

**BrainLMM: A Label-Free Framework for Mapping Multi-Semantic
Representation in the Human Visual Cortex**

Anonymous submission

ResNet50_{clip}



Figure 1: **Cakes with text patterns.** We employ ResNet50_{CLIP} as the backbone and select the top 5 images exhibiting the highest activations from a subset of voxels within the VWFA region, where the label with the highest score corresponds to “cake”. As illustrated in the images, the majority of these cakes are characterized by distinct textual patterns, highlighting the strong correlation between the VWFA region and textual features.

ViT-32_{clip}

Voxel: 53



Voxel: 54



Voxel: 117



Voxel: 149



Voxel: 157



Voxel: 165



Voxel: 166



Voxel: 167



Voxel: 168



Voxel: 177



Voxel: 178



Voxel: 186



Voxel: 221



Voxel: 237



Voxel: 248



Figure 2: **Cakes with text patterns.** We employ ViT-32_{CLIP} as the backbone and select the top 5 images exhibiting the highest activations from a subset of voxels within the VWFA region, where the label with the highest score corresponds to “cake”. As illustrated in the images, the majority of these cakes are characterized by distinct textual patterns, highlighting the strong correlation between the VWFA region and textual features.

NSD ResNet50_{clip}

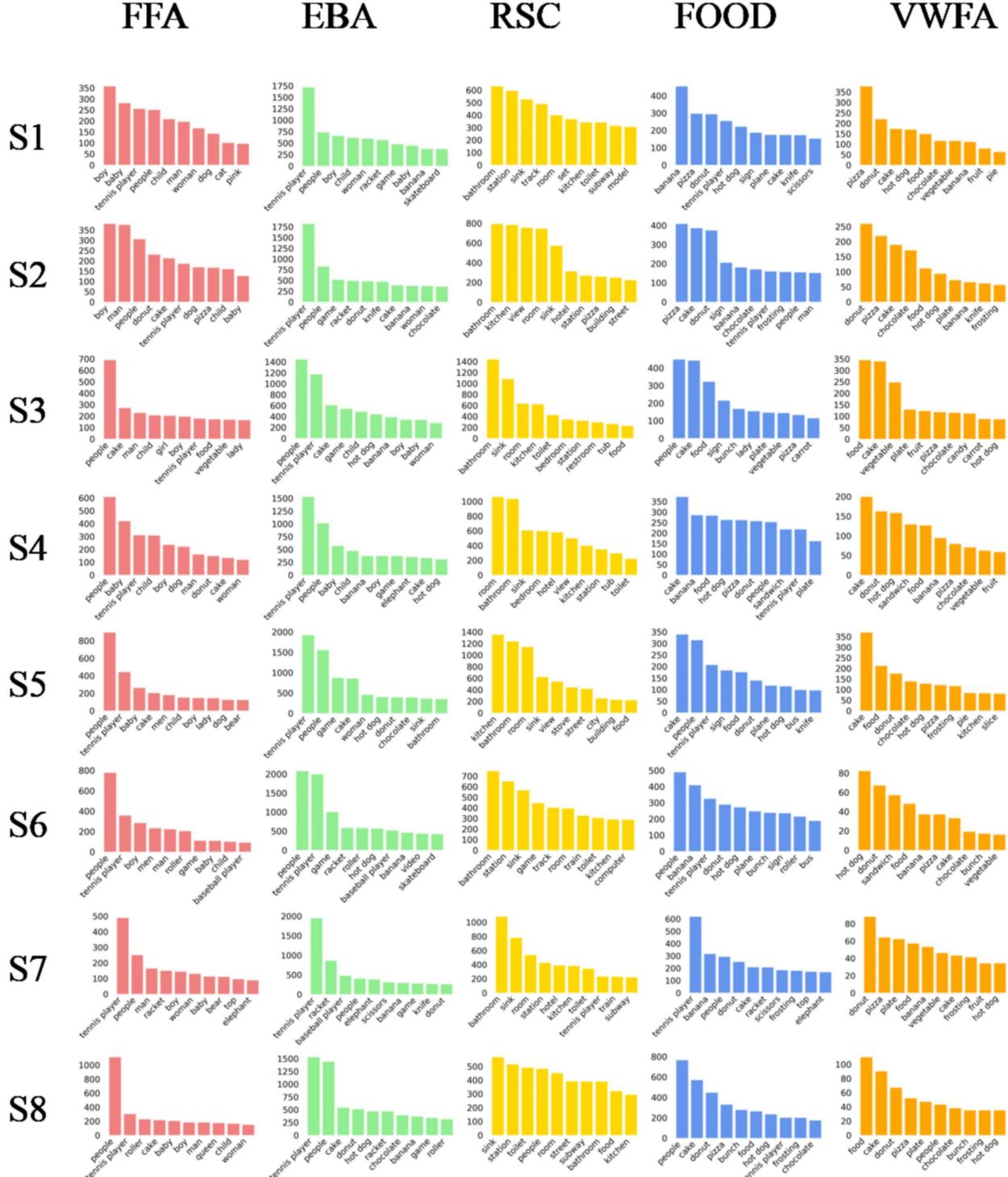


Figure 3: **The BrainLMM enables precise dissection of high-level visual categories in the selective regions.** Single semantic mapping captured via BrainLMM of all subjects in the NSD for ResNet50_{CLIP} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

