# Predicting dementia types using machine learning approach on aggregate images
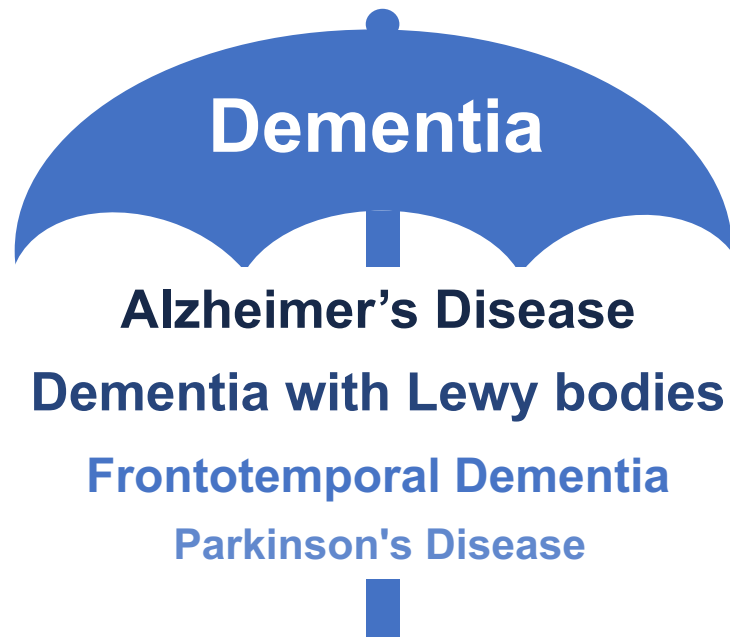
Yuting Gu

MSc HDAML Research Project

Supervisor: Dr. Yu Ye (Ye Lab)
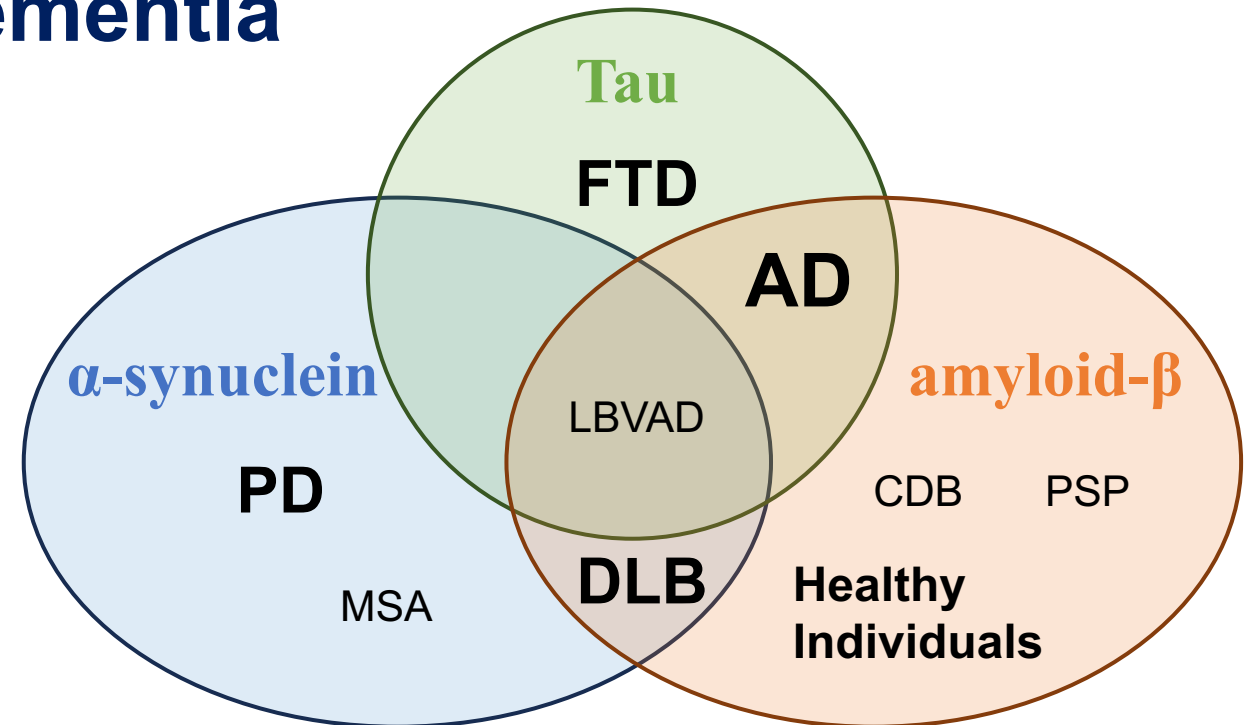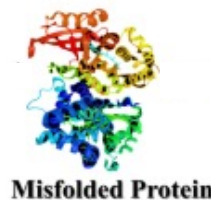
19th September 2024

