

E-ISSN: 2583-7141

## International Journal of Scientific Research in Technology & Management



# Visual Effects (VFX) Using Deep Learning: A Comprehensive Review

Arun Pratap Singh

Dept. of Computer Science & Engineering Samrat Ashok Technological Institute, Vidisha, Madhya Pradesh, India singhprataparun@gmail.com

Abstract— Visual effects (VFX) have evolved into a crucial component of modern entertainment, enabling filmmakers, game developers, and content creators to achieve visuals that transcend physical constraints. Traditional pipelines, grounded in computer graphics and manual artistry, often demand extensive effort and resources. The emergence of deep learning has introduced a shift, allowing for data-driven automation, photorealistic rendering, and intelligent scene manipulation. Deep learning models such as convolutional neural networks (CNNs), generative adversarial networks (GANs), transformers, diffusion models, and neural radiance fields (NeRFs) have reshaped workflows in areas such as object removal, background replacement, motion capture, super-resolution, style transfer, and text-to-video synthesis. This paper provides a comprehensive review of deep learning in VFX, consolidating advances in architectures, datasets, evaluation methods, and real-world applications. Key challenges—such as temporal consistency, computational overhead, dataset scarcity, and ethical concernsare analyzed, while emerging research directions including multimodal control, efficient generative modeling, and real-time deployment are highlighted.

Keywords— Deep learning, visual effects, VFX, generative adversarial networks, transformers, diffusion models, neural rendering, NeRF, computer vision.

#### I. Introduction

Visual effects (VFX) refer to the integration of computergenerated imagery (CGI) with live-action footage to create visual content that would otherwise be impractical or impossible to capture on camera. Over the past decades, VFX has grown into a multi-billion-dollar industry spanning cinema, gaming, advertising, and immersive media such as augmented reality (AR) and virtual reality (VR). Traditional VFX pipelines rely on manual rotoscoping, chroma keying, physics-based simulations, and handcrafted modeling, which, while effective, are resource-intensive and timeSanjay Kumar Sharma

Dept. of Computer Science & Engineering Technocrats Institute of Technology Excellence, Bhopal, Madhya Pradesh, India sanjaysharmaemail@gmail.com

consuming [1]. Deep learning has revolutionized this landscape by automating core processes through data-driven methods. For instance, CNNs facilitate accurate foregroundbackground segmentation [2], GANs produce photorealistic textures and synthetic imagery [3], and transformers enhance temporal coherence across video sequences [4]. Diffusion models, which iteratively refine images through stochastic denoising, are emerging as state-of-the-art for generating high-quality and diverse imagery [5]. The motivation for this paper is to provide researchers and practitioners with a unified review of how deep learning has transformed VFX pipelines. We emphasize the relationship between foundational computer vision tasks and their direct integration into creative workflows. This includes not only technical contributions but also practical adoption in industries such as cinema (e.g., The Irishman for de-aging [6]) and streaming media. Furthermore, we examine the societal implications of AI-driven VFX, such as deepfakes and intellectual property challenges, highlighting the necessity of ethical frameworks [7].



Fig.1. VFX Example [7]

#### II. RELATED WORKS

#### A. Traditional VFX Pipelines

Conventional VFX pipelines relied heavily on a combination of computer graphics, procedural modeling, and manual compositing techniques. Tools such as Autodesk Maya, 3ds Max, and Adobe After Effects became industry standards for generating visual effects, enabling the design of particle systems, fluid dynamics, and physically based simulations [8]. Techniques like rotoscoping and chroma keying (green/blue screen compositing) were central to isolating subjects and integrating them into synthetic environments. However, while effective, these methods demanded significant manual labor, often requiring frameby-frame annotation by skilled artists. Moreover, scalability was limited, as every new project required custom tuning of parameters, and rendering times were often prohibitively expensive [9]. Despite incremental improvements in hardware acceleration and procedural workflows, traditional pipelines lacked adaptability to real-world variability such as motion blur, occlusion, and dynamic lighting, creating bottlenecks for large-scale productions [10].