NSD ViT-32_{clip}

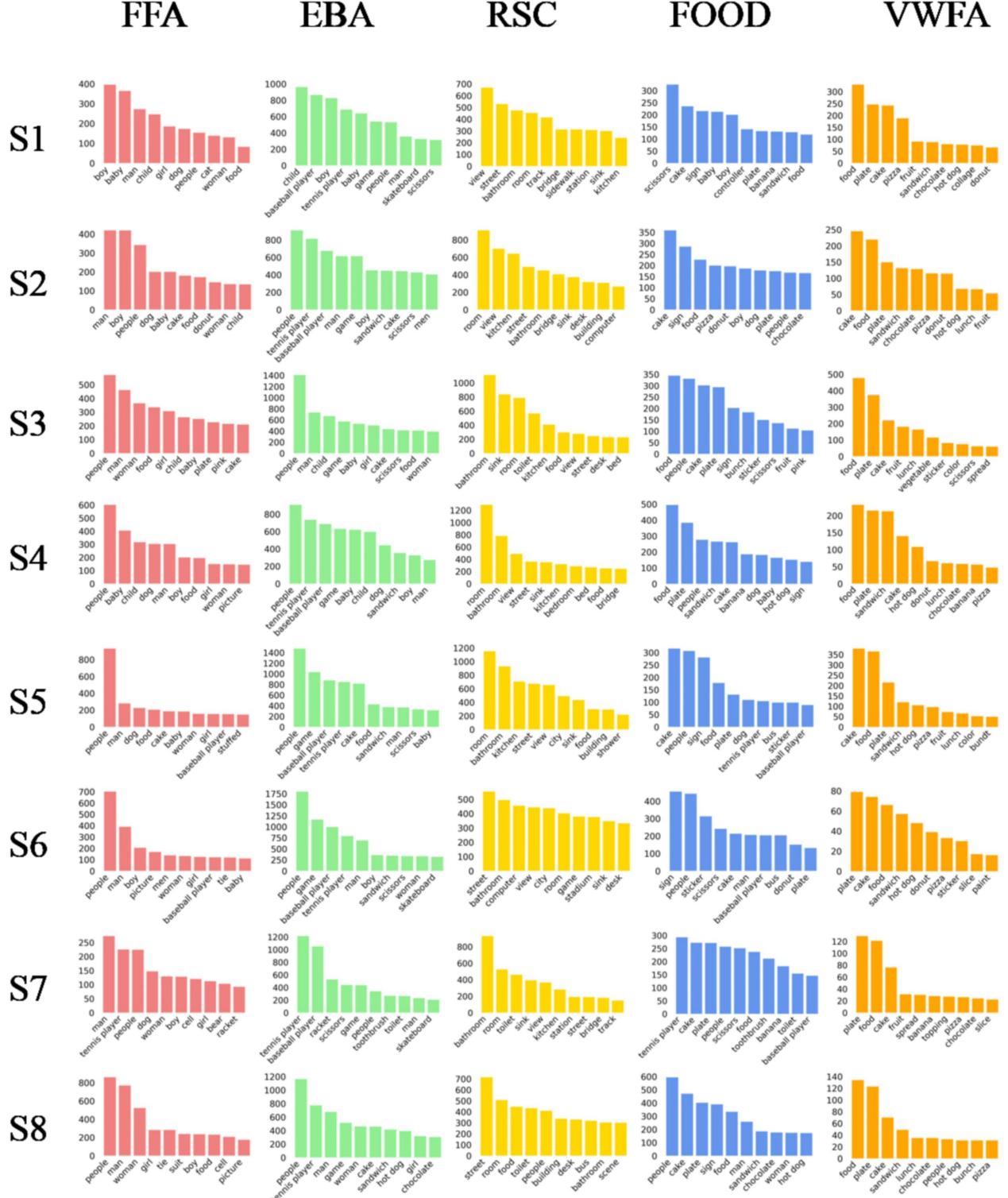


Figure 4: **The BrainLMM enables precise dissection of high-level visual categories in the selective regions.** Single semantic mapping captured via BrainLMM of all subjects in the NSD for ViT-32_{CLIP} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

NSD AlexNet_{ImageNet}

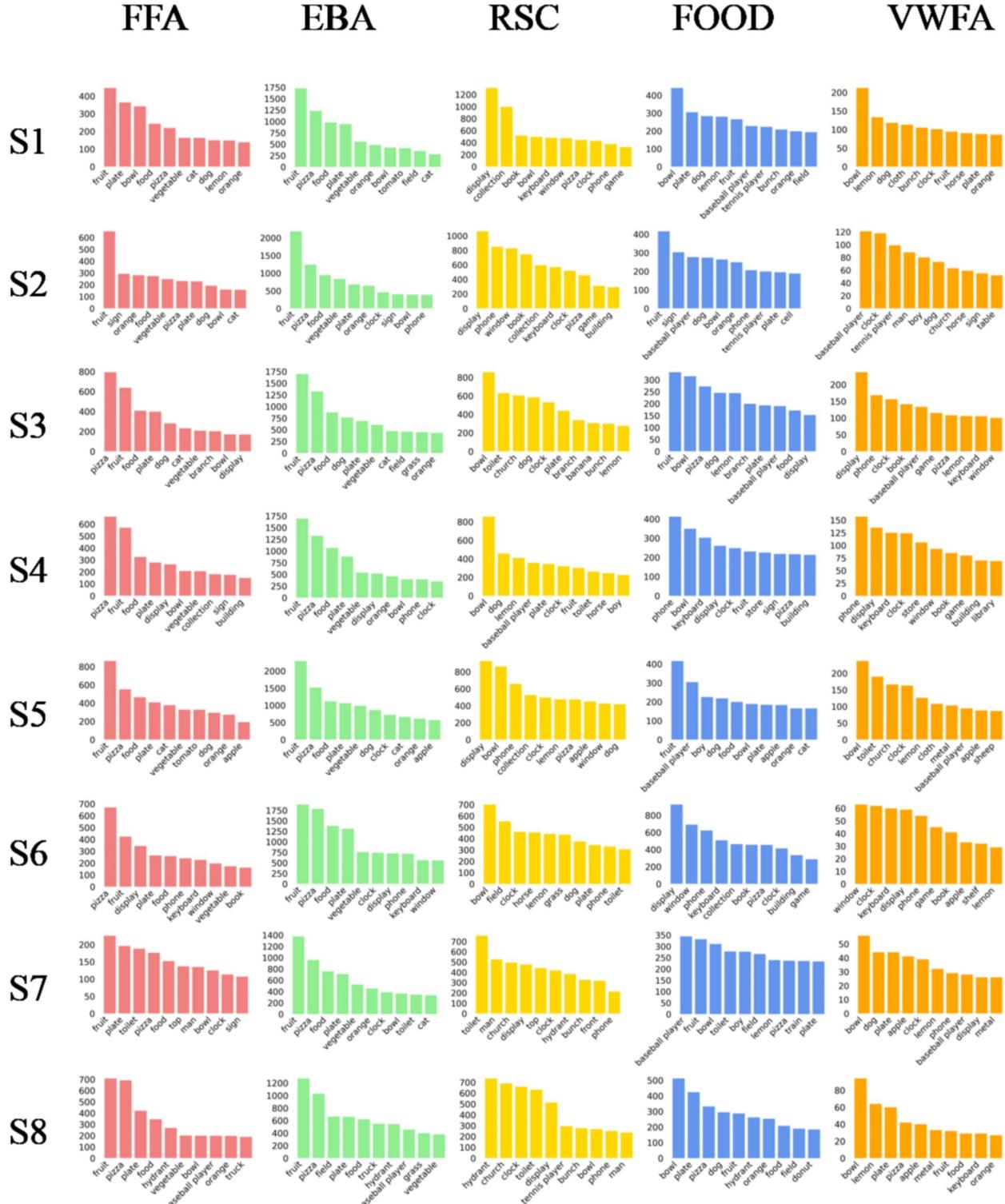


Figure 5: The BrainLMM enables precise dissection of high-level visual categories in the selective regions. Single semantic mapping captured via BrainLMM of all subjects in the NSD for AlexNet_{ImageNet} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

