takes "rewards"
204         loss_kwargs = {}
205         if "rewards" in loss_fn.__code__.co_varnames:
206             loss_kwargs["rewards"] = batch["rewards"]
207         if "ratings" in loss_fn.__code__.co_varnames:
208             loss_kwargs["ratings"] = batch["ratings"]
209         final_loss = loss_fn(noise_pred, noise, **loss_kwargs)
210         final_loss.backward()
211         optimizer.step()
212         optimizer.zero_grad()
213         scheduler.step()
214         wandb.log({"loss": final_loss.item()})
215
216     end_time = time.time() # End timing
217     total_time = end_time - start_time
218     output_dir = Path(f"./outputs_10_epochs/{run_name}") # Changing
219         # the folder's name for each epoch set
220     output_dir.mkdir(parents=True, exist_ok=True)
221     adapter_path = output_dir / "./lora_adapters"
222     pipe.unet.save_pretrained(str(adapter_path))
223
224     # Save training metadata
225     result_log = {
226         "run_name": run_name,
227         "loss_type": loss_type,
228         "num_epochs": num_epochs,
229         "total_time_sec": round(total_time, 2),
230         "final_loss": final_loss.item(),
231         "timestamp": datetime.now().isoformat()
232     }
233
234     with open(output_dir / "results.jsonl", "a") as f:
235         f.write(json.dumps(result_log) + "\n")
```

```
235
236     print(f"Training completed. Adapter saved to {adapter_path}")
237     print(f"Log written to {output_dir / 'results.jsonl'}")
238
239     # -----
240     # 4. Main
241     # -----
242     def main():
243         run_name = input("Enter run name (no spaces): ").strip()
244         jsonl_path = Path("./training_dataset.jsonl")
245         pairs = load_image_pairs(jsonl_path)
246         if not pairs:
247             raise ValueError("No valid pairs found.")
248         dataset = Dataset.from_list(pairs, features=dataset_features)
249         pipe = StableDiffusionImg2ImgPipeline.from_pretrained(
250             MODEL_PATH,
251             torch_dtype=torch.float32
252         ).to(DEVICE)
253
254         pipe.unet = prepare_unet_for_lora(pipe.unet)
255         lora_config = LoraConfig(
256             r=8,
257             lora_alpha=32,
258             target_modules=["to_q", "to_k", "to_v"],
259             layers_to_transform=list(range(16)),
260             lora_dropout=0.1,
261             bias="none",
262             fan_in_fan_out=False
263         )
264         pipe.unet = get_peft_model(pipe.unet, lora_config)
265         pipe.unet = pipe.unet.to(device=DEVICE, dtype=torch.float32)
266         pipe.vae.requires_grad_(False)
267         pipe.text_encoder.requires_grad_(False)
268         loss_type = "margin_reward_loss" # change loss function here.
269         # options: "mse", "clip_weighted_mse", "reward_weighted_mse",
270             "inverse_reward_loss", "margin_reward_loss"
271         custom_train_loop(pipe, dataset, reward_model, run_name=
272             run_name, loss_type=loss_type)
273
274 if __name__ == "__main__":
275     main()
```

Listing A.1: train.py

```
"test.py"
```

```
1 import torch
2 from diffusers import StableDiffusionImg2ImgPipeline
3 from peft import PeftModel
4 from transformers import CLIPProcessor, CLIPModel
5 from PIL import Image
6 from pathlib import Path
7 import json
8 import csv
9 import os
10 from tqdm import tqdm
11
12 # test.py
13
14 # --- Config ---
15 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
16 IMAGE_SIZE = (512, 512)
17 SEED = 42
18
19 torch.manual_seed(SEED)
20
21 # --- Paths ---
22 MODELS = {
23     "baseline": None,
24     "mse": "/home/student/khan/youssif/reward_model/
25           outputs_10_epochs/mse_e10/lora_adapters",
26     "clip_weighted_mse": "/home/student/khan/youssif/reward_model/
27           outputs_10_epochs/clip_weighted_mse_e10/lora_adapters",
28     "reward_weighted_mse": "/home/student/khan/youssif/
29           reward_model/outputs_10_epochs/reward_weighted_mse_e10/
30           lora_adapters",
31     "inverse_reward_loss": "/home/student/khan/youssif/
32           reward_model/outputs_10_epochs/inverse_reward_loss_e10/
33           lora_adapters",
34     "margin_reward_loss_10": "/home/student/khan/youssif/
35           reward_model/outputs_10_epochs/margin_reward_loss_e10/
```

