

Dynamic Scene Modeling and Rendering: A Survey of Methods and Applications

Yuwei Zhao^{1,2}, Kaiyuan Zhang^{1,2}, Yuxiang Liu^{1,2}, Yilin Zhang^{1,2}, and Keqin Zhang^{1,2}

Abstract—Dynamic scene modeling and reconstruction from video streams have witnessed a paradigm shift with the advent of Neural Radiance Fields (NeRF) and 3D Gaussian Splatting (3DGS). While NeRF-based methods offer photorealistic quality via implicit neural representations, they often suffer from prohibitive training and rendering costs. Recently, 3DGS has emerged as a powerful alternative, enabling real-time rendering through explicit splatting techniques, though it introduces new challenges regarding storage overhead and temporal consistency. In this paper, we present a comprehensive survey of video-based 3D reconstruction methods, categorizing them into implicit and explicit approaches. We systematically evaluate their trade-offs in terms of rendering quality, speed, and memory consumption across diverse scenarios. Furthermore, we summarize the limitations of current state-of-the-art methods and discuss future directions. Our compiled resources and code are available at <https://github.com/8barbatos/Academic-English-Group-Paper>.

Index Terms—3D Reconstruction, Dynamic Scene Modeling, Neural Radiance Fields (NeRF), Gaussian Splatting, Novel View Synthesis, Survey.

I. INTRODUCTION

WITH the rapid advancement of virtual reality (VR), augmented reality (AR), and the Metaverse, the demand for photorealistic 3D content creation has surged exponentially. Dynamic scene modeling, which aims to reconstruct 3D geometry and appearance from 2D video streams, serves as a fundamental technology for these applications, enabling immersive telepresence, free-viewpoint video, and digital human avatars [1]. Unlike static scene reconstruction, modeling dynamic scenes from video presents a highly ill-posed inverse problem due to the entanglement of object motion, topology changes, and time-variant lighting conditions.

Traditionally, 3D reconstruction relied on Structure-from-Motion (SfM) and Multi-View Stereo (MVS) algorithms [2]. While these methods, such as COLMAP, provide robust camera pose estimation, they struggle to capture thin structures and view-dependent effects (e.g., reflections). More importantly, traditional pipelines typically assume a static world, making them brittle when applied to dynamic video sequences where geometry consistency is violated over time.

The field witnessed a paradigm shift with the introduction of **Neural Radiance Fields (NeRF)** [3]. By representing scenes as implicit continuous functions parameterized by Multi-Layer Perceptrons (MLPs), NeRF achieved unprecedented rendering quality. Subsequent works, such as D-NeRF [4] and Nerfies [5], extended this implicit paradigm to dynamic domains by

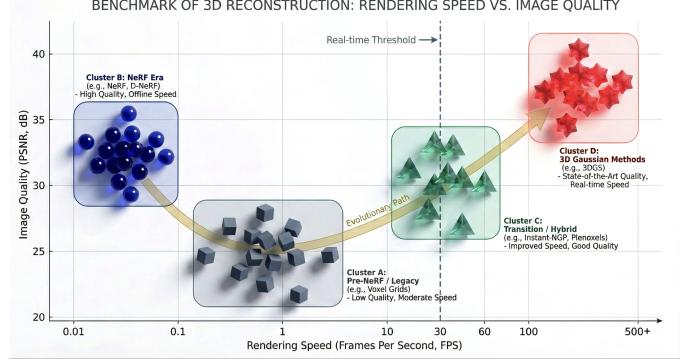


Fig. 1: Performance landscape of video-based 3D reconstruction methods. We visualize the trade-off between rendering speed (FPS, log scale) and reconstruction quality (PSNR). Implicit methods achieve high quality but slow inference. Explicit methods (3DGS) achieve real-time performance while maintaining high fidelity.

introducing deformation fields to handle non-rigid motion. However, implicit methods suffer from prohibitive computational costs due to the extensive ray-marching sampling required during both training and inference, limiting their deployment in real-time applications.