NSD ResNet50_{ImageNet}

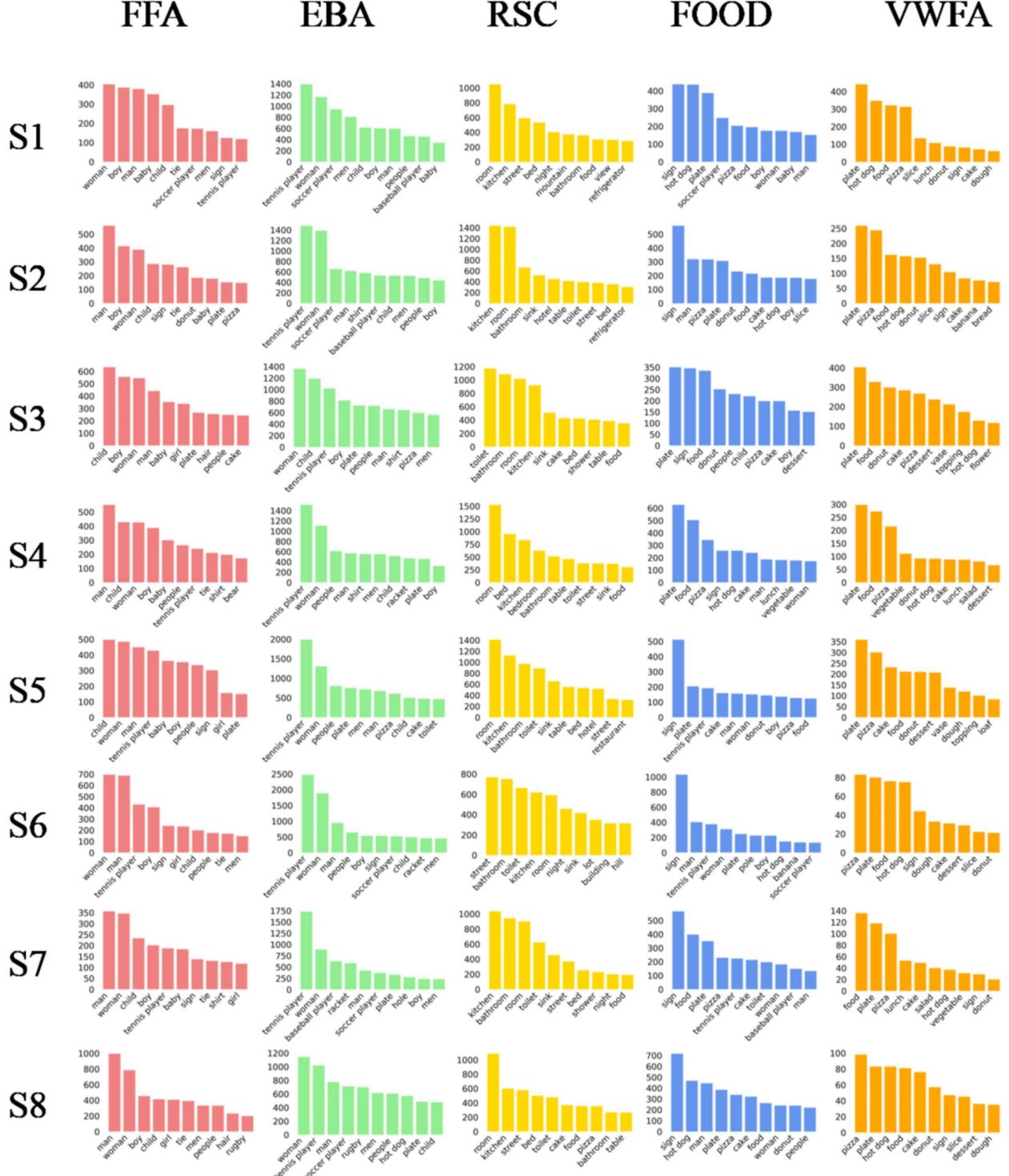


Figure 6: **The BrainLMM enables precise dissection of high-level visual categories in the selective regions.** Single semantic mapping captured via BrainLMM of all subjects in the NSD for ResNet50_{ImageNet} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

NOD ResNet50_{clip}

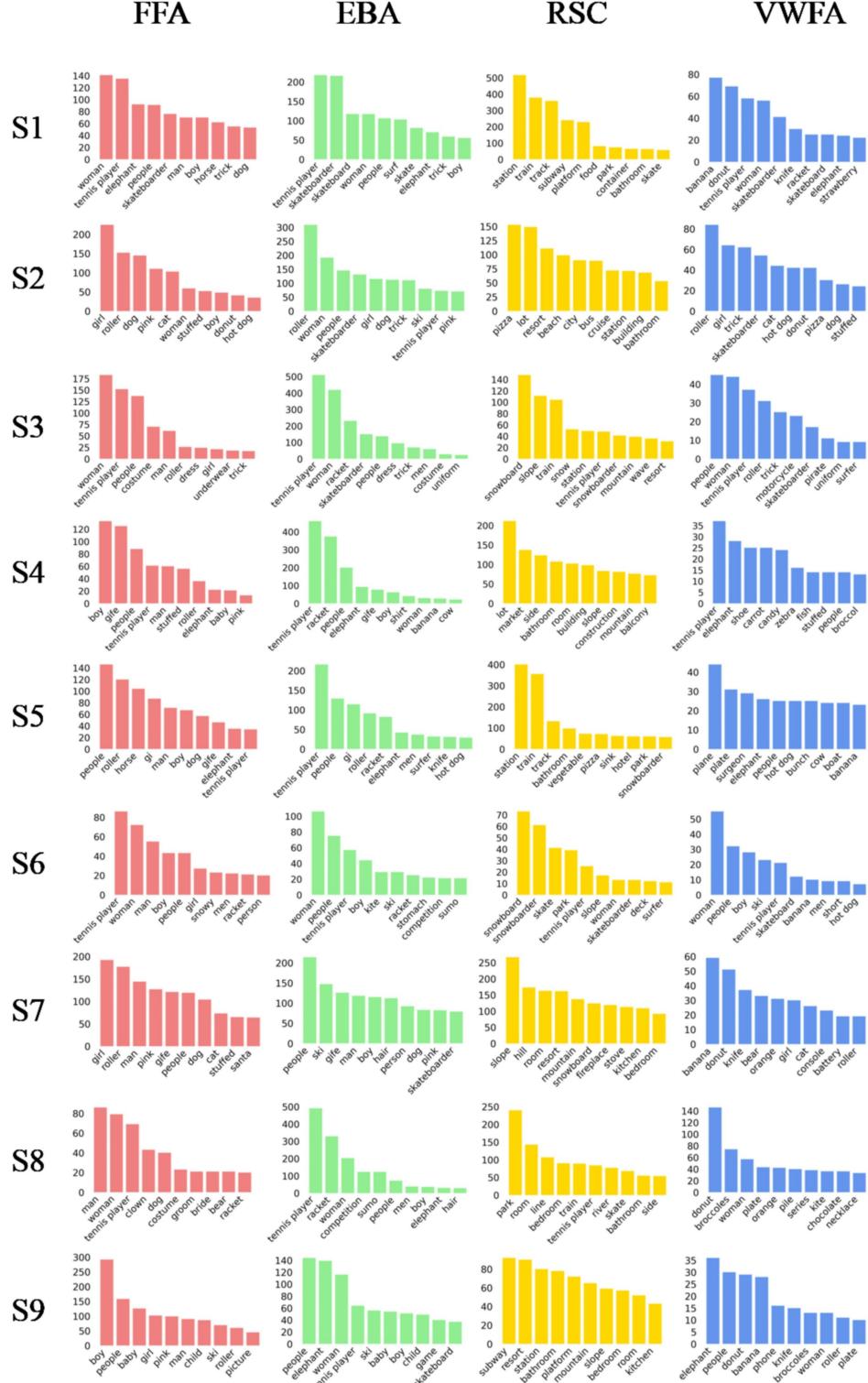


Figure 7: **The BrainLMM enables precise dissection of high-level visual categories in the selective regions.** Single semantic mapping captured via BrainLMM of all subjects in the NOD for ResNet50_{CLIP} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