```
        lora_adapters"
29 }
30
31 MODEL_PATH = ""
32 TEST_JSONL = "./test_dataset.jsonl"
33 OUTPUT_DIR = Path("./misc_test")
34 OUTPUT_DIR.mkdir(parents=True, exist_ok=True)
35
36 # --- CLIP Reward Model ---
37 class CLIPRewardModel:
38     def __init__(self):
39         self.model = CLIPModel.from_pretrained("openai/clip-vit-
40             base-patch32").to(DEVICE)
41         self.processor = CLIPProcessor.from_pretrained("openai/
42             clip-vit-base-patch32")
43
44     def score(self, images, texts):
45         inputs = self.processor(text=texts, images=images,
46             return_tensors="pt", padding=True, truncation=True).to(
47             DEVICE)
48         with torch.no_grad():
49             outputs = self.model(**inputs)
50             return outputs.logits_per_image.squeeze().tolist()
51
52 # --- Load Test Data ---
53 def load_test_set(path):
54     samples = []
55     with open(path) as f:
56         for line in f:
57             data = json.loads(line)
58             samples.append({
59                 "input_image": data["input_image"],
60                 "prompt": data["feedback_text"],
61                 "category": data.get("Category", "unknown")
62             })
63     return samples
64
65 # --- Inference Function ---
66 def generate_image(pipe, image_path, prompt):
67     image = Image.open(image_path).convert("RGB").resize(
68         IMAGE_SIZE)
```

```
64     generator = torch.Generator(DEVICE).manual_seed(SEED)
65     with torch.no_grad():
66         result = pipe(prompt=prompt, image=image, strength=0.8,
67                       guidance_scale=7.5, num_inference_steps=30, generator=
68                           generator)
69     return result.images[0].convert("RGB")
70
71 # --- Main Evaluation ---
72 def main():
73     test_samples = load_test_set(TEST_JSONL)
74     scorer = CLIPRewardModel()
75
76     for name, lora_path in MODELS.items():
77         print(f"\nEvaluating model: {name}")
78         model_output_dir = OUTPUT_DIR / name
79         model_output_dir.mkdir(exist_ok=True)
80
81         pipe = StableDiffusionImg2ImgPipeline.from_pretrained(
82             MODEL_PATH, torch_dtype=torch.float32).to(DEVICE)
83         if lora_path:
84             pipe.unet = PeftModel.from_pretrained(pipe.unet,
85                                                  lora_path).to(DEVICE)
86         pipe.unet.eval()
87         pipe.vae.eval()
88         pipe.text_encoder.eval()
89
90         results = []
91
92         for idx, sample in enumerate(tqdm(test_samples)):
93             try:
94                 gen_img = generate_image(pipe, sample["input_image"]
95                                         ], sample["prompt"])
96                 score = scorer.score([gen_img], [sample["prompt"]
97                                         ]])
98                 out_path = model_output_dir / f"{idx:04d}.png"
99                 gen_img.save(out_path)
100                results.append({
101                    "idx": idx,
102                    "category": sample["category"],
103                    "prompt": sample["prompt"],
104                    "score": score,
```

```
99         "image_path": str(out_path)
100     })
101     except Exception as e:
102         print(f"Error on sample {idx}: {e}")
103
104     # Save CSV
105     csv_path = model_output_dir / "results.csv"
106     with open(csv_path, "w", newline="") as csvfile:
107         writer = csv.DictWriter(csvfile, fieldnames=results
108                                 [0].keys())
109         writer.writeheader()
110         writer.writerows(results)
111
112 if __name__ == "__main__":
113     main()
```

Listing A.2: test.py

"analyzer.py"

```
1 import pandas as pd
2 import os
3 from pathlib import Path
4
5 # analyzer.py
6 # Path to the directory where model outputs are stored
7 OUTPUT_DIR = Path("//home/student/khan/youssif/
8     instruct_reward_model/test_results/30_epochs_test") # change
9     the input for analysis
10
11 def load_all_results(output_dir):
12     records = []
13     for model_dir in output_dir.iterdir():
14         csv_path = model_dir / "results.csv"
15         if not csv_path.exists():
16             print(f"Skipping {model_dir.name}, no results.csv
17                 found.")
18             continue
19         df = pd.read_csv(csv_path)
20         df["model"] = model_dir.name
21         records.append(df)
```

