# The 1st International Workshop on Disentangled Representation Learning for Controllable Generation (DRL4Real): Methods and Results

Qiuyu Chen[†]  Xin Jin[†,*]  Yue Song[†]  Xihui Liu[†]  Shuai Yang[†]  Tao Yang[†]
Ziqiang Li[†]  Jianguo Huang[†]  Yuntao Wei[†]  Ba'ao Xie[†]  Nicu Sebe[†]  Wenjun (Kevin) Zeng[†]

Jooyeol Yun  Davide Abati  Mohamed Omran  Jaegul Choo  Amir Habibian  Auke Wiggers
Masato Kobayashi  Ning Ding  Toru Tamaki  Marzieh Gheisari  Auguste Genovesio  Yuheng Chen
Dingkun Liu  Xinyao Yang  Xinping Xu  Baicheng Chen  Dongrui Wu  Junhao Geng
Lexiang Lv  Jianxin Lin  Hanzhe Liang  Jie Zhou  Xuanxin Chen  Jinbao Wang
Can Gao  Zhangyi Wang  Zongze Li  Bihan Wen  Yixin Gao  Xiaohan Pan
Xin Li  Zhibo Chen  Baorui Peng  Zhongming Chen  Haoran Jin

*Abstract*—*This paper reviews the 1st International Workshop on Disentangled Representation Learning for Controllable Generation (DRL4Real), held in conjunction with ICCV 2025. The workshop aimed to bridge the gap between the theoretical promise of Disentangled Representation Learning (DRL) and its application in realistic scenarios, moving beyond synthetic benchmarks. DRL4Real focused on evaluating DRL methods in practical applications such as controllable generation, exploring advancements in model robustness, interpretability, and generalization. The workshop accepted 9 papers covering a broad range of topics, including the integration of novel inductive biases (e.g., language), the application of diffusion models to DRL, 3D-aware disentanglement, and the expansion of DRL into specialized domains like autonomous driving and EEG analysis. This summary details the workshop's objectives, the themes of the accepted papers, and provides an overview of the methodologies proposed by the authors.*

## I. INTRODUCTION

Disentangled Representation Learning (DRL) is a critical area of research aimed at enabling AI systems to decompose observed data into underlying, interpretable factors of variation. By decoupling complex entities into independent latent factors, DRL holds the potential to address fundamental challenges in AI, such as enhancing the controllability and interpretability of generative systems and improving model generalization.

Despite significant academic interest and progress in DRL methodologies, primarily based on Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs), the field has largely remained confined to synthetic datasets. The transition to realistic scenarios has been hindered by the complex nature of real-world data and the lack of robust benchmarks and unified evaluation metrics. Traditional DRL methods often struggle when faced with the weaker inductive biases present in real-world environments.

The ICCV 2025 DRL4Real Workshop was organized to address this gap. It aimed to foster the development of novel,

realistic datasets and comprehensive benchmarks for evaluating DRL methods in practical applications. The workshop encouraged submissions exploring how DRL can advance model capabilities, with a focus on key areas including controllable generation. This summary provides an overview of the workshop and the contributions of the accepted papers.

## II. WORKSHOP OVERVIEW

The DRL4Real workshop aimed to achieve two primary goals: (1) to provide a comprehensive review of recent developments in applying DRL to realistic scenarios, and (2) to serve as a forum for researchers to explore the challenges and opportunities in controllable generation using disentangled representations.

The workshop attracted diverse submissions spanning various modalities and applications. Following a rigorous review process, 9 papers were accepted for presentation.

## III. WORKSHOP PAPERS AND THEMES

The accepted papers showcased several innovative techniques and highlighted emerging trends in the application of DRL to realistic challenges. We summarize the principal themes observed across the submissions:

1) **DRL for Precise Controllable Generation and Editing.** A major focus was on leveraging DRL to achieve fine-grained control in generation and editing tasks. This included using pre-trained DRL models to extract semantic priors that explicitly constrain edits [8] and modeling spatial reasoning for plausible object placement [1].

