# ASCII-Bench: Evaluating Language-Model-Based Understanding of Visually-Oriented Text

### **Anonymous Author(s)**

Affiliation Address email

### **Abstract**

Large language models (LLMs) have demonstrated several emergent behaviors with scale, including reasoning and fluency in long-form text generation. However, they continue to struggle with tasks requiring precise spatial and positional reasoning. ASCII art, a symbolic medium where characters encode structure and form, provides a unique probe of this limitation. We introduce ASCIIBench, a novel benchmark for evaluating both the generation and classification of ASCII-text images. ASCIIBench consists of a filtered dataset of 5,315 class-labeled ASCII images and is, to our knowledge, the first publicly available benchmark of its kind. Alongside the dataset, we release weights for a fine-tuned CLIP model adapted to capture ASCII structure, enabling the evaluation of LLM-generated ASCII art. Our analysis shows that cosine similarity over CLIP embeddings fails to separate most ASCII categories, yielding chance-level performance even for low-variance classes. In contrast, classes with high internal mean similarity exhibit clear discriminability, revealing that the bottleneck lies in representation rather than generational variance. These findings position ASCII art as a stress test for multimodal representations and motivate the development of new embedding methods or evaluation metrics tailored to symbolic visual modalities. All resources are available at https://github.com/ASCIIBench/ASCIIBench.

### 19 1 Introduction

2

3

5

6

7

8

9

10

11 12

13

14

15

16

17

18

Scaling language models has been shown to induce emergent capabilities [Wei et al., 2022], includ-20 ing those involving positional understanding, such as the generation and editing of TikZ drawings 21 [Bubeck et al., 2023]. We define ASCII art as the intersection of text and vision. The generation 22 and classification of ASCII art introduces challenges that are distinct from conventional NLP and 23 multimodal benchmarks: characters function as visual primitives rather than semantic tokens, necessitating strict structural regularity seen in other forms of structured data like tables [Chen, 2022]. 25 26 In contrast to natural images, ASCII art is both present in the pretraining distribution of unimodal 27 language models and natively aligned with their tokenization schemes, enabling direct evaluation without additional adaptation. 28

### 29 **2** The ASCIIBench Dataset

- 30 We introduce ASCIIBench, a high-quality benchmark for ASCII art understanding and generation.
- 31 Sourced ethically from ascii.co.uk, the data underwent a rigorous multi-stage curation pipeline.
- The final dataset contains 5,315 unique ASCII art pieces across 752 classes (e.g., aircraft, birds).
- 33 All art is credited to the original creators on ascii.co.uk. In the absence of explicit licensing, we
- adhered to standard research practices described in Appendix A.

#### 35 2.1 Data Curation & Analysis

- 36 Raw ASCII art contains pervasive noise like signatures and tags. An in-depth description of our data
- cleaning methodology and dataset analysis can be found in Appendix B and Appendix C.

# 38 3 Classification

- 39 **Models** We evaluated multiple models on classification and generation tasks, including Llama 3-8B,
- 40 Llama 3-8B-Instruct, GPT-3.5, GPT-4o, GPT-4o-mini, GPT-5-mini, and Claude 3.5 Sonnet, testing
- 41 text-only, vision-only, and text-vision prompts to compare performance across modalities.

# 42 3.1 Model Testing Procedure



Figure 1: Example classification prompt with result

- 43 **Prompt** ASCII images are preprocessed based on input modalities. Image preprocessing is de-
- 44 scribed in Appendix F. The model is then prompted to select one of four choices in the format shown
- in Figure 1.
- 46 Evaluation Metrics Performance was measured by the model's macro and micro accuracy

### 47 4 Generation

- 48 **Models** We prompted GPT-3.5, GPT-4, and GPT-40 to generate 5 ASCII images for each class.
- 49 **Approach** To evaluate fidelity, we require an image-to-image metric that captures both the visual
- and textual characteristics of ASCII art. We use CLIP [Radford et al., 2021], which aligns images
- 51 and text through large-scale contrastive training. By comparing embeddings of generated images to
- 52 reference embeddings derived from ground-truth data, we assess generation accuracy.

