



# Animate Your Thoughts: Reconstruction of Dynamic Natural Vision from Human Brain Activity

• Yizhuo Lu<sup>1,2</sup>, Changde Du<sup>\*1</sup>, Chong Wang<sup>4</sup>, Xuanliu Zhu<sup>5</sup>, Liuyun Jiang<sup>1,2</sup>, Xujin Li<sup>1,2</sup>, Huiguang He<sup>1,2,3</sup>,

• <sup>1</sup> State Key Laboratory of Brain Cognition and Brain-inspired Intelligence Technology, CASIA,

• <sup>2</sup> School of Future Technology, University of Chinese Academy of Sciences, <sup>3</sup> School of Artificial Intelligence, UCAS,

• <sup>4</sup> School of Computer and Artificial Intelligence, Zhengzhou University, <sup>5</sup> Beijing University of Posts and Telecommunications

• Email: [luyizhuo2023@ia.ac.cn](mailto:luyizhuo2023@ia.ac.cn), [huiguang.he@ia.ac.cn](mailto:huiguang.he@ia.ac.cn)



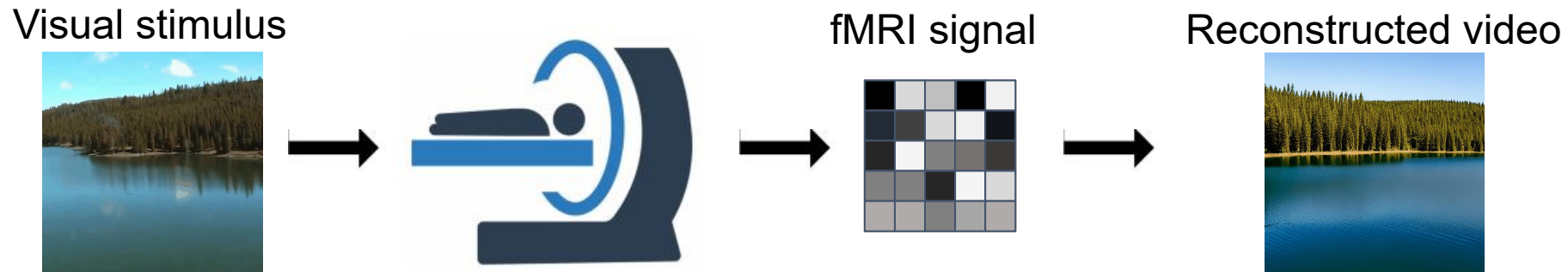
**中国科学院自动化研究所**<sup>1</sup>  
INSTITUTE OF AUTOMATION  
CHINESE ACADEMY OF SCIENCES



**中国科学院大学**<sup>2</sup>  
University of Chinese Academy of Sciences

# Introduction

---

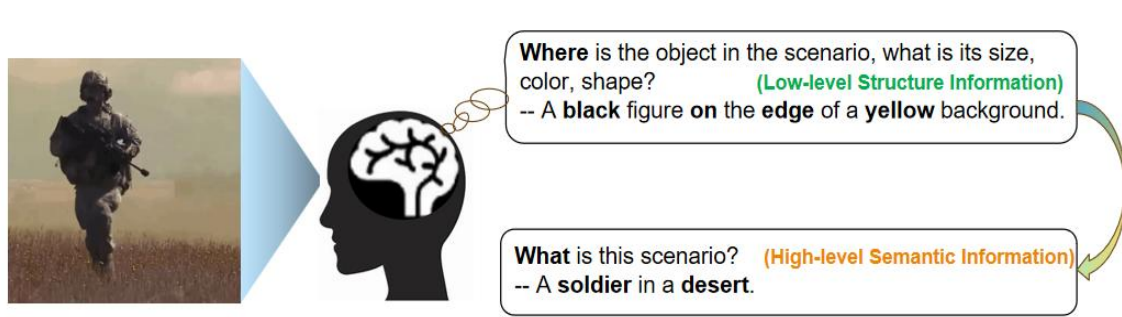


**(1) Method:** We propose Mind-Animator, which enables video reconstruction by decoupling semantic, structural, and motion information from fMRI data for the first time.

