

HIERARCHICAL REPRESENTATIONAL TRANSFORMATIONS OF WORKING MEMORY IN BRAINS AND MACHINES

Anonymous authors

Paper under double-blind review

ABSTRACT

Working memory (WM) maintains past inputs while processing new ones, yet how representations transform between encoding and retrieval remains unclear. Clarifying whether these representations are sustained through stable coding formats, dynamically updated subspaces, or their interplay is key to uncovering the mechanisms of WM. To address this, we combined high-resolution 7T fMRI from the Natural Scenes Dataset with recurrent neural networks (RNNs) trained on a naturalistic 1-back task. Using representational similarity, cross-decoding, and subspace geometry analyses, we directly compared rotational and non-rotational transformations between WM encoding and retrieval phases in brain regions and model layers. Our analyses revealed convergent evidence for a mixture mechanism of WM coding for encoding and retrieval information: early visual regions (V1–hV4) underwent large representational changes across encoding to retrieval phases, including both rotational and non-rotational transformations. Whereas higher-order regions in the prefrontal cortex (FEF, dlPFC) were more stable. Applying the same analyses to models showed a similar mechanism across layers, but critically depended on the learning objective and the recurrent architecture. We examined two different encoder architectures, ResNet and Vision Transformer (ViT), each trained with supervised and self-supervised learning objectives. Models with supervised encoders preserved a hierarchical layer dissociation paralleling the cortical gradient in both rotational and non-rotational transformations, while models with self-supervised encoders diverged in the rotational transformation. Among recurrent architectures, gated architectures (GRU, LSTM) better reproduced the brain-like mixture of subspace rotational transformation. Taken together, these results established hierarchical shifts between flexibility and stability in WM representational transformation in both humans and machines, with supervised learning objectives combined with gated recurrent dynamics most closely resembling human WM mechanisms.

1 INTRODUCTION

Working memory (WM) is the neural and cognitive process that temporarily stores and manipulates sensory information (Baddeley, 1992; 2003). It supports a wide range of higher-order cognitive functions such as learning, reasoning and decision-making (Collins & Frank, 2012; Daneman & Carpenter, 1980; Süß et al., 2002; Wagner, 1999; Engle, 2010; Cools & D’Esposito, 2011). Classic work has highlighted the limitations of WM in capacity and precision: its precision decreases with the number of items stored (Luck & Vogel, 1997; Ma et al., 2014) and with longer delays (Pertzov et al., 2017; Magnussen et al., 1998; Shin et al., 2017). Despite its limited capacity, WM exhibits a high degree of flexibility. People can maintain WM content even in the presence of new incoming stimuli. Such flexibility indicates that WM representations are not just mere traces of the original inputs, but undergo dynamic transformation to serve diverse cognitive demands.

The representational format of WM content remains an ongoing debate in neuroscience and cognitive science. WM-related signals observed in early visual cortex have been interpreted as evidence that WM and perception share a similar representational format (Serences, 2016; Christophel et al., 2017; Harrison & Tong, 2009). However, a shared format across different stages of information pro-

054 cessing could lead to interference between types of information that should be kept separate. Consistent with this notion, studies have shown that WM representations can differ from their perceptual counterparts. Such differences may naturally emerge during memory delays (Kwak & Curtis, 2022; Li & Curtis, 2023; Spaak et al., 2017; Murray et al., 2017), in the presence of sensory distractions (Xu, 2024; 2025; Libby & Buschman, 2021; Degutis et al., 2025), or when the WM content is assigned with different levels of priority (Wan et al., 2020; 2022). Reconciling these observations, some electrophysiological evidence suggests that a dynamic neural code and a stable neural code may coexist when maintaining WM information (Stokes et al., 2013; Murray et al., 2017).

062 In this study, we investigate how WM representations are transformed from encoding to retrieval, the stage where stored information is accessed to guide behaviors, as in a 1-back task requiring comparison between past items held in memory and the current input. Retrieval is particularly important because it is the stage when stored information is accessed and used for comparison, such as in an n-back task where past items must be retrieved to evaluate against the current input. We aim to characterize these transformations across brain regions along the cortical hierarchy. Specifically, we ask whether different stages of the cortical hierarchy implement distinct transformations: some regions may represent retrieved information in a format distinct from currently encoded signals, thereby reducing interference, whereas other regions may preserve a more coherent code between encoding and retrieval. By analyzing representational similarity and geometry, we aim to evaluate both rotational and non-rotational representational transformation. In addition, we compare the transformations from WM encoding to retrieval between human neural recordings and recurrent neural networks (RNNs) trained to perform continuous WM tasks. By testing neural networks with varying architectures and learning objectives, we aim to further identify the factors in models that support human-like representational transformation in neural networks.

076 To address this, we combine high-resolution 7T fMRI from the Natural Scenes Dataset (NSD) (Allen et al., 2022; Gifford et al., 2025) with RNNs trained on a 1-back WM task. We analyze representational similarity, cross-phase decoding, and geometric subspace rotation angles to test whether memory representations are best explained by stable or dynamically updated subspace mechanisms across the cortical hierarchical biological and artificial systems. Our contributions are as follows:

- 082 • **Mixture of dynamic and stable WM subspaces across the hierarchy:** Encoding–retrieval representations are partially overlapping but systematically transformed along the hierarchy, with greater transformation in early visual regions (e.g. V1-hV4), and more stable subspaces in higher-order regions (FEF, dlPFC) in both rotational and non-rotational transformation.
- 087 • **Model–brain alignment in WM encoding–retrieval transformation depends on both the encoder learning objectives and the recurrence modules:** Supervised models and gated architectures (GRU, LSTM) better captured the brain-like mixture of WM subspace transformation along the layers.

091 Together, these findings provide new insights into how the brain and machines using dynamic transformations to minimize interference in early regions while maintaining stable subspaces in higher-order regions to ensure reliability in WM. Furthermore, this efficient mixture of coding strategies is supported by the supervised learning objective and gated recurrent architectures in machines.

096 2 RELATED WORKS

098 RNNs are widely used in neuroscience to model WM due to their ability to maintain information over time and generate dynamics resembling cortical activity observed in humans and non-human primates (Wang, 1999; Wimmer et al., 2014; Compte et al., 2000; Bouchacourt & Buschman, 2019; Wang, 2021; Yang & Wang, 2020; Esnaola-Acebes et al., 2022). By training with WM delayed-response tasks, RNNs can reproduce hallmarks of WM phenomena such as persistent activity, attractor dynamics, and flexible subspace organization. These capabilities make RNNs a powerful computational framework for probing how WM representations are formed, maintained, and adapted under varying task demands.

106 Despite extensive research on neural representations of WM, studies linking representational transformations in neural networks and brains remain rare. Prior work has compared RNNs to post-cue WM tasks in non-human primates (Piwek et al., 2023) and human whole-brain EEG in n-back tasks

(Wan et al., 2022), but these studies relied on simple artificial stimuli (e.g. color and orientation) in predefined subspaces. As such, they could not capture representational dynamics for naturalistic stimuli, or dissociate coding strategies across brain regions given the low spatial resolution of EEG signals. Some studies analyzed representational transformations in human fMRI with simple artificial stimuli (Kwak & Curtis, 2022; Li & Curtis, 2023; Degutis et al., 2025) or high-dimensional naturalistic stimuli (Xu, 2024; 2025; Nakamura et al., 2025), but without comparing the neural representations of the brain to those of neural networks. More recently, transformation dynamics have been quantified in RNNs trained on n-back tasks with naturalistic stimuli (Lei et al., 2024), but without comparing RNNs with human brain data. These efforts face common limitations: reliance on simplified stimuli, scalar subspace metrics ill-suited to high-dimensional geometry, and isolated analyses of either neural networks or brain data. We address these gaps by directly comparing WM encoding-retrieval transformation matrices between artificial and biological systems, by using representational similarity, cross-decoding, and geometric alignment analyses.

3 METHODS

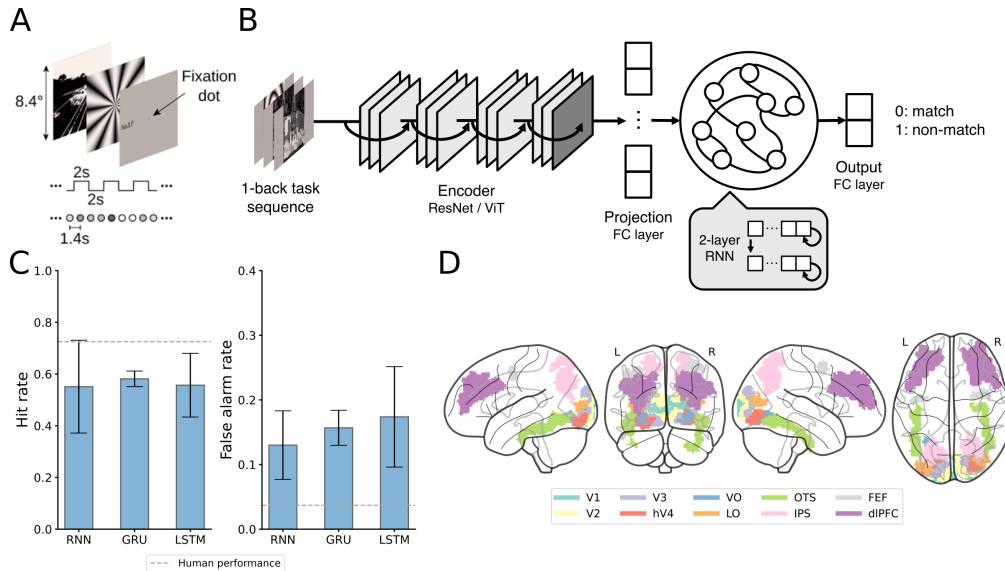


Figure 1: Tasks, model architecture, performance, and brain ROIs. (A) Trial design of the task. For 1-back task, subjects were asked to judge whether the currently presented image was identical to the previously one. Image from Gifford et al. (2025). (B) Model architecture. The model largely consists of these parts: encoder (ResNet, ViT) with different training objective (supervised, self-supervised learning), and RNN (vanilla RNN, LSTM, GRU). The input was the 1-back task stimuli sequences, and the encoder part extracted the representation, which was passed to a project layer in order to match the RNN’s input dimension. Then RNN module was linked to a fully connected layer to produce the decision. (C) Model performances on the one-back task. Error bars show the mean ± 1 s.e.m across averaged testing performance over all the testing sets of models across 4 encoder types. (D) Regions of interest. Subject 1 from NSD was used as the sample subject for visualization.

3.1 NEURAL RECORDINGS AND TASKS

The NSD-synthetic dataset extends the Natural Scenes Dataset (NSD-core; Allen et al., 2022) with 7T fMRI recordings of fixation and 1-back tasks using new synthetic stimuli (Gifford et al., 2025). Data were collected from 8 participants viewing 284 synthetic images spanning diverse formats and semantic information, including 8 subclasses such as natural and manipulated scenes, noise, words, gratings, and chromatic variations.

Participants alternated between two tasks: a fixation task measuring perceptual representation and a 1-back task probing WM encoding and retrieval. Each task was performed in 4 runs. Each run

162 contained 93 trials, with each trial consisting of a 2-s image presentation followed by a 2-s inter-trial
 163 interval. In the fixation task, participants responded to color changes in the central dot, while in the
 164 1-back task they indicated whether the current image matched with the stimulus before it (Fig. 1A).

165 The analyses were focused on 10 regions of interest (ROI, Fig. 1D), ranging from early visual areas
 166 (V1–hV4) to mid-level regions (LO, VO, IPS, OTS) and higher-level regions such as FEF and dIPFC,
 167 defined by the mask provided by the NSD (Allen et al., 2022), a probabilistic map (Wang et al.,
 168 2015), and an established parcellation atlas (Glasser et al., 2016). The more detailed description of
 169 the fMRI dataset and the tasks is in Appendix A.2.

171 3.2 MODEL

172 We modeled the 1-back task with a two-stage architecture (Fig. 1B). The frozen visual encoder en-
 173 coded each image to a feature vector, which was passed through a fully connected project layer
 174 before being processed by a recurrent module. Note that the project layer and recurrent modules
 175 were trained in this work. The recurrent module consumed sequences of projected feature vectors,
 176 and emitted binary decisions, indicating whether the current stimuli matched the last one. For repre-
 177 sentational analyses, we extracted activations from the the project layer and recurrent hidden layers
 178 (layers 1 and 2). Implementation details are provided in Appendix A.3.

