

MonoMobility: Zero-Shot 3D Mobility Analysis from Monocular Videos

Supplementary Material

001 Summary

002 The supplementary material consists of the following two
 003 parts: (1) A detailed introduction to the dataset, includ-
 004 ing both the synthetically simulated scenarios and the real-
 005 world captured scenarios; see Section 1. (2) Additional vi-
 006 sualization results, including qualitative comparison ex-
 007 perimental results with state-of-the-arts methods [3, 5], as well
 008 as more qualitative experimental results of our method; see
 009 Section 2.

010 1. Details of Dataset

011 Our goal is to analyze motion parts and their motion at-
 012 tributes from monocular videos. For effective evaluation
 013 of our algorithm, we have constructed a Motion Pars-
 014 ing Dataset that primarily comprises virtual simulation
 015 and real-world scenarios. It includes various common
 016 articulated object categories such as drawers, wardrobes,
 017 laptops, staplers, and liftchairs, which cover three main
 018 types of articulated motion: translation, rotation, and ro-
 019 tation+translation. Some scenes contain multiple motion
 020 parts with different motion types, aimed at verifying the ef-
 021 fectiveness of relevant algorithms in solving complex tasks.
 022 The statistical details of the dataset are presented in Tab.1.

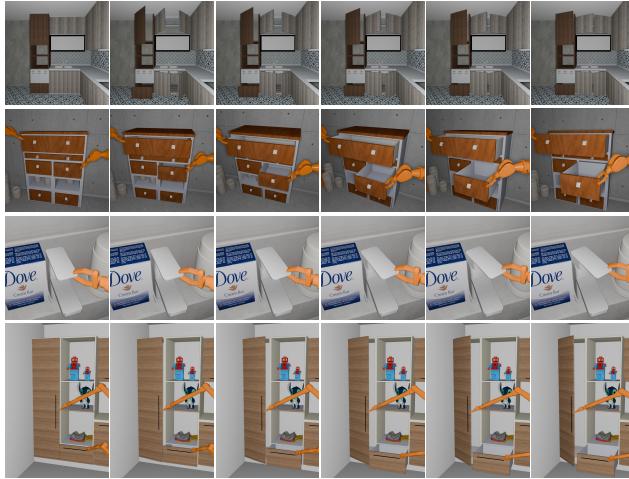


Figure 1. Virtual simulation scenarios. Each row represents a sequence of video frames captured in a virtual scene.

023 1.1. Virtual Simulation Scenarios

024 For virtual simulation scenarios, we first collect 3D models
 025 from 3D Warehouse [2] and annotate the motion parts and
 026 motion attributes using Blender [1]. Subsequently, robots



Figure 2. Real-world scenarios. Each row represents a sequence of video frames captured in a real-world scene.

027 are introduced into the scenes to simulate the interaction
 028 operations of articulated object. Finally, the constructed
 029 dynamic scenes are rendered using Blender [1], capturing
 030 motion videos sequences. To quantitatively analyze the re-
 031 sults of the motion parts and motion attributes parsing, we
 032 also captured 3D point clouds with annotations of motion
 033 parts and motion attributes. A total of 15 virtual motion
 034 simulation scenarios were built, some examples are shown
 035 in Fig.1.

036 1.2. Real-world scenarios

037 For real-world scenarios, we use the camera directly capture
 038 motion videos of articulated objects. The intrinsic parame-
 039 ters of camera are calibrated using COLMAP [4]. Unlike
 040 virtual simulation scenarios, we cannot obtain geometric
 041 and annotation information for real-world scenarios. There-
 042 fore, we only perform qualitative analysis on the real-world
 043 data. A total of 11 real-world scenarios were constructed,
 044 some examples are shown in Fig.2.

2. Additional Experimental Results

045 To further substantiate the superiority and effectiveness of
 046 our method, we conducted qualitative comparisons in mo-
 047 tion parameter prediction with the state-of-the-arts algo-
 048 rithms (PARIS-scene, PARIS-obj [3], DGMarbles* [5], and
 049 ours(w/o optim)) on additional data, as shown in Fig. 3;
 050 as well as more qualitative analysis results of motion parts
 051 and motion attributes of our method from more video se-
 052 quences, detailed in Fig. 4–17, where each figure includes
 053 the analysis results of two scenes. For each scene, the first
 054 row represents the input (video sequence) to the algorithm,
 055 while the second to fourth rows show the continuous mo-
 056 tion visualization of motion parts based on motion attributes
 057 from different viewpoints.

Data Type	Number of Scenes	Number of Motion Parts	Distribution of Motion Types	Object Categories
Virtual	15	18	10×rotation	Fridge(3), Door(1), Cupboard(4), Faucet(1), Laptop(1)
			7×translation	Drawer(6), Flatdoor(1)
			1×rotation-translation	Liftchair(1)
Real-world	11	13	8×rotation	Cupboard(3), Laptop(3), Wrench(1), Stapler(1)
			4×translation	Drawer(4)
			1×rotation-translation	Pumpbottle(1)

Table 1. The statistical of the dataset. '10×rotation' indicates that the rotation motion type contains 10 motion parts, 'Fridge(3)' indicates that the articulated object category ‘‘Fridge’’ contains 3 motion parts, and similar symbols apply accordingly.

