

# MultiDiffNet: A Multi-Objective Diffusion Framework for Generalizable Brain Decoding

Meng-Chun Zhang<sup>1\*</sup>, Kateryna Shapovalenko<sup>2\*</sup>, Yucheng Shao<sup>2</sup>, Eddie Guo<sup>2</sup>, Parusha Pradhan<sup>1</sup>,

<sup>1</sup>University of Pittsburgh

<sup>2</sup>Carnegie Mellon University

MEZ141@pitt.edu, kshapova@alumni.cmu.edu, yshao3@andrew.cmu.edu, yuzhiguo@andrew.cmu.edu, pap203@pitt.edu

## Abstract

Neural decoding from electroencephalography (EEG) remains fundamentally limited by poor generalization to unseen subjects, driven by high inter-subject variability and the lack of large-scale datasets to model it effectively. Existing methods often rely on synthetic subject generation or simplistic data augmentation, but these strategies fail to scale or generalize reliably. We introduce *MultiDiffNet*, a diffusion-based framework that bypasses generative augmentation entirely by learning a compact latent space optimized for multiple objectives. We decode directly from this space and achieve state-of-the-art generalization across various neural decoding tasks using subject and session disjoint evaluation. We also curate and release a unified benchmark suite spanning four EEG decoding tasks of increasing complexity (SSVEP, Motor Imagery, P300, and Imagined Speech) and an evaluation protocol that addresses inconsistent split practices in prior EEG research. Finally, we develop a statistical reporting framework tailored for low-trial EEG settings. Our work provides a reproducible and open-source foundation for subject-agnostic EEG decoding in real-world BCI systems. Project code: <https://github.com/eddieguo-1128/DualDiff>

## Introduction

Electroencephalography (EEG) is a widely used modality in brain-computer interfaces (BCIs), supporting applications from assistive communication to cognitive monitoring. Deep learning has improved decoding across motor imagery, SSVEP, and speech tasks (Gu et al. 2025; Ahmadi and Mesin 2025; Lee and Lee 2022), yet generalizing to unseen subjects remains challenging due to high inter-subject variability and limited data (Huang et al. 2023; Barmpas et al. 2023).

Subject-specific models require extensive per-user calibration (Hartmann, Schirrmester, and Ball 2018; Luo and Cai 2024), while multi-subject models struggle to generalize (Rommel et al. 2022; Liu et al. 2022; Wu 2016). The alternative is to use two-stage pipelines that generate EEG via GANs or diffusion and then train decoders (Hartmann, Schirrmester, and Ball 2018; Torma and Szegletes 2025), but they suffer from low realism, artifact transfer, and inefficiencies.

\*These authors contributed equally.

We propose *MultiDiffNet*, a unified multi-objective diffusion framework that learns a shared latent space, eliminating the need for synthetic augmentation and enhancing generalization. To benchmark progress, we release a curated suite spanning SSVEP, Motor Imagery, P300, and Imagined Speech tasks, with standardized subject- and session-disjoint evaluation. We also develop a statistical reporting protocol tailored for low-trial EEG research, addressing a persistent gap in reproducibility.

## Related work

**EEG Decoding and Generalization** EEG decoding has evolved from handcrafted features to deep architectures, with EEGNet emerging as a widely adopted baseline due to its efficient depthwise-separable convolutions and lightweight design (Lawhern et al. 2018). Recent models explore transformers (Liao, Liu, and Wang 2025; Song et al. 2022a) and graph neural networks (Tang et al. 2024; Hu et al. 2023), but EEGNet remains favored for its robustness and simplicity. A key limitation is poor cross-subject generalization, with 20-40% accuracy drops despite strong within-subject performance (Huang et al. 2023; Barmpas et al. 2023). Attempts to address this require expensive calibration (Rommel et al. 2022; Liu et al. 2022; Wu 2016). Scalable BCIs require subject-agnostic models that generalize without per-user retraining.

**Diffusion Models for EEG** Denoising Diffusion Probabilistic Models (DDPMs) model data distributions via iterative denoising and outperform GANs in EEG synthesis by avoiding mode collapse (Tosato, Dalbagno, and Fumagalli 2023; Ho, Jain, and Abbeel 2020). Recent enhancements, such as reinforcement learning (An et al. 2024) and progressive distillation (Torma and Szegletes 2025), have further improved realism and sampling speed. Diff-E (Kim et al. 2023) extended diffusion to imagined-speech decoding via joint reconstruction and classification, but remained task-specific and did not address cross-subject generalization. Broader research suggests that combining generative and discriminative objectives yields stronger representations (Chow et al. 2024; Grathwohl et al. 2019), yet EEG models typically optimize only one. We explore this joint learning paradigm across diverse EEG tasks, aiming to learn generalizable representations that capture both signal structure and task-relevant information.