181 **The frozen encoder and projection module.** We used two different frozen encoder architectures:
 182 a convolutional neural network (ResNet-50; He et al., 2016) and a Vision Transformer (ViT-B/16;
 183 Dosovitskiy et al., 2020). We tested four frozen image encoder architecture–learning objective set-
 184 ups: ResNet-50 and ViT-B/16, each pretrained with ImageNet-1K supervised learning (SL) or self-
 185 supervised learning (SSL). We used contrastive learning objectives (He et al., 2020; Chen* et al.,
 186 2021) for self-supervised learning (see Appendix A.3 for details). These encoders remained frozen
 187 during training, and their final-layer feature outputs were processed by the project layer to reduce
 188 the dimension before the recurrent module. We extracted 4 encoder layers of each setup: one per
 189 block from ResNet-50’s 4 blocks, and layers 3, 6, 9, 12 from ViT-B/16.

190 **Recurrent modules.** We evaluated 3 recurrent architectures—vanilla RNN (Elman, 1990),
 191 LSTM (Hochreiter & Schmidhuber, 1997), and GRU (Chung et al., 2014)—each with two lay-
 192 ers. Our design crossed encoder types (ResNet vs. ViT), learning objectives (supervised vs. self-
 193 supervised), and recurrent modules (RNN, LSTM, GRU), enabling systematic comparison of archi-
 194 tectures and learning objectives on human-like sequence processing. For brevity, we denote the
 195 4 encoder–objective setups as ResNet-SL (ResNet-50, supervised), ResNet-SSL (ResNet-50, self-
 196 supervised), ViT-SL (ViT-B/16, supervised), and ViT-SSL (ViT-B/16, self-supervised).

197 4 RESULTS

198 4.1 MODEL PERFORMANCE

199 We trained 12 model configurations with 2 different learning objects \times 2 encoder architectures
 200 \times 3 recurrent architectures (details in A.3.5) with NSD-core images (Allen et al., 2022). Across
 201 all configurations, training accuracy exceeded 90%, and validation accuracy exceeded 85%. Next,
 202 these trained models were tested on the NSD-synthetic dataset (Gifford et al., 2025), which consists
 203 of 284 out-of-distribution (OOD) images. We trained the models with varying levels of “match”
 204 events, while the test set NSD-synthetic contained a highly imbalanced distribution of match events
 205 (11% targets). In the imbalanced environment, quantifying performance with hit rate and false alarm
 206 rate is more reliable than raw accuracy.

207 Here we reported the test performance with NSD-synthetic (Fig. 1C), where each of 12 models
 208 was trained with 5 random seeds (details in A.3.6). All recurrent architectures (vanilla RNN, GRU,
 209 LSTM) achieved moderate hit rates in the 1-back task, while lower than the human observers. Addi-
 210 tionally, models exhibited higher false alarm rates than humans, suggesting a higher bias. Together,
 211 these results indicate that while models could perform the 1-back task, they do not fully reproduce
 212 the human decision strategy due to a mismatch between the match frequency of the training and
 213 testing environment.

4.2 NEURAL MECHANISM OF CONCURRENT WM ENCODING AND RETRIEVAL

After confirming that our models successfully learned the 1-back task, we next turned to our central question: How are representations transformed between encoding and retrieval? Specifically, we tested three competing assumptions about the mechanisms of processing concurrent WM encoding and retrieval information: (I) Stimuli held in WM maintain a stable representational format preserved across encoding and retrieval. (II) Representations are dynamically updated as the trials unfold, with systematic transformations between the encoding and retrieval phases. (III) A mixture of stable and transformed memory subspaces, and different cortical areas may implement different strategies.

To adjudicate between these accounts, we analyzed human fMRI and model activations using representational similarity, cross-decoding, and subspace geometry. These analyses revealed consistent evidence for the mixture mechanism, with early visual regions (V1–hV4) exhibiting stronger dynamic transformations and higher-order regions (dlPFC) maintaining more stable subspaces.

4.2.1 REPRESENTATIONAL SIMILARITY EVIDENCE

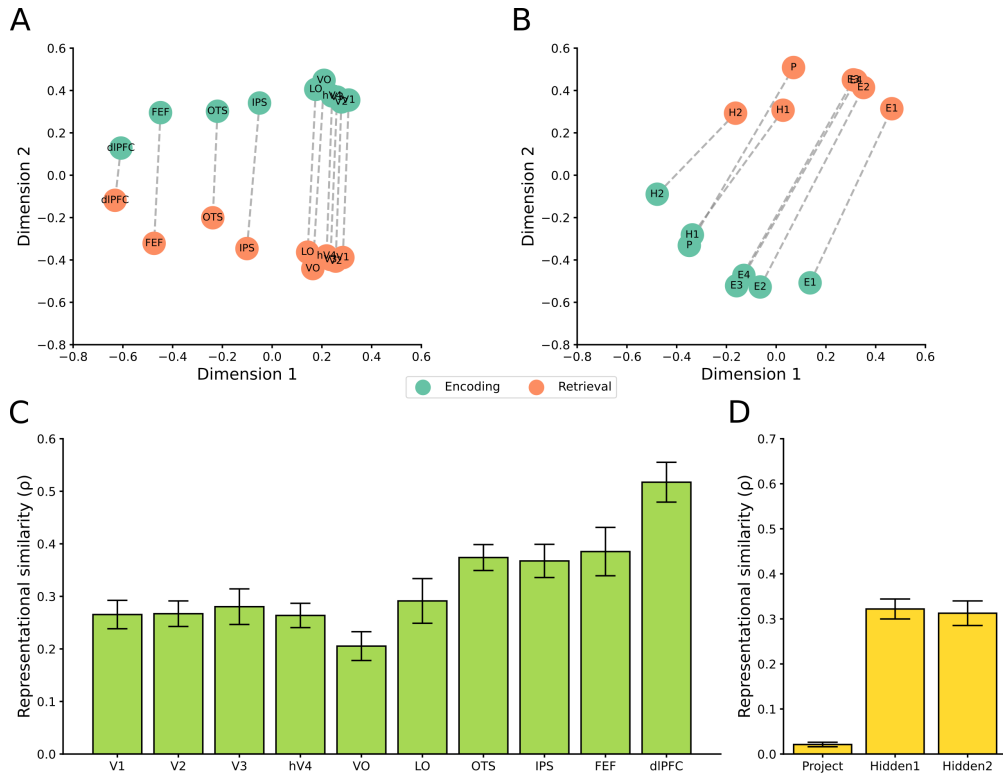


Figure 2: Representational similarity results. (A) MDS visual-path of the dissimilarity matrix based on RDMs of encoding and retrieval representations for each human ROI. (B) MDS visual-path of the averaged dissimilarity matrix across 12 models based on RDMs of encoding and retrieval representations for each model layer. E denotes frozen weight encoder layers, P denotes the project layer, and H denotes hidden layers. The number following the letter denotes the order of that layer. (C) Similarity scores (Spearman’s ρ) of encoding and retrieval RDMs for each human ROI. Error bars represent ± 1 s.e.m across 8 subjects. (D) Similarity scores (Spearman’s ρ) of encoding and retrieval RDMs for each model layer. Error bars represent ± 1 s.e.m across 12 models

With representational similarity analysis on brains and models, we consistently observed an increasing similarity scores between encoding and retrieval representations from low-level to high-level cortical regions or model layers. For the brain data, we constructed the identity-based representational structures (RDMs) for each ROI in both encoding and retrieval phases. We then calculated the dissimilarity ($1 - \text{Spearman’s } \rho$) of the RDMs from both phases and all ROIs. Finally, we ap-

plied multidimensional scaling (MDS) to project these relationships into a two-dimensional space (Fig. 2A, with methodology details in Appendix A.7.1).

Within each phase, the representations of low-level (V1-hV4) and mid-level areas (LO and VO) clustered close to each other, while the representations shifted a lot from low-level to the higher-level cortical areas (dlPFC). Interestingly, the WM encoding representations in the lower areas (V1-hV4) were very distinct from their representations in the retrieval phase, while the representations in the higher-level areas (dlPFC) exhibited larger similarity across the phases. This similarity between the encoding to retrieval phase increased gradually along the hierarchy. RNNs exhibited the same hierarchical pattern, as the representations in two phases were closer in the higher-level hidden layers compared to the lower untrained encoder layers and trained project layer (Fig. 2B; Fig. A6 for model-wise visualizations). In the following analyses, we mainly focused on the project and 2 hidden layers trained on the 1-back WM task.

We further confirmed our observation in the MDS of RDM-dissimilarity by directly computing the similarity score Spearman’s ρ for each brain ROI or model layer between RDM of encoding v.s. retrieval phase (Fig. 2C). A permutation based one-way ANOVA test showed a significant main effect of ROI on the similarity score ($p < .001$), where the score was highest in dlPFC compared to all the other regions (permutation post-hoc test, $p < .001$, FDR corrected). For the models, a one-way ANOVA revealed a significant main effect of layer on encoding–retrieval similarity ($p < .001$). FDR post-hoc paired t-tests showed higher similarity for both hidden layers than for the project layer ($p < .001$, Fig. 2D). To assess the reliability of the RSA, we shuffled RDM cells and ran the same set of analyses (Fig. A7). A permutation based repeated-measures ANOVA revealed significant main effects of observed or shuffled condition and ROI ($p < .001$), as well as a significant interaction between condition and ROI ($p < .001$). FDR post-hoc paired t-test between the observed pattern and shuffled conditions (random baseline) confirmed the robust effects across all ROIs ($p < .001$).

Overall, the evidence of relative representational distance along stimuli suggested hierarchical changes from dynamically updated memory subspaces in lower regions/layers to more stable representational formats in high-level regions/layers, in both biological and artificial organizations.

4.2.2 DECODING EVIDENCE

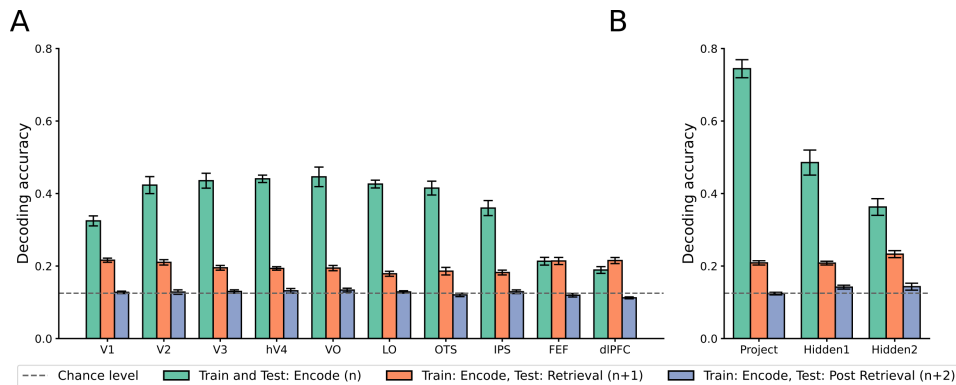


Figure 3: Decoding accuracies of the decoder trained in the encoding (n) phase, tested in the encoding (n), retrieval ($n+1$), and post-retrieval ($n+2$) phases. (A) Results of each human ROI. Error bars represent ± 1 s.e.m. across subjects. (B) Results of each model’s layer. Error bars represent ± 1 s.e.m. across models.

To complement RSA, we next used cross-decoding to investigate representational changes between the encoding and retrieval phases (Harrison & Tong, 2009; Xu, 2025). For each ROI or model layer, we trained linear decoders to classify the stimulus identity across 8 subclasses (Gifford et al., 2025). Decoders were trained during the encoding phase with the current stimuli (step n) and evaluated either within the same phase (n), during the retrieval trial ($n+1$) when a match/no-match action is required, or during the post-retrieval trial ($n+2$) when the stimuli information is irrelevant (details in Appendix A.7.2).