2) **Leveraging Diffusion Models and Novel Inductive Biases.** Recognizing the limitations of traditional regularization, many authors integrated DRL with Diffusion Probabilistic Models (DPMs) and introduced novel inductive biases. This included using textual semantics as a regularization prior [5], leveraging inherent diffusion properties like time-varying bottlenecks [3], and incorporating structural prompts [7].

3) **3D-Aware and Sequential Disentanglement.** Several papers addressed the challenge of disentangling factors in complex spatial and temporal data. Methods explored 3D-aware generation for autonomous driving [9], diffusion-based video factorization [3], reducing static bias in action recognition [2], and semantic isolation theory for 3D anomaly detection [6].

4) **Expanding DRL to Specialized Domains.** The workshop demonstrated the broadening scope of DRL into specialized, realistic domains. Contributions included applications in autonomous driving [9], 3D anomaly detection [6], and EEG analysis for Brain-Computer Interfaces [4].

5) **Foundation Models for Compact Representations.** The interplay between large foundation models and representation learning was also explored, notably in using multimodal LLMs (GPT-4o) to generate high-fidelity images from highly compact textual representations for compression [7].

## IV. ACCEPTED PAPERS

This section details the methodologies and contributions of the papers accepted at the DRL4Real workshop, organized thematically.

### A. DRL for Controllable Image Generation and Editing

*1) A Guided Fine-tuning Framework for Diffusion Models with Disentangled Semantic Priors for Multi-Factor Image Editing [8]:* This paper addresses the challenge of unintended alterations in complex, multi-factor image editing using diffusion models. They propose a Guided Fine-tuning Framework that incorporates disentangled semantic priors as structural constraints.

**Description.** The framework (Figure 1) introduces a dual-conditioning approach. While text prompts guide what to change, a disentangled semantic prior guides what to preserve.

1) **Semantic Prior Extraction:** A pre-trained disentanglement model (EncDiff) is used as a Semantic Encoder ($\tau_\phi$) to extract a set of independent semantic concept tokens (S) from the input image.

2) **Guidance Adapter:** A lightweight, trainable adapter module is integrated into a pre-trained editing model (InfEdit). This adapter fuses the text prompt (P) and the semantic prior (S) using Multi-Head Cross-Attention, where $E_p$ serves as Query and $E_s$ serves as Key/Value.

The resulting structurally-aware prompts ($P^{refined}$) guide the editing model more precisely.

**Implementation Details.** The framework is validated on the 'DRL for Real' competition Multi-Factor Track dataset. By fine-tuning only the adapter module, the method reduces computational costs. Evaluations using CLIP Score, LPIPS, and FID demonstrated that the framework significantly reduces unintended attribute changes while maximizing desired edits.
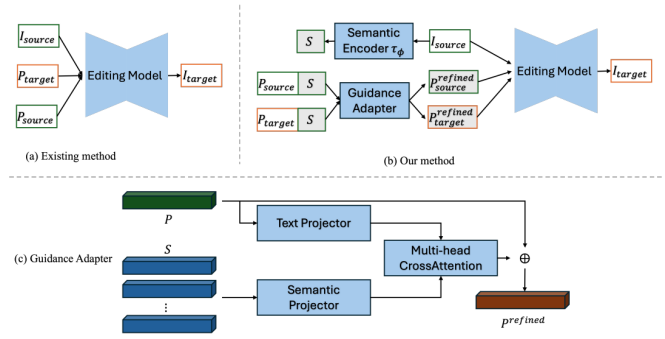


Fig. 1. Proposed Guided Fine-tuning Framework. (a) Existing methods. (b) Our method using Semantic Encoder and Guidance Adapter. (c) Detailed architecture of the Guidance Adapter. (Source: Paper 9 [8])

*2) Textual Semantics Matters: Unsupervised Representation Disentanglement in Realistic Scenarios with Language Inductive Bias (TA-Dis) [5]:* This paper proposes Text-Aided Disentanglement (TA-Dis), a framework that leverages the inherently disentangled nature of textual semantics to regularize DRL in the visual domain, addressing the limitations of purely visual constraints in realistic scenarios.