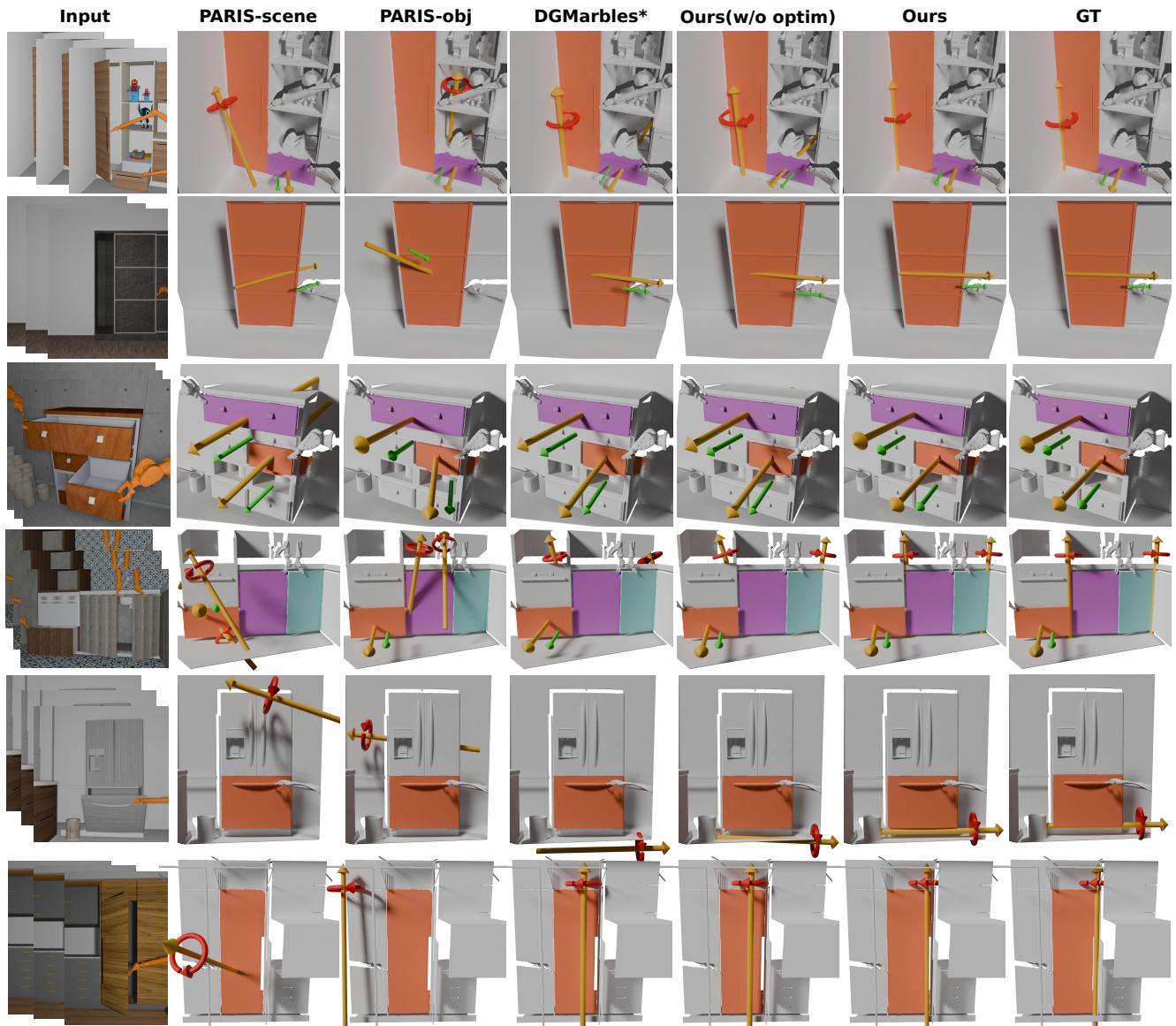


Figure 3. More qualitative comparison with state-of-the-arts methods (PARIS-scene, PARIS-obj, DGMarbles*) and Ours(w/o optim).

Real-world Data Results

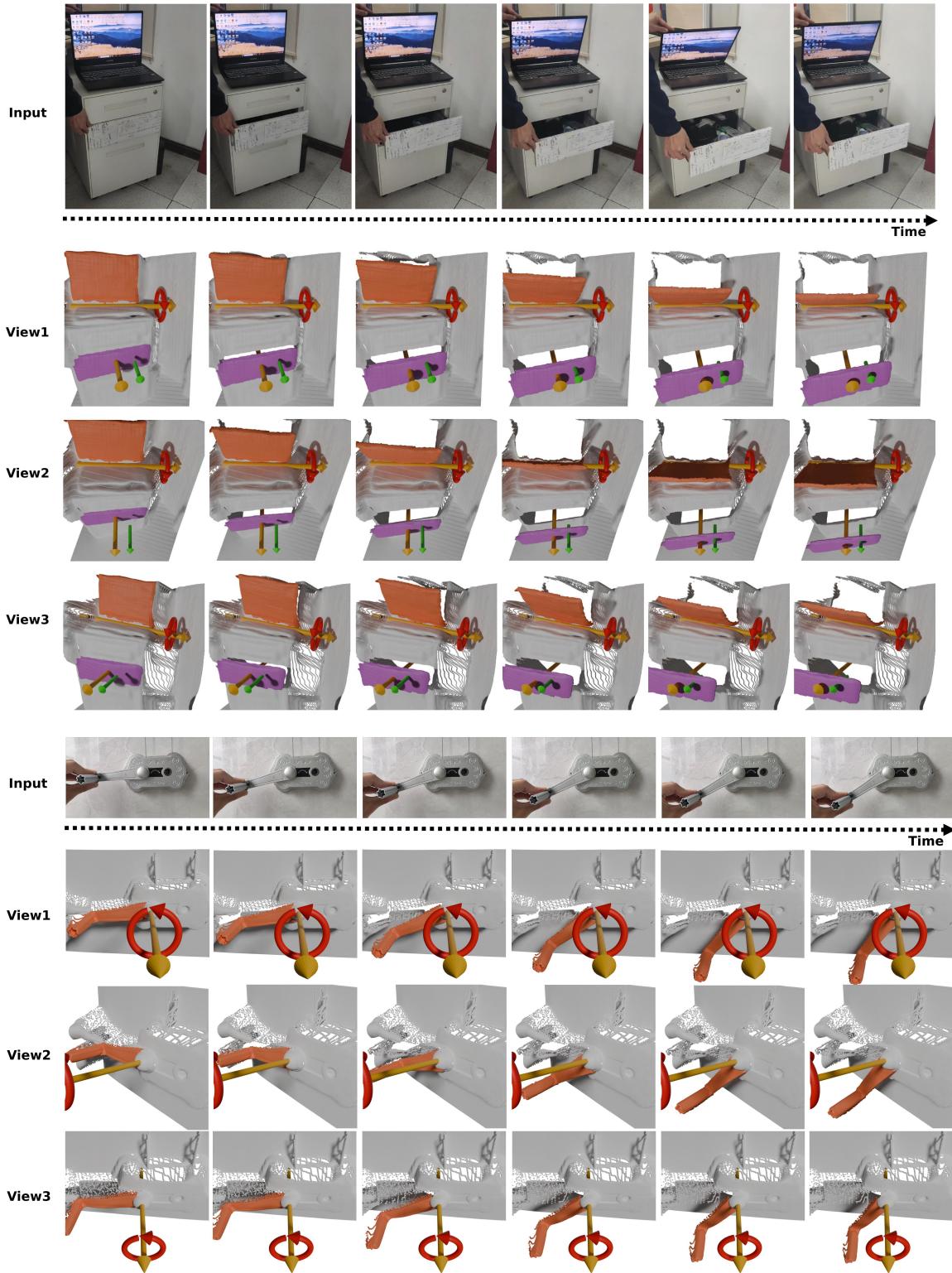


Figure 4. Analysis results of motion parts and their motion attributes for real-world scenes 1-2.

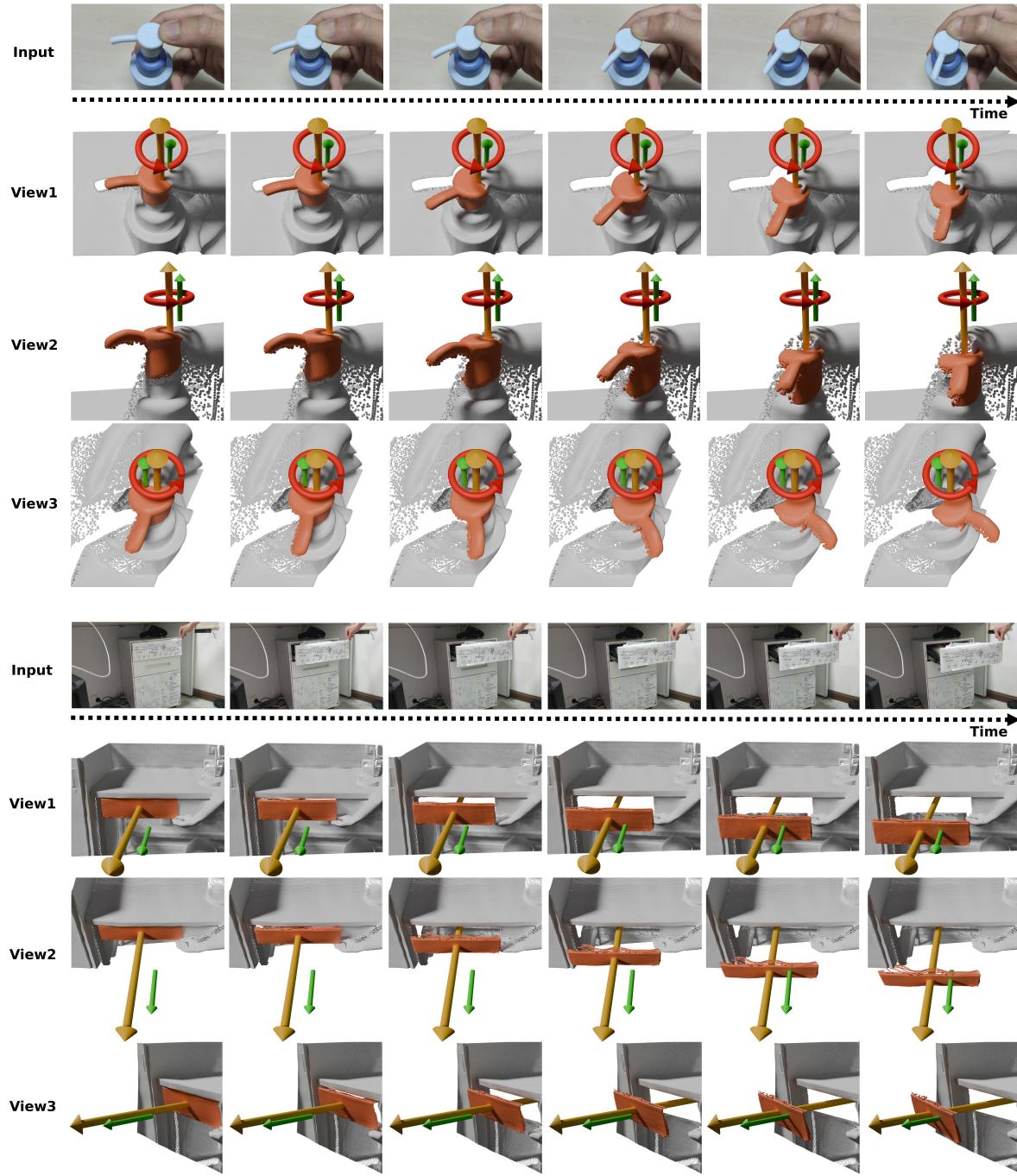


Figure 5. Analysis results of motion parts and their motion attributes for real-world scenes 3-4.

