
🎮ByteCraft: Generating video games and animations through bytes

Alexia Jolicoeur-Martineau
Samsung SAIL Montréal
alexia.j@samsung.com

Emy Gervais
Independent

Abstract

Generating new video games using AI has potential to be the next holy grail of the video game industry. Current AI efforts have focused on two directions: i) controllable video generation and ii) code generated by Large Language Models (LLMs). The first direction is limited due to short-term memory and increasing corruption (blur, noise) over time. The second direction is promising, but it requires a lot of human hand-holding with human-provided assets.

Generating hours of coherent interactive video content is infeasible. In this paper, we instead attempt the overly ambitious problem of end-to-end generation of Small Web Format (SWF) games and animations through bytes. By modeling bytes, one does not need code or assets to potentially obtain full games with title screen, narrative, text, graphics, music, and sounds.

We make a first attempt by fine-tuning a 7-billion-parameter LLM at 32K context length to generate the bytes of video games and animations conditional on a text description. Our model (ByteCraft) can generate up to 32K tokens, each containing at most 4-5 bytes (generating files as big as 140 KB). Some of the generated files are partially working (4.8-12%), or fully working (0.4-1.2%). ByteCraft is a proof-of-concept highlighting what could be possible given more scaling and engineering effort.

We open-source our model and inference code alongside a dataset of 10K synthetic prompts for use with ByteCraft.



Figure 1: Screenshots of files generated by ByteCraft.

1 Introduction

Current methods for video game generation. While a lot of research is done on generating videos (Runaway, 2023; OpenAI, 2024; DeepMind, 2024), very little effort is done on the video game front. Current state-of-the-art models for generating video games focus mainly on controllable video generation (Menapace et al., 2021; Yang et al., 2024b; Valevski et al., 2024; Che et al., 2024; Yu et al., 2025; Kanervisto et al., 2025), which effectively consists of walking simulators with a few seconds of memory and increasingly blurry images over time. More recently, Large Language Models (LLMs) (Todd et al., 2024; Hu et al., 2024; Anjum et al., 2024; Rosebud AI, 2024; X, 2025) have started being used to generate video games through code. This approach is promising, but generating complex games requires many rounds of human-AI interactions and user-provided graphics and audio assets.

A promising direction: directly generating bytes. A promising new line of direction is the generation of bytes from various types of files (e.g., images, videos, text, programs, etc.) using language models (Horton et al., 2023; Wu et al., 2024; Pérez et al., 2024; Han et al., 2024; Pérez et al., 2024). By staying in the byte world, one can generate any type of file found on a computer.

Why is byte generation not widespread? This direction has received very little attention for many reasons. First, simple files can take an enormous amount of bytes, leading to extremely large context lengths. For example, a simple 1Mb file requires 1 million byte tokens. Second, most modalities that people care about have specific neural network architecture and methods that have been developed and iterated over many years by researchers to make them excellent (e.g., images using 2D convolutions). Treating these modalities as bytes is challenging and may require many improvements in order to beat existing methods.

What are the problems associated with using bytes and how do we deal with them? The main difficulties when generating bytes of games and animations are 1) scaling (due to high-context length and limited data) and 2) overfitting. To handle longer context length than byte-level models, we tokenize bytes into 108K tokens, leading to approximately 2.29 bytes per token (with some extreme cases at 4-5 bytes per token). With a 32K sequence length, we can generate games of around 73 KB in size. To reduce the risks of overfitting, we produced multiple prompts per sequence of bytes.

The first-of-a-kind. In this work, we make the first attempt at building a generator of Small Web Format (SWF) (Systems, 2010) games and animations through bytes. SWF is a complex multi-modal format containing images, videos, code, fonts, and more data types. A single incorrectly placed byte could break the generated file. We are the first to tackle such a challenging task. We do so by fine-tuning an LLM (Qwen2.5-7B) to generate bytes conditional on a text prompt describing the game/animation, making our approach fully end-to-end. Our model, ByteCraft, was trained on limited resources (4 GPUs) for many months. It can generate up to 32K bytes (games/animations smaller than 73KB on average, with the largest we have seen at around 140KB).