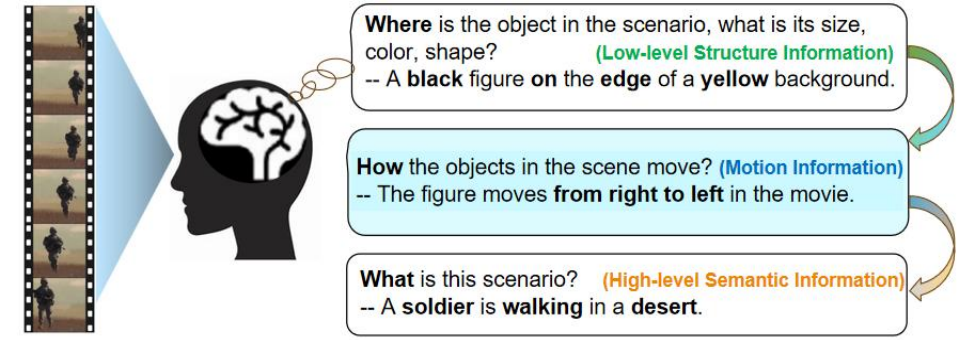
**(2) Interpretability:** We use voxel-wise and ROI-wise visualization techniques to elucidate the interpretability of our proposed model from a neurobiological perspective.

**(3) Comprehensive evaluation:** We introduce eight evaluation metrics that comprehensively assess the reconstruction results of our model and all previous models across three dimensions—semantic, structure, and spatiotemporal consistency—on three publicly available video-fMRI datasets. This establishes our work as the first unified benchmark for subsequent researchers. We will release all data and code to facilitate future research.

