# ComicScene154: A Scene Dataset for Comic Analysis

**Sandro Paval,  Ivan P. Yamshchikov,  Pascal Meißner**

CAIRO, THWS, Technical University of Applied Sciences Würzburg-Schweinfurt

`sandro.paval@study.thws.de, ivan.yamshchikov@thws.de,pascal.meissner@thws.de`

## Abstract

Comics offer a compelling yet under-explored domain for computational narrative analysis, combining text and imagery in ways distinct from purely textual or audiovisual media. We introduce *ComicScene154*, a manually annotated dataset of scene-level narrative arcs derived from public-domain comic books spanning diverse genres. By conceptualizing comics as an abstraction for narrative-driven, multimodal data, we highlight their potential to inform broader research on multi-modal storytelling. To demonstrate the utility of *Comic-Scene154*, we present a scene segmentation baseline, providing an initial benchmark for future studies to build upon. Our results indicate that *ComicScene154* constitutes a valuable resource for advancing computational methods in multimodal narrative understanding and expanding the scope of comic analysis within the Natural Language Processing community.

## 1 Introduction

Advancements in Computer Vision (CV) and Natural Language Processing (NLP) have enabled the analysis and processing of diverse media types in both standalone and multimodal settings (Bayoudh et al., 2022). However, one medium that remains underexplored is comics, which uniquely interweaves text and images into panels, creating a sequence of discontinuous frames. The narrative link between these frames is implicit, requiring a cognitive process called *closure* to bridge the gaps. This characteristic makes comics particularly intriguing not only as a distinct medium, but also as a compact abstraction for more continuous, multimodal data such as videos or movies (Rao et al., 2020; Cohn, 2013). By sampling frames from data streams and compiling them in a comic-like structure, one can preserve narrative structure while reducing the data's overall complexity.

Progress in computational comic analysis largely falls into four categories: *object detection*

*tasks* (Ogawa et al., 2018; Yanagisawa et al., 2018), *multi-modal understanding* (Li et al., 2024; Rigaud et al., 2024), *image generation* (Proven-Bessel et al., 2021; Yang et al., 2021; Wang et al., 2012) and *dataset proposals* (Iyyer et al., 2017; Dunst et al., 2017). In comparison, the problem of *narrative understanding* (Pratt, 2009) in comics remains relatively underdeveloped.

Comic analysis tasks, such as character re-identification (Sachdeva and Zisserman, 2024) and closure inference (Iyyer et al., 2017), would rely on narrative context but have to resort to operating at broad levels (e.g., entire comics) or strict physical divisions (e.g., single pages). While page-wise processing may be practical from a layout point-of-view, the internal narrative structure is more nuanced (Zehe et al., 2021; Cohn, 2010). Segmenting comics according to narrative arcs could better support tasks like story summarization (Huang et al., 2016) and entity tracking (Kim and Schuster, 2023) by ensuring that only relevant information is included. Conversely, using an entire comic might introduce superfluous content, whereas strict physical divisions risk omitting key story details.

This take on segmentation could also serve as a simplified approach to scene segmentation in video data (Rao et al., 2020), where select frames are extracted to create a comic-like structure conducive to smaller-scale analysis. Beyond comptutational problems, our take on segmentation could benefit tasks in digital humanities: identifying recurring storytelling patterns, comparing artistic styles, or performing linguistic analyses on segmented narrative arcs (Cohn, 2013; Zehe et al., 2021).

