

Figure 1. **Adaptive Gabor formulation (a) Smooth transition between Gaussian and Gabor kernels.** Our method (rightmost column,  $S_{adap}(x)$ ) uses a compensation term  $b$  to maintain energy stability while transitioning from pure Gaussian ( $\omega = 0$ , top) to frequency-modulated Gabor ( $\omega = 1$ , bottom). Naive combination  $1 + S(x)$  (third column) suffers from intensity artifacts. **(b) Frequency weight combinations.** Different  $(\omega_0, \omega_1)$  pairs generate diverse spatial patterns, from smooth (low  $\omega$ ) to high-frequency textures (high  $\omega$ ), enabling adaptive detail capture in different scene regions.

#### 072 B.4. Adaptive Gabor Formulation.

073 We provide additional visualization of the proposed Adaptive Gabor design in Fig. 1. Figure 1(a) shows that a naive  
074 Gaussian-to-Gabor replacement can lead to intensity changes  
075 and visible artifacts, since the modulation term alters the kernel energy as frequency content increases. In contrast, our  
076 compensation term  $b$  stabilizes this transition, allowing the primitive to move smoothly from Gaussian-like behavior to  
077 frequency-enhanced Gabor modeling. Figure 1(b) further  
078 illustrates that different combinations of wave coefficients  
079 produce diverse spatial-frequency patterns, ranging from  
080 smooth low-frequency responses to richer high-frequency  
081 textures. These visualizations help explain why the proposed  
082 formulation can adapt its frequency behavior to different  
083 scene regions while maintaining stable rendering quality.  
086

#### 087 B.5. Conclusion

088 This proof demonstrates that our Adaptive Gabor representation gracefully degrades to a standard Gaussian when frequency content is not needed ( $\omega_i \rightarrow 0$ ), while smoothly transitioning to frequency-enhanced Gabor modes when high-frequency details are required ( $\omega_i > 0$ ). This adaptive behavior is crucial for maintaining energy stability across diverse scene regions with varying frequency characteristics.  
094

#### 095 C. Additional Visual Comparisons and Results

096 For comprehensive visual comparisons with baseline methods across various dynamic scenes, please refer to Figs. 2,  
097

5 and 6. These figures demonstrate our method’s superior performance in preserving high-frequency texture details and maintaining temporal consistency across challenging scenarios including fast motion, occlusions, and complex deformations.  
102

For interactive visualization of downstream application results, including frame interpolation, video editing, and stereo view synthesis, please refer to the supplementary HTML page (`index.html`). The interactive viewer allows frame-by-frame inspection and video playback to better appreciate the temporal coherence and visual quality of our method.  
108

#### D. Additional Applications

**Depth Consistency.** We further evaluate whether the reconstructed representation preserves temporally stable geometry in depth space. As shown in Fig. 3, our method maintains consistent depth estimates for static scene elements across frames, while the per-frame baseline exhibits noticeable temporal flickering and boundary inconsistency. This improvement comes from our explicit 3D primitive representation together with smooth temporal motion modeling, which encourages coherent geometry over time rather than frame-wise independent predictions. Such temporal stability is particularly important for downstream applications that rely on consistent depth, such as video editing, view synthesis, and geometry-aware video processing.  
122

**Stereo View Synthesis.** Our representation also supports stereo view synthesis from monocular video. As shown  
123  
124