As shown in Fig. 3A, decoding accuracies were highest when trained and tested within encoding, of which all ROIs showed above chance accuracy (permutation test, $p < .01$). The decoder trained in the WM encoding phase can reliably generalize to the retrieval ($n+1$) phase in all ROIs, achieving above chance level decoding accuracy (permutation test, $p < .01$). However, the decoder trained in WM encoding phase tested in post-retrieval ($n+2$) phase was not above chance level in all ROIs (permutation test, $p > .05$). Therefore, the decoders in frontal areas are reliably generalized from encoding to retrieval phase, but not the post-retrieval phase when the encoded information is irrelevant. More importantly, there is a significant 2-way interaction of ROI and test condition ($p < .001$): the decoders of most the ROIs showed higher decoding accuracies when tested in encoding phase compared to the retrieval phase (FDR corrected, permutation based post-hoc paired t-test, $p < .001$), except for frontal areas (FEF and dlPFC $p > .05$). This significant cross-decoding performance drop could signal a large representational transformation (Xu, 2025) between encoding and retrieval phase in the low- to mid-level visual regions (V1-IPS). Additionally, all ROIs showed above chance accuracy when the decoders were trained and tested in the retrieval ($n+1$) phase (Fig. A1, permutation test, $p < .01$), indicating an active engagement of the all ROIs, including FEF and dlPFC, during the memory retrieval.

Decoding accuracy in the fixation task for decoders trained and tested in the encoding phase didn't achieve above chance accuracy across all ROIs, and lower than that in the 1-back task (permutation test, $p < .01$), indicating that 1-back decoding reflected WM processes beyond sensory responses.

For RNNs, decoding showed a pattern similar to humans, as shown in Fig. 3B. The decoder trained in encoding could reliably generalize to the retrieval but not post-retrieval phase. A 2-way ANOVA showed an interaction between layer and test condition ($p < .001$) and a main effect of layer and test condition ($p < .001$). For decoder train and test in encoding phase, the decoding accuracy was highest in the project layer, and then the 1st hidden layer, with lowest accuracy in the 2nd hidden layer (FDR post-hoc paired t-test $p < .001$). When tested in the retrieval phase, the accuracy was higher in the 2nd hidden layer than the project layer (FDR post-hoc paired t-test $p < .05$).

Hence, consistent with RSA results, decoding results also suggest that stimulus information is well-represented at encoding but undergoes substantial transformation from encoding to retrieval. Importantly, RSA is insensitive to subspace rotations since it is built from pairwise representational distances (Fig. 4D). This motivates our geometric rotation analyses, where we directly quantify the rotation angles between encoding and retrieval subspaces to test for stable versus dynamic subspace mechanisms.

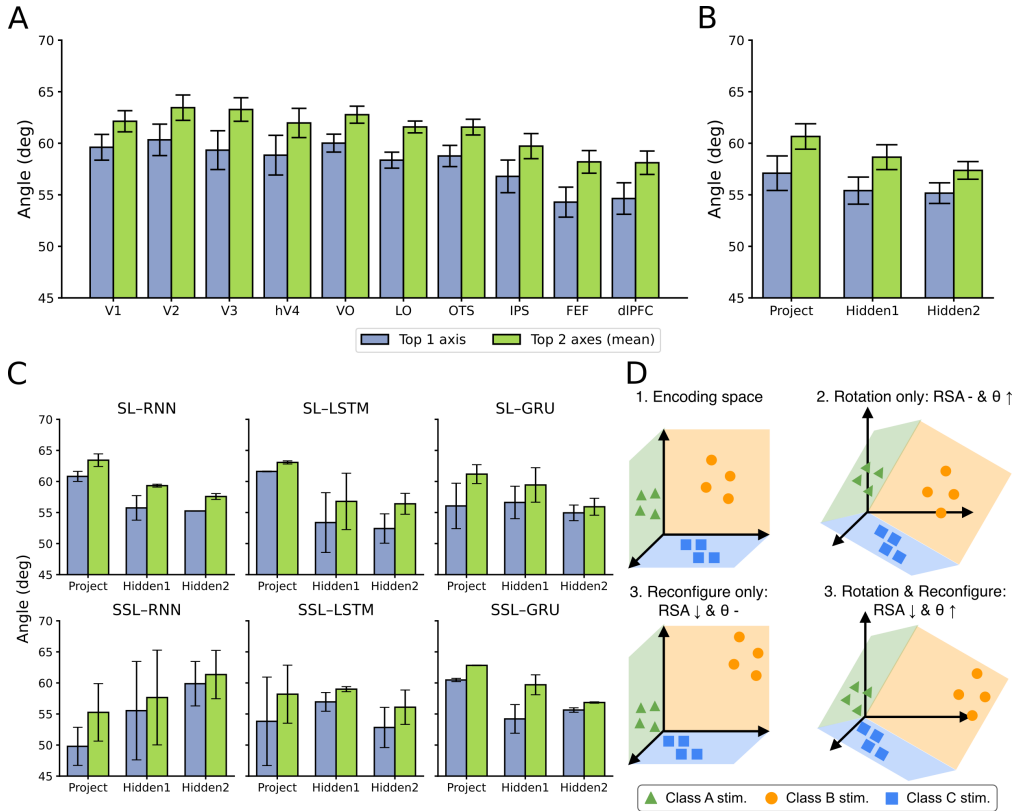
4.2.3 GEOMETRY ROTATION EVIDENCE

To directly assess the rotational geometric relationship between the encoding and retrieval representations, we computed the principal angles between the subspaces spanned by class-selective decoder weight vectors in each ROI (Fig. 4A) (Gower, 1975). We first used the top-1 principal angle to quantify the degree of rotations of the neural subspaces between the encoding and retrieval phases. This top-1 angle represents the neural dimension that is most aligned across the two phases. We observed a significant main effect of ROIs in the top-1 angle (permutation test, $p < .05$), and there was a reduction in rotations when moving from low-level visual areas to higher-level brain regions (FEF, dlPFC). We conducted the same analysis using the average of the top-2 angles, which is a more robust measure of rotation, and found similar results: there was a significant main effect of ROIs (permutation test, $p < .01$), and the rotation angle between encoding and retrieval was larger in the low-level regions (V2, V3) compared to the high-level frontal areas (FEF, dlPFC; post-hoc paired t-tests, FDR corrected $p < .05$). The permutation based 2-way ANOVA (2 angles x 10 ROIs) showed no significant interactions ($p > .05$). Therefore the pattern of the rotation angles across ROIs remains across the angle dimensions.

As a validation of whether decision geometry across the WM phases can be approximated by a rigid rotation, for each ROI we trained a classifier in the encoding subspace and tested it on retrieval with vs. without applying the rotation derived from principal angles (details see Appendix A.7.2). Across ROIs, the rotation improved cross-phase decoding when we kept the most aligned Top 1–2 axes (one-sample t-test $p < .05$). In other words, only smaller portion of well-aligned dimension of subspaces carry transferable signals with rotation, whereas full-dimensional representational transformations are explained by non-rotational reconfigurations (such as the relative representational distance).

378 Combined with the RSA results, these findings indicate that, while encoding and retrieval share some
 379 stable representational subspaces in higher-order regions, the representational transformation in the
 380 low- to mid-level regions undergoes systematic transformation, in both rotational and non-rotational
 381 manners. Together, these findings indicate a WM coding strategy that leverages dynamic transfor-
 382 mations in lower-level regions to minimize interference, while maintaining stable representations in
 383 higher-order regions to preserve reliability.

4.2.4 THE EFFECT OF TASK OBJECTIVES AND RECURRENT ARCHITECTURES IN ROTATION



414 Figure 4: Geometry analysis results. (A) Rotation angles of the top-1 axis and top-2 axes (mean)
 415 from encoding to retrieval in each ROI. Error bars represent ± 1 s.e.m across 8 subjects.(B) Rota-
 416 tion angles of the top-1 axis and top-2 axes (mean) from encoding to retrieval in each layer of the
 417 models. Error bars represent ± 1 s.e.m across 12 models. (C) Group average rotation angles of the
 418 top-1 and top-2 axes (mean) from encoding to retrieval in each layer of the models. SL denotes
 419 supervised, SSL denotes self-supervised. We averaged the rotation angles of Resnet-50 and Vit-
 420 B/16, since the patterns were consistent across encoder architectures. Error bar represent ± 1 s.e.m
 421 across models with varying encoder architectures.(D) Example illustrations of the representation
 422 transformation. The example representation space in the encoding phase (1), and how rotational
 423 or/and non-rotational transformation from encoding space to retrieval space would influence RSA
 424 and rotation angle θ (2-4). - denotes no difference, \uparrow denotes increase, \downarrow denotes decrease.

426 To have a direct comparison between the human brains and artificial systems, we next applied the
 427 same subspace rotation analysis on the RNNs. Interestingly, the models also exhibited smaller
 428 rotation angles (greater subspace alignment) in the higher hidden layer compared to the lower layers
 429 (Fig. 4B) similar to the brain data. The ANOVA revealed a significant main effect of the layer on
 430 the rotation angle ($p < .001$), without an interaction between layer and the angle types ($p > .05$).
 431 However, there was no significant difference between layers in post-doc paired t-tests ($p < .05$),
 suggesting the hierarchical pattern of rotation in model might not as consistent as humans.

432 We further investigate the factors under the divergence of the rotation angles in human and machines
433 by separately looking into the results of models with different learning objectives of the encoder
434 (Supervised learning -SL vs Self-Supervised learning -SSL), encoder architectures (ResNet vs ViT),
435 and recurrent architectures (vanilla RNN, GRU and LSTM) (Fig. 4C).

436 **The effect of learning objectives.** Overall, models that used supervised learning (SL) encoder
437 showed shared mechanism with humans than those that used self-supervised learning (SSL): For the
438 SL models, the project layer showed larger rotational angles which parallel with early to mid-level
439 ROIs (V1-hV4, VO-IPS), while the hidden layer exhibited smaller rotations similar to the higher-
440 order regions (FEF, dIPFC; upper row of Fig. 4C compared to the lower row). This layer-dependent
441 shift parallels the hierarchical gradient observed in the brain (Fig. 4A). By contrast, self-supervised
442 models did not display such gradient across layers: The rotation angle in the project layer was not
443 consistently higher than the hidden layers as in humans, suggesting that self-supervised pretraining
444 induces representational rotations that diverge more strongly from the transformations measured in
445 cortex.

446 **The effect of architectures.** Among the recurrent architectures, we also observed systematic but
447 distinct rotation patterns. GRU and LSTM exhibited the layer-dependent shift (Fig. 4C columns 2-
448 3): its project layer showed larger rotation angles which is similar to human early and mid-level vi-
449 sual ROIs (V1-hV4, VO, LO, OTS), while hidden layers showed smaller rotational angles mapping
450 comparably onto higher-order regions (FEF, dIPFC). By contrast, vanilla RNN models (especially
451 the self supervised ones) exhibited a different profile: the project layer showed smaller rotation
452 angles compared to the hidden layers. Additionally, we didn't observe the effect of the encoder
453 architecture (ResNet vs ViT) in the rotational transformation between WM phases.

454 5 DISCUSSION

455 **Hierarchical mixture of stable and dynamically updated memory subspaces.** Our findings sup-
456 port the view that WM representations combine both stable and dynamic components. Lower-level
457 visual areas tend to reformat representations across phases, whereas higher-order regions such as
458 prefrontal cortex exhibit more stable coding. The co-existence of dynamic and stable neural code
459 echoes the electrophysiological studies suggesting monkey's prefrontal cortex could adapt to behav-
460 ioral demands (Stokes et al., 2013; Murray et al., 2017), whereas the current work consider various
461 human ROIs, which parallels the layers in neural networks. The increasing stability of WM repre-
462 sentations along the cortical hierarchy is consistent with recent fMRI work on WM representations
463 using simpler stimuli (Li & Curtis, 2023). Additionally, merging evidence in comparing perception
464 to visual imaginary or long-term memory have also shown a stable representations across phases in
465 high-level cortical areas (Favila et al., 2022; Breedlove et al., 2020). Together, these findings sug-
466 gest that stable coding may emerge prominently in higher-order cortex for higher reliability, whereas
467 low-level areas undergo larger transformations for lower interference.

468 **The impact learning objectives in rotational WM transformation.** By training artificial neural
469 networks on the same WM tasks performed by humans, we observed greater rotational transforma-
470 tion of encoding-retrieval representations in the early project layer than in the late hidden layers,
471 paralleling the results in human brains. However, this shared mechanism between human and ma-
472 chine strongly depended on the learning objectives of the encoders. Supervised encoders preserved a
473 hierarchical project-hidden dissociation that paralleled the cortical gradient, whereas self-supervised
474 encoders diverged more strongly from this hierarchical structure. This suggests that supervised train-
475 ing enforces categorical abstraction across layers, which in turn constrains how memory subspaces
476 are rotated across time. However, self-supervised contrastive objectives prioritize invariance over
477 categorical structure, which may distort the temporal transformations relevant for WM (Konkle &
478 Alvarez, 2022).

