

## INTRODUCTION

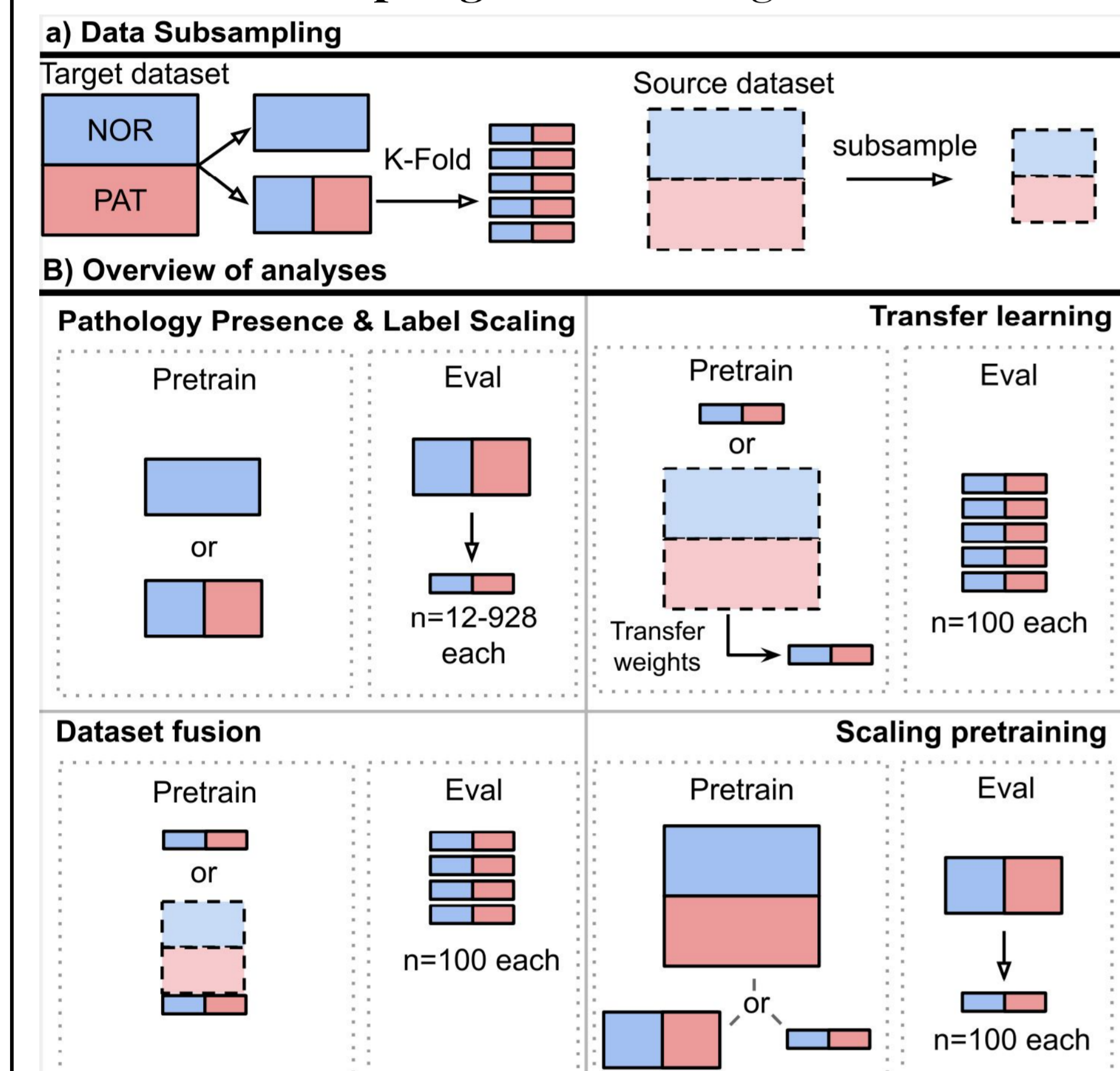
The use of neural signals such as those measured by EEG for pathology detection has been an enduring goal with the potential for high clinical relevance. Deep learning has enabled considerable progress in various fields, but requires large, labeled datasets, which are unavailable in clinical neuroimaging. Self-supervised learning allows for pretraining with unlabeled data and may therefore be a promising approach for label-scarce pathology data. Initial applications to EEG data are promising [1-3], but important questions remain, which we aim to address in the present work:

1. If only between-subject pathological information is of interest, can we improve existing SSL approaches?
2. Can SSL methods effectively differentiate between healthy and diseased populations using unlabeled data?
3. Can SSL address the issue of small pathological datasets via transfer learning or data fusion?