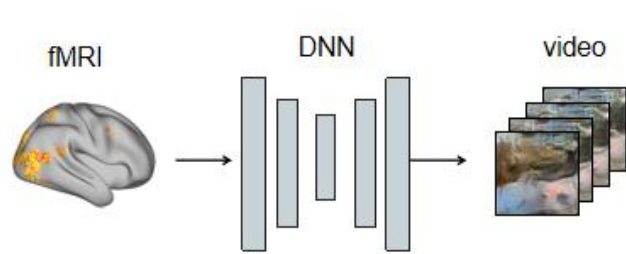
# Motivation and Related works



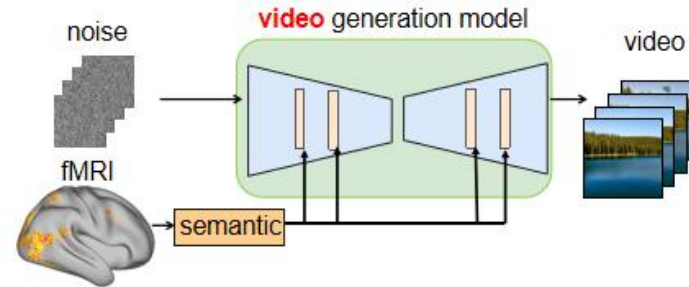
(a) Static image stimulus



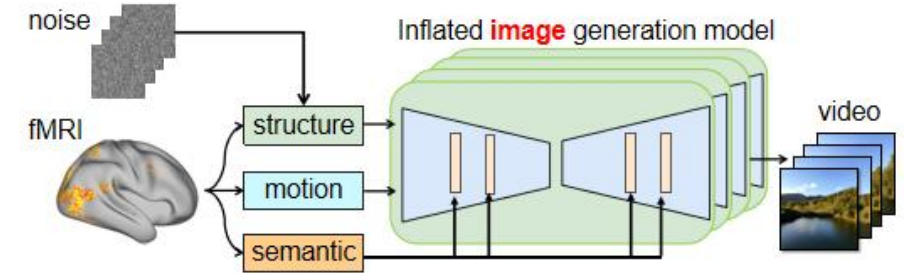
(b) Dynamic video stimulus



(c) End-to-end models



(d) Video generation model based



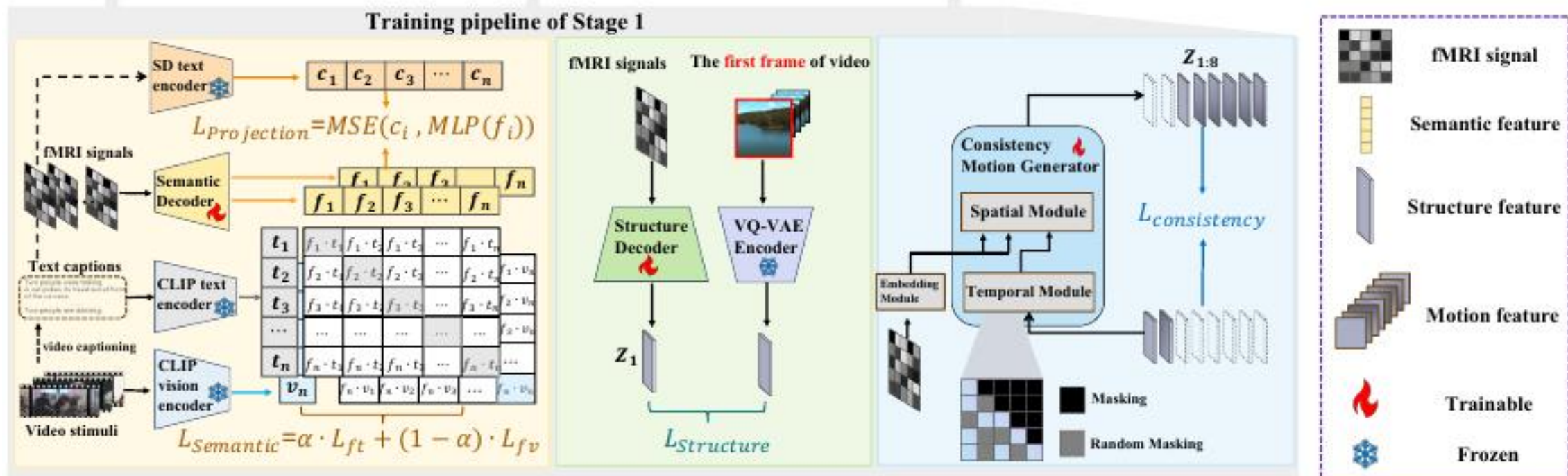
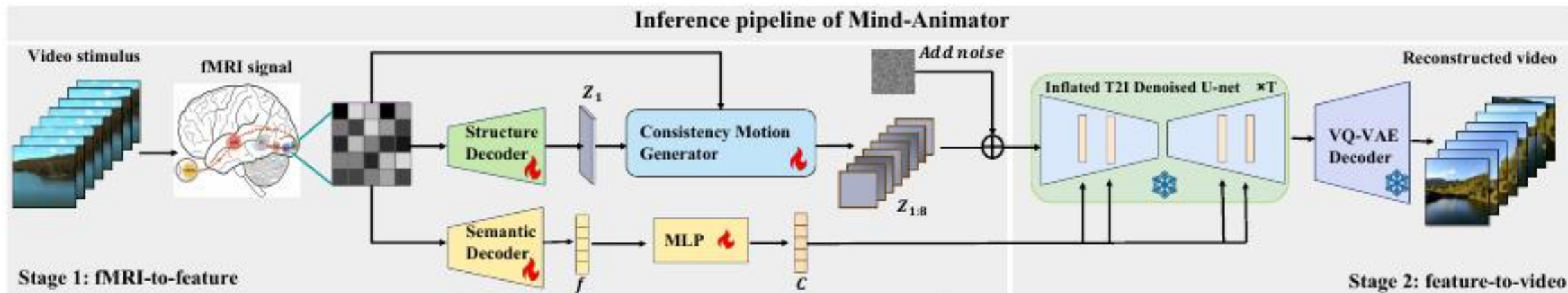
(e) Our model