NOD ViT-32_{clip}

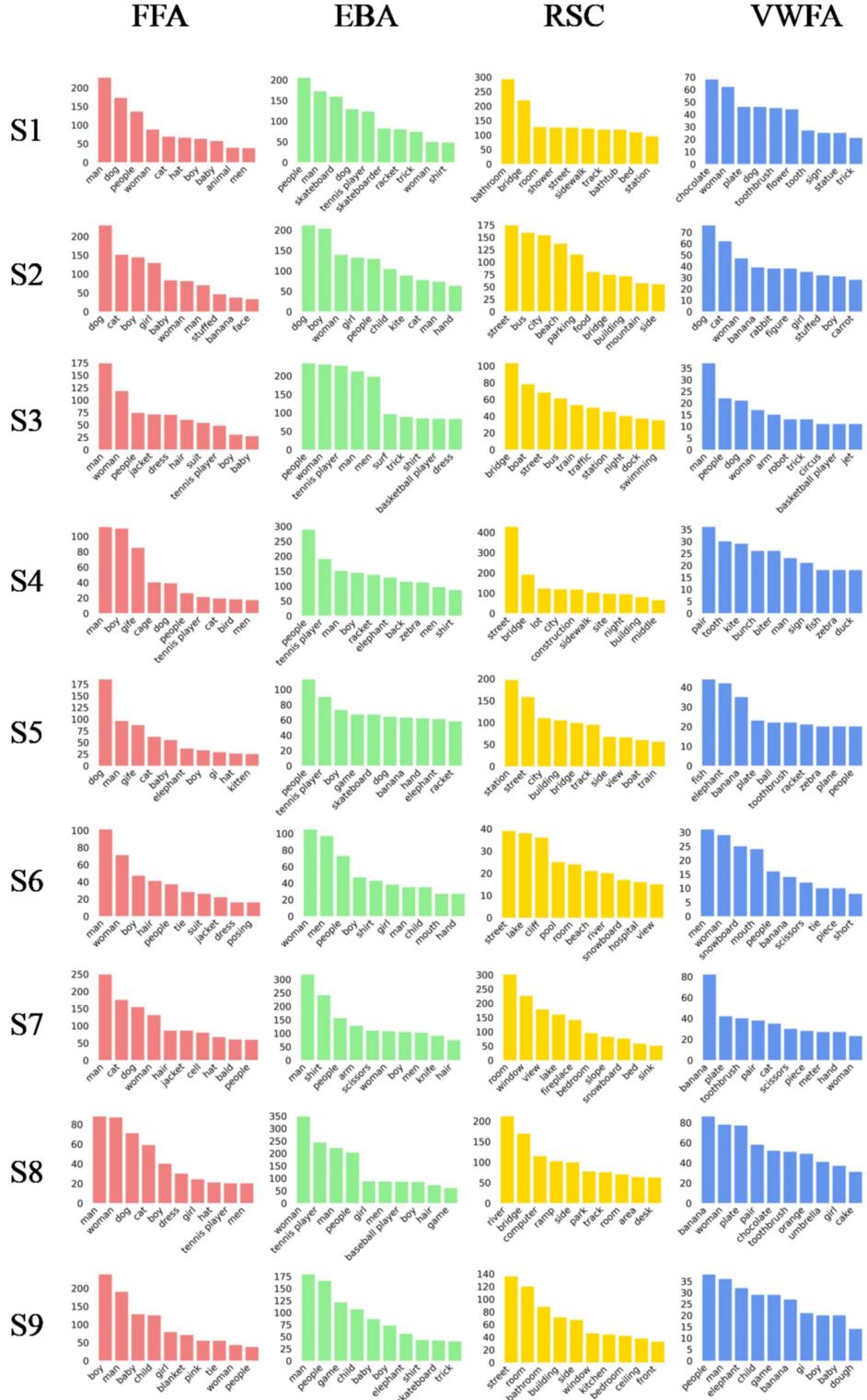


Figure 8: **The BrainLMM enables precise dissection of high-level visual categories in the selective regions.** Single semantic mapping captured via BrainLMM of all subjects in the NOD for ViT-32_{CLIP} across selective regions. The X-axis is the semantic label from the label set, and the Y-axis is the number of voxels mapped to each semantic label.

ResNet50_{clip}

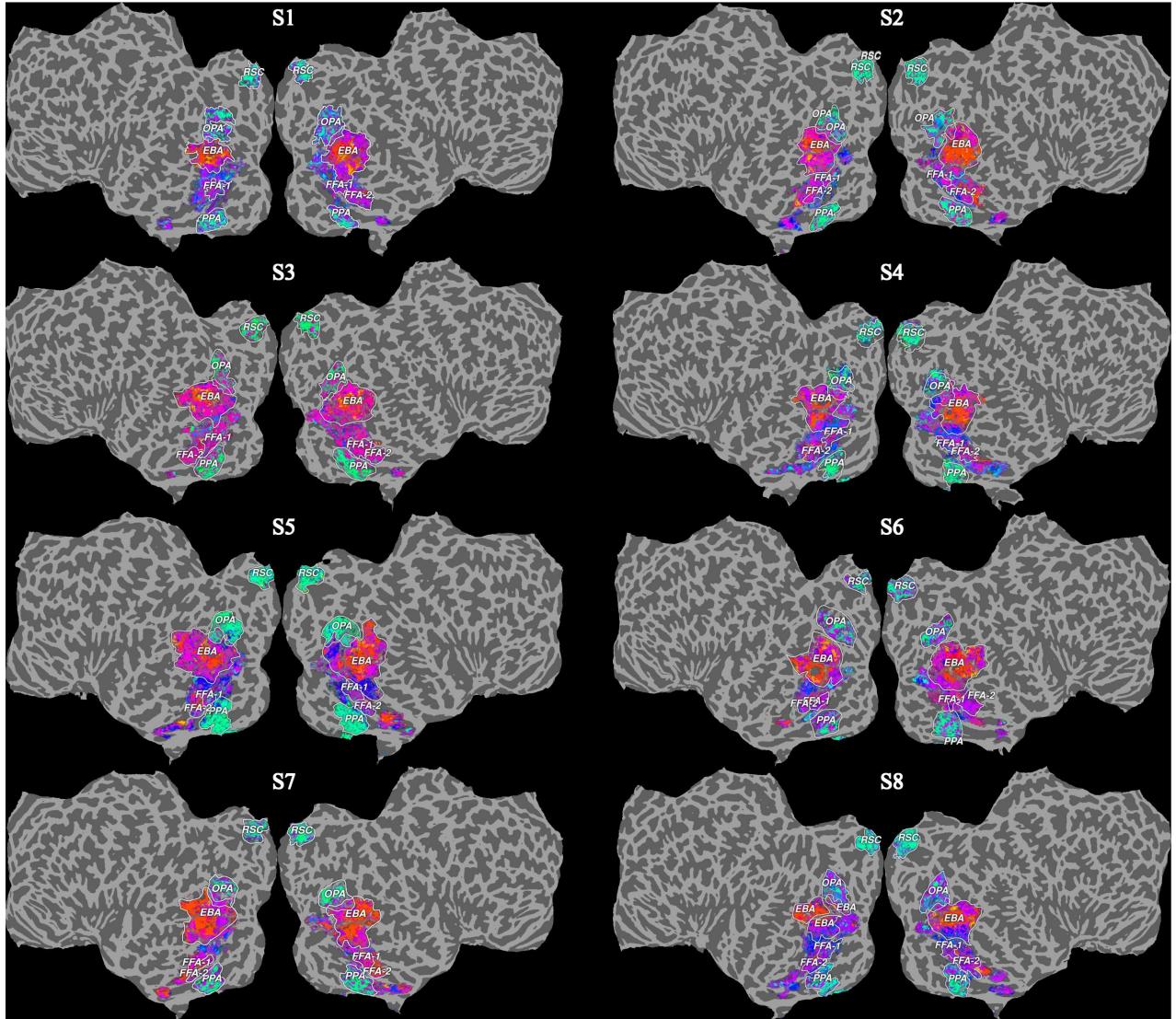


Figure 9: Voxel-level dissection results for all subjects in the NSD dataset are obtained using ResNet50_{CLIP} as the visual backbone. Image embeddings for all stimulus images are computed through the CLIP image encoder, followed by a UMAP projection into a three-dimensional space. For each voxel, the highest-scoring label from its dissection result is embedded using the CLIP text encoder, then projected into the same UMAP space and normalized to generate RGB color values. We present a flatmap of all subjects with labeled ROIs.

