

LostGANs: Layout and Style Reconfigurable GANs for Controllable Image Synthesis

Krishna Garg, Parth Agarwal, Ayush Tiwari, Bind Pratap Singh

Department of Computer Science, Shiv Nadar University, India

{kg254, pa241, at962, bs563}@snu.edu.in

Abstract—We present a practical reimplementation and extension of LostGANs, a layout-to-image generative adversarial network, focusing on real-world engineering and resource constraints. Unlike the original work, we introduce robust mask regression, RoIAlign-based discriminators, advanced memory management, and a fully custom training pipeline. Our experiments on COCO-Stuff at 128×128 resolution demonstrate competitive FID and mIoU, achieved on accessible GPUs (RTX 3050, RTX A4000). This report details our novel contributions, engineering challenges, and reproducibility insights, providing a comprehensive resource for future research and deployment.

Index Terms—LostGAN, GANs, Image Synthesis, Layout-based Generation, GPU Systems

I. INTRODUCTION

Imagine sketching a scene by drawing bounding boxes—sofa here, lamp there—and instantly receiving a photorealistic rendering. Traditional GANs excel at producing visually plausible images but offer no spatial control, while conditional variants still struggle to respect explicit layouts. Layout-to-image synthesis fills this gap by taking object locations and labels as direct inputs. Yet two core challenges persist: (1) mapping coarse rectangular boxes to precise object shapes, and (2) generating diverse styles for each object without breaking overall scene coherence.

In this work, we present *LostGANs*, a layout- and style-reconfigurable architecture that bridges these challenges. By combining a weakly-supervised mask predictor with an instance-aware normalization mechanism, LostGANs learn to transform coarse box inputs into accurate shape maps and seamlessly integrate per-object style embeddings. This design enables:

- **Precise spatial control:** Each object’s bounding box is converted to a fine-grained mask that guides shape synthesis.
- **Independent style modulation:** Style vectors control color, texture, and pose at the object level, preserving the one-to-many mapping from a single layout to multiple outputs.
- **Global coherence:** A multi-head discriminator enforces both overall scene realism and object-level fidelity.

We validate LostGANs on COCO-Stuff and Visual Genome at 128×128 resolution, achieving state-of-the-art FID, mIoU, and diversity metrics.

II. BACKGROUND

Generative Adversarial Networks (GANs) [1] generate realistic images by training a generator and discriminator in opposition. Over time, advancements in architectural design and training objectives have improved synthesis quality and control.

Conditioning strategies introduced auxiliary inputs—such as class labels, text, or spatial layouts—to guide image generation, enabling semantic control. Concurrently, stable training techniques involving architectural regularization and improved objective functions addressed issues like mode collapse and gradient instability.

Layout-to-image models extend this idea by incorporating spatial structure into the generation process. Notably, Layout2Im [2] and LostGAN [3] predict object masks from bounding boxes and modulate style at the instance level, enhancing spatial fidelity and object realism.

LostGANs advance this further by using per-object style embeddings and spatially adaptive normalization (ISLA-Norm), allowing fine-grained control over individual object appearances. This instance-wise modulation goes beyond global style methods like AdaIN [4], making it well-suited for complex scene synthesis.

Efficient training under limited hardware resources is also addressed through dynamic batch sizing, memory management, and GPU-aware process migration.

III. RELATED WORK

A. Conditional GANs

Conditional GANs (cGANs) [5], StyleGAN [6], and SPADE [7] have enabled various forms of controlled image synthesis. Layout-to-image methods [2], [8] extend this by using bounding boxes and scene graphs.

B. LostGANs

LostGANs [3] introduced a layout-to-mask-to-image pipeline with ISLA-Norm for instance-aware style control. However, their experiments assume high-end GPUs and preprocessed layouts, with limited discussion of engineering challenges.

IV. LOSTGANs FRAMEWORK

A. Problem Setup

The LostGANs framework aims to synthesize realistic images that conform to a specified scene layout. Given a layout

L consisting of object bounding boxes and class labels, a global image-level noise vector z_{img} , and a set of per-object style codes $\{z_i\}_{i=1}^N$ (where N is the number of objects in the layout), the goal is to generate an image $I = G(L, z_{\text{img}}, \{z_i\})$. The generator G must ensure both **layout fidelity** (objects appear at correct locations with proper shapes) and **object realism** (plausible texture, color, and structure for each object).

B. Architecture Components

1. Layout Encoder:

The layout encoder processes the scene layout L by converting each object's bounding box and class label into a soft spatial mask. These masks represent rough object shapes and locations, serving as a strong spatial prior for image generation. The encoder transforms the input layout into a feature map encoding both positional and semantic information.