```
19     return pd.concat(records, ignore_index=True) if records else
20         pd.DataFrame()
21
22 def main():
23     df_all = load_all_results(OUTPUT_DIR)
24
25     if df_all.empty:
26         print("No results found. Make sure the evaluation outputs
27             exist.")
28         return
29
30     # --- Summary Logs ---
31     print("=" * 60)
32     print("CLIP Score Interpretation")
33     print(" - Each score measures how well the generated tactile
34         image aligns with the instruction.")
35     print(" - Higher is better (stronger visual-text alignment).")
36     print(" - Scores in your setting typically range 20-35.")
37     print(" - This helps compare how well each model follows human
38         feedback.")
39     print("=" * 60)
40
41     # --- Mean Score per Model ---
42     mean_scores = df_all.groupby("model")["score"].mean().
43         sort_values(ascending=False)
44     max_score = df_all["score"].max()
45     min_score = df_all["score"].min()
46
47     print("\n=== Mean CLIP Score per Model (Higher is Better) ==="
48         )
49     print(mean_scores.round(2))
50     print(f"\n Score Range Across All Models: min={min_score:.2f},
51         max={max_score:.2f}")
52
53     # --- Category-wise Mean Scores ---
54     print("\n=== Category-wise Mean Scores per Model ===")
55     pivot = df_all.pivot_table(index="category", columns="model",
56         values="score", aggfunc="mean")
57     print(pivot.round(2))
58
59     # --- Optional Save ---
```

```
52     pivot.to_csv("/home/student/khan/youssif/instruct_reward_model
        /analysis_results/30_epochs_test_by_category.csv") #change
        the name of the file
53     mean_scores.to_csv("/home/student/khan/youssif/
        instruct_reward_model/analysis_results/30
        _epochs_test_overall.csv") #change the name of the file
54
55 if __name__ == "__main__":
56     main()
```

Listing A.3: analyzer.py

"train2.py"

```
1 import torch
2 from torchvision import transforms
3 import json
4 from pathlib import Path
5 from typing import List, Dict
6 from PIL import Image
7 from datasets import Dataset, Features, Image as ImageFeature,
    Value
8 from transformers import CLIPProcessor, CLIPModel, get_scheduler
9 from diffusers import StableDiffusionImg2ImgPipeline
10 from peft import LoraConfig, prepare_model_for_kbit_training,
    get_peft_model
11 from torch.nn import functional as F
12 from torch.utils.data import DataLoader
13 import wandb
14 from tqdm import tqdm
15 import random
16 from torchvision.transforms import Compose, ToTensor, Normalize
17 import time
18 from datetime import datetime
19 import os
20 import lpips
21 import piq
22 import itertools
23
24
25 # Config
```

```
26 MODEL_PATH = "/home/student/khan/youssif/pix2pix_baseline/15
    k_steps"
27 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
28 BATCH_SIZE = 4
29 IMAGE_SIZE = (512, 512)
30 NUM_EPOCHS = 30
31 LEARNING_RATE = 5e-5 #1e-5
32
33 # Reward setup
34 REWARD_TYPE = "clip+lpips"      # "clip", "lpips", "ssim", "pixel
    ", "clip+lpips", "clip+ssim", "clip+pixel"
35 REWARD_WEIGHTS = {
36     "clip": 0.6, # adjust this back to 0.6
37     "lpips": 0.4, #adjsut this back to 0.4
38     "ssim": 0.4,
39     "pixel": 0.4,
40 }
41
42 wandb.init(project="tactile-dpo-training")
43
44 # -----
45 # 1. Load Dataset
46 # -----
47 def load_image_pairs(jsonl_path: Path) -> List[Dict]:
48     pairs = []
49     with open(jsonl_path) as f:
50         for line in f:
51             try:
52                 data = json.loads(line)
53                 rating = 1 if data["rating"].lower() == "accept"
                    else 0
54                 input_img_path = Path(data["input_image"])
55                 gen_img_path = Path(data["generated_image"])
56
57                 if not input_img_path.exists() or not gen_img_path
                    .exists():
58                     continue
59
60                 pairs.append({
61                     "input_image": Image.open(input_img_path).
                        convert("RGB").resize(IMAGE_SIZE),
```

```
62         "generated_image": Image.open(gen_img_path).
           convert("RGB").resize(IMAGE_SIZE),
63         "instruction": data["feedback_text"] if data["
           feedback_text"].lower() != "nan" else "",
64         "rating": rating
65     })
66     except Exception as e:
67         print(f"Error: {e}")
68     return pairs
69
70 dataset_features = Features({
71     "input_image": ImageFeature(),
72     "generated_image": ImageFeature(),
73     "instruction": Value("string"),
74     "rating": Value("int32")
75 })
76
77 vae_preprocess = Compose([
78     ToTensor(),
79     Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5])
80 ])
81
82 # -----
83 # 2. Reward Model
84 # -----
85 class CLIPRewardModel:
86     def __init__(self):
87         self.model = CLIPModel.from_pretrained("openai/clip-vit-
           base-patch32").to(DEVICE)
88         self.processor = CLIPProcessor.from_pretrained("openai/
           clip-vit-base-patch32")
89
90     def calculate_reward(self, images: List[Image.Image], texts:
           List[str]) -> torch.Tensor:
91         inputs = self.processor(text=texts, images=images,
           return_tensors="pt", padding=True, truncation=True).to(
           DEVICE)
92         with torch.no_grad():
93             outputs = self.model(**inputs)
94             sims = outputs.logits_per_image
95             diagonal = sims.diagonal()
```