# 53 4.1 Evaluation Metrics

- We leverage CLIP cosine similarity between generated and reference images and representation
- 55 quality of the embedding space using alignment and uniformity [Wang and Isola, 2022]. We report
- 56 ROC-AUC for same-class retrieval in Section 5.2. ROC-AUC (Receiver Operating Characteristic –
- 57 Area Under the Curve) quantifies how well a model separates positive from negative pairs, with 0.5
- indicating random performance and 1.0 indicating perfect discrimination.

# 4.2 CLIP Cosine Similarity

59

- Purpose CLIP cosine similarity is a metric used to evaluate how similar two images are in the context of their high-level features extracted by the CLIP model.
- 62 **Implementation** ASCII art is rendered following the steps in Appendix F and then embedded with
- 63 CLIP. The CLIP model, known for its ability to understand high-level visual concepts through natural
- 64 language supervision, is used to process these images [Radford et al., 2021]. The model extracts
- 65 feature vectors representing the semantic content of each image. The cosine similarity score ranges
- 66 from -1 (completely different) to 1 (exactly the same), with higher scores indicating greater similarity
- 67 [Radford et al., 2021].

# 88 4.3 Alignment & Uniformity

Alignment measures intra-class compactness, while uniformity quantifies dispersion in the embedding space [Wang and Isola, 2022]. Out-of-the-box CLIP shows alignment of 5.85 (squared 34.20). Fine-tuning increases alignment to 8.90 (squared 79.16) and improves uniformity from baseline to -7.61 ( $t\!=\!1$ ), -8.09 ( $t\!=\!5$ ), and -8.21 ( $t\!=\!10$ ). Together with stable cosine similarities, these results confirm that CLIP is not experiencing representation collapse.

# 74 5 Results

### 5.1 Classification Results

We evaluate the performance of various models when classifying ASCII art using the methods in Section 3.1 with a maximum of 50 output tokens. We report results across three modalities: T (text-only), V (vision-only), and T+V (text+vision). Responses were filtered for possible string parsing errors, resulting in a <2% average removal. Unfiltered and filtered results are shown in Table 1.

Table 1: Model performance comparison on raw (left) and filtered (right) datasets.

Raw (Unfiltered) Dataset					Filtered Dataset				
Model	Mod.	Micro acc. (%)	Macro acc. (%)	Pass rate (%)	Model	Mod.	Micro acc. (%)	Macro acc. (%)	Pass rate (%)
LLaMA3.1-8B-Inst	Т	34.27	31.89	91.78	LLaMA3.1-8B-Inst	T	34.50	32.01	91.69
LLaMA3.1-8B	T	29.00	25.07	82.67	LLaMA3.1-8B	T	29.39	25.40	82.69
GPT-5-mini	T	61.60	62.39	99.38	GPT-5-mini	T	61.36	61.97	99.35
	V	77.25	84.13	99.24		V	77.25	84.13	99.24
	T+V	73.27	73.84	99.01		T+V	73.27	73.84	99.01
GPT-4o-mini	T	73.61	77.60	95.23	GPT-4o-mini	T	73.52	77.27	95.19
	V	75.72	77.77	97.27		V	75.72	77.77	97.27
	T+V	76.02	77.55	96.67		T+V	76.02	77.55	96.67
GPT-4o	T	75.44	80.23	96.63	GPT-4o	T	75.64	80.26	96.61
	V	77.49	82.16	98.75		V	77.49	82.16	98.75
	T+V	76.56	79.74	98.52		T+V	76.02	77.55	96.67
GPT-3.5-turbo	T	39.05	33.54	91.34	GPT-3.5-turbo	T	39.98	33.77	91.31
Claude-3.5-Sonnet	T	59.55	56.98	98.54	Claude-3.5-Sonnet	T	59.84	57.23	98.65
	V	76.40	76.92	99.08		V	76.40	76.92	99.08
	T+V	76.48	76.89	99.08		T+V	76.48	76.89	99.08

# 5.1.1 Interpretation

83

84 85

86

87

88

89

90

91

Our results align with those of Jia et al. [2024]. Larger models had greater performance, and all accuracy values were over 25%, indicating that models did not choose arbitrarily. Across both raw and filtered datasets, we find that vision-only models consistently outperform text-only and text+vision counterparts, with GPT-40 achieving the highest macro accuracy at 82.2%. Text-only performance lags significantly, especially for LLaMA and GPT-3.5, underscoring the difficulty of modeling ASCII art as pure text. Surprisingly, adding text to vision does not improve performance and in some cases degrades it, suggesting that current multimodal fusion strategies do not capture ASCII structure effectively. Filtering has little effect on overall trends, indicating robustness of the observed modality gaps.