**Description.** TA-Dis is built upon the Latent Diffusion Model (LDM) in a two-stage approach (Figure 2). This novel use of language as a regularizer builds upon recent trends in unsupervised disentanglement that leverage diffusion model properties [14], [15], sparse transformations [16], and graph-based reasoning with large language models [17], all aiming for more robust and meaningful latent representations [18]. The first stage establishes a semantic projector $P_{sem}$ to obtain a semantic code $Z_{sem}$ with primary disentanglement capability. The core innovation is the second stage, introducing Text-Aided Regularization based on CLIP scores to further enhance the disentanglement capability of $Z_{sem}$. Three language inductive biases are designed: Image-Text Alignment (using $\mathcal{L}_{pull}$ and $\mathcal{L}_{push}$), Cycle Consistency (Order Loss $\mathcal{L}_o$), and Shift Equivariance (Exchangeable Loss $\mathcal{L}_e$). These losses regularize the disentanglement process by enforcing constraints in the image-text space.
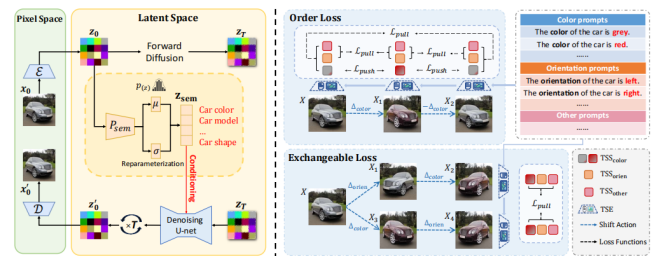


Fig. 2. Pipeline of the TA-Dis framework. (a) LDM backbone with a semantic projector $P_{sem}$. (b) Text-aided regularization via Order Loss and Exchangeable Loss. (Source: Paper 5 [5])

**Implementation Details.** The framework utilizes a pre-trained VAE and U-Net, optimizing the semantic projector

using AdamW. Experiments demonstrated superior disentanglement (measured by TAD) compared to VAE- and other Diffusion-based DRL methods on realistic datasets.

*3) Imagining the Unseen: Generative Location Modeling for Object Placement [1]:* This paper tackles the problem of location modeling—determining plausible locations for non-existing objects in a scene. They propose a generative approach to handle the inherent ambiguity and data sparsity of the task.

**Description.** The authors reframe location modeling as a generative task, $P(Y|X, C)$, using an autoregressive transformer (Figure 3, left). The input image and target object class condition the model, which sequentially generates bounding box coordinates. To utilize negative annotations (implausible locations), they incorporate Direct Preference Optimization (DPO), as shown in Figure 3 (right). DPO fine-tunes the model by maximizing the likelihood that positive locations ($Y^+$) are preferred over negative locations ($Y^-$) based on the Bradley-Terry model.
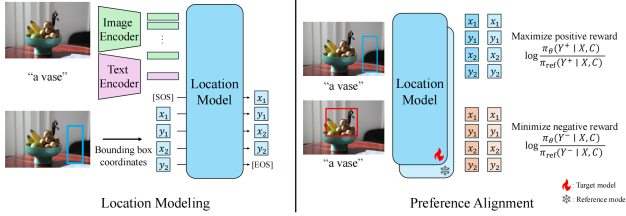


Fig. 3. Overview of the generative location model. Left: The autoregressive transformer generates bounding box coordinates. Right: Direct Preference Optimization (DPO) refines the model by aligning with positive and negative location preferences. (Source: Paper 1 [1])

**Implementation Details.** The model uses a small GPT-2 architecture, pretrained on the PIPE dataset and fine-tuned on OPA. The generative model achieved superior placement accuracy on OPA and improved visual coherence in object insertion tasks compared to instruction-tuned editing methods.