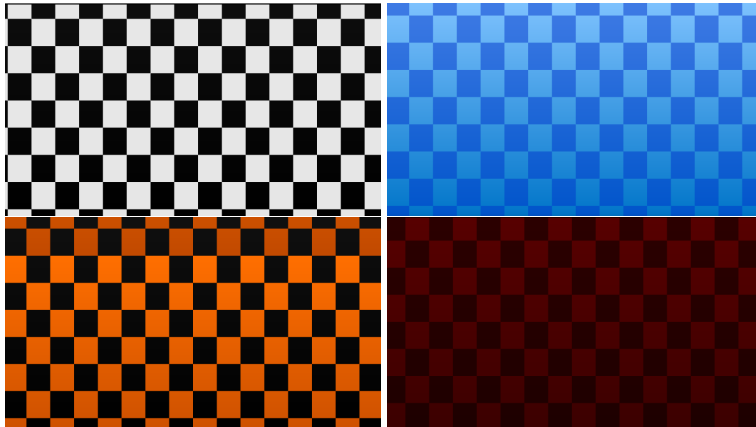


Figure 2: Checkered patterns in motion generated in different colors by ByteCraft.

2 ByteCraft

2.1 Architecture

We need to be able to condition on both text description and bytes. Since no model exists for generating bytes of video games and animations, this part has to be learned. LLMs are particularly good with unstructured text prompts.

We considered the following open-source LLMs as base model: Llama-3.1-8B (Dubey et al., 2024), Ministral-8B (Mistral AI, 2024), Qwen2-VL-7B (Wang et al., 2024), and Qwen2.5-7B (Yang et al., 2024a). In preliminary experiments, we found Qwen2.5-7B to be the best; thus, we used it as our base model.

2.2 Tokenization

We used Byte Pair Encoding (BPE) (Gage, 1994) to encode the bytes of video games and animations into 108K new tokens containing on average 2.29 bytes per token (going up to 4-5 bytes per token in rare cases). This allowed us to scale to larger games than would be possible with byte-level generation. These new tokens were added to the tokens of the pre-trained model.

2.3 Data augmentation on prompts

Using Qwen2.5-7B, we generated multiple prompts per sequence of bytes in order to produce more diversity and reduce risks of overfitting.

2.4 Training

We fine-tuned ByteCraft with progressively bigger context length (4K, 8K, 32K) using AdamW (Kingma, 2014; Loshchilov, 2017). For the early stages of training, we used the Muon optimizer (Jordan et al., 2024), which accelerated training. Modern techniques such as Fully Sharded Data Parallel(FSDP) (Zhao et al., 2023) and fused kernels (Hsu et al., 2024) were used to accelerate training and reduce memory cost. Training took around 4 months in total using 4 GPUs. The model was trained until it reached a cross-entropy loss of 0.15 (Perplexity (PPL) of 1.16).

2.5 Usage

The user provides a text prompt to describe the video game or animation that they want to be generated. It can be written in any language. We provide some prompt examples in Figure 3. Once the prompt is given, ByteCraft generates k SWF files per prompt, where k is set by the user. These files can then be opened with the Ruffle player (Ruffle, 2025).

Using vLLM v0.7.3 (Kwon et al., 2023), the model can simultaneously generate 25 files at 32K context length in around 10 minutes on a single A100 with 80Gb memory.

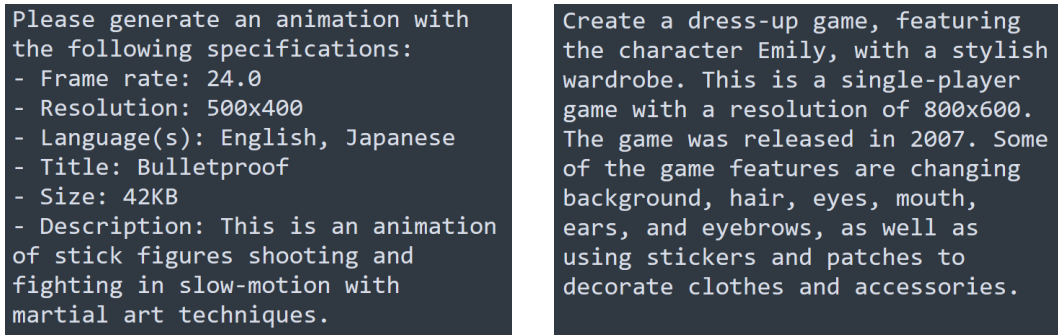


Figure 3: Examples of prompts (Left: Structured, Right: Unstructured)

3 Results

We tested ByteCraft qualitatively and quantitatively after generating 250 files from 250 prompts from a held-out set.