### 5.2 Generation Results

On unfiltered generations, CLIP showed weak class separation (ROC-AUC  $\approx 0.55$ ; silhouette -0.46), and t-SNE revealed no clear clusters. After filtering inconsistent generations (std > 0.15, mean similarity < 0.3), ROC-AUC rose to 0.83, demonstrating that CLIP can discriminate effectively when ASCII generations are semantically consistent. This indicates that the bottleneck lies in the quality of LLM-generated ASCII rather than in the evaluator.

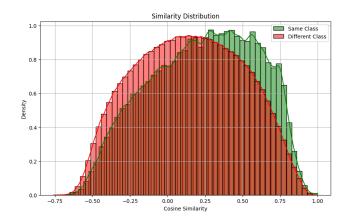


Figure 2: Cosine similarity distributions. Green indicates positive (intra-class) pairs, red indicates negative (inter-class) pairs.

97 Shown in Figure 2, while there is a separation between inter- and intra-class distributions, the 98 separation is not complete, with substantial overlap remaining.

# 99 5.2.1 CLIP Representation Analysis

We examined cosine similarities, silhouette scores, and t-SNE visualizations of CLIP embeddings 100 (Figure 5). Out-of-the-box CLIP produced weak intra- and inter-class separation (AUC  $\approx 0.558$ , 101 silhouette -0.46), and fine-tuning with triplet loss yielded only modest gains. Filtering noisy genera-102 tions did not resolve this, as even the lowest-variance subsets approached near chance performance 103 (AUC = 0.641). However, when restricting analysis to classes with high mean similarity, AUC 104 increased to 0.83, indicating that CLIP can represent ASCII structure only for a subset of well-formed 105 categories. These results show that the primary limitation lies in CLIP's representational capacity for 106 ASCII art, along with variance in model generations. 107

# 108 6 Limitations

Our findings show that evaluation quality depends strongly on input consistency. CLIP performs well only when generations are visually coherent and semantically aligned, but typical LLM outputs are noisy and inconsistent, especially for vague categories. This highlights a dual bottleneck: the instability of ASCII generation and the limited ability of a broad, general-purpose model like CLIP to represent ASCII structure. Filtering demonstrates an upper bound of performance but is not a sustainable evaluation strategy, as it amounts to testing on inputs already close to the training distribution. Future work should explore specialized, smaller models, which may capture ASCII-specific patterns more effectively than CLIP.

# 7 Conclusion

117

We introduce ASCIIBench, a benchmark for evaluating ASCII art on classification and generation 118 tasks, and used it to probe how multimodal models represent symbolic visual inputs. Empirically, 119 vision-only models consistently outperform text-only and text+vision settings on classification, while 120 CLIP-based evaluation of generations provides limited class separation on unfiltered outputs and 121 improves primarily for classes with high internal similarity. These trends position ASCII art as 122 a stringent stress test for multimodal reasoning: performance hinges on both the consistency of 123 generations and the representational suitability of the embedding model for ASCII structure. Looking 124 ahead, we advocate standardized rendering and preprocessing protocols to enable fair cross-model 125 comparisons, improved prompting and training strategies for ASCII generation, and exploration of structure and variance-aware metrics to better capture and evaluate symbolic layout.