```
96         return diagonal
97
98 reward_model = CLIPRewardModel()
99
100 augment_ref = transforms.Compose([
101     transforms.RandomHorizontalFlip(),
102     transforms.ColorJitter(brightness=0.1, contrast=0.1,
103         saturation=0.1),
104     transforms.GaussianBlur(kernel_size=3),
105 ])
106
107 augment_tgt = transforms.Compose([
108     transforms.RandomHorizontalFlip(),
109     transforms.ColorJitter(brightness=0.1, contrast=0.1), # no
110         saturation jitter here
111 ])
112
113 def collate_fn(batch: List[Dict]) -> Dict:
114     augmented = []
115     for item in batch:
116         augmented.append(item)
117         augmented.append({
118             "input_image": augment_ref(item["input_image"]),
119             "generated_image": augment_tgt(item["generated_image"]
120                 ),
121             "instruction": item["instruction"],
122             "rating": item["rating"]
123         })
124
125     return {
126         "input_images": [x["input_image"] for x in augmented],
127         "generated_images": [x["generated_image"] for x in
128             augmented],
129         "instructions": [x["instruction"] for x in augmented],
130         "ratings": [x["rating"] for x in augmented],
131     }
132
133 class RewardSuite:
134     def __init__(self):
135         self.clip = CLIPRewardModel()
```

```
133     self.lpiips = lpips.LPIPS(net='vgg').to(DEVICE).eval()
134     self.ssim = piq.ssim
135
136     @staticmethod
137     def pil_to_tensor_01(img: Image.Image) -> torch.Tensor:
138         return transforms.functional.to_tensor(img).unsqueeze(0).
139             to(DEVICE)
140
141     @staticmethod
142     def pil_to_tensor_vae(imgs: List[Image.Image]) -> torch.Tensor
143         :
144         t = [vae_preprocess(im) for im in imgs]
145         return torch.stack(t).to(DEVICE)
146
147     def get_clip_reward(self, gen_imgs, texts) -> torch.Tensor:
148         return self.clip.calculate_reward(gen_imgs, texts)
149
150     def get_lpips_reward(self, inp_imgs, gen_imgs) -> torch.Tensor
151         :
152         with torch.no_grad():
153             d = []
154             for a, b in zip(inp_imgs, gen_imgs):
155                 A = transforms.functional.to_tensor(a).unsqueeze
156                     (0).to(DEVICE) # 0..1
157                 B = transforms.functional.to_tensor(b).unsqueeze
158                     (0).to(DEVICE)
159                 A = A * 2 - 1 # -> [-1,1]
160                 B = B * 2 - 1
161                 d.append(self.lpiips(A, B).squeeze())
162             d = torch.stack(d).to(DEVICE)
163             d = torch.clamp(d, 0.0, 1.0)
164             reward = 1.0 - d
165         return reward
166
167     def get_pixel_reward(self, inp_imgs, gen_imgs) -> torch.Tensor
168         :
169         with torch.no_grad():
170             A = self.pil_to_tensor_vae(inp_imgs)
171             B = self.pil_to_tensor_vae(gen_imgs)
```

```
167         mse = F.mse_loss(A, B, reduction='none').mean(dim
168             =(1,2,3))
169         mse = torch.clamp(mse, 0.0, 1.0)
170         reward = 1.0 - mse
171     return reward
172
173 def get_ssim_reward(self, inp_imgs, gen_imgs) -> torch.Tensor:
174     with torch.no_grad():
175         vals = []
176         for a, b in zip(inp_imgs, gen_imgs):
177             A = self.pil_to_tensor_01(a)
178             B = self.pil_to_tensor_01(b)
179             vals.append(self.ssim(A, B, data_range=1.0).
180                 squeeze())
181         return torch.stack(vals).to(DEVICE)
182
183 def get_saturation_reward(self, imgs) -> torch.Tensor:
184     # High when images are closer to grayscale
185     with torch.no_grad():
186         vals = []
187         for im in imgs:
188             t = transforms.functional.to_tensor(im).to(DEVICE)
189             # [3,H,W], 0..1
190             mx, mn = t.max(dim=0).values, t.min(dim=0).values
191             delta = mx - mn + 1e-6
192             v = mx + 1e-6
193             s = (delta / v).mean()           # mean HSV-like
194             saturation proxy
195             vals.append(1.0 - s)           # higher reward for
196             lower saturation
197         return torch.stack(vals).to(DEVICE)
198
199 def get_combined_reward(self, inp_imgs, gen_imgs, texts, kind:
200     str, weights: Dict[str, float]) -> torch.Tensor:
201     kind = kind.lower()
202     if "+" not in kind:
203         if kind == "clip": return self.get_clip_reward(
204             gen_imgs, texts)
205         if kind == "lpips": return self.get_lpips_reward(
206             inp_imgs, gen_imgs)
```

