

Collaborative Comic Generation: Integrating Visual Narrative Theories with AI Models for Enhanced Creativity

Yi-Chun Chen, Arnav Jhala

North Carolina State University, Computer Science, Raleigh, NC 27606, USA

Abstract

This study presents a theory-inspired visual narrative generative system that integrates conceptual principles—comic authoring idioms—with generative and language models to enhance the comic creation process. Our system combines human creativity with AI models to support parts of the generative process, providing a collaborative platform for creating comic content. These comic-authoring idioms, derived from prior human-created image sequences, serve as guidelines for crafting and refining storytelling. The system translates these principles into system layers that facilitate the creation of comics through sequential decision-making, addressing narrative elements such as panel composition, story tension changes, and panel transitions. Key contributions include the integration of machine learning models into the human-AI cooperative comic generation process, the deployment of abstract narrative theories into AI-driven comic creation, and a customizable tool for narrative-driven image sequences. This approach improves narrative elements in generated image sequences and brings engagement of human creativity in an AI-generative process of comics. We open-source the code at https://github.com/Rimichen/Collaborative_Comic_Generation.

Keywords

Generative AI System, Human-AI interaction, Comic Generation, Visual Narrative Theories

1. Introduction

Despite their various names, comics, manga, and visual stories represent a dominant form of storytelling that spans cultures and age groups. Authors combine their creative storytelling ideas with textual expressions and graphical representations to convey intricate narratives through multi-modal panel sequences. Recently, generative AI, a trending topic in AI creativity, has explored the automatic generation of narratives and their synthesis with visual representations, simulating an activity traditionally rooted in human creativity. These studies lead to an exciting research question: how can AI collaborate with human creativity to create image sequences, such as comics and visual stories?

Nowadays, most generative AI models handle almost the entire process of creating image sequences, leaving little room for human authors to modify the details during generation. Traditionally, however, creating visual stories or comics relies heavily on authors' familiarity

CREAI 2024: International Workshop on Artificial Intelligence and Creativity, ECAI, 2024, Santiago de Compostela, Spain

✉ ychen74@ncsu.edu (Y. Chen); ahjhala@ncsu.edu (A. Jhala)

🌐 <https://sites.google.com/ncsu.edu/rimichen/home> (Y. Chen); <https://www.csc.ncsu.edu/people/ahjhala> (A. Jhala)

🆔 0009-0003-4035-9894 (Y. Chen)

© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



with narrative idioms [1], principles of visual storytelling [2], and their skills in translating narrative content into visual representations [3].

We propose a solution that balances these approaches, allowing AI models to assist in generating visual stories while also providing space for human creativity. Narrative theories and idioms are important tools for authors, helping them convey stories with clarity and making visual narratives more engaging for readers. These elements, derived from existing comics, serve as underlying schemes. We aim to leverage these narrative elements by integrating them into the generative process, enhancing the results. Current studies on generating image sequences or visual storytelling include text-to-image synthesis [4], character consistency algorithms [5], and narrative structuring models [6, 7]. Although most synthesis methods allow users to provide some input as seeds, such as captions or base images, authors have limited flexibility in the creative process. Furthermore, while narrative idioms like grammar and patterns of visual storytelling are widely discussed in the analysis of image sequences, they are rarely incorporated into AI-driven comic generation. Even in studies that integrated rule-based methodologies to allow narrative theories to influence visual storytelling [8, 9, 10, 11], authors' creativity remained largely excluded from the generative process.

Building on prior studies, we propose a comic-generating system that integrates a human-in-the-loop process, enabling authors to customize the generation by leveraging multiple AI models. Our system combines AI-driven narrative development and art generation, reducing the effort needed to create dynamic visuals. A sentiment analysis model also guides the story arc, serving as a plotline framework.