In this paper, we aim to advance all these fields by introducing the scene-annotated dataset for comics *ComicScene154*. We devised a prototypical approach to display the limitations and challenges associated with such a task. This dataset can facilitate new segmentation methods while providing a benchmark for existing approaches. Addition-

31574

| Dataset | Tasks | Years | Style | Books | Pages |
|---|---|---|---|---|---|
| eBDtheque | d,t | 1905-2012 | mix | 28 | 100 |
| COMICS | c | 1938-1954 | comics | 3948 | 198000 |
| GCN | d,t | 1978-2013 | comics | 253 | 38000 |
| DCM772 | d | 1938-1954 | comics | 27 | 772 |
| Manga109 | d,t,r | 1970-2010 | manga | 109 | 10000 |
| BCBId | - | - | bangla | 64 | 3000 |
| PopManga | d,t,r | 2010-2023 | manga | 25 | 1800 |
| CoMix | d,t,r,N,D | 1938-2023 | mix | 100 | 3800 |
| **Ours** | d,S | 1942-1962 | comics | 4 | 154 |

Table 1: Overview of comic datasets including ours, adapted from (Vivoli et al., 2024). **Tasks:** Classification (c), Detection (d), Text-Character Association (t), Character Re-Identification (r), Character Naming (N), Dialog (D), and Scene Segmentation (S).

ally, it offers a consistent resource for tasks such as story summarization and character identification on a more fine-grained, narrative-based scale. Section 3 will detail the dataset construction and the evaluation, while Section 4 will present an example to illustrate the inherent challenges of the task.

## 2 Related Work

### 2.1 Scenes in Comics and Beyond

We follow the definition of scenes from (Rao et al., 2020), and describe a scene as a plot-based semantic unit, in which an overarching task is pursued by a certain cast of characters. In most cases, these semantic units maintain temporal and spatial coherence. Additionally, we treat a scene analogous to narrative arcs as defined by (Cohn, 2013)

In movies, different shots—sequences of images captured from the same viewpoint—form the visual narrative. Similarly, comics consists of panels that depict imagery from diverse viewpoints. By selecting a representative image from each shot, along with its corresponding narrative data (e.g., spoken text or script), the resulting sequence can be compacted. By recreating a comic-like structure via representative frames, image data can be significantly reduced beyond comics.

### 2.2 Analogy to Semantic Text Segmentation

In NLP, semantic text segmentation involves breaking text down into semantically coherent segments or representations, often focusing on higher-level logical forms[1]. This task shares conceptual parallels with scene segmentation in comics: both seek

to identify coherent boundaries within a data stream - be it textual sentences or visual-narrative frames.

To evaluate segmentation quality, text-based approaches often rely on the $p_k$ metric, which quantifies the proportion of sentences (or segments) that are incorrectly "cut" or "joined" (Glavaš et al., 2016) via segmentation. We adopt a similar perspective for evaluating scene segmentation: just as $p_k$ measures segmentation consistency in text, we devise and adapt an analogous metric for assessing how effectively consecutive frames (or panels) are grouped into scenes within comics and movies.

### 2.3 Comic Datasets

Obtaining ethically sourced data for comics is non-trivial due to the commercial nature of the medium. The diverse datasets for manga—most notably *Manga109* (Fujimoto et al., 2016) and *PopManga*(Sachdeva and Zisserman, 2024)—are not a suitable replacement, as mangas often differ considerably from classic Western comics in terms of storytelling conventions, structural layout, and linguistic features (Cohn, 2011). These differences highlight the need for building datasets tailored to distinct comic traditions, enabling more targeted research in both CV and multimodal NLP tasks.

Western datasets on the other hand consist almost exclusively of public domain comics, which restricts them to older sources. Notable datasets compiled in Table 1 include DCM772(Nguyen et al., 2018), focusing on object detection and association, or the combined manga and comic dataset CoMix(Sachdeva and Zisserman, 2024) which draws from already existing datasets for benchmarking purposes. In terms of cultural diversity, eBDtheque (Guérin et al., 2013), offers not only a mix of manga and American comics, but

---

[1]For a detailed survey on semantic text segmentation, we refer the reader to (Kamath et al., 2019)

| Comic, Volume | Pages | Panels | Scenes | Genre |
|---|---|---|---|---|
| Alley Oop,1 | 35 | 191 | 37 | Humor |
| Champ Comics,24 | 38 | 263 | 55 | Heroes |
| Treasure Comics,6 | 41 | 221 | 41 | Fantasy |
| Western Love,4 | 40 | 279 | 52 | Love |

Table 2: Details on of our *ComicScene154* dataset.

also French bande dessinée, while BCBId(Dutta et al., 2022) offers Bangladeshi comics.

