

# GPT4Motion: Scripting Physical Motions in Text-to-Video Generation via Blender-Oriented GPT Planning

Jiayi Lv<sup>1,2\*</sup> Yi Huang<sup>1,2\*</sup> Mingfu Yan<sup>1,2\*</sup> Jiancheng Huang<sup>1,2</sup> Jianzhuang Liu<sup>1</sup>  
Yifan Liu<sup>1</sup> Yafei Wen<sup>3</sup> Xiaoxin Chen<sup>3</sup> Shifeng Chen<sup>1†</sup>

<sup>1</sup>Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences,

<sup>2</sup>University of Chinese Academy of Sciences, <sup>3</sup>VIVO AI Lab

{jx.lv1, yi.huang, mf.yan, jc.huang, jz.liu, yf.liu2, shifeng.chen}@siat.ac.cn

{yafei.wen, xiaoxin.chen}@vivo.com

## Supplementary Material

*“A white T-shirt flutters in light wind”*

*“A white T-shirt flutters in the wind”*

*“A white T-shirt flutters in strong wind”*

GPT4Motion

AnimateDiff [1]

ModelScope [8]

Text2Video-Zero [4]

DiracT2V [2]

Figure 1. Comparison of the video results generated by different text-to-video models under different physical conditions. *Best viewed with Acrobat Reader for animation.*

In this supplement material, we further compare the generation ability of GPT4Motion and four baselines for different physical conditions in Sec. 1. In Sec. 2, we introduce the quantitative metrics we used. In Sec. 3, we describe the details of our settings for Blender. In Sec. 4 and Sec. 5, we show how GPT-4 generates encapsulated Python functions via prompts, and how GPT-4 incorporates its own knowledge of physics to control the motion of objects, respectively.

## 1. More Comparison with Baselines

In the main paper, we have compared GPT4Motion with four baselines (AnimateDiff [1], ModelScope [8], Text2Video-Zero [4], and DirecT2V [2]) on three scenarios (rigid object drop and collision, cloth draping and swinging, and liquid flow). Here, we further conduct an experiment on dynamic effects of a T-shirt being blown by the wind under three wind strengths. The results are shown in Figure 1, where the seed is randomly chosen and fixed in all the generations. We can see that these baselines all fail to generate videos that match the user prompts and are unable to control the intensity of physical phenomena solely based on the linguistic descriptions. In contrast, our GPT4Motion not only precisely designs the parameters of Blender encapsulated functions (such as wind strength) through GPT-4, but also leverages Blender’s physics engine to simulate the complex flapping and twisting dynamics of the T-shirt in the wind.

## 2. Quantitative Evaluation Metrics

Here, we introduce the metrics employed in the main paper:

1. *Motion Smoothness* [3]. This metric evaluates the smoothness of motion in generated videos, ensuring it conforms to the physical laws of the real world. The evaluation utilizes motion priors from the video frame interpolation model [5] to assess the smoothness of generated motions.
2. *Temporal Flickering* [3]. This metric identifies imperfections in temporal consistency within generated videos, especially in local and high-frequency details. The method involves analyzing static frames and computing the mean absolute difference across frames. These values are then subtracted from 255 and normalized between 0 and 1.
3. *CLIP-Score* [6]. This metric is designed to assess Text-Video Consistency. It leverages a pretrained ViT-B/32 CLIP model [7] as a feature extractor to calculate the cosine similarity between each video frame and the associated text.

## 3. Blender Settings

We use Blender to generate two sequences of scene depth maps and edge maps. The edge maps are obtained by

Blender’s built-in Freestyle<sup>1</sup> feature, which is an advanced rendering technique for non-photorealistic line drawings from 3D scenes. It not only allows for various line styles, such as hand-drawn or sketch, but also serves as an independent rendering pass without losing edge information of the scene and additional post-processing. Moreover, we normalize the scene’s Z-channel and render it as an independent depth map. Since the rendering of edge maps and depth maps is unrelated to ray tracing, we choose Blender’s “Workbench” rendering mode. This mode avoids the substantial computational resources required for ray tracing, significantly enhancing rendering efficiency. In our experiments, completing the entire rendering process using Blender scripts takes only 2 to 3 minutes, with most of the time spent on calculating physical effects.