Second, our system dynamically applies narrative idioms during the generation process. The layer-wised customizations allow authors to edit the image sequence iteratively, selecting different rules to refine the results. This human-in-the-loop process ensures that authors can adjust based on the generated image sequences. Third, rather than treating the comic panel as a single entity, we decompose it into multiple layers: background, foreground, compositional layer, and symbol layer. This approach enables authors to customize specific details without altering the entire panel.

Finally, our system provides a Graphical User Interface (GUI) and an Application Programming Interface (API). These interfaces allow authors to extend and customize the comic generation process, creating engaging and varied content. Additionally, the system supports testing various machine learning models in comic generation. This approach balances AI and human creativity, making the comic creation process more flexible and efficient.

2. Methodology

The system is structured to a human-in-loop workflow as Figure 1.

The system comprises six modules: the graphical interface, a container of models, the pool of visual sets and data, an image sequence model, a generator, and a renderer.

- **Graphical Interface:** This is the primary means of interaction for human authors. It allows them to input base images, launch and apply ML models to the current image sequence, import scripts for customized editing layers, and perform simple operations like selection and dragging to modify images.

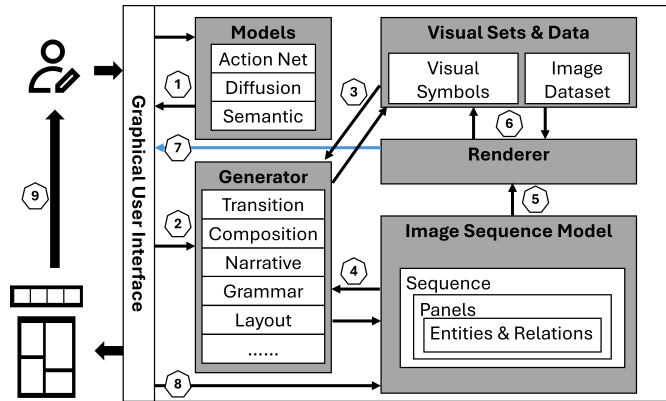


Figure 1: Overall system architecture and the author-in-loop workflow.

- **Container of Models:** This module houses various ML models that can be applied during generation.
- **Pool of Visual Sets and Data:** This repository contains visual elements and data sets the system can draw upon to create the image sequences.
- **Image Sequence Model:** This model organizes the graph model for linking the elements of the image sequence, including sequence information, panel transitions, panel content, characters, and narrative elements.
- **Generator:** This component integrates the customized editing layers into a pipeline that iteratively edits the image sequence by applying the narrative goals of each layer.
- **Renderer:** This module finalizes the visual representation of the generated sequences, ensuring they are ready for presentation or further editing.

Due to the complexity of applying ML models, some customizations of the generating process rely on scripts that deploy API functions across different system modules. The process typically follows nine steps between these modules:

1. **Register Models:** After receiving user scripts as input, the *Container of Models* registers all imported ML models and initializes them.
2. **Register Editing Layers:** This step links the interface buttons with class scripts, dynamically imports user-defined classes, and executes them upon button click.
3. **Retrieve Data:** If user-customized scripts include functions to retrieve information, the *Generator* communicates with the *Pool of Visual Sets and Data* to obtain necessary data.
4. **Apply Editing:** The customized editing results reflect the corresponding changes to the data nodes in the *Image Sequence Model*.
5. **Compose Visual Panels:** After updating the *Image Sequence Model* with the desired changes, the *Renderer* is triggered to start composing the visual results of comic panels. The updated graph model will also be passed to the *Renderer*.
6. **Map Semantic to Visual:** The *Renderer* maps the data nodes to visual elements in the *Pool of Visual Sets and Data* and forms the multi-layered panels. These layers include background, foreground, composition, symbol, and any other user-customized layers.