# References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
   Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan
- Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian
- Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo
- Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language
- model for few-shot learning, 2022.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, John A. Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuan-Fang Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi,
- Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with
- gpt-4. ArXiv, abs/2303.12712, 2023. URL https://api.semanticscholar.org/CorpusID:
- 139 257663729.
- Wenhu Chen. Large language models are few(1)-shot table reasoners. ArXiv, abs/2210.06710, 2022.
   URL https://api.semanticscholar.org/CorpusID:252872943.
- Moonjun Chung and Taesoo Kwon. Fast text placement scheme for ascii art synthesis. *IEEE Access*, 10:40677–40686, 2022. doi: 10.1109/ACCESS.2022.3167567.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. A neural algorithm of artistic style, 2015.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
   Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- Qi Jia, Xiang Yue, Shanshan Huang, Ziheng Qin, Yizhu Liu, Bill Yuchen Lin, and Yang You. Visual perception in text strings, 2024. URL https://arxiv.org/abs/2410.01733.
- Fengqing Jiang, Zhangchen Xu, Luyao Niu, Zhen Xiang, Bhaskar Ramasubramanian, Bo Li, and Radha Poovendran. Artprompt: Ascii art-based jailbreak attacks against aligned llms, 2024.
- Yongcheng Jing, Yezhou Yang, Zunlei Feng, Jingwen Ye, Yizhou Yu, and Mingli Song. Neural style transfer: A review, 2018.
- Rémi Kazmierczak, Gianni Franchi, Nacim Belkhir, Antoine Manzanera, and David Filliat. A study of deep perceptual metrics for image quality assessment, 2022.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language, 2019.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguisticrepresentations for vision-and-language tasks, 2019.
- Kazuyuki Matsumoto, Akira Fujisawa, Minoru Yoshida, and Kenji Kita. Ascii art classification based on deep neural networks using image feature of characters. *J. Softw.*, 13:559–572, 2018. URL https://api.semanticscholar.org/CorpusID:53281518.
- Katsunori Miyake, Henry Johan, and Tomoyuki Nishita. An interactive system for structure-based
   ascii art creation. 01 2011.
- Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models, 2020.
- Yingxue Pang, Jianxin Lin, Tao Qin, and Zhibo Chen. Image-to-image translation: Methods and applications, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
  Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
  Learning transferable visual models from natural language supervision, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michael Swedrowski, Michael Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

203

204

205

206

207

208

209 210

211

212

213

214

215

218

219

220

221

222

223

224

225

226

227

228

229

230

Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sa-232 jant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, 233 Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan 234 Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, 235 Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, 236 Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima, 237 Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, 238 Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano 239 Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, 240 Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, 241 Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas 242 Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Ger-243 stenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair 247 Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan 248 Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. 249 Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the 250 capabilities of language models, 2023. 251

Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere, 2022.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
 Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals,
 Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022.

Xuemiao Xu, Linling Zhang, and Tien-Tsin Wong. Structure-based ascii art. In *ACM SIGGRAPH* 2010 Papers, SIGGRAPH '10, New York, NY, USA, 2010. Association for Computing Machinery.
 ISBN 9781450302104. doi: 10.1145/1833349.1778789. URL https://doi.org/10.1145/1833349.1778789.

Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan,
 William Yang Wang, and Linda Ruth Petzold. Gpt-4v(ision) as a generalist evaluator for vision language tasks, 2023.

Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2021. doi: 10.1109/JPROC.2020.3004555.

# 267 Appendix

# 268 A Data Sourcing

We express our gratitude to the ASCII artists. We made slight modifications to the original ASCII art and provide the URL to the source of our data. Our dataset is licensed under CC BY NC 4.0.

# B Data Curation

We developed a custom web crawler and an 11-step automated pipeline to remove these artifacts, followed by a multi-stage manual review described in Appendix E. Abstract or ambiguous categories (e.g. "small") were excluded, retaining only well-defined classes. Three annotators then applied a strict rubric to eliminate pieces with: (1) inappropriate content, (2) excessive intra-class variation, (3) overly complex structures, or (4) low quality. This conservative process, requiring strong annotator agreement, removed over 13,000 low-quality images and 1,800 ambiguous classes, resulting in a focused, high-quality benchmark.