#### B. Transition to Machine Learning

Prior to the deep learning revolution, machine learning approaches began to augment traditional VFX workflows. Early applications included optical flow estimation for motion tracking [11], active contours (snakes) for object boundary refinement, and handcrafted descriptors such as SIFT and SURF for feature matching and scene reconstruction [12]. These methods offered a degree of automation but were brittle under complex scenarios, particularly in the presence of heavy occlusions, illumination changes, or non-rigid object motion. Probabilistic graphical models such as Markov Random Fields (MRFs) and Conditional Random Fields (CRFs) were also used for segmentation and matting, but their reliance on manually engineered features limited scalability [13]. By the late 2000s, shallow learning methods, including support vector machines (SVMs) and random forests, were introduced for tasks like scene classification and face tracking, leading to modest improvements in semiautomated VFX [14]. However, these models lacked the representational power to capture high-dimensional visual patterns, and their performance degraded on large-scale, uncurated datasets. These limitations paved the way for the adoption of deep learning approaches, which demonstrated superior generalization by leveraging massive data and hierarchical feature extraction [15].

### C. Rise of Deep Learning in VFX

The introduction of AlexNet in 2012 [16] marked a watershed moment for computer vision, demonstrating the power of deep convolutional networks on large-scale datasets like ImageNet. This success rapidly influenced VFX pipelines, particularly in automating tasks such as semantic segmentation, matte extraction, and style transfer. Soon after, Generative Adversarial Networks (GANs), introduced by Goodfellow et al. in 2014 [17], transformed visual content generation by producing photorealistic

textures and completing missing regions in images. This was especially significant for object removal, scene completion, and super-resolution, all core components of VFX workflows [18]. In parallel, advances in 3D deep learning extended applications beyond static imagery. Works on volumetric CNNs [19], point cloud networks [20], and later Neural Radiance Fields (NeRFs) [21] enabled the reconstruction of 3D geometry and dynamic character animation from minimal inputs such as sparse images or monocular video. These methods facilitated realistic virtual cinematography, relighting, and camera re-projection, key requirements for modern VFX. Recent studies further consolidated the transformative role of deep learning in creative media. For example, Berthelot et al. [22] applied GAN-based adversarial training for video frame synthesis, producing temporally coherent imagery, while Chan et al. [23] demonstrated deep networks for facial reenactment, enabling realistic lip-syncing and facial animation. More recently, diffusion models [24], transformers [25], and CNN-transformer architectures hybrid [26] significantly improved temporal consistency, semantic accuracy, and generation fidelity, further bridging the gap between AI-driven automation and artist-directed creativity. Collectively, these milestones illustrate the growing synergy between computer vision and the creative industries, where deep learning is no longer an auxiliary tool but a central driver of scalability, photorealism, and accessibility in VFX pipelines.

#### III. DEEP LEARNING APPROACHES IN VFX

#### A. Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) form the backbone of most computer vision-driven VFX applications due to their ability to hierarchically extract spatial features. In the context of VFX, CNNs have been widely applied to semantic segmentation, object recognition, depth estimation, and background subtraction, tasks that are crucial for automating matte extraction and scene compositing [15]. Architectures such as U-Net [16], originally designed for biomedical image segmentation, have been successfully adapted for VFX to perform pixel-level separation of foreground and background, thus replacing the laborintensive process of rotoscoping. Similarly, Mask R-CNN [17] extended region-based CNNs to instance segmentation, enabling object-level manipulations such as selective removal, replacement, or augmentation in VFX workflows. Recent works integrate CNNs with temporal consistency modules for video effects, where maintaining coherence across frames is essential for preventing visual artifacts [18]. Despite their effectiveness, CNN-based methods struggle with capturing global dependencies and handling large occlusions, which has motivated the adoption of more advanced architectures such as GANs, transformers, and diffusion models.

