# Towards Scalable, Flexible, and Interpretable Self-Supervised Learning for Multiview Biomedical Data
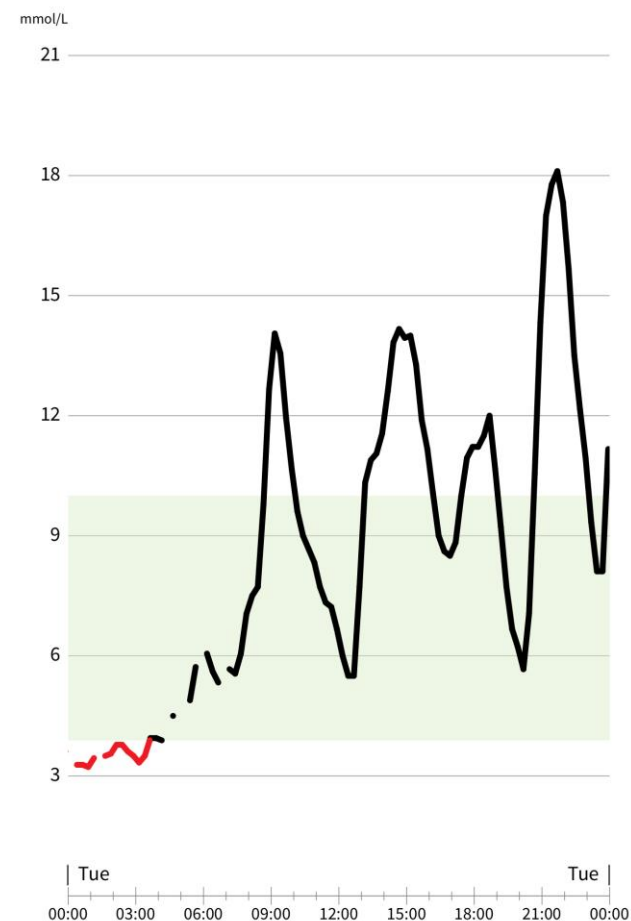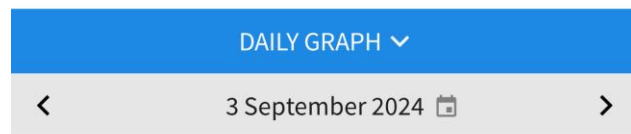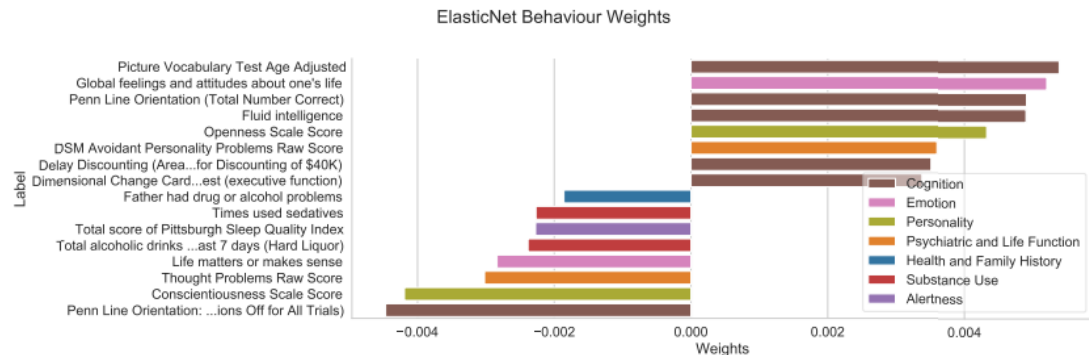
James Chapman

# Context

- Started working on *Deep Learning for Brain Behaviour Associations*

- Discovered that neuroimaging data is *really* high-dimensional and hard to work with

- And sample sizes are still small... but growing

- ... and most of the relevant code was in MATLAB

# Research Objectives

- Develop regularized Canonical Correlation Analysis (CCA) methods for improved interpretability

- Create efficient algorithms for large-scale datasets

- Extend CCA to deep learning and self-supervised settings

- Provide open-source implementations for the research community

# FRALS Framework

- Addresses a known weakness of existing sparse CCA (Partial Least Squares) methods