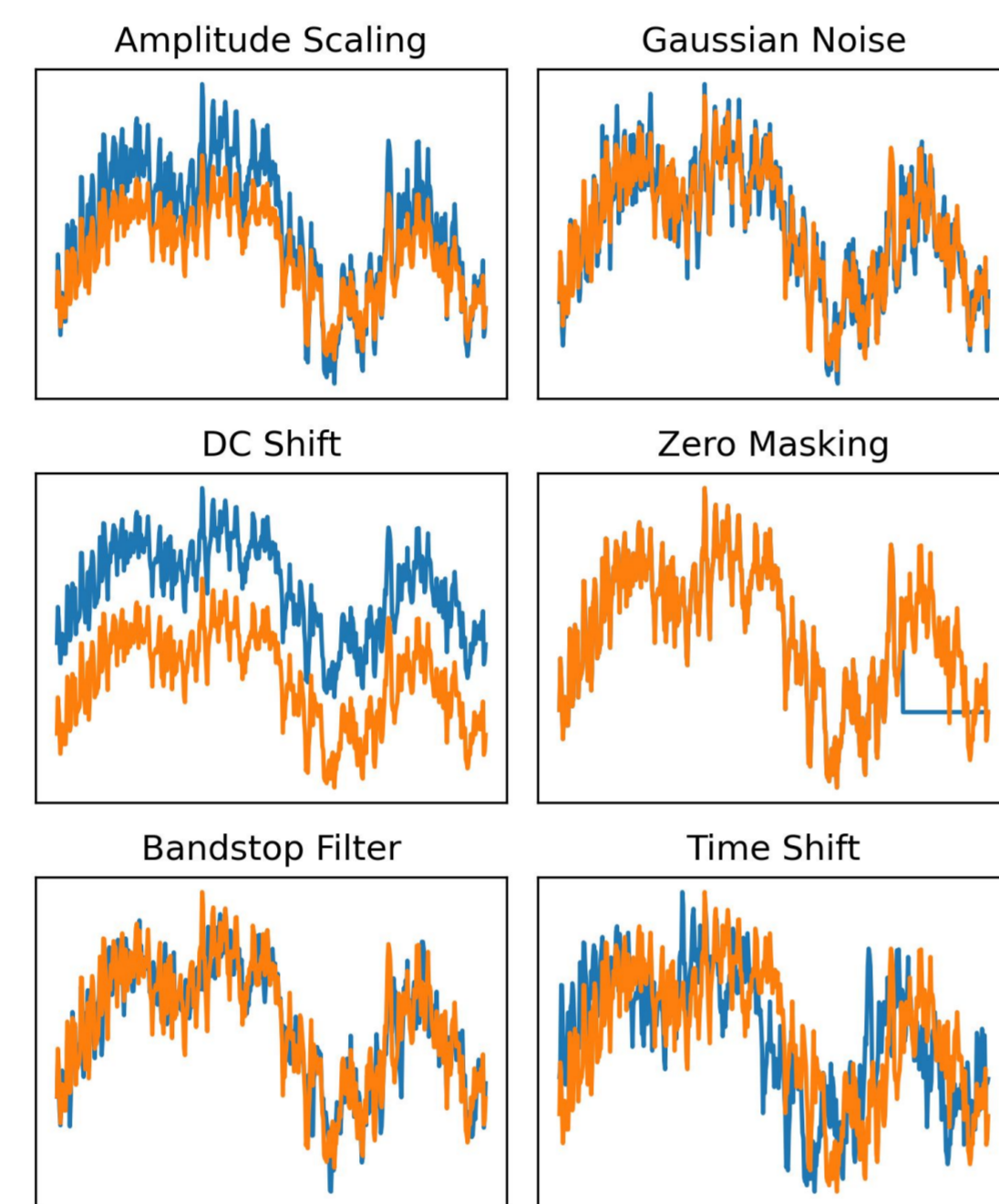
## METHODS

Datasets	n train/test	Binary target
TUH Abnormal Corpus (TUAB)	2711 / 274	Pathological abnormality
Healthy Brain Network (HBN)	2300 / 411	Low or High Functioning (CGAS)