## 3 The ComicScene154 Dataset

### 3.1 Data Source

For the sake of reproducibility, only freely available public-domain data were used in this study. While large-scale datasets exist for Japanese manga, we opted not to include them here due to the format discrepancies discussed in Section 2.3. All data was collected from *Comic Book Plus*[2], a repository offering a diverse range of public-domain comics. The comics used skew toward older storytelling and artistic styles, reflecting their origins in the "Golden Age" of comics, roughly 1940–1960.

### 3.2 Dataset

Our dataset, *ComicScene154*, comprises four public-domain comic magazines containing 34 distinct stories that span a total of 154 pages across various genres and publication years[3]. These magazines were selected to ensure diversity in storytelling, thus covering multiple narrative styles. Table 2 outlines the chosen comics and their associated metadata. Due to the focus on narrative, none of the existing datasets in Table 1 were used, as their data is not categorized by genre or consist of samples of comics. The latter is particularly important because being able to segment the dataset at the level of full stories is crucial, as scene construction relies on the broader comic narrative. Randomly sampled panels or pages are therefore not suitable for the task of scene segmentation, we address.

Before distributing the dataset to the annotators, all panels were extracted from comic pages along with their coordinates, and then numbered in reading order. Annotating suchlike ensures the consistency of the dataset, as any discrepancies in panel numbering (e.g., missing or extra panels) would corrupt the reading order. Additionally, each panel was tagged with a Boolean value indicating

[2]https://comicbookplus.com
[3]Here is the link to the dataset and code https://github.com/Knorrsche/ComicScene154



Figure 1: Panels with 2 scenes visualized as intervals.

whether it marks the start of a new scene. An example for visualization can be seen at Figure 1, where the scene boundaries are marked blue.

### 3.3 Reliability

Given the inherently subjective nature of scene segmentation, evaluating the reliability of our annotations is critical. To this end, one-third of the dataset was independently labeled by three different groups of two annotators (Tester 1 & Tester 2), and their annotations were compared against our own. This triple annotation process not only ensures data consistency but also gauges whether different annotators share a conceptual understanding of scene boundaries. Annotators were tasked with marking each panel that signals the start of a new narrative arc. Following the definition in Section 2.1, they received a brief introduction to the task and a sample annotated example for guidance.

To quantify agreement, we employed the $p_k$ metric from the semantic text segmentation literature (Glavaš et al., 2016). This measure compares two segmentations by assessing the proportion of sliding windows that are inconsistently segmented. In our context, $p_k = 0$ indicates perfect alignment of scene boundaries, while $p_k = 1$ indicates total misalignment. The choice of window size $k$ is typically informed by half of the average length of segments. Lacking an external reference for segment length, we computed $k$ by averaging the scene lengths from our own annotations and those produced by the other annotators, yielding $k = 3$.

With an average tester $p_k$ score of $0.17 = 0.15 + 0.19/2$, there is some notable agreement in the interpretation of what a scene is yet as with other narrative related tasks there is certain noise associated with the subjectivity of a given assessor. This suggests that though the overall concepts align rea-

| Excerpts | Tester 1 | Tester 2 | In-between |
|---|---|---|---|
| Alley Oop (1) | 0.07 | 0.00 | 0.07 |
| Champs (1) | 0.27 | 0.16 | 0.22 |
| Champs (2) | 0.21 | 0.07 | 0.28 |
| Treasure C. (1) | 0.06 | 0.19 | 0.14 |
| Treasure C. (2) | 0.17 | 0.33 | 0.26 |
| Western L. (1) | 0.12 | 0.37 | 0.29 |
| **Average** | 0.15 | 0.19 | 0.21 |

Table 3: Agreement scores ($p_k$) for 6 randomly chosen excerpts. Each excerpt was tested by 2 of 6 tester.

sonably well in Table 3, a major challenge in scene segmentation remains the lack of an intersubjective definition of a scene. The varying scores illustrate this issue—most notably in the excerpt extracted from **Western Love**, where Tester 2 exhibited the highest disagreement with a score of 0.37. When analyzing the annotation data, it became clear that Tester 2 segmented the comic into much shorter scenes. This behavior can be attributed to the inherent subjectivity of the interpretation of narratives.

