

Bridging the Vision-Brain Gap with an Uncertainty-Aware Blur Prior

Supplementary Material

A. Experimental details

A.1. Datasets details

THINGS-EEG [21] is a large scale EEG dataset included 10 subjects with the Rapid Serial Visual Presentation (RSVP) paradigm [23, 30, 32]. The EEG data are collected using 64-channel EASYCAP equipment with the standard 10-10 system [49]. The training set includes 1654 concepts with each concept 10 images, and each image repeats 4 times (1654 concepts \times 10 images/concept \times 4 trials/image) per subject. The test set includes 200 concepts with each concept 1 image, and each image repeats 80 times (200 concepts \times 1 image/concept \times 80 trials/image) per subject.

For data preprocessing, we follow the method detailed in [60]. Raw EEG data filtered to [0.1, 100] Hz has 63 channels and a sample rate of 1000 Hz. EEG data is epoched into trials ranging from 0 to 1000 ms after stimuli onset with baseline correction using the prior 200 ms average. EEG data is down-sampled to 250 Hz and 17 channels are selected overlying occipital and parietal cortex related to visual¹. For the purpose of high Signal-to-Noise Ratio (SNR), EEG repetitions are averaged, resulting in total of 16540 training samples and 200 test samples per subject. Additionally, we store EEG data in float16 format to enable faster reading speeds and reduce storage requirements.

THINGS-MEG [26] dataset involves four participants and is characterized by 271 channels. The experimental design incorporates a relatively long stimulus duration of 500 ms, followed by a blank screen with a duration of 1000 ± 200 ms. It consists of 1854 concepts \times 12 images \times 1 repetitions in the training stage and 200 concepts \times 1 image \times 12 repetitions in the test stage.

We follow the settings described in [60]. During the data processing phase, 200 test concepts are discarded from the training set to construct the zero-shot task, mirroring the procedures in that study. Subsequently, the MEG data are epoched into trials covering the period from 0 to 1000 ms after the stimuli onset. For preprocessing, a band-pass filter within the range of [0.1, 100] Hz is utilized, and baseline correction is carried out after down-sampling the data to 200 Hz. Additionally, we average all MEG repetitions of one image to ensure the signal-to-noise ratio. Additionally, we store EEG data in float16 format to enable faster reading speeds and reduce storage requirements.

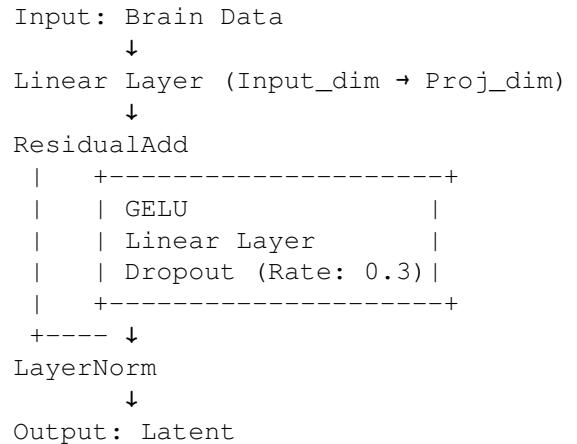
¹P7, P5, P3, P1, Pz, P2, P4, P6, P8, PO7, PO3, POz, PO4, PO8, O1, Oz, O2

A.2. Implementation details

Environment. Our method is implemented with Python 3.8.19, CUDA 12.0, and PyTorch 2.4.1. The required libraries are specified in the requirements.txt file provided in the repository. The experiments are performed on a machine equipped with an Intel Xeon Platinum 8352V CPU, four V100 GPUs, and 256 GB of RAM.

Training Configuration. We use a batch size of 1024 and train the model for 50 epochs. The learning rate is set to 1e-4 for intra-subject setting and 1e-5 for inter-subject setting. Gradient updates are performed using the AdamW[42] optimizer with weight decay set to 1e-4. Early stopping is employed to monitor training loss and validation performance, concluding the training process to mitigate overfitting when improvements stabilize. Notably, we use the softplus function instead of the exponential function to ensure that temperature parameter τ remains positive and continuous, as softplus offers a smoother and more stable transition, avoiding the numerical instability of the exponential function. For all above experiments, the hyperparameter r_0 is set to 0.25 and c is set to 10.