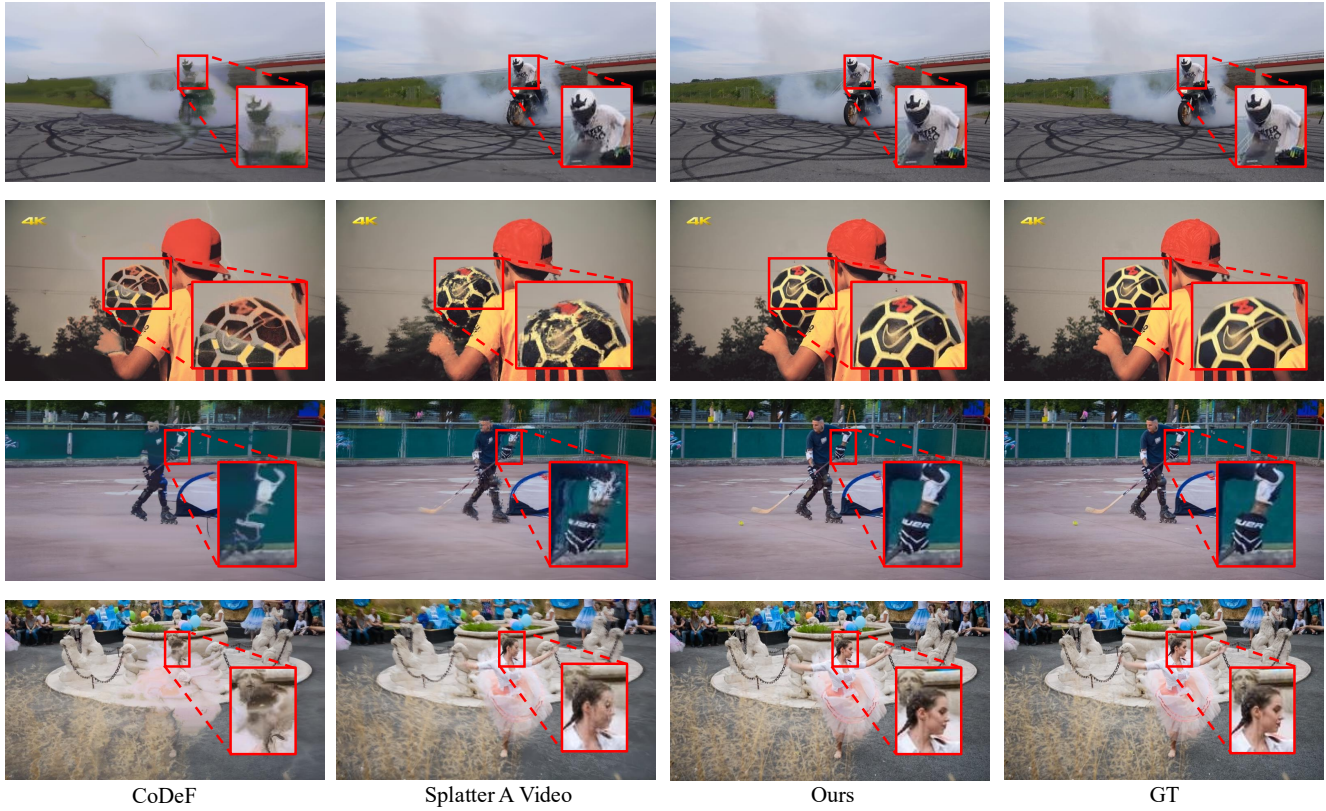


Figure 2. Visual comparison on the DAVIS dataset.

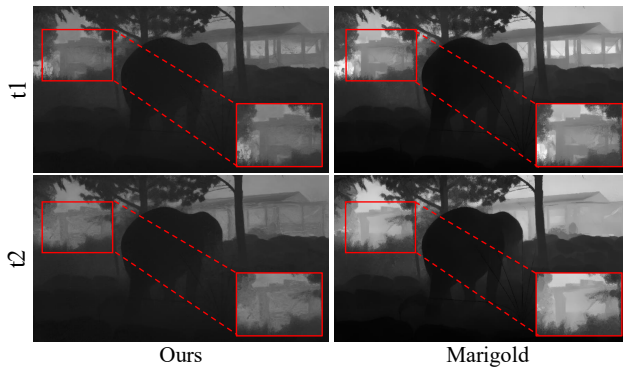


Figure 3. **Depth consistency across time.** (Left) Our 3D primitive representation maintains consistent depth for static elements across frames. (Right) Per-frame estimation (Marigold [2]) shows temporal flickering (red boxes). Explicit 3D geometry with smooth motion modeling ensures temporal coherence essential for depth-based video applications.

