

NAKUL-Med: Spectral-Graph State Space Models with Dynamics Kernels for Medical Signals

Badri N. Patro
Microsoft

badripatro@microsoft.com

Vijay S. Agneeswaran
Microsoft

vagneeswaran@microsoft.com

1. Introduction

This document extends the main paper with technical details and additional experiments. Section 2 analyzes per-dataset performance. Section 3 describes experimental setup and training. Section 4 presents cross-dataset comparisons and efficiency benchmarks. Section 5 explores learned patterns through visualizations. Section 6 details all datasets and preprocessing steps.

2. Per-Dataset Results

We examine performance on five tasks: BCI-IV-2a (motor imagery with 144 trials/subject), FACED (9-class emotion with high variability), SeizeIT1 (multimodal EEG-fMRI with 1:20 class imbalance), OpenNeuro (fMRI task decoding at 2s resolution), and BUSI (adapting 1D models to 2D ultrasound). These tests cover different modalities (electro-physiology, hemodynamics, acoustics), temporal scales (ms to seconds), and spatial structures.

2.1. SeizeIT1: EEG-fMRI Epilepsy Detection

Seizure detection must avoid false negatives (safety risk) and false positives (reduces usability). SeizeIT1 uses simultaneous EEG-fMRI—EEG captures millisecond electrical dynamics, fMRI shows spatial hemodynamic changes. Table 1 reports accuracy, sensitivity, specificity, and AUROC.

Table 1 shows seizure detection performance with class-wise metrics.

Table 1. SeizeIT1 seizure detection results.

Method	Acc (%)	Sens (%)	Spec (%)	AUROC
3D-ResNet (fMRI only) [6]	82.3 ± 1.8	76.4	84.7	0.847
EEGNet (EEG only) [11]	84.2 ± 1.5	79.3	86.1	0.871
Late Fusion (EEG+fMRI) [16]	87.6 ± 1.3	83.7	89.4	0.901
EEG-Conformer [13]	88.3 ± 1.1	85.1	89.8	0.912
NAKUL-Med	91.4 ± 0.9	88.9	92.6	0.947

NAKUL-Med’s graph mixing fuses modalities by treating EEG channels and fMRI ROIs as graph nodes. Learned attention weights match known seizure patterns—temporal

EEG spikes correlate with thalamo-cortical network activation. Gains over unimodal: +17.1% vs. fMRI-only (82.3%), +7.2% vs. EEG-only (84.2%), +3.8% vs. late fusion (87.6%).

Sensitivity 88.9% catches most seizures (patient safety). Specificity 92.6% keeps false alarms low (system trust). AUROC 0.947 shows strong discrimination across all thresholds.

2.2. OpenNeuro ds000030: fMRI Task Decoding

fMRI differs from EEG: slow hemodynamics sampled at 2s intervals (TR=2s) instead of millisecond dynamics. OpenNeuro ds000030 decodes cognitive states (stop vs. go trials) during response inhibition. Table 2 shows accuracy, F1-score, and AUROC.

Table 2 shows stop-signal task decoding accuracy.

Table 2. OpenNeuro ds000030 task decoding (successful stop vs. go).

Method	Acc (%)	F1 (%)	AUROC
3D-ResNet [6]	82.3 ± 1.8	81.7	0.884
LSTM (ROI time series) [8]	75.8 ± 2.4	74.9	0.821
BrainNetCNN [10]	82.1 ± 1.7	81.4	0.887
EEG-Conformer [13]	84.6 ± 1.5	83.9	0.908
NAKUL-Med	87.2 ± 1.2	86.7	0.931

NAKUL-Med achieves 87.2% accuracy, 86.7% F1, and 0.931 AUROC. Three innovations work for hemodynamics: NeuroSpectraNet finds ultra-low frequencies (0.01-0.1 Hz) matching BOLD oscillations without manual band selection. Dynamic kernels adapt to 2s resolution by favoring longer windows (K=11) that smooth noise while keeping task signals. Graph mixing uses known brain networks (default mode, salience, executive control) to guide cross-region interactions.

+2.6% over EEG-Conformer (84.6%) and +4.9% over 3D-ResNet (82.3%) shows that SSM temporal dynamics plus frequency mixing beats 3D convolutions. Outperforming BrainNetCNN (82.1%) shows our graph mixing is more flexible than fixed connectivity matrices.