Anger	Quick moving	Slow moving	Anxious	Collision	Relieved	Shock	Big shock
							

Table 1

Examples of metaphor symbols of action and characters' emotions from Manga109. AisazuNihalarrenai© Yoshimasako, AkkeraKanjinchou© Kobayashiyuki, Akuhamu© Araisatoshi

7. **Render:** The multi-layered image panel is passed to the interface to support user selection and drag functions.
8. **Modify Results:** The interface inputs user interactions to update the *Image Sequence Model* and repeats the rendering steps.
9. **Get Result:** The generated and modified image sequence is presented on the interface for users to view.

2.1. Build-in and Extendable Elements:

To demonstrate the system's capabilities, we incorporated a built-in model, the action causal graph, and a visual symbol set as the default sources to support the plotline of the visual representations. These components are extendable through the APIs of the *Container of Models* and the *Pool of Visual Sets and Data*, respectively. Detailed documentation will be provided in the API subsections.

2.1.1. Common Symbols

In human-created comics, authors often use abstract symbols to visualize ideas such as atmosphere, motion, and emotions, thereby exaggerating characters' reactions. These symbols include emojis to emphasize emotions, speed lines to show movement, explosion shapes to represent collisions, cross shapes to denote anger, and many others. Table 1 presents examples from the Manga109 dataset, a collection of 109 Japanese manga titles published in commercial magazines [12, 13]. These examples serve as references for the symbols or emojis used in our default set.

2.1.2. Action Causal Network

The *Action Causal Graph* is a directed graph model where each node represents an action that a character might perform, and the links indicate the causal relationships between actions and possible reactions. For example, "Fall" links to "Fly," "Jump," and "Run," while "Dizzy," "Collide," and "Hit" link to "Run." Consecutive nodes form action pools for plot planning. The current version of the default action set built into the system includes a small group of common daily

actions. The detailed method for expanding the graph and its scalability will be discussed in subsequent subsections.

2.2. AI-Assisted Editing Layers and Narrative Theories:

We implement three ML-driven editing layers in the generator to demonstrate the system’s image sequence pipeline. The first layer uses a stable diffusion model for modifying visual elements, while the second layer combines the PAD emotion model with a semantic analysis language model to guide plotline decisions.

2.2.1. Diffusion Model for Visual Elements

We employed a pre-trained stable diffusion model developed by the CompVis group to implement one of the editing layers [14, 15]. This high-resolution image synthesis model transforms input text descriptions into detailed and coherent images based on latent diffusion. The model allows users to adjust visual elements like characters or scenes.

The generated panel sequences are rendered through our pipeline, where information is first updated in the Image Sequence Model before being rendered. Any changes in the visual representation of an entity will be reflected in its semantic nodes, ensuring character consistency throughout the panel sequence. For example, if the diffusion model alters the visual representation of character_x, this change will be mirrored in its semantic mapping. Similar rules apply to other visual elements. Furthermore, our system supports multi-layer rendering, dividing panel images into background, foreground, compositional, and symbol layers. This feature enables partial redrawing of the generated panel, ensuring that changes applied to one entity do not interfere with others.

2.2.2. Plotline and Narrative Theories

To form a simple plotline, we incorporate narrative idioms and theories. This editing layer aims to generate content for the image sequence according to a specific narrative arc. We use Cohn’s narrative grammar [16, 17, 18] to estimate the narrative arc, as the grammar categories indicate plot changes throughout the comic sequence. After establishing the narrative arc, we target the characters’ actions to predict story tension.

We use the arousal level concept from the PAD emotion model [19, 20, 21] and sentiment labels from the language model to predict arousal scores for character actions, where higher scores indicate greater story tension. Differences in arousal scores between consecutive actions (based on our action causal graph) form a probability distribution for subsequent actions. This mapping enables the editing layer to select actions probabilistically while maintaining narrative arc alignment. Detailed explanations of each component are provided in the following subsections.

Narrative Grammar

We formalize the plot generation of new comic sequences using Cohn’s Visual Narrative Grammar (VNG). Starting with the narrative structure to determine the content’s global reason-






E	I	L	P	R
				

Table 2

Examples of a comic sequence that followed grammar categories.

ing, we adopt Cohn’s theory, which proposes that coherent image sequences follow a grammar, organizing their global structure into five categories.