### Dataset Subsampling to balance age and sex



### Data Augmentations [1]



### Baseline methods

Riemannian Filterbank  
Handcrafted features  
Supervised deep learning

### SSL Pretraining Methods

Augmentation-based: SimCLR, BYOL, VICReg, ContraWR [3]  
Subject-based: SubCLR

for embeddings  $z$  create mask given subject identities  $p$  compute loss for batch with size  $N$

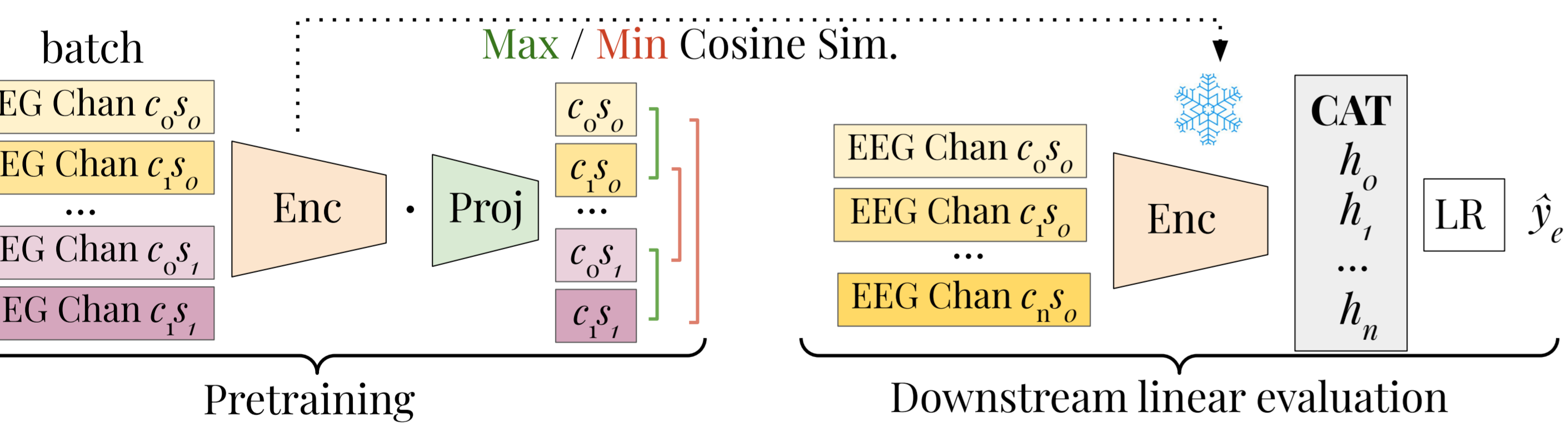
$$z = \frac{z}{\|z\|_2}$$

$$S = \frac{zz^T}{\tau}$$

$$M_{ij} = \begin{cases} 1, & \text{if } p_i = p_j \\ 0, & \text{otherwise} \end{cases}$$

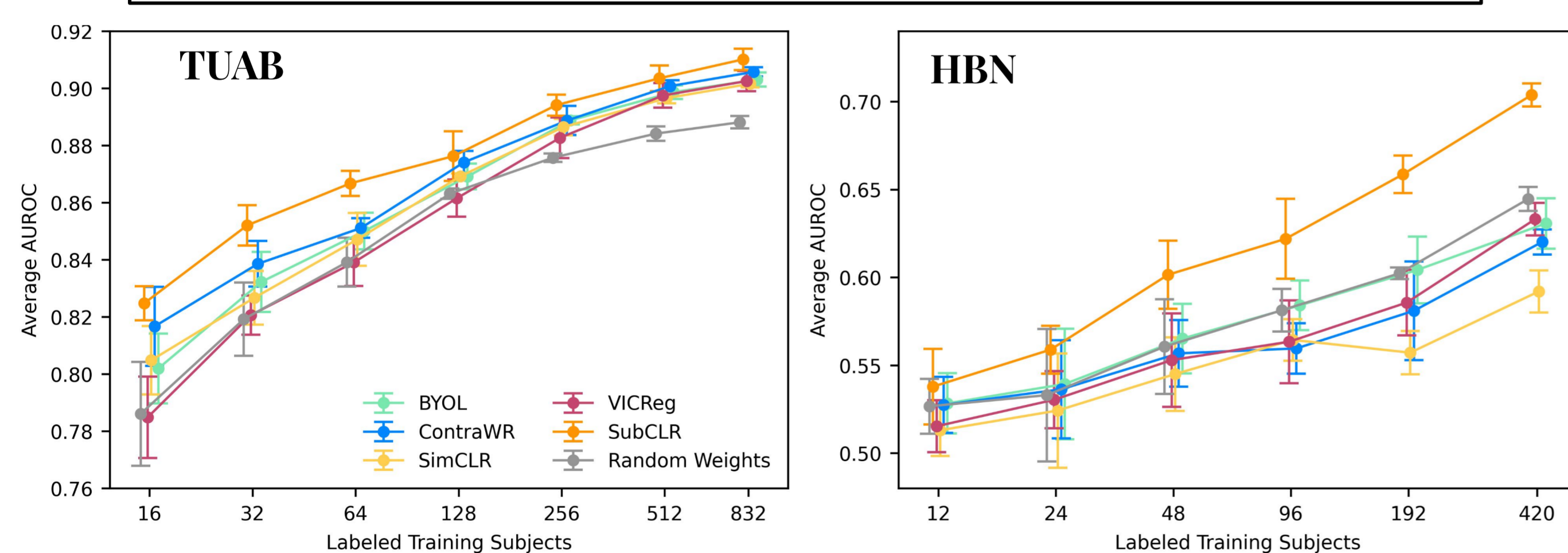
$$M_{ii} = 0, \text{ for all } i \in \{1, \dots, N\}$$