**Mixup Methods** Signal-level augmentation has evolved from basic jittering and filtering to temporal, spectral, and channel-wise mixup (Luo and Cai 2025; Liu, Lu, and Zheng 2025; Kim, Han, and Ko 2021; Pei et al. 2021; Zhang et al. 2017), but many variants introduce unrealistic artifacts that hinder generalization. This motivates our systematic evaluation of weighted and temporal input mixup across encoder layers, along with latent-space mixing

**Evaluation Strategies** Effective cross-subject EEG decoding requires both rigorous training strategies and standardized evaluation. Leave-one-subject-out (LOSO) validation remains common but is computationally intensive and impractical for real-time deployment (Del Pup et al. 2025; Chen et al. 2025; Zhao et al. 2024; Barmpas et al. 2023; Kunjan et al. 2021), while simpler subject splits often neglect session independence and true seen/unseen separation (Zhang et al. 2023). We address it in our work by introducing a standardized subject- and session-disjoint evaluation.

## Methodology

### MultiDiffNet architecture

*MultiDiffNet* is a modular architecture designed to jointly optimize classification, reconstruction, and contrastive structure learning from EEG signals. It consists of a Denoising Diffusion Probabilistic Model (DDPM), a discriminative encoder, a generative decoder, and a classifier (Figure 1).

Given a raw EEG signal  $x \in \mathbb{R}^{C \times T}$ , where  $C$  is the number of EEG channels and  $T$  is the number of timepoints, the model processes the input in two parallel paths. First, the DDPM denoises the signal via a learned reverse diffusion process, producing a refined version  $\hat{x} \in \mathbb{R}^{C \times T}$ . Simultaneously, the same input  $x$  is passed through an EEGNet-based encoder (See Section ) to extract a latent representation  $z \in \mathbb{R}^D$ , where  $D$  is the embedding dimension. The latent vector  $z$  is then used for two purposes: (1) it is passed to a lightweight decoder to reconstruct the denoised signal  $\hat{x}$ , resulting in a reconstruction  $x_{\text{dec}} \in \mathbb{R}^{C \times T}$ ; and (2) it is passed to a fully connected classification head to predict class logits  $\hat{y} \in \mathbb{R}^K$ , where  $K$  is the number of classes.

To further structure the latent space,  $z$  is locally normalized (Section ) and then projected to  $z_{\text{proj}} \in \mathbb{R}^{D'}$ , which is optimized with a supervised contrastive loss. All classification and reconstruction are performed directly from  $z$ , without relying on generated augmentations.

We performed an extensive ablation study across architectural variants, modifying the presence of DDPMs, encoder inputs, decoder pathways, classifier heads, and loss terms. The configuration described here reflects the best-performing combination.

### EEGNet-style encoder with attention pool

Given EEGNet’s demonstrated effectiveness across multiple EEG decoding tasks, we adapt its architecture as our discriminative encoder, hypothesizing that its proven feature extraction capabilities can produce powerful latent representations  $z$  for our multi-objective framework. Our encoder

extracts multi-scale features  $(dn_1, dn_2, dn_3)$  from different layers and applies attention pooling:

$$z = \text{AttentionPool}(dn_3) \in \mathbb{R}^D,$$

### Subject-wise latent normalization

To mitigate inter-subject variability, we apply subject-wise normalization on the encoder output  $z$ :

$$z_{\text{norm}} = \frac{z - \mu_s}{\sigma_s},$$

where  $\mu_s$  and  $\sigma_s$  denote the mean and standard deviation computed per subject  $s$  using a subset of training trials. During evaluation, we adopt a two-mode strategy: for seen subjects, normalization uses pre-computed statistics from their training data; for unseen subjects, statistics are estimated on-the-fly using their own calibration trials, simulating realistic deployment scenarios.

### Mixup strategies

Mixup strategies can improve robustness in low-trial EEG decoding. However, standard mixup techniques may not fully exploit the structure of neural time series. We therefore explore two complementary strategies: *Weighted Average Mixup* and a novel *Temporal Masked Mixup*. *Weighted Average Mixup* performs linear interpolation between the original EEG input  $x$ , the DDPM-denoised output  $\hat{x}$ , and the decoder reconstruction  $x_{\text{dec}}$ . We investigate multiple integration points in the model: **(0)** Input-level mixup, **(1-3)** Mixup after encoder layers 1, 2, or 3, respectively, **(4)** Mixup after the final attention pooling layer. To address the limitations of global interpolation, we propose *Temporal Masked Mixup*, which perturbs only localized segments of the input time series while preserving surrounding structure. See Algorithm 1 for pseudocode.

---

#### Algorithm 1: Temporal Masked Mixup

---

```

0: Initialize a binary mask  $M \in \{0, 1\}^{C \times T}$  with all zeros.
0: Flip each 0 in  $M$  to 1 with probability  $p = 0.01$ .
0: for each position in  $M$  with value 1 do
0:   Expand to a temporal window of random length (uniform
     between min and max size).
0: end for
0: Flip each 1 in  $M$  to  $-1$  with:
  • Fixed probability 0.5 (fixed ratio), or
  • Probability drawn from Beta(0.2, 0.2) each epoch
    (random ratio).
0: Apply the final mask:
  •  $0 \rightarrow x$  (original input)
  •  $1 \rightarrow \hat{x}$  (DDPM output)
  •  $-1 \rightarrow x_{\text{dec}}$  (decoder output)
=0

```

---

### Loss functions

*MultiDiffNet* is trained using a weighted sum of three objectives:

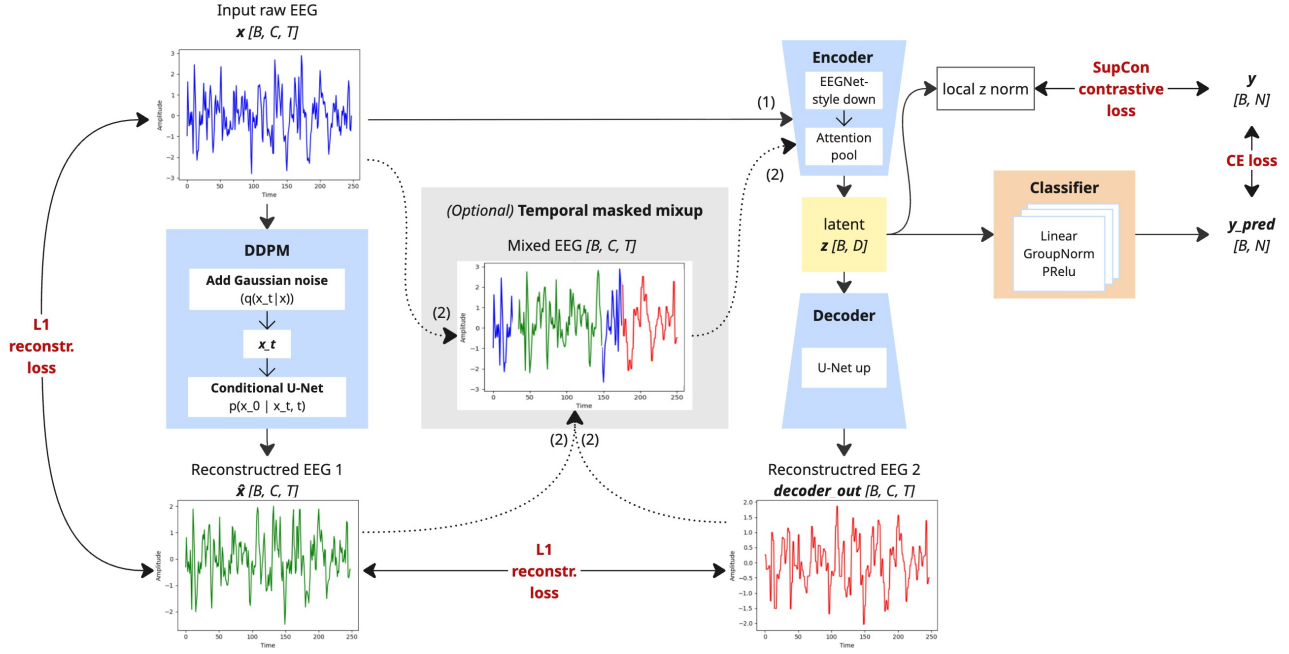


Figure 1: Overview of the *MultiDiffNet* that jointly optimizes a conditional DDPM, a contrastive encoder, and a generative decoder through a shared latent space  $z$ . The encoder produces discriminative features used for both classification and contrastive learning, while the decoder and DDPM reconstruct the input signal. An optional *temporal masked mixup* module stochastically blends the original, DDPM-denoised, and decoder-reconstructed EEG to improve representation quality.

$$\mathcal{L}_{\text{total}} = \underbrace{\alpha \mathcal{L}_{\text{CE/MSE}}(\hat{y}, y)}_{\text{classification}} + \underbrace{\beta \mathcal{L}_{\text{L1}}(x_{\text{dec}}, \hat{x})}_{\text{reconstruction}} + \underbrace{\gamma \mathcal{L}_{\text{SupCon}}(z_{\text{proj}}, y)}_{\text{contrastive}}$$

We fix  $\alpha = 1.0$  and progressively scale  $\beta$  and  $\gamma$  to stabilize training:

$$\beta = \min\left(1.0, \frac{\text{epoch}}{100}\right) \cdot 0.05, \quad \gamma = \min\left(1.0, \frac{\text{epoch}}{50}\right) \cdot 0.2$$

Details on loss formulation and weighting strategies are provided in the Appendix.

### Evaluation metrics

We evaluate model performance primarily using downstream classification accuracy, which quantifies the proportion of correctly classified EEG samples. Accuracy is defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

where  $TP$ ,  $TN$ ,  $FP$ , and  $FN$  denote true positives, true negatives, false positives, and false negatives, respectively. In addition, we report F1 score, precision, recall, and AUC for a more comprehensive evaluation; detailed formulas and results are provided in the Appendix.

### Trend-level statistical reporting framework

Conventional  $p$ -values often fail under the high-variance, low-trial, subject-disjoint conditions of EEG decoding. To address this, we introduce a robust trend-level statistical framework (detailed in the Appendix) that synthesizes effect sizes, cross-seed consistency, and Bayesian posterior probabilities. This allows us to detect systematic, reproducible gains even when classical significance tests return null results. Our approach represents a principled shift toward reproducible, evidence-based model evaluation in brain decoding.