- **Establisher(E)**: Sets the objects and scenes without involving any action.
- **Initial(I)**: Marks the beginning of a story arc—the starting point of a sequence of actions or events.
- **Prolongation(L)**: Represents the middle state of the story arc, extending an action.
- **Peak(P)**: Indicates the highest story tension—the climax of an action.
- **Release(R)**: Releases the tension—the outcome or result of an action.

The five categories form basic phases through linear ordering:

Phase (Establisher) - Initial(Prolongation) - Peak - (Release)

The use of parentheses indicates that categories are optional when forming phases. The categories have different levels of importance, ranked from highest to lowest as follows: Peak, Initial, Release, Establisher, and Prolongation. Additionally, more complex combinatorial structures can be created through the conjunction of embeddings. Our editing layer generates the narrative structure using center-embedding, expanding a new tree structure by replacing a single category with a phase.

Table 2 provides an example of comic sequences (illustrated with simplified icons) following this grammar structure. In the same scene, the circle character appears in the first panel and then begins performing actions in the subsequent panels. In the fourth panel, a significant event occurs, creating a small climax, which leads to the resolution in the final panel. This sequence follows the narrative structure *E-I-L-P-R*.

The implementation follows the algorithm flow below:

The input is a comic sequence, either a generated result or an empty sequence with a certain length. Then, the editing layer creates an object dictionary for grammar phases and then expands the tree structure. The structure then decides the narrative arc of the comic sequence. It assigned a grammar phase to a sequence, and we then formed the possible narrative arc based on the structure.

Narrative Arc

The importance ranking of narrative grammar categories mirrors the narrative arc, reflecting a story’s progression from a calm beginning through a peak of tension in the middle to conflict

Algorithm 1 Grammar Structure

```
1: procedure GRAMMER STRUCTURE SCRIPT
2:    $P \leftarrow$  input the current panel sequence.
3:    $VNG \leftarrow$  Create object list for VNG basic phases.
4:    $S \leftarrow$  Expend a center-embedded tree with  $VNG$ , get the reference narrative structure.
5:   if then  $Length(P) \neq Length(S)$ 
6:     Add or Subtract empty panels from  $P$ 
7:   end if
8:   loop: Assign each phases in  $VNG$  to  $P$ 
9:   return  $P$ 
10: end procedure
```

resolution at the end. The editing layer projects these categories onto a curve that illustrates changes in story tension.

Algorithm 2 Narrative Arc

```
1: procedure MAPPING NARRATIVE ARC WITH NARRATIVE STRUCTURE
2:    $P \leftarrow$  input the current panel sequence.
3:    $REF \leftarrow$  create value dictionary with VNG phases with assigned tension scores.
4:    $REF = \{E : 0I : 2L : 4P : 6R : 2\}$ 
5:   if panels in  $P$  have assigned grammar phase then
6:     loop: over  $P$ , apply  $REF$  with the grammar phase
7:   elseuse default narrative arc scores
8:   end if
9:   return  $P$ 
10: end procedure
```

To capture the abstract concept of tension, we assigned each grammar phase a score between one and ten based on its narrative function. For example, the Peak(P) phase has the highest tension, while the Release(R) phase somewhat alleviates the tension. We begin by creating a value dictionary for each grammar phase. The next step is to map these phases to their respective scores and generate the curve of the narrative arc.

PAD Emotion State Model and Semantic Analysis

The PAD emotion state model, developed by Albert Mehrabian and James A. Russell, describes emotions through three dimensions: Pleasure, Arousal, and Dominance. This model quantifies emotional states, making it valuable in psychology, user experience design, and AI. Specifically, the Arousal dimension measures how energized or calm one feels, reflecting the activation level of an emotion. We adapt this concept to model narrative momentum—story tension.