# Protein Aggregates Associated with Different Types of Dementia

**Dementia**

**Alzheimer's Disease**

**Dementia with Lewy bodies**

**Frontotemporal Dementia**

**Parkinson's Disease**

**hallmark**

Progressive accumulation
**Misfolded protein aggregates**

Misfolded Protein

Tau

FTD

AD

α-synuclein

amyloid-β

LBVAD

PD

CDB     PSP
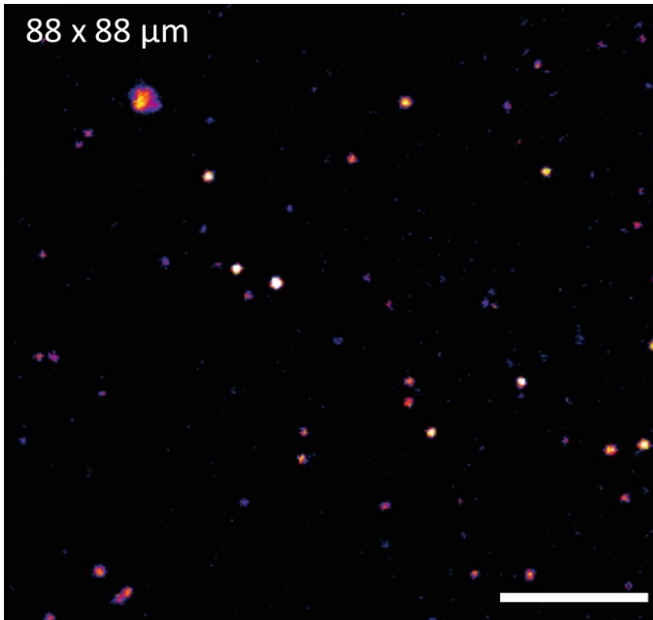
MSA

DLB

Healthy Individuals

**Challenge**

Hard to distinguish different disease associated with same protein aggregates
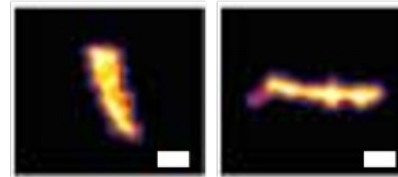
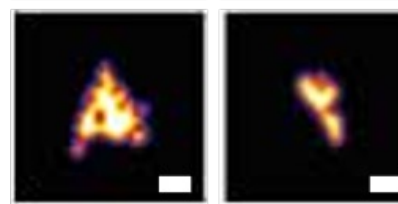# Morphological Difference of Aggregates

**Super-resolved SMLM images**

88 x 88 µm

**down to ~25 nm**
Beyond the diffraction limit of
a typical optical microscope

**PD**

**DLB**

Previous work from the lab

Super-resolved imaging technique

Some shape more prevalent to
one disease than another

**Use images capturing
aggregate morphology to
classify different types of
dementia**

Morten, M.J. et al. (2022) 'Quantitative super-resolution imaging of pathological aggregates reveals distinct toxicity profiles in different synucleinopathies', Proceedings of the National Academy of Sciences, 119(41), p. e2205591119.
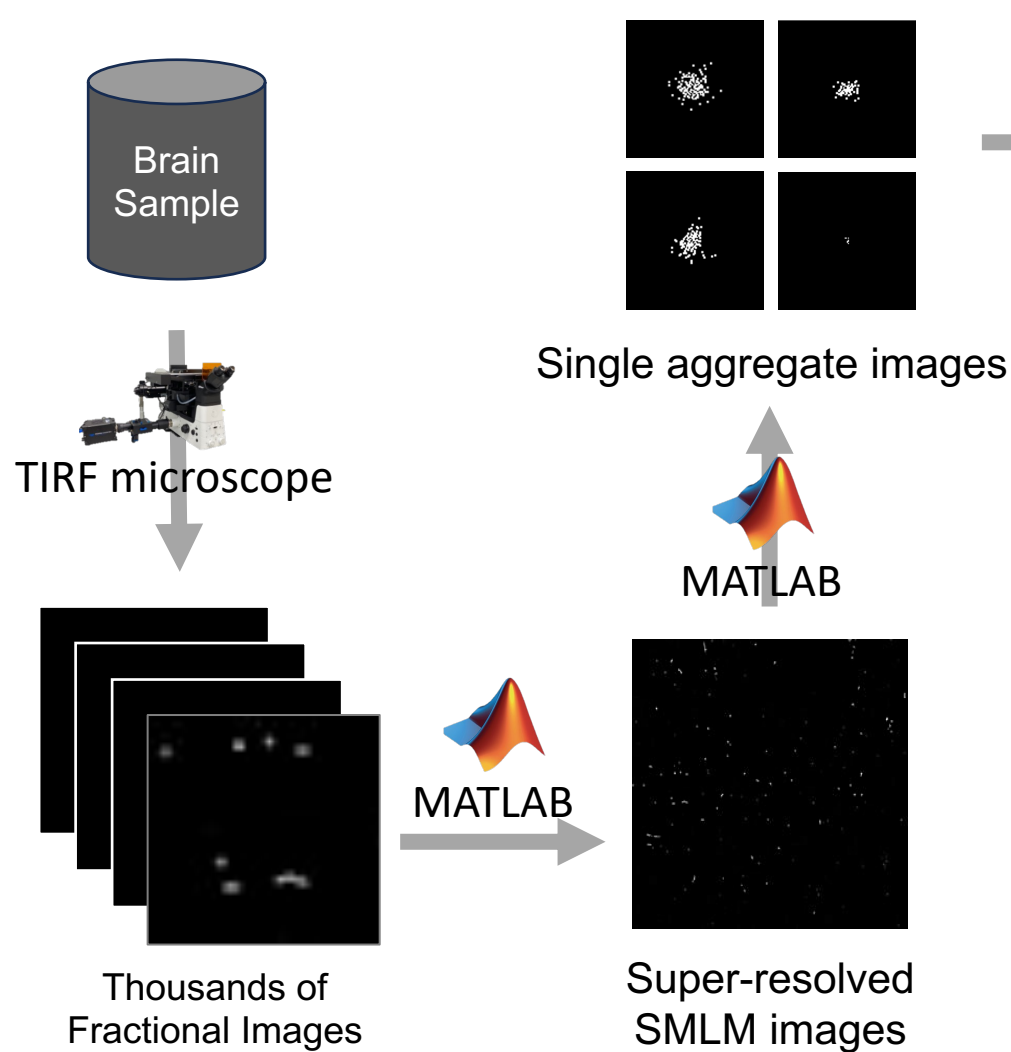
2