- Incorporates structured priors into CCA models through regularization (e.g., elastic net)

- Enhances interpretability of CCA models



$$\underset{u^{(i)}}{\mathrm{argmin}} \left\{ \| X^{(i)} u^{(i)} - t \|_2^2 + \lambda_i P_i(u^{(i)}) \right\}$$

# Weights & Loadings

• Categorized methods into explicit and implicit latent variable models

• Demonstrated the theoretical and practical robustness of loadings

• Produced meaningful multiview simulated data with known ground truth at the scale of neuroimaging studies
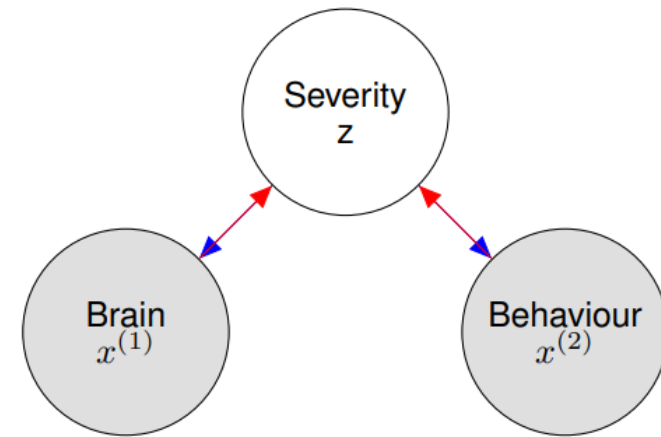


**Figure IV.1: *Forward and Backward Multiview Models:*** The generative/forward and discriminative/backward approaches in CCA.

# GEP-EY Algorithm

- Efficient solution for generalized eigenvalue problems

- Enables application of CCA and PLS to large-scale datasets (e.g., UK Biobank)

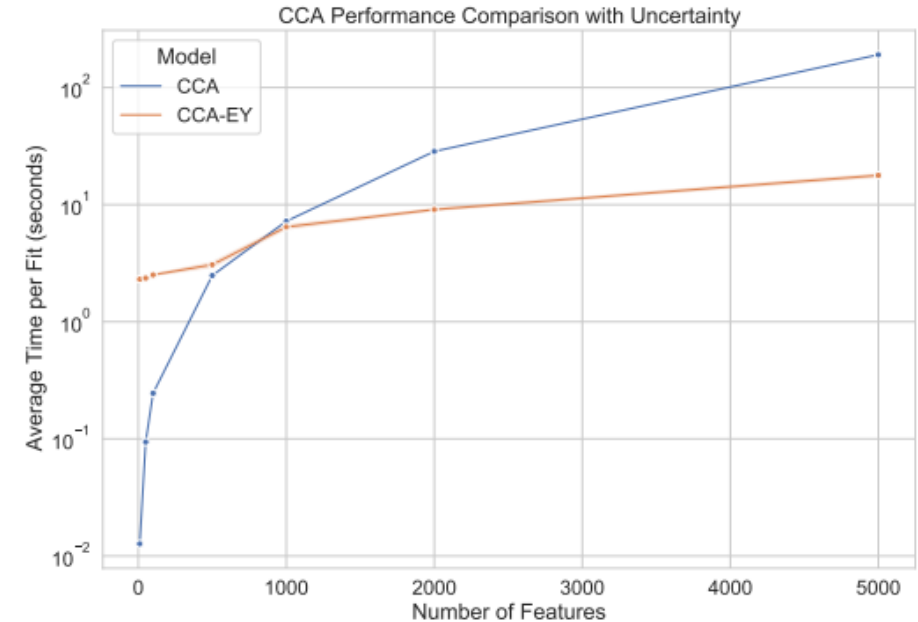- Outperforms existing methods in terms of scalability and convergence