479 Pervious studies in brain-model alignment have centered on assessing models based on its brain
480 activation predictability, and reported mixed results comparing SL and SSL models. Several studies
481 have found advantages of certain SSL models (Zhuang et al., 2021), while others highlights simi-
482 larities or even SL advantage (Conwell et al., 2024; Konkle & Alvarez, 2022; Khaligh-Razavi &
483 Kriegeskorte, 2014). The current study did not evaluate models based on the brain predictability,
484 whereas we investigated whether models exhibit shared transformation across WM phases, leading
485 to a different alignment criterion in our work. On the other hand, when we consider the emer-

486 gence of human-like representational transformation mechanism as a subset of the above topic, our
487 conclusion aligns with previous studies showing SL advantage.

488
489 **Gated recurrent mechanism supports human-like WM transformation.** Across different types
490 of recurrent modules, GRUs and LSTMs showed decreasing rotational transformation from WM
491 encoding to retrieval phases across layers, similar to humans. In contrast, vanilla RNNs exhibited
492 inconsistent patterns to humans, suggesting that insufficient control over information flow may fail
493 to capture the balance of stability and dynamics observed in the cortex during WM. Overall, these
494 results indicate that gated recurrent dynamics in GRUs and LSTMs are essential for modeling the
495 coexistence of stable and dynamic subspaces in the brain.

496 **Out-of-distribution generalization.** In the current study, the models were trained and validated
497 on a subset of the NSD-core, while tested on NSD-synthetic, an out-of-distribution (OOD) dataset.
498 Model performance was lower than on in-distribution NSD-core, highlighting a distribution shift.
499 This likely reflects two potential limitations: (I) the frozen encoder may not fully capture task-
500 relevant features for novel stimuli, and (II) recurrent backbones tend to maintain features learned
501 during training, which can be hard to adapt to unseen inputs. In contrast, humans readily generalize
502 WM across diverse visual inputs, reflecting the brain’s ability to flexibly abstract stimulus features
503 while preserving task-relevant geometry. Prior work shows that human vision achieves robust OOD
504 generalization by leveraging hierarchical abstractions and flexible context adaptation (Fang & Sims,
505 2025). These comparisons underscore that while models can approximate in-distribution represen-
506 tational dynamics, OOD generalization remains a challenge for both supervised and self-supervised
507 vision models (Geirhos et al., 2021). Future work should explore adaptive or hierarchical encoders
508 that can dynamically preserve task-relevant features across domains.

509 6 LIMITATIONS

510
511 The current work focused on the 1-back task, which is widely used to investigate the WM neural
512 mechanisms, combined with human fMRI (Malisza et al., 2005; Ricciardi et al., 2006; Lee et al.,
513 2013; Ateş et al., 2017) and EEG (Audrain et al., 2020; Gjini et al., 2007). Including multiple
514 WM tasks would further strengthen the generalization of our findings in WM transformation across
515 encoding to retrieval phases. However, NSD-synthetic 1-back dataset (Gifford et al., 2025) is cur-
516 rently the only available WM benchmark that provides high-resolution, single-trial fMRI data on
517 naturalistic stimuli. Because the focus of the present study is on shared mechanisms between mod-
518 els and humans, we are therefore constrained by the availability of such biological ground truth. To
519 partially address this limitation on the modeling side, we analyzed models with fine-tuned encoders
520 trained on the 1-back task (Fig. A5) and additionally trained new models on a 2-back task (Fig. A4).
521 The results converged to the same central conclusion: hierarchical representational transformations
522 remain consistent regardless of task difficulty or whether encoders are frozen or fine-tuned.

523 The NSD-core (Allen et al., 2022) contains over 73,000 images, we sampled only 600 to build
524 the 1-back training sequences. This choice kept training computationally feasible across multiple
525 architecture-objective combinations. The subset spans all 80 categories, ensuring diversity while
526 reducing training time. While more samples might improve performance, we expect relative com-
527 parisons (e.g., gated vs. non-gated recurrent modules, supervised vs. self-supervised encoders) to
528 hold since all models used the same dataset, while larger training and testing datasets would strength
529 our conclusions. Additionally, though models perform the 1-back task well, their error patterns dif-
530 ferred from humans. Several models showed a higher bias toward predicting “match” compared to
531 human participants. Thus, while the models captured aspects of WM representational transforma-
532 tion across phases, their ability to regulate decisions was weaker than humans. Though the bias
533 result wasn’t link to the model-brain alignment in the WM representation transformation, future
534 work would benefit from comparing computational principles between humans and models exhibit
535 more human-like behavioral patterns.

540 ETHICS AND REPRODUCIBILITY STATEMENT

541
542 **Ethics statement.** This work uses a publicly available fMRI dataset released by Allen et al. (2022)
543 and Gifford et al. (2025). As reported in the original dataset publication, informed written consent
544 was obtained from all participants, and the experimental protocol was approved by the University
545 of Minnesota Institutional Review Board. Our study did not involve any new data collection with
546 human subjects. We adhere to the ICLR Code of Ethics in conducting and presenting this research.

547 **Reproducibility statement.** We provide detailed descriptions of the data preparation procedures
548 and model training settings to ensure transparency and reproducibility. In addition, we plan to
549 publicly release the processed data and the code used for analyses for reproducibility.

550
551 REFERENCES

- 552
553 Emily J Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Jacob S Prince, Logan T Dowdle,
554 Matthias Nau, Brad Caron, Franco Pestilli, Ian Charest, et al. A massive 7t fmri dataset to bridge
555 cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, 2022.
- 556 Fatma Ebru Ateş, Banu Cangöz, Erguvan Tuğba Özel Kızıl, Bora Baskak, Zeynel Baran, and
557 Halise Devrimci Özgüven. Frontal activity during a verbal emotional working memory task in
558 patients with alzheimer’s disease: A functional near-infrared spectroscopy study. *Psychiatry Re-*
559 *search: Neuroimaging*, 261:29–34, 2017. doi: 10.1016/j.psychresns.2016.12.013.
- 560
561 Samantha P Audrain, Charline M Urbain, Veronica Yuk, Rachel C Leung, Simeon M Wong, and
562 Margot J Taylor. Frequency-specific neural synchrony in autism during memory encoding, main-
563 tenance and recognition. *Brain Communications*, 2(2):fcaa094, 2020. doi: 10.1093/braincomms/
564 fcaa094.
- 565 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*
566 *arXiv:1607.06450*, 2016.
- 567
568 Alan Baddeley. Working memory. *Science*, 255(5044):556–559, 1992. doi: 10.1126/science.
569 1736359.
- 570
571 Alan Baddeley. Working memory: looking back and looking forward. *Nature reviews neuroscience*,
572 4(10):829–839, 2003. doi: 10.1038/nrn1201.
- 573 Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful
574 approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*,
575 57(1):289–300, 1995.
- 576
577 Flora Bouchacourt and Timothy J Buschman. A flexible model of working memory. *Neuron*, 103
578 (1):147–160, 2019. doi: 10.1016/j.neuron.2019.04.020.
- 579 Jesse L Breedlove, Ghislain St-Yves, Cheryl A Olman, and Thomas Naselaris. Generative feedback
580 explains distinct brain activity codes for seen and mental images. *Current Biology*, 30(12):2211–
581 2224, 2020. doi: 10.1016/j.cub.2020.04.014ExternalLink.
- 582
583 Xinlei Chen*, Saining Xie*, and Kaiming He. An empirical study of training self-supervised vision
584 transformers. *arXiv preprint arXiv:2104.02057*, 2021.
- 585
586 Thomas B Christophel, P Christiaan Klink, Bernhard Spitzer, Pieter R Roelfsema, and John-Dylan
587 Haynes. The distributed nature of working memory. *Trends in cognitive sciences*, 21(2):111–124,
588 2017. doi: 10.1016/j.tics.2016.12.007.
- 589 Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of
590 gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- 591
592 Anne GE Collins and Michael J Frank. How much of reinforcement learning is working memory,
593 not reinforcement learning? a behavioral, computational, and neurogenetic analysis. *European*
Journal of Neuroscience, 35(7):1024–1035, 2012. doi: 10.1111/j.1460-9568.2011.07980.x.

- 594 Albert Compte, Nicolas Brunel, Patricia S Goldman-Rakic, and Xiao-Jing Wang. Synaptic mech-
595 anisms and network dynamics underlying spatial working memory in a cortical network model.
596 *Cerebral cortex*, 10(9):910–923, 2000. doi: 10.1093/cercor/10.9.910.
597
- 598 Colin Conwell, Jacob S Prince, Kendrick N Kay, George A Alvarez, and Talia Konkle. A large-scale
599 examination of inductive biases shaping high-level visual representation in brains and machines.
600 *Nature communications*, 15(1):9383, 2024. doi: 10.1038/s41467-024-53147-y.
- 601 Roshan Cools and Mark D’Esposito. Inverted-u-shaped dopamine actions on human working
602 memory and cognitive control. *Biological psychiatry*, 69(12):e113–e125, 2011. doi: 10.1016/
603 j.biopsych.2011.03.028.
604
- 605 Meredyth Daneman and Patricia A Carpenter. Individual differences in working memory and
606 reading. *Journal of verbal learning and verbal behavior*, 19(4):450–466, 1980. doi: 10.1016/
607 S0022-5371(80)90312-6.
- 608 Jonas Karolis Degutis, Simon Weber, Joram Soch, and John-Dylan Haynes. Neural dynamics of
609 visual working memory representation during sensory distraction. *Elife*, 13:RP99290, 2025. doi:
610 <https://doi.org/10.7554/eLife.99290.4>.
611
- 612 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
613 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
614 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
615 *arXiv:2010.11929*, 2020.
- 616 Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
617
- 618 Randall W Engle. Role of working-memory capacity in cognitive control. *Current anthropology*,
619 51(S1):S17–S26, 2010. doi: 10.1086/650572.
- 620 Jose M Esnaola-Acebes, Alex Roxin, and Klaus Wimmer. Flexible integration of continuous sensory
621 evidence in perceptual estimation tasks. *Proceedings of the National Academy of Sciences*, 119
622 (45):e2214441119, 2022. doi: 10.1073/pnas.2214441119.
623
- 624 Zeming Fang and Chris R Sims. Humans learn generalizable representations through efficient cod-
625 ing. *Nature Communications*, 16(1):3989, 2025. doi: 10.1038/s41467-025-58848-6.
626
- 627 Serra E Favila, Brice A Kuhl, and Jonathan Winawer. Perception and memory have distinct spatial
628 tuning properties in human visual cortex. *Nature communications*, 13(1):5864, 2022. doi: 10.
629 1038/s41467-022-33161-8.
- 630 Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge,
631 Felix A Wichmann, and Wieland Brendel. Partial success in closing the gap between human and
632 machine vision. *Advances in Neural Information Processing Systems*, 34:23885–23899, 2021.
633 doi: 10.48550/arXiv.2106.07411.
- 634 Alessandro T Gifford, Radoslaw M Cichy, Thomas Naselaris, and Kendrick Kay. A 7t fmri dataset
635 of synthetic images for out-of-distribution modeling of vision. *arXiv preprint arXiv:2503.06286*,
636 2025. doi: 10.48550/arXiv.2503.06286.
637
- 638 Klevest Gjini, Takashi Maeno, Keiji Iramina, and Shoogo Ueno. A multichannel whole-head eeg
639 study on visual working memory processing of spatiality in the human brain. In *World Congress*
640 *on Medical Physics and Biomedical Engineering 2006: August 27–September 1, 2006 COEX*
641 *Seoul, Korea “Imaging the Future Medicine”*, pp. 2752–2755. Springer, 2007. doi: 10.1007/
642 978-3-540-36841-0_694.
- 643 Matthew F Glasser, Timothy S Coalson, Emma C Robinson, Carl D Hacker, John Harwell, Essa
644 Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F Beckmann, Mark Jenkinson, et al. A
645 multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171–178, 2016.
646
- 647 John C Gower. Generalized procrustes analysis. *Psychometrika*, 40(1):33–51, 1975. doi: 10.1007/
BF02291478.