### B. Sequential and 3D Disentanglement

*1) DiViD: Disentangled Video Diffusion for Static-Dynamic Factorization [3]:* This work introduces **DiViD**, the first end-to-end video diffusion framework designed explicitly for static-dynamic factorization, aiming to overcome the information leakage common in VAE/GAN approaches.

**Description.** DiViD incorporates several key inductive biases within a DDPM framework (Figure 4). The sequence encoder employs an **Architectural Bias** by extracting the static token (s) from the first frame ($f_1$) and dynamic tokens ($d_i$) from the residuals ($f_i - f_1$), explicitly removing static content from the motion code. The decoder utilizes **Diffusion-driven Inductive Biases**: a Time-Varying Information Bottleneck inherent to the diffusion process, and Cross-Attention Interaction in the U-Net to route global static and local dynamic information appropriately.
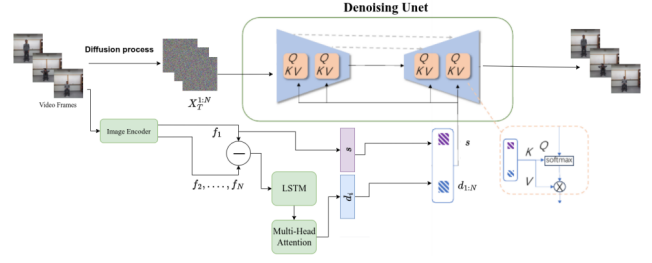


Fig. 4. Overview of DiViD. The sequence encoder uses residual encoding (subtracting $f_1$) to separate static (s) and dynamic ($d_{1:N}$) tokens, which condition the Denoising U-Net. (Source: Paper 3 [3])

**Implementation Details.** DiViD is trained end-to-end using the standard DDPM loss augmented by an orthogonality regularization term between static and dynamic tokens. Evaluations on MHAD and MEAD showed superior performance compared to SOTA methods.

*2) Disentangling Static and Dynamic Information for Reducing Static Bias in Action Recognition [2]:* This paper addresses the problem of static bias in action recognition, where models rely excessively on static cues rather than dynamic motion.

**Description.** A two-stream architecture is proposed to separate unbiased features ($f_u$, dynamics) from biased features ($f_b$, static cues), as illustrated in Figure 5. The biased stream is designed to be inaccessible to temporal information (either via architecture or input manipulation). Disentanglement is achieved through two mechanisms: (1) Statistical Independence Loss ($\mathcal{L}_{ind}$) using the Hilbert-Schmidt Independence Criterion (HSIC) to minimize dependence between $f_u$ and $f_b$. (2) Adversarial Scene Prediction Loss ($\mathcal{L}_S$), using a Gradient Reversal Layer (GRL) on the unbiased stream to force it to fail at scene prediction, thereby removing background information from $f_u$.

**Implementation Details.** Experiments on datasets emphasizing temporal information demonstrated that the method effectively reduces static bias metrics.

*3) Controllable Generation with Disentangled Representative Learning of Multiple Perspectives in Autonomous Driving [9]:* This paper presents a framework for controllable multi-view image generation in autonomous driving scenarios, focusing on disentangling key semantic factors from the scene representation and viewpoint.

**Description.** The proposed method employs a structured latent representation that decomposes the generative process into three factors: scene content ($z_{sc}$), weather ($z_w$), and speed ($z_s$). These codes are derived via variational encoding of semantic labels. The architecture utilizes a triplane-based 3D generator conditioned on these latent codes. The triplane representation efficiently maps the latent codes into a compact 3D scene representation. A factor-aware decoder then estimates color and density ($c, \sigma$) for sampled 3D locations, which are synthesized into 2D views via differentiable volumetric rendering (similar to NeRF). This formulation allows independent control over weather, motion, and viewpoint.