ViT-32_{clip}

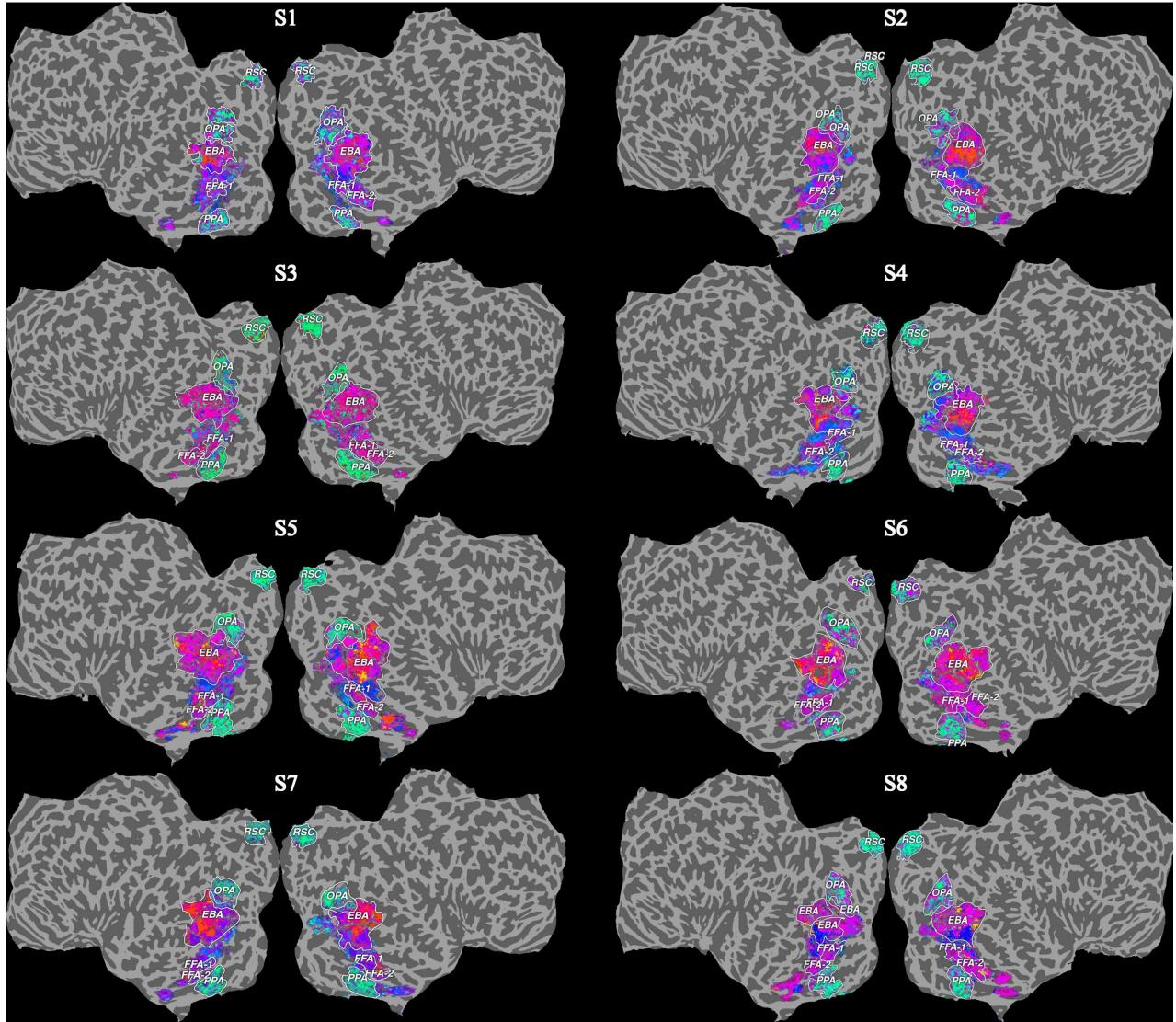


Figure 10: Voxel-level dissection results for all subjects in the NSD dataset are obtained using ViT-32_{CLIP} as the visual backbone. Image embeddings for all stimulus images are computed through the CLIP image encoder, followed by a UMAP projection into a three-dimensional space. For each voxel, the highest-scoring label from its dissection result is embedded using the CLIP text encoder, then projected into the same UMAP space and normalized to generate RGB color values. We present a flatmap of all subjects with labeled ROIs.

ResNet50_{ImageNet}

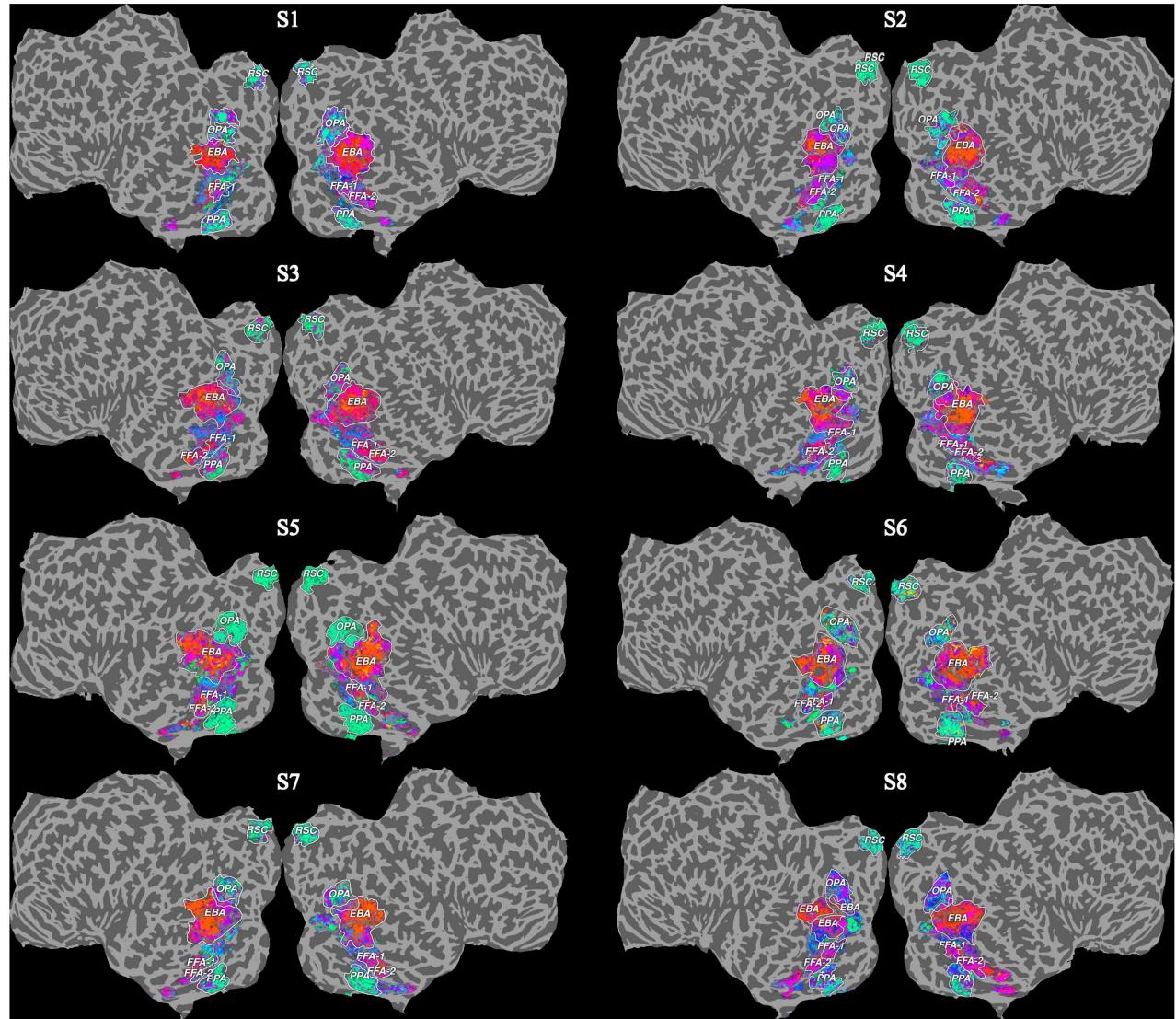


Figure 11: Voxel-level dissection results for all subjects in the NSD dataset are obtained using ResNet50_{ImageNet} as the visual backbone. Image embeddings for all stimulus images are computed through the CLIP image encoder, followed by a UMAP projection into a three-dimensional space. For each voxel, the highest-scoring label from its dissection result is embedded using the CLIP text encoder, then projected into the same UMAP space and normalized to generate RGB color values. We present a flatmap of all subjects with labeled ROIs.