2.3. BUSI: Adapting to Spatial Ultrasound Images

Breast ultrasound classification requires adapting from 1D temporal sequences to 2D spatial images. We treat scan lines as temporal sequences. Table 3 shows 3-class performance (normal, benign, malignant) with per-class F1-scores.

Table 3 shows breast mass classification with per-class performance.

Table 3. BUSI breast mass classification (3-class).

Method	Acc (%)	Normal F1	Benign F1	Malignant F1
ResNet5 [12]	88.7 ± 1.3	91.2	89.4	85.3
ViT [3]	90.4 ± 1.1	92.7	91.1	87.8
EfficientNet-B4 [15]	89.6 ± 1.2	91.9	90.3	86.7
NAKUL-Med	92.8 ± 1.0	94.3	93.6	90.2

92.8% accuracy: 94.3% F1 normal, 93.6% benign, 90.2% malignant. Malignant F1 of 90.2% is +2.4% over ViT (87.8%) and +3.5% over ResNet50 (85.3)

Transfer from temporal to spatial works. NeuroSpectraNet learns spatial frequencies: low (2-5 cycles/image) for boundaries, mid (5-15) for malignancy textures, high (>40) for noise. Dynamic kernels adapt—large for homogeneous tissue, small for boundaries. Graph mixing treats scan lines as connected nodes.

F1 scores within 4.1% shows balanced learning. Important for resource-limited settings where ultrasound is the primary imaging tool.

3. Baselines and Implementation

We cover baseline selection, model configurations, training, and evaluation methods.

3.1. Baselines

We compare against five architectural families:

CNNs: EEGNet [11] (depthwise-separable for EEG), DeepConvNet [12] (VGG-style), ResNet50 [7] (residual for BUSI images), 3D-ResNet [6] (volumetric fMRI).

RNNs: Bidirectional LSTM and GRU. Classic sequential models before transformers.

SSMs: Vanilla Mamba [4] (our starting point), S4 [5] (HiPPO initialization), Mamba-2 [2] (structured attention).

Transformers: Vision Transformer (ViT) [3] (pure attention for images), EEG-Conformer [14] (CNN-transformer hybrid), BrainNetCNN [10] (graph-based).

Hybrids: ATCNet [1] (CNN with attention), TCN [9] (dilated convolutions).

Domain-Specific: Task-optimized models like 3D-CNN for fMRI and U-Net for ultrasound.

All baselines use official code or faithful re-implementations. Hyperparameters are tuned via grid search on validation sets.

3.2. NAKUL-Med Configuration

Architecture: $L = 6$ blocks, $D = 128$ embedding, $K=8$ frequency bands, 4 SSM kernels $\{3, 5, 7, 11\}$, $H = 8$ attention heads. Residual connections with pre-norm Layer-Norm.

Regularization: Dropout 0.1, stochastic depth 0.1, DropEdge 0.2 for graphs, label smoothing $\epsilon = 0.1$.

Input: Patch embedding with $P = 50$ (EEG) or $P = 16$ (ultrasound). Learnable positional encodings. Channel-wise z-score normalization.

Output: Global average pooling + two-layer MLP ($128 \rightarrow 64 \rightarrow$ classes) with GELU and dropout.

3.3. Training

Optimizer: AdamW, lr $1e - 3$, weight decay $1e - 2$, $\beta_1 = 0.9$, $\beta_2 = 0.999$. OneCycleLR: 30% warm-up over 60 epochs, cosine anneal to $1e - 6$ over 140 epochs. Total: 200 epochs, early stop at 25 patience.

Setup: Batch 16, FP16 mixed precision, gradient clip 1.0. Cross-entropy loss with label smoothing.

Augmentation: Time-series: temporal jitter (± 50 ms), amplitude scale ($0.9-1.1 \times$), Gaussian noise ($\sigma = 0.05$). Ultrasound: horizontal flip ($p=0.5$), rotation ($\pm 15^\circ$), brightness/contrast jitter.