# Research Rationale

## Hypothesis

Each type of dementia has a subset of **morphological disease-specific** aggregates
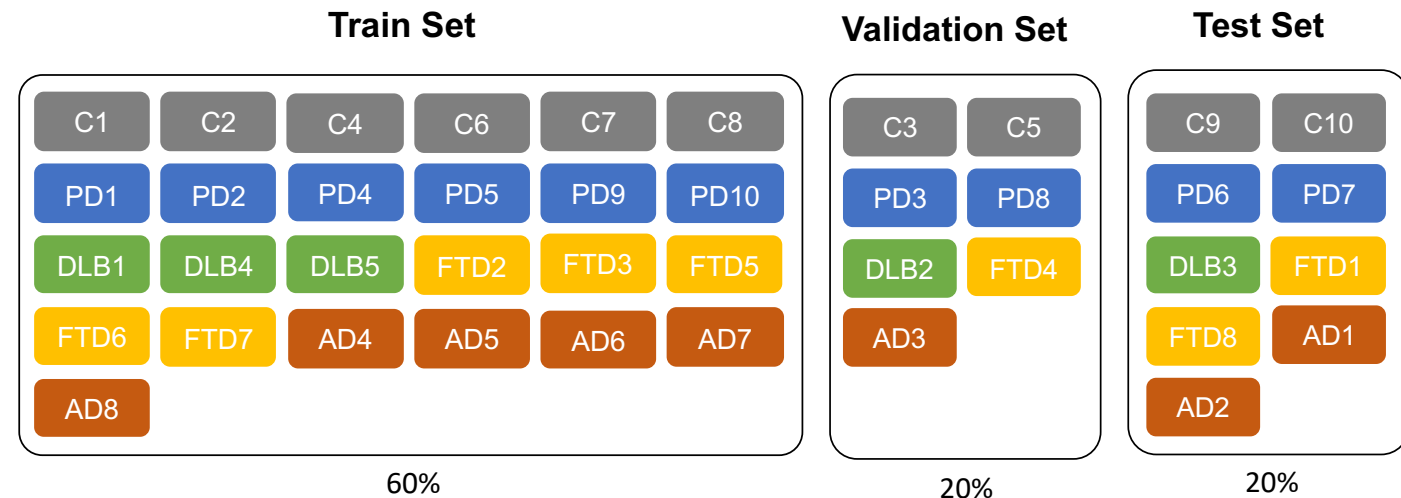
## Aims

- **Aggregate Morphology Analysis** - Use machine learning/deep learning approaches to analyse aggregate images to understand the differences in aggregate morphology between different diseases.

- **Prediction Pipeline** - Develop a prediction pipeline for patient-level dementia classification based on aggregate populations from patient samples.

# Data Acquisition and Preprocessing

Brain Sample

TIRF microscope

Thousands of Fractional Images

MATLAB

Super-resolved SMLM images

MATLAB

Single aggregate images

Single aggregate images (preprocessed)
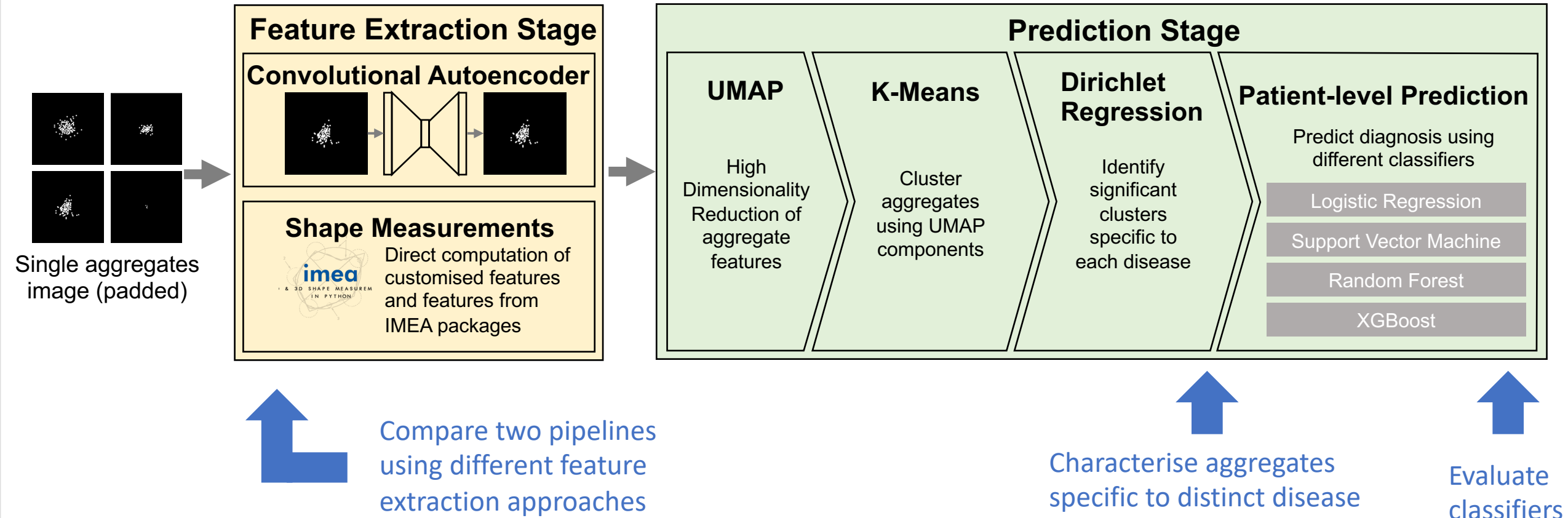
120nm          120nm