(a) with extractor-based biased stream
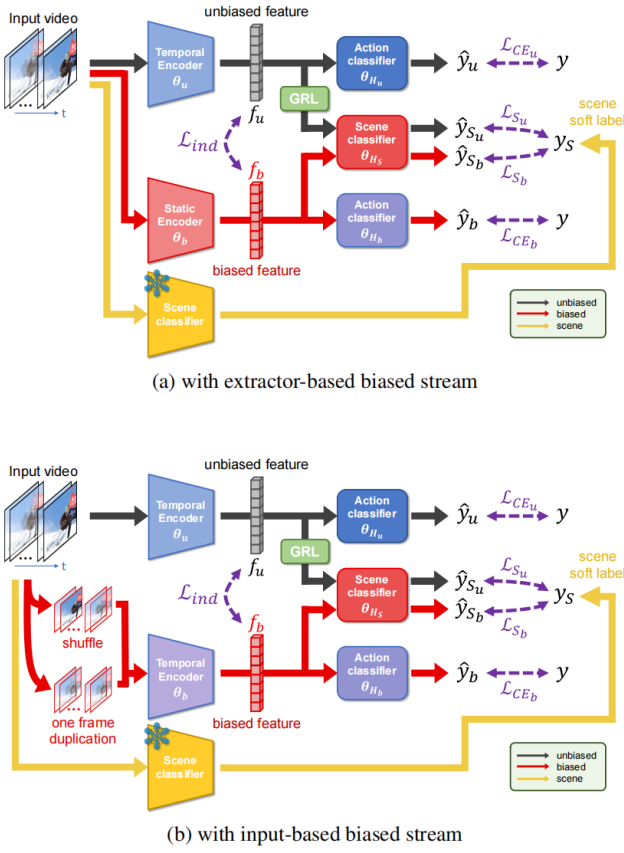


(b) with input-based biased stream

Fig. 5. The proposed two-stream architecture for disentanglement. An unbiased stream processes temporal dynamics, while a biased stream focuses on static cues. Disentanglement is enforced via an independence loss ($\mathcal{L}_{ind}$) and an adversarial scene prediction loss. (Source: Paper 2 [2])

**Implementation Details.** The model is trained with a hybrid objective including reconstruction loss, semantic consistency loss across latent-modified samples, and volumetric rendering-based regularization. The approach was validated on a custom multi-view driving dataset, demonstrating high-fidelity, semantically controllable view synthesis, including view completion.

*4) Fence Theorem: Towards Dual-Objective Semantic-Structure Isolation in Preprocessing Phase for 3D Anomaly Detection [6]:* This paper addresses the lack of a unified theoretical foundation for preprocessing design in 3D Anomaly Detection (AD). It establishes the **Fence Theorem** and proposes Patch3D as an implementation.

**Description.** The Fence Theorem formalizes preprocessing as a dual-objective semantic isolator: (1) mitigating cross-semantic interference, and (2) confining anomaly judgments to aligned semantic spaces to establish intra-semantic comparability. The theorem posits that preprocessing aims to divide the point cloud into mutually non-interfering (orthogonal) semantic spaces. Guided by this theorem, the authors implement **Patch3D** (Figure 6):

1) **Patch-Cutting:** Uses FPS and K-Means to segment a single point cloud into multiple independent semantic spaces based on its structure.

2) **Patch-Matching:** Merges similar semantic spaces across different point clouds to align their meanings.