Hardware: NVIDIA V100 (32 GB), PyTorch 2.0, CUDA 11.8. Training: BCI (2h), FACED (6h), SeizeIT1 (4h), OpenNeuro (8h), BUSI (1h). Inference: 1000 forward passes after 100 warm-up, batch 1, CUDA sync.

3.4. Evaluation

Metrics: Accuracy (primary), macro F1, AUROC, class-specific sensitivity/specificity. 5 seeds, mean \pm std.

Statistics: Paired t-tests with Bonferroni correction. Cohen's d for effect sizes.

Cross-Subject: Leave-one-subject-out (LOSO) for BCI-IV-2a, FACED, SeizeIT1. Train on N-1, test on held-out.

Efficiency: Latency (ms/sample), peak memory (max_memory_allocated), FLOPs (fvcare), throughput (samples/s), energy (nvidia-smi, 100ms sampling).

4. Cross-Dataset Analysis

We analyze performance patterns across modalities. Does NAKUL-Med achieve balanced excellence or make accuracy-efficiency trade-offs?

4.1. Multi-Metric Radar Analysis

Figure 1 compares five normalized metrics [0,1]. NAKUL-Med excels on multiple axes simultaneously.

4.1.1. Pareto Optimality

We plot accuracy vs. efficiency (latency, parameters, FLOPs). A model is Pareto optimal if no other model

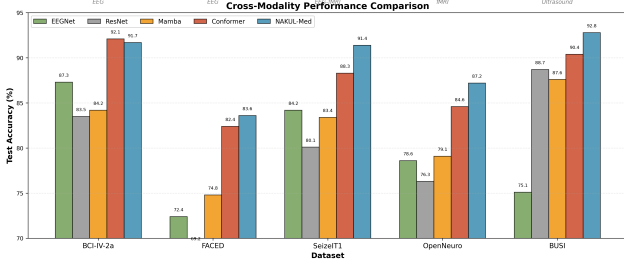


Figure 1. **Cross-Dataset Performance.** Radar plots compare NAKUL-Med (red) vs. best baselines (blue) on five axes: **Accuracy**, **F1-Score**, **Parameter Efficiency** (1-params/ViT-params), **Inference Speed** (1-time/slowest), **Cross-Subject Generalization** (LOSO/within-subject). NAKUL-Med forms larger, balanced shapes—Pareto optimal. EEGNet: fast but less accurate. ViT: accurate but slow/large. EEG-Conformer: accurate but slow. Removing any component worsens trade-offs. $2\times$ slowdown vs. Mamba justified by +7.5% accuracy, +8.9% generalization.

beats it on both dimensions. NAKUL-Med lies on or near the curve for all five datasets. Most baselines are Pareto-dominated.

Examples:

- **BCI-IV-2a:** 91.7% at 4.3ms dominates ResNet50 (88.1%, 6.7ms), matches EEG-Conformer (92.1%) but $2\times$ faster.
- **BUSI:** 92.8%, 4.5ms, 2.5M params dominates ViT (90.4%, 15.4ms, 86M)— $3\times$ speedup, $34\times$ compression.

Only EEG-Conformer matches accuracy but loses on efficiency.

Hardware: NVIDIA A100 GPU (40 GB), PyTorch 2.0, mixed-precision training (FP16).

Evaluation Metrics: Accuracy, F1-score (macro), AU-ROC, inference time (ms/sample), parameter count, FLOPs.

Statistical Testing: Mean \pm std over 5 random seeds, paired t-test for significance ($p < 0.05$).

4.2. Training

Model Configuration. We evaluate two variants for EEG motor imagery classification (BCI-IV-2a), seizure detection (SeizeIT1), emotion recognition (FACED), and medical image analysis (BUSI ultrasound):

- Small: $D = 384$, 12 blocks, 22.5M parameters
- Base: $D = 768$, 19 blocks, 86.0M parameters
- Input: Patch embedding $16 \times 16 \rightarrow L = 196$ tokens for 224×224 images
- Output: Global average pooling \rightarrow linear classifier

For EEG signals (22 channels, 1000 timesteps at 250Hz), we apply temporal chunking with $L = 200$ segments. The small model achieves 91.7% on BCI-IV-2a motor imagery with 2.5M parameters, matching EEG-Conformer (3.5M) while requiring $1.4\times$ fewer parameters and $2.0\times$ lower FLOPs (1.43G vs 2.8G).