```
199         if kind == "pixel": return self.get_pixel_reward(  
200             inp_imgs, gen_imgs)  
201         if kind == "ssim": return self.get_ssim_reward(  
202             inp_imgs, gen_imgs)  
203         if kind == "sat": return self.get_saturation_reward(  
204             gen_imgs)  
205         raise ValueError(f"Unknown REWARD_TYPE: {kind}")  
206  
207     parts = [k.strip() for k in kind.split("+")]  
208     w_sum = sum(weights.get(p, 0.0) for p in parts)  
209     if w_sum <= 0:  
210         raise ValueError("Combined reward has zero total  
211             weight.")  
212     norm = {p: weights.get(p, 0.0)/w_sum for p in parts}  
213  
214     total = None  
215     for p in parts:  
216         if p == "clip": r = self.get_clip_reward(gen_imgs,  
217             texts)  
218         elif p == "lpips": r = self.get_lpips_reward(inp_imgs,  
219             gen_imgs)  
220         elif p == "pixel": r = self.get_pixel_reward(inp_imgs,  
221             gen_imgs)  
222         elif p == "ssim": r = self.get_ssim_reward(inp_imgs,  
223             gen_imgs)  
224         elif p == "sat": r = self.get_saturation_reward(  
225             gen_imgs)  
226         else: raise ValueError(f"Unknown combo part: {p}")  
227         total = r*norm[p] if total is None else total + r*norm  
228         [p]  
229     return total  
230  
231 reward_suite = RewardSuite()  
232  
233 # -----  
234 # 3. Custom Training Loop  
235 # -----  
236 def prepare_unet_for_lora(unet):  
237     return prepare_model_for_kbit_training(unet)  
238  
239 # -----
```

```
230 # Loss Variants
231 # -----
232 def mse_loss_fn(pred, target, **kwargs):
233     return F.mse_loss(pred, target)
234
235 def clip_weighted_mse_loss(pred, target, ratings, **kwargs):
236     accepted_mask = torch.tensor(ratings, device=DEVICE, dtype=
237         torch.bool)
238     loss_weights = torch.where(
239         accepted_mask,
240         torch.tensor(0.5, device=DEVICE),
241         torch.tensor(2.0, device=DEVICE)
242     )
243     base_loss = F.mse_loss(pred, target)
244     return base_loss * loss_weights.mean()
245
246 def reward_weighted_mse_loss(pred, target, rewards=None, **kwargs)
247 :
248     rewards = rewards.to(device=pred.device, dtype=pred.dtype)
249     base_loss = F.mse_loss(pred, target, reduction='none')
250     weight = rewards.reshape(-1, 1, 1, 1)
251     return (base_loss * weight).mean()
252
253 def inverse_reward_loss(pred, target, rewards=None, **kwargs):
254     rewards = rewards.to(device=pred.device, dtype=pred.dtype)
255     base_loss = F.mse_loss(pred, target, reduction='none')
256     inverse_weight = (1.0 - rewards).unsqueeze(1).unsqueeze(2).
257         unsqueeze(3)
258     return (base_loss * inverse_weight).mean()
259
260 def margin_reward_loss(pred, target, rewards=None, **kwargs):
261     rewards = rewards.to(device=pred.device, dtype=pred.dtype)
262     base_loss = F.mse_loss(pred, target, reduction='none')
263     margin = 0.2
264     clipped = torch.clamp(margin - rewards.unsqueeze(1).unsqueeze
265         (2).unsqueeze(3), min=0)
266     return (base_loss * clipped).mean()
267
268 # Update map:
269 loss_fn_map = {
```

```

267     "mse": mse_loss_fn,
268     "clip_weighted_mse": clip_weighted_mse_loss,
269     "reward_weighted_mse": reward_weighted_mse_loss,
270     "inverse_reward_loss": inverse_reward_loss,
271     "margin_reward_loss": margin_reward_loss
272 }
273
274 def custom_train_loop(pipe, dataset, run_name, loss_type="
clip_weighted_mse", num_epochs=NUM_EPOCHS):
275     dataloader = DataLoader(dataset, batch_size=BATCH_SIZE,
        shuffle=True, collate_fn=collate_fn)
276     trainable_params = itertools.chain(
277         (p for p in pipe.unet.parameters() if p.requires_grad),
278         (p for p in pipe.text_encoder.parameters() if p.
        requires_grad),
279     )
280     optimizer = torch.optim.AdamW(trainable_params, lr=
        LEARNING_RATE)
281
282     scheduler = get_scheduler("cosine", optimizer,
        num_warmup_steps=10, num_training_steps=len(dataloader) *
        num_epochs)
283     pipe.unet.train()
284     pipe.text_encoder.train()
285
286
287     start_time = time.time() # Start timing
288     for epoch in range(num_epochs):
289         for step, batch in enumerate(tqdm(dataloader, desc=f"Epoch
        {epoch+1}/{num_epochs}")):
290             input_imgs_pil = batch["input_images"]
291             gen_imgs_pil = batch["generated_images"]
292             instructions = batch["instructions"]
293             ratings = batch["ratings"]
294
295             rewards = reward_suite.get_combined_reward(
296                 input_imgs_pil, gen_imgs_pil, instructions,
297                 kind=REWARD_TYPE, weights=REWARD_WEIGHTS
298             )
299             assert rewards.shape[0] == len(gen_imgs_pil), f"
        rewards {rewards.shape} vs batch {len(gen_imgs_pil)}

```