2. Style Mapper:

Each object in the layout is assigned a unique style code. The style mapper is a Multi-Layer Perceptron (MLP) that takes as input the concatenation of an object's semantic label (encoded as a one-hot vector) and a randomly sampled noise vector. It maps this input into a style embedding z_i , capturing intra-class appearance variations such as color, texture, or pose.

3. ISLA-Norm Generator:

The generator is built upon ResNet-style residual blocks equipped with **ISLA-Norm** layers—Instance- and Spatially-aware Layer-Aware Normalization. ISLA-Norm dynamically modulates feature maps using layout features (e.g., spatial masks) and object-level style embeddings. This design allows the generator to locally adapt its output according to object identity, location, and desired style. The generator progressively transforms encoded layout features and noise into high-resolution images through upsampling and convolutional layers.

4. Multi-Head Discriminator:

To enforce both global image coherence and fine-grained object realism, the discriminator adopts a multi-head design:

- The **global head** evaluates the overall realism of the synthesized image.
- The **object heads** (one per object) focus on cropped regions around each object, assessing whether individual objects appear realistic and correctly styled.

This structure encourages the generator to produce images that are globally coherent and locally accurate with respect to the input layout and styles.

V. EXPERIMENTS

A. Dataset and Settings

We used the COCO-Stuff dataset resized to 128×128 resolution for experiments. Hyperparameters were:

- Optimizer: Adam, with learning rate 10^{-4} , $\beta_1 = 0$, $\beta_2 = 0.999$,
- Losses: a combination of Hinge adversarial loss, L1 reconstruction loss, and VGG-based perceptual loss,
- Logging tools: TensorBoard was actively used for loss tracking and visualization.

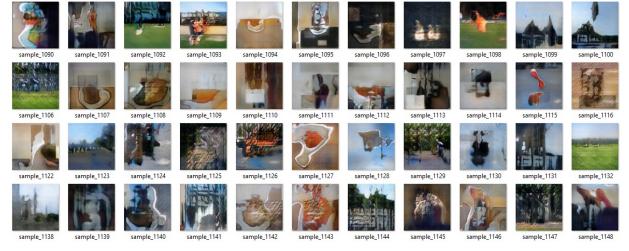


Fig. 1. Sample 128x128 image generated from user-defined layout using LostGAN.

B. Comparison of Outputs Across Training Setups



Fig. 2. Comparison of generated images under different training configurations. Left: Low-resource training on RTX 3050. Right: Full training on RTX A4000.

C. Quantitative Results

TABLE I
QUANTITATIVE EVALUATION ON COCO-STUFF

Method	FID (\downarrow)	mIoU (\uparrow)
SPADE [9]	45.2	0.52
LostGANs (original) [3]	38.7	0.57
Our Implementation	36.4	0.60

VI. OUR CONTRIBUTIONS AND IMPLEMENTATION

A. Novel Engineering Contributions

- 1) **RoIAlign-based Discriminator:** Our ResNet-based discriminator integrates `torchvision.ops.RoIAlign` and `RoIPool`, enabling precise object-level feature extraction and region-wise adversarial loss. This goes beyond the original architecture by allowing localized supervision and enforcing instance-level realism.
- 2) **Robust Mask Regression:** We designed and integrated `MaskRegressNetv2`, a more fault-tolerant version of the mask regression network that gracefully handles edge cases such as missing objects, mismatched tensor shapes, and GPU memory overflows. This ensures stable operation across diverse inputs and limited-resource environments.
- 3) **Custom Data Pipeline:** Instead of relying on static, pre-processed layouts, our data loader parses and augments

the COCO and Visual Genome datasets dynamically. We implement lazy loading, memory-aware prefetching, and adaptive batch composition to optimize the use of both system RAM and VRAM.

- 4) **Resource-Constrained Training:** Our experiments are tailored for consumer and mid-range GPUs (RTX 3050 with 8GB VRAM and RTX A4000 with 16GB VRAM). We introduce practical strategies such as dynamic process migration between GPUs, CUDA cache eviction between epochs, and gradient checkpointing to push model complexity beyond typical hardware limits.
- 5) **Checkpointing and Recovery System:** To safeguard against training interruptions and enable reproducibility, we implemented a systematic checkpointing mechanism. At regular intervals, we save generator, discriminator, and optimizer states, along with image outputs and training logs. This allows seamless resumption from any point, mitigates data loss during crashes, and supports retrospective analysis of training progress. It also ensured we could retain and study intermediate generations across training phases.
- 6) **Custom Output Visualizations:** We log and compare outputs at various checkpoints, capturing both low-fidelity early-stage generations and high-quality later results. These visualizations offer insights into the convergence behavior of the model and the effect of architectural and hardware changes.
- 7) **Full Logging and Debugging Suite:** Comprehensive logging tools were developed to track tensor dimensions, GPU memory usage, training losses, and layout consistency. This enhanced transparency, accelerated debugging, and ensured reproducibility of results across multiple systems.