Recently, **3D Gaussian Splatting (3DGS)** [10] has emerged as a compelling alternative, marking a return to explicit volumetric representations. By combining the differentiability of deep learning with the efficiency of rasterization-based rendering, 3DGS enables real-time rendering (100+ FPS) and fast training speeds. This breakthrough has triggered a new wave of research focused on extending Gaussian primitives to 4D spatiotemporal modeling [11], [12], aiming to combine the efficiency of explicit representations with the flexibility of neural fields.

Despite the explosion of research papers in this domain, a systematic comparison between implicit (NeRF-based) and explicit (Gaussian-based) approaches for video modeling is lacking. In this survey, we provide a comprehensive review of the state-of-the-art methods for video-based 3D reconstruction. The main contributions of this paper are summarized as follows:

- We propose a structured taxonomy of dynamic scene modeling methods, categorizing them into *Deformation-based*, *Spacetime-based*, and *Hybrid* approaches across both implicit and explicit representations.
- We provide an in-depth analysis of the transition from NeRF to 3D Gaussian Splatting, highlighting the trade-

¹School of Computer Science and Technology, Ocean University of China, Qingdao 266100, China.

²School of MPs & EPs, Heriot-Watt University, Edinburgh EH14 4AS, UK.

THE EVOLUTION OF VIDEO-BASED 3D RECONSTRUCTION

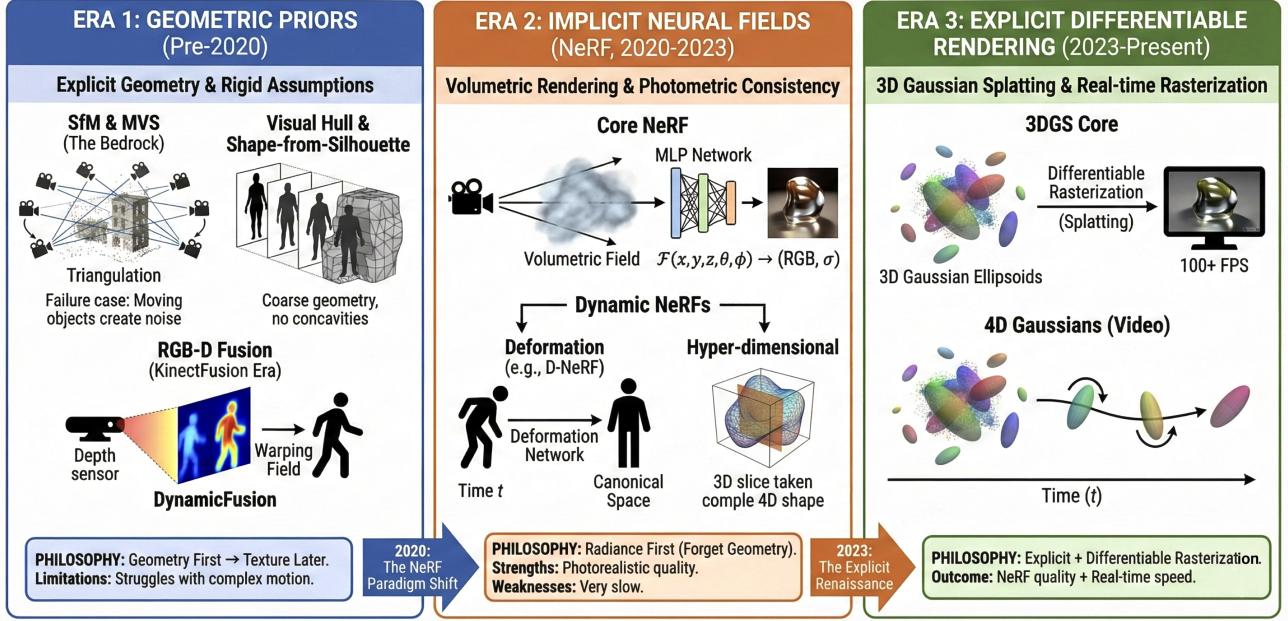


Fig. 2: Visualization of the evolution of video-based 3D reconstruction

offs between rendering quality, training efficiency, and storage overhead.