3) **Separation Modeling:** Independently models normal features within each aligned space for anomaly detection.
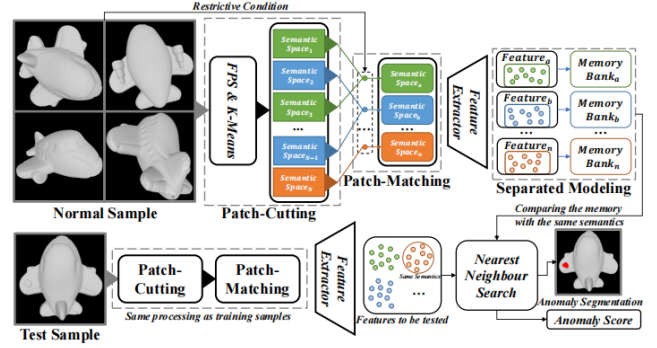


Fig. 6. Pipeline of Patch3D, illustrating the implementation of the Fence Theorem via Patch-Cutting, Patch-Matching, and Separation Modeling. (Source: Paper 6 [6])

**Implementation Details.** The feature extractor utilizes Fast Point Feature Histogram (FPFH) features. Experiments on Anomaly-ShapeNet and Real3D-AD demonstrated that the progressively finer-grained semantic alignment achieved by Patch3D enhances point-level anomaly detection accuracy, validating the Fence Theorem. The proposed theorem and its implementation in Patch3D align with a broader research direction focusing on novel reconstruction techniques [20] and the perception of internal spatial modalities [19] to advance the state-of-the-art in 3D anomaly detection.

*C. DRL in Specialized Domains and Modalities*

*1) FusionGen: Feature Fusion-Based Few-Shot EEG Data Generation [4]:* This paper addresses data scarcity in EEG-based BCIs by proposing **FusionGen**, a data generation framework based on disentangled representation learning and feature fusion for few-shot scenarios.

**Description.** FusionGen utilizes a U-Net-shaped encoder-decoder architecture (Figure 7). The core innovation is the **Feature Matching Fusion** module. This module integrates features across different samples (source and target) in the latent semantic space by randomly sampling target embeddings and replacing them with the most similar source feature. This injects diversity while preserving semantics.

**Implementation Details.** The network is trained as a denoising autoencoder minimizing MSE. Experiments on three public Motor Imagery (MI) and one steady-state visual evoked (SSVEP) EEG datasets showed that FusionGen significantly outperforms existing augmentation techniques in few-shot scenarios. This approach complements other recent advancements in the EEG domain, such as pipelines for high-quality data construction [11], foundation models for classification [12], and techniques for cross-headset distribution alignment [10], alongside related work in multi-graph adversarial networks [13].
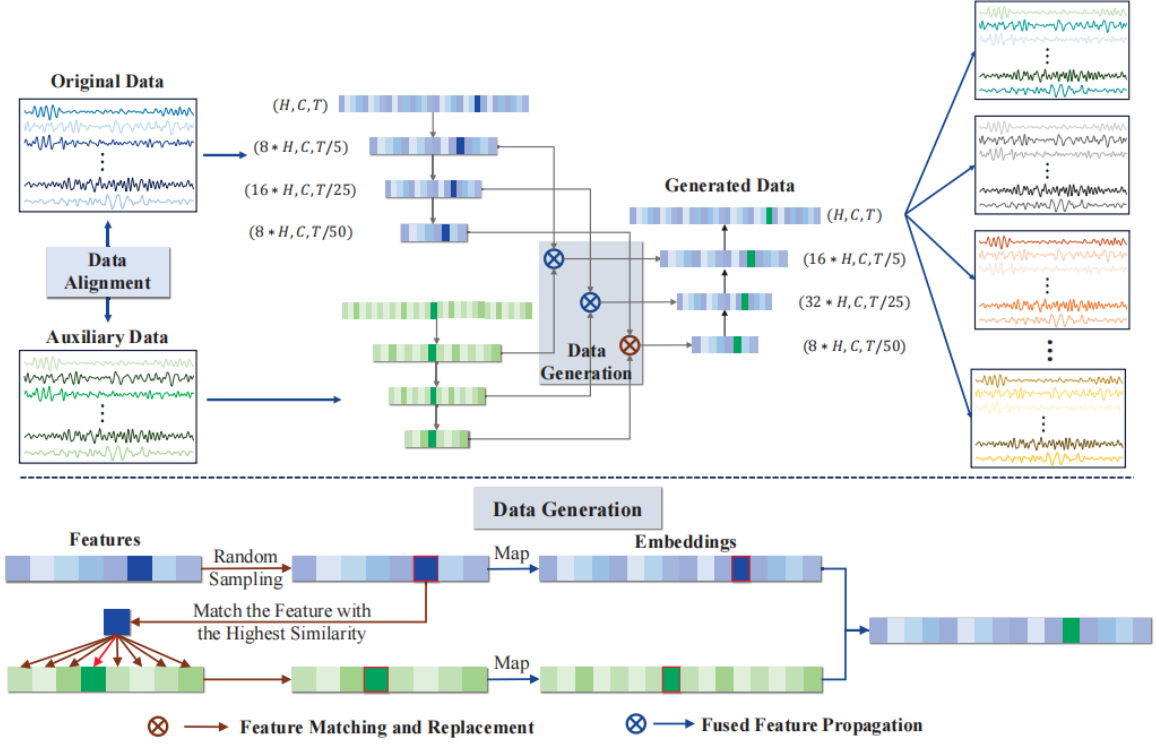
Fig. 7. Architecture of FusionGen. Aligned raw and auxiliary trials are encoded, fused via the Feature Matching and Replacement mechanism in the latent space, and propagated through the decoder to generate EEG trials. (Source: Paper 4 [4])