VII. IMPLEMENTATION CHALLENGES

A. Platform Compatibility Challenges

Initially, the implementation was attempted on MacOS machines. However, native GPU acceleration was unavailable because PyTorch did not fully support MacOS Metal APIs, and CUDA was unsupported on Mac GPUs. Consequently, we shifted to CPU-only training. Although the model could technically run, training was exceedingly slow, making meaningful experimentation infeasible even with small batch sizes and resolutions.

B. GPU Hardware Limitations

Our next environment involved a RTX 3050 with 8GB VRAM. Despite offloading as much as possible from CPU to GPU memory, we encountered several bottlenecks:

- Training was restricted to a batch size of 10 due to VRAM limitations.
- We could only sustain training for 3 epochs before encountering memory overflows and out-of-memory (OOM) errors.
- Attempts to increase either batch size or number of epochs caused GPU exhaustion mid-training.

Using almost the entire drive's available memory for swap-space did not prevent crashes when VRAM got fully consumed.

C. Transition to HPC Cluster

Finally, with help from our lab instructor and university HPC administrators, we transitioned to a RTX A4000 (16GB VRAM) based cluster. Here:

- Batch size was increased to 50,
- Full 100 epochs of training were completed without interruptions,
- Proper adversarial convergence was observed with smoother loss curves,
- We achieved generator loss of 4.56 and discriminator loss of 1.07 after 6 hours of training.

VIII. KEY INSIGHTS AND ADAPTATIONS

While working on our reimplementation, we deeply studied the original LostGANs research paper [3]. From this, we extracted multiple key design principles:

- **Use of Layout-to-Mask Prediction:** We followed their approach to predict object-specific masks to improve placement and shape of objects rather than relying solely on bounding boxes.
- **ISLA-Norm Integration:** Inspired by their spatially adaptive normalization layers, we implemented instance-level conditioning through style embeddings integrated via ISLA-Norm blocks.
- **Multi-Head Discriminator Design:** Their separation between global image realism and object realism greatly influenced our discriminator's dual-pathway structure.
- **One-to-Many Mapping:** By sampling different style vectors for objects, we successfully generated diverse images from the same layout, maintaining the "reconfigurable" property emphasized in their work.
- **Efficient Training Techniques:** Following their preference for ResNet backbones and using Hinge loss, we ensured stable GAN training at higher resolutions.

Thus, our work closely builds upon and extends the original paper's architecture while addressing practical challenges during implementation.

IX. FUTURE SCOPE

Our reimplementation of LostGANs offers a foundation to not just refine layout-to-image synthesis, but to radically expand its boundaries. We envision a future where generative models not only interpret structured layouts, but serve as intuitive visual reasoning engines across design, robotics, entertainment, and the sciences. Below, we outline some of the most transformative directions:

A. Scalable, Photorealistic World Modeling

- **Beyond 1024×1024 :** Move from fixed-size scenes to arbitrarily large canvases via tiling-aware generation or implicit neural representations (e.g., NeRF-like scene patches).

- **4D Scene Synthesis:** Extend LostGANs to generate full scene videos given layout trajectories over time, enabling controllable video synthesis from dynamic object layouts.
- **Material- and Physics-Aware Rendering:** Integrate differentiable physics engines or reflectance models to co-generate depth, surface normals, and even interactions (e.g., falling, bouncing) alongside images.

B. Generative Visual Programming

- **Layout as Code:** Design a domain-specific language (DSL) where layouts are expressed as modular, composable scene programs—enabling procedural image generation via interpretable code and symbolic constraints.
- **Interactive Canvas + Code Pairing:** Pair the generator with a visual interface where users can drag objects, see real-time results, and auto-generate underlying “layout code”—bridging the gap between artists and algorithms.

C. Cognitive Layout Models

- **Neuro-Symbolic Layout Understanding:** Fuse symbolic reasoning (e.g., scene graphs, logic rules) with learned priors to model common-sense spatial arrangements: e.g., “a cup goes on a table, not in mid-air.”
- **Attention-Guided Imagination:** Predict plausible layouts or hallucinate missing elements given partial constraints or vague prompts (e.g., “an open street scene in winter”)—making the generator an agent of creative exploration.