#### B. Generative Adversarial Networks (GANs)

Generative Adversarial Networks (GANs) have emerged as a powerful tool in VFX for texture synthesis, style transfer, inpainting, and resolution enhancement. GANs operate on a generator-discriminator paradigm, where the generator creates realistic outputs and the discriminator enforces perceptual fidelity [19]. Conditional GANs (cGANs) enable controlled editing, where outputs can be guided by semantic labels, sketches, or text prompts, greatly expanding creative flexibility [20]. StyleGAN [21] has become a benchmark in high-resolution facial generation, producing photorealistic avatars and enabling advanced techniques such as digital deaging, facial reenactment, and identity-preserving synthesis for film and television. GAN-based super-resolution models, including ESRGAN [22], have been employed for legacy film restoration, where low-quality archival footage is enhanced to modern viewing standards. More recent developments incorporate temporal GANs for video frame interpolation and cycle-consistent GANs for style adaptation across cinematic domains, such as converting live-action sequences into stylized animations [23]. While GANs produce sharp and detailed textures, challenges such as mode collapse, training instability, and lack of diversity persist, motivating hybridization with diffusion or transformer-based frameworks.

#### C. Transformers and Attention Mechanisms

Transformers, originally designed for sequential data in natural language processing, have been adapted for vision tasks through the Vision Transformer (ViT) architecture [24]. In VFX, transformers offer the advantage of capturing long-range dependencies across both spatial and temporal domains, enabling applications in scene reconstruction, motion tracking, and temporal consistency for video editing [25]. For example, Video Transformers [26] extend this framework to sequential frames, providing robust handling of occlusions and maintaining global coherence across complex scenes. Transformers have also been integrated into hybrid CNN-transformer pipelines, leveraging CNNs for local detail extraction and attention mechanisms for global semantic reasoning [27]. These models are particularly impactful in video inpainting, motion retargeting, and neural style transfer for film editing, where maintaining frame-to-frame coherence is critical to avoiding flicker. Recent innovations also involve temporal multimodal transformers, which align video, text, and audio modalities to support creative control in VFX—such as generating effects guided by natural language prompts [28].

#### D. Diffusion Models

Diffusion probabilistic models have recently emerged as a dominant paradigm in generative modeling, surpassing GANs in terms of sample diversity and photorealistic quality. These models iteratively transform noise into structured images through a denoising process, which has proven especially effective in video inpainting, creative content generation, and photorealistic synthesis [29]. Denoising Diffusion Probabilistic Models (DDPMs) [30] and their accelerated variants, such as Denoising Diffusion Implicit Models (DDIMs) [31], allow the generation of

high-fidelity, temporally consistent visual content, which is critical in film post-production. Applications in VFX include object removal, scene completion, relighting, and cinematic style transfer, where diffusion-based methods outperform GANs in maintaining global semantic coherence [32]. Moreover, diffusion models are inherently more stable during training and less prone to mode collapse. Their integration into interactive editing workflows has further accelerated adoption, allowing artists to generate multiple variations of effects and select the most appropriate outcome for a given scene. Despite their computational cost due to iterative sampling, recent advances in distilled diffusion models [33] have made them increasingly practical for production environments requiring real-time previews.

#### E. Neural Radiance Fields (NeRFs)

Neural Radiance Fields (NeRFs) represent a breakthrough in reconstructing 3D scenes from sparse 2D inputs, making them invaluable in VFX for virtual cinematography, relighting, and scene re-rendering [34]. By learning a volumetric representation of a scene, NeRFs allow artists to generate novel views and simulate realistic camera motions without requiring full 3D asset modeling. This has transformative applications in film production, where it can drastically reduce the need for physical set construction or complex motion capture systems [35]. Extensions such as Dynamic NeRFs [36] enable the modeling of non-rigid and time-varying scenes, making it possible to capture human motion, facial expressions, and environmental changes directly from monocular video. Moreover, recent works have combined NeRFs with diffusion and transformer-based pipelines to allow text-guided 3D scene editing [37], bridging the gap between creative intent and technical execution. While NeRFs are computationally intensive and face challenges with real-time rendering, advances such as Instant-NGP [38] have significantly accelerated inference, making NeRFs increasingly viable for integration into production-grade VFX pipelines.