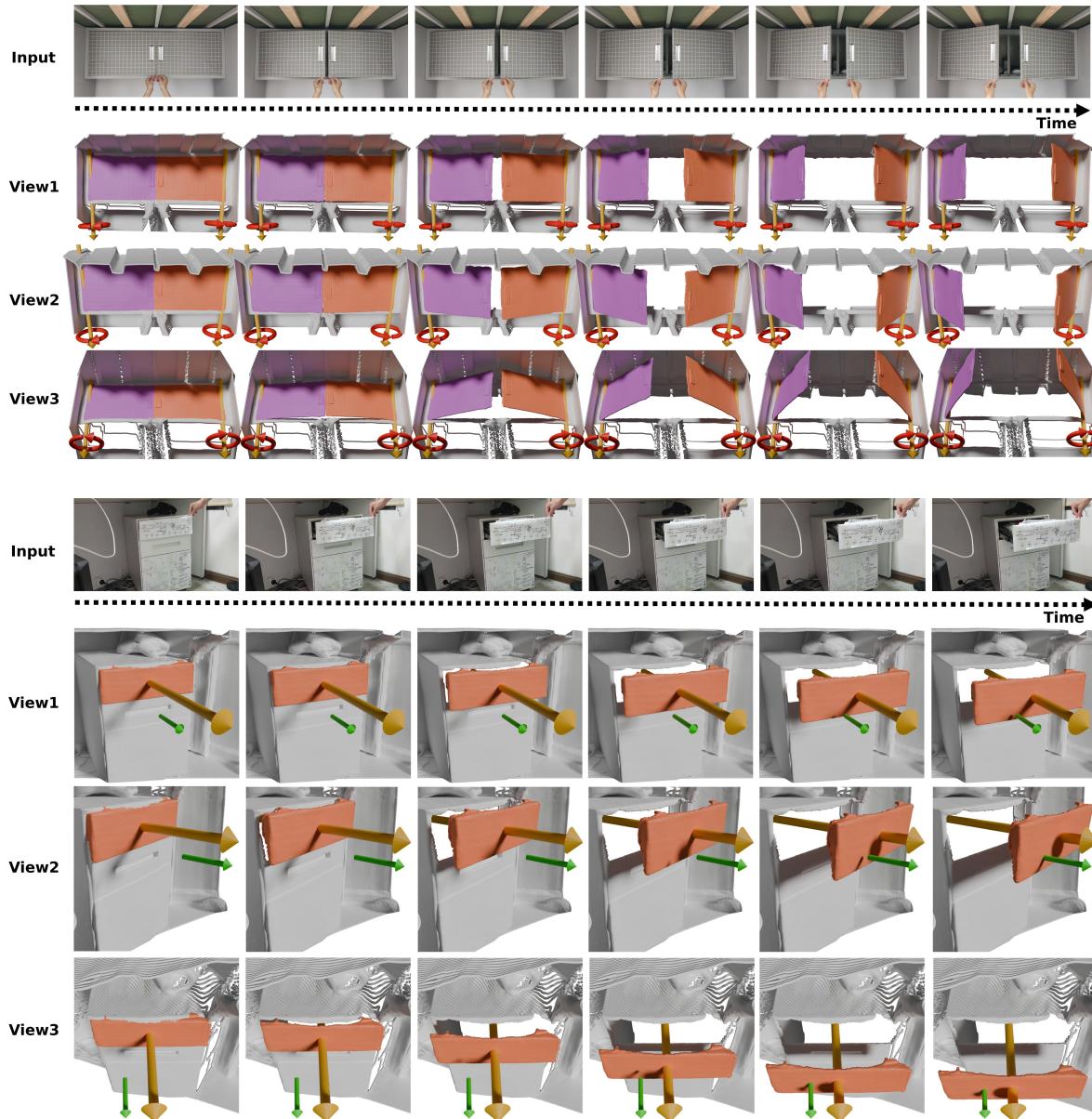


Figure 6. Analysis results of motion parts and their motion attributes for real-world scenes 5-6.

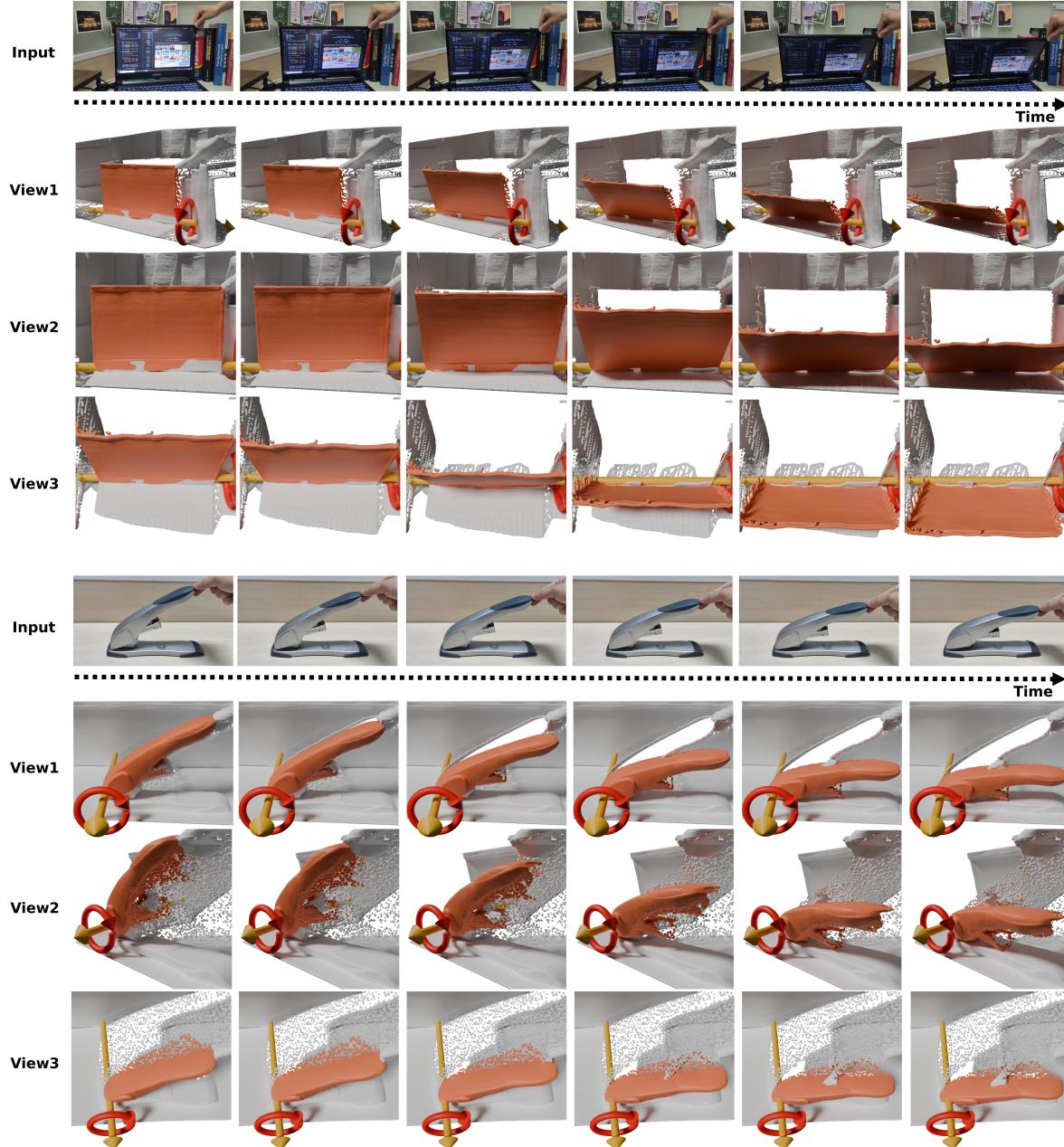


Figure 7. Analysis results of motion parts and their motion attributes for real-world scenes 7-8.