Architectures. We use EEGProject as the brain encoder, detailed as follows:



We provided the number of parameters and embedding dimension within different EEG encoders [34, 56, 60], in Tab. 21. Compared to other models, EEGProject achieves its performance through a simple yet effective architecture, while remaining lightweight with 5.154M parameters, especially in comparison to the vision branch.

We also provide the parameter counts for various CLIP vision branch models [29] to offer a comprehensive comparison across architectures in Tab. 22.

B. Results details

B.1. Retrieval Case Analysis

We present the top-5 retrieval results on THINGS-EEG dataset, including both good cases and bad cases, as shown in Fig. 9 and Fig. 10, respectively. Good cases demonstrate the model's capability to effectively align with the target stimuli and retrieve relevant results. An intriguing retrieval result is that the model not only retrieves items with similar materials but also demonstrates **associations with the orientation and quantity of objects**. These observations warrant further investigation in future studies.

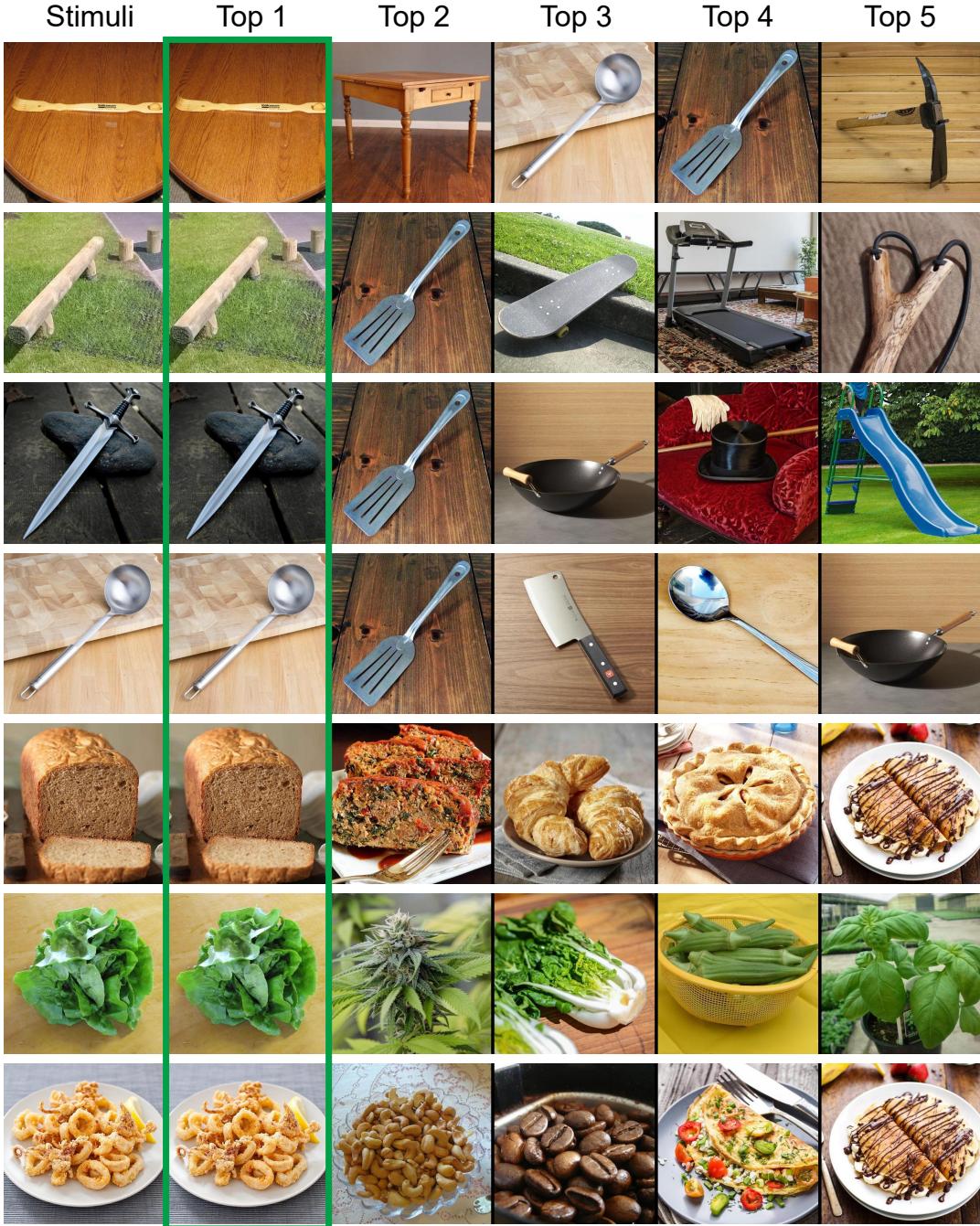


Figure 9. Good Cases: Top-5 Retrieval Results for Various Stimuli.

In contrast, bad cases reveal limitations in distinguishing fine-grained features or addressing semantic inconsistencies. It is challenging to distinguish highly similar stimuli due to the limited information contained in brain signals.