#### IV. DATASETS AND BENCHMARKS

The success of deep learning in VFX is closely tied to the availability of large-scale, domain-specific datasets and standardized benchmarks. These resources provide diverse visual content, enabling models to generalize across tasks such as object segmentation, motion tracking, facial reenactment, texture synthesis, and 3D scene reconstruction. Moreover, benchmarks ensure reproducibility and fair evaluation, which are critical for comparing competing methods in research and production environments. COCO (Common Objects in Context). The COCO dataset [25] contains over 330,000 images with pixel-level annotations for object detection, instance segmentation, and keypoint estimation. In VFX, COCO-trained models are frequently used to automate foreground-background separation, prop manipulation, and scene composition. Its diverse object categories also make it suitable for training models used in cinematic object replacement or augmentation. ImageNet. ImageNet [26], a large-scale dataset with over 14 million labeled images across 21,000 categories, serves as a

foundation for pretraining deep networks. Transfer learning from ImageNet-trained models significantly accelerates convergence and improves generalization for VFX tasks, ranging from style transfer to scene classification and domain adaptation in visual storytelling. DAVIS (Densely Annotated VIdeo Segmentation). The DAVIS benchmark [27] provides densely annotated video sequences for evaluating video object segmentation, a task central to VFX where maintaining temporal consistency across frames is vital. Methods trained on DAVIS have been used to develop tools for automated rotoscoping, dynamic background removal, and temporal-aware video inpainting, all of which reduce manual effort in post-production. CelebA-HQ. CelebA-HQ [28] is a high-resolution facial image dataset employed for generative modeling, reenactment, and digital humans. VFX applications include de-aging actors, identity-preserving facial synthesis, facial replacement, and digital stunt doubles. GAN-based and diffusion-based models pretrained on CelebA-HQ have already been integrated into commercial post-production pipelines. ScanNet and ShapeNet. For 3D scene understanding, ScanNet provides richly annotated RGB-D video data of indoor environments, while ShapeNet offers a repository of over 50,000 3D models across various categories [29]. These datasets underpin tasks like neural rendering, volumetric reconstruction, and asset generation, which are increasingly crucial in virtual production, digital set design, and hybrid CGI-live action integration. LAION-5B and Objaverse. More recently, multimodal datasets such as LAION-5B and Objaverse [30] have enabled the training of text-to-image and text-to-3D generative models. These resources align natural language descriptions with visual and 3D assets, empowering novel workflows where VFX artists can generate scenes or characters through text prompts. Objaverse, in particular, provides millions of 3D assets that enhance training of NeRFs and diffusion models, bridging the gap between procedural asset design and datadriven content creation. Together, these datasets and benchmarks provide the training corpus and evaluation standards necessary for advancing deep learning in VFX. While early datasets like ImageNet and COCO enabled foundational breakthroughs, the shift toward video, 3D, and multimodal corpora reflects the evolving demands of production-grade VFX pipelines.

#### A. Evaluation Metrics

Evaluation in VFX must balance fidelity, perceptual realism, temporal consistency, and computational efficiency. Unlike purely scientific imaging tasks, VFX is judged not only on accuracy but also on aesthetic quality and viewer experience. Pixel-level similarity. Traditional measures such as Peak Signal-to-Noise Ratio (PSNR) and the Structural Similarity Index (SSIM) [31] quantify reconstruction fidelity against ground truth. While widely reported, these metrics often fail to capture perceptual realism, since visually plausible but different reconstructions may score poorly. Perceptual realism. To align better with human vision, Learned Perceptual Image Patch Similarity (LPIPS)