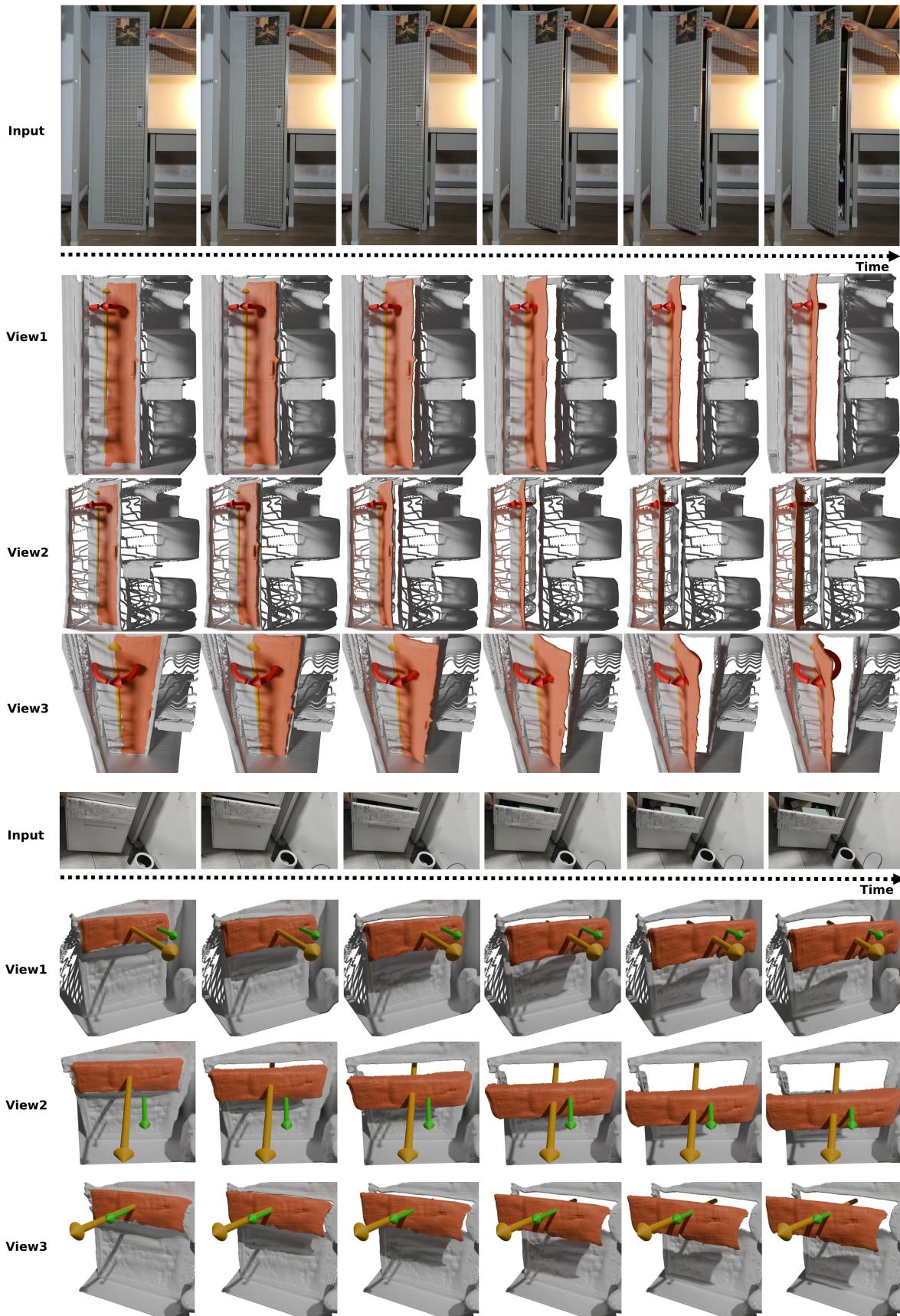


Figure 8. Analysis results of motion parts and their motion attributes for real-world scenes 9-10.

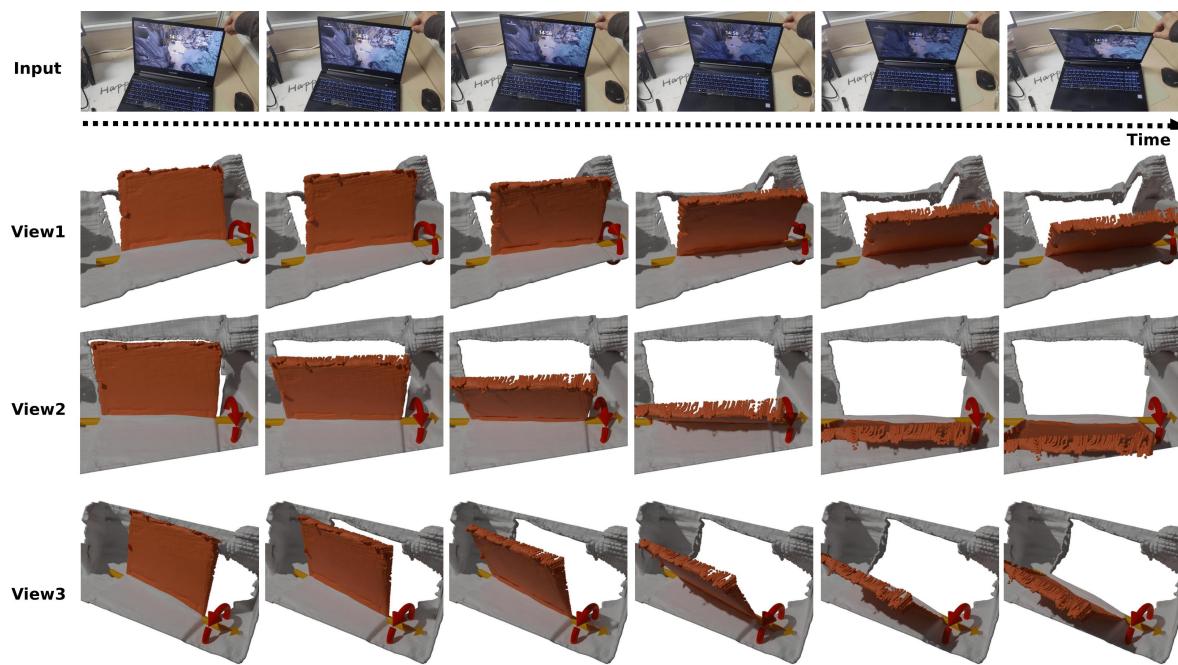


Figure 9. Analysis results of motion parts and their motion attributes for real-world scene 11.