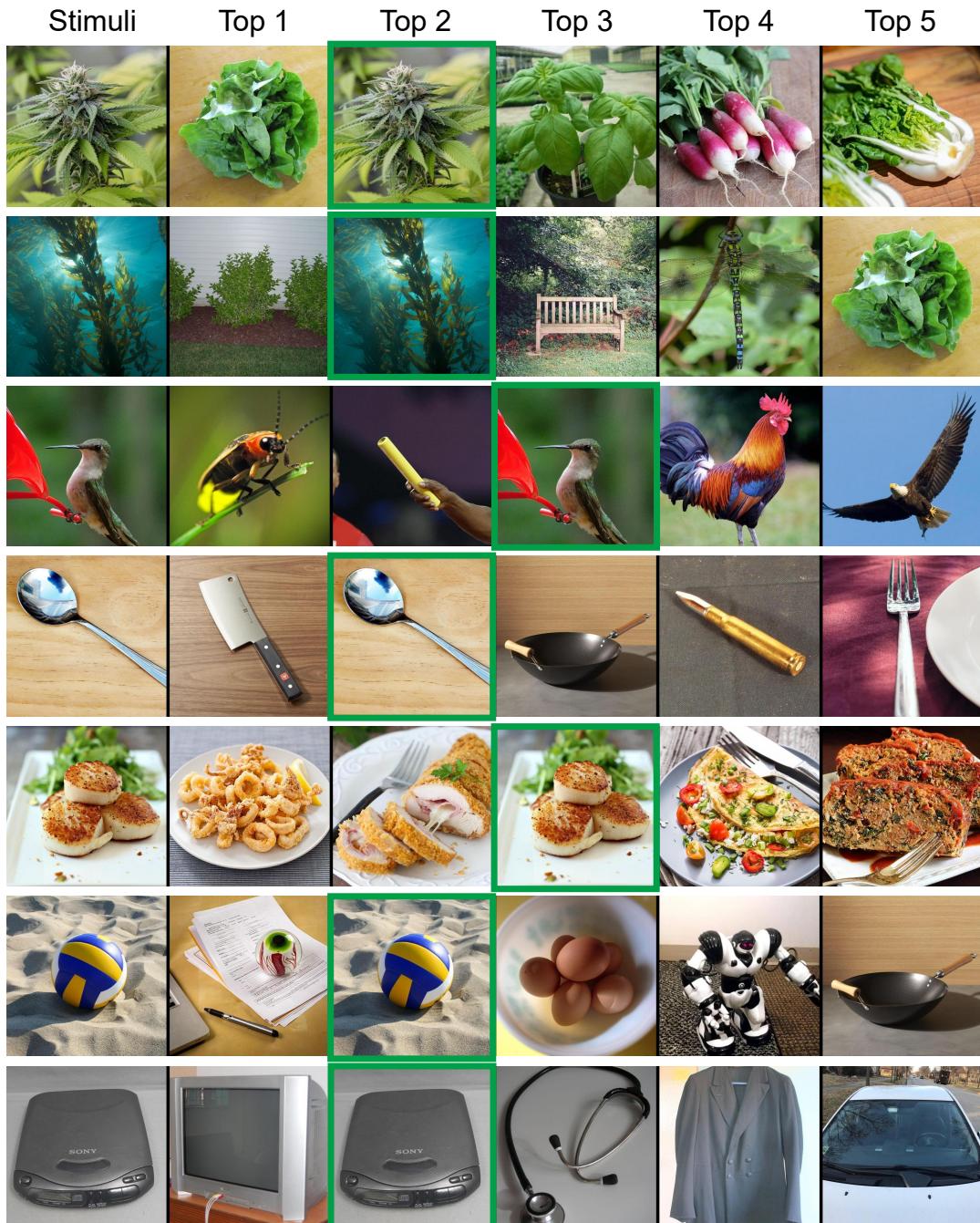
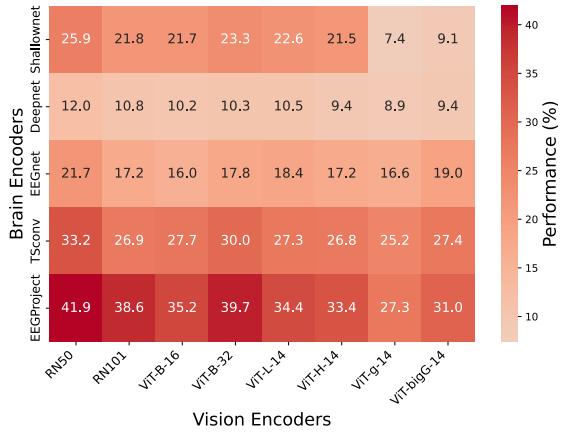
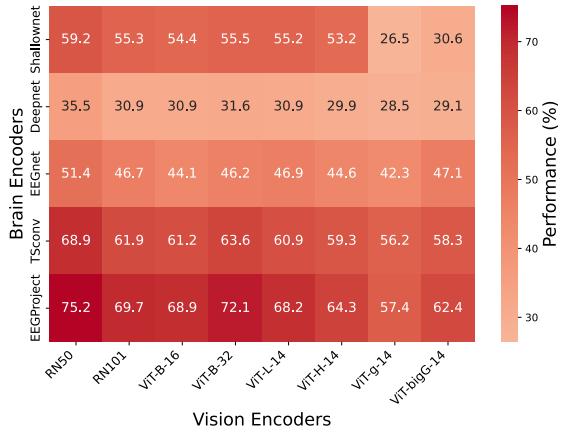


Figure 10. Bad Cases: Top-5 Retrieval Results for Various Stimuli.