## D. Generative Models for Compact Representations

*1) Why Compress What You Can Generate? When GPT-4o Generation Ushers in Image Compression Fields [7]:* This work explores the potential of large foundation models (GPT-4o) for ultra-low bitrate image compression, advocating for generating pixels rather than compressing them.

**Description.** The framework (Figure 8) decouples the image signal into textual descriptions and an optional low-resolution visual prior. The key challenge is maintaining consistency during generation. To address this, the authors introduce **structural raster-scan prompt engineering**. This method involves designing a prompt that instructs the MLLM (GPT-4o) to describe visual elements in a specific order (top-to-bottom, left-to-right), explicitly preserving the spatial arrangement, object identities, and stylistic properties. The textual description is losslessly compressed and transmitted. If used, the visual condition (downsampled and compressed using MS-ILLM codec) provides base structural information. GPT-4o then reconstructs the image guided by these inputs without any additional fine-tuning.

**Implementation Details.** Text compression uses Lempel-Ziv coding. Evaluations on the DIV2K validation set showed that the method achieves competitive performance compared to recent generative compression approaches at ultra-low bitrates (e.g., 0.001 bpp).
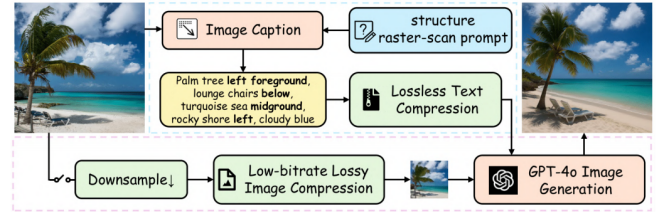


Fig. 8. Overall pipeline of the multimodal image compression framework based on GPT-4o image generation. (Source: Paper 8 [7])

## V. CONCLUSION

The 1st DRL4Real Workshop highlighted significant progress in moving Disentangled Representation Learning from synthetic benchmarks to realistic applications. The accepted papers demonstrated a clear trend towards leveraging powerful generative architectures, particularly diffusion models, and incorporating novel inductive biases, such as language and explicit structural priors, to achieve controllable generation. Furthermore, the expansion of DRL principles into specialized domains like autonomous driving and EEG analysis underscores the growing impact of the field. The innovations presented at the workshop provide valuable insights and pave the way for more robust, interpretable, and controllable AI systems.