- Exclude larger than 128 pixels
- Exclude aggregate with sizes <= 4 pixels
- Padded all to 128 x 128

- Randomly assign subset of patients to Training/Validation/Testing Set

**Train Set**

| C1 | C2 | C4 | C6 | C7 | C8 |
|----|----|----|----|----|----|
| PD1 | PD2 | PD4 | PD5 | PD9 | PD10 |
| DLB1 | DLB4 | DLB5 | FTD2 | FTD3 | FTD5 |
| FTD6 | FTD7 | AD4 | AD5 | AD6 | AD7 |
| AD8 | | | | | |

60%

**Validation Set**

| C3 | C5 |
|----|----|
| PD3 | PD8 |
| DLB2 | FTD4 |
| AD3 | |

20%

**Test Set**

| C9 | C10 |
|----|----|
| PD6 | PD7 |
| DLB3 | FTD1 |
| FTD8 | AD1 |
| AD2 | |

20%

# Convolutional Autoencoder(CAE)

# Convolutional Autoencoder(CAE)

**Use Mean Squared Error (MSE) to select best dimensionality of reduction**
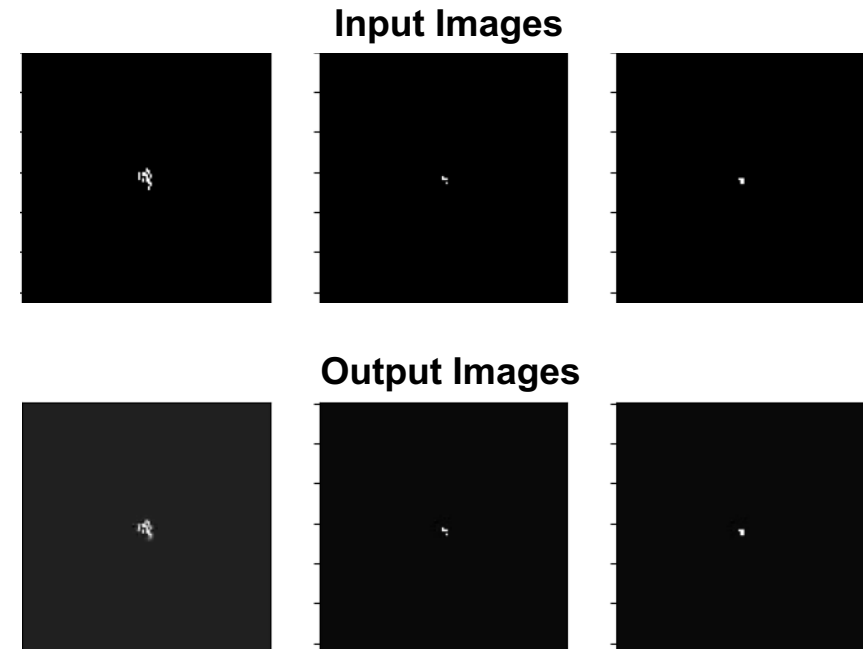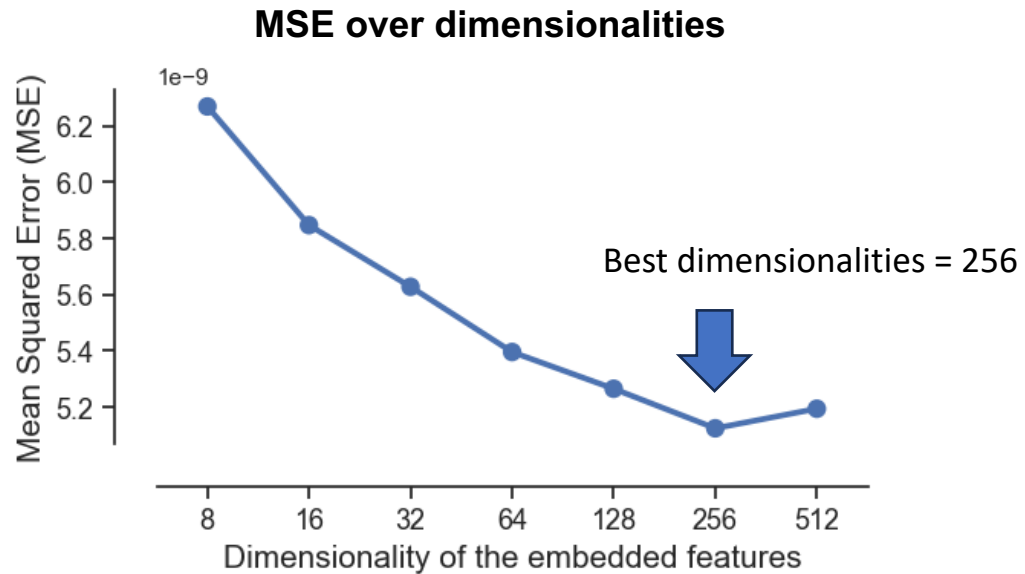
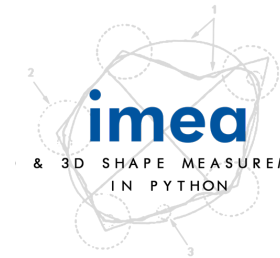Example of CAE performance at Best dim=256