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \left( - \frac{\sum_{j=1}^N \mathbf{LSM}_{ij} M_{ij}}{\sum_{j=1}^N M_{ij}} \right)$$

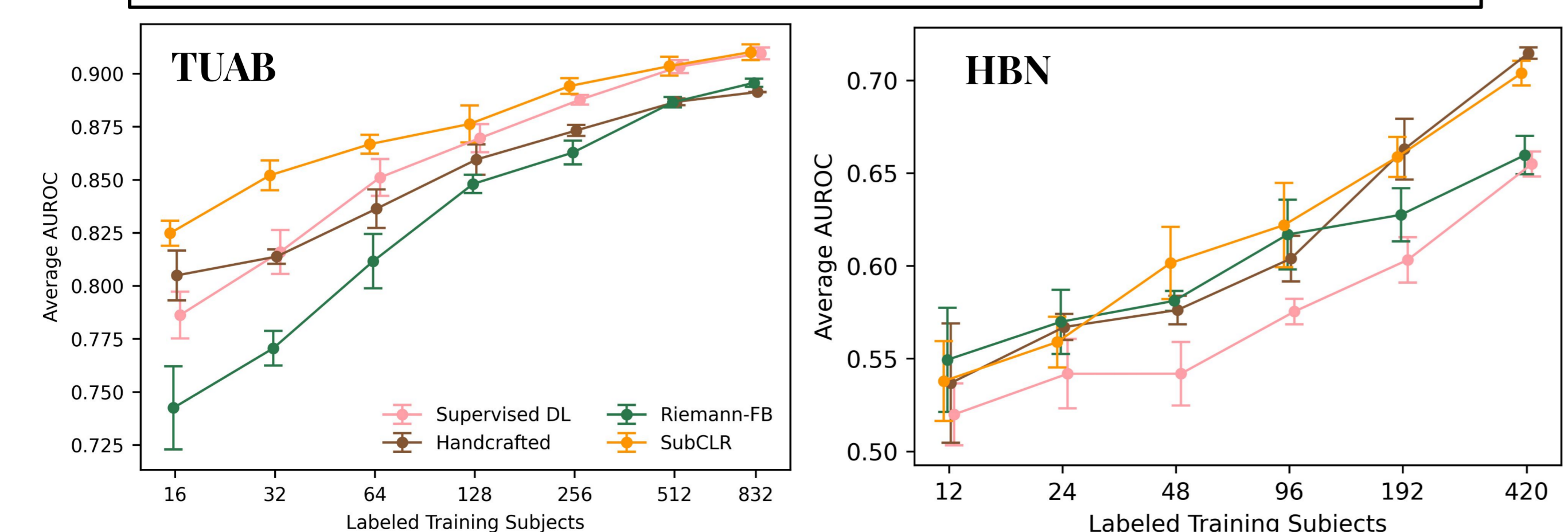


## RESULTS

### Comparison of SSL methods for pathology

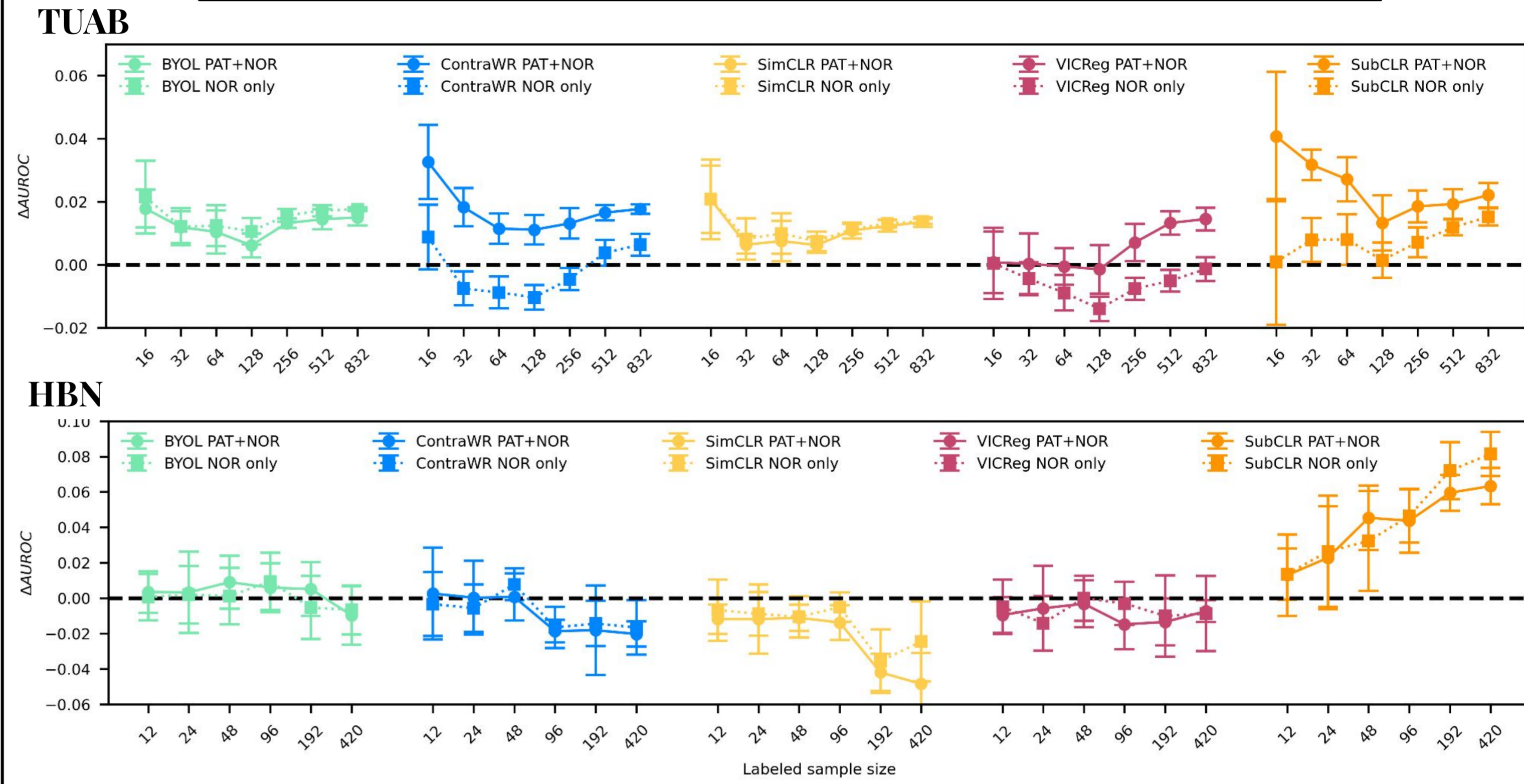


### Comparison with baseline methods

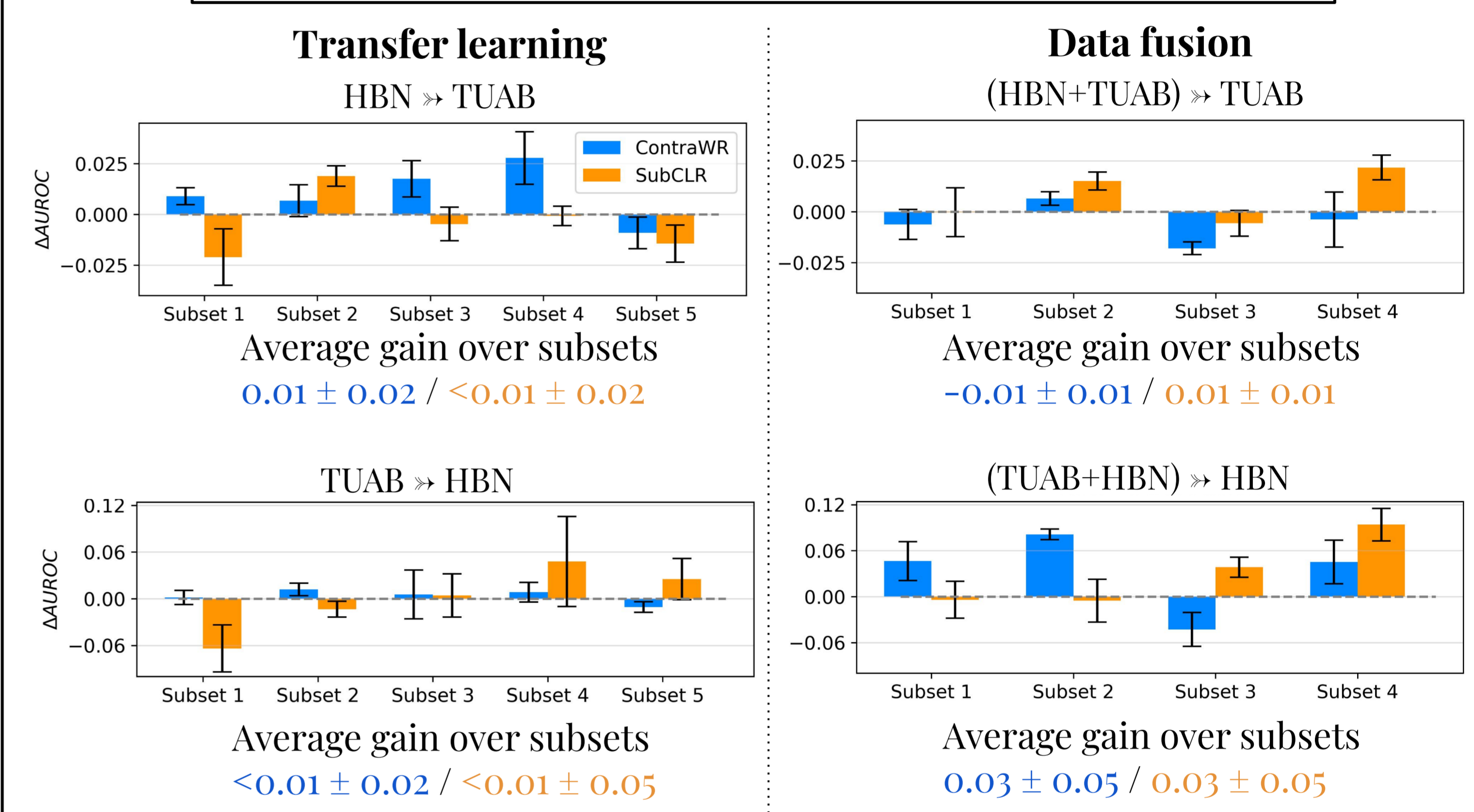


## RESULTS CONT.

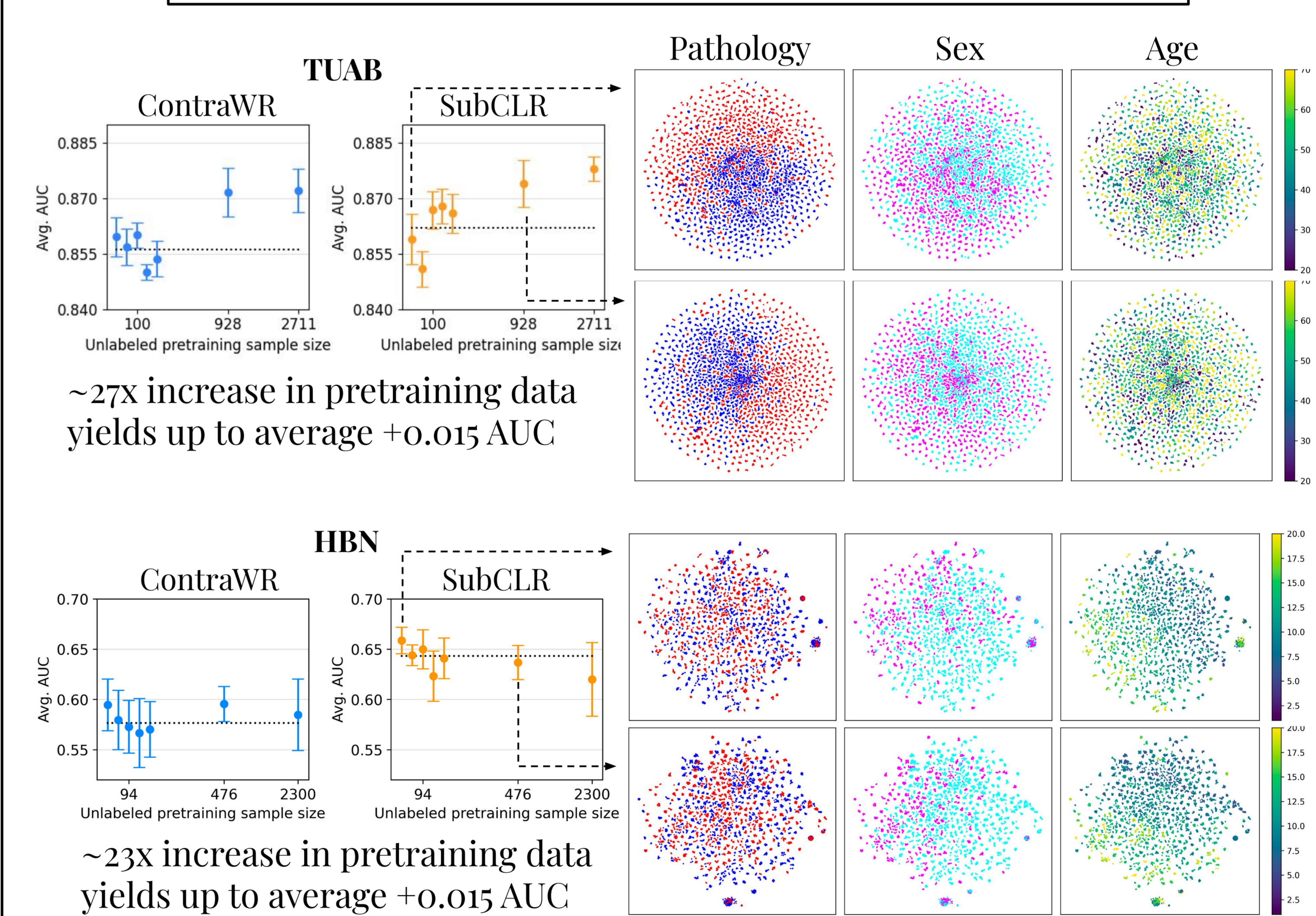
### Effect of including pathological data when pretraining



### Unreliable Effects of