```
    }"
300
301
302     # input_images = [vae_preprocess(img).to(device=DEVICE
      , dtype=torch.float32) for img in batch["
      input_images"]]
303     # input_images = torch.stack(input_images) # shape: [
      B, 3, 512, 512]
304
305     # VAE encode both; NOISE THE TARGET
306     ref_images = torch.stack([vae_preprocess(x) for x in
      input_imgs_pil]).to(DEVICE)
307     tgt_images = torch.stack([vae_preprocess(x) for x in
      gen_imgs_pil]).to(DEVICE)
308
309     with torch.no_grad():
310         ref_latents = pipe.vae.encode(ref_images).
          latent_dist.sample() * 0.18215
311         tgt_latents = pipe.vae.encode(tgt_images).
          latent_dist.sample() * 0.18215
312
313     # latents = latents.to(dtype=torch.float32)
314     timesteps = torch.randint(
315         0,
316         pipe.scheduler.config.num_train_timesteps,
317         (tgt_latents.shape[0],),
318         device=DEVICE,
319         dtype=torch.long
320     )
321     noise = torch.randn_like(tgt_latents).to(dtype=torch.
      float32, device=DEVICE)
322     noisy_tgt = pipe.scheduler.add_noise(tgt_latents,
      noise, timesteps)
323     # latents = latents.to(device=DEVICE, dtype=torch.
      float32)
324     # noisy_latents = pipe.scheduler.add_noise(latents,
      noise, timesteps)
325     text_inputs = pipe.tokenizer(batch["instructions"],
      padding="max_length", max_length=pipe.tokenizer.
      model_max_length, truncation=True, return_tensors="
      pt").to(DEVICE)
```

```
326         encoder_hidden_states = pipe.text_encoder(**
327             text_inputs).last_hidden_state
328         encoder_hidden_states = encoder_hidden_states.to(
329             device=DEVICE, dtype=torch.float32)
330
331     # --- choose correct UNet input (4 or 8 chans)
332     in_ch = pipe.unet.config.in_channels
333     if in_ch == 4:
334         latent_model_input = noisy_tgt
335     elif in_ch == 8:
336         latent_model_input = torch.cat([noisy_tgt,
337             ref_latents], dim=1)
338     else:
339         raise ValueError(f"Unexpected in_channels={in_ch}"
340             )
341
342     noise_pred = pipe.unet(
343         latent_model_input,
344         timesteps,
345         encoder_hidden_states=encoder_hidden_states
346     ).sample
347
348     # --- one-time debug (avoids log spam)
349     if epoch == 0 and step == 0:
350         tqdm.write(f"UNet in_channels: {in_ch}")
351         tqdm.write(f"latent_model_input: {tuple(
352             latent_model_input.shape)}")
353         tqdm.write(f"noise_pred/noise: {tuple(noise_pred.
354             shape)} {tuple(noise.shape)}")
355     assert noise_pred.shape == noise.shape, "UNet output
356         must match noise target."
357
358     loss_fn = loss_fn_map[loss_type]
359     # Dynamically check if loss function takes "rewards"
360     loss_kwargs = {}
361     if "rewards" in loss_fn.__code__.co_varnames:
362         loss_kwargs["rewards"] = rewards.detach().to(
363             noise_pred.dtype).to(noise_pred.device)
364
365     if "ratings" in loss_fn.__code__.co_varnames:
```

