
Synthaesthetic Art: Human-Machine Creative Collaboration

Justin Baird
Tesseract.art
justin@tesseract.art

Ivy Chen
Tesseract.art
ivy@tesseract.art

Jookyung Song
Seoul National University
chsjk9005@gmail.com

MooKyoung Kang
Xorbis
nostone@gmail.com

Richard Savery
University of Canberra
saveryrichard0@gmail.com

Abstract

1 *Synthaesthetic Art* is a human-machine collaborative framework combining live
2 perception, generative AI, and robotic execution. Presented as a live installation
3 at the **Super AI Conference 2025**, the system captured portraits of attendees,
4 transformed them into stylised caricatures using a diffusion model that was fine-
5 tuned on artist-in-residence Ivy Chen’s works, and rendered them on canvas with
6 acrylic paint via a robotic arm. Unlike conventional AI-driven art systems, our
7 Synthaesthetic Art system was designed to preserve and amplify artist authorship
8 Artist authorship is preserved through stylistic training, selective curation, and real-
9 time control. This work introduces the conceptual foundations of Synthaesthetic
10 Art and explores how hybrid creativity can expand agency, re-frame authorship,
11 and deepen emotional expression through machine augmentation.

12 1 Introduction

13 What happens when a robot doesn’t just assist the artist, but becomes part of the artwork?

14 *Synthaesthetic Art* is a contemporary movement that explores how generative algorithms, sensory
15 inputs, and robotic systems can be integrated into live, co-creative processes. Rather than replacing
16 the artist, Synthaesthetic Art positions AI and robotics as collaborators, amplifying human intent
17 while challenging traditional notions of authorship and aesthetic form. This paper is grounded in
18 a collaborative project first demonstrated at the **Super AI Conference 2025 in Singapore**, where
19 portraits of attendees were captured, stylised via a diffusion model fine-tuned on artist-in-residence
20 Ivy Chen’s caricatures, and painted by a robotic arm. The process introduced Ivy’s distinctive style
21 and showcased a new form of hybrid creativity in which human authorship and machine execution
22 are intertwined.

23 The movement is grounded in six core principles: **Machine-Augmented Creativity**, which em-
24 phasises the machine’s role as an extension of human expression; **Hybrid Intelligence**, wherein
25 decisions are shaped by both human intuition and computational inference; **Temporal and Spatial**
26 **Fluidity**, which embraces adaptive, real-time creation; **Sensory Fusion**, blending modalities like
27 sound, vision, and gesture into unified compositions; **Algorithmic Aesthetics**, exploring emergent
28 patterns and computational beauty; and **Decentralised Authorship**, which questions the notion of a
29 singular artistic origin.

30 This paper presents a system that embodies these principles, transforming facial images into acrylic
31 portraits while preserving and extending the artist’s authorship. We situate this work as both a

32 technical contribution and an artistic provocation, demonstrating how AI systems can be embedded
33 in creative workflows in ways that remain expressive, ethical, and deeply human.

34 2 Early Explorations in Sound and Vision

35 Our initial studies connected live music with robotic painting through a real-time audio-robot pipeline.
36 Max/MSP analysed features such as pitch and onset density, triggering brush gestures via OSC
37 messages to a Python-controlled robotic arm. Figure 1 illustrates examples from these formative
38 experiments, showing different ways we connected music and robotic drawing. This early setup raised
39 the central question of our practice: where does authorship lie when sound becomes image? (Savery
40 et al., 2022)

41 A key turning point came in a performance with vocalist *Nyali*, where we introduced a vision layer.
42 A camera tracked her facial position and micro-expressions, while audio features continued to shape
43 brushstroke timing and type. OSC sliders modulated speed, jitter, and stroke behavior. The result was
44 a responsive, multimodal system that performers described as “the painting singing back.”

45 This collaboration revealed two core lessons: (1) mapping embodied cues such as gaze enables clearer
46 performer agency; (2) maintaining human supervision keeps the robot an expressive instrument, not
47 an autonomous painter. These insights shaped our current Synthaesthetic pipeline. (Savery and Baird,
48 2024)

49 Earlier prototypes included violin-controlled sketches, a large-canvas system with automated brush
50 swapping, and group performances where musicians signed the resulting artwork—early attempts
51 that explored the boundaries of shared creative authorship.