- We conduct a comparative evaluation of representative frameworks and discuss open challenges, including storage optimization and long-duration video modeling, to guide future research directions.

II. PRELIMINARIES

In this section, we first provide a unified mathematical formulation for the problem of video-based 3D reconstruction. We then detail the fundamental principles of the two dominant paradigms: the implicit neural representation introduced by NeRF and the explicit differentiable representation proposed by 3D Gaussian Splatting.

A. Mathematical Formulation

Video-based 3D reconstruction can be fundamentally formulated as an inverse rendering problem. Given a sequence of N observed images $\mathcal{I} = \{I_1, I_2, \dots, I_N\}$ captured from a video stream, along with their corresponding camera poses $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$ (typically estimated via Structure-from-Motion, SfM), the goal is to optimize a 3D scene representation \mathcal{S} such that it can synthesize novel views consistent with the input observations.

Formally, we seek to minimize the photometric reconstruction error between the rendered images and the ground truth images:

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \sum_{k=1}^N \mathcal{L}(\mathcal{R}(\mathcal{S}, P_k), I_k) \quad (1)$$

where $\mathcal{R}(\cdot)$ denotes the differentiable rendering function that projects the 3D representation \mathcal{S} onto the 2D image plane

given a camera pose P_k , and \mathcal{L} represents the loss function (typically ℓ_1 , ℓ_2 , or D-SSIM) measuring the discrepancy between the synthesized and observed pixels.

B. Fundamentals of Neural Radiance Fields

Neural Radiance Fields (NeRF) represent a scene implicitly as a continuous volumetric function parameterized by a Multi-Layer Perceptron (MLP).

1) *Scene Representation*: The core function, denoted as F_Θ , maps a continuous 5D coordinate input—consisting of a 3D spatial location $\mathbf{x} = (x, y, z)$ and a 2D viewing direction $\mathbf{d} = (\theta, \phi)$ —to a volume density σ and a view-dependent emitted color $\mathbf{c} = (r, g, b)$:

$$F_\Theta : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma) \quad (2)$$

By conditioning the color on the viewing direction \mathbf{d} , NeRF is able to model complex view-dependent effects such as specular reflections.

2) *Volume Rendering Equation*: To synthesize an image, NeRF employs classical volume rendering techniques. The color $C(\mathbf{r})$ of a camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is computed by integrating the radiance along the ray from near bound t_n to far bound t_f :

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt \quad (3)$$

where $T(t)$ denotes the accumulated transmittance, representing the probability that the ray travels from t_n to t without hitting any other particles:

$$T(t) = \exp \left(- \int_{t_n}^t \sigma(\mathbf{r}(s)) ds \right) \quad (4)$$

3) *Discretization*: In practice, the continuous integral is approximated using stratified sampling. The ray is divided into N bins, and the integral is reduced to a finite sum:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \quad (5)$$

where $T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$

Here, δ_i is the distance between adjacent samples. This differentiable rendering pipeline allows the MLP parameters Θ to be optimized via stochastic gradient descent.

C. Fundamentals of 3D Gaussian Splatting

To address the computational inefficiency of volumetric ray-marching in NeRF, 3D Gaussian Splatting (3DGS) introduces an explicit representation based on anisotropic 3D Gaussians, enabling real-time rendering via rasterization.