MSE over dimensionalities

Best dimensionalities = 256

# Shape Measurements

## 56 shape measurement features collected for each image

5 commonly used shape measurement features

| Feature Name | Definition |
|---|---|
| Area | the total number of white pixels of an aggregate |
| Solidity | the ratio of the area to the area of a convex hall (i.e. the smallest polygon that aggregates region) of an aggregate, representing the density of this shape |
| Eccentricity | the ratio of the distance between the foci of the best-fit ellipse to its major axis length, measuring how much the shape deviates from being a perfect circle |
| Number of branches | the number of pixels in the skeletonised aggregate that are surrounded by three or four other pixels |
| Skeleton size | the number of pixels of the skeleton of an aggregate |

IMPERIAL

imea
& 3D SHAPE MEASUREM
IN PYTHON

51 computed using IMEA package

**macro descriptors**
geometric features

perimeter, area, diameter and etc.

**meso descriptors**
intermediate details

like erosions

**micro descriptors**
finer details

roughness of particle contours, specific diameter measurements like Feret and etc.

**statistical lengths**
distribution and variation in lengths

various chord lengths

# Prediction Stage: UMAP and K-Means

256 CAE image features

56 IMEA image features

Cluster Percentage

$$P_i = \frac{C_i}{\sum_{j=1}^{K} C_j}$$

$C_i$ is the number of aggregates affiliated to cluster $i$, and $K$ is the total number of clusters

Reduce dimensionalities to 20 features to enable the feasibility of clustering tools

UMAP

K-Means

20 CAE-UMAP features

20 IMEA-UMAP features

# Clustering aggregate morphologies using K-Means

**Explore best choice of K**

**Repeating** the clustering for different **6 splits** and K ranging from **10 to 150**.
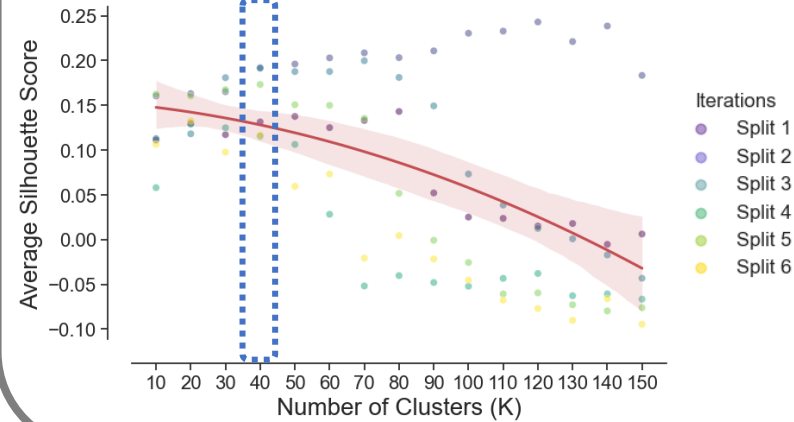
Evaluate clustering performance by **average silhouette score.** Higher score is better.