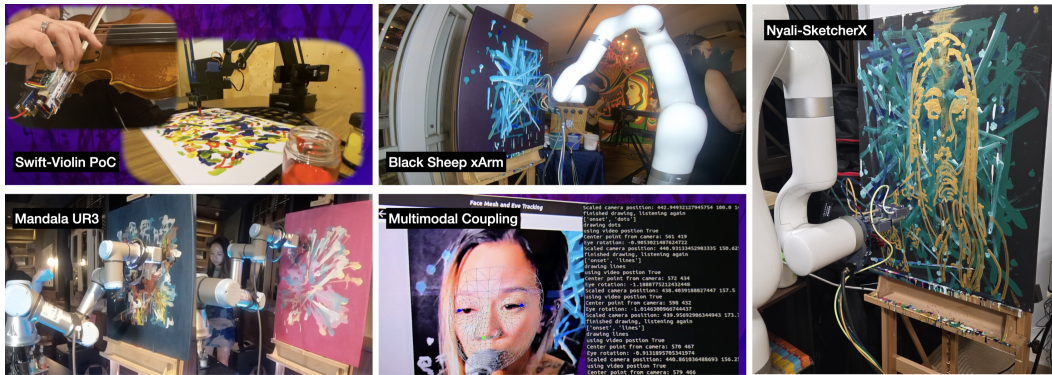


Figure 1: Examples of early experiments linking music and robotic drawing. Left and middle: music-driven robotic sketches. Right: performance with vocalist *Nyali*, where her singing generated the painted background and her live portrait was drawn at the end of the performance.

52 3 System Design and Technical Pipeline

53 3.1 Image Preprocessing and Prompt Generation

54 To produce stylistically consistent portraits, we developed a streamlined capture and preprocessing
55 pipeline aligned with our artist-adapted diffusion model. Portraits were taken using a fixed DSLR
56 setup, then automatically aligned and cropped to match the forward-facing composition typical of Ivy
57 Chen’s caricatures.

58 We applied Meta’s Segment Anything Model v2 (SAM2) (Ravi et al., 2024) to remove background
59 and clothing details, reducing stroke complexity and emphasising facial features for improved stylistic
60 fidelity and faster robotic execution.

61 Because our fine-tuned model builds on Stable Diffusion, high-quality text prompt generation is
62 essential. To automate this, we used Microsoft’s Florence-2 (Xiao et al., 2024) to extract visual
63 features—such as age, gender, hairstyle, and accessories—which were rephrased by a LLaMA-based

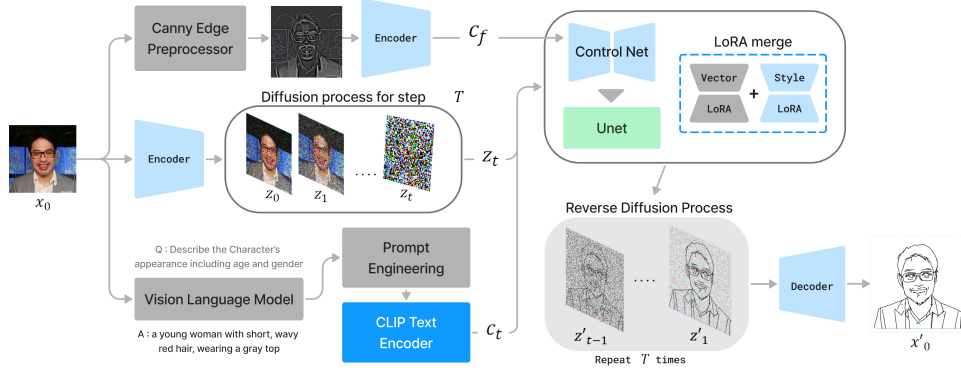


Figure 2: Technical pipeline showing portrait capture, image preprocessing feature extraction with canny edges, and diffusion-based stylisation.