1) *3D Gaussian Representation*: The scene is represented by a set of 3D Gaussians. Each primitive is defined by a center position μ , a covariance matrix Σ , an opacity α , and spherical harmonics (SH) coefficients for view-dependent color. The influence of a Gaussian at a point \mathbf{x} is defined as:

$$G(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right) \quad (6)$$

To ensure the covariance matrix Σ remains positive semi-definite during optimization, it is factorized into a scaling matrix S and a rotation matrix R :

$$\Sigma = RSS^T R^T \quad (7)$$

2) *Projection and Splatting*: Unlike NeRF's ray-marching, 3DGS projects the 3D Gaussians directly onto the 2D image plane. Given a viewing transformation W , the covariance matrix Σ' in camera coordinates is approximated as:

$$\Sigma' = JW\Sigma W^T J^T \quad (8)$$

where J is the Jacobian of the affine approximation of the projective transformation. This step effectively "splats" the 3D ellipsoids into 2D ellipses.

3) *Tile-based Blending*: The final pixel color is computed using α -blending of the N sorted Gaussians overlapping the pixel:

$$C = \sum_{i \in \mathcal{N}} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (9)$$

Here, α_i is obtained by evaluating the 2D Gaussian multiplied by the learned opacity. This tile-based rasterization approach allows for highly efficient backpropagation and real-time rendering speeds (> 100 FPS), making it particularly suitable for high-resolution video reconstruction.

III. TAXONOMY OF DYNAMIC SCENE MODELING

A. Implicit Deformation Fields

This category of methods addresses dynamic scene modeling by decoupling the scene's geometry into two components: a static structure and a time-variant deformation.

1) *Core Assumption and Architecture*: The fundamental assumption of this paradigm is the existence of a time-independent **Canonical Space**. The dynamic scene observed at any time instance t is modeled as a warped version of this canonical configuration [4]. This approach simplifies the learning objective by separating motion from appearance.

Architecturally, these frameworks typically employ two joint Multilayer Perceptrons (MLPs):

- **Deformation Network**: This network learns the mapping from the observation space to the canonical space. It takes the current 3D position \mathbf{x} and a temporal indicator (time t or a latent code ω_i) as input, and outputs a displacement vector $\Delta\mathbf{x}$ (or a dense vector field).
- **Canonical Network**: This is a standard static NeRF that represents the scene in its "rest pose". It takes the corrected position $\mathbf{x}_{canon} = \mathbf{x} + \Delta\mathbf{x}$ and viewing direction \mathbf{d} to predict color \mathbf{c} and density σ .

Mathematically, this process can be formulated as:

$$(\mathbf{x}_{canon}, \Delta\mathbf{x}) = \mathcal{D}(\mathbf{x}, t), \quad (\mathbf{c}, \sigma) = \mathcal{F}_{NeRF}(\mathbf{x}_{canon}, \mathbf{d}) \quad (10)$$

where \mathcal{D} represents the deformation network and \mathcal{F}_{NeRF} represents the canonical radiance field [4].

2) *Representative Frameworks*: The evolution of implicit deformation fields can be traced through three milestone works, each addressing specific limitations of the previous:

D-NeRF: Pumarola et al. established the baseline for this category by introducing the translation field $\mathbf{x}(t) \rightarrow \mathbf{x}_{canon} + \Delta\mathbf{x}$. D-NeRF optimizes the scene using a sparse set of monocular images, successfully reconstructing non-rigid objects (e.g., a moving mechanical arm or a jumping human) by learning a continuous displacement field conditioned on time t [4].

Nerfies: While D-NeRF works well for synthetic or controlled data, applying this concept to casual selfie videos introduces challenges such as camera movement and under-constrained deformations. Park et al. introduced *Nerfies*, which replaces the explicit time input t with learned latent deformation codes ω_i to handle irregular timestamps. Crucially, Nerfies introduces **elastic regularization** (based on the Jacobian of the deformation field) to impose rigidity priors. This prevents the "background" or "rigid parts" of a face from distorting unrealistically, ensuring that the deformation field remains plausible even with sparse views [5].