[32] compares deep feature embeddings extracted from pretrained networks, correlating strongly with subjective human judgments of image similarity. This is particularly relevant in texture synthesis, facial generation, and background replacement. Generative quality. For evaluating realism at a distributional level, metrics such as Fréchet Inception Distance (FID) and Kernel Inception Distance (KID) [33] are widely adopted. These metrics assess how closely the distribution of generated images matches that of real-world datasets, making them critical for GAN and diffusion-based methods in VFX. Video consistency. Temporal coherence is essential for avoiding artifacts like flicker in moving sequences. Metrics such as temporal optical flow (tOF) and temporal LPIPS (tLPIPS) [34] measure how smoothly generated content evolves across frames, ensuring consistent object appearance and background continuity in long video segments. Efficiency. Since VFX workflows often process large volumes of data, practical evaluation must include computational cost metrics such as floating-point operations per second (FLOPs), runtime latency, and GPU/CPU memory usage [35]. These determine whether models are suitable for real-time production (e.g., virtual sets) or offline rendering (e.g., film post-production). Finally, human subjective studies remain indispensable, as perceptual quality in VFX is inherently tied to viewer experience. Many studios conduct side-byside user evaluations or blind preference tests to validate realism in applied contexts.

### B. Applications in VFX

Deep learning is transforming nearly every stage of the visual effects production pipeline, providing automation, scalability, and novel creative tools. Object removal and inpainting. GANs and diffusion models automate contentaware fill, enabling seamless removal of props, wires, or even entire characters while maintaining background consistency [36]. Background replacement. CNN-based semantic segmentation models allow green-screen-free compositing, reducing the need for controlled chroma setups and enabling virtual production workflows [37]. Character animation. Deep learning-based pose estimation and motion transfer techniques provide markerless motion capture [38], lowering costs and enabling realistic digital stunt doubles and crowd animation. Facial synthesis. GAN-based models support de-aging, face-swapping, lip-sync generation, and digital doubles [39], increasingly used in cinema to create photorealistic performances without extensive manual postprocessing. Super-resolution and restoration. Deep learningdriven super-resolution models enhance the visual quality of legacy films, restoring archival footage and adapting content for 4K and beyond [40]. Text-to-video generation. Multimodal models such as text-conditioned diffusion and transformers allow script-driven previsualization, enabling directors to rapidly generate scene drafts directly from textual descriptions [41]. Collectively, these applications demonstrate how AI not only automates tedious tasks but also introduces new creative possibilities, reshaping both production efficiency and artistic freedom in VFX.

#### V. CHALLENGES AND LIMITATIONS

Despite rapid progress, deep learning in VFX faces several technical, practical, and ethical challenges. Dataset scarcity. Training high-capacity models demands large, domainspecific datasets, yet annotated cinematic-quality video data is limited [42]. While synthetic data can supplement training, it may not capture the complexity of real-world scenes. Computational expense. State-of-the-art GANs, transformers, and diffusion models are resource-intensive, requiring powerful GPUs and long training times [43]. This limits adoption in smaller studios and real-time environments. Generalization limits. Models trained on specific datasets often fail to generalize to new environments, lighting conditions, or artistic styles [44], requiring domain adaptation and fine-tuning for production. Temporal consistency. Maintaining seamless realism across thousands of frames is still difficult. Even state-of-the-art models suffer from flicker and inconsistent object persistence, a major barrier for film-quality output. Integration with traditional pipelines. Deep learning tools must integrate with existing VFX software (e.g., Autodesk Maya, Houdini, Nuke). Lack of standardized interfaces slows adoption in industry. Ethical risks. The same methods enabling de-aging or face replacement can be misused for deepfakes and misinformation [45]. Establishing ethical guidelines and detection mechanisms is critical to safeguard responsible use. In sum, while deep learning offers powerful advances for VFX, scalability, robustness, and ethics remain active areas of research before widespread, frictionless industry adoption becomes possible.