Virtual Data Results

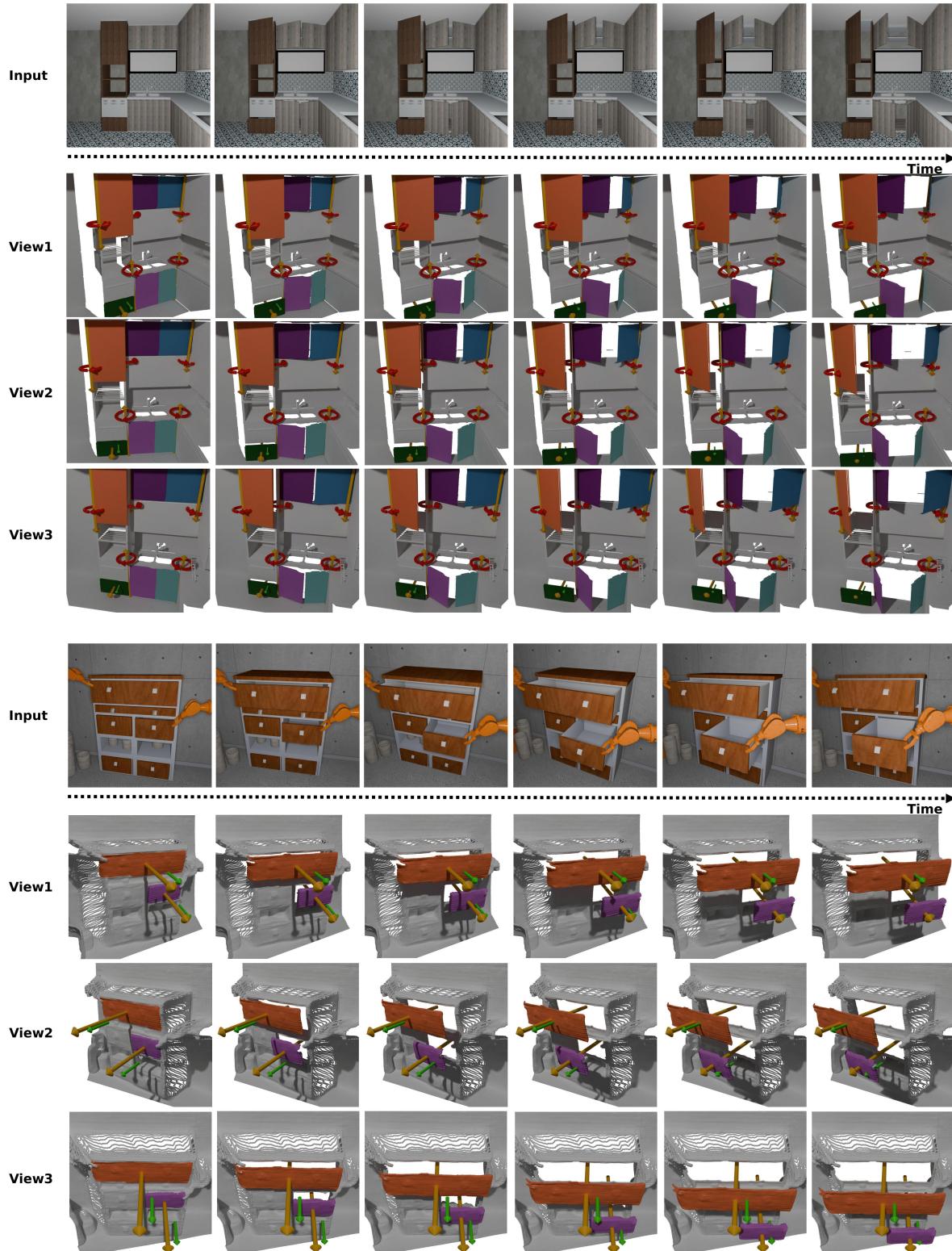


Figure 10. Analysis results of motion parts and their motion attributes for virtual simulation scenes 1-2.

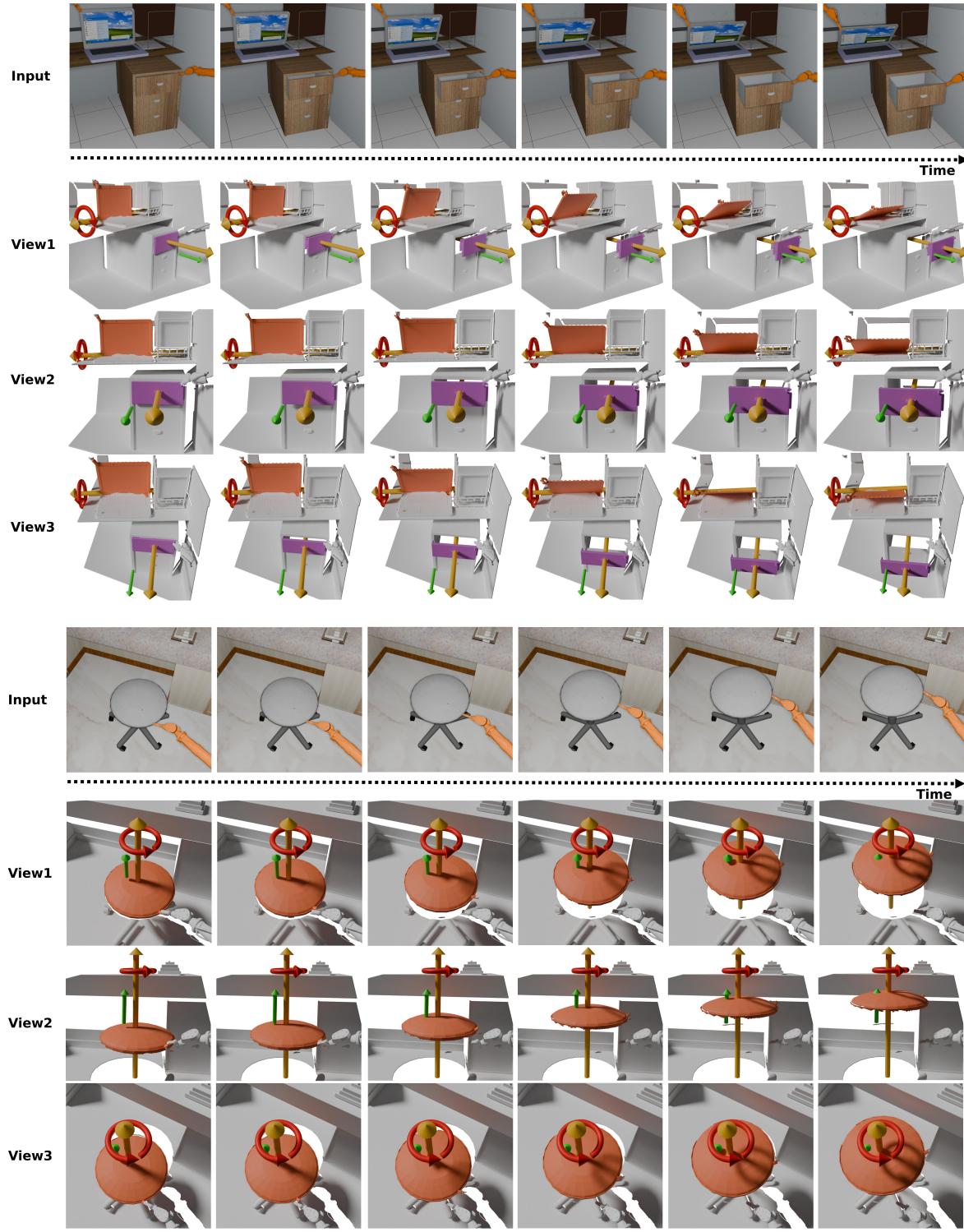


Figure 11. Analysis results of motion parts and their motion attributes for virtual simulation scenes 3-4.