# 279 C Data Analysis

280

281

282

283

284

285

286

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

The curated dataset exhibits a natural long-tail class distribution (Figure 3). The largest categories are aircraft (13.3%), land transportation (11.1%), and birds (10.4%) (Figure 4). A t-SNE visualization of class embeddings (Figure 5) confirms semantic coherence, showing clear clustering of related concepts (e.g. animals), demonstrating that ASCII art encodes learnable semantic structures. Character frequency analysis (Figure 7) reveals the artistic "vocabulary": the space character is dominant (>1.6M occurrences), followed by structural elements like -, |, and \_. Alphanumeric characters are used sparingly as accents.

# D Dataset Analysis Figures

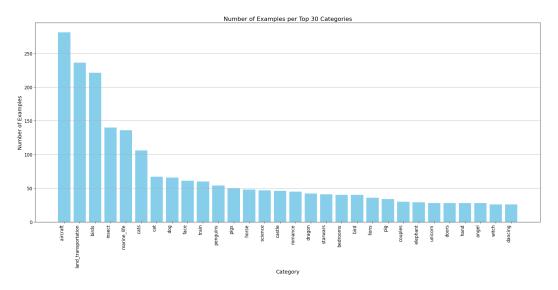


Figure 3: Top 30 Class Histogram

# 288 E Noise Removal Pipeline

- 1. Tags consisting of three or fewer alphabetic characters are replaced with whitespace.
- 2. Creator tags in the last three lines of the artwork are detected and removed while preserving structural spacing.
- 3. Non-visible Unicode control characters are filtered out.
- 4. Alphanumeric creator tags appended to the end of the ASCII art are removed.
- 5. Full names and abbreviated creator tags (e.g., 'Matthew Kenner' or 'jps') positioned on the right margin are eliminated.
- 6. Tags labeled as 'unknown' are discarded.
- 7. Left-aligned creator signatures are identified and removed.
- 8. Common date formats (e.g., 12/21/2023, 21nov2023, 12.21.2023) are detected via expression matching and subsequently stripped.
- 9. Contact information, such as email addresses, is localized and pruned.
- 10. Three-letter creator tags enclosed in dashes (-) or brackets ([]) are filtered out.
  - 11. Known problematic signatures, maintained in a blacklist, are systematically removed.

# F Preprocessing

We follow Jia et al. [2024], using a black monospaced font (DejaVu Sans Mono) on a white background. No blur is added to preserve structural integrity.

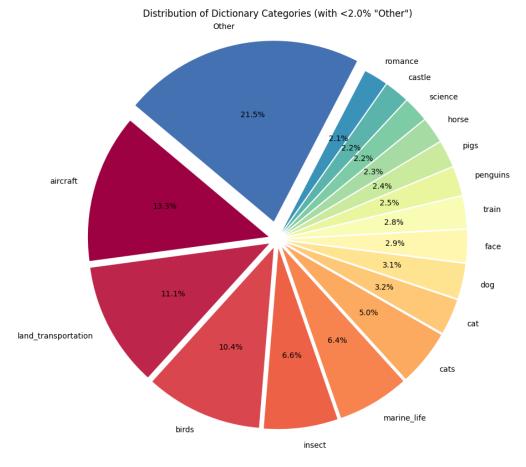


Figure 4: Class Distribution Pie Chart

# **G** Related Works

**Emergent Behaviors:** LLMs have been studied for their emergent properties as they scale. Wei et al. [2022] highlights that larger models, ones with more parameters and diverse data, possess emergent abilities such as improved reasoning and fluency in text generation. Our work extends this to ASCII-text image generation, which requires textual understanding and visual creativity, skills not typically emphasized in standard LLM evaluations.

**LLM ASCII Word Recognition:** The Beyond the Imitation Game benchmark (BIG-bench) introduced by Srivastava et al. [2023] addresses the need to understand the capabilities of LLMs across various tasks. Our research focused on its ASCII word recognition dataset. Similarly, the ArtPrompt jailbreak attack highlights the need for the improvement of LLM performance in identifying ASCII text, an ability crucial to prevent the circumvention of safeguards and elicitation of unintended behaviors [Jiang et al., 2024].