While this framework enhances reproducibility, it is not meant to substitute conventional  $p$ -value testing. Instead, it addresses a well-documented limitation: in low-trial, high-variance EEG decoding, even systematic improvements may fail to reach arbitrary significance thresholds. Such small yet consistent gains—for instance, 2–3% accuracy in imagined speech or SSVEP—can substantially affect usability in BCI systems. By combining effect sizes, cross-seed consistency, and Bayesian evidence, the framework provides a principled way to surface these domain-relevant improvements, while remaining fully compatible with classical and non-parametric statistical tests.

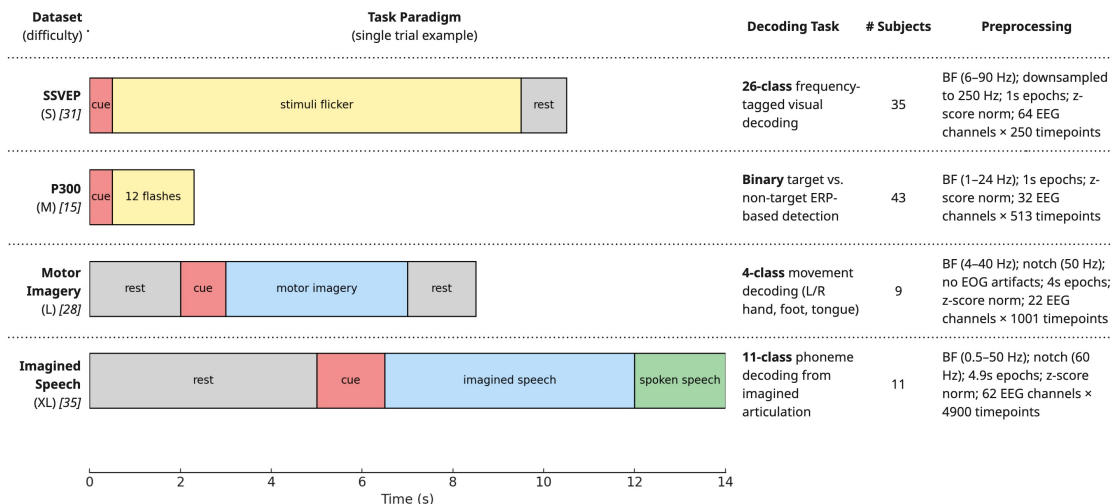


Figure 2: Overview of four EEG datasets ranked by task difficulty from easiest (top) to hardest (bottom). Task paradigms and preprocessing details are adapted from the original publications: SSVEP (Wang et al. 2017), P300 (Korcowski et al. 2019), Motor Imagery (Tangermann et al. 2012), and Imagined Speech (Zhao and Rudzicz 2015).

## Experiments and results

### Benchmark dataset suite

We curated four diverse EEG benchmarks (SSVEP, P300, Motor Imagery, and Imagined Speech), spanning increasing decoding difficulty. Each dataset is split into train, val, and two test sets: a seen-subject (intra-subject) split and an unseen-subject (cross-subject) split. This standardized protocol enables rigorous evaluation of both personalization and generalization, addressing the inconsistent and often unrealistic split practices prevalent in prior EEG research, where models are evaluated on mixed subject data or using computationally expensive LOSO.

### Baselines

We benchmarked our model against a diverse set of carefully selected baselines to ensure robust and fair comparisons. Our selection criteria were twofold: (i) prioritize architectures that are widely used for generalization to unseen subjects or sessions, and (ii) cover the main inductive biases found in EEG decoding, such as spatial filtering, temporal modeling, and attention mechanisms.

Specifically, we include: (1) **EEGNet** (Lawhern et al. 2018), a compact depthwise-separable CNN that is widely adopted for cross-subject generalization due to its strong accuracy–efficiency trade-off; (2) **ShallowFBCSPNet** (Schirrmester et al. 2017), which implements learnable filter-bank Common Spatial Patterns (CSP) to extract frequency–spatial features; (3) **TIDNet** (Kostas and Rudzicz 2020), which introduces dilated convolutions and residual connections to improve robustness under subject shift; (4) **EEGConformer** (Song et al. 2022b), which combines a convolutional front-end with self-attention to model both local spatial structure and global temporal context; and (5) **EEGTCNet** (Ingolfsson et al. 2020), a temporal convolu-

tional network tailored for EEG that emphasizes causal and dilated temporal modeling, offering complementary inductive bias to purely spatial–spectral models.

All models are evaluated using identical input windows of shape  $(C, T)$ , and trained with a unified global training schedule to ensure comparability. Public implementations and recommended hyperparameters are used where available, with no method-specific tuning.

### Generalization performance

*MultiDiffNet* helps with generalization. Unlike raw EEG representations, where class boundaries blur due to subject-specific noise, our learned latent space forms clearly separable, label-aligned clusters (Figure 3). This structured representation enables robust decoding across subjects. As shown in Table 1, *MultiDiffNet* consistently reduces the seen–unseen accuracy gap across all tasks. In SSVEP, it lifts cross-subject accuracy from 81.08% (EEGNet baseline) to 84.72%, further boosted to 85.25% with Temporal Masked Mixup. For comparison, other representative architectures such as ShallowFBCSPNet (58.87%), EEGConformer (51.92%), TIDNet (25.96%), and EEGTCNet (49.57%) fall well behind, highlighting the robustness of our latent-space design.