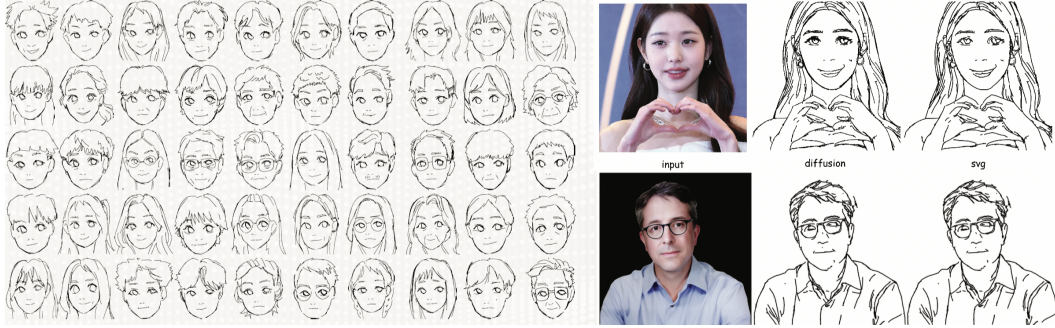


Figure 3: Left: hand-drawn caricatures by Ivy Chen used to train the model. Right: model inference pipeline – input photograph of a person (left), diffusion-generated stylised image (middle), and final vectorised SVG output (right) prepared for robotic execution.

language model into descriptive prompts. This combination enabled high-quality, personalised generation without manual intervention.

3.2 Diffusion Model and LoRA Training

To adapt Stable Diffusion to Ivy Chen’s caricature style, we fine-tuned it using two Low-Rank Adaptation (LoRA) modules: a *Style LoRA* trained on 50 of Ivy’s hand-drawn portraits (examples shown in Figure 3, left), and a *Vector LoRA* introduced in Song et al. (Song et al., 2024). It is optimized for clean, continuous strokes suited to robotic execution. This dual-LoRA approach preserved stylistic fidelity while ensuring vectorization compatibility.

The dataset included diverse demographics to support generalisation. Style images were paired with regularisation images to maintain structural coherence. To enhance alignment, we incorporated context consistency loss and ControlNet with canny-edge conditioning (Zhang et al., 2023).

Prompt embeddings used CLIP (Radford et al., 2021) to align textual inputs with generated features. The resulting model reliably produced stylised, vector-friendly caricatures, enabling efficient SVG conversion and robotic rendering via continuous Bézier strokes—crucial for fluid, human-like execution.

3.3 Vectorisation and Robotic Rendering

The selected diffusion output is converted to scalable vector paths (SVG). We approximate primitives (paths, lines, curves, polygons) as continuous Bézier segments and normalise them to the physical canvas.

83 Each path is then executed stroke-by-stroke by an xArm robot. We approach, draw, and retreat with
84 controlled lift motions to avoid smearing, and expose real-time parameters over OSC from Max/MSP.
85 These modifiers let sound or performance cues adjust stroke density and continuity without breaking
86 geometric fidelity.

87 Brush changes, pressure offsets, and per-stroke safety margins are handled in software, so the arm
88 behaves like a responsive prosthetic. The artist can pause, skip, or replay individual strokes, keeping
89 curatorial control over every mark. In practice, this makes the robot an expressive executor, not an
90 autonomous painter.

91 4 Artistic Agency and Direction

92 The creative foundation of the Synthaesthetic system is grounded in the work of Ivy Chen, the
93 artist-in-residence. Her hand-drawn caricatures formed the primary training corpus for the diffusion
94 model, and her aesthetic decisions shaped both the stylistic parameters of the generation process and
95 the behavior of the robotic rendering system.

96 To construct the dataset, Chen produced over fifty black-and-white line drawings, each reflecting
97 her distinct visual language. These caricatures encompassed a broad range of demographic fea-
98 tures—including variations in facial structure, hairstyle, age, and accessories—to promote stylistic
99 generalisation while maintaining authorial coherence. Technical constraints for robotic painting
100 informed the drawings: shading was omitted, filled regions avoided, and features such as eyes were
101 rendered with open lines or hatching. The dataset evolved through iterative cycles of testing with
102 close attention to line clarity, demographic balance, and stroke feasibility.

103 This preparation reflects the principle of **Algorithmic Aesthetics**, wherein visual form is shaped
104 through computational constraint and artistic abstraction, and **Machine-Augmented Creativity**, as
105 the artist’s own style becomes both the training signal and the creative boundary within which the
106 machine operates.

