

Hypothesis Report Practice

Max Matkovski¹

I. QUESTION 1

The authors identify a significant gap in generative models' ability to synthesize realistic, long-term scene motion from a single, still image. It is true that recent advances in generative models—specifically diffusion models—have allowed for the creation of diverse and realistic-seeming images, even short video clips. However, these models typically operate on overall pixel values. These models struggle to capture the fundamental, low-dimensional structure underlying natural dynamics. This can be seen in examples such as trees swaying or clothes fluttering in the wind. Existing approaches require either multiple input frames or are unable to generalize the complex, temporally coherent motions which we see in real-world scenes. The gap, therefore, lies in the absence of a model that can learn a prior over long-range pixel motion trajectories. The authors address this shortcoming by proposing a novel framework which learns from a generative prior in the Fourier domain, which enables the synthesis of spectral volumes that represent the motion dynamics of a scene. This can later be used for applications like seamlessly looping video generation with interactive dynamic simulations.

II. QUESTION 2

The hypothesis of the article is mainly that learning a generative prior over dense, long-term scene motion in the Fourier domain, using motion trajectories which are extracted from real videos, allows a model to synthesize more realistic and coherent motion dynamics from one single image. This approach leverages a frequency coordinated diffusion process that can produce spectral volumes which capture the essential dynamics of natural scenes, allowing for convincing animations and interactive manipulations of still images.

III. QUESTION 3

To fully understand the claims, I need to look up:

Spectral Volume Representation: The paper relies on representing motion as spectral volumes in the Fourier domain. Understanding the mathematical formulation and prior uses of spectral volumes in motion analysis is crucial, as it underpins the model's ability to capture oscillatory dynamics. For context, modal analysis and spectral methods have been used in graphics and engineering to model vibrations and motions (Pentland and Williams, 1989).

Diffusion Models in Frequency Domain: The model predicts motion in the frequency domain using a diffusion

process. I need to review how diffusion models are adapted for non-RGB data and how coordination across frequency bands is achieved, as this impacts the model's ability to generate coherent motion (Ho et al., 2020).

Looking up these concepts would help me clarify the effectiveness of the approach, providing further context for the model architecture and help me evaluate overall efficacy and validity of the results stated.

IV. QUESTION 4

One limitation is the model's reliance on training data that consists primarily of natural, oscillatory motions (e.g., trees, flowers, fabrics). As a result, its ability to generalize to scenes with non-oscillatory or highly nonlinear dynamics is limited. The approach used in the paper assumes that the dominant motions in a scene can be effectively captured in the frequency domain, which may not hold for all types of motion. For example, abrupt or non-repetitive interactions involving multiple objects may not work well. Further, the model's performance is based on the quality and variety of the motion trajectories extracted during training. Errors or biases in this data or training period could lead to erroneous outputs. This method also presumes accurate segmentation and motion extraction from videos, which can be challenging in cluttered or dynamic environments. Lastly, while the model enables interactive dynamics, the plausibility of these interactions is limited by abstractions inherent in the spectral volume representation, potentially lowering realism in complex scenarios.

V. QUESTION 5

Next steps can possibly include expanding the model to handle a wider range of motion types, including non-oscillatory, abrupt, or multi-object dynamics. Integrating semantic understanding or object-level reasoning could potentially enhance the model's ability to animate scenes with complex interactions or articulated motion. Another direction is improving the physical plausibility of interactive dynamics by incorporating explicit physics-based constraints or learning from physically simulated data. Additionally, we can explore the application of this approach to different domains such as medical imaging or robotics. This realistic motion synthesis from limited data is valuable and could be impactful.

Hypothesis: Incorporating object segmentation and physics-informed priors into the generative model will improve the realism and diversity of synthesized motion, enabling the generation of plausible dynamics for a wider variety of scenes. This includes those with multiple

¹Max Matkovski is with the Georgia Institute of Technology. Contact: max.matkovski@gatech.edu

interacting objects or non-periodic motion. This could be tested by training the model on datasets with annotated object boundaries and comparing the quality of generated animations in both oscillatory and non-oscillatory scenarios.

REFERENCES

- [1] A. Pentland and J. Williams, “Good vibrations: Modal dynamics for graphics and animation,” ACM SIGGRAPH Comput. Graph., vol. 23, no. 3, pp. 207–214, 1989.
- [2] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” arXiv preprint arXiv:2006.11239, 2020.