AlexNet_{ImageNet}

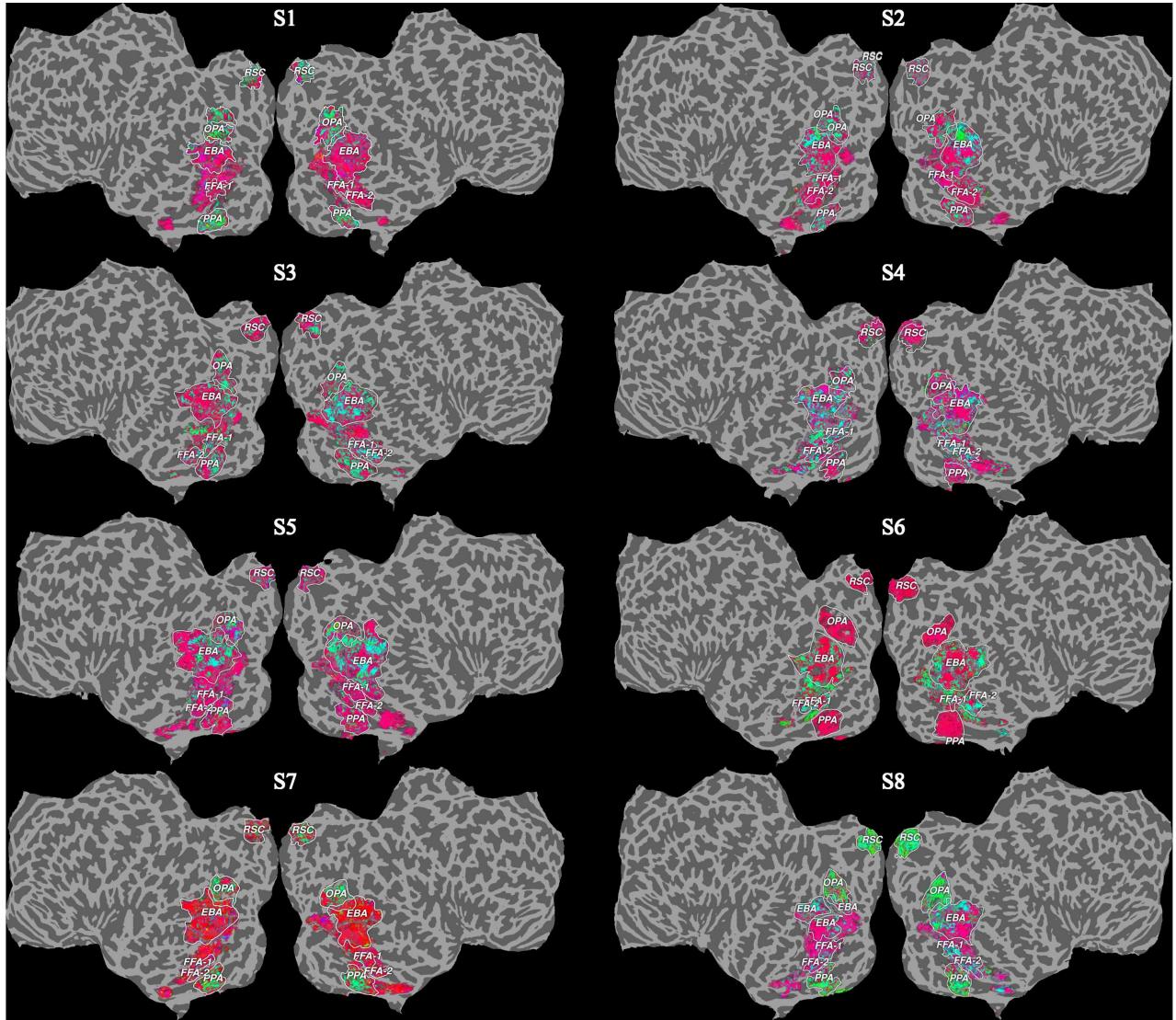


Figure 12: Voxel-level dissection results for all subjects in the NSD dataset are obtained using AlexNet_{ImageNet} as the visual backbone. Image embeddings for all stimulus images are computed through the CLIP image encoder, followed by a UMAP projection into a three-dimensional space. For each voxel, the highest-scoring label from its dissection result is embedded using the CLIP text encoder, then projected into the same UMAP space and normalized to generate RGB color values. We present a flatmap of all subjects with labeled ROIs.

NSD ResNet50_{clip}

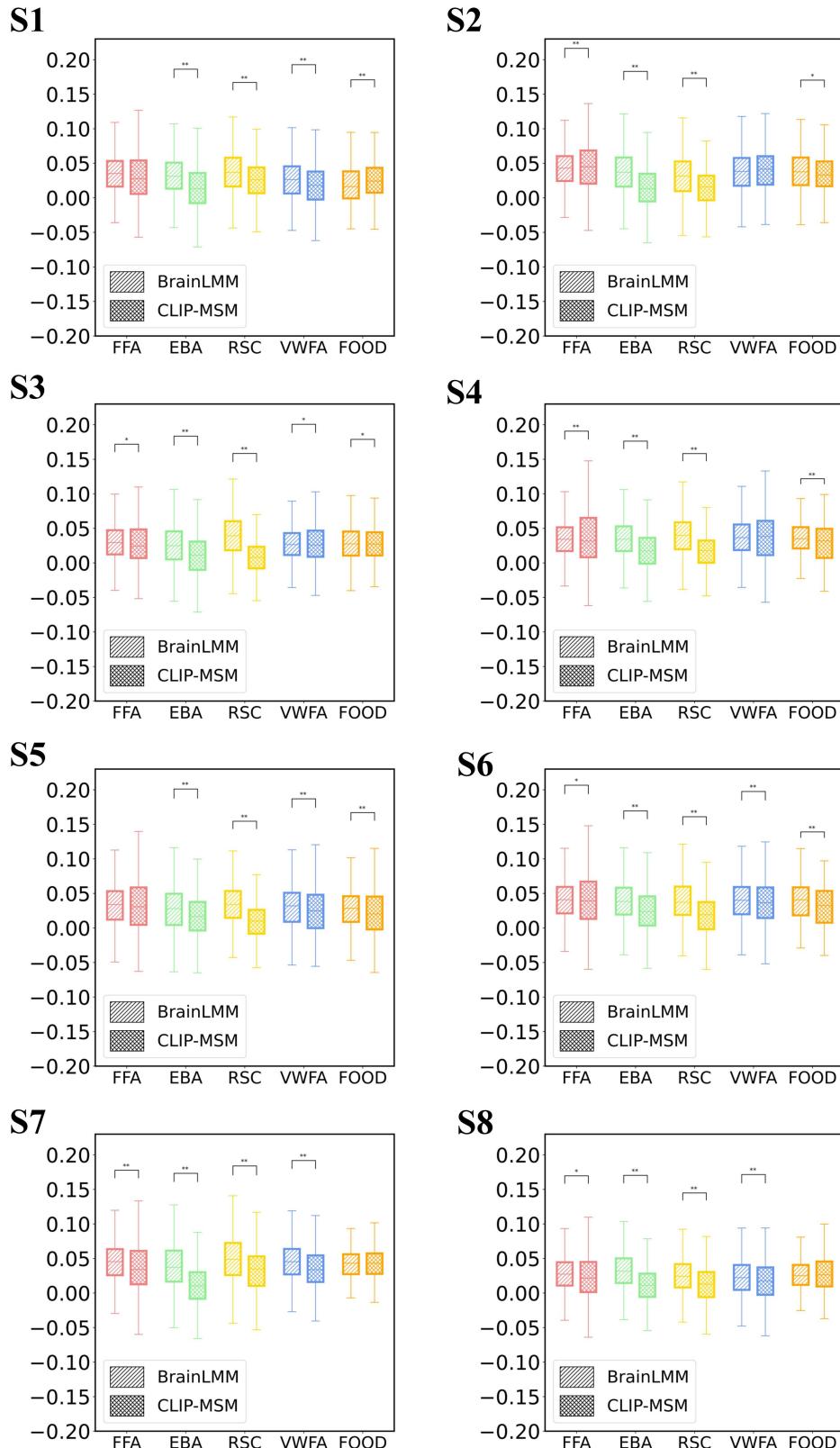
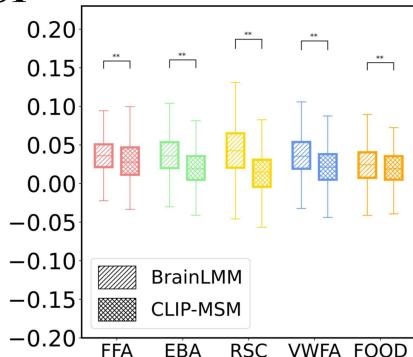


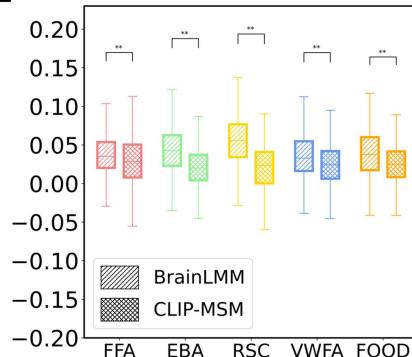
Figure 13: Performance evaluation of BrainLMM and CLIP-MSM for semantic map capture using the NSD dataset. We evaluate the performance of both methods by mapping the label of each voxel to its corresponding text embedding using the text encoder. ResNet50_{CLIP} utilizes its respective text encoders, in line with the original CLIP implementation. A higher cosine similarity between the text embeddings and voxel-wise weights signifies a stronger semantic alignment with the selective regions. As demonstrated, BrainLMM provides a more precise semantic mapping compared to CLIP-MSM (* $P < 0.05$, ** $P < 0.001$, paired t -test).