125 in Fig. 4, the learned 3D structure enables plausible view-  
 126 point variation and consistent scene geometry, producing a  
 127 convincing stereo effect from a single input video. This result  
 128 suggests that Adaptive Gabor primitives in the orthographic  
 129 camera coordinate space capture not only appearance details  
 130 but also sufficiently stable spatial structure for downstream



Figure 4. **Stereo view synthesis.** Our 3D representation enables novel view synthesis for stereo visualization from monocular video. This demonstrates that Adaptive Gabor primitives in orthographic camera coordinate space capture accurate 3D geometry, enabling immersive applications.

view synthesis. The ability to generate stereo views further  
 demonstrates the practical utility of our explicit dynamic  
 representation beyond standard reconstruction.

**References**

[1] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 1  
 [2] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estima-

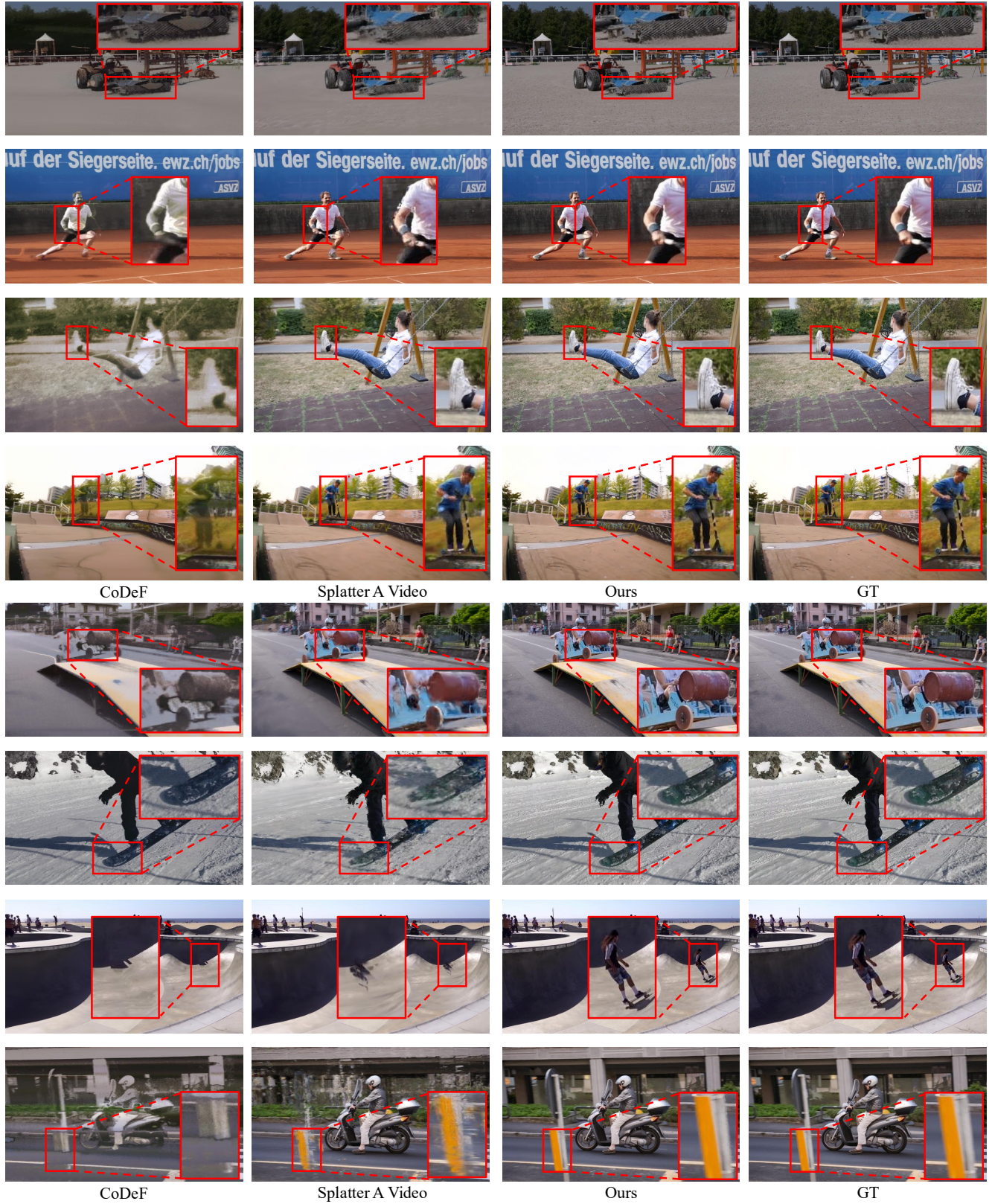


Figure 5. Visual comparison on the DAVIS dataset.