## 4. GPT-4’s Generation of Encapsulated Python Functions for Blender

Blender incorporates an independent Python library, enabling the transformation of various manual processes within the software to automate tasks through its Python API. Leveraging GPT-4’s capabilities of grasping Blender’s Python API, we can rapidly generate encapsulated Python functions for Blender. This involves inputting specific functional goals and the associated Blender operation steps into GPT-4. Consequently, GPT-4 can autonomously create Python code that corresponds to these steps, effectively encapsulating them into Python functions. One example is given below.



**User:**

You are an expert in Blender and Python. Please generate encapsulated Blender Python functions by following our step-by-step instructions, writing the Blender Python API corresponding to each step, designing each function’s name, and explaining its functionality and the meaning of each parameter in the docstring.



**GPT-4:**

Sure, I’d be happy to help with that. Please go ahead and provide me with the step-by-step instructions for the Blender Python functions you need, and I will write the corresponding code for you.

<sup>1</sup><https://docs.blender.org/manual/en/latest/render/freestyle/introduction.html>

**User:**

Please build a function to remove all objects, cameras and lights from the initial scene.

**GPT-4:**

Certainly! To create a Blender Python function that removes all objects, cameras, and lights from the initial scene, you can define a function like this:

```

1 # bpy is the Python module used
  for scripting within Blender
2 import bpy
3
4 def clear_scene():
5     """
6     Clears all objects from the
7     current Blender scene.
8
9     This function selects all
10    objects in the scene and
11    deletes them.
12
13    It is useful when starting a
14    new scene setup or
15    resetting the scene to a
16    blank state.
17
18    No inputs or outputs.
19
20    Typically called at the
21    beginning of a script
22    when starting a new scene
23    setup.
24    """
25    bpy.ops.object.select_all(
26        action='SELECT')
27    bpy.ops.object.delete()

```

You can call this function whenever you need to clear the scene of all types of objects. Just make sure to be in the correct context where the scene is accessible.

**User:**

Please now create two physics functions, one that adds collision properties to the inputted object, and another that adds rigid body properties to the object. Then, construct a function to create a floor that scales it up 50 times in

all directions, with physics set to collision and rigid body (passive, with elasticity set to 1).

**GPT-4:**

To accomplish this task, we'll create three functions in Blender's Python API:

1. Add Collision Properties: This function will add collision properties to a specified object.
2. Add Rigid Body Properties: This function will add rigid body properties to a specified object.
3. Create a Floor with Physics: This function will create a floor, scale it, and set its physics properties to collision and rigid body.

```

1 import bpy
2
3 def add_collision(obj):
4     """
5     Adds a collision modifier to
6     a Blender object if it
7     does not already have one
8     .
9
10    Parameters:
11    - obj (Blender Object): The
12    object to which the
13    collision modifier will
14    be added.
15
16    Returns:
17    - None, but the object will
18    have a collision modifier
19    added to it if it wasn't
20    present before.
21
22    This function is used in
23    physics simulations where
24    it's necessary for
25    objects to interact with
26    each other,
27    such as rigid body or soft
28    body simulations. The
29    collision modifier makes
30    the object a collider in
31    the physics
32    simulation, allowing other

```

```

16         objects to bounce off or
17         slide along its surface.
18     """
19     if "Collision" not in obj.
20         modifiers:
21             bpy.ops.object.
22                 modifier_add(type='
23                     COLLISION')
24
25 def add_rigid_body(obj, mass=1,
26 elasticity=0.5,
27 rigid_body_type='ACTIVE'):
28     """
29     Adds a rigid body physics
30     characteristic to a
31     Blender object.
32
33     Args:
34     - obj (Blender Object): The
35       object to which the rigid
36       body physics will be
37       applied.
38     - mass (float): The mass of
39       the object, default is 1.
40     - elasticity (float): The
41       restitution (bounciness)
42       of the object, default is
43       0.5.
44     - rigid_body_type (str): The
45       type of rigid body, can
46       be 'ACTIVE' or 'PASSIVE',
47       default is 'ACTIVE'.
48
49     Returns:
50     - None, but the object is now
51       a rigid body with the
52       specified physics
53       properties.
54
55     This function is crucial for
56     physics simulations,
57     defining how the object
58     behaves under physical
59     forces.
60     """
61     bpy.ops.object.select_all(
62         action='DESELECT')
63     obj.select_set(True)
64     bpy.context.view_layer.
65         objects.active = obj
66     bpy.ops.rigidbody.object_add
67         ()
68     obj.rigid_body.type =
69         rigid_body_type
70     obj.rigid_body.mass = mass
71     obj.rigid_body.restitution =
72         elasticity
73     obj.rigid_body.