$$S(i) = \frac{b(i) - a(i))}{max(a(i), b(i))}$$
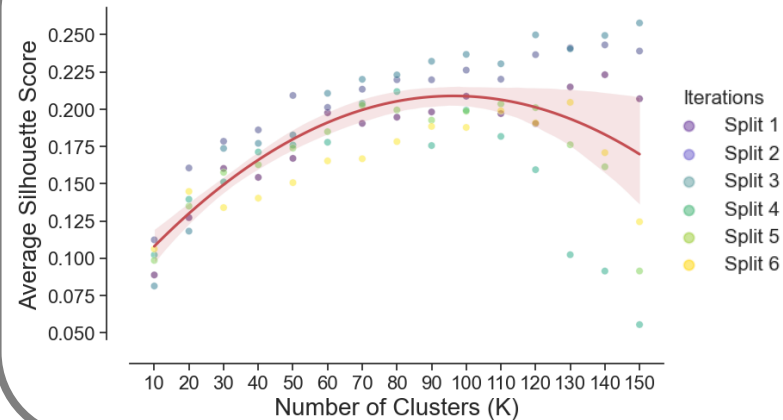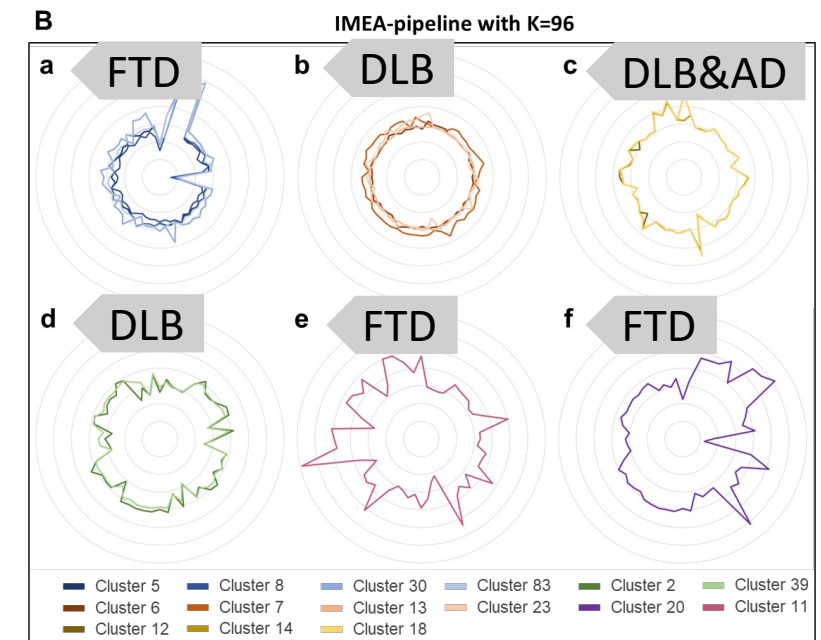
Solve best K of a fitted quadratic line



**CAE-pipeline**

No feasible solution of the fitted quadratic line

Best K = 40

**IMEA-pipeline**

Fitted quadratic line

$-1.4 \times 10^{-5}x^2 + 0.0026x + 0.083 = 0$
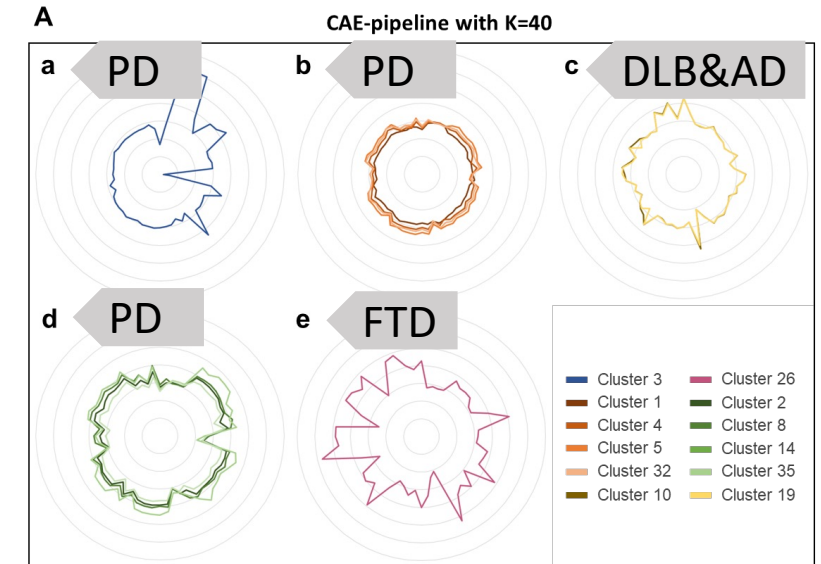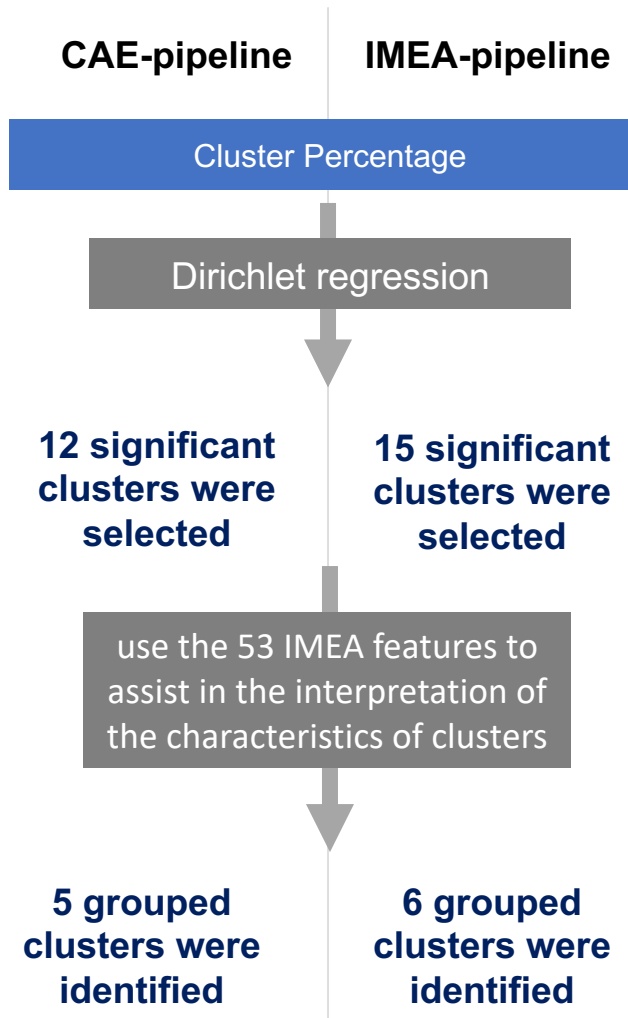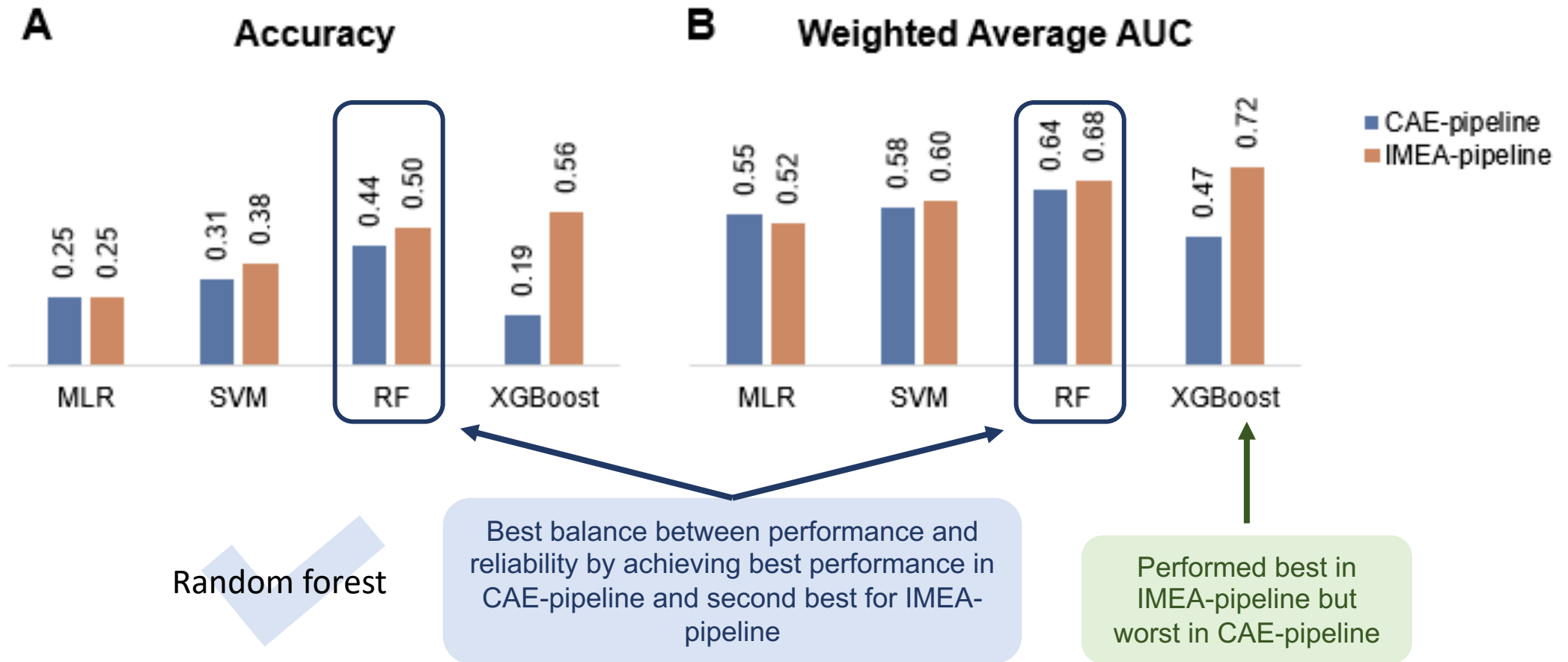
Best K = 96