Vision-Language Integration and Multimodal Models: GPT-4V(ision) produces human-aligned scores with detailed explanations, showing promise as a universal automatic evaluator despite some limitations [Zhang et al., 2023]. Flamingo models demonstrate strong few-shot learning capabilities, showcasing potential to give LLMs adaptive abilities and decreased dependence on large task-specific datasets [Alayrac et al., 2022]. Ramesh et al. [2021] introduces an autoregressive transformer-based approach for text-to-image generation with competitive zero-shot performance compared to domain-specific models. These capabilities lead into ASCII art generation, which poses unique challenges due to merging textual and visual information. However, the advanced multimodal reasoning of these models also introduces new vulnerabilities.[Jia et al., 2024] demonstrate that Large Vision-Language Models (LVLMs) like LLaVA and GPT-4V are highly susceptible to *Self*-

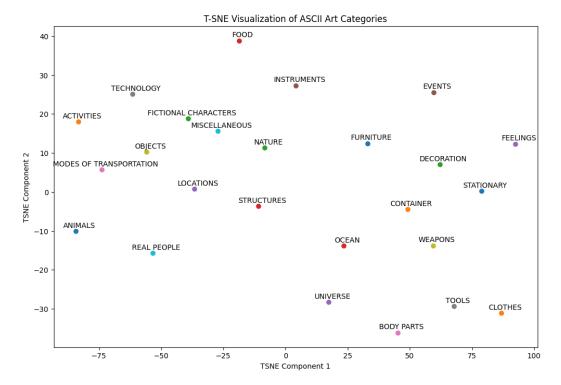


Figure 5: T-SNE visualization of class embeddings

Generated Typographic Attacks, where the model itself is leveraged to create deceptive text and descriptions that cause misclassification, reducing accuracy by up to 60%. This underscores a critical weakness in how LVLMs fuse and weight visual and textual signals. Prior works have explored the generation of stylized textual outputs using neural networks, with some focusing specifically on the artistic transformations of text and image data [Matsumoto et al., 2018]. Models like ViLBERT and VisualBERT use multimodal pre-training to boost performance in vision-language tasks [Li et al., 2019] [Lu et al., 2019]. They guide our fine-tuning of the CLIP model for effective ASCII art generation from text descriptions.

**Evaluation Metrics and Methodologies:** Finally, our evaluation of LLM outputs uses standard metrics used in both language and image processing domains. Metrics such as FID scores, typically used to assess image quality, were adapted to assess the uniqueness and clarity of ASCII-text images produced by our model. We employ CLIP for its ability to effectively bridge text and image representations. Research by Radford et al. [2021] demonstrates its robust performance in zero-shot classification tasks. This makes CLIP ideal for our needs, as our models must both generate and classify ASCII images from minimal prompts.

# **G.1** Other ASCII generation methods

The exploration of AI in the context of ASCII art has witnessed growing interest in recent years, with researchers exploring methods to optimize conversion accuracy. There are many notable contributions in this field. [Goodfellow et al., 2014] proposed a new framework for estimating generative models via an adversarial process. Researchers have delved into the use of GANs to generate realistic and well-designed ASCII art. By leveraging the adversarial training paradigm, these models can produce a large range of high-quality ASCII representations. Similarly, [Gatys et al., 2015] explored Convolutional Neural Networks and their ability to create artistic imagery by experimenting with the style and content of an image, also known as Neural Style Transfer. [Jing et al., 2018] provides an extension of this idea, comparing different Neural Style Transfer qualitatively and quantitatively. They discuss applications of NST and problems to be addressed in future research. Studies have investigated how deep neural networks can be trained to transfer artistic styles onto ASCII images, showcasing the potential for creative synthesis.