Even in the low-SNR regime of Imagined Speech, *MultiDiffNet* improves cross-subject accuracy from 10.61% (EEGNet) to 12.12%, while simultaneously achieving a much larger gain on seen-subject accuracy (11.26% → 17.57%). Other baselines such as ShallowFBCSPNet (10.48/13.78%), EEGConformer (9.21/10.62%), TIDNet (10.35/9.10%), and EEGTCNet (10.10/12.64%) hover close to chance level on both splits, further highlighting the robustness of our approach. For such a challenging task, even modest absolute gains are meaningful, as they can indicate more reliable signal extraction under extreme noise

| Task        | Model                       | Subj. | Classes | Seen Acc. (%)          | Unseen Acc. (%)        |
|-------------|-----------------------------|-------|---------|------------------------|------------------------|
| SSVEP       | ShallowFBCSPNet             | 35    | 26      | 69.58 ± 1.30*          | 58.87 ± 9.37*          |
|             | EEGConformer                | 35    | 26      | 66.98 ± 2.83           | 51.92 ± 9.06           |
|             | TIDNet                      | 35    | 26      | 28.01 ± 4.12           | 25.96 ± 5.29           |
|             | EEGTCNet                    | 35    | 26      | 58.31 ± 4.02           | 49.57 ± 9.14           |
|             | EEGNet                      | 35    | 26      | 89.16 ± 0.57***        | 81.08 ± 9.16**         |
|             | <b>MultiDiffNet</b>         | 35    | 26      | 85.08 ± 1.53**         | <b>84.72 ± 6.03***</b> |
|             | <b>MultiDiffNet + Mixup</b> | 35    | 26      | 86.79 ± 1.75***        | <b>85.25 ± 6.94***</b> |
| P300        | ShallowFBCSPNet             | 43    | 2       | 87.72 ± 0.33           | 86.20 ± 1.45           |
|             | EEGConformer                | 43    | 2       | 88.54 ± 0.54**         | 86.30 ± 1.73           |
|             | TIDNet                      | 43    | 2       | 88.24 ± 0.31*          | 85.63 ± 0.58**         |
|             | EEGTCNet                    | 43    | 2       | 88.69 ± 0.59***        | 87.02 ± 1.62***        |
|             | EEGNet                      | 43    | 2       | 88.79 ± 0.67***        | 87.24 ± 2.01***        |
|             | <b>MultiDiffNet</b>         | 43    | 2       | 85.35 ± 1.12           | 79.47 ± 0.54*          |
|             | <b>MultiDiffNet + Mixup</b> | 43    | 2       | 85.61 ± 0.52           | 79.56 ± 4.43           |
| MI          | ShallowFBCSPNet             | 9     | 4       | 64.34 ± 3.61***        | 36.46 ± 6.60           |
|             | EEGConformer                | 9     | 4       | 59.57 ± 5.60**         | 36.49 ± 7.72           |
|             | TIDNet                      | 9     | 4       | 44.27 ± 2.60           | 34.42 ± 3.60           |
|             | EEGTCNet                    | 9     | 4       | 58.85 ± 4.54           | 32.99 ± 6.94           |
|             | EEGNet                      | 9     | 4       | 67.01 ± 5.38***        | 46.18 ± 7.20***        |
|             | <b>MultiDiffNet</b>         | 9     | 4       | 55.85 ± 2.80           | 39.24 ± 8.00***        |
|             | <b>MultiDiffNet + Mixup</b> | 9     | 4       | 57.69 ± 3.27*          | 36.78 ± 5.23           |
| Img. Speech | ShallowFBCSPNet             | 14    | 11      | 13.78 ± 1.55**         | 10.48 ± 0.64           |
|             | EEGConformer                | 14    | 11      | 10.62 ± 0.82           | 9.21 ± 3.00            |
|             | TIDNet                      | 14    | 11      | 9.10 ± 0.54            | 10.35 ± 0.18           |
|             | EEGTCNet                    | 14    | 11      | 12.64 ± 1.58           | 10.10 ± 0.64           |
|             | EEGNet                      | 14    | 11      | 11.26 ± 2.01*          | 10.61 ± 0.93*          |
|             | <b>MultiDiffNet</b>         | 14    | 11      | <b>15.55 ± 0.62***</b> | <b>11.62 ± 1.29***</b> |
|             | <b>MultiDiffNet + Mixup</b> | 14    | 11      | <b>17.57 ± 1.16***</b> | <b>12.12 ± 0.38***</b> |

Table 1: Final results across tasks and models. Accuracy is reported for both seen-subject (intra-subject) and unseen-subject (cross-subject) test splits. Tasks are ranked by task difficulty. Stars denote win percentage: \*\*\*  $\geq 80\%$ , \*\*  $\geq 60\%$ , \*  $\geq 40\%$ . Detailed results are in the Appendix.

conditions. On Motor Imagery, *MultiDiffNet* also surpasses most baselines on unseen accuracy, e.g., outperforming TIDNet (34.42%) and EEGTCNet (32.99%), while maintaining competitive seen accuracy (57.69% vs. 44.27% for TIDNet and 58.85% for EEGTCNet). Although it remains slightly below EEGNet (46.18/67.01%), this is likely due to ceiling effects and dataset scale.