## 4 Scene Segmentation Benchmark

For benchmarking the dataset, a two step scene segmentation pipeline was developed: First, utilizing a multi-modal model to predict scene boundaries and then refining the predictions using a reasoning-based large language model (LLM).

For both steps, Gemini's reasoning model, `gemini-2.0-flash-thinking-exp`, was used. In the initial iteration, each comic page, along with the coordinates of the panels and their reading order, was provided as context alongside a prompt. The model's task was to generate descriptions of narrative arcs and identify the panel where each arc begins, and hence scenes. While these outputs were useful, their quality remained limited.

Due to the non-deterministic nature of LLM, the performance was evaluated using the $p_k$ metric, based on ten iterations per comic. Table 4 presents the average scores in these ten iterations, compared to scenes we defined randomly. The results indicate that the performance of the multi-modal model is only marginally better than random definition, highlighting the current limitations. Additionally, the agreement between the ten iterations was analyzed, revealing high consistency. This suggests that while the model struggles to fully align with human-annotated scenes, it does capture some underlying patterns. It also highlights the present challenge of subjectiveness in scene segmentation.

| Comic | Dataset | Random | In-between |
|---|---|---|---|
| Alley Oop | 0.43 | 0.47 | 0.10 |
| Champs | 0.40 | 0.43 | 0.09 |
| Treasure Comics | 0.43 | 0.46 | 0.05 |
| Western Love | 0.43 | 0.46 | 0.07 |
| **Average** | 0.42 | 0.46 | 0.06 |

Table 4: Agreement scores ($p_k$) of multi-modal, random and in-between different multi-modal iterations

To further refine the results, the initial model outputs were processed using a reasoning-based LLM. Since this refinement step depends on the output of the multi-modal model, each of the ten iterations was further benchmarked ten times, resulting in a total of 100 evaluations per comic. Table 5 compares the average score across all 100 iterations to the best average score from a single multi-modal output (based on ten iterations).

While the overall output still exhibits a relatively high $p_k$ score, and the improvement from averaging 100 iterations is marginal, the best-performing iterations for each comic show noticeable gains. Although these best results represent the most favorable outcomes, it is important to note that they still consist of ten independent iterations based on the same multi-modal model output. These large performance variations further illustrate the inherent challenges and subjectivity in scene segmentation.

| Comic | All Iter. | Best Iter. | In-between |
|---|---|---|---|
| Alley Oop | 0.36 | 0.27 | 0.04 |
| Champs | 0.40 | 0.37 | 0.05 |
| Treasure Comics | 0.41 | 0.34 | 0.05 |
| Western Love | 0.39 | 0.36 | 0.04 |
| **Average** | 0.39 | 0.34 | 0.05 |

Table 5: Refined agreement scores ($p_k$)

## 5 Discussion

Comparing the results of both human (Table 3) and AI (Tables 4, 5) benchmarks reveals a significant difference. While the human benchmark showed notable improvement compared to randomly defined scenes, the multi-modal approach demonstrated almost no improvement. The refined method did improve results, but the largest performance gains were observed only in specific iterations rather than across all outputs. This highlights the present challenge of defining scenes, as even human annotators disagree. Nevertheless, while the refined performance relied on the multi-modal output, the annotations themselves showed low differences in between, demonstrating the model's ability to detect certain patterns—though not the

semantic units we aimed to capture, i.e., the scenes.

# 6 Conclusion

We introduce ComicScene154, a novel dataset designed to facilitate the study of scene segmentation in comics. By addressing the unique challenges (e.g. closure, narrative segmentation) of this medium, ComicScene154 provides a robust foundation for developing algorithms capable of understanding complex narrative structures like scenes (narrative arcs). Furthermore we propose the leveraging of comic analysis techniques on formats other than comics, like movies, to reduce the data scale and complexity of tasks such as narrative summarizations or object clustering.