```

```

43         collision_shape = 'MESH'
44 def create_floor(elasticity=1):
45     """
46     Creates a floor plane in
47     Blender, scales it, and
48     sets it up with collision
49     and rigid body physics.
50     The created floor is scaled
51     to be large enough to act
52     as a ground plane for
53     most scenes.
54
55     Parameters:
56     - elasticity (float): The
57       restitution or '
58       bounciness' of the floor.
59       A value of 1 means
60       perfectly elastic, while
61       0 means no elasticity.
62       Default is 1.
63     """
64     bpy.ops.mesh.
65         primitive_plane_add(size
66         =1, enter_editmode=False,
67         align='WORLD', location
68         =(0, 0, 0))
69     floor = bpy.context.
70         active_object
71     floor.scale = (50, 50, 50)
72     floor.name = 'GROUND'
73     add_collision(floor)
74     add_rigid_body(floor,
75         rigid_body_type='PASSIVE'
76         , elasticity=elasticity)

```

These functions together allow you to dynamically add physics properties to objects and create a large, physical floor in your Blender scene. Remember to run these functions in a Blender Python scripting context where the bpy module is accessible.

After generating the function, we enhance its usability for GPT-4 by incorporating instructions (such as the strength interval of the wind) in the function's docstring. This ensures that GPT-4's usage aligns more closely with our expectations.

## 5. Combination of Blender and GPT-4's Physics Knowledge

We discover that GPT4Motion can not only use GPT-4's semantic understanding and coding capabilities to generate Blender Python scripts, but it can also leverage GPT-4's



physical knowledge to control object motion. For example, when prompted with “a basketball thrown towards the camera.”, we find that GPT-4 can calculate the time it takes for the basketball to fall, based on the height difference between the initial position of the basketball and the camera. This calculation helps to determine the initial velocity required for the basketball to fly to the camera. The speed calculation code generated by GPT-4 is as follows:

```
1 # Setup the basketball
2 basketball = setup_object(
3     os.path.join(ASSETS_PATH, "
4         basketball.obj"),
5     "Basketball",
6     size=0.24, # Standard basketball
7     diameter in meters
8     mass=0.625, # Standard basketball
9     mass in kilograms
10    position=(0, 0, 4) # Initial
11    position for the basketball
12 )
13
14 # Calculate the initial velocity for the
15 # basketball to hit the camera
16 g = 9.81 # Acceleration due to gravity
17 # in m/s^2
18 camera_height = 1.8521
19 height_difference = 4 - camera_height
20 time_to_fall = math.sqrt(2 *
21     height_difference / g)
22 initial_horizontal_velocity = 13.665 /
23     time_to_fall # Distance to camera
24 # along Y-axis
25
26 # Set the initial velocity and rotation
27 # for the basketball
28 add_initial_velocity_for_rigid_body(
29     basketball,
30     initial_velocity=(0,
31         initial_horizontal_velocity, 0),
32     initial_rotation=(0, 0, 0) # No
33     initial rotation
34 )
```

In the above script, GPT-4 first calculates the height difference between the initial position of the basketball and the camera, which determines the time required for the basketball to fall to the same height as the camera. Subsequently, GPT-4 calculates the distance between the basketball and the camera along the Y-axis to determine the required initial velocity of the basketball. This process effectively integrates basic principles of physics, such as the equations of motion, to solve a practical problem in a simulated environment like Blender.

## References

- [1] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 1, 2
- [2] Susung Hong, Junyoung Seo, Sunghwan Hong, Heeseong Shin, and Seungryong Kim. Large language models are frame-level directors for zero-shot text-to-video generation. *arXiv preprint arXiv:2305.14330*, 2023. 1, 2
- [3] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023. 2
- [4] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. In *ICCV*, 2023. 1, 2
- [5] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun-Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-field transforms for efficient frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9801–9810, 2023. 2
- [6] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tiejong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. 2023. 2
- [7] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763, 2021. 2
- [8] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modellscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 1, 2