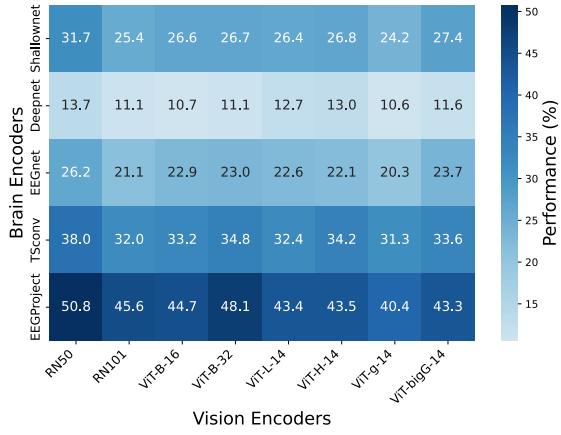
B.2. THINGS-EEG Results



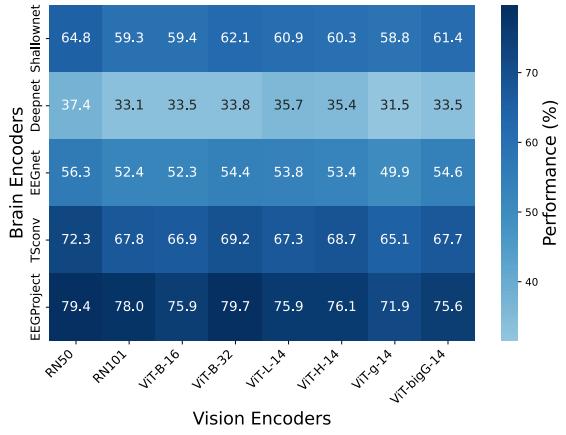
(a) Top-1 accuracy (%) of **Vanilla** on the THINGS-EEG dataset.



(b) Top-5 accuracy (%) of **Vanilla** on the THINGS-EEG dataset.



(c) Top-1 accuracy (%) of **UBP** on the THINGS-EEG dataset.



(d) Top-5 accuracy (%) of **UBP** on the THINGS-EEG dataset.



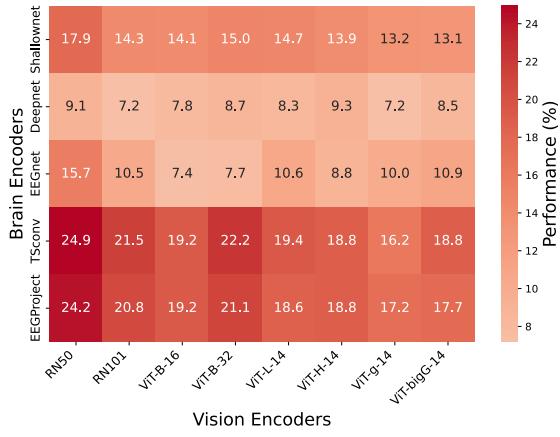
(e) Top-1 accuracy improvement (%) on the THINGS-EEG dataset.



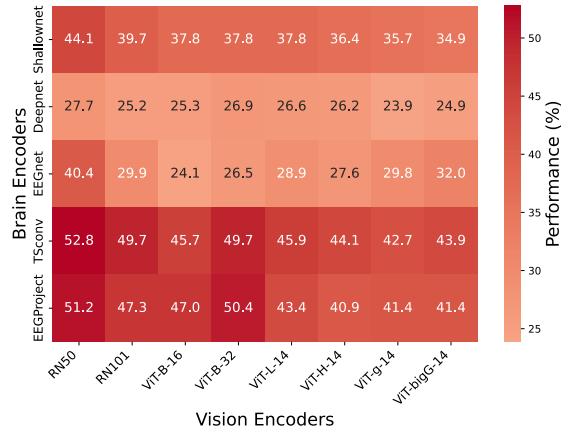
(f) Top-5 accuracy improvement (%) on the THINGS-EEG dataset.

Figure 11. Results on the THINGS-EEG dataset.

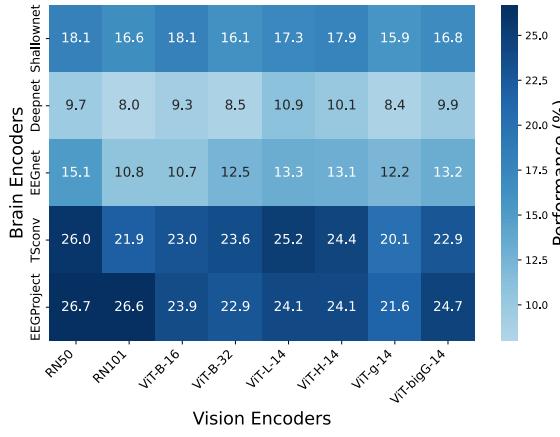
B.3. THINGS-MEG Results



(a) **Top-1** accuracy (%) of **Vanilla** on the THINGS-MEG dataset.



