

## Appendix Materials

### Manuscript Title

**PIFiA: Self-supervised Approach for Protein Functional Annotation from Single-Cell Imaging Data**

### Authors

Anastasia Razdaibiedina<sup>1,2,4</sup>, Alexander Brechalov<sup>1,2</sup>, Helena Friesen<sup>2</sup>, Mojca Mattiazzi Usaj<sup>2\*</sup>, Myra Paz David Masinas<sup>2</sup>, Harsha Garadi Suresh<sup>2</sup>, Kyle Wang<sup>1,2</sup>, Charles Boone<sup>1,2,5†</sup>, Jimmy Ba<sup>3,4,†</sup>, Brenda Andrews<sup>1,2,†</sup>

<sup>1</sup> Department of Molecular Genetics, University of Toronto, Toronto ON, Canada

<sup>2</sup> The Donnelly Centre, University of Toronto, Toronto ON, Canada

<sup>3</sup> Department of Computer Science, University of Toronto, Toronto ON, Canada

<sup>4</sup> Vector Institute for Artificial Intelligence, Toronto ON, Canada

<sup>5</sup> RIKEN Center for Sustainable Resource Science, 2-1 Hirosawa, Wako, Saitama, Japan.

\*Current address: Department of Chemistry and Biology, Toronto Metropolitan University, Toronto, ON, Canada

### Corresponding authors

Brenda Andrew

Email: [brenda.andrews@utoronto.ca](mailto:brenda.andrews@utoronto.ca)

Jimmy Ba

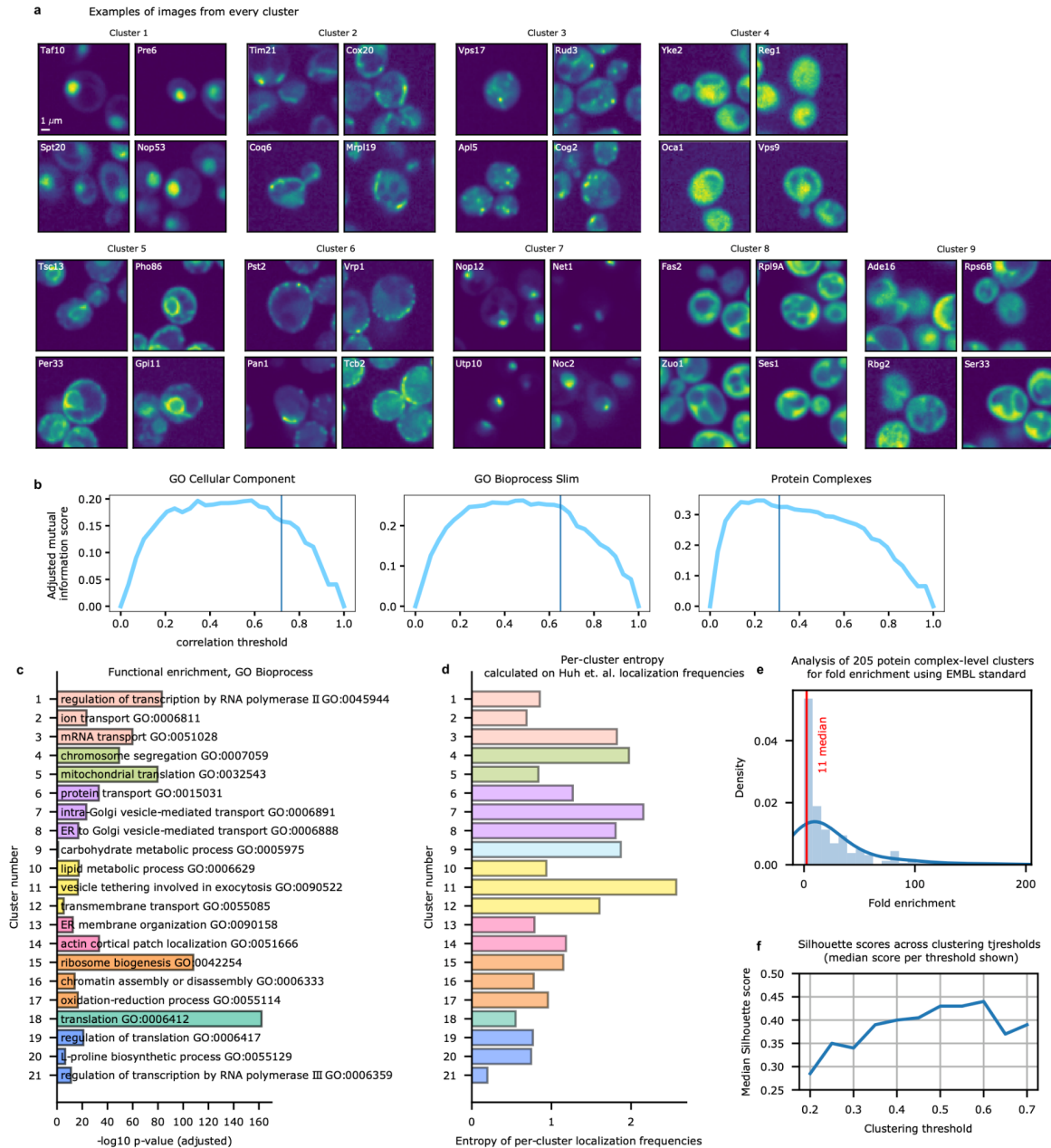
Email: [jba@cs.toronto.edu](mailto:jba@cs.toronto.edu)