**Figure V.1:** Comparison of the time taken to solve CCA using `eigh` and our CCA-EY method.

# Deep Learning Extensions

•DCCA-EY: Novel formulation of Deep CCA for stochastic settings

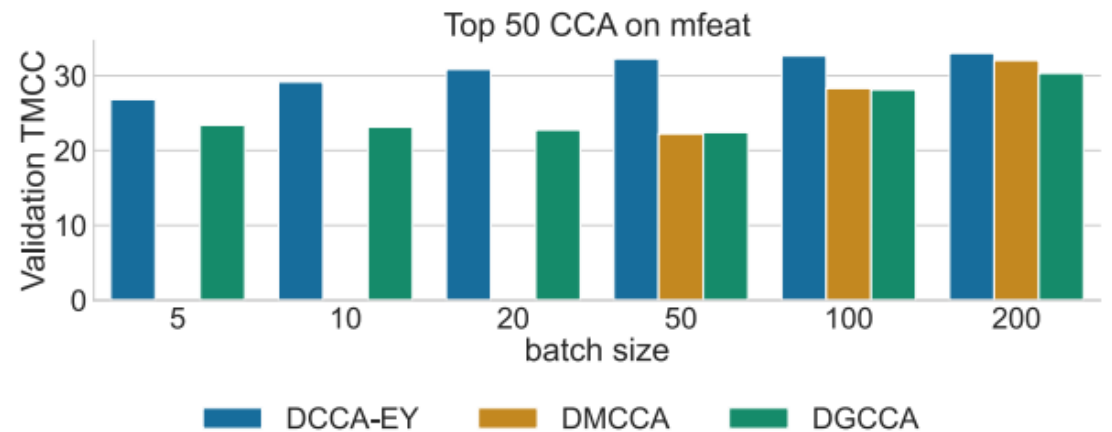•SSL-EY: New self-supervised learning method competitive with state-of-the-art with less tuning
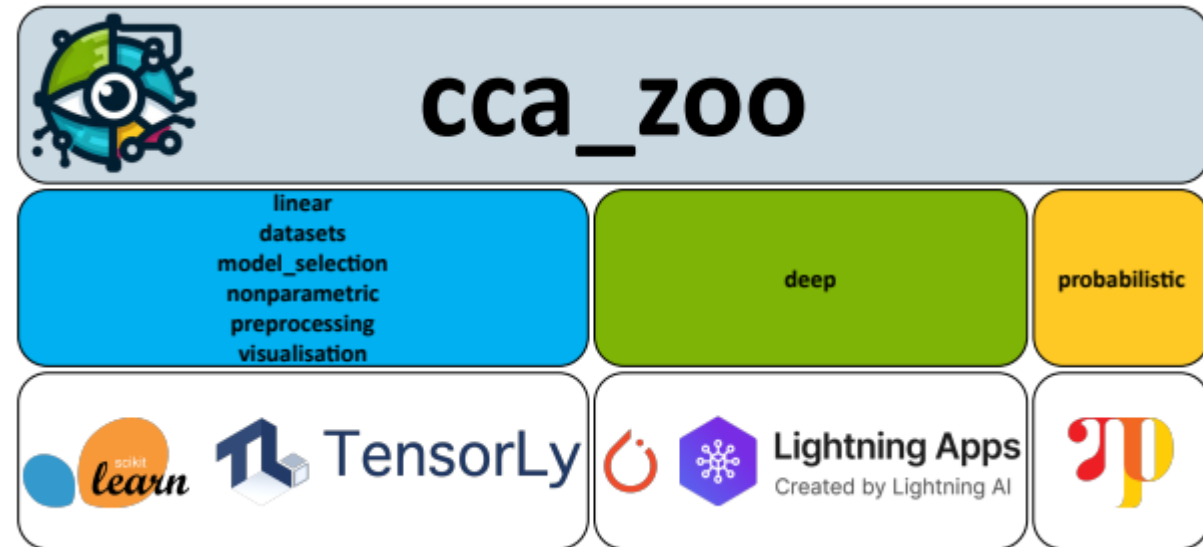


**Figure VI.8:** Deep Multi-view CCA on mfeat: Comparison across various mini-batch sizes using the Validation TMCC metric.

# CCA-Zoo Software

- Comprehensive Python library for multiview learning methods

- Implements various CCA, PLS, and related techniques

- Facilitates broader adoption and innovation in the research community

# Impact

- Can run (Multiview) CCA on UK Biobank on your laptop

- Can run Deep (Multiview) CCA on a single GPU

- 190 stars on GitHub from people who have built on my small contributions