# Identifying significant clusters using Dirichlet regression

- Different grouped clusters present very **distinct characteristics** by their polygon shape

- Both feature extraction approaches (CAE and IMEA) result in almost the **same grouped clusters**, except the IMEA pipeline separated one more grouped cluster (f) than the CAE pipeline

- Some grouped clusters were consistently identified and **selected** as significant clusters **for specific dementia types.**

CAE-pipeline | IMEA-pipeline

Cluster Percentage

Dirichlet regression

**12 significant clusters were selected** | **15 significant clusters were selected**

use the 53 IMEA features to assist in the interpretation of the characteristics of clusters

**5 grouped clusters were identified** | **6 grouped clusters were identified**



A CAE-pipeline with K=40
a PD  b PD  c DLB&AD
d PD  e FTD

Cluster 3 — Cluster 26
Cluster 1 — Cluster 2
Cluster 4 — Cluster 8
Cluster 5 — Cluster 14
Cluster 32 — Cluster 35
Cluster 10 — Cluster 19

B IMEA-pipeline with K=96
a FTD  b DLB  c DLB&AD
d DLB  e FTD  f FTD

Cluster 5 — Cluster 8 — Cluster 30 — Cluster 83 — Cluster 2 — Cluster 39
Cluster 6 — Cluster 7 — Cluster 13 — Cluster 23 — Cluster 20 — Cluster 11
Cluster 12 — Cluster 14 — Cluster 18

# Patient-level Prediction

**Tree based method, RF and XGBoost, tended to show superior results**



**A  Accuracy**

| | CAE-pipeline | IMEA-pipeline |
|---|---|---|
| MLR | 0.25 | 0.25 |
| SVM | 0.31 | 0.38 |
| RF | 0.44 | 0.50 |
| XGBoost | 0.19 | 0.56 |