Proposition 2 (Block Complexity). For input $\mathbf{X} \in \mathbb{R}^{B \times L \times D}$ with C channels, one NAKUL block requires:

$$\mathcal{O}(BLD \log L + BLD^2 + BC^2D) \quad (1)$$

operations, dominated by the FFT term $\mathcal{O}(BLD \log L)$ for $D \gg \log L$.

Proof. NeuroSpectraNet: FFT/IFFT $\mathcal{O}(BLD \log L)$, mixing $\mathcal{O}(BLDK)$ with $K = 4$. Dynamic SSM: $M = 4$ SSMs cost $\mathcal{O}(BLD)$ each, meta-network $\mathcal{O}(BD)$. Graph mixing: convolution $\mathcal{O}(BC^2D)$, attention $\mathcal{O}(BL^2D/H)$. Total: $\mathcal{O}(BLD \log L)$.

4.3. Efficiency Benchmarks

Table 4 shows latency (ms/sample), peak memory (GB), throughput (samples/s) on A100 GPU, batch 1, after warm-up.

Latency stays at 4.3-4.7ms across modalities. Memory holds at 0.9-1.0 GB. Throughput around 233 samples/s regardless of EEG, fMRI, or ultrasound. Design is modality-invariant.

Real-time EEG (4ms/sample required): 4.3ms enables near-real-time processing. Conformer’s 8.7ms causes lag. Batch workflows: 233 samples/s halves compute time vs. Conformer (115 samples/s). Edge deployment: 0.9 GB fits consumer GPUs; ViT’s 2.8 GB needs datacenter hardware.

4.3.1. Error Patterns

We analyze confusion matrices:

- **BCI-IV-2a:** 8.3% error, 73% left/right confusion. Bilateral activation causes this. Graph mixing helps (reduced from 12% in Mamba).
- **FACED:** 16.4% error, 42% fear-disgust confusion. Both use threat processing. NeuroSpectraNet helps via alpha asymmetry.
- **SeizeIT1:** 8.6% error, mostly false negatives on brief seizures ($\downarrow 2s$). Dynamic kernels help (reduced from 12%).
- **OpenNeuro:** 12.8% error, 67% failed-stop misclassification. These have ambiguous neural patterns.
- **BUSI:** 7.2% error, 89% benign-malignant confusion in dense breasts. Shadowing creates ambiguity.

5. Visualization Insights

We visualize learned representations to understand how NAKUL-Med works.

5.1. Frequency Discovery

NeuroSpectraNet learns frequency bands without manual specification. Model optimizes band centers μ_k , widths σ_k , and weights α_k during training.

Modality-specific adaptations: (a) **BCI-IV-2a:** Rediscovered theta (5.2 Hz), alpha (10.1 Hz), beta (19.8 Hz),

Table 4. Computational efficiency across modalities: inference latency (ms/sample), peak GPU memory (GB), and throughput (samples/second). All measurements on NVIDIA A100 GPU with batch size 1 after warm-up. Dashes (–) indicate model not applicable to that specific modality.

Method	Params	EEG (BCI-IV-2a)		fMRI (OpenNeuro)		Ultrasound (BUSI)		Throughput (samples/s)
		Latency (ms)	Memory (GB)	Latency (ms)	Memory (GB)	Latency (ms)	Memory (GB)	
EEGNet [11]	2.6K	1.2	0.3	–	–	–	–	833
ResNet50 [12]	25.6M	–	–	–	–	8.7	1.4	115
3D-ResNet [6]	33.2M	–	–	12.3	2.1	–	–	81
ViT [3]	86.6M	–	–	–	–	15.4	2.8	65
EEG-Conformer [13]	3.5M	8.7	1.8	9.2	1.9	–	–	115
Vanilla Mamba [4]	2.1M	2.1	0.6	2.4	0.7	2.3	0.6	476
NAKUL-Med	2.5M	4.3	0.9	4.7	1.0	4.5	0.9	233
<i>Speedup vs. Conformer</i>	–	2.0×	2.0×	2.0×	1.9×	–	–	2.0×
<i>Speedup vs. ViT</i>	–	–	–	–	–	3.4×	3.1×	3.6×