For quantitative measures, we 1) calculate the maximum Jaro-Winkler string similarity (Jaro, 1989) between the generated file and training files ("Max similarity"), 2) the Jaro-Winkler string similarity (Jaro, 1989) between the generated file and the true file with the associated prompt ("Truth similarity"), and 3) the percentage of "fully broken" files as determined by verifying its header and metadata; if parsing fails, the file is likely invalid.

Since this is the first attempt at generating such files, assessing the quality of the files is tricky. We thus rely on qualitative measures by manually opening each file using the Ruffle player (Ruffle, 2025) (Nightly 2025-03-14). We categorize games that are not "fully broken" as i) "blank canvas" (with a flat background color), ii) "stuck on a loading screen", iii) "showing/hearing something" (this is subjective; it means that something is shown or heard and it is not a loading screen or blank canvas), iv) "Fully Working" (the file is working and playable from a quick assessment).

Table 1: Results from 250 prompts.

Classification	$\min_p = 0.05$ $T = 0.1$	$\min_p = 0.05$ $T = 0.3$	$\min_p = 0.10$ $T = 0.3$	$\min_p = 0.20$ $T = 0.3$
Fully broken (wrong format)	20.4%	26.4%	20.4%	26.0%
Blank canvas (flat color)	69.6%	62.8%	67.2%	68.0%
Stuck on loading screen	4.4%	4.4%	4.0%	0.8%
Showing/hearing something	4.8%	5.6%	8.0%	4.0%
Fully working	0.4%	0.8%	0.4%	1.2%
Max similarity	0.790	0.790	0.789	0.791
Truth similarity	0.627	0.628	0.628	0.626

The results are shown in Table 1. We see that most games are broken, but a few show something visually/audibly interesting, and a tiny percentage functions properly. On average, around 21% of the bytes in the files are novel, as determined by the maximum Jaro-Winkler similarity. The results show that the model can learn some patterns about bytes enough to produce some semi-working or working files.

4 Potential future improvements

ByteCraft is a first attempt at generating open-web games and animations through bytes. However, it generates many broken files. In this section, we propose some improvements that could improve the performance of this method.

4.1 Scaling

ByteCraft was trained on extremely limited resources (4 GPUs over 4 months). Given large-scale resources, ByteCraft could improve in performance and be extended to generate larger games at higher context length.

4.2 Reinforcement learning

To reduce the number of broken files, one could i) generate many files from various prompts with the current model, ii) have an automatic visual evaluator (possibly using a visual LLM such as the latest Qwen2.5-VL (Bai et al., 2025)) from a screenshot giving a reward between 0 and 1 for broken to fully working, and iii) use reinforcement learning to push the model toward generating valid working files or more aligned with the prompt.

4.3 Test-time compute

Test-time scaling methods (Snell et al., 2024; Kim et al., 2025) have the potential to increase quality and diversity of the generated samples at inference without retraining.

4.4 Better generalization on small data using data augmentations on bytes

Data augmentation is currently applied only to prompts: we generate a different prompt for each sequence of bytes. However, to truly enable the model to understand and generalize given a limited amount of training examples, one could i) randomly permute the order in which file components are stored (e.g., image 1, font1, code1, image2 \rightarrow image2, image1, code1, font1). Other data augmentations are also possible, such as ii) having an AI rewrite the programming language code differently or iii) adding noise to the image assets contained in the file or changing their format (e.g., PNG (Roelofs, 2010) \rightarrow JPEG (Collins English Dictionary, 2013)). The difficulty with all the data augmentation strategies proposed is that they require a good understanding of the file format, as one would need to determine where each asset starts and ends with their format and possibly where they are referred to in the code contained in the file. Ideally, this could be inferred directly from bytes, but it is possible that this would instead require expensive reverse engineering.