```
359         loss_kwargs["ratings"] = ratings
360
361         final_loss = loss_fn(noise_pred, noise, **loss_kwargs)
362         final_loss.backward()
363         optimizer.step()
364         optimizer.zero_grad()
365         scheduler.step()
366         wandb.log({"loss": final_loss.item()})
367
368     end_time = time.time() # End timing
369     total_time = end_time - start_time
370     output_dir = Path(f"/home/student/khan/youssif/
371         instruct_reward_model/train_results/{run_name}") #Changing
372         the folder's name for each epoch set
373     output_dir.mkdir(parents=True, exist_ok=True)
374
375     adapter_root = output_dir / "lora_adapters"
376     (adapter_root / "unet").mkdir(parents=True, exist_ok=True)
377     (adapter_root / "text_encoder").mkdir(parents=True, exist_ok=
378         True)
379
380     pipe.unet.save_pretrained(str(adapter_root / "unet"))
381     pipe.text_encoder.save_pretrained(str(adapter_root / "
382         text_encoder"))
383
384     # Save training metadata
385     result_log = {
386         "run_name": run_name,
387         "loss_type": loss_type,
388         "reward_type": REWARD_TYPE,
389         "num_epochs": num_epochs,
390         "total_time_sec": round(total_time, 2),
391         "final_loss": final_loss.item(),
392         "timestamp": datetime.now().isoformat()
393     }
394
395     with open(output_dir / "results.jsonl", "a") as f:
396         f.write(json.dumps(result_log) + "\n")
397
398     print(f"Training completed. Adapter saved to {adapter_root}")
```

```
396     print(f"Log written to {output_dir / 'results.jsonl'}")
397
398     # -----
399     # 4. Main
400     # -----
401     def main():
402         run_name = input("Enter run name (no spaces): ").strip()
403         jsonl_path = Path("./training_dataset.jsonl")
404         pairs = load_image_pairs(jsonl_path)
405         if not pairs:
406             raise ValueError("No valid pairs found.")
407         dataset = Dataset.from_list(pairs, features=dataset_features)
408         pipe = StableDiffusionImg2ImgPipeline.from_pretrained(
409             MODEL_PATH,
410             torch_dtype=torch.float32
411         ).to(DEVICE)
412
413         pipe.unet = prepare_unet_for_lora(pipe.unet)
414         lora_config = LoraConfig(
415             r=8,
416             lora_alpha=32,
417             target_modules=["to_q", "to_k", "to_v"],
418             layers_to_transform=list(range(16)),
419             lora_dropout=0.1,
420             bias="none",
421             fan_in_fan_out=False
422         )
423         pipe.unet = get_peft_model(pipe.unet, lora_config)
424         pipe.unet = pipe.unet.to(device=DEVICE, dtype=torch.float32)
425
426         lora_config_te = LoraConfig(
427             r=8,
428             lora_alpha=32,
429             target_modules=["q_proj", "k_proj", "v_proj", "out_proj"],
430                 # CLIPTextModel attn projections
431             lora_dropout=0.1,
432             bias="none",
433             fan_in_fan_out=False,
434         )
435         pipe.text_encoder = get_peft_model(pipe.text_encoder,
436             lora_config_te)
```

```
435     pipe.text_encoder = pipe.text_encoder.to(device=DEVICE, dtype=
         torch.float32)
436
437     pipe.vae.requires_grad_(False)
438     # pipe.text_encoder.requires_grad_(True)
439     loss_type = "reward_weighted_mse" # change loss function here
         .
440     # options: "mse", "clip_weighted_mse", "reward_weighted_mse",
         "inverse_reward_loss", "margin_reward_loss"
441     custom_train_loop(pipe, dataset, run_name=run_name, loss_type=
         loss_type)
442
443 if __name__ == "__main__":
444     main()
```

Listing A.4: "train2.py" for final adapter creation

"test2.py"

```
1 import torch
2 from diffusers import StableDiffusionInstructPix2PixPipeline
3 from peft import PeftModel
4 from transformers import CLIPProcessor, CLIPModel
5 from PIL import Image
6 from pathlib import Path
7 import json
8 import csv
9 import os
10 from tqdm import tqdm
11
12 # --- Config ---
13 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
14 IMAGE_SIZE = (512, 512)
15 SEED = 42
16
17 torch.manual_seed(SEED)
18
19 def load_lora_adapters(pipe, lora_root: str, merge: bool = True,
        device: str = "cuda"):
20     """
21     lora_root should contain:
```

```
22     - lora_root/unet
23     - lora_root/text_encoder
24     """
25     root = Path(lora_root)
26     unet_dir = root / "unet"
27     te_dir = root / "text_encoder"
28
29     if unet_dir.exists():
30         pipe.unet = PeftModel.from_pretrained(pipe.unet, str(
31             unet_dir)).to(device)
32     else:
33         print(f"[WARN] UNet adapter not found at: {UNET_DIR}")
34
35     if te_dir.exists():
36         pipe.text_encoder = PeftModel.from_pretrained(pipe.
37             text_encoder, str(te_dir)).to(device)
38     else:
39         print(f"[WARN] Text-encoder adapter not found at: {TE_DIR}")
40
41     if merge:
42         # Merge LoRA weights into the base modules for faster
43         # inference
44         if hasattr(pipe.unet, "merge_and_unload"):
45             pipe.unet = pipe.unet.merge_and_unload()
46         if hasattr(pipe.text_encoder, "merge_and_unload"):
47             pipe.text_encoder = pipe.text_encoder.merge_and_unload()
48
49     return pipe
50
51 LORA_PATH = "/home/student/khan/youssif/instruct_reward_model/
52     train_results/30_epochs_clip+lpips/lora_adapters"
53 # --- Paths ---
54 RUNS = {
55     "baseline_only": {"use_lora": False, "merge": False, "
56         lora_path": None},
57     "adapter_only": {"use_lora": True, "merge": False, "lora_path
58         ": LORA_PATH},
59     "both_together": {"use_lora": True, "merge": True, "lora_path
60         ": LORA_PATH},
```