## Limitations

Despite its contributions, ComicScene154 has certain limitations. The primary challenge remains the subjectivity of scene segmentation, as interpretations can vary across annotators. This was also seen in the performance of scene segmentation, where the average annotations of the models displayed lacking results. Additionally, the dataset is largely composed of "Golden Age" comics, which may limit its applicability to modern comics due to differences in artistic style, storytelling techniques, and narrative complexity.

## Ethics Statement

This paper complies with the ACL Ethics Policy. The comics included in ComicScene154 are in the public domain, ensuring compliance with copyright regulations. However, we acknowledge the potential for annotator bias, which could impact dataset consistency. To mitigate this, we implemented a standardized annotation framework and conducted inter-annotator agreement evaluations to enhance reliability.

## Acknowledgments

We extend our gratitude to the annotators for their meticulous work and to the reviewers for their valuable feedback.

## References

Khaled Bayoudh, Raja Knani, Fayçal Hamdaoui, and Abdellatif Mtibaa. 2022. A survey on deep multimodal learning for computer vision: advances, trends, applications, and datasets. *The Visual Computer*, 38(8):2939–2970.

Neil Cohn. 2010. The limits of time and transitions: Challenges to theories of sequential image comprehension. *Studies in Comics*, 1(1):127–147.

Neil Cohn. 2011. A different kind of cultural frame: An analysis of panels in american comics and japanese manga. *Image & Narrative*, 12(1):120–134.

Neil Cohn. 2013. *The Visual Language of Comics: Introduction to the Structure and Cognition of Sequential Images*. Bloomsbury Academic, New York.

Alexander Dunst, Rita Hartel, and Jochen Laubrock. 2017. The graphic narrative corpus (gnc): design, annotation, and analysis for the digital humanities. In *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, volume 3, pages 15–20. IEEE.

Arpita Dutta, Samit Biswas, and Amit Kumar Das. 2022. Bcbid: first bangla comic dataset and its applications. *International Journal on Document Analysis and Recognition (IJDAR)*, 25(4):265–279.

Azuma Fujimoto, Toru Ogawa, Kazuyoshi Yamamoto, Yusuke Matsui, Toshihiko Yamasaki, and Kiyoharu Aizawa. 2016. Manga109 dataset and creation of metadata. In *Proceedings of the 1st international workshop on comics analysis, processing and understanding*, pages 1–5.

Goran Glavaš, Federico Nanni, and Simone Paolo Ponzetto. 2016. Unsupervised text segmentation using semantic relatedness graphs. In *Proceedings of the Fifth Joint Conference on Lexical and Computational Semantics*, pages 125–130. Association for Computational Linguistics.

Clément Guérin, Christophe Rigaud, Antoine Mercier, Farid Ammar-Boudjelal, Karell Bertet, Alain Bouju, Jean-Christophe Burie, Georges Louis, Jean-Marc Ogier, and Arnaud Revel. 2013. ebdtheque: a representative database of comics. In *2013 12th International Conference on Document Analysis and Recognition*, pages 1145–1149. IEEE.

Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239.

Mohit Iyyer, Varun Manjunatha, Anupam Guha, Yogarshi Vyas, Jordan Boyd-Graber, Hal Daume, and Larry S Davis. 2017. The amazing mysteries of the gutter: Drawing inferences between panels in comic book narratives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern recognition*, pages 7186–7195.

Apurva Kamath, Rishiraj Das, et al. 2019. A survey on semantic parsing. *ACM Computing Surveys*.

Najoung Kim and Sebastian Schuster. 2023. Entity tracking in language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3835–3855.

Yingxuan Li, Ryota Hinami, Kiyoharu Aizawa, and Yusuke Matsui. 2024. Zero-shot character identification and speaker prediction in comics via iterative multimodal fusion. *arXiv preprint arXiv:2404.13993*.

