

Gene-Brain CCA Analysis: Concise Report

Author: Allie Date: January 14, 2026 Dataset: UK Biobank (N=4,218)

Executive Summary

This study investigated whether combining genetic embeddings (from DNABERT-2 foundation model) with brain imaging (fMRI) data improves Major Depressive Disorder (MDD) prediction.

Key Findings:

- * Gene-only prediction achieves AUC 0.759 (holdout)
 - * fMRI adds no predictive value (early fusion AUC 0.762, +0.003)
 - * Unsupervised CCA/SCCA underperforms direct supervised learning by 17-23 AUC points
 - * Full 768-D embeddings improve performance by +29% vs scalar pooling
-

Korean Summary (한글 요약)

DNABERT-2 foundation model의 유전적 임베딩과 뇌 영상(fMRI)은 주요 우울 장애(MDD) 예측에 개선을 가져온다.

주요 발견:

- * 유전만 예측 AUC 0.759 (홀드아웃)
 - * fMRI는 예측 가치를 더해지지 않음 (초기 결합 AUC 0.762, +0.003)
 - * UCCA/SCCA는 직접 감독 학습보다 17-23 AUC 점수 차이로 저작된다
 - * 전체 768-D 임베딩은 스칼라 폴링에 비해 +29% 개선을 가져온다
-

Dataset Overview

Metric	Value
Total subjects	4,218
MDD cases	1,735 (41.1%)
Controls	2,483 (58.9%)
Gene features	111 genes x 768-D
fMRI features	180 brain ROIs

Methods Summary

Experiment 1: Two-Stage CCA/SCCA

Stage 1 (Unsupervised): CCA/SCCA finds gene-brain correlations
Stage 2 (Supervised): Predict MDD from canonical variates

Gene reduction strategies:

- * Mean pooling: Average of 768 dimensions
- * Max pooling: Maximum of 768 dimensions

Experiment 2: Leakage-Safe Pipelines

Pipeline A: Interpretable SCCA on scalar genes Pipeline B: Supervised prediction with full 768-D embeddings

Key Results

Experiment 1: Mean vs Max Pooling

Metric	Mean Pooling	Max Pooling
Stage 1 p-value	0.040 (sig)	0.995 (n.s.)
Gene-only AUC	0.588	0.505

Mean pooling preserves more predictive information. Max pooling destroys the genetic signal.

Experiment 2: Pipeline B Results

Model	Holdout AUC	Note
gene_only_logreg	0.759	Best
early_fusion_logreg	0.762	Marginal
fmri_only_logreg	0.559	Chance
cca_joint_logreg	0.546	Weak
scca_joint_logreg	0.566	Poor

Master Comparison

Experiment	Best AUC	Key Insight
Exp1 Mean Pool	0.588	Mean preserves signal
Exp1 Max Pool	0.522	Max loses signal
Exp2 Pipeline B	0.762	Full embeddings best
Yoon et al.	0.851	Reference (N=29k)

Scientific Conclusions

Finding	Evidence
Gene-brain coupling weak	$r=0.368, p=0.04$

CCA/SCCA hurts prediction	0.566 vs 0.759 AUC
fMRI adds no value	AUC 0.50-0.56
Full embeddings essential	+29% improvement

Clinical Implications

English

- * Brain imaging does not improve genetic prediction of MDD
- * Foundation model embeddings must be preserved (not pooled)
- * Gene-brain alignment is statistically real but clinically irrelevant

Korean (한국어)

- * 뇌 이미징은 MDD의 유전 예측을 개선하지 않습니다
- * 기본 모델 임베딩은 보존되어야 합니다 (복합화X)
- * 유전자-뇌 정렬은 통계학적으로 실제지만 임상적으로는 관련이 없습니다

Recommendations

Immediate Next Steps

- * Gene curation - Filter to Yoon's 38 genes (Expected AUC: 0.80-0.84)
- * Remove PCA bottleneck - Use LASSO on full 85K features (Expected AUC: 0.78-0.82)
- * Match methodology - Implement 10-fold nested CV

Future Directions

- * Expand sample size (target N>10,000)
- * Test alternative brain features (network-specific)
- * Explore fMRI foundation models (BrainLM)
- * Supervised feature selection for interpretability

한국어 (Korean)

- * 표본 크기 확장 (N>10,000)
- * 다른 뇌 특성 테스트 (네트워크 특수)
- * fMRI 기본 모델 탐색 (BrainLM)
- * 감독형 특성 선택을 위한 해석 가능성

Technical Glossary

Term	Definition	Korean
AUC	Area Under ROC Curve	ROC 면적

CCA	Canonical Correlation Analysis	Canonical Correlation Analysis
SCCA	Sparse CCA	Sparse CCA
Foundation Model	Pre-trained neural network	Pre-trained neural network
fMRI	Functional MRI	Functional MRI
PCA	Principal Component Analysis	PCA
Holdout Set	Fixed test set	Holdout Set
MDD	Major Depressive Disorder	MDD

End of Report / 报告结束