including transfer learning or data fusion



### ...Possibly due to poor scaling with pretraining sample size



## DISCUSSION & CONCLUSION

- Using subject identities for SSL can simplify the methodology and improve pathology detection.
- SSL can benefit from pretraining on pathological data yet improvements are modest.
- Directly transferring pretrained models does not consistently help in identifying pathology, which may be due to investigated methods scaling poorly with pretraining sample sizes.

The similar performance of diverse methods, the surprising relative performance of untrained CNNs, and poor scaling despite dramatic increases in dataset sizes, indicate that learned features are of relatively low complexity. Future work may benefit from focusing on either clinically-relevant features or improving the scaling of models, as opposed to more sophisticated data fusion methodology.

## REFERENCES

1. Mostafa Neo Mohsenvand, Mohammad Rasool Izadi, and Pattie Maes. Contrastive representation learning for electroencephalogram classification. In Machine Learning for Health, pages 238-253. PMLR, 2020.
2. Hubert Banville, Omar Chehab, Aapo Hyvärinen, Denis-Alexander Engemann, and Alexandre Gramfort. Uncovering the structure of clinical eeg signals with self-supervised learning. Journal of Neural Engineering, 18(4):046020, 2021.
3. Chaoqi Yang, Danica Xiao, M Brandon Westover, and Jimeng Sun. Self-supervised eeg representation learning for automatic sleep staging. arXiv preprint arXiv:2110.15278, 2021