# REFERENCES

[1] J. Yun, D. Abati, M. Omran, J. Choo, A. Habibian, and A. Wiggers, "Imagining the unseen: Generative location modeling for object placement," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[2] M. Kobayashi, N. Ding, and T. Tamaki, "Disentangling static and dynamic information for reducing static bias in action recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[3] M. Gheisari and A. Genovesio, "DiViD: Disentangled video diffusion for static–dynamic factorization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[4] Y. Chen, D. Liu, X. Yang, X. Xu, B. Chen, and D. Wu, "FusionGen: Feature fusion-based few-shot EEG data generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[5] J. Geng, L. Lv, and J. Lin, "Textual semantics matters: Unsupervised representation disentanglement in realistic scenarios with language inductive bias," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[6] H. Liang, J. Zhou, X. Chen, J. Wang, and C. Gao, "Fence theorem: Towards dual-objective semantic-structure isolation in preprocessing phase for 3d anomaly detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[7] Y. Gao, X. Pan, X. Li, and Z. Chen, "Why compress what you can generate? when GPT-4o generation ushers in image compression fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[8] B. Peng and Z. Chen, "A guided fine-tuning framework for diffusion models with disentangled semantic priors for multi-factor image editing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[9] H. Jin, "Controllable generation with disentangled representative learning of multiple perspectives in autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2025.

[10] D. Liu, S. Li, Z. Wang, W. Li, and D. Wu, "Spatial distillation based distribution alignment (SDDA) for cross-headset EEG classification," *arXiv preprint arXiv:2503.05349*, 2025.

[11] D. Liu, Z. Chen, and D. Wu, "CLEAN-MI: A scalable and efficient pipeline for constructing high-quality neurodata in motor imagery paradigm," *arXiv preprint arXiv:2506.11830*, 2025.

[12] D. Liu, Z. Chen, J. Luo, S. Lian, and D. Wu, "MIRepNet: A pipeline and foundation model for EEG-based motor imagery classification," *arXiv preprint arXiv:2507.20254*, 2025.

[13] D. Liu, H. Zhou, Y. Qu, H. Zhang, and Y. Xu, "UMMAN: Unsupervised multi-graph merge adversarial network for disease prediction based on intestinal flora," *IEEE Transactions on Computational Biology and Bioinformatics*, 2025.

[14] X. Jin, B. Li, B. Xie, W. Zhang, J. Liu, Z. Li, T. Yang, and W. Zeng, "Closed-loop unsupervised representation disentanglement with $\beta$-vae distillation and diffusion probabilistic feedback," in *European Conference on Computer Vision*. Springer, 2024, pp. 270–289.

[15] T. Yang, C. Lan, Y. Lu, and N. Zheng, "Diffusion model with cross attention as an inductive bias for disentanglement," *Advances in Neural Information Processing Systems*, vol. 37, pp. 82 465–82 492, 2024.

[16] Y. Song, T. A. Keller, Y. Yue, P. Perona, and M. Welling, "Unsupervised representation learning from sparse transformation analysis," *arXiv preprint arXiv:2410.05564*, 2024.

[17] B. Xie, Q. Chen, Y. Wang, Z. Zhang, X. Jin, and W. Zeng, "Graph-based unsupervised disentangled representation learning via multimodal large language models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 103 101–103 130, 2024.

[18] J. Liu, R. Feng, Y. Qi, Q. Chen, Z. Chen, W. Zeng, and X. Jin, "Rate-distortion-cognition controllable versatile neural image compression," in *European Conference on Computer Vision*. Springer, 2024, pp. 329–348.

[19] H. Liang, G. Xie, C. Hou, B. Wang, C. Gao, and J. Wang, "Look inside for more: Internal spatial modality perception for 3d anomaly detection," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 5, pp. 5146–5154, Apr. 2025. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/32546

[20] H. Liang, J. Zhang, T. Dai, L. Shen, J. Wang, and C. Gao, "Taming anomalies with down-up sampling networks: Group center preserving reconstruction for 3d anomaly detection," *arXiv preprint arXiv:2507.03903*, 2025.