Charles Boone

Email: [charlie.boone@utoronto.ca](mailto:charlie.boone@utoronto.ca)

## Table of Contents

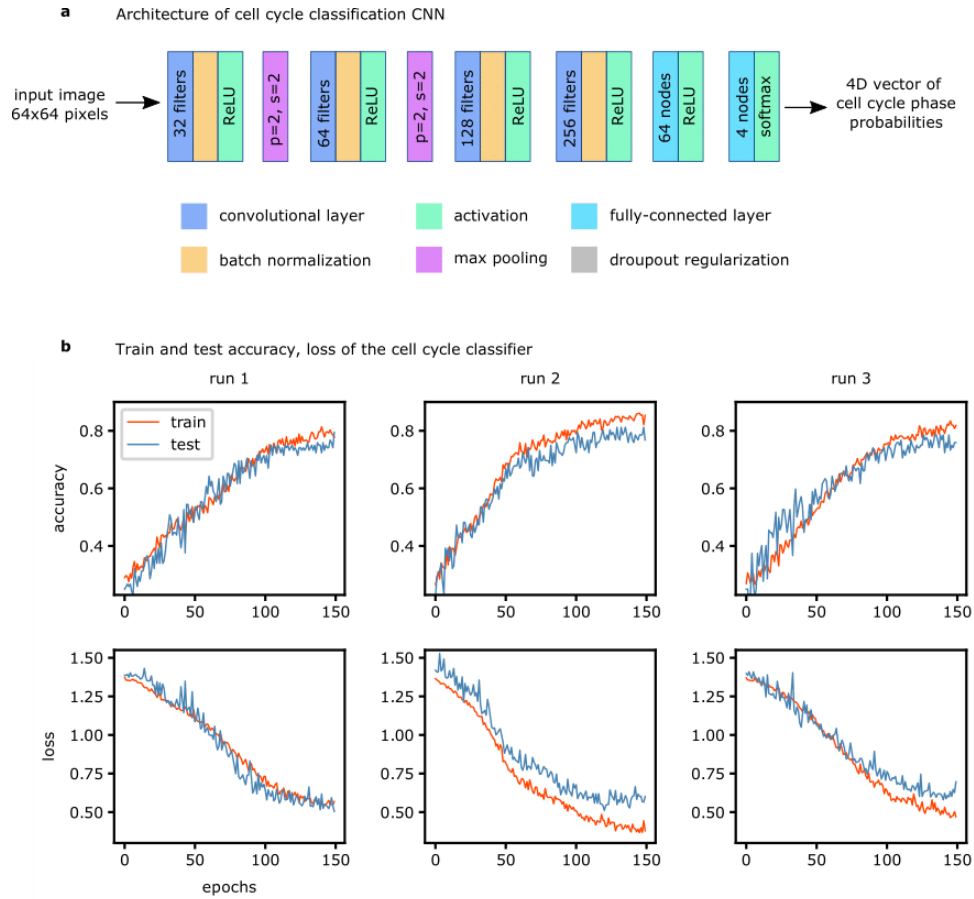
Appendix Contents	Page Number
Appendix Figure S1.....	3
Appendix Figure S2.....	5
Appendix Figure S3.....	6



**Appendix Figure S1. Examples of AMI cutoffs and cell images from nine clusters. Related to Fig 3.**