Nhu-Van Nguyen, Christophe Rigaud, and Jean-Christophe Burie. 2018. Digital comics image indexing based on deep learning. *Journal of Imaging*, 4(7):89.

Toru Ogawa, Atsushi Otsubo, Rei Narita, Yusuke Matsui, Toshihiko Yamasaki, and Kiyoharu Aizawa. 2018. Object detection for comics using manga109 annotations. *arXiv preprint arXiv:1803.08670*.

Henry John Pratt. 2009. Narrative in comics. *The Journal of Aesthetics and Art Criticism*, 67(1):107–117.

Ben Proven-Bessel, Zilong Zhao, and Lydia Chen. 2021. Comicgan: Text-to-comic generative adversarial network. *arXiv preprint arXiv:2109.09120*.

Anyi Rao, Linning Xu, Yu Xiong, Guodong Xu, Qingqiu Huang, Bolei Zhou, and Dahua Lin. 2020. A local-to-global approach to multi-modal movie scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Christophe Rigaud, Jean-Christophe Burie, and Samuel Petit. 2024. Toward accessible comics for blind and low vision readers. In *International Conference on Document Analysis and Recognition*, pages 198–215. Springer.

Ragav Sachdeva and Andrew Zisserman. 2024. The manga whisperer: Automatically generating transcriptions for comics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12967–12976.

Emanuele Vivoli, Marco Bertini, and Dimosthenis Karatzas. 2024. Comix: A comprehensive benchmark for multi-task comic understanding. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Meng Wang, Richang Hong, Xiao-Tong Yuan, Shuicheng Yan, and Tat-Seng Chua. 2012. Movie2comics: Towards a lively video content presentation. *IEEE Transactions on Multimedia*, 14(3):858–870.

Hideaki Yanagisawa, Takuro Yamashita, and Hiroshi Watanabe. 2018. A study on object detection method from manga images using cnn. In *2018 International Workshop on Advanced Image Technology (IWAIT)*, pages 1–4. IEEE.

Xin Yang, Zongliang Ma, Letian Yu, Ying Cao, Baocai Yin, Xiaopeng Wei, Qiang Zhang, and Rynson WH Lau. 2021. Automatic comic generation with stylistic multi-page layouts and emotion-driven text balloon generation. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 17(2):1–19.

Albin Zehe, Leonard Konle, Lea Katharina Dümpelmann, Evelyn Gius, Andreas Hotho, Fotis Jannidis, Lucas Kaufmann, Markus Krug, Frank Puppe, Nils Reiter, et al. 2021. Detecting scenes in fiction: A new segmentation task. In *Proceedings of the 16th conference of the European chapter of the association for computational linguistics: Main volume*, pages 3167–3177.

# 7 Appendix

## 7.1 Annotators

The annotators were friends or colleagues of the authors who followed the provided guidelines without any additional input from the authors. They volunteered without compensation.

## 7.2 Guideline

**Dear Participant,**

I warmly invite you to participate in a survey designed to evaluate the accuracy of a dataset I have developed. The focus of this survey is to assess how effectively the dataset can identify the start and end points of scenes in comics.

**Purpose of the Survey:**

The segmentation of comic scenes is often subjective and can vary from reader to reader. Sometimes, it is difficult to pinpoint exactly where one scene ends and the next begins. In such cases, deviations are completely acceptable. What is most important is whether there is significant agreement on the core panels of a scene—the key panels that form the heart of the storyline.

**Instructions:**

1. **Review the comic images:** You will be provided with a selection of comic images representing different scenes.

2. **Mark the start and end points:** Indicate or describe where you perceive the transition from one scene to the next.

3. **Reference example:** A sample comic is attached for your orientation, but please note that it is merely a reference—you do not need to strictly follow it.

**Important Note:**

Comics are a visual art form, and often, there is no definitive answer to where a scene begins or ends. Your feedback will help me identify patterns and commonalities that will be valuable in improving the dataset.

## 7.3 Additional Statistics