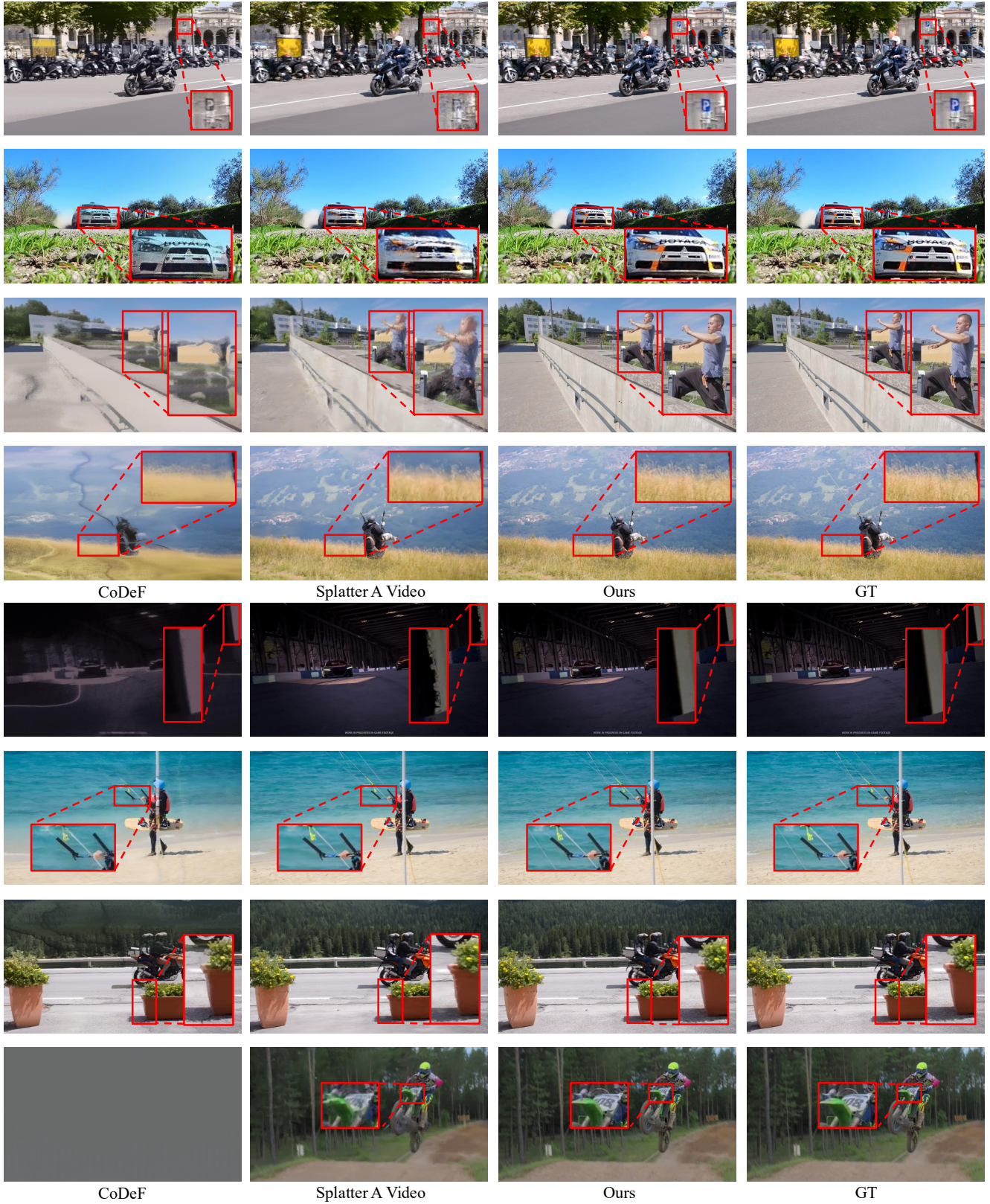


Figure 6. Visual comparison on the DAVIS dataset.

142  
143

tion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9492–9502, 2024. 3