- 648 Stephenie A Harrison and Frank Tong. Decoding reveals the contents of visual working memory in
649 early visual areas. *Nature*, 458(7238):632–635, 2009. doi: 10.1038/nature07832.
- 650
- 651 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
652 human-level performance on imagenet classification. In *Proceedings of the IEEE international
653 conference on computer vision*, pp. 1026–1034, 2015.
- 654
- 655 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
656 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition
657 (CVPR)*, June 2016.
- 658
- 659 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
660 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on
661 Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 662
- 663 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
664 1735–1780, 1997.
- 665
- 666 Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsupervised,
667 models may explain it cortical representation. *PLoS computational biology*, 10(11):e1003915,
668 2014. doi: 10.1371/journal.pcbi.1003915.
- 669
- 670 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
671 arXiv:1412.6980*, 2014.
- 672
- 673 Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for
674 human ventral stream representation. *Nature communications*, 13(1):491, 2022.
- 675
- 676 Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. Representational similarity analysis-
677 connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:249, 2008.
678 doi: 10.3389/neuro.06.004.2008.
- 679
- 680 Yuna Kwak and Clayton E Curtis. Unveiling the abstract format of mnemonic representations.
681 *Neuron*, 110(11):1822–1828, 2022. doi: 10.1016/j.neuron.2022.03.016.
- 682
- 683 Tien-Wen Lee, Ho-Ling Liu, Yau-Yau Wai, Han-Jung Ko, and Shwu-Hua Lee. Abnormal neural
684 activity in partially remitted late-onset depression: an fmri study of one-back working memory
685 task. *Psychiatry Research: Neuroimaging*, 213(2):133–141, 2013. doi: 10.1016/j.pscychresns.
686 2012.04.010.
- 687
- 688 Xiaoxuan Lei, Takuya Ito, and Pouya Bashivan. Geometry of naturalistic object representations in
689 recurrent neural network models of working memory. *Advances in Neural Information Processing
690 Systems*, 37:100604–100629, 2024. doi: 10.48550/arXiv.2411.02685.
- 691
- 692 Hsin-Hung Li and Clayton E Curtis. Neural population dynamics of human working memory. *Cur-
693 rent Biology*, 33(17):3775–3784, 2023. doi: 10.1016/j.cub.2023.07.067.
- 694
- 695 Alexandra Libby and Timothy J Buschman. Rotational dynamics reduce interference between sensory
696 and memory representations. *Nature neuroscience*, 24(5):715–726, 2021. doi: 10.1038/
697 s41593-021-00821-9.
- 698
- 699 Steven J Luck and Edward K Vogel. The capacity of visual working memory for features and
700 conjunctions. *Nature*, 390(6657):279–281, 1997. doi: 10.1038/36846.
- 701
- 702 Wei Ji Ma, Masud Husain, and Paul M Bays. Changing concepts of working memory. *Nature
703 neuroscience*, 17(3):347–356, 2014. doi: 10.1038/nn.3655.
- 704
- 705 Svein Magnussen, Espen Idås, and Steinar Holst Myhre. Representation of orientation and spatial
706 frequency in perception and memory: a choice reaction-time analysis. *Journal of Experimental
707 Psychology: Human perception and performance*, 24(3):707, 1998. doi: [https://doi.org/10.1037/
708 0096-1523.24.3.707](https://doi.org/10.1037/0096-1523.24.3.707).

- 702 Krisztina L Malisza, Ava-Ann Allman, Deborah Shiloff, Lorna Jakobson, Sally Longstaffe, and
703 Albert E Chudley. Evaluation of spatial working memory function in children and adults with fetal
704 alcohol spectrum disorders: a functional magnetic resonance imaging study. *Pediatric research*,
705 58(6):1150–1157, 2005. doi: 10.1203/01.pdr.0000185479.92484.a1.
- 706
707 John D Murray, Alberto Bernacchia, Nicholas A Roy, Christos Constantinidis, Ranulfo Romo, and
708 Xiao-Jing Wang. Stable population coding for working memory coexists with heterogeneous
709 neural dynamics in prefrontal cortex. *Proceedings of the National Academy of Sciences*, 114(2):
710 394–399, 2017. doi: <https://doi.org/10.1073/pnas.1619449114>.
- 711 Tomoya Nakamura, Seng Bum Michael Yoo, Kendrick Kay, Hakwan Lau, and Ali Moharramipour.
712 Representational geometries of perception and working memory: A pilot study. *bioRxiv*, pp.
713 2025–09, 2025. doi: <https://doi.org/10.1101/2025.09.07.674590>.
- 714
715 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
716 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
717 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 718 Yoni Pertzov, Sanjay Manohar, and Masud Husain. Rapid forgetting results from competition over
719 time between items in visual working memory. *Journal of Experimental Psychology: Learning*,
720 *Memory, and Cognition*, 43(4):528, 2017. doi: 10.1037/xlm0000328.
- 721
722 Michael Petrides. Lateral prefrontal cortex: architectonic and functional organization. *Philosophical*
723 *Transactions of the Royal Society B: Biological Sciences*, 360(1456):781–795, 2005.
- 724
725 Emilia P Piwek, Mark G Stokes, and Christopher Summerfield. A recurrent neural network model
726 of prefrontal brain activity during a working memory task. *PLoS Computational Biology*, 19(10):
727 e1011555, 2023. doi: 10.1371/journal.pcbi.1011555.
- 728
729 Jacob S Prince, Ian Charest, Jan W Kurzwaski, John A Pyles, Michael J Tarr, and Kendrick N Kay.
730 Improving the accuracy of single-trial fmri response estimates using glmsingle. *Elife*, 11:e77599,
731 2022. doi: 10.7554/eLife.77599.
- 732
733 Emiliano Ricciardi, Daniela Bonino, Claudio Gentili, Lorenzo Sani, Pietro Pietrini, and Tomaso
734 Vecchi. Neural correlates of spatial working memory in humans: a functional magnetic resonance
735 imaging study comparing visual and tactile processes. *Neuroscience*, 139(1):339–349, 2006. doi:
736 10.1016/j.neuroscience.2005.08.045.
- 737
738 John T Serences. Neural mechanisms of information storage in visual short-term memory. *Vision*
739 *research*, 128:53–67, 2016. doi: 10.1016/j.visres.2016.09.010.
- 740
741 Hongsup Shin, Qijia Zou, and Wei Ji Ma. The effects of delay duration on visual working memory
742 for orientation. *Journal of vision*, 17(14):10–10, 2017. doi: <https://doi.org/10.1167/17.14.10>.
- 743
744 Eelke Spaak, Kei Watanabe, Shintaro Funahashi, and Mark G Stokes. Stable and dynamic coding
745 for working memory in primate prefrontal cortex. *Journal of neuroscience*, 37(27):6503–6516,
746 2017. doi: 10.1523/JNEUROSCI.3364-16.2017.
- 747
748 Mark G Stokes, Makoto Kusunoki, Natasha Sigala, Hamed Nili, David Gaffan, and John Duncan.
749 Dynamic coding for cognitive control in prefrontal cortex. *Neuron*, 78(2):364–375, 2013. doi:
750 10.1016/j.neuron.2013.01.039.
- 751
752 Heinz-Martin Süß, Klaus Oberauer, Werner W Wittmann, Oliver Wilhelm, and Ralf Schulze.
753 Working-memory capacity explains reasoning ability—and a little bit more. *Intelligence*, 30(3):
754 261–288, 2002. doi: 10.1016/S0160-2896(01)00100-3.
- 755
756 Anthony D Wagner. Working memory contributions to human learning and remembering. *Neuron*,
757 22(1):19–22, 1999. doi: [https://doi.org/10.1016/s0896-6273\(00\)80674-1](https://doi.org/10.1016/s0896-6273(00)80674-1).
- 758
759 Quan Wan, Ying Cai, Jason Samaha, and Bradley R Postle. Tracking stimulus representation across
760 a 2-back visual working memory task. *Royal Society open science*, 7(8):190228, 2020. doi:
761 10.1098/rsos.190228.

756 Quan Wan, Jorge A Menendez, and Bradley R Postle. Priority-based transformations of stimulus
757 representation in visual working memory. *PLoS Computational Biology*, 18(6):e1009062, 2022.
758 doi: 10.1371/journal.pcbi.1009062.
759

760 Liang Wang, Ryan EB Mruzek, Michael J Arcaro, and Sabine Kastner. Probabilistic maps of visual
761 topography in human cortex. *Cerebral cortex*, 25(10):3911–3931, 2015.

762 Xiao-Jing Wang. Synaptic basis of cortical persistent activity: the importance of nmda receptors to
763 working memory. *Journal of Neuroscience*, 19(21):9587–9603, 1999. doi: 10.1523/JNEUROSCI.
764 19-21-09587.1999.

765 Xiao-Jing Wang. 50 years of mnemonic persistent activity: quo vadis? *Trends in Neurosciences*, 44
766 (11):888–902, 2021. doi: 10.1016/j.tins.2021.09.001.
767

768 Klaus Wimmer, Duane Q Nykamp, Christos Constantinidis, and Albert Compte. Bump attractor
769 dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nature*
770 *neuroscience*, 17(3):431–439, 2014. doi: 10.1038/nn.3645.

771 Yaoda Xu. The human posterior parietal cortices orthogonalize the representation of different
772 streams of information concurrently coded in visual working memory. *PLoS Biology*, 22(11):
773 e3002915, 2024. doi: 10.1371/journal.pbio.3002915.
774

775 Yaoda Xu. Transformed visual working memory representations in human occipitotemporal and
776 posterior parietal cortices. *eneuro*, 12(7), 2025. doi: 10.1523/ENEURO.0162-25.2025.

777 Guangyu Robert Yang and Xiao-Jing Wang. Artificial neural networks for neuroscientists: a primer.
778 *Neuron*, 107(6):1048–1070, 2020. doi: 10.1016/j.neuron.2020.09.005.
779

780 Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J Di-
781 Carlo, and Daniel LK Yamins. Unsupervised neural network models of the ventral visual
782 stream. *Proceedings of the National Academy of Sciences*, 118(3):e2014196118, 2021. doi:
783 10.1073/pnas.2014196118.
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A APPENDIX

A.1 LLM USAGE DISCLOSURE

Large language models (LLMs) were employed solely as a tool to refine wording, improve clarity, and enhance the readability of the manuscript. LLMs were not used to generate original ideas, analyses, or conclusions. All substantive content, arguments, and interpretations presented in this work are from authors' own.

A.2 NEURAL RECORDINGS AND TASKS

The NSD-synthetic dataset extends the Natural Scenes Dataset (NSD-core) with high-resolution 7T fMRI recordings, providing single-trial beta estimates across fixation and 1-back tasks with new synthetic stimuli (Gifford et al., 2025).

Participants and sequences. The NSD-synthetic dataset consisted of data from eight subjects who had completed the NSD-core experiment (Gifford et al., 2025) and subsequently participated in the NSD-synthetic experiment, in which 284 synthetic images were tested. The whole-brain BOLD was acquired at 1.8 mm isotropic resolution (TR = 1.6 s) using gradient-echo EPI sequences on a Siemens 7 T scanner.

Stimuli. 284 synthetic images were generated to vary in visual format and semantic content, including 8 subclasses - noise, natural scenes, manipulated scenes, contrast modulation, phase-coherence modulation, words, spiral gratings, chromatic noise (Gifford et al., 2025). The manipulated scenes are the transformed versions of the natural scenes, such as upside-down, Mooney and lie-drawing scenes. This design enabled systematic manipulation of stimulus identity and format.

Tasks. Two tasks were performed in alternating runs, consisting of 4 runs (93 trials per run). Each trial started with the presentation of an image for 2 s, followed by a 2-s inter-trial interval. (1) In the Fixation task, subjects detected color changes in the fixation dot and immediately pressed button, providing a baseline measure of pure perceptual representation. (2) 1-back WM task, in which subjects pressed a button to indicate whether the current image matched the immediately preceding image in identity, thus dissociating the encoding and retrieval stages.

Preprocessing and GLM. We used the preprocessed single-trial GLM beta estimates provided by NSD-synthetic (1.8mm with TR = 1.33s). The beta weights were estimated using GLMs—employing optimized HRFs, GLMdenoise, and ridge-regularized regression—to yield percent signal change estimates for each image under both task conditions (Prince et al., 2022).