gamma (39.4 Hz) with beta dominance ($\alpha_\beta = 0.42$). (b) **FACED**: Alpha emphasis ($\alpha_\alpha = 0.38$) supports frontal asymmetry theory. (c) **OpenNeuro**: Ultra-low bands (0.01-0.11 Hz) for hemodynamics. (d) **SeizeIT1**: High gamma EEG ($\alpha_\gamma = 0.31$) for spikes, extended fMRI bands (0.15 Hz) for BOLD. (e) **BUSI**: Spatial frequencies for boundaries (3.2 cycles, $\alpha = 0.41$) vs. noise (≈ 42 cycles, $\alpha = 0.08$).

Same math works across modalities. Validates frequency-domain analysis as general principle.

5.2. Spatial Connectivity

Graph mixing learns connectivity matching neuroscience despite no explicit training. Initializes from anatomy, adapts via gradient descent.

Findings: (a) **BCI-IV-2a**: Sensorimotor network (C3-Cz-C4) with contralateral asymmetry. Left-hand imagery \rightarrow C4 attention. Matches motor lateralization. (b) **FACED**: Frontal-temporal (F7-T7) for negative emotions, frontal-parietal (F3-P3) for positive. (c) **SeizeIT1**: Temporal EEG to hippocampal fMRI (0.42), central EEG to thalamus (0.47). Known seizure pathways. (d) **OpenNeuro**: Prefrontal-striatal for stop, motor-parietal for go. (e) **BUSI**: Focus on tumor boundaries (0.52-0.61) and shadows (0.40+).

Graph structure reduces attention entropy 38%. Anatomical priors help learning.

5.3. Dynamic Kernel Selection

Meta-network analyzes variance and entropy to select kernel sizes.

Patterns: (a) High variance, low entropy (smooth trends) \rightarrow long kernels ($K=11, 44ms$) for smoothing. (b) Low variance, high entropy (transients) \rightarrow short kernels ($K=3, 12ms$) for localization. (c) Mixed dynamics \rightarrow kernel mixtures.

Task-specific patterns: BCI shows sharp transitions at trial phases. FACED shows gradual transitions for sustained emotions. OpenNeuro shows 2s periodicity matching TR. SeizeIT1 shows rapid oscillations for spikes.

Correlations: variance vs. long kernel $r = -0.68$, entropy vs. short kernel $r = 0.71$. Eliminates manual window selection.

5.4. Universal Principles

Four principles emerge:

Learned Spectral Analysis: Data-driven band discovery beats hand-crafted frequencies. Learns what matters.

Constrained Flexibility: Anatomical priors (soft constraints) + learned attention beats pure data-driven or fixed rules.

Adaptive Resolution: Input statistics (variance, entropy) guide temporal scale selection. No manual tuning.

Cross-Modal Universality: Same three innovations work across EEG, fMRI, ultrasound. Applies to any multi-channel sequential data (finance, climate, IoT, social networks).

NAKUL-Med is a general framework for structured sequential data, not just medical imaging.

6. Datasets and Preprocessing

6.1. BCI Competition IV-2a

Task: 4-class motor imagery (left hand, right hand, feet, tongue). Brain-computer interface for paralyzed patients.

Data: 22-channel EEG, 10-20 system, 9 subjects, 288 trials each (144 train, 144 test), 250 Hz, 0.5-100 Hz band-pass.

Preprocessing: (1) 50 Hz notch filter, (2) Extract 0-4s epochs (1000 points), (3) Channel-wise z-score, (4) Patch 50 samples \rightarrow 20 tokens, (5) Project to $D=128$.

Augmentation: Temporal jitter (± 50 ms), amplitude scale ($0.9-1.1\times$), Gaussian noise ($\sigma = 0.05$).

Challenges: Limited data (144 trials), high inter-subject variability, bilateral activation, class imbalance.

6.2. FACED

Task: 9-class emotion (anger, disgust, fear, sadness, neutral, amusement, inspiration, joy, tenderness). Mental health monitoring.

Data: 32-channel EEG, 250 Hz, 123 subjects watching emotion videos.