### Ablation studies

To understand what drives generalization in *MultiDiffNet*, we ran extensive ablation experiments, over 100 controlled configs. All results are reported for both seen- and unseen-subject accuracy, with statistical evidence matrices and trend-level effect sizes in the Appendix.

**Decoder input.** Feeding only  $z$  to the decoder often matches or exceeds more complex fusion variants. For example, SSVEP unseen accuracy reaches 84.72% with  $z$  alone, further boosted to 85.25% with mixup, while more elaborate fusions ( $z + x$ ,  $x_{\text{hat}} + \text{skips}$ ) show no consistent gains. These findings validate our architectural decision to decode primarily from  $z$ . For completeness, the best accuracies achieved in this ablation are 85.86/84.72 on SSVEP,

85.88/81.41 on P300, 56.89/40.36 on Motor Imagery, and 18.58/12.88 on Imagined Speech (seen/unseen).

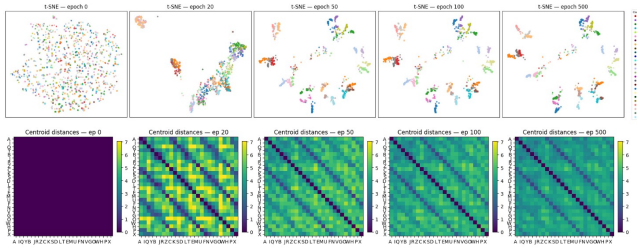
**Classifier head.** A lightweight FC head on  $z$  delivers state-of-the-art generalization with minimal complexity. It rivals or outperforms EEGNet classifiers trained on  $x$ , especially in low-SNR tasks. This supports our choice to use FC as the default classification head. For completeness, the best accuracies achieved in this ablation are 85.08/84.72 on SSVEP, 85.35/84.12 on P300, 55.85/39.24 on Motor Imagery, and 17.95/11.61 on Imagined Speech (seen/unseen).

**Encoder and decoder.** Using raw  $x$  as encoder input consistently outperforms  $\hat{x}$ , showing that denoising is useful for regularization. Interestingly, removing the decoder entirely sometimes improves generalization, suggesting that reconstruction may introduce noise if overemphasized. For completeness, the best accuracies in this ablation are 90.95/85.58 on SSVEP, 85.71/80.93 on P300, 55.85/40.16 on Motor Imagery, and 19.22/13.76 on Imagined Speech (seen/unseen).

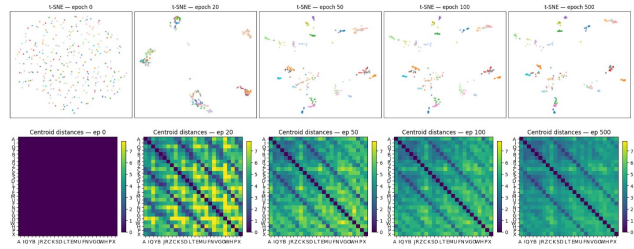
**Loss combinations.** Combining CE with mild MSE or contrastive losses improves stability, particularly when auxiliary weights are gently annealed. The best results use  $\beta = 0.05$ ,  $\gamma = 0.2$ —balancing reconstruction as a regularizer

(A) Latent space evolution across epochs via t-SNE and class centroid distances, showing how the model learns to better separate different classes over time (SSVEP task).

(A.1) Test seen subjects

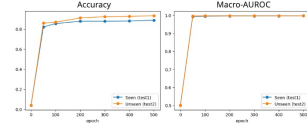


(A.2) Test unseen subjects

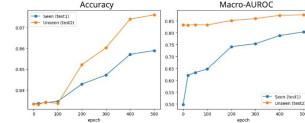


(B) Classification performance from frozen latent space across epochs via accuracy and macro-AUROC on seen/unseen test sets, showing improved generalization over time.

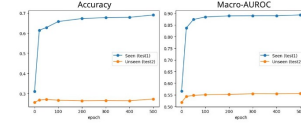
(B.1) SSVEP task



(B.2) P300 task



(B.3) Motor imagery task



(B.4) Imagined speech task

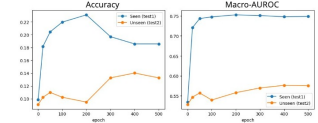


Figure 3: (A) Visualization of latent space across training epochs. (B) Downstream classification performance from frozen latent representations.

without overpowering the classification objective. For completeness, the best accuracies in this ablation are 86.40/85.58 on SSVEP, 85.69/80.18 on P300, 59.67/41.44 on Motor Imagery, and 19.60/13.51 on Imagined Speech (seen/unseen).

**Mixup strategies.** Mixup effects are task-specific. For SSVEP, *Temporal Masked Mixup* outperforms all variants. Motor Imagery benefits from *Weighted Average Mixup*, while P300 and Imagined Speech show limited sensitivity, highlighting that mixup is most impactful in high-SNR regimes. For completeness, the best accuracies in this ablation are 87.84/85.26 on SSVEP, 85.78/79.56 on P300, 63.44/38.83 on Motor Imagery, and 19.47/12.12 on Imagined Speech (seen/unseen).