Considering the image sequence-generating system’s future expansion and possible integration with a more extended narrative, we employed a sentiment analysis model for sentences—Roberta base model for emotion classification, fine-tuned on the GoEmotions dataset [22]. We

then project the sentiment labels in the RoBERTa model to the PAD model’s emotion labels, dividing the emotions into high, medium, and low arousal levels. By computing the feature vectors through the BERT model [23] of the two label sets, we can measure the Euclidian distance between the labels, estimating the possible arousal level scores for the sentiment labels. The mapping process follows the flow below:

Algorithm 3 Mapping Sentiment Labels with Emotion Labels

- 1: **procedure** ESTIMATING AROUSAL LEVELS
 - 2: $E \leftarrow$ emotion labels from the PAD emotion state model, where [high, medium, low] arousal level maps to [1, 0, -1].
 - 3: $S \leftarrow$ sentiment class labels from the RoBERTa model.
 - 4: $Sdis \leftarrow$ Distance matrix for distance between labels in S and labels in E .
 - 5: loop: over S , compute the distance to each element in E , and get $Sdis$
 - 6: loop: over each row in $Sdis$, flat the vector except the minimum two elements in the row.
 - 7: loop: over each row in $Sdis$, use $Sdis/sum(Sdis)$ as the weight multiply with E to get $S_arousal$.
 - 8: normalize $S_arousal$ to $[-1, 1]$
 - 9: **return** $S_arousal$
 - 10: **end procedure**
-

Figure 2 shows the results, where the blue points represent the emotion labels from the PAD model, and the orange points represent the sentiment classes in the language model.

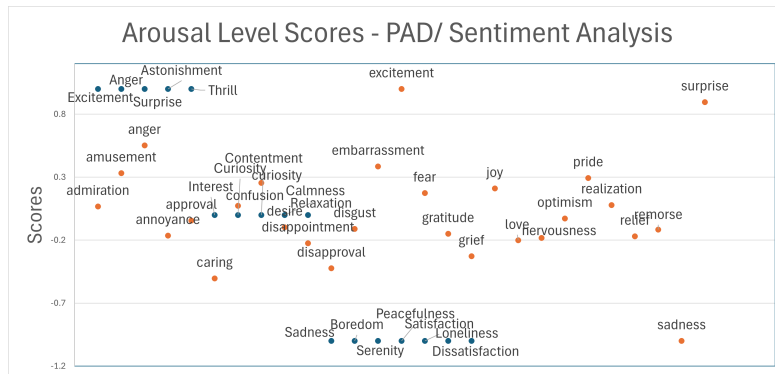


Figure 2: The arousal level scores are estimated using label set mapping.

The sentiment analysis model predicts probabilities across various sentiment classes when using a sentence or word as input. We use these probabilities as weights and then measure the distance between the input action and the sentiment classes, predicting the possible score of the input action. This score indicates the story’s tension. We further compute the slope between two consecutive actions and divide it by the sum of slopes for all subsequent actions. Finally, we normalize this value to determine the probability of the following action.

Action Mapping

By mapping arousal scores with actions, the editing layer integrates the narrative arc and actions by referencing these scores. Using the curve from the narrative arc, the layer calculates changes along the curve and selects actions that best fit these changes. Additionally, it sets a likelihood tolerance with the probability of actions, expanding potential narrative diversity. The process is described below:

Algorithm 4 Narrative Arc Mapping

```
1: procedure MAPPING ACTIONS TO FIT NARRATIVE ARC
2:    $P \leftarrow$  input the current panel sequence.
3:    $NET \leftarrow$  get the action causal graph network.
4:    $ACT \leftarrow$  creates a value dictionary for actions according to the arousal scores.
5:    $ARC \leftarrow$  get referenced Narrative Arc.
6:   loop: over  $P$ , check the characters' actions in the next panel and compute the score
   difference, according to  $ARC$ 
7:   Select actions in likelihood from  $ACT$  and  $NET$ 
8:   loop: revise action selection other panels according to  $NET$ 
9:   return  $P$ 
10: end procedure
```