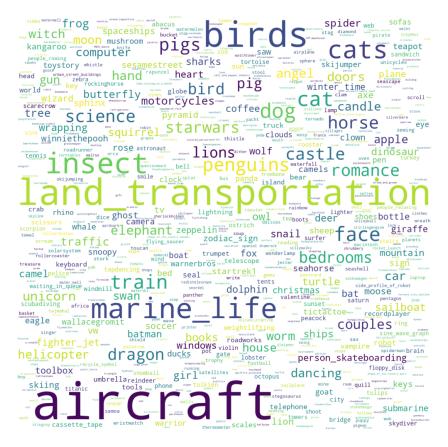


Figure 6: Word cloud of class labels

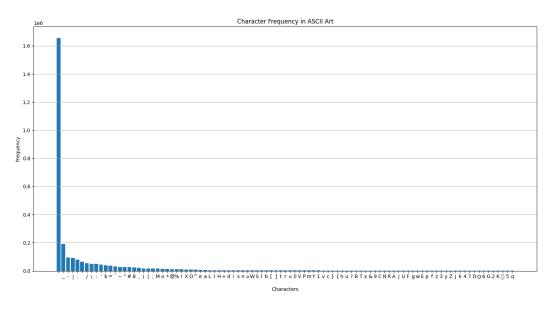


Figure 7: Character Frequency Histogram

# 356 G.1.1 Transfer Learning

Naturally, it is also extremely important for the models to efficiently produce accurate results. [Zhuang et al., 2021] reviews more than forty representative transfer learning approaches from a data and

model perspective. Their paper provides more than twenty experiments of different learning models' performances. The exploration of the efficacy of transfer learning in training models for ASCII art generation and leveraging pre-trained models on large datasets allow researchers to enhance the efficiency and artistic quality of AI-generated ASCII images.

### G.1.2 Methods & Metrics

363

371

385

397

Our research paper took inspiration from methods mentioned in [Pang et al., 2021], which explored developments in I2I translations and analyzed key techniques to elaborate on the effect of I2I on the research and industry community. The paper introduced the problem setting of the image-to-image translation task, introduced the generative models used for I2I methods, and discussed the work and applications of multi-domain I2I tasks. Many of these methods and metrics mentioned in this paper are parallel to evaluation metrics we implemented, but while this paper mainly evaluated methods on I2I, our research drew comparisons between these methods and ASCII image generation standards.

# **G.1.3** Frechet Inception Distance

More specifically, Frechet Inception Distance, which produces an FID score, has been the most 372 widely used metric for measuring the similarity between real and generated images. Muhammad 373 [Naeem et al., 2020] focuses on the reliability of certain methods. They concluded that while variants 374 of precision and recall metrics are generally unreliable methods, density and coverage metrics will 375 provide more interpretable and reliable comparisons. While precision metrics can overestimate the 376 manifold around real outliers, density fixes this issue by more accurately representing the distribution around real samples. The objective of coverage is to improve upon the recall metric. When models generate many unrealistic and diverse samples, this can skew the data and lead to a false increase 379 in the recall measure. Coverage addresses this by building the manifolds around real samples as 380 opposed to fake ones. This approach is less prone to overestimation since real samples tend to have 381 fewer outliers compared to generated samples. The goal for our purposes would be to capture how 382 well the generated ASCII art (fake samples) represents the original ASCII art (real samples) in both 383 details (density) and overall composition (coverage). 384

# G.1.4 Interactive System for Structure-based ASCII Art Creation

Interactive structure-based systems [Xu et al., 2010] invite users to actively participate in the creation 386 of ASCII art; the interactive paradigm empowers users to foster a collaborative synergy between 387 human creativity and computational assistance. While it is an earlier paper, [Miyake et al., 2011] proposes to input images divided into grids for glyph matching using four metrics employed for 389 converting images into ASCII art: template matching which considers pixel positions for dissimilarity 390 measure, normalized cross-correlation which minimizes the influence of line width differences 391 using histograms, Histogram of Oriented Gradients (HOG) representing line directions, and distance 392 transformation indicating line positions. In comparison, our work focuses mainly on AI generation, so 393 while it is not as collaborative it focuses on the optimization and comparison of methods for comparing 394 the output of ASCII art utilizing either tone-based style, which is a detailed and comprehensive image, 395 or structure-based style, which is a simple outline of the image using less characters. 396

# **G.1.5** Perceptual Metrics for Image Quality Assessment

The recent success of perceptual messages based on deep neural networks in regards to the Image Quality Assessment (IQA) task [Kazmierczak et al., 2022] has led to a growing interest in new metrics that outperform previous metrics to develop perceptual information at different resolutions. Whereas the IQA metric is generally easily perceivable for humans, it is more difficult to set a metric for a computational algorithm. Our work investigates the model's abilities to generate accurate images by comparing it to Euclidean distance and the SSIM index, groundwork laid by [Chung and Kwon, 2022].