NSD ViT-32_{clip}

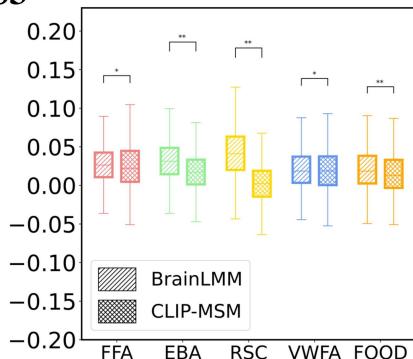
S1



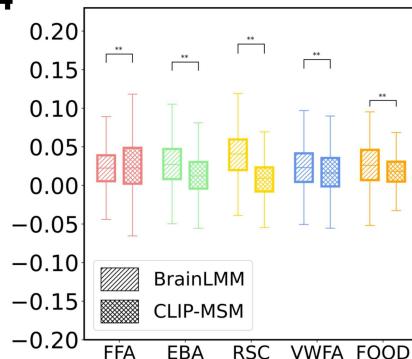
S2



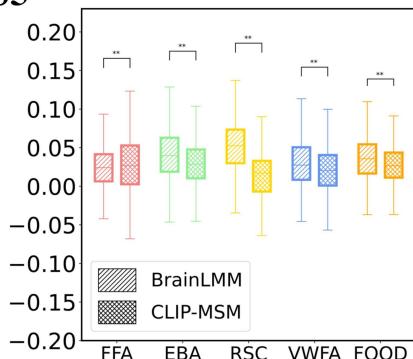
S3



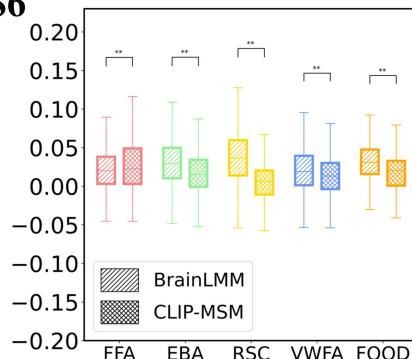
S4



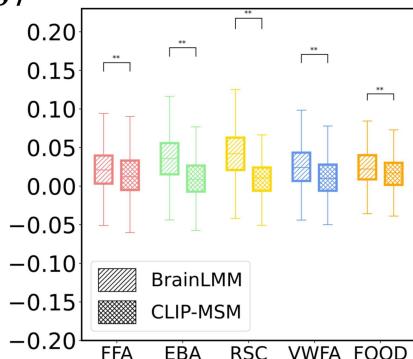
S5



S6



S7



S8

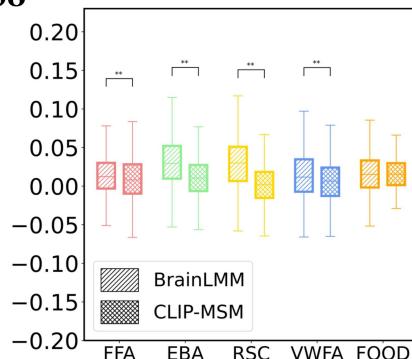


Figure 14: Performance evaluation of BrainLMM and CLIP-MSM for semantic map capture using the NSD dataset. We evaluate the performance of both methods by mapping the label of each voxel to its corresponding text embedding using the text encoder. ViT-32_{CLIP} utilizes its respective text encoders, in line with the original CLIP implementation. A higher cosine similarity between the text embeddings and voxel-wise weights signifies a stronger semantic alignment with the selective regions. As demonstrated, BrainLMM provides a more precise semantic mapping compared to CLIP-MSM (* $P < 0.05$, ** $P < 0.001$, paired t -test).

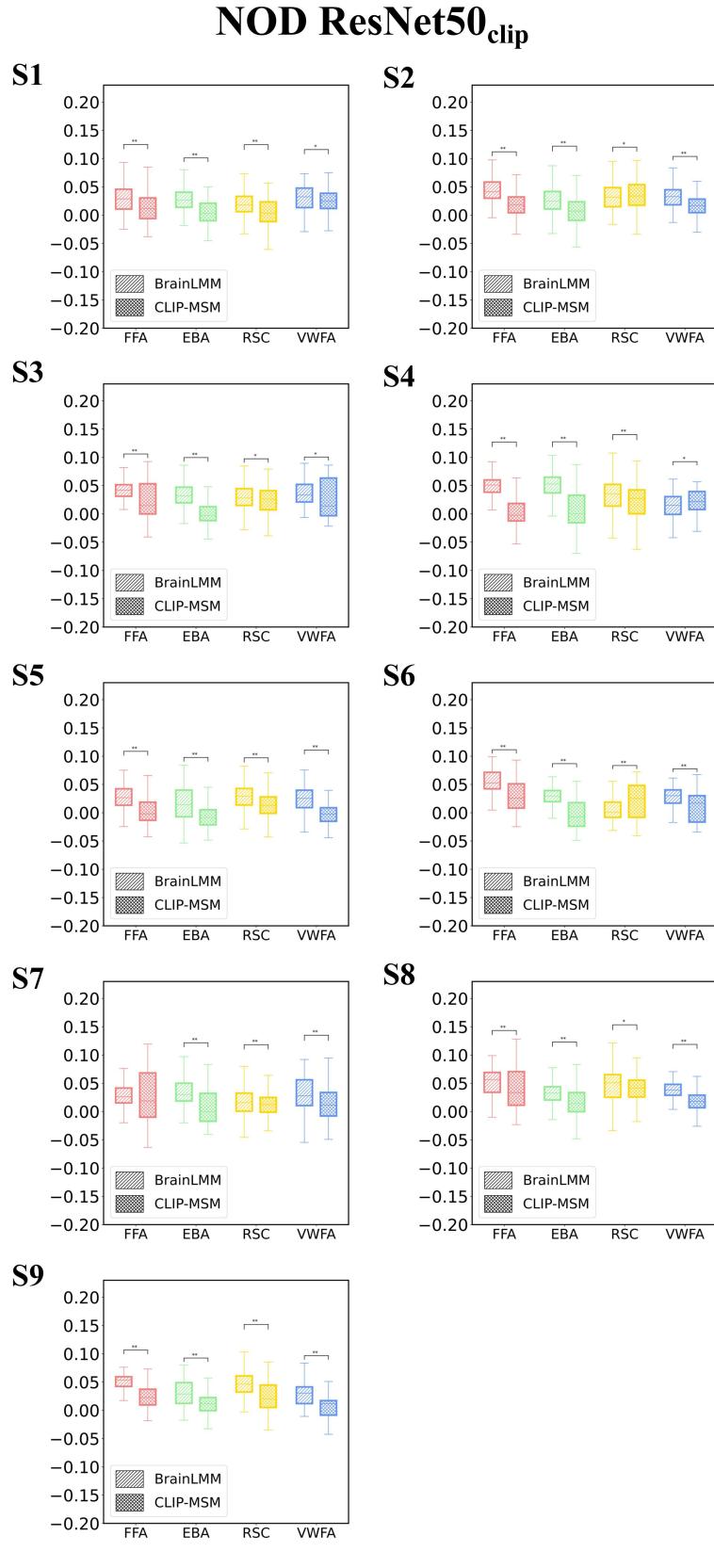


Figure 15: Performance evaluation of BrainLMM and CLIP-MSM for semantic map capture using the NOD dataset. We evaluate the performance of both methods by mapping the label of each voxel to its corresponding text embedding using the text encoder. ResNet50_{CLIP} utilizes its respective text encoders, in line with the original CLIP implementation. A higher cosine similarity between the text embeddings and voxel-wise weights signifies a stronger semantic alignment with the selective regions. As demonstrated, BrainLMM provides a more precise semantic mapping compared to CLIP-MSM (* $P < 0.05$, ** $P < 0.001$, paired t -test).