HyperNeRF: A critical limitation of pure deformation fields is their continuity; they cannot model **topological changes** (e.g., a mouth opening/closing or separating objects) because mapping a continuous space to a discontinuous one requires a break in the vector field [6]. Park et al. proposed *HyperNeRF* to address this by lifting the scene into a higher-dimensional "hyperspace". Instead of simply deforming points, HyperNeRF learns a **slicing surface** in this ambient hyperspace. Topological variations are effectively modeled as continuous deformations of the slicing manifold, allowing the method to handle complex level-set changes that D-NeRF and Nerfies cannot [6].

3) Summary:

- **Pros**: These methods excel at reconstructing non-rigid, organic objects (such as human faces and bodies) with

photorealistic quality, inheriting the high-frequency detail capabilities of the underlying NeRF representation [3].

- **Cons:** The primary drawback is the computational cost; training is extremely slow due to the MLP queries required for ray marching.

Furthermore, pure deformation-based methods (excluding HyperNeRF) fundamentally struggle with large topological changes due to the continuity constraints of neural networks [6].

B. Spacetime Neural Fields

Unlike deformation-based methods that warp a static template, Spacetime Neural Fields treat time as an intrinsic dimension of the scene. These methods directly learn a continuous function over the spatiotemporal domain, effectively slicing the 4D volume to retrieve the state of the scene at any specific moment.

1) *Core Concepts and Architecture*: The fundamental shift in this paradigm is the expansion of the input domain. While a standard NeRF maps a 5D coordinate (\mathbf{x}, \mathbf{d}) to color and density, Spacetime Neural Fields extend this to a **Higher-dimensional Input** by incorporating time t . The mapping function evolves into $F_\Theta : (\mathbf{x}, \mathbf{d}, t) \rightarrow (\mathbf{c}, \sigma)$. To handle complex dynamics, two primary strategies are employed:

Latent Code Modulation: Instead of using a raw scalar t , methods like DyNeRF learn a set of compact, time-variant *latent codes* \mathbf{z}_t for each frame. These codes are fed into the MLP alongside spatial coordinates, acting as a switch that modulates the network's activations to represent topological changes, varying illumination, and transient objects.

Physical Consistency (Scene Flow): To enforce temporal coherence, some approaches incorporate physical constraints. A notable example is **NSFF (Neural Scene Flow Fields)** [7], which predicts not only the radiance and density but also the 3D *Scene Flow* vector \mathbf{f} . By enforcing consistency between the backward/forward flow and the geometric positions across frames, the model ensures that the dynamic representation adheres to physical motion laws.

Mathematically, the general formulation for this category can be expressed as:

$$(\mathbf{c}, \sigma) = F_\Theta(\mathbf{x}, \mathbf{d}, \mathbf{z}_t) \quad \text{or} \quad (\mathbf{c}, \sigma, \mathbf{f}) = F_\Theta(\mathbf{x}, t) \quad (11)$$

where \mathbf{z}_t represents the learned latent code and \mathbf{f} denotes the predicted scene flow vector.

2) *Representative Frameworks*: In the following, we discuss two representative frameworks that address dynamic modeling under distinct input constraints: **DyNeRF** for synchronized multi-view videos and **Video-NeRF** for casual monocular captures.

DyNeRF (Dynamic NeRF): Li et al. proposed DyNeRF to address the challenge of reconstructing dynamic scenes from multi-view videos [8]. Instead of a deformation field, DyNeRF conditions the radiance field on a learnable latent embedding \mathbf{z}_t . This allows the model to capture complex dynamics that are difficult to model with continuous deformations, such as flame, pouring liquid (topology changes), and view-dependent lighting effects.

Video-NeRF: Focusing on the challenging monocular setting, Video-NeRF learns a direct mapping from spacetime coordinates (\mathbf{x}, t) to an irradiance field (omitting view dependence to handle single-video ambiguity). To prevent the “motion-appearance ambiguity” (where geometry changes are incorrectly explained as texture changes), it incorporates depth supervision from pre-trained monocular estimators and introduces a “static scene loss” to constrain unobserved regions by propagating information across time [9].