#### VI. CONCLUSION & FUTURE SCOPE

Deep learning has emerged as a transformative force in VFX, fundamentally reshaping how visual narratives are crafted, from object removal and facial synthesis to scene reconstruction and text-driven generation. By automating labor-intensive tasks, enabling novel creative workflows, and bridging the gap between technical complexity and artistic vision, AI-driven models have introduced unprecedented scalability and realism into visual effects pipelines. Nonetheless, challenges persist, including dataset scarcity, computational expense, temporal consistency, and ethical concerns. Yet, the field's rapid progress indicates that these obstacles are likely to be addressed through a combination of efficient architectures, multimodal integration, and stronger evaluation protocols. convergence of artistry and artificial intelligence promises not only to streamline workflows but also to unlock new forms of storytelling, making immersive, high-quality VFX more accessible than ever before. The trajectory of deep learning in VFX suggests several promising avenues for innovation. One critical direction is multimodal content creation, where models integrate text, sketches, audio, or exemplar images to provide artists with more intuitive and controllable editing tools [46]. Early work in text-to-image generation has already demonstrated this capability, and its extension to full video pipelines could enable script-driven or storyboard-based previsualization. Parallel to this,

advances in efficient transformers and distilled diffusion models are anticipated to reduce computational overhead, thereby supporting real-time applications such as live virtual production and interactive editing [47]. A second research front involves self-supervised and unsupervised learning, which can significantly reduce reliance on labeled datasets [48]. Since high-quality annotated video data for cinematic effects is scarce, leveraging unlabeled content for representation learning could accelerate the development of scalable models for production environments. Another frontier lies in Neural Radiance Fields (NeRFs) and their extensions to dynamic, time-varying scenes [49]. While NeRFs currently excel at static reconstructions, extending them to dynamic scenarios would allow 3D-consistent camera re-lighting, and immersive scene manipulation directly from 2D inputs—ushering in new paradigms for virtual cinematography. Finally, as deep learning-based VFX tools become mainstream, the community must address ethical, legal, and creative challenges. Issues such as misuse of generative models for deepfakes, lack of transparency in AI-assisted artistry, and dataset bias must be mitigated through standardized benchmarks, fairness guidelines, and regulatory frameworks [50]. This balance between innovation and responsibility will be critical to ensuring that the transformative potential of AI in VFX benefits both industry professionals and society at large.

#### REFERENCES

- [1] Lin, T. Y., et al. (2014). Microsoft COCO: Common Objects in Context. ECCV.
- [2] Deng, J., et al. (2009). ImageNet: A Large-Scale Hierarchical Image Database. CVPR.
- [3] Perazzi, F., et al. (2016). A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation. CVPR.
- [4] Karras, T., et al. (2018). Progressive Growing of GANs for Improved Quality, Stability, and Variation. ICLR.
- [5] Chang, A. X., et al. (2015). ShapeNet: An Information-Rich 3D Model Repository. arXiv:1512.03012.
- [6] Deitke, M., et al. (2023). Objaverse: A Universe of Annotated 3D Objects. CVPR.
- [7] Berthelot, D., et al. (2017). BEGAN: Boundary Equilibrium Generative Adversarial Networks. arXiv:1703.10717.
- [8] Autodesk Maya, Adobe After Effects (Software). Industry-standard VFX and compositing tools.
- [9] Horn, B. K., & Schunck, B. G. (1981). Determining Optical Flow. Artificial Intelligence.
- [10] Krizhevsky, A., Sutskever, I., & Hinton, G. (2012). ImageNet Classification with Deep Convolutional Neural Networks (AlexNet). NeurIPS
- [11] Goodfellow, I., et al. (2014). Generative Adversarial Nets. NeurIPS.
- [12] Wu, Z., et al. (2015). 3D ShapeNets: A Deep Representation for Volumetric Shapes. CVPR.
- [13] Berthelot, D., et al. (2017). Unsupervised Learning for Image Synthesis Using GANs. arXiv.
- [14] Chan, C., et al. (2019). Everybody Dance Now. ICCV.
- [15] Long, J., Shelhamer, E., & Darrell, T. (2015). Fully Convolutional Networks for Semantic Segmentation. CVPR.
- [16] Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image Segmentation. MICCAI.
- [17] He, K., et al. (2017). Mask R-CNN. ICCV.