Regions of interest. We included 10 regions of interest (ROIs) in various regions, from retinotopic visual regions to higher-order cognition regions. For V1, V2, V3 (grouping ventral and dorsal subregions for each corresponding ROI) and hV4, we used regions provided from NSD, which were manually drawn for each subject based on pRF experiment (Allen et al., 2022). Lateral occipital (LO; grouping LO1 and LO2), ventral occipital (VO; grouping VO1 and VO2) and intraparietal sulcus (IPS; grouping IPS0 to IPS5) were defined using maximum probability map (Wang et al., 2015). Occipito-temporal sulcus (OTS) was defined using “corticalsulc” ROIs provided by NSD (Allen et al., 2022). Frontal eye field (FEF) and dorsolateral prefrontal cortex (dlPFC) were defined using the HCP-MMP1.0 atlas by Glasser et al. (2016). To be specific, the dlPFC was defined as a grouped region of Brodmann areas 9, 46 and 9/46 (Petrides, 2005).

A.3 1-BACK TASK MODELS DETAILS

In this section, we provide additional details to ensure full reproducibility of our experiments. We first expand on the design choices for the encoder and recurrent modules, followed by dataset construction, preprocessing, and training hyperparameters. Unless otherwise noted, all models were implemented in PyTorch (Paszke et al., 2019).

A.3.1 ENCODERS

ResNet-50. We used the ResNet-50 implementation (He et al., 2016) from `torchvision.models`, with pretrained weights from two sources. (1) Supervised weights

864 trained on ImageNet-1K (`weights=ResNet50.Weights.IMAGENET1K_V2`) were down-
 865 loaded from PyTorch Hub¹. (2) Self-supervised weights were obtained from MoCo v2 (He et al.,
 866 2020), released in the official GitHub repository², where we selected the 800-epoch pretrained
 867 checkpoint. For ResNet-50, features were extracted from the output of the last residual block
 868 (`layer4.2.relu_2`), which yields a 2048-channel representation.

870 **ViT-B/16.** We used the ViT-B/16 implementation (Dosovitskiy et al., 2020) from
 871 `torchvision.models`, again with pretrained weights from two sources. (1) Supervised weights
 872 trained on ImageNet-1K were provided in PyTorch Hub³. (2) Self-supervised weights were obtained
 873 from MoCo v3 (Chen* et al., 2021), released in the official GitHub repository⁴, where we used
 874 the ViT checkpoint pretrained for 300 epochs. For ViT, features were extracted from the output of
 875 the final (12th) transformer encoder block (`encoder.layers.encoder_layer.11.add.1`),
 876 which yields a 768-dimensional representation.

877 A.3.2 RECURRENT MODULES

879 **Vanilla RNN.** We implemented a two-layer vanilla RNN (Elman, 1990) with input size 512 and
 880 hidden size 256, using PyTorch’s `nn.RNN`. The hidden state at each time step h_t was updated by
 881 combining the projected encoder feature f_t and the previous hidden state h_{t-1} through a linear
 882 transformation followed by \tanh activation. Dropout with probability $p = 0.3$ was applied to the
 883 hidden states of each layer. Parameters were initialized using Kaiming uniform initialization (He
 884 et al., 2015). At each time step, the hidden state from the last RNN layer was passed through a fully
 885 connected layer and sigmoid activation to yield a binary prediction \hat{y}_t .

887 **LSTM.** The long short-term memory (LSTM) (Hochreiter & Schmidhuber, 1997) was imple-
 888 mented with two layers and hidden size 256 with input size 512, using PyTorch’s `nn.LSTM`. We
 889 adopted the standard gating mechanism (input, forget, and output gates) as in PyTorch’s `nn.LSTM`
 890 implementation. The same dropout ($p = 0.3$) and initialization scheme as the vanilla RNN were
 891 used.

892 **GRU.** The gated recurrent unit (GRU) (Chung et al., 2014) was similarly implemented with two
 893 layers and hidden size 256 with input size 512, using PyTorch’s `nn.GRU`. The GRU simplifies the
 894 LSTM by merging input and forget gates into a single update gate, thereby reducing parameter
 895 count. As with the other recurrent models, dropout was applied to hidden states and predictions
 896 were generated via a fully connected layer followed by sigmoid activation.

898 A.3.3 ENCODER-RNN CONNECTION

899 To bridge the encoder and recurrent modules, we extracted image features from the final representa-
 900 tional layer of each encoder (ResNet-50: `layer4.2.relu_2`; ViT-B/16: `encoder_layer.11`).
 901 These features were high-dimensional (2048 channels for ResNet, 768 for ViT) and were adapted
 902 for recurrent processing through a projection module. Specifically, we first applied a pointwise
 903 (1×1) convolution to reduce the feature dimensionality to match the RNN input size. The resulting
 904 features were then passed through a fully connected layer and subsequently normalized with layer
 905 normalization (Ba et al., 2016). We denote this processed representation as the *projected feature*
 906 *layer* (project layer), which served as the input to the recurrent module.

908 A.3.4 DATASET GENERALIZATION

909 Images were sampled from the NSD-core dataset (Allen et al., 2022) and preprocessed following
 910 NSD-synthetic conventions (Gifford et al., 2025). Each image was center-cropped to a square using
 911 the smallest dimension and resized to 224×224 pixels. Synthetic images were incorporated by

913 ¹[https://docs.pytorch.org/vision/main/models/generated/torchvision.
 914 models.resnet50.html#torchvision.models.resnet50](https://docs.pytorch.org/vision/main/models/generated/torchvision.models.resnet50.html#torchvision.models.resnet50)

915 ²<https://github.com/facebookresearch/moco>

916 ³[https://docs.pytorch.org/vision/main/models/generated/torchvision.
 917 models.vit_b_16.html#torchvision.models.ViT_B_16_Weights](https://docs.pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html#torchvision.models.ViT_B_16_Weights)

⁴<https://github.com/facebookresearch/moco-v3>

concatenating the official NSD-synthetic with an additional collection of chromatic noise images. Invalid indices (e.g., corrupted crops) were removed prior to concatenation.

Training and validation sequences were generated using a 1-back trial generator. Each run consisted of 93 image presentations, with repeat events ($y_t = 1$ if $x_t = x_{t-1}$, otherwise $y_t = 0$) introduced at predefined probabilities $\{0.02, 0.1, 0.3, 0.5, 0.7\}$. The training set contained 40 runs (8 per probability condition), while the validation set contained 10 runs. During the data generation, we set random seed as 42.

rapped in PyTorch Dataset objects and loaded with the DataLoader API (batch size 1, shuffled). Each batch yielded sequences of images and labels suitable for recurrent training.

A.3.5 MODEL TRAINING

To evaluate generalization, we trained and validated models on a subset of the Natural Scenes Dataset (NSD-core) (Allen et al., 2022), ensuring no overlap with the NSD-synthetic test set. From the 73,000 available images, we randomly sampled 600 spanning all 80 categories, which were used to generate 1-back training sequences. The training set comprised 40 runs and the validation set 10 runs, each with 93 image presentations. Repeat events (label = 1, when $x_t = x_{t-1}$) were inserted at controlled probabilities. All images were preprocessed following NSD-synthetic conventions (Gifford et al., 2025): center-cropped to a square using the shortest dimension and resized to 224×224 pixels.

Model parameters were optimized using Adam (Kingma & Ba, 2014). We used an initial learning rate of 1×10^{-3} for GRU and LSTM parameters, 5×10^{-4} for vanilla RNN parameters, and 2×10^{-3} for the projection and classifier layers. A weight decay of 1×10^{-4} was applied to all trainable parameters. Encoders remained frozen with pretrained weights, while recurrent modules were initialized with Kaiming initialization (He et al., 2015). Dropout ($p = 0.3$) was applied to RNN hidden states, the pointwise convolution, and the fully connected layers. All models were optimized with binary cross-entropy loss. Training was performed with batch size 1 on a single NVIDIA A100 GPU. We trained all the models for 90 epochs, and models with the best validation accuracy were selected for the following analyses of our work.

A.3.6 RANDOM SEED TESTING

To ensure the stability of our 1-back models, we trained all architectures under multiple random initializations. Specifically, each model was trained using a distinct random seed (random seed set to 100, 200, 300, 400, 500) that controlled weight initialization, data shuffling, and three different RNN types state initialization. This procedure allowed us to assess whether the learned working-memory representations were consistent across training replicas.

For evaluation, we generated a new random test dataset using NSD synthetic images and a fixed master seed (random seed 2025). All trained replicas were tested on this identical test set, ensuring that performance variability reflected differences in the learned representations rather than stochasticity in test sequence construction. The test set followed the same structure as in the main 1-back protocol: 4 runs of continuous visual streams (sequence length $T = 93$), matched in stimulus composition across all model variants. We report both the mean and standard error of hit rate and false alarm rate across seeds for each architecture.

A.4 2-BACK TASK MODELS

To evaluate whether our architecture supports higher-order working memory demands, we extended the framework to a 2-back variant of the task. In this paradigm, the model must determine whether the current stimulus x_t matches the stimulus presented two steps earlier, x_{t-2} . The overall training and testing procedures closely mirror those used in the main 1-back experiment, with modifications described below.

A.5 MODEL TRAINING

We trained hybrid visual-temporal architectures that combine frozen visual encoders (ResNet-50, MoCo v2, ViT-B/32, and MoCo v3) with a GRU-based recurrent head. Unlike the main 1-back

972 setup—where training was performed solely on NSD-core images—the 2-back extension employed
 973 a *mixed N-back training curriculum*. Specifically, each training batch interleaved 1-back and 2-back
 974 trials, requiring the GRU to flexibly maintain and update memory states across variable temporal
 975 horizons.

976 Training used the AdamW optimizer with a cosine-annealing learning-rate schedule. A masked
 977 cross-entropy loss was applied to accommodate variable-length label sequences generated by the
 978 mixed-N curriculum. Consistent with the original human behavioral experiment, all models were
 979 trained with the same target-to-true ratio (0.02, 0.1, 0.3, 0.5, 0.7).
 980

981 A.6 MODEL TESTING

982 Testing procedures were designed to be directly comparable across all architectures. As in the main
 983 1-back experiment, we used NSD synthetic images as the evaluation stimuli. However, because
 984 no dedicated 2-back test split exists, we generated test sequences using fixed random seeds while
 985 maintaining the same evaluation structure used in the 1-back experiments.
 986

987
 988 **Protocol.** For each model, we conducted 10 independent test runs, each initialized with a dis-
 989 tinct master seed (1000–1009). Every run consisted of three blocks of continuous visual streams
 990 (sequence length $T = 93$), ensuring robustness to random variation in trial composition.
 991

992 **Stimuli.** Stimuli were constructed from NSD synthetic images following the same procedure as in
 993 the 1-back evaluation, but paired with 2-back labels.
 994

995 **Metrics.** For each model, we computed task accuracy and extracted activations from all layers
 996 and task phases. These activations were used to construct representational dissimilarity matrices
 997 (RDMs) for later representational-geometry analyses.
 998

999 A.6.1 2-BACK TASK PERFORMANCE

1000 Across all four architectures, models achieved over 95% training accuracy and roughly 75% val-
 1001 idation accuracy on the mixed 1-back/2-back training curriculum. When evaluated on randomly
 1002 generated NSD-synthetic test sequences. For testing, all model performance accuracy higher than
 1003 65% accuracy.
 1004

1005 A.7 ANALYSIS

1006 A.7.1 REPRESENTATIONAL SIMILARITY ANALYSIS

1007
 1008 **Representational dissimilarity matrices (RDMs).** RDMs enable the comparison of representa-
 1009 tion across systems. The dissimilarity between two representations X_i and X_j can be expressed as
 1010 $D(X_i, X_j) = 1 - \frac{\text{cov}(X_i, X_j)}{\sigma(X_i)\sigma(X_j)}$, where $\text{cov}(X_i, X_j)$ denotes the covariance between the two repre-
 1011 sentations, with $\sigma(X_i)$ being the standard deviation of X_i . This formula denotes a $1 - \text{Pearson}$
 1012 distance between the representations of stimuli i and j .
 1013

1014 The RDM is a symmetric $n \times n$ matrix R where R_{ij} reflects the dissimilarity between the represen-
 1015 tations of stimuli, resulting in a 284 by 284 matrix in the current study. Mathematically, the RDM is
 1016 given by $R_{ij} = D(X_i, X_j) \quad \forall i, j \in \{1, 2, \dots, n\}$.
 1017

1018 We constructed RDMs for each subject and each ROI in both encoding and retrieval phases of the
 1019 fMRI recordings. Similarly, we constructed RDMs for each model and each layer in two phases
 1020 based on the layer activations.
 1021

1022 **Representational Similarity Analysis (RSA).** RSA provides a common framework to quanti-
 1023 tatively compare relative representational distances as in RDMs across different modalities and
 1024 phases of representations (Kriegeskorte et al., 2008). To assess the similarity of neural and
 1025 model representations, we calculate the similarity score (Spearman’s ρ) between the RDMs: $\rho =$
 $\text{Spearman}(\text{vec}(R_A), \text{vec}(R_B))$.