2.2.3. Panel Relations

Panel transitions refer to the changes in content between consecutive panels. McCloud proposed six categories of transition types to model various aspects of content change[2], while Cohn introduced conjunction schemes to capture more complex panel transitions [16]. Our system combines these theories to modify the visual composition of comic panels—how elements are arranged—and reflect content changes according to the transition types. The narrative goal of this editing layer is to arrange the panel transitions in the generated comic sequence to create dynamic viewport changes and increase tension.

Here is how we map the transitions with panel content changes:

- **Action:** McCloud's action-to-action transition indicates changes in actions between consecutive panels. We use this transition to guide the selection of different character actions in the next panel.
- **Scene:** McCloud's scene-to-scene transition indicates changes in scenes between consecutive panels. We use this transition to guide the change of location where the character's action occurs in the next panel.
- **Object:** The object-to-object transition indicates a shift in focus from one object to another. We use this transition to guide the focus to different objects.
- **Addition:** Cohn's additive conjunction involves panels that add information or detail to the ongoing story. We use this transition to introduce new objects into the panels.
- **Alternation:** Cohn's alternating conjunction presents alternative scenarios or actions, offering possibilities within the narrative. Compared to scene transitions, we use this transition to guide panels to alter most elements while maintaining consistent characters.

In our system’s current version, some complex transitions or conjunctions are not yet supported but can be added through customizations. For example, the Temporal Conjunction depicts the progression of time, a rather abstract concept for visual representation. Another example is the Contrasting Conjunction, which depicts opposing ideas, actions, and emotions; achieving this requires deep semantic analysis of the narrative. However, we have integrated causal conjunctions to some extent by using action causal graphs.

3. Usage Explanations and Showcases

This section introduces our system’s application programming interface (API), and graphical user interface (GUI). The API functions are in Table 3. The system was constructed using several core classes, with the *Parameter*, *Layer*, and *AttributeNode* classes being the most crucial. The *Parameter* class manages all the registers of models and scripts. The *AttributeNode* class serves as the fundamental component of the graph model, representing the entire sequence. The *Layer* class is the parent class for all editing layer scripts. To execute modifications, the *apply* methods within the *Layer* class must be overridden and then executed by the Generator module.

The GUI screenshot is Figure 3. The two columns on the left display the image inputs for the character and scene of the comic sequence. Users can generate these images using the diffusion model or import them from their work. The column on the right side contains buttons that link to the imported scripts, triggering functions to apply modifications and generate results. The large area in the middle shows the currently generated result and allows users to select comic panels and the elements within them.

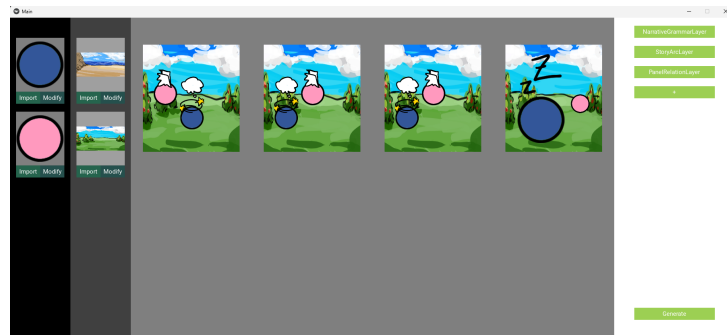


Figure 3: The graphical user interface of the generating system.