5 Conclusion

We built ByteCraft, the first generator of games and animations through bytes. Contrary to existing approaches, ByteCraft is fully end-to-end. Given our limited resources (4 GPUs), it is limited to small files (32K context-length, around 73KB on average, but can be as big as 140KB). ByteCraft will require more scaling and could be extended to more types of games, including compressed versions of larger video games. This approach of generating files from text is not limited to video games and animations; it could be used to generate any kind of file end-to-end on a computer.

Parallel with early molecule generation. A parallel exists between ByteCraft and autoregressive molecule generation. Molecules can be represented as SMILES strings (Weininger, 1988) and their context length is generally small (around 10-250 tokens without BPE). We show below some of the progress of molecule generation over time on the Zinc-250K dataset (Sterling and Irwin, 2015):

1. GVAE (Gómez-Bombarelli et al., 2016): 0.7% valid molecules (\leftarrow ByteCraft is here)
2. CVAE (Kusner et al., 2017): 7.2% valid molecules
3. RVAE (Ma et al., 2018): 34.9% valid molecules
4. GFVAE (Ma and Zhang, 2021), STGG (Ahn et al., 2021), and many others: 100% valid molecules, but not necessarily realistic and synthesizable
5. STGG+AL (Jolicoeur-Martineau et al., 2025): 100% valid molecules with high synthesizability and out-of-distribution properties (\leftarrow a future ByteCraft version could be here)

In our case, we tackle the much harder problem of autoregressively generating SWF files with up to 32K tokens. Currently, some of the generated files are partially working (4.8-12%) or fully working (0.4-1.2%). We are thus at step 1, the equivalent of GVAE for molecule generation in 2016. Our end goal is generating 100% valid files with high novelty. We propose some potential directions of improvements in Section 4 on how to get there. We hope this crazy project inspires researchers and hobbyists toward the lofty goal of generating games through bytes.