We used Spearman correlation, between the encoding and the retrieval phase, of each brain ROI or model layer, to investigate how the relative representational distances across stimuli change from the encoding to the retrieval phase. Furthermore, we analyzed a path-based RSA to inspect the representational similarity changes along the visual hierarchy or model layers. Specifically, we treated the set of ROI-level or layer-level RDMs from each phase (encoding or retrieval) as a trajectory in representational space. RDMs were vectorized, and dissimilarities of RDMs across all ROIs (or layers) were computed as $1 - \rho$. The resulting dissimilarity matrix was then embedded in two dimensions using multidimensional scaling (MDS), allowing us to visualize each phase as a continuous path through representational geometry.

Shuffling RSA for reliability. To evaluate the reliability of the RSA, we performed the same set of RSA with randomly shuffling the rows and columns of the RDMs within each subject. For each ROI, permuted distributions were constructed by repeatedly shuffling the indices of both encoding and retrieval RDMs (100 iterations per subject), correlating the permuted RDM pairs, and averaging across permutations to yield a subject-level null estimate. Both observed and permuted values were compared using paired t-tests across subjects ($p < .001$; corrected for multiple comparisons using Benjamini–Hochberg method (Benjamini & Hochberg, 1995)).

A.7.2 DECODING AND GEOMETRY ANALYSIS

Decoding analysis. For decoding analysis, we evaluated how well the obtained fMRI beta weights from different WM phases could classify the stimulus identity (8 subclasses (Gifford et al., 2025)). For each ROI and subject, we selected the top $k = 100$ voxels using univariate ANOVA (f -test) of stimuli class on the encoding data, and used the same subset in retrieval. In parallel, we extracted activations from models, separately for encoding and retrieval phases, respectively for the project and hidden layer.

For each ROI or model layer, we trained linear decoders to classify 8 image subclasses from extracted activity. To prevent the classifier from being dominated by high-frequency classes, we re-weighted the loss inversely to class frequency during the training, and tested the decoding performance using the average of per-class recall accuracy. The chance level of the decoding accuracy is 0.125.

We used one-vs-rest linear SVMs with L2 regularization ($C=1.0$). The fitted decoders yield weight matrices $W_{\text{enc}}, W_{\text{ret}} \in \mathbb{R}^{k \times C}$, where $k = 100$ is the number of selected features and $C = 8$ is the number of stimulus classes. Each column of W represents the normal vector defining the separating hyperplane for a class.

Procrustes analysis and subspace alignment. We normalized decoder weights column-wise and extracted orthonormal bases via SVD: $W'_{\text{enc}}, W'_{\text{ret}} \in \mathbb{R}^{k \times r}$, where $r = \min(\text{rank}(W_{\text{enc}}), \text{rank}(W_{\text{ret}}))$. In our case, $r = 8$.

We then estimated the orthogonal matrix $R \in \mathbb{R}^{r \times r}$ that best aligns encoding to retrieval subspaces by solving the orthogonal Procrustes problem (Gower, 1975):

$$R^* = \arg \min_{R \in O(r)} \|W'_{\text{enc}}R - W'_{\text{ret}}\|_F.$$

The solution is given by $R^* = UV^\top$, where $U\Sigma V^\top$ is the singular value decomposition of $W'_{\text{ret}}{}^\top W'_{\text{enc}}$. We used this singular values σ_i of the $W'_{\text{ret}}{}^\top W'_{\text{enc}}$ to define the cosines of principal angles:

$$\cos \theta_i = \sigma_i, \quad \theta_i = \arccos(\sigma_i),$$

yielding an ordered sequence of angles with increasing degrees $0 \leq \theta_1 \leq \dots \leq \theta_r \leq 90^\circ$ that quantify alignment between subspaces.

For each ROI, we averaged principal angles across subjects to obtain an averaged canonical spectrum of rotation $\{\bar{\theta}_i\}$. For each model layer, we computed the corresponding spectrum of rotation from decoder weights. We then quantified brain–model similarity using the mean absolute deviation (MAD) between brain and model angles over the first n axes:

$$\Delta_{\text{ROI}, \ell} = \frac{1}{n} \sum_{i=1}^n |\bar{\theta}_i^{(\text{ROI})} - \theta_i^{(\ell)}|,$$

where ℓ indexes the model layer. Setting n gives the the average rotation degree of the top n shared axes between encoding and retrieval subspaces. We report the differences in the rotation degree $\Delta_{\text{ROI},\ell}$ between brain and model when $n = 1$ and $n = 3$, as in the top one and top three axes of space. Lower values indicate greater alignment between encoding–retrieval rotations in the brain ROI and the model layer ℓ .

Cross-decoding validation of rotation. To test whether retrieval decision geometry can be explained as a rigid rotation of the encoding geometry, we performed a cross-decoding analysis in the learned subspaces. For each ROI and subject, we formed orthonormal bases for the encoding and retrieval decision subspaces, $U_e \in \mathbb{R}^{p \times d}$ and $U_r \in \mathbb{R}^{p \times d}$. Axes were ranked by increasing principal angle (largest cosine first), and the top- k axes were retained to give $U_{e,k}$ and $U_{r,k}$ ($k = 1, \dots, d$). The orthogonal Procrustes rotation aligning these k -dimensional subspaces was

$$R_k = \arg \min_{R \in O(k)} \|U_{e,k}R - U_{r,k}\|_F.$$

All preprocessing (feature selection, standardization) and the linear multi-class classifier were fit *only on encoding* to avoid information leakage. Let $X_e, X_r \in \mathbb{R}^{n \times p}$ be encoding and retrieval data. We trained a classifier f on encoding-subspace coordinates

$$Z_e = X_e U_{e,k},$$

and then evaluated the *same* classifier on retrieval under two mappings:

$$\text{No rotation: } Z_r^{\text{naive}} = X_r U_{e,k} \quad \text{Rotated: } Z_r^{\text{rot}} = X_r U_{r,k} R_k^\top.$$

The outcome for each k was retrieval accuracy with and without rotation; we defined the gain as

$$\Delta_k = \text{acc}_{\text{rot},k} - \text{acc}_{\text{naive},k}.$$

For each ROI, Δ_k was computed per subject and then averaged across subjects to yield an ROI-level mean gain. To summarize “on average across ROIs,” we took the mean of the ROI-level Δ_k values at each k . Statistical inference tested whether the across-ROI gain was above 0 using one-sided one-sample t-tests (a priori directional hypothesis that rotation improved cross-decoding).

1134 A.8 SUPPLEMENTARY RESULTS

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

A.8.1 DECODERS TRAINED AND TESTED ON THE RETRIEVAL PHASE

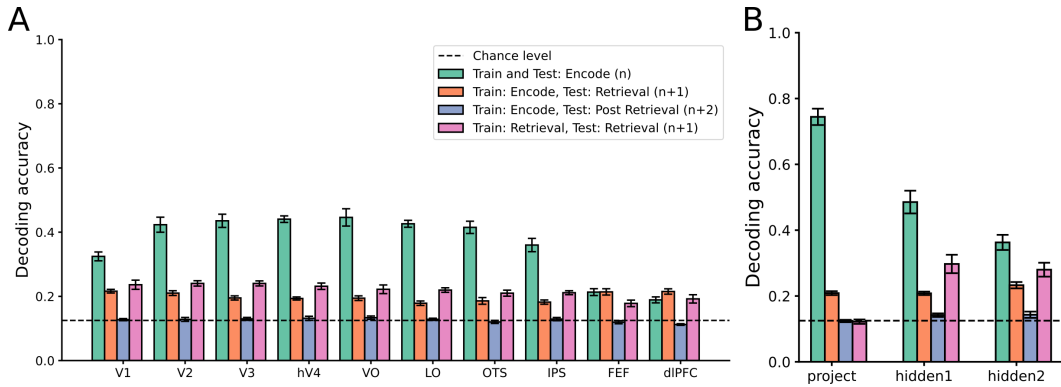


Figure A1: Decoding accuracies with decoders trained and tested on the retrieval phase. (A) Results of each human ROI. Error bars represent \pm s.e.m. across 8 subjects. (B) Results of each model's layer. Error bars represent \pm 1 s.e.m. across 12 models.

A.8.2 RSA AND DECODING ANALYSIS RESULTS WITH CONTROL REGION PRIMARY AUDITORY CORTEX (A1)

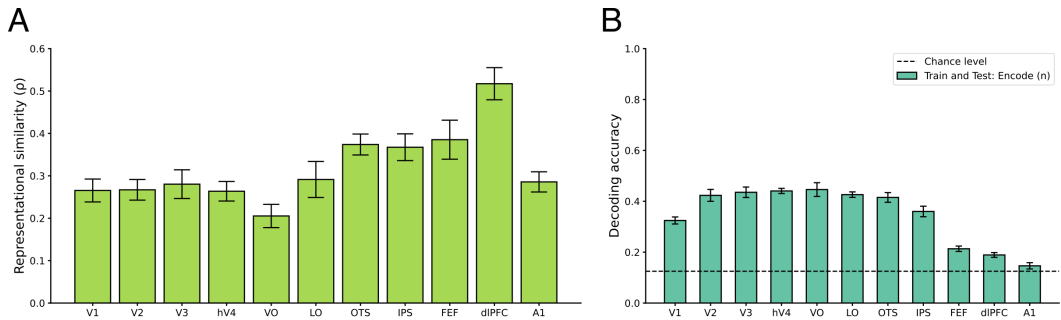


Figure A2: Results with control ROI A1. (A) Similarity scores (Spearman's ρ) of encoding and retrieval RDMs for each human ROI. Error bars represent \pm 1 s.e.m across 8 subjects. (B) Decoding accuracies of the decoder trained and tested on the encoding phase. Error bars represent \pm 1 s.e.m across 8 subjects.

A.8.3 CENTERED KERNEL ALIGNMENT (CKA) RESULTS WITH RBF KERNEL

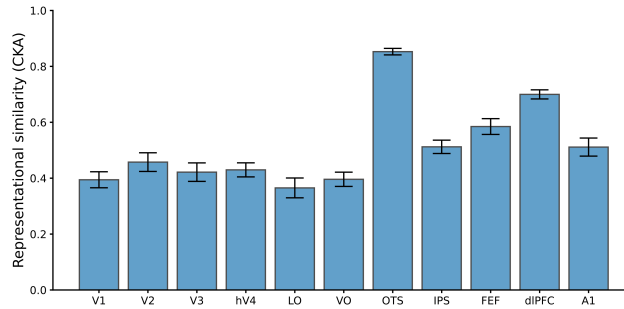


Figure A3: RBF kernel-based CKA results with control ROI A1. Error bars represent ± 1 s.e.m across subjects. The CKA RBF kernel parameter was set as 0.5. The parameters reflect the fraction of the median Euclidean distance used as σ .