3.1. Showcases

Table 4 compares two sets of generated results, one before and one after applying the changes implemented in the current version. The first set, generated with user-imported images and all editing layers turned off, features content with default panel composition and randomly chosen characters’ actions. After partial redrawing with the stable diffusion model, the second set demonstrates the results, followed by an editing layer that uses the sentiment language model to achieve narrative planning. It includes a partially redrawn icon and scene using a diffusion model from an Einstein head icon and a WindowsXP desktop photo as inputs. In addition, based

Function Name	Parameters	Description	Return Type
Parameter()	Basic parameter for GUI settings <ul style="list-style-type: none"> win_w, win_h menu_w sequence_len 	The class that registers and manages all the modules.	None
addModel	Selected ML model class <ul style="list-style-type: none"> module_name 	Register the ML models for use in the generating process.	None
importModels()	None	Import all the registered models.	None
Models()	None	Initialize the Container of Models and initialize all the registered models.	None
addModule()	Customized editing layer class <ul style="list-style-type: none"> layer_name 	Register script for a new editing layer.	None
importModules()	None	Import all the registered editing layers.	None
Sequence()	Initialized Parameter object and attribute node type <ul style="list-style-type: none"> attribute_type Parameter() 	Inherit from the AttributeNode class and initialize the root of the graph model for the Image Sequence Model.	None
Generator()	None	Initialize the Generator Module.	None
addVisuals()	Name of target visual set, file path, or folder path for image input <ul style="list-style-type: none"> set_name path 	Expand or create the visual set with the assigned name.	None
GUIInterface()	Initialized Parameter object <ul style="list-style-type: none"> Parameter() 	Initilize the GUI and the Renderer.	None
Layer()	None	The Parent class of all the editing layers.	None
Layer.apply()	The graph model of a generated sequence <ul style="list-style-type: none"> Sequence() 	Apply the modifications to the comic sequence and return the results.	Sequence()
AttributeNode()	The attribute type and the initialized Parameter class <ul style="list-style-type: none"> attribute_type Parameter() 	The Parent class of all the attribute nodes, including Sequence, Panel, Character, etc. It is the Node class for the graph model in the Image Sequence Model.	None
addAttribute()	The parent node's name, attribute type, and the child node (self). <ul style="list-style-type: none"> parent_node attribute_type AttributeNode() 	Add an attribute node to the graph model as a child node of the assigned parent node.	None

Table 3
API Functions









Set#				
1				
2				

Table 4
Examples of before and after applying the changes.

on the action causal network, it shows a short narrative in which the two characters ate apples and felt dizzy after eating, then shocked and finally rested in the garden.

3.2. Data Availability

The data and code used in this study are openly available in a GitHub repository. The repository includes the raw data, processed data, and all scripts necessary to reproduce the analyses presented in this paper. You can access the repository at https://github.com/RimiChen/Collaborative_Comic_Generation. The repository is licensed under MIT License, allowing for reuse and modification with appropriate attribution.

4. Conclusion and Future Work

This paper presents an extensible system for generating comic-style visual narratives, integrating narrative theory with human-AI collaboration. We address challenges in visual modification and plotline planning, balancing automation with user control for customization.

While the current system effectively integrates abstract narrative theories, it has limitations in coordinating visual components and diversifying narratives. Separating visual layers, though beneficial for user modifications, reduces scene and character interaction, limiting cohesive artwork and rich actions. Additionally, using symbols to represent actions restricts narrative diversity, requiring user customization to expand content. Integrating advanced language models could enhance narrative richness.

Future work will focus on refining the integration of narrative and visual components and exploring advanced models to enhance the system’s capability to handle diverse narratives. We also plan to conduct user studies to evaluate effectiveness and usability across different groups. In conclusion, our work advances human-AI cooperative visual narrative generation, offering a versatile platform for creating engaging experiences.