107 Following model fine-tuning, Chen retained curatorial authority over all system outputs. She con-
108 ducted stylistic evaluations of generated portraits, selecting those that aligned with her visual intent
109 and rejecting or refining those that did not. During robotic execution, she exercised real-time control
110 to modulate stroke attributes. In this manner, authorship was not merely pre-configured but remained
111 dynamic and embodied during the act of rendering.

112 This interactive control loop exemplifies **Hybrid Intelligence**, whereby the system’s generative and
113 mechanical components are continuously shaped by human intuition and aesthetic judgment. It also
114 enacts **Temporal and Spatial Fluidity**, enabling the artwork to evolve over time and across layers of
115 abstraction—from image to vector to brushstroke.

116 Moreover, the multimodal integration of visual, physical, and performance-based inputs reflects the
117 principle of **Sensory Fusion**, allowing audio features, gesture, and live input to co-shape the final
118 composition.

119 Chen’s role extended across the entire creative pipeline: dataset curation, style definition, generative
120 evaluation, and real-time robotic orchestration. This expansive authorship challenges conventional
121 notions of singular authorship in computational art and directly aligns with the principle of **De-**
122 **centralised Authorship**, wherein the final artefact emerges from a distributed network of human,
123 algorithmic, and mechanical contributors, with the artist acting as a conductor of this hybrid ensemble.

124 This framework defines the function of the **Synthaesthetic Artist**: one who authors the work, but
125 the generative conditions, the operational logic, and the expressive dynamics of the system itself.
126 Chen’s influence is embedded not only in line, form, and style, but in the computational grammar and
127 mechanical gestures that bring each portrait into being. Her authorship is both visible and systemic; a
128 signature of intent encoded across every layer of the creative process.

129 5 Foundations of Human-Machine Collaboration

130 The development of our real-time human-machine collaborative art system is fundamentally shaped
131 by the core principles of the Synthaesthetic Art movement. These principles—Machine-Augmented
132 Creativity, Hybrid Intelligence, Temporal and Spatial Fluidity, Sensory Fusion, Algorithmic Aes-



Figure 4: Ivy Chen holding a completed portrait painted, while in the same scene the robot continues to generate a new portrait at SuperAI conference in Singapore.

thetics, and Decentralised Authorship—provide both philosophical and technical scaffolding for our design. Below, we articulate how each principle is practically embodied within our system and workflow.

5.1 Machine-Augmented Creativity

Our system positions artificial intelligence not as a replacement for human creativity, but as a means of augmenting it. Ivy Chen, the artist-in-residence, retains creative control throughout the process—from dataset curation and style definition to output validation and robot supervision. The generative diffusion model, trained on her caricatures, functions as an extension of her visual language. The robot acts not as an autonomous painter but as a prosthetic executor of her artistic decisions. In this way, the machine serves to expand Ivy’s expressive reach, enabling high-speed, stylistically consistent production while remaining subordinate to her authorship.

5.2 Hybrid Intelligence

The system exemplifies a dynamic interplay between human intuition and machine computation. The generative pipeline requires both algorithmic inference (e.g. vision-language embeddings, prompt synthesis, vectorisation) and real-time human judgment. For example, while Florence-2 and LLaMA construct initial prompts from visual features, the final output is shaped through human-led curation, allowing Ivy to reject or refine undesirable results. Real-time robotic rendering is further modulated via OSC controls exposed in Max/MSP, enabling intuitive performance-like interaction during execution. This co-creative structure embodies the hybrid intelligence model at the heart of Synthaesthetic Art.

5.3 Temporal and Spatial Fluidity

Unlike fixed, linear creative pipelines, our system embraces fluidity in both time and space. Paintings evolve through iterative processes—feedback loops between generation, human evaluation, and robotic execution—rather than a one-directional flow. The use of real-time inputs such as audio features or facial tracking data allows the artwork to respond temporally to live stimuli. Spatial fluidity is achieved through vector-based rendering and parameterised robotic motion, enabling adaptive composition adjustments without rigid constraints. These qualities align with the Synthaesthetic ideal of living, adaptable artworks that exist across multiple temporal and representational states.