**B  Weighted Average AUC**

| | CAE-pipeline | IMEA-pipeline |
|---|---|---|
| MLR | 0.55 | 0.52 |
| SVM | 0.58 | 0.60 |
| RF | 0.64 | 0.68 |
| XGBoost | 0.47 | 0.72 |

Random forest

Best balance between performance and reliability by achieving best performance in CAE-pipeline and second best for IMEA-pipeline

Performed best in IMEA-pipeline but worst in CAE-pipeline
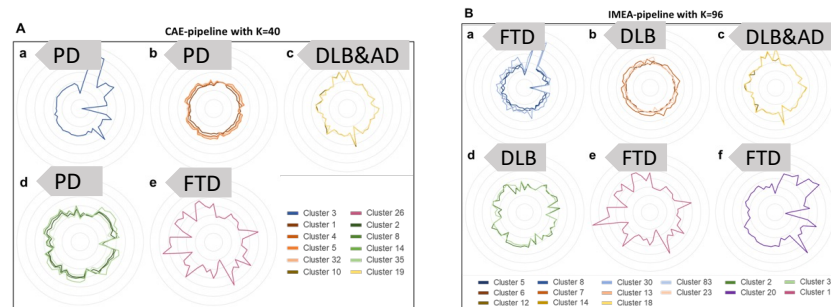
# Strengths, Limitations and Future Work

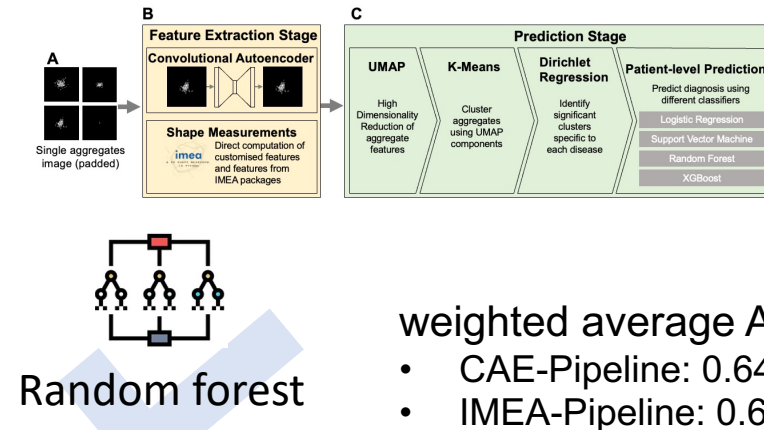| Strengths | <ul><li>**Identified clusters of disease-specific aggregates**<ul><li>Observed resemblance in characteristics of these aggregate subsets in both pipelines, indicating stability of this aggregate identification method</li></ul></li><li>**Developed end-to-end prediction pipeline**<ul><li>unsupervised feature extraction method (CAE)</li><li>fine-tune the models</li><li>practically more efficient</li></ul></li><li>**Evaluation of different feature extraction approaches and classifiers**</li></ul> |
|---|---|
| Limitations | <ul><li>Limited sample size (41 donors)</li><li>Binary (black-and-white) images, limit the data's richness.</li><li>One split of dataset, need more evaluation of generalisation</li></ul> |
| Future Work | <ul><li>Collect more donors' sample</li><li>Non-binary images with pixel intensity or coloured images</li><li>Optimise models and replicate pipeline</li></ul> |

# Conclusion

## Aim1: Aggregates morphology

**Identified and characterised the morphological differences of disease-specific aggregates**
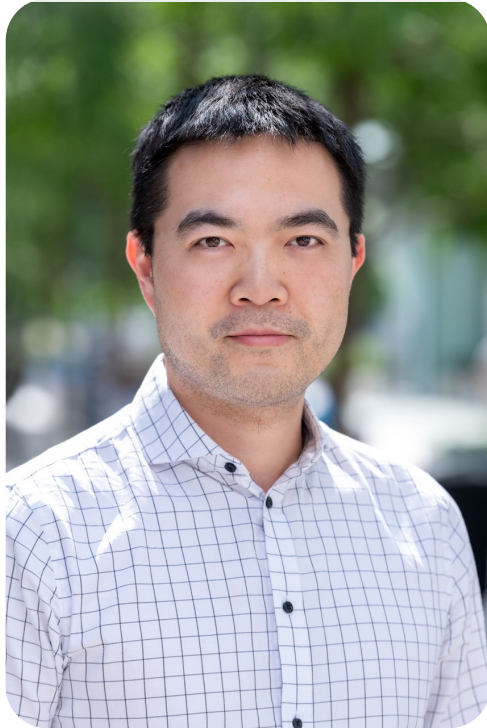


## Aim2: Prediction pipeline

**Novel development of the comprehensive end-to-end machine learning pipeline for dementia prediction using aggregate images**



Random forest

weighted average AUC
- CAE-Pipeline: 0.64
- IMEA-Pipeline: 0.68

**Important foundational framework for developing a scalable, high-throughput diagnostic tool using aggregate morphology**

**Thank you for your attention!**

# Acknowledgement

**Dr. Yu Ye**
Main Supervisor
Principal Investigator
Senior Lecturer in Molecular Neuroscience
at Imperial College London

**Hailey Gu**
HDA alumni
Research Technician in
Computational Biology

**Dr. Michael Morten**
Secondary Supervisor
Research Associate in Biophysics