References

- Sungsoo Ahn, Binghong Chen, Tianzhe Wang, and Le Song. Spanning tree-based graph generation for molecules. In *International Conference on Learning Representations*, 2021.
- Asad Anjum, Yuting Li, Noelle Law, Megan Charity, and Julian Togelius. The ink splotch effect: A case study on chatgpt as a co-creative game designer. In *Proceedings of the 19th International Conference on the Foundations of Digital Games*, pages 1–15, 2024.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibong Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin, and Hao Chen. Gamegen-x: Interactive open-world game video generation. *arXiv preprint arXiv:2411.00769*, 2024.
- Collins English Dictionary. Definition of JPEG, 2013. Archived from the original on 21 September 2013. Retrieved 23 May 2013.
- DeepMind. Veo 2: A High-Definition Generative Model for Video, 2024. URL <https://deepmind.google/technologies/veo/veo-2/>.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Philip Gage. A new algorithm for data compression. *The C Users Journal*, 12(2):23–38, 1994.
- Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *arXiv e-prints*, pages arXiv–1610, 2016.
- Xiaochuang Han, Marjan Ghazvininejad, Pang Wei Koh, and Yulia Tsvetkov. Jpeg-lm: Llms as image generators with canonical codec representations. *arXiv preprint arXiv:2408.08459*, 2024.
- Maxwell Horton, Sachin Mehta, Ali Farhadi, and Mohammad Rastegari. Bytes are all you need: Transformers operating directly on file bytes. *arXiv preprint arXiv:2306.00238*, 2023.
- Pin-Lun Hsu, Yun Dai, Vignesh Kothapalli, Qingquan Song, Shao Tang, Siyu Zhu, Steven Shimizu, Shivam Sahni, Haowen Ning, and Yanning Chen. Liger kernel: Efficient triton kernels for llm training. *arXiv preprint arXiv:2410.10989*, 2024. URL <https://arxiv.org/abs/2410.10989>.
- Chengpeng Hu, Yunlong Zhao, and Jialin Liu. Game generation via large language models. In *2024 IEEE Conference on Games (CoG)*, pages 1–4. IEEE, 2024.
- Matthew A Jaro. Advances in record-linkage methodology as applied to matching the 1985 census of tampa, florida. *Journal of the American Statistical association*, 84(406):414–420, 1989.
- Alexia Jolicœur-Martineau, Yan Zhang, Boris Knyazev, Aristide Baratin, and Cheng-Hao Liu. Generating π -functional molecules using stgg+ with active learning, 2025. URL <https://arxiv.org/abs/2502.14842>.
- Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://web.archive.org/web/20250122060345/https://kellerjordan.github.io/posts/muon/>.
- Anssi Kanervisto, Dave Bignell, Linda Yilin Wen, Martin Grayson, Raluca Georgescu, Sergio Valcarcel Macua, Shan Zheng Tan, Tabish Rashid, Tim Pearce, Yuhao Cao, et al. World and human action models towards gameplay ideation. *Nature*, 638(8051):656–663, 2025.
- Hyeonah Kim, Sanghyeok Choi, Jiwoo Son, Jinkyoo Park, and Changhyun Kwon. Neural genetic search in discrete spaces. *arXiv preprint arXiv:2502.10433*, 2025.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. In *International conference on machine learning*, pages 1945–1954. PMLR, 2017.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626, 2023.
- I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Changsheng Ma and Xiangliang Zhang. Gf-vae: a flow-based variational autoencoder for molecule generation. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pages 1181–1190, 2021.
- Tengfei Ma, Jie Chen, and Cao Xiao. Constrained generation of semantically valid graphs via regularizing variational autoencoders. *Advances in neural information processing systems*, 31, 2018.
- Willi Menapace, Stephane Lathuiliere, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci. Playable video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10061–10070, 2021.
- Mistral AI. Un ministral, des ministraux, 2024. URL <https://mistral.ai/news/ministraux/>.
- OpenAI. Sora: Creating video from text, 2024. URL <https://openai.com/sora>.
- Juan C Pérez, Alejandro Pardo, Mattia Soldan, Hani Itani, Juan Leon-Alcazar, and Bernard Ghanem. Compressed-language models for understanding compressed file formats: a jpeg exploration. *arXiv preprint arXiv:2405.17146*, 2024.
- Greg Roelofs. History of png. libpng, May 2010. Retrieved 20 October 2010.
- Rosebud AI. AI Game Creator | AI-Powered Game Dev Platform, 2024. URL <https://lab.rosebud.ai/ai-game-creator>.
- Ruffle. Ruffle - Flash Emulator, 2025. URL <https://ruffle.rs>.
- Runway. Runway Research | Gen-2: Generate novel videos with text, images or video clips, 2023. URL <https://runwayml.com/research/gen-2>.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Teague Sterling and John J Irwin. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling*, 55(11):2324–2337, 2015.
- Adobe Systems. *Adobe Flash Player Administration Guide for Flash Player 10.1*, June 2010. URL <https://www.adobe.com>. Archived from the original (PDF) on 2010-11-21. Retrieved 2011-03-10.
- Graham Todd, Alexander G Padula, Matthew Stephenson, Éric Piette, Dennis Soemers, and Julian Togelius. Gavel: Generating games via evolution and language models. *Advances in Neural Information Processing Systems*, 37:110723–110745, 2024.
- Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.

- Shangda Wu, Xu Tan, Zili Wang, Rui Wang, Xiaobing Li, and Maosong Sun. Beyond language models: Byte models are digital world simulators. *arXiv preprint arXiv:2402.19155*, 2024.
- X. Grok 3 Beta — The Age of Reasoning Agents | xAI, 2025. URL <https://x.ai/news/grok-3>.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024a.
- Mingyu Yang, Junyou Li, Zhongbin Fang, Sheng Chen, Yangbin Yu, Qiang Fu, Wei Yang, and Deheng Ye. Playable game generation. *arXiv preprint arXiv:2412.00887*, 2024b.
- Jiwen Yu, Yiran Qin, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Gamefactory: Creating new games with generative interactive videos. *arXiv preprint arXiv:2501.08325*, 2025.
- Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023.