## Conclusions and future work

We presented *MultiDiffNet*, a diffusion-based neural decoder that learns a compact, multi-objective latent space for EEG decoding without synthetic augmentation. Through unified benchmarks and rigorous cross-subject evaluation, we showed that *MultiDiffNet* achieves strong generalization across diverse BCI paradigms, particularly in challenging low-signal settings such as SSVEP and Imagined Speech. Our statistical analysis framework further addresses reproducibility challenges in low-trial EEG research. Future work will explore scaling *MultiDiffNet* to larger and more diverse EEG datasets and extending the architecture to other neural modalities.

For completeness, we note that our trend-level statistical framework is intended only as a complementary tool for low-trial EEG research; detailed rationale is provided in Section 3.7, with Bayesian and non-parametric validations reported in the Appendix.

## Acknowledgments

The authors would like to thank Professor Bhiksha Raj of Carnegie Mellon University for his guidance and support.

## References

- Ahmadi, H.; and Mesin, L. 2025. Universal semantic feature extraction from EEG signals: a task-independent framework. *Journal of Neural Engineering*, 22(3): 036003.
- An, Y.; Tong, Y.; Wang, W.; and Su, S. W. 2024. Enhancing EEG Signal Generation through a Hybrid Approach Integrating Reinforcement Learning and Diffusion Models. *arXiv preprint arXiv:2410.00013*.
- Barmpas, K.; Panagakis, Y.; Bakas, S.; Adamos, D. A.; Laskaris, N.; and Zafeiriou, S. 2023. Improving generalization of CNN-based motor-imagery EEG decoders via dynamic convolutions. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31: 1997–2005.
- Chen, Z.; Wang, P. T.; Ibrahim, M.; Baveja, S.; Mu, R.; Do, A. H.; and Nenadic, Z. 2025. Leveraging Transfer Learning and User-Specific Updates for Rapid Training of BCI Decoders. *arXiv preprint arXiv:2506.14120*.
- Chow, W.; Li, J.; Yu, Q.; Pan, K.; Fei, H.; Ge, Z.; Yang, S.; Tang, S.; Zhang, H.; and Sun, Q. 2024. Unified generative and discriminative training for multi-modal large language models. *Advances in Neural Information Processing Systems*, 37: 23155–23190.
- Del Pup, F.; Zanolà, A.; Tshimanga, L. F.; Bertoldo, A.; Finos, L.; and Atzori, M. 2025. The role of data partitioning on the performance of EEG-based deep learning models in supervised cross-subject analysis: a preliminary study. *Computers in Biology and Medicine*, 196: 110608.
- Grathwohl, W.; Wang, K.-C.; Jacobsen, J.-H.; Duvenaud, D.; Norouzi, M.; and Swersky, K. 2019. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*.
- Gu, H.; Chen, T.; Ma, X.; Zhang, M.; Sun, Y.; and Zhao, J. 2025. CLTNet: A Hybrid Deep Learning Model for Motor Imagery Classification. *Brain Sciences*, 15(2): 124.
- Hartmann, K. G.; Schirrmester, R. T.; and Ball, T. 2018. EEG-GAN: Generative adversarial networks for electroencephalographic (EEG) brain signals. *arXiv preprint arXiv:1806.01875*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Hu, W.; Jiang, G.; Han, J.; Li, X.; and Xie, P. 2023. Regional-asymmetric adaptive graph convolutional neural network for diagnosis of autism in children with resting-state EEG. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 32: 200–211.
- Huang, G.; Zhao, Z.; Zhang, S.; Hu, Z.; Fan, J.; Fu, M.; Chen, J.; Xiao, Y.; Wang, J.; and Dan, G. 2023. Discrepancy between inter- and intra-subject variability in EEG-based motor imagery brain-computer interface: Evidence from multiple perspectives. *Frontiers in neuroscience*, 17: 1122661.
- Ingolfsson, T. M.; Hersche, M.; Wang, X.; Kobayashi, N.; Cavigelli, L.; and Benini, L. 2020. EEG-TCNet: An accurate temporal convolutional network for embedded motor-imagery brain-machine interfaces. *arXiv preprint arXiv:2006.00622*. EEGTCNet.
- Kim, G.; Han, D. K.; and Ko, H. 2021. Specmix: A mixed sample data augmentation method for training withtime-frequency domain features. *arXiv preprint arXiv:2108.03020*.
- Kim, S.; Lee, Y.-E.; Lee, S.-H.; and Lee, S.-W. 2023. Diff-E: Diffusion-based learning for decoding imagined speech EEG. *arXiv preprint arXiv:2307.14389*.
- Korczowski, L.; Cederhout, M.; Andreev, A.; Cattan, G.; Rodrigues, P. L. C.; Gautheret, V.; and Congedo, M. 2019. *Brain Invaders calibration-less P300-based BCI with modulation of flash duration Dataset (bi2015a)*. Ph.D. thesis, GIPSA-lab.
- Kostas, D.; and Rudzicz, F. 2020. Thinker invariance: Enabling deep neural networks for BCI across more people. *Journal of Neural Engineering*, 17(5): 056008. TIDNet.
- Kunjan, S.; Grummett, T. S.; Pope, K. J.; Powers, D. M.; Fitzgibbon, S. P.; Bastiampillai, T.; Battersby, M.; and Lewis, T. W. 2021. The necessity of leave one subject out (LOSO) cross validation for EEG disease diagnosis. In *International conference on brain informatics*, 558–567. Springer.
- Lawhern, V. J.; Solon, A. J.; Waytowich, N. R.; Gordon, S. M.; Hung, C. P.; and Lance, B. J. 2018. EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces. *Journal of neural engineering*, 15(5): 056013.
- Lee, Y.-E.; and Lee, S.-H. 2022. EEG-transformer: Self-attention from transformer architecture for decoding EEG of imagined speech. In *2022 10th International winter conference on brain-computer interface (BCI)*, 1–4. IEEE.
- Liao, W.; Liu, H.; and Wang, W. 2025. Advancing BCI with a transformer-based model for motor imagery classification. *Scientific Reports*, 15(1): 23380.
- Liu, S.; Zhang, J.; Wang, A.; Wu, H.; Zhao, Q.; and Long, J. 2022. Subject adaptation convolutional neural network for EEG-based motor imagery classification. *Journal of Neural Engineering*, 19(6): 066003.
- Liu, X.-H.; Lu, B.-L.; and Zheng, W.-L. 2025. mix-EEG: Enhancing EEG Federated Learning for Cross-subject EEG Classification with Tailored mixup. *arXiv preprint arXiv:2504.07987*.
- Luo, T.-j.; and Cai, Z. 2024. Diffusion models-based motor imagery EEG sample augmentation via mixup strategy. *Expert Systems with Applications*, 235: 125585.
- Luo, T.-j.; and Cai, Z. 2025. Diffusion models-based motor imagery EEG sample augmentation via mixup strategy. *Expert Systems with Applications*, 262: 125585.
- Pei, Y.; Luo, Z.; Yan, Y.; Yan, H.; Jiang, J.; Li, W.; Xie, L.; and Yin, E. 2021. Data augmentation: Using channel-level recombination to improve classification performance for motor imagery EEG. *Frontiers in Human Neuroscience*, 15: 645952.
- Rommel, C.; Paillard, J.; Moreau, T.; and Gramfort, A. 2022. Data augmentation for learning predictive models on EEG: a systematic comparison. *Journal of Neural Engineering*, 19(6): 066020.