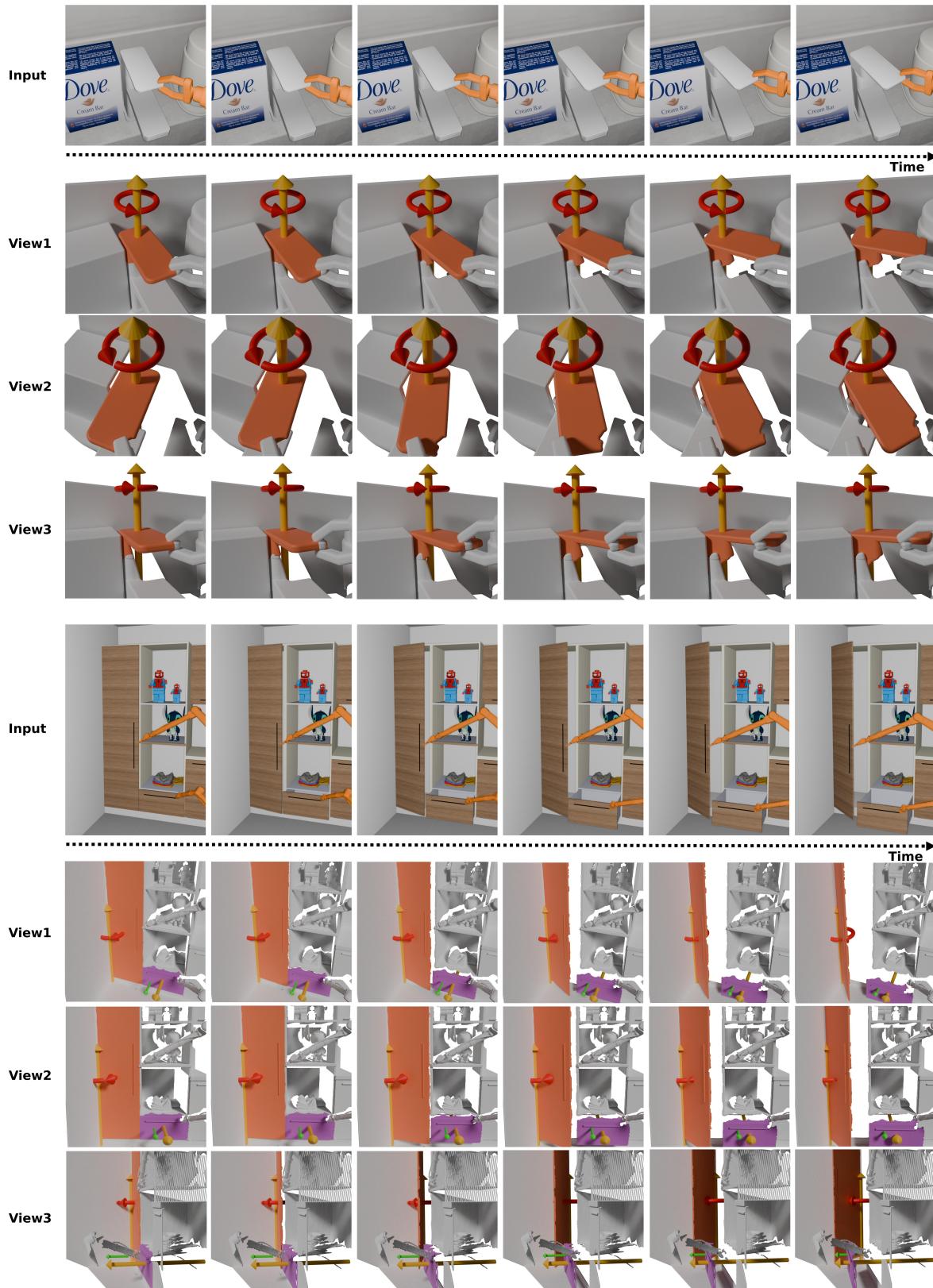


Figure 12. Analysis results of motion parts and their motion attributes for virtual simulation scenes 5-6.

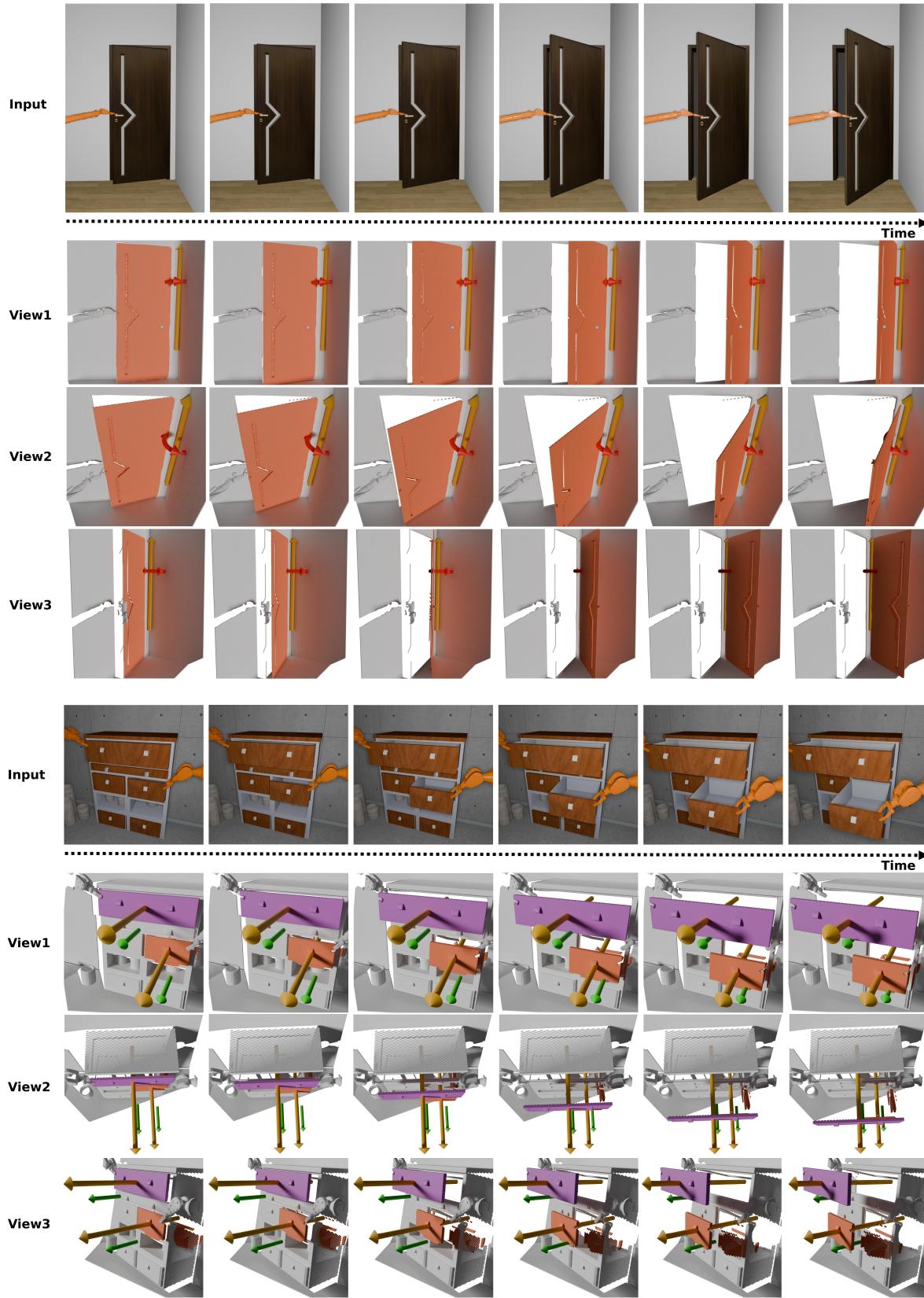


Figure 13. Analysis results of motion parts and their motion attributes for virtual simulation scenes 7-8.