# 405 G.1.6 ASCII Representation Learning

While prior work has explored ASCII conversion from images [Matsumoto et al., 2018], little attention has been given to understanding how models internally represent ASCII structures. Our probing

of CLIP embeddings extends this line of inquiry, revealing architecture and how ASCII images are represented in the model.

# NeurIPS Paper Checklist

418

419

420

421

436

437

438

439

440

441

442

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

- The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count
- towards the page limit.

  Hease read the checklist guidelines carefully for information on how to answer these questions.
  - Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:
    - You should answer [Yes], [No], or [NA].
    - [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
    - Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. 426 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally 428 expensive" or "we were unable to find the license for the dataset we used"). In general, answering 429 "[No] " or "[NA] " is not grounds for rejection. While the questions are phrased in a binary way, we 430 acknowledge that the true answer is often more nuanced, so please just use your best judgment and 431 write a justification to elaborate. All supporting evidence can appear either in the main paper or the 432 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found. 434

### 435 IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- · Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

# 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state the contributions of our paper such as a benchmark to measure the quality of language-model-generated ASCII-text images and classification ability, a dataset of 20k images, a high-quality evaluation set of 320 images, and a fine-tuned CLIP model.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

# Answer: [Yes]

Justification: We discuss the limitations of our experiments, specifically the methods used. Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach.
   For example, a facial recognition algorithm may perform poorly when image resolution
   is low or images are taken in low lighting. Or a speech-to-text system might not be
   used reliably to provide closed captions for online lectures because it fails to handle
   technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Our paper does not present any theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Our paper discloses and describes in great detail the dataset creation process and model testing procedure. Additionally, our model weights and evaluation set are to be released for those to reproduce the main experimental results.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
  well by the reviewers: Making the paper reproducible is important, regardless of
  whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We link our github

### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
  proposed method and baselines. If only a subset of experiments are reproducible, they
  should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

566

567

568

569

570

571

573

574

575

576

577

578

579

580

581

582

583

584

585

586 587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607 608

609

610

611

612

613

615

Justification: Our paper specifies the model architecture, number of epochs, optimizer (e.g. Adam), full set of hyperparameters, which hyperparameters we chose and why.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail
  that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We report error bars for some figures, but do not report them in the classification accuracy table. The numbers in that table were rounded from the number we get by dividing the number of correct answers by the number of total answers.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Our paper explicitly states the type of compute used for different stages such as the free-tier CPU for noise removal and a paid-tier Google Colab A100 GPU for model training. Additionally in our paper we state how long it took for our training loops to execute.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our research involves creating a benchmark from publicly available ASCII art for model evaluation. Our dataset was filtered for inappropriate and harmful content, additionally our work does not involve human subjects, weapons research, or other prohibited activities, thus conforming to the NeurIPS Code of Ethics.

### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We specifically mention ArtPrompt, an ASCII art-based LLM jailbreak and how a dataset is needed to develop a solution in that field.

### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal
  impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.

• If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: Each ASCII image went through an extensive filtering pipeline, and was manually checked over by a human to prevent the release of unsafe images. Our data is licensed under CC BY NC 4.0, which permits only non-commercial use and is intended exclusively for research purposes.

# Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We explicitly credit the source, ascii.co.uk, in our dataset section. In the absence of a clear license, our use is strictly non-commercial and transformative for scholarly research.

# Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

# 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: All new datasets, model weights, and code are released with clear licenses and comprehensive documentation to ensure full reproducibility.

### Guidelines:

720

721

722

723

724

725

726

727

728

729

730

731

732 733

734

735

736

737

738

739

740

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

769

770

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: : Our research did not involve crowdsourcing or human subjects.

### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our research did not involve human subjects so IRB approval was not required.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent)
  may be required for any human subjects research. If you obtained IRB approval, you
  should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: In our research LLMs were utilized for benchmarking experiments, not a component of the core methodology used to conduct our research.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
  - Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.