5.4 Sensory Fusion

Our early explorations integrated sound (pitch, onset density) and visual cues (facial expression, gaze) into a single synthaesthetic canvas, translating auditory and visual stimuli into gestural strokes. The

ongoing use of OSC signals to perturb robotic behaviors—altering speed, stippling, or pressure based on sound or movement—blurs traditional boundaries between modalities. This multimodal mapping reflects the movement’s aim to blend senses in ways inspired by neurological synesthesia, where inputs in one modality produce responses in another. Our system literalises this fusion, turning song into line, emotion into form, and presence into composition.

5.5 Algorithmic Aesthetics

The artwork generated by the system embraces complexity, imperfection, and the emergent properties of machine learning. The LoRA-trained diffusion model captures not just Ivy’s stylistic features but also introduces subtle variations—emergent quirks, abstraction, and character—that challenge normative standards of beauty. Rather than seeking photorealism or technical polish, we explore an algorithmically-mediated visual language that values stylisation, abstraction, and the uncanny. These properties resonate with the algorithmic aesthetic sensibility of the Synthaesthetic movement, where meaning and form are discovered through generative processes and computational play.

5.6 Decentralised Authorship

The system challenges the traditional notion of singular authorship. While Ivy’s role is primary, the system’s outcomes are shaped by contributions from model developers, performers (e.g. Nyali), technical collaborators, and even the viewers whose presence influences real-time execution. The curated dataset, open-source components (e.g. SAM2, ControlNet), and community-driven design also reflect an ethos of shared creative ownership. This decentralisation of authorship is a key tenet of Synthaesthetic Art, which embraces the plurality of creators, both human and non-human, in shaping the final work.

6 Conclusion and Future Directions

Synthaesthetic Art is not a single system or technique, but a broader conceptual movement, one that reimagines creativity as a distributed, co-authored process between humans and intelligent machines. It challenges the traditional hierarchies of artistic production by blending sensory modalities, decentralising authorship, and treating generative algorithms and robotic systems as active participants in the creative act.

The system we present in this paper is a concrete realisation of these ideals. It demonstrates how an artist’s intent can be embedded into a generative model, how human oversight can guide machine execution, and how live performance data can be translated into physical expression through code and motion. Yet this system is only one instantiation, only one articulation of a much larger set of questions posed by Synthaesthetic Art: What does it mean to co-create with a machine? How do we preserve artistic voice in computational pipelines? And how can technology be harnessed to deepen, rather than dilute, human expression?

Looking ahead, we aim to explore more diverse forms of real-time interaction including gesture, gaze, and audience input to broaden the system’s adaptability to multiple artists, styles, and modalities. As machine capabilities evolve, so too must our frameworks for collaboration, ethics, and creative agency. Synthaesthetic Art offers one such evolving framework: not to replace the artist, but to extend their reach and re-imagine what art can become when human and machine think, see, and feel together.

Acknowledgements

We thank Enya "Nyali" Lim, Eivind Lodemel and Tristan Seow for their live music performance. We thank the Super AI Singapore team for supporting the system demonstration.

References

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmlR, 2021.

- 211 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr,
212 Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. arXiv
213 preprint arXiv:2408.00714, 2024.
- 214 Richard Savery and Justin Baird. Redefining artistic boundaries: A real-time interactive painting robot for
215 musicians. In 38th Conference on Neural Information Processing Systems (NeurIPS 2024) Creative AI
216 Track, 2024.
- 217 Richard Savery, Anna Savery, and Justin Baird. Robotic arm generative painting through real-time analysis
218 of music performance. In Proceedings of the 10th International Conference on Human-Agent Interaction,
219 pages 253–255, 2022.
- 220 Jookyung Song, Mookyoung Kang, and Noju Kwak. Sketcherx: Ai-driven interactive robotic drawing220 with
221 diffusion model and vectorization techniques. In 38th Conference on Neural Information Processing Systems
222 (NeurIPS 2024) Creative AI Track, 2024.
- 223 Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan.
224 Florence-2: Advancing a unified representation for a variety of vision tasks. In Proceedings of the IEEE/CVF
225 Conference on Computer Vision and Pattern Recognition, pages 4818–4829, 2024.
- 226 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models.
227 In Proceedings of the IEEE/CVF international conference on computer vision, pages 3836–3847, 2023.