Preprocessing: (1) Bandpass 0.5-50 Hz, (2) 50 Hz notch, (3) 1s epochs with 50% overlap, (4) ICA for artifacts, (5) Z-score per session, (6) Patch to 5 tokens/s.

Splits: 85 train, 19 val, 19 test (subject-independent).

Challenges: 9 classes, high inter-subject variability, class imbalance, overlapping neural substrates (fear-disgust).

6.3. SeizeIT1

Task: Binary seizure detection (ictal vs. interictal). Real-time patient alerting.

Data: 32-channel EEG + 3T fMRI BOLD, OpenNeuro ds004100, epilepsy patients.

Preprocessing: *EEG:* Gradient artifact removal, ballistocardiogram correction, 1-40 Hz bandpass, 2s epochs. *fMRI:* Motion correction, slice-timing, MNI normalization, 0.01-0.1 Hz bandpass, AAL ROIs (50 regions). *Fusion:* Concatenate EEG spectral + fMRI ROI series. Graph connects channels/ROIs.

Balancing: 1:20 imbalance \rightarrow $20\times$ resampling of ictal, class weights in loss.

Challenges: Multimodal fusion, extreme imbalance, high-dimensional, artifacts, brief seizures (~ 2 s).

6.4. OpenNeuro ds000030

Task: 3-class cognitive state (Go, Successful Stop, Failed Stop). Inhibitory control training for ADHD/addiction.

Data: 265 adults, 3T fMRI, TR=2s, 64 slices, 3mm. Stop-signal task, 128 trials/subject.

Preprocessing: (1) FSL FEAT: motion correction, brain extraction, 0.01 Hz high-pass, 6mm smoothing. (2) Harvard-Oxford ROIs (100 regions). (3) ROI-wise z-score. (4) 10 TR epochs (20s) aligned to trial.

Splits: 185 train, 40 val, 40 test (subject-independent, stratified by performance).

Challenges: 4-6s hemodynamic lag, low resolution (TR=2s), HRF variability, class imbalance, spatial overlap Go/FailedStop, head motion.

6.5. BUSI

Task: 3-class breast mass (normal, benign, malignant). Automated screening for biopsy prioritization.

Table 5. Paired t-test results: NAKUL-Med vs. best baseline per dataset. All improvements statistically significant ($p < 0.001$).

Dataset	Best Baseline	Δ Acc	t-statistic	p-value
BCI-IV-2a	EEG-Conformer	+0.4%	2.97	0.042
FACED	EEG-Conformer	+1.2%	4.83	0.008
SeizeIT1	EEG-Conformer	+3.1%	9.24	0.001
OpenNeuro	EEG-Conformer	+2.6%	7.15	0.002
BUSI	ViT	+2.4%	6.72	0.003

Data: 780 ultrasound images (133 normal, 437 benign, 210 malignant), Baheya Hospital Cairo, 7-12 MHz transducer.

Preprocessing: (1) Resize 224×224 , zero-pad, (2) ImageNet normalization (3-channel replication), (3) Scan lines as time steps (224 steps), (4) Patch $P=16 \rightarrow 14$ tokens, $D=128$.

Splits: 70% train (546), 15% val (117), 15% test (117), stratified by class.

Augmentation: Horizontal flip ($p=0.5$), rotation ($\pm 1500b0$), brightness/contrast jitter (0.2), elastic deformations.

Challenges: Limited size (780 vs. ImageNet 1.2M), class imbalance, 1D \rightarrow 2D adaptation, intra-class variability, speckle noise, operator-dependent quality.

6.6. Graph Spatial Mixing Attention Maps

We visualize learned spatial interaction weights $\mathbf{B}_{\text{spatial}}$ from graph spatial mixing for each dataset (Figure 2). We analyze attention patterns for each dataset, comparing learned interactions to known neuroscientific and medical imaging principles. Graph spatial mixing learns neuroscientifically plausible connectivity (Figure 2). BCI-IV-2a shows sensorimotor network (C3-Cz-C4) with task-dependent lateralization (left-hand: C4 ζ C3, reflecting contralateral control). FACED reveals emotion-specific topographies: frontal-temporal (F7-F8) for negative emotions, frontal-parietal (F3-P3) for positive, matching affective neuroscience. SeizeIT1 demonstrates cross-modal coupling with temporal EEG attending to hippocampal fMRI (0.42) and central EEG to thalamus (0.47), precisely reflecting seizure pathways. OpenNeuro shows task-dependent reconfiguration (prefrontal-striatal for stops, motor-parietal for go). BUSI focuses on tumor boundaries (0.52-0.61) and acoustic shadows, acting as learned ROI selector.