NOD ViT-32_{clip}

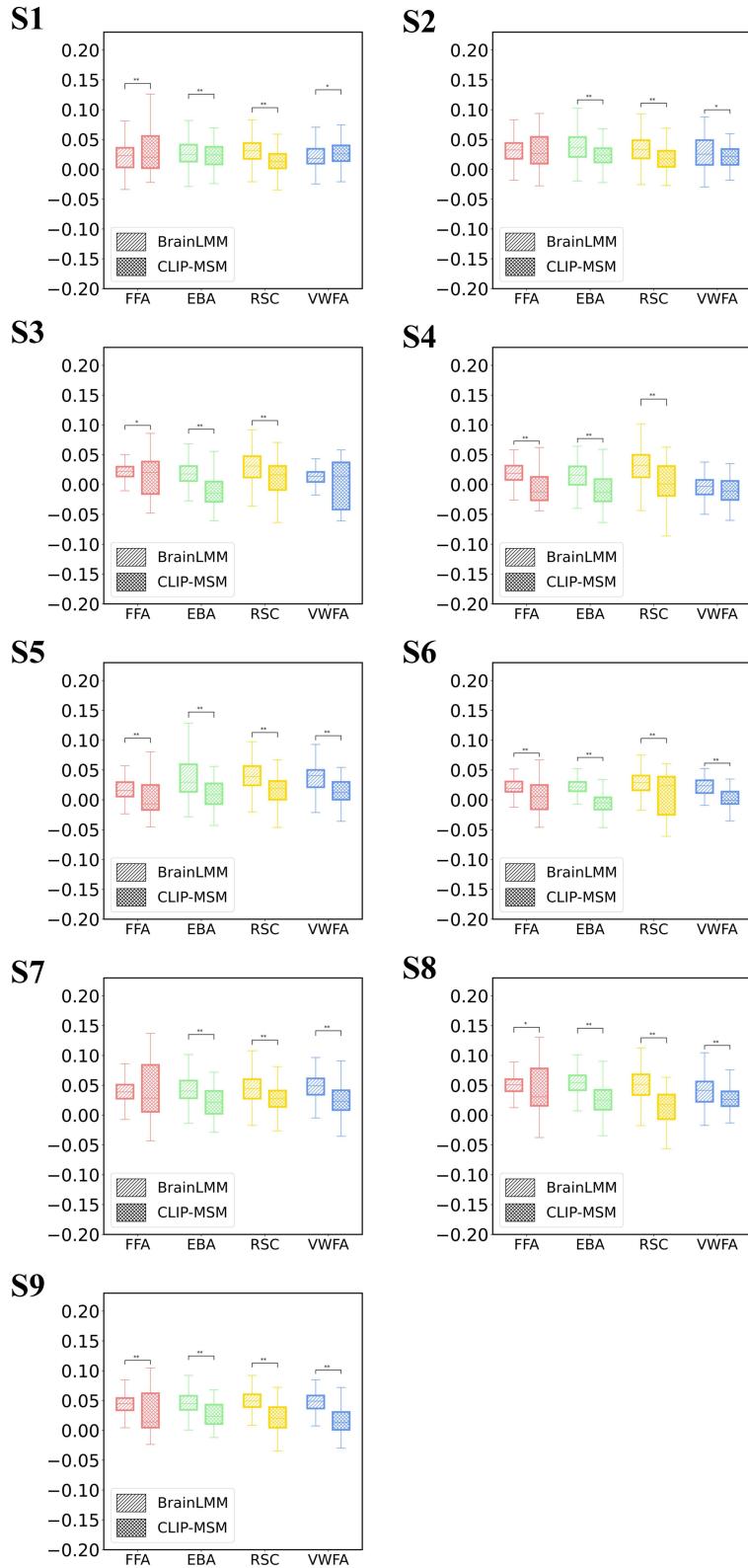


Figure 16: Performance evaluation of BrainLMM and CLIP-MSM for semantic map capture using the NOD dataset. We evaluate the performance of both methods by mapping the label of each voxel to its corresponding text embedding using the text encoder. ViT-32_{CLIP} utilizes its respective text encoders, in line with the original CLIP implementation. A higher cosine similarity between the text embeddings and voxel-wise weights signifies a stronger semantic alignment with the selective regions. As demonstrated, BrainLMM provides a more precise semantic mapping compared to CLIP-MSM (* $P < 0.05$, ** $P < 0.001$, paired t -test).

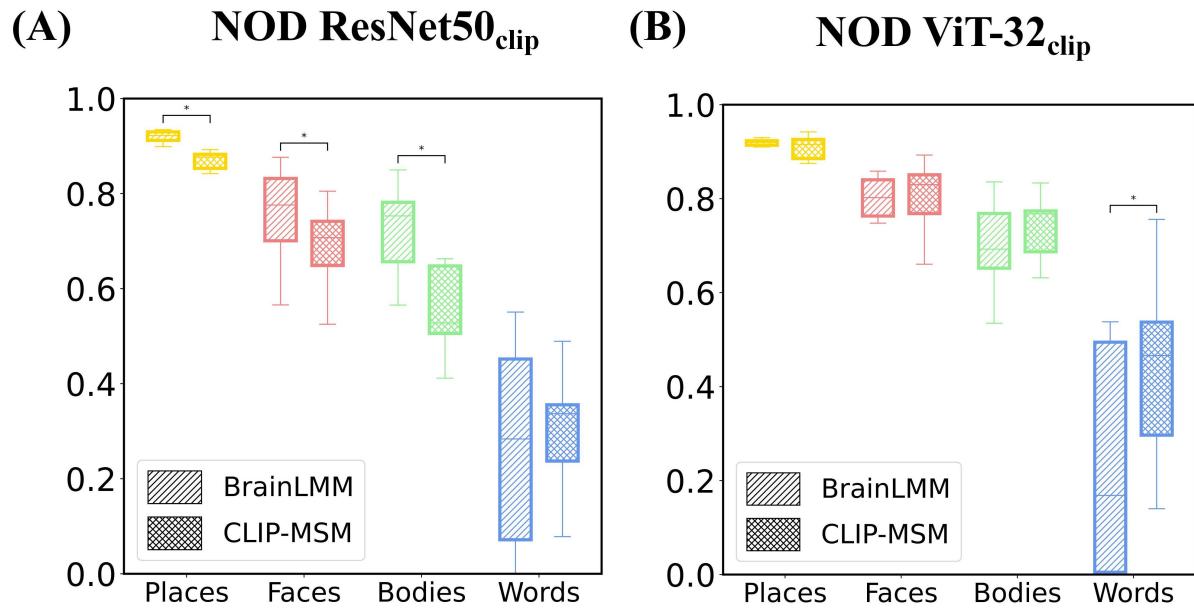


Figure 17: **Comparison of multi-semantic mapping between BrainLMM and CLIP-MSM.** Pearson correlation coefficients quantifying the alignment between multi-semantic mapping derived from both BrainLMM and CLIP-MSM and the ground-truth voxel responses across selective cortical regions. **(A)** Using ResNet50_{CLIP} as the visual encoding backbone, BrainLMM achieves significantly higher prediction accuracy than CLIP-MSM for visual responses associated with places, food, faces, bodies, and words. **(B)** Using ViT-32_{CLIP} as the visual encoding backbone, BrainLMM also outperforms CLIP-MSM in predicting responses for places, food, bodies, and words (* $P < 0.05$, ** $P < 0.001$, paired t -test).