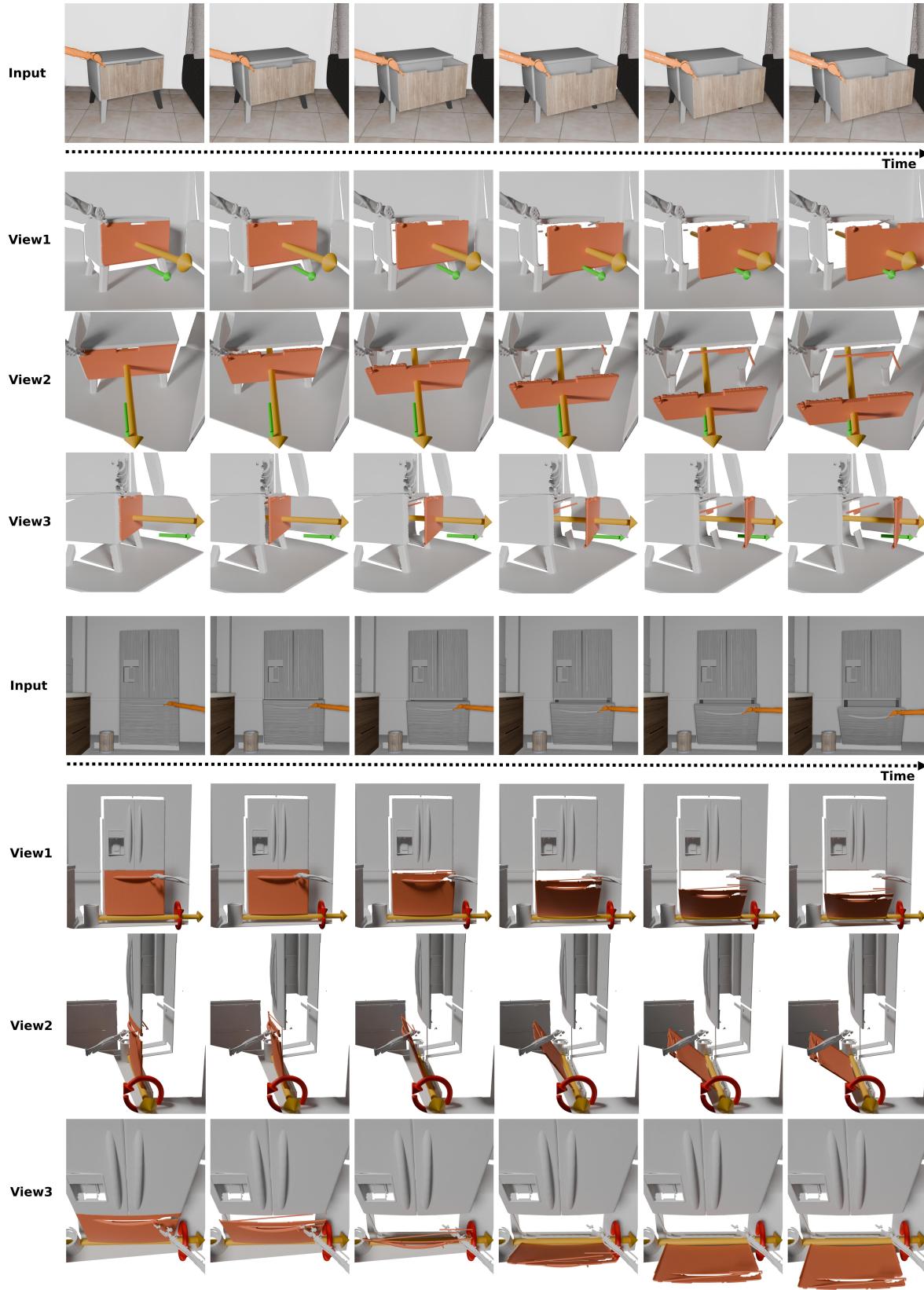


Figure 14. Analysis results of motion parts and their motion attributes for virtual simulation scenes 9-10.

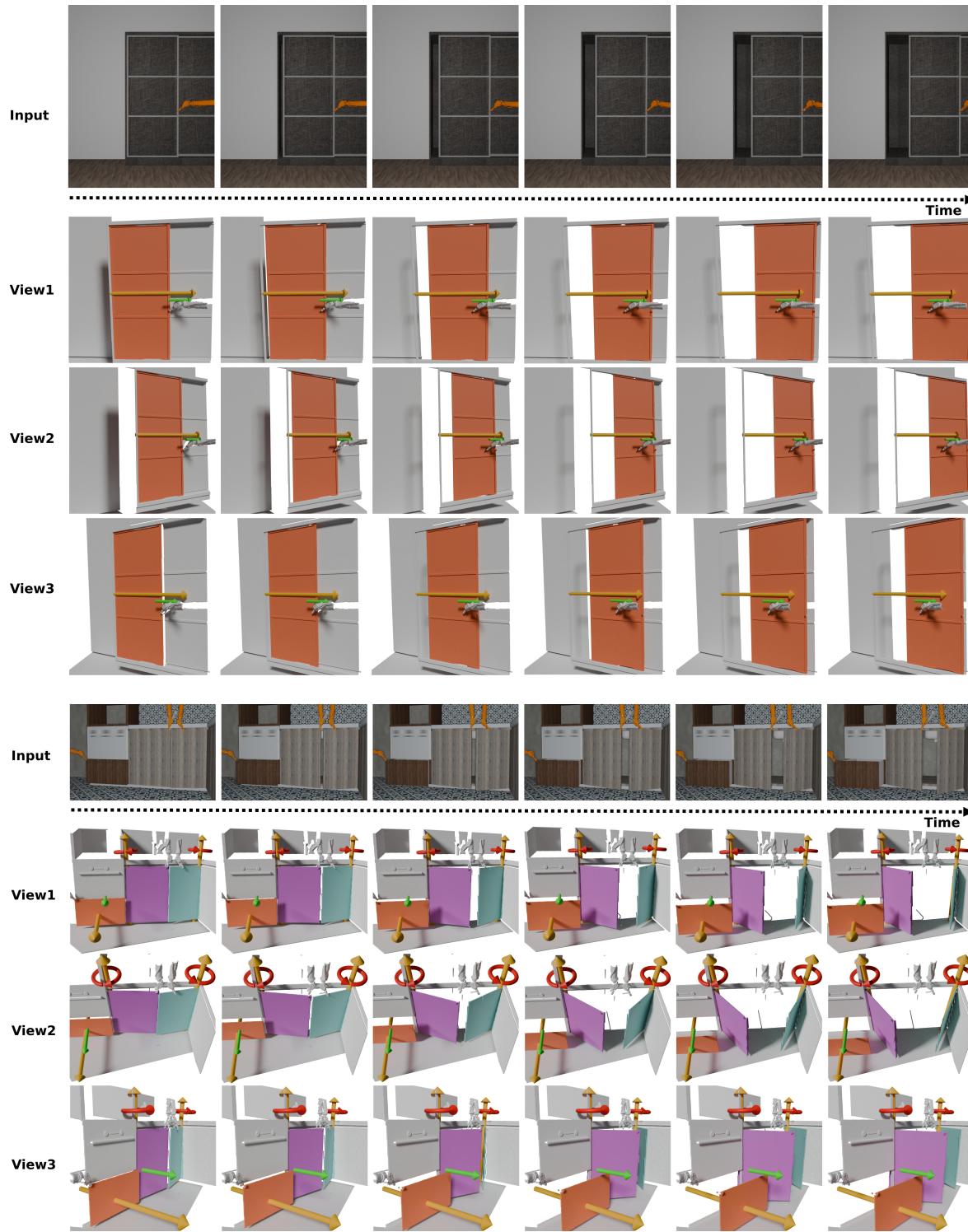


Figure 15. Analysis results of motion parts and their motion attributes for virtual simulation scenes 11-12.