```
54 }
55
56 MODEL_PATH = "/home/student/khan/youssif/pix2pix_baseline/15
    k_steps"
57 TEST_JSONL = "/home/student/khan/youssif/instruct_reward_model/
    test/test.jsonl"
58 OUTPUT_DIR = Path("/home/student/khan/youssif/
    instruct_reward_model/test_results/30_epochs_test/")
59 OUTPUT_DIR.mkdir(parents=True, exist_ok=True)
60
61 # --- CLIP Reward Model ---
62 class CLIPRewardModel:
63     def __init__(self):
64         self.model = CLIPModel.from_pretrained("openai/clip-vit-
            base-patch32").to(DEVICE)
65         self.processor = CLIPProcessor.from_pretrained("openai/
            clip-vit-base-patch32")
66
67     def score(self, images, texts):
68         inputs = self.processor(text=texts, images=images,
            return_tensors="pt", padding=True, truncation=True).to(
            DEVICE)
69         with torch.no_grad():
70             outputs = self.model(**inputs)
71         return outputs.logits_per_image.squeeze().tolist()
72
73 # --- Load Test Data ---
74 def load_test_set(path):
75     samples = []
76     with open(path) as f:
77         for line in f:
78             data = json.loads(line)
79             samples.append({
80                 "input_image": data["input_image"],
81                 "prompt": data["feedback_text"],
82                 "category": data.get("Category", "unknown")
83             })
84     return samples
85
86 # --- Inference Function ---
87 def generate_image(pipe, image_path, prompt):
```

```
88     image = Image.open(image_path).convert("RGB").resize(  
89         IMAGE_SIZE)  
90     generator = torch.Generator(DEVICE).manual_seed(SEED)  
91     with torch.no_grad():  
92         result = pipe(prompt=prompt, image=image, strength=0.8,  
93             guidance_scale=10, num_inference_steps=30, generator=  
94                 generator)  
95     return result.images[0].convert("RGB")  
96  
97 # --- Main Evaluation ---  
98 def main():  
99     test_samples = load_test_set(TEST_JSONL)  
100     scorer = CLIPRewardModel()  
101  
102     for name, cfg in RUNS.items():  
103         print(f"\nEvaluating: {name}")  
104         model_output_dir = OUTPUT_DIR / name  
105         model_output_dir.mkdir(exist_ok=True)  
106  
107         pipe = StableDiffusionInstructPix2PixPipeline.  
108             from_pretrained(  
109                 MODEL_PATH,  
110                 torch_dtype=torch.float32,  
111                 safety_checker=None  
112             ).to(DEVICE)  
113  
114         if cfg["use_lora"]:  
115             pipe = load_lora_adapters(pipe, cfg["lora_path"],  
116                 merge=cfg["merge"], device=DEVICE)  
117  
118         pipe.unet.eval()  
119         pipe.vae.eval()  
120         pipe.text_encoder.eval()  
121  
122         results = []  
123         for idx, sample in enumerate(tqdm(test_samples)):  
124             try:  
125                 gen_img = generate_image(pipe, sample["input_image"  
126                     ], sample["prompt"])  
127                 score = scorer.score([gen_img], [sample["prompt"  
128                     ]])
```

```
122         out_path = model_output_dir / f"{idx:04d}.png"
123         gen_img.save(out_path)
124         results.append({
125             "idx": idx,
126             "category": sample["category"],
127             "prompt": sample["prompt"],
128             "score": score,
129             "image_path": str(out_path)
130         })
131     except Exception as e:
132         print(f"Error on sample {idx}: {e}")
133
134     if results:
135         csv_path = model_output_dir / "results.csv"
136         with open(csv_path, "w", newline="") as csvfile:
137             writer = csv.DictWriter(csvfile, fieldnames=
138                                     results[0].keys())
139             writer.writeheader()
140             writer.writerows(results)
141     else:
142         print(f"No results to save for {name}")
143
144 if __name__ == "__main__":
145     main()
```

Listing A.5: "test2.py" for final adapter and instruct-tuned img2img baseline testing