References

- [1] N. Cohn, R. Jackendoff, P. J. Holcomb, G. R. Kuperberg, The grammar of visual narrative: Neural evidence for constituent structure in sequential image comprehension, *Neuropsychologia* 64 (2014) 63–70.
- [2] S. McCloud, M. Martin, *Understanding comics: The invisible art*, volume 106, Kitchen sink press Northampton, MA, 1993.
- [3] C. Martens, N. EDU, R. E. Cardona-Rivera, U. EDU, The visual narrative engine: A computational model of the visual narrative parallel architecture, in: *8th Annual Conference on Advances in Cognitive Systems*, 2020.
- [4] M. C. Team, Image sequence creation with microsoft copilot, <https://microsoft.com/copilot>, 2023. Accessed: 2024-05-25.
- [5] Y. Zhou, D. Zhou, M. Cheng, J. Feng, Q. Hou, Storydiffusion: Consistent self-attention for long-range image and video generation, *arXiv preprint arXiv:2405.01434* (2024). URL: <https://arxiv.org/abs/2405.01434>.
- [6] G. Jing, Y. Hu, Y. Guo, Y. Yu, W. Wang, Content-aware video2comics with manga-style layout, *IEEE Transactions on Multimedia (TMM)* 17 (2015) 2122–2133.
- [7] P. P. Gunasekara, P. M. Perera, C. D. Adhihetty, D. D. Kollure, N. Kodagoda, A. Caldera, Generate comic strips using ai, in: *Proceedings of Conference on Transdisciplinary Research in Engineering*, volume 1, 2024.
- [8] T. Alves, A. McMichael, A. Simões, M. Vala, A. Paiva, R. Aylett, Comics2d: Describing and creating comics from story-based applications with autonomous characters, in: *Proceedings of the International Conference on Computer Animation and Social Agents (CASA)*, Springer, 2007, pp. 67–74.
- [9] C. Martens, R. E. Cardona-Rivera, Discourse-driven comic generation, in: *Proc. International Conference on Computational Creativity*, 2016.
- [10] C. Martens, R. E. Cardona-Rivera, Generating abstract comics, in: *International Conference on Interactive Digital Storytelling*, Springer, 2016, pp. 168–175.
- [11] Y.-C. Chen, A. Jhala, A customizable generator for comic-style visual narrative, *arXiv preprint arXiv:2401.02863* (2023). URL: <https://arxiv.org/abs/2401.02863>.
- [12] R. Narita, K. Tsubota, T. Yamasaki, K. Aizawa, Sketch-based manga retrieval using deep features, in: *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 3, IEEE, 2017, pp. 49–53.
- [13] Y. Matsui, K. Ito, Y. Aramaki, A. Fujimoto, T. Ogawa, T. Yamasaki, K. Aizawa, Sketch-based manga retrieval using manga109 dataset, *Multimedia Tools and Applications* 76 (2017) 21811–21838.
- [14] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, 2021. *arXiv:2112.10752*.
- [15] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, <https://github.com/CompVis/stable-diffusion>, 2022. Accessed: 2024-05-24.
- [16] N. Cohn, Visual narrative structure, *Cognitive science* 37 (2013) 413–452.
- [17] N. Cohn, How to analyze visual narratives: A tutorial in visual narrative grammar, Online: http://www.visuallanguagelab.com/P/VNG_Tutorial.pdf [last accessed: 1 March 2016]

(2015).

- [18] M. I. N. COHN, From visual narrative grammar to filmic narrative grammar: The narrative structure of static and moving images, in: *Film Text Analysis*, Routledge, 2016, pp. 106–129.
- [19] A. Mehrabian, J. A. Russell, Basic dimensions for a general psychological theory: Implications for personality, social, environmental, and developmental studies, *Journal of Comparative and Physiological Psychology* 55 (1974) 439–449.
- [20] J. A. Russell, L. M. Ward, G. Pratt, Affective quality attributed to environments: A factor analytic study, *Environment and behavior* 13 (1981) 259–288.
- [21] J. A. Russell, A circumplex model of affect, *Journal of Personality and Social Psychology* 39 (1980) 1161–1178.
- [22] S. Lowe, Roberta base model fine-tuned on goemotions for emotion classification, <https://huggingface.co/SamLowe/roberta-base-go_emotions>, 2021. Accessed: 2024-05-25.
- [23] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 4171–4186.