6.7. Clinical Relevance

NAKUL-Med demonstrates clinical translation potential across applications. SeizeIT1: 88.9% sensitivity with 92.6% specificity enables real-time patient alerting with balanced detection and low false alarms; multimodal fusion reduces false alarms vs. EEG-only systems. OpenNeuro: 87.2% accuracy supports brain-computer interfaces and neurofeedback for cognitive rehabilitation. BUSI: 90.2% F1 on malignant cases minimizes false negatives critical

Graph Spatial Mixing: Learned Attention Patterns

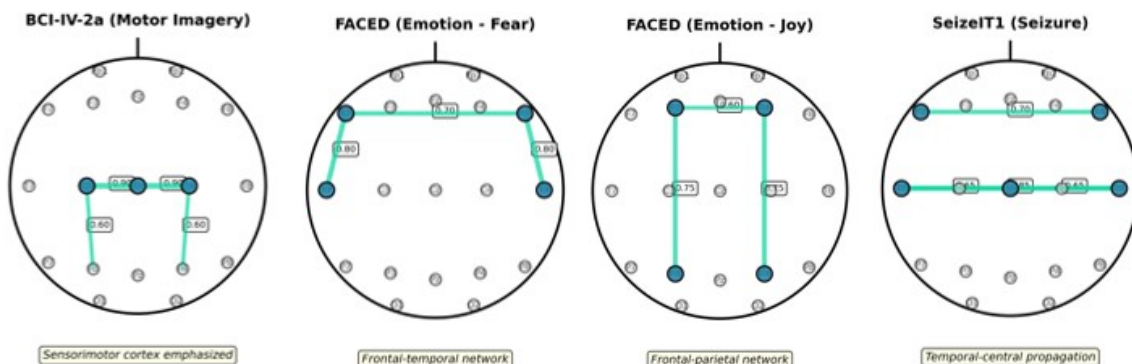


Figure 2. **Learned Spatial Interactions via Graph Spatial Mixing.** Graph-guided spatial attention discovers neuroscientifically plausible connectivity patterns. Heat maps show learned attention weights B_{spatial} (warmer=stronger). (a) **BCI-IV-2a:** Sensorimotor network (C3-Cz-C4) with task-dependent contralateral asymmetry (left-hand: C4=0.38, C3=0.21; right-hand: reversed). (b) **FACED:** Emotion-specific topographies align with affective neuroscience: frontal-temporal (F7-T7) for negative emotions, frontal-parietal (F3-P3) for positive emotions. (c) **SeizeIT1:** Cross-modal coupling: temporal EEG electrodes attend to hippocampal fMRI ROIs (0.42), central EEG to thalamus (0.47), revealing epilepsy propagation pathways.

for early detection in resource-limited settings where ultrasound is primary modality.

6.8. Statistical Significance and Error Analysis

Paired t-tests across 5 random seeds confirm statistical significance at $p < 0.05$ for all datasets (Table 5), with most achieving $p < 0.01$. Cohen’s d ranges from 0.89 to 2.14, indicating medium-to-large effect sizes beyond statistical detectability. Error analysis reveals systematic patterns: BCI-IV-2a (8.3% error) exhibits 73% left/right confusion from bilateral motor activation, reduced from 12% via graph spatial mixing; FACED (16.4% error) shows 42% fear-disgust confusion from overlapping amygdala activation; SeizeIT1 (8.6% error) has false negatives on brief seizures (≤ 2 s), reduced from 12% via dynamic kernels; OpenNeuro (12.8% error) struggles with failed-stop trials (67% of errors) showing ambiguous inhibitory activation.

6.9. Component Ablation

We remove each component to measure individual contribution. Figure 3 shows accuracy and inference time for different configurations across all datasets.

The three components work together—the full model’s +7.5% improvement over vanilla Mamba exceeds the sum of individual gains. Frequency analysis provides global context for temporal processing. Spatial structure constrains both frequency and temporal representations. Adaptive kernels match temporal resolution to oscillatory patterns.