## Existing Issues:

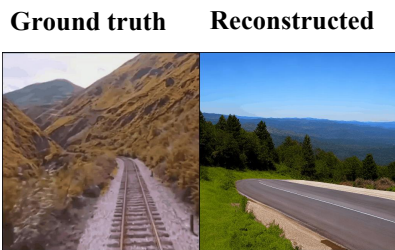
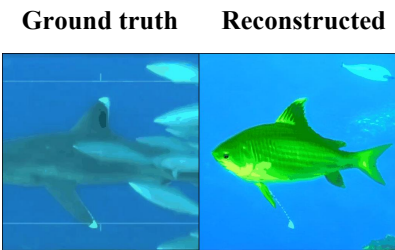
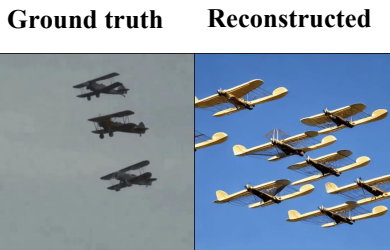
- (1) Unable to accurately model semantic, structural, and motion features.
- (2) External video data (e.g., motion priors learned by video generation models) interfere with the motion features in the reconstructed video.



# Methodology

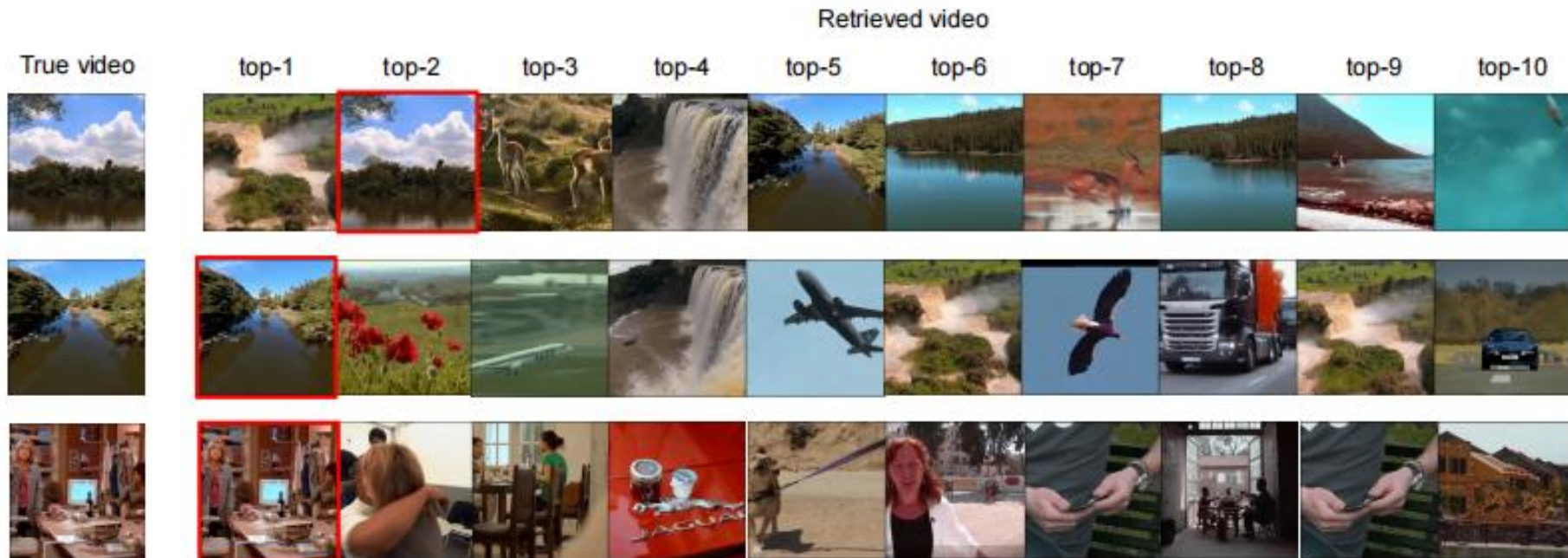


# Results (reconstruction task)





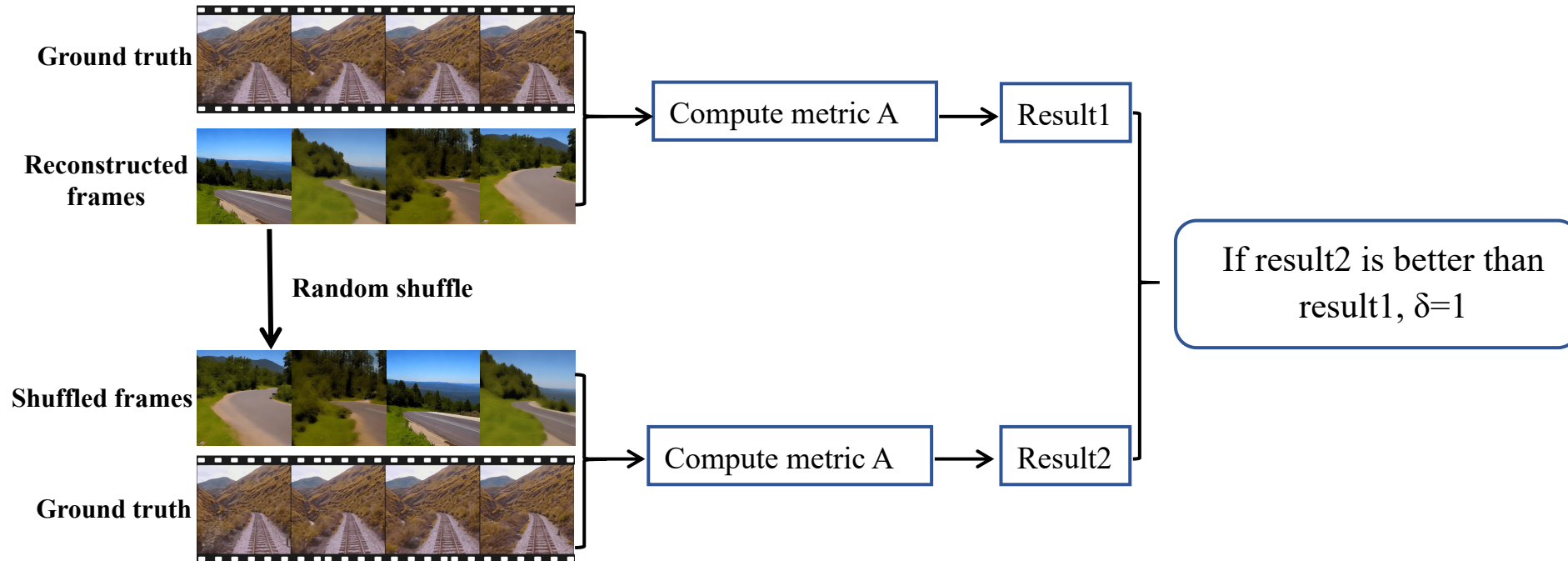
# Results (retrieval task)



Dataset		CC2017							
			Subjet1		Subjet2		Subjet3		Average
Model	Test set	top-10	top-100	top-10	top-100	top-10	top-100	top-10	top-100
Wen (Wen et al. (2018))	Small	2.17 <sub>*</sub>	19.50 <sub>*</sub>	3.33 <sub>*</sub>	19.17 <sub>*</sub>	—	—	2.75 <sub>*</sub>	19.33 <sub>*</sub>
Kupershmidt (Kupershmidt et al. (2022))	Small	1.09 <sub>*</sub>	8.57 <sub>*</sub>	0.92 <sub>*</sub>	8.24 <sub>*</sub>	0.84 <sub>*</sub>	8.24 <sub>*</sub>	0.95 <sub>*</sub>	8.35 <sub>*</sub>