- [18] Karras, T., et al. (2019). StyleGAN: A Style-Based Generator Architecture for Generative Adversarial Networks. CVPR.
- [19] Wang, X., et al. (2018). ESRGAN: Enhanced Super-Resolution Generative Adversarial Networks. ECCV Workshops.
- [20] Dosovitskiy, A., et al. (2020). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale (ViT). ICLR.
- [21] Arnab, A., et al. (2021). ViViT: A Video Vision Transformer. ICCV.
- [22] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. NeurIPS.
- [23] Song, J., et al. (2021). Denoising Diffusion Implicit Models. ICLR.
- [24] Mildenhall, B., et al. (2020). NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. ECCV.
- [25] Lin, T. Y., et al. (2014). Microsoft COCO: Common Objects in Context. ECCV.
- [26] Deng, J., et al. (2009). ImageNet: A Large-Scale Hierarchical Image Database. CVPR.
- [27] Perazzi, F., et al. (2016). DAVIS: A Benchmark Dataset and Evaluation Methodology for Video Object Segmentation. CVPR.
- [28] Karras, T., et al. (2018). Progressive Growing of GANs for Improved Quality, Stability, and Variation. ICLR.
- [29] Chang, A. X., et al. (2015). ShapeNet: An Information-Rich 3D Model Repository. arXiv.
- [30] Deitke, M., et al. (2023). Objaverse: A Universe of Annotated 3D Objects. CVPR.
- [31] Wang, Z., et al. (2004). Image Quality Assessment: From Error Visibility to Structural Similarity (SSIM). IEEE Transactions on Image Processing.
- [32] Zhang, R., et al. (2018). The Unreasonable Effectiveness of Deep Features as a Perceptual Metric (LPIPS). CVPR.
- [33] Heusel, M., et al. (2017). GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium (FID). NeurIPS.
- [34] Zhou, T., et al. (2018). Temporal Consistency Metrics for Video Prediction. ECCV.

- [35] Howard, A., et al. (2017). MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. arXiv:1704.04861.
- [36] Yu, J., et al. (2019). Free-Form Image Inpainting with Gated Convolution. ICCV.
- [37] Chen, L. C., et al. (2018). DeepLab: Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs. TPAMI.
- [38] Cao, Z., et al. (2017). Realtime Multi-Person 2D Pose Estimation Using Part Affinity Fields. CVPR.
- [39] Thies, J., et al. (2016). Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. CVPR.
- [40] Ledig, C., et al. (2017). Photo-Realistic Single Image Super-Resolution Using a GAN (SRGAN). CVPR.
- [41] Singer, A., et al. (2022). Make-A-Video: Text-to-Video Generation without Text-Video Data. arXiv:2209.14792.
- [42] Torralba, A., & Efros, A. A. (2011). Unbiased Look at Dataset Bias. CVPR.
- [43] Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and Policy Considerations for Deep Learning in NLP. ACL.
- [44] Geirhos, R., et al. (2020). Shortcut Learning in Deep Neural Networks. Nature Machine Intelligence.
- [45] Chesney, R., & Citron, D. (2019). Deep Fakes: A Looming Challenge for Privacy, Democracy, and National Security. California Law Review.
- [46] Ramesh, A., et al. (2022). Hierarchical Text-Conditional Image Generation with CLIP Latents (DALL·E 2). arXiv:2204.06125.
- [47] Meng, C., et al. (2023). Distillation of Diffusion Models for Fast Sampling. ICLR.
- [48] Chen, T., et al. (2020). A Simple Framework for Contrastive Learning of Visual Representations (SimCLR). ICML.
- [49] Pumarola, A., et al. (2021). D-NeRF: Neural Radiance Fields for Dynamic Scenes. CVPR.
- [50] Floridi, L., et al. (2018). AI4People—An Ethical Framework for a Good AI Society. Minds and Machines.