A.8.4 2-BACK TASK RESULTS

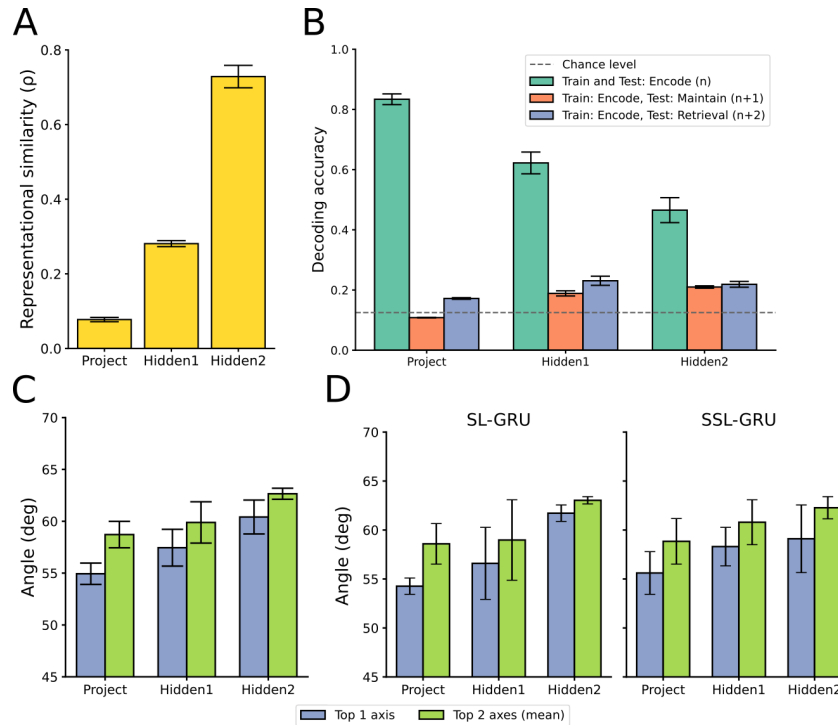


Figure A4: Analysis results with 2-back task models. Note that GRU was used as the recurrent module. (A) Similarity scores (Spearman's ρ) of encoding and retrieval RDMs for each model layer. Error bars represent ± 1 s.e.m across 4 models. (B) Decoding accuracies of each model's layer. Error bars represent ± 1 s.e.m across models. (C) Rotation angles of the top-1 axis and top-2 axes (mean) from encoding to retrieval in each layer of the models. Error bars represent ± 1 s.e.m across 4 models. (D) Group average rotation angles of the top-1 and top-2 axes (mean) from encoding to retrieval in each layer of the models. SL denotes supervised, SSL denotes self-supervised. Error bars represent ± 1 s.e.m across models with varying encoder architectures.

A.8.5 MODELS WITH FINE-TUNED ENCODERS RESULTS

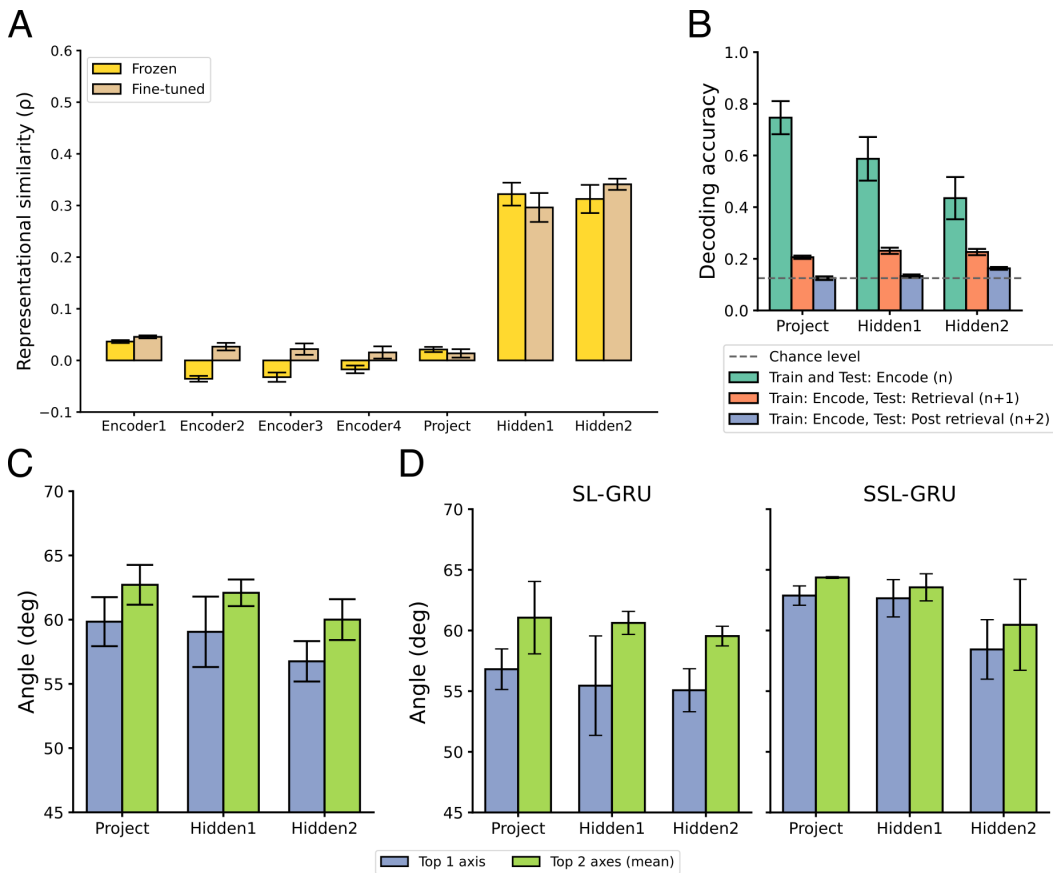


Figure A5: Analysis results of models with fine-tuned encoders. Note that GRU was used as the recurrent module for fine-tuned models. (A) Similarity scores (Spearman's ρ) of encoding and retrieval RDMs per model layer for both models with frozen encoders and with fine-tuned encoders. Error bars represent ± 1 s.e.m across models. (B) Decoding accuracies of each model's layer. Error bars represent ± 1 s.e.m across models. (C) Rotation angles of the top-1 axis and top-2 axes (mean) from encoding to retrieval in each layer of the models. Error bars represent ± 1 s.e.m across 4 models. (D) Group average rotation angles of the top-1 and top-2 axes (mean) from encoding to retrieval in each layer of the models. SL denotes supervised, SSL denotes self-supervised. Error bars represent ± 1 s.e.m across models with varying encoder architectures.

A.8.6 ADDITIONAL RSA RESULTS

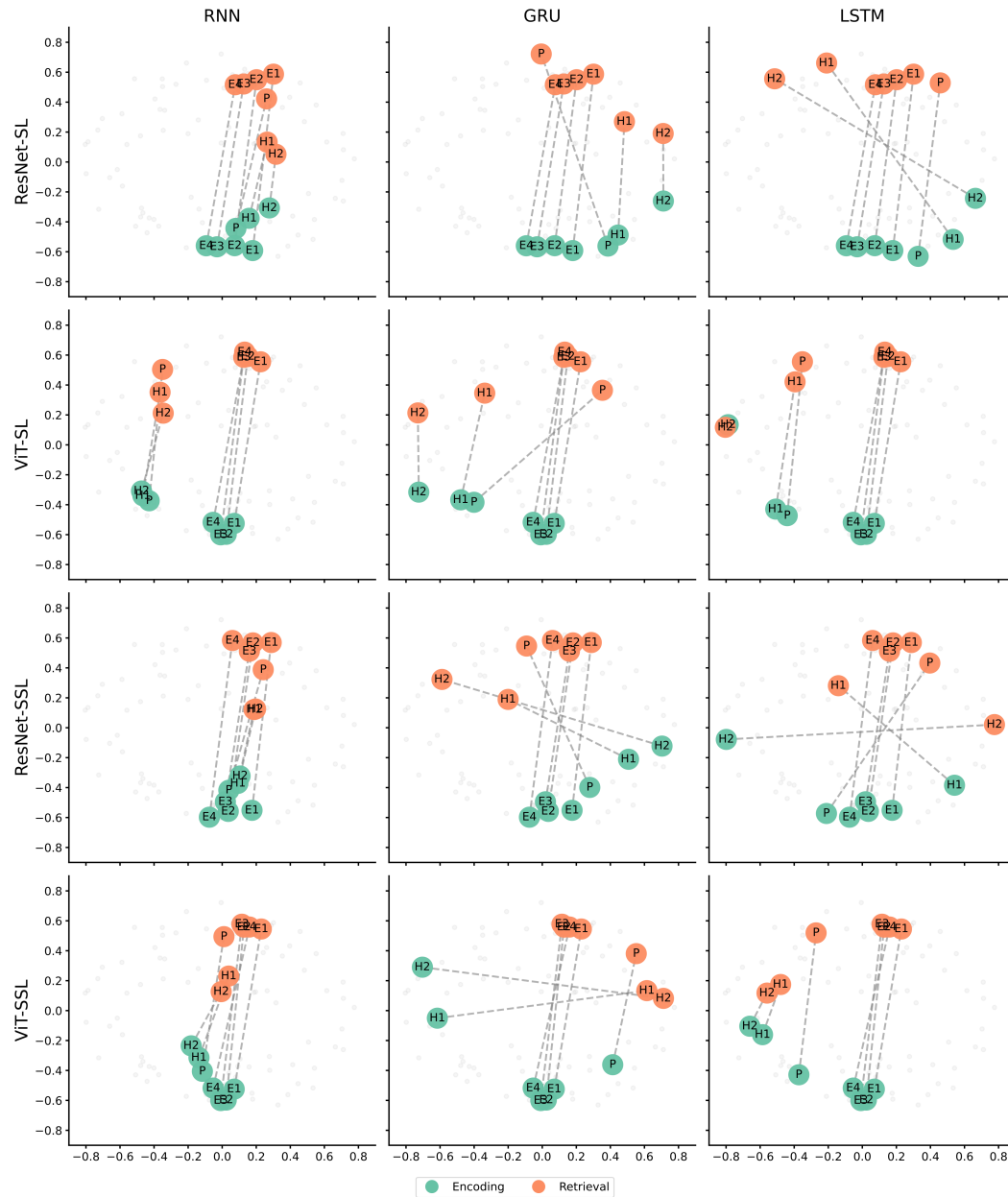


Figure A6: MDS visual-path of the dissimilarity matrix based on RDMS of encoding and retrieval representations for each model layer. E denotes frozen weight encoder layers, P denotes the project layer, and H denotes hidden layers. The number following the letter denotes the order of that layer. Note that across all the models, the similarity scores across encoding to retrieval phases were higher in 2 hidden layers compared to the project layer (FDR post-hoc paired t-tests, $p < .001$).

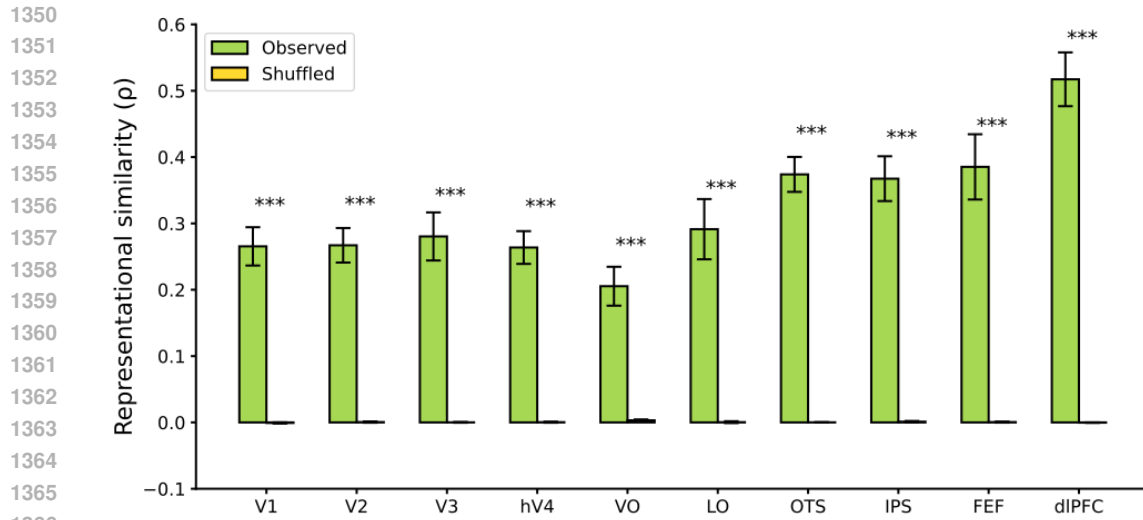


Figure A7: Shuffled RDM representational similarity results. Repeated-measures ANOVA showed significant effects of condition and ROI, and their interaction ($p < .001$). FDR post-hoc paired t-tests confirmed robust differences between observed and shuffled conditions across all ROIs. $p < .05$: *; $p < .01$: **; $p < .001$: ***.