Mind-video (Chen et al. (2024))	Small	<b>3.22<sub>*</sub></b>	19.08 <sub>*</sub>	2.75 <sub>*</sub>	16.83 <sub>*</sub>	3.58 <sub>*</sub>	22.08 <sub>*</sub>	3.18 <sub>*</sub>	19.33 <sub>*</sub>
Ours	Small	3.08	<b>22.58</b>	<b>4.75</b>	<b>26.90</b>	<b>4.50</b>	<b>24.67</b>	<b>4.11</b>	<b>24.72</b>
Wen (Wen et al. (2018))	Large	1.41 <sub>*</sub>	11.58 <sub>*</sub>	2.08 <sub>*</sub>	9.58 <sub>*</sub>	—	—	1.75 <sub>*</sub>	10.58 <sub>*</sub>
Kupershmidt (Kupershmidt et al. (2022))	Large	0.17 <sub>*</sub>	2.94 <sub>*</sub>	0.17 <sub>*</sub>	2.77 <sub>*</sub>	0.25 <sub>*</sub>	2.18 <sub>*</sub>	0.19 <sub>*</sub>	2.63 <sub>*</sub>
Mind-video (Chen et al. (2024))	Large	1.75 <sub>*</sub>	7.17 <sub>*</sub>	0.83 <sub>*</sub>	5.17 <sub>*</sub>	1.25 <sub>*</sub>	9.00 <sub>*</sub>	1.28 <sub>*</sub>	7.11 <sub>*</sub>
Ours	Large	<b>2.17</b>	<b>12.50</b>	<b>2.25</b>	<b>17.00</b>	<b>2.75</b>	<b>16.42</b>	<b>2.39</b>	<b>15.31</b>