3) Summary:

- **Pros:** These methods are well-suited for large-scale scenes and complex dynamics involving drastic topological changes (e.g., fluid dynamics, fire) where defining a canonical template is infeasible.
- **Cons:** The addition of the time dimension significantly expands the sampling space, leading to the “curse of dimensionality.” Without sufficient data (e.g., multi-view rigs) or strong regularization (e.g., depth/flow priors), these models are prone to overfitting or producing temporally incoherent artifacts.

C. Dynamic Gaussian Splatting

...

D. Hybrid Representations

...

IV. DATASETS AND METRICS

A. Common Datasets (*N3V, D-NeRF, etc.*)

...

B. Evaluation Metrics (*PSNR, SSIM, LPIPS, FPS*)

...

V. COMPARATIVE ANALYSIS

A. Quantitative Benchmarks

...

B. Qualitative Visual Comparisons

...

C. Efficiency Analysis (*Time & Memory*)

...

VI. OPEN CHALLENGES AND FUTURE DIRECTIONS

A. Storage Efficiency and Compression

...

B. Handling Long-Duration Videos

...

C. Integration with Generative Models

...

VII. CONCLUSION

...

VIII. BIOGRAPHY SECTION

REFERENCES

- [1] Apple Inc., *Expression Estimation for Headsets Using Low-Profile Antenna and Impedance Characteristic Sensing*.
- [2] J. Schönberger, J.-M. Frahm, M. Pollefeys, P.-E. Sarlin, and S. Liu, *COLMAP*, [Online]. Available: <https://colmap.github.io/index.html>
- [3] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis,” *arXiv preprint*, arXiv:2003.08934, 2020. [Online]. Available: <https://arxiv.org/abs/2003.08934>
- [4] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer, “D-NeRF: Neural Radiance Fields for Dynamic Scenes,” *arXiv preprint*, arXiv:2011.13961, 2020. [Online]. Available: <https://arxiv.org/abs/2011.13961>
- [5] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla, “Nerfies: Deformable Neural Radiance Fields,” *arXiv preprint*, arXiv:2011.12948, 2021. [Online]. Available: <https://arxiv.org/abs/2011.12948>
- [6] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz, “HyperNeRF: A Higher-Dimensional Representation for Topologically Varying Neural Radiance Fields,” *arXiv preprint*, arXiv:2106.13228, 2021. [Online]. Available: <https://arxiv.org/abs/2106.13228>
- [7] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang, “Neural Scene Flow Fields for Space-Time View Synthesis of Dynamic Scenes,” *arXiv preprint*, arXiv:2011.13084, 2021. [Online]. Available: <https://arxiv.org/abs/2011.13084>
- [8] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv, “Neural 3D Video Synthesis from Multi-view Video,” *arXiv preprint*, arXiv:2103.02597, 2021. [Online]. Available: <https://arxiv.org/abs/2103.02597>
- [9] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim, “Space-time Neural Irradiance Fields for Free-Viewpoint Video,” *arXiv preprint*, arXiv:2011.12950, 2020. [Online]. Available: <https://arxiv.org/abs/2011.12950>
- [10] J. Kerbl, A. Wang, F. Rousselle, and V. Koltun, “3D Gaussian Splatting for Real-Time Radiance Field Rendering,” *arXiv preprint*, arXiv:2304.08914, 2023. [Online]. Available: <https://arxiv.org/abs/2304.08914>
- [11] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang, “4D Gaussian Splatting for Real-Time Dynamic Scene Rendering,” *arXiv preprint*, arXiv:2310.08528, 2024. [Online]. Available: <https://arxiv.org/abs/2310.08528>
- [12] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan, “Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis,” *arXiv preprint*, arXiv:2308.09713, 2023. [Online]. Available: <https://arxiv.org/abs/2308.09713>

Yuwei ZHAO Use \begin{IEEEbiographynophoto} and the author name as the argument followed by the biography text.