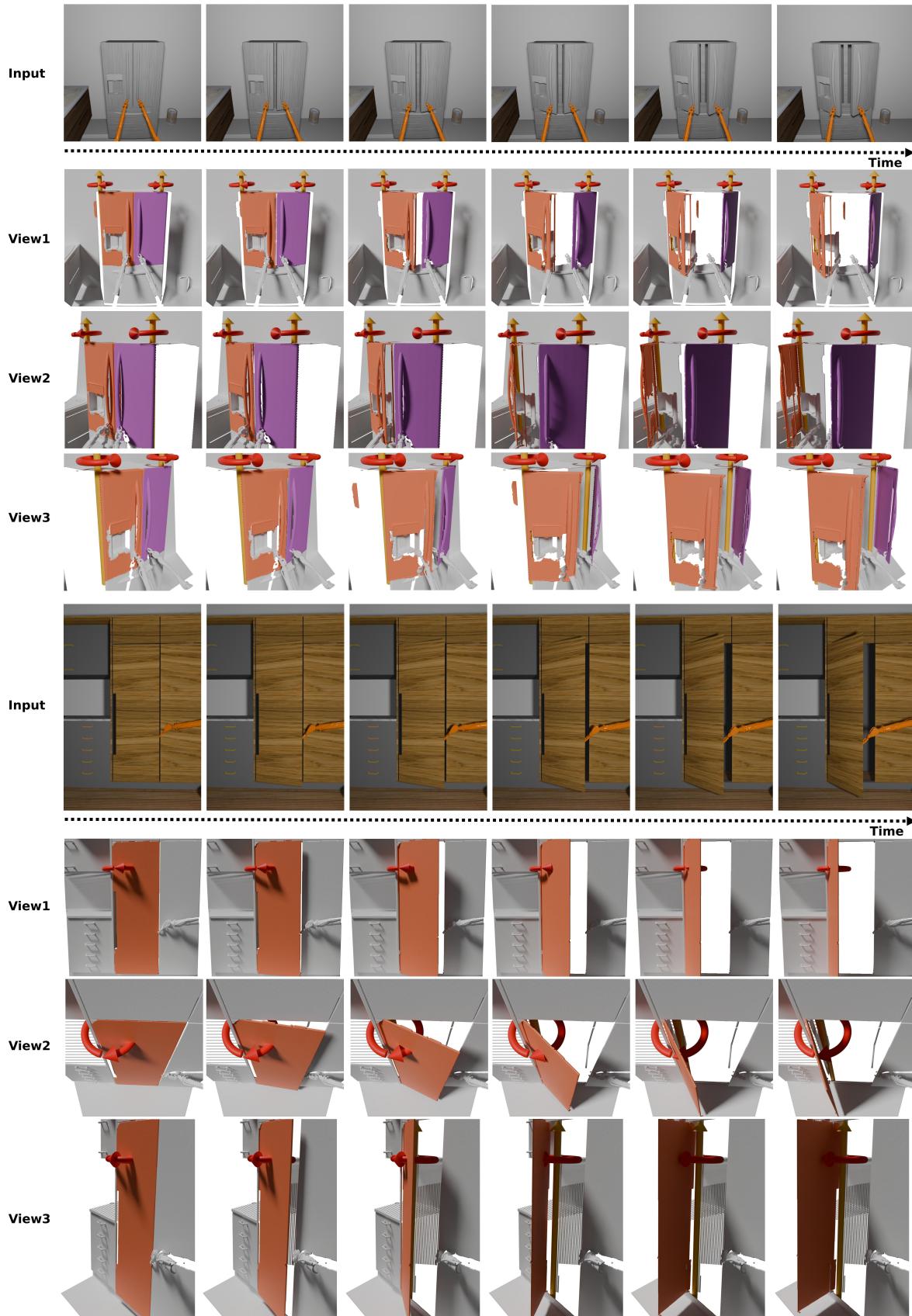


Figure 16. Analysis results of motion parts and their motion attributes for virtual simulation scenes 13-14.

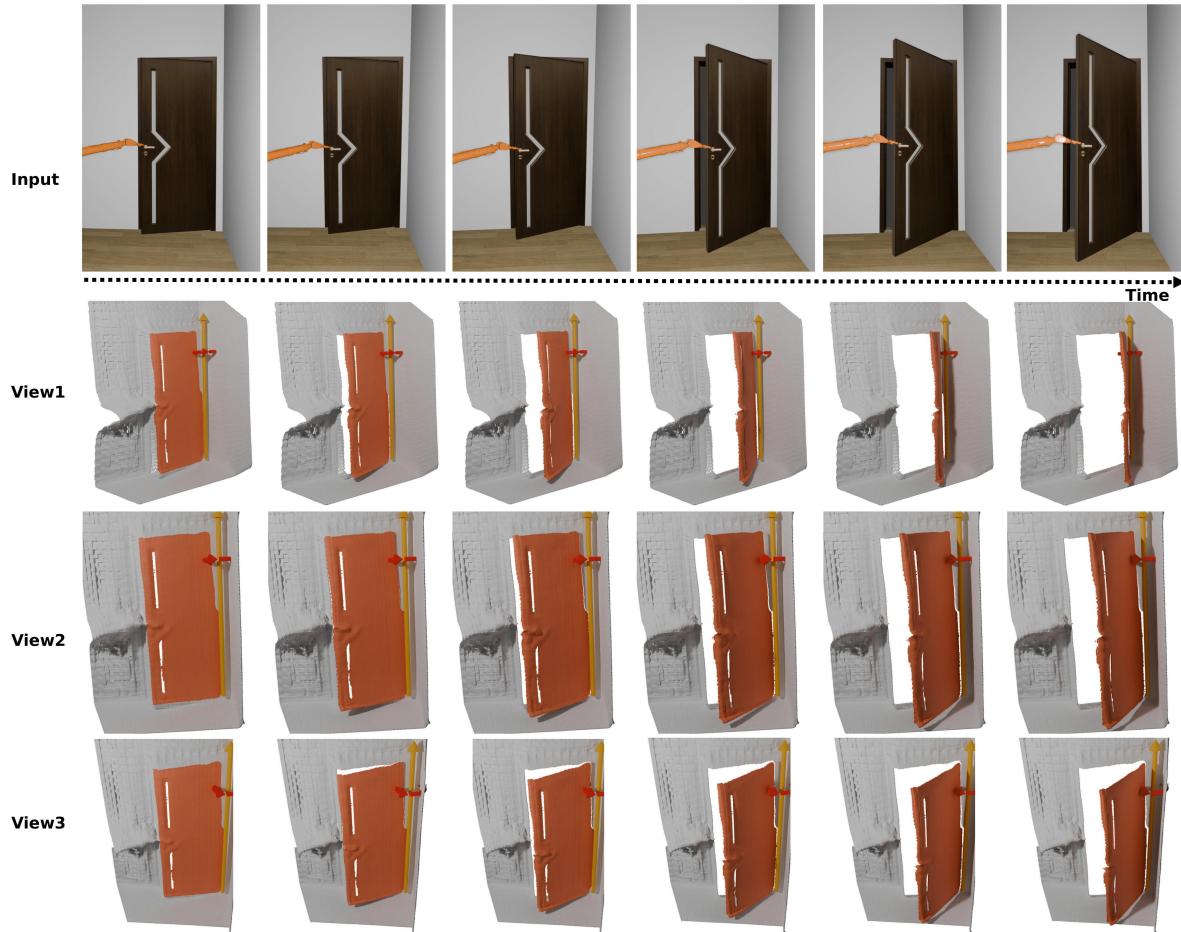


Figure 17. Analysis results of motion parts and their motion attributes for virtual simulation scenes 15.

059

References

- 060 [1] Blender Foundation. Blender: The free and open source 3d
061 creation suite, 2024. [1](#)
- 062 [2] Trimble Inc. 3d warehouse: The largest library of free 3d
063 models, 2024. [1](#)
- 064 [3] Jiayi Liu, Ali Mahdavi-Amiri, and Manolis Savva. Paris: Part-
065 level reconstruction and motion analysis for articulated ob-
066 jects. In *Proceedings of the IEEE/CVF International Confer-
067 ence on Computer Vision*, pages 352–363, 2023. [1](#)
- 068 [4] Johannes L. Schoenberger and Jan-Michael Frahm. Colmap:
069 A general-purpose structure-from-motion and multi-view
070 stereo pipeline, 2024. [1](#)
- 071 [5] Colton Stearns, Adam Harley, Mikaela Uy, Florian Dubost,
072 Federico Tombari, Gordon Wetzstein, and Leonidas Guibas.
073 Dynamic gaussian marbles for novel view synthesis of casual
074 monocular videos. *arXiv preprint arXiv:2406.18717*, 2024. [1](#)