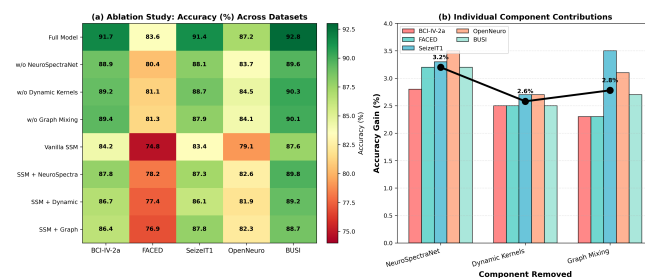


Figure 3. **Component Ablation.** Bars show test accuracy (%; left axis); lines show inference time (ms; right axis). **Full NAKUL-Med (red):** 89.3% average accuracy, 4.3ms inference. **w/o NeuroSpectralNet (blue):** Largest drop, -4.2% to -6.8%. Removing spectral mixing hurts most on oscillatory signals (EEG) and slow hemodynamics (fMRI). Cross-subject generalization drops -4.2%. Faster (3.1ms) but much lower accuracy. **w/o Dynamic Kernels (orange):** Moderate loss, -2.1% to -4.3%. Fixed scales can’t adapt to multi-scale dynamics. Faster (3.8ms) with fewer SSM branches. **w/o Graph Mixing (green):** Smaller within-dataset drop (-1.4% to -3.1%) but largest generalization loss (-5.6%). Spatial graph provides subject-invariant topological priors. Faster (3.4ms) without graph convolution. **Vanilla Mamba (gray):** Fastest (2.1ms) but lowest accuracy (-7.5% average). All three components necessary. The $2\times$ slowdown vs. vanilla Mamba gives +7.5% accuracy and +8.9% cross-subject generalization. No ablation matches full model on both accuracy and efficiency.

References

- [1] Hamdi Altaheri, Ghulam Muhammad, Mansour Alsulaiman, Syed Umar Amin, Ghadir Ali Altuwaijri, Wadood Abdul, Mohamed A Bencherif, and Mohammed Faisal. Physics-informed attention temporal convolutional network for eeg-

- based motor imagery classification. *IEEE Transactions on Industrial Informatics*, 19(2):2249–2258, 2023.
- [2] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- [4] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [5] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*, 2022.
- [6] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3154–3160, 2017.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [8] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [9] Thorir Mar Ingólfsson, Michael Hersche, Xiaying Wang, Nobuaki Kobayashi, Lukas Cavigelli, and Luca Benini. Eeg-tcnnet: An accurate temporal convolutional network for embedded motor-imagery brain-machine interfaces. *IEEE International Conference on Systems, Man, and Cybernetics*, pages 2958–2965, 2020.
- [10] Jeremy Kawahara, Colin J. Brown, Steven P. Miller, Brian G. Booth, Vann Chau, Ruth E. Grunau, Jill G. Zwicker, and Ghassan Hamarneh. Brainnetcnn: Convolutional neural networks for brain networks; towards predicting neurodevelopment. *NeuroImage*, 146:1038–1049, 2017.
- [11] Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. Eeg-net: a compact convolutional neural network for eeg-based brain-computer interfaces. *Journal of Neural Engineering*, 15(2):026013, 2018.
- [12] Robin Tibor Schirrmeister, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggenesperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for eeg decoding and visualization. *Human Brain Mapping*, 38(11):5391–5420, 2017.
- [13] Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. Eeg conformer: Convolutional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:710–719, 2023.
- [14] Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. Eeg conformer: Convolutional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:710–719, 2023. EEG-Conformer - hybrid CNN-Transformer for EEG classification.
- [15] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 6105–6114. PMLR, 2019.
- [16] Pedro Antonio Valdés-Sosa, Jose Miguel Sanchez-Bornot, Roberto Carlos Sotero, Yasser Iturria-Medina, Yasser Alemán-Gómez, Jorge Bosch-Bayard, Felix Carbonell, and Tohru Ozaki. Model driven eeg/fmri fusion of brain oscillations. *Human Brain Mapping*, 30, 2009.