# Interpretability

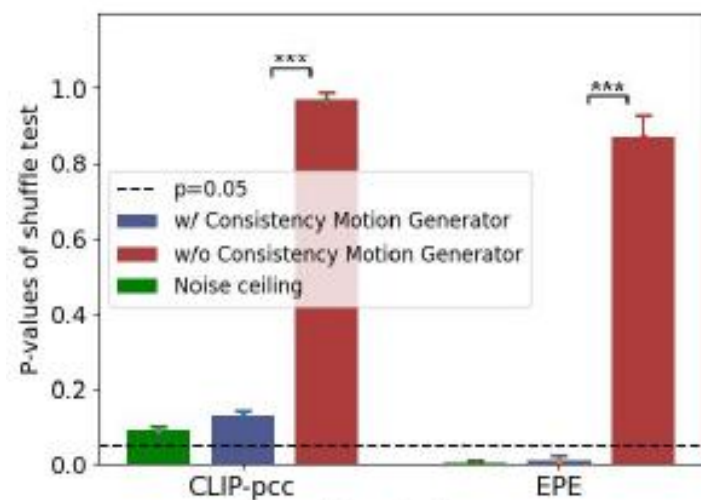
- Have we truly decoded motion information from fMRI?



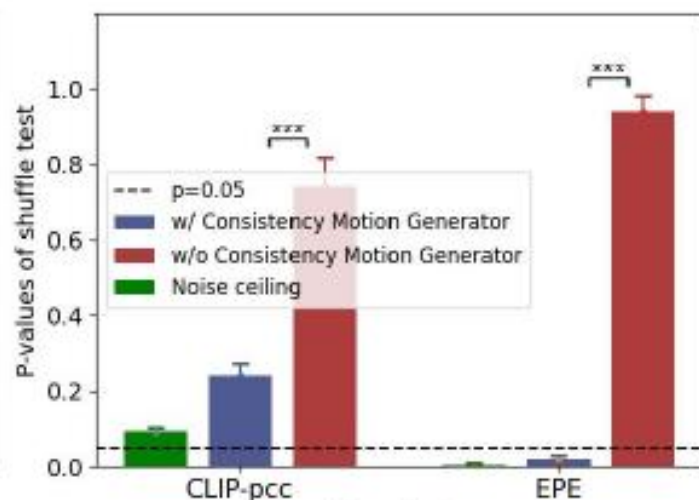
Repeating the aforementioned procedure 100 times, the P-value of the shuffle test can be estimated as  $P = \sum_{i=1}^{100} \delta_i / 100$ .

A lower P-value indicates a higher consistency between the reconstructed video frames and the ground truth before shuffling.