Schirrneister, R. T.; Springenberg, J. T.; Fiederer, L. D.; Glasstetter, M.; Eggenberger, K.; Tangermann, M.; Hutter, F.; and Ball, T. 2017. Deep learning with convolutional neural networks for EEG decoding and visualization. *Human Brain Mapping*. ShallowFBCSPNet.

Song, Y.; Zheng, Q.; Liu, B.; and Gao, X. 2022a. EEG conformer: Convolutional transformer for EEG decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31: 710–719.

Song, Y.; Zheng, Q.; Liu, B.; and Gao, X. 2022b. EEG Conformer: Convolutional transformer for EEG decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31: 710–719. EEGConformer.

Tang, X.; Zhang, J.; Qi, Y.; Liu, K.; Li, R.; and Wang, H. 2024. A spatial filter temporal graph convolutional network for decoding motor imagery EEG signals. *Expert Systems with Applications*, 238: 121915.

Tangermann, M.; Müller, K.-R.; Aertsen, A.; Birbaumer, N.; Braun, C.; Brunner, C.; Leeb, R.; Mehring, C.; Miller, K. J.; Müller-Putz, G. R.; et al. 2012. Review of the BCI competition IV. *Frontiers in neuroscience*, 6: 55.

Torma, S.; and Szegletes, L. 2025. Generative modeling and augmentation of EEG signals using improved diffusion probabilistic models. *Journal of Neural Engineering*, 22(1): 016001.

Tosato, G.; Dalbagno, C. M.; and Fumagalli, F. 2023. EEG Synthetic Data Generation Using Probabilistic Diffusion Models. arXiv:2303.06068.

Wang, Y.; Chen, X.; Gao, X.; and Gao, S. 2017. A Benchmark Dataset for SSVEP-Based Brain–Computer Interfaces. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 25(10): 1746–1752.

Wu, D. 2016. Online and offline domain adaptation for reducing BCI calibration effort. *IEEE Transactions on Human-Machine Systems*, 47(4): 550–563.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.

Zhang, H.; Ji, H.; Yu, J.; Li, J.; Jin, L.; Liu, L.; Bai, Z.; and Ye, C. 2023. Subject-independent EEG classification based on a hybrid neural network. *Frontiers in Neuroscience*, 17: 1124089.

Zhao, S.; and Rudzicz, F. 2015. Classifying phonological categories in imagined and articulated speech. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 992–996. IEEE.

Zhao, W.; Jiang, X.; Zhang, B.; Xiao, S.; and Weng, S. 2024. CTNet: a convolutional transformer network for EEG-based motor imagery classification. *Scientific reports*, 14(1): 20237.