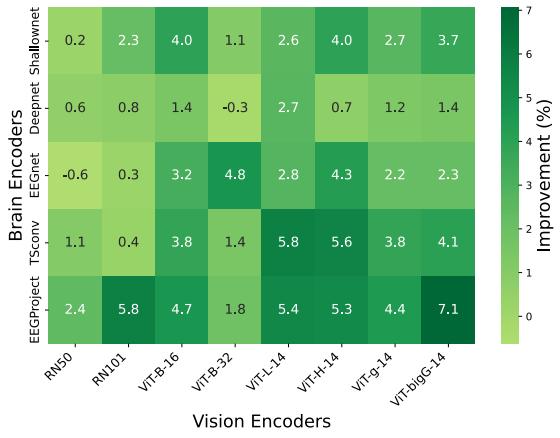
(b) **Top-5** accuracy (%) of **Vanilla** on the THINGS-MEG dataset.



(c) **Top-1** accuracy (%) of **UBP** on the THINGS-MEG dataset.



(d) **Top-5** accuracy (%) of **UBP** on the THINGS-MEG dataset.



(e) **Top-1** accuracy **improvement** (%) on the THINGS-MEG dataset.



(f) **Top-5** accuracy **improvement** (%) on the THINGS-MEG dataset.

Figure 12. Results on the THINGS-EEG dataset.

Table 21. Details of different EEG encoders with Emb dimension of 1024

Brain Encoder	Params
ShallowNet	2.56 M
DeepNet	2.76 M
EEGNet	2.34 M
TSConv	2.56 M
EEGProject	5.40 M

Table 22. Details of different Vision encoders

Vision Encoder	Params	Emb dim
RN50	38.32 M	1024
RN101	56.26 M	512
ViT-B-16	86.19 M	512
ViT-B-32	87.85 M	512
ViT-L-14	303.97 M	768
ViT-H-14	632.08 M	1024
ViT-g-14	1012.65 M	1024
ViT-bigG-14	1844.91 M	1280

B.4. EEG Feature Selection Analysis

Unless otherwise specified, EEG data spanning 1000 ms were selected, focusing on 17 visual-related (O+P) channels out of a total of 63, following previous work [60], where ablation studies are conducted on both channels and epochs. To verify the influence of the selection of EEG features, we conducted relevant ablation experiments. Tab. 23 and Fig. 13 show consistent improvements across various settings.

Table 23. Top-1 ACC (%) comparation with different channels.

Method	Occipital	Parietal	O+P (Our)	Others	All
Vanilla	21.0	33.7	42.0	10.2	35.6
UBP	26.9 (+5.9)	40.7 (+7.0)	50.9 (+8.9)	12.1 (+1.9)	42.2 (+6.6)

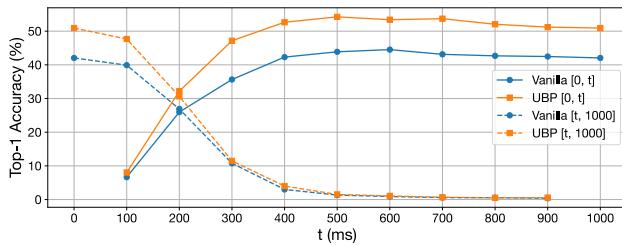


Figure 13. Top-1 ACC (%) comparation with different epoching.

B.5. More Evaluation Metrics

For a more comprehensive assessment of the method’s performance, we provided additional metric values to validate UBP’s effectiveness, specifically two new metrics: **mAP** and **Similarity Score** of paired samples. We present the performance averaged over 10 subjects. Tab. 24 demonstrates the results of the proposed method evaluated across different metrics. The calculation of metrics can be found in the repo.

Table 24. Retrieval Performance on THINGS-EEG.

Method	Top-1 Acc	Top-5 Acc	mAP*	Similarity*
Vanilla	42.0	75.2	56.6	0.160
UBP	50.9 (+8.9)	79.7 (+4.5)	63.8 (+7.2)	0.199 (+0.039)