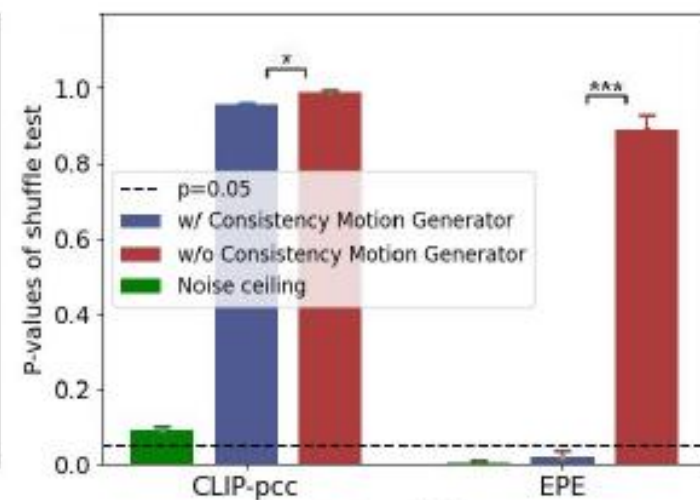
# Interpretability



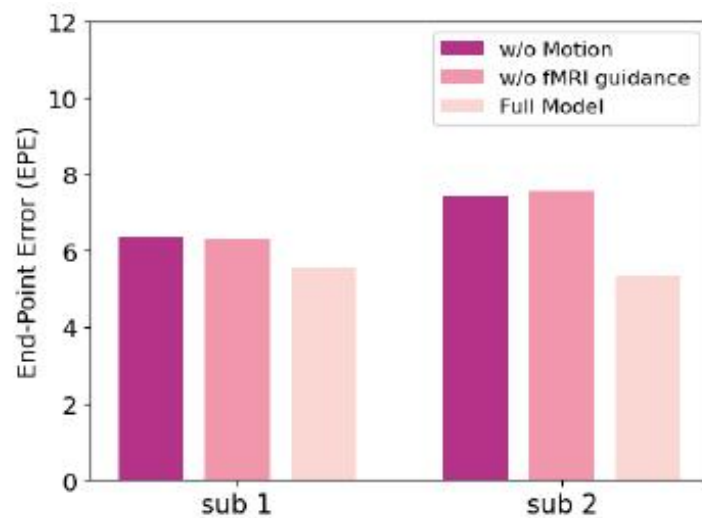
(a) sub 1



(b) sub 2



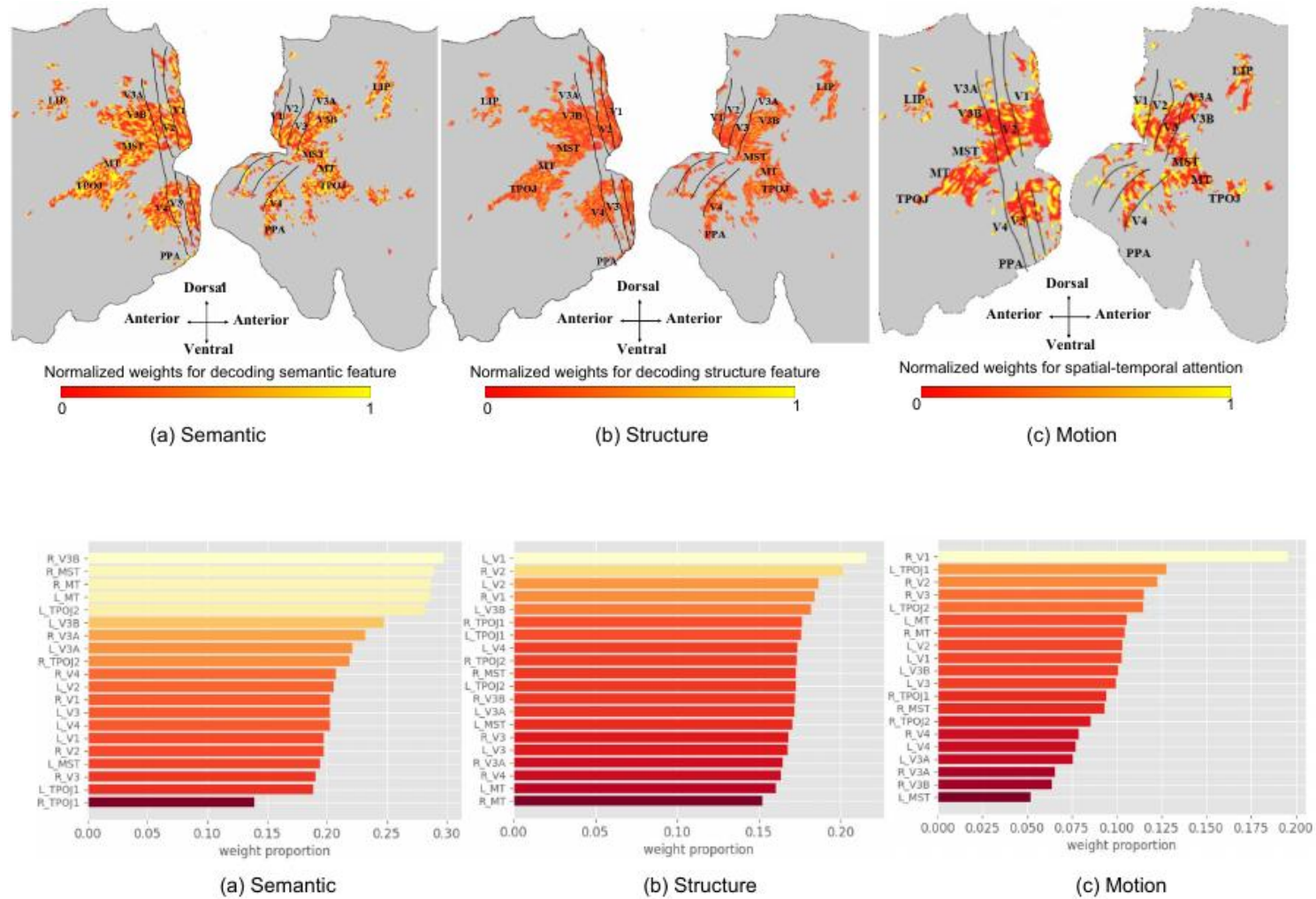
(c) sub 3





# Interpretability

- Which brain regions are responsible for decoding different features, respectively?



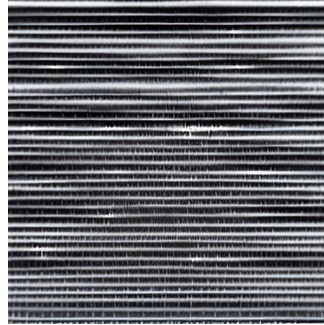
# Ablation Study

---

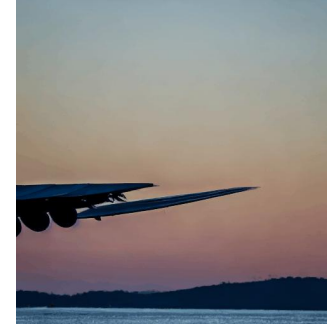
**Ground truth    Reconstructed**



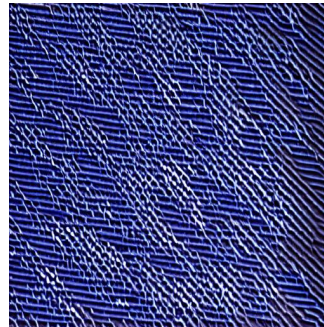