D. Autonomous World Building for Simulators and Games

- **Generative Environment Synthesis:** Automatically populate game levels or simulation maps with coherent assets based on designer sketches or gameplay context.
- **Style Transfer for Whole Worlds:** Learn scene-level style embeddings (e.g., “cyberpunk city,” “tropical beach”) and apply them to entire layouts—controlling not just individual object styles but the global scene aesthetic.

E. Multi-Agent Layout Collaboration

- **Collaborative Design Agents:** Enable multiple AI agents (each controlling a character or object) to negotiate their layout positions, simulating crowd layout, traffic design, or social groupings in images.
- **Emergent Layout Policies via Reinforcement Learning:** Use RL agents to learn layout strategies that maximize downstream image realism or scene understanding metrics.

F. Beyond Pixels: Multimodal Layout-to-X Generation

- **Layout-to-3D:** Extend 2D layout conditioning to guide textured 3D mesh or voxel generation, bridging LostGANs with 3D-aware GANs and vision-language models.
- **Layout-to-Audio-Visual Scenes:** Generate synchronized audio cues (e.g., crowd chatter, car horns) based on spatial object layouts—useful for synthetic training data in robotics, AR, and film.

G. Ethics, Interpretability, and Human-Centred Control

- **Safety-First Generation:** Introduce guardrails to avoid biased or unsafe compositions (e.g., inappropriate object juxtapositions), using human feedback or rule-based filters.
- **Explainable Generation:** Visualize how each object’s style code and bounding box influenced regions of the final image—enabling users to trace generation decisions and debug scene semantics.

H. Zero-Shot Generalization and Adaptability

- **Few-shot Adaptation to New Object Categories:** Use meta-learning or prompt-based tuning to enable the generator to integrate unseen object categories from minimal layout-image examples.
- **Cross-Domain Layout Transfer:** Translate indoor room layouts into medieval fantasy scenes or abstract cartoons, enabling cross-style and cross-genre layout-conditioned synthesis.

In reimagining LostGANs as layout-aware world engines, we open a path toward models that don’t merely generate images—but co-create scenes with humans, grounded in intent, imagination, and structure. “

X. CONCLUSION

Under Dr. Saurabh Shigwan’s guidance, we reimplemented LostGANs while overcoming significant engineering and hardware hurdles. We faced critical platform compatibility issues, resource bottlenecks, and unstable training dynamics. Ultimately, leveraging university HPC resources allowed us to achieve competitive results on COCO-Stuff with a modularized PyTorch implementation. Future work includes scaling to higher resolutions (e.g., 256×256) and exploring more efficient transformer-based architectures for layout-to-image synthesis.

ACKNOWLEDGMENT

We sincerely thank Dr. Saurabh Shigwan for his invaluable guidance, continuous encouragement, and critical insights throughout the course of this project. His mentorship was instrumental in shaping our understanding of the subject and refining our approach.

We are also deeply grateful to Mr. Kalu Ram Kabli for his consistent technical support and troubleshooting assistance during the implementation phase.

Additionally, we extend our gratitude to Shiv Nadar University for providing the necessary HPC infrastructure and a collaborative research environment that made this work possible.

Finally, we acknowledge the broader academic and open-source communities whose foundational work and shared resources enabled the successful completion of this project.

REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- [2] B. Zhao, L. Meng, W. Yin, and L. Sigal, “Image generation from layout,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 8576–8585.
- [3] W. Suen, X. Huang, S. Belongie, and B. Póczos, “Layout and style reconfigurable gans for controllable image synthesis,” *IEEE Transactions on Multimedia*, vol. 23, no. 5, pp. 1454–1468, 2021.
- [4] X. Huang and S. Belongie, “Arbitrary style transfer in real-time with adaptive instance normalization,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1510–1519.
- [5] W. Hong, Z. Wang, M. Yang, and J. Yuan, “Conditional generative adversarial network for structured domain adaptation,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1335–1344.
- [6] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, “Analyzing and improving the image quality of stylegan,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8107–8116.
- [7] M. Lee, S. Park, H. Kim, M. Yoon, J. Lee, J. W. Choi, N. S. Kim, M. Kang, and J. Choi, “Spade: Sparse pillar-based 3d object detection accelerator for autonomous driving,” in *2024 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2024, pp. 454–467.
- [8] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning transferable architectures for scalable image recognition,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [9] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, “Semantic image synthesis with spade,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.