**Without Semantic**



**Without Structural**



**Without Motion**



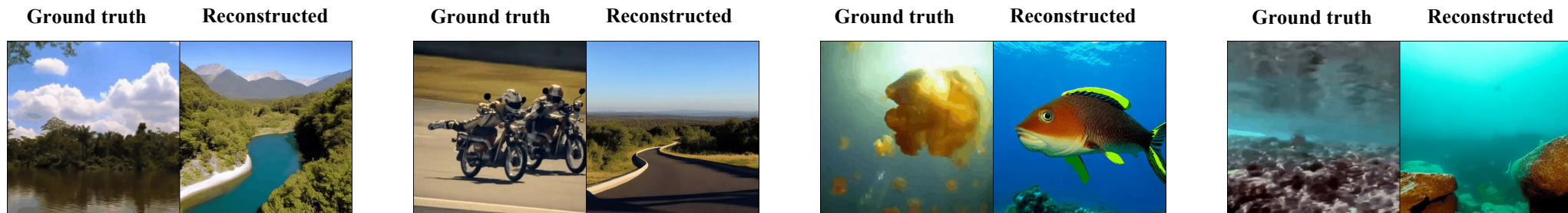
# Fail Cases

---

**Case 1:** Decoding errors in semantic, structural, or motion features due to low decoding accuracy.



**Case 2:** The data acquisition paradigm causes abrupt content transitions at the boundaries of video clips, which are uniformly segmented from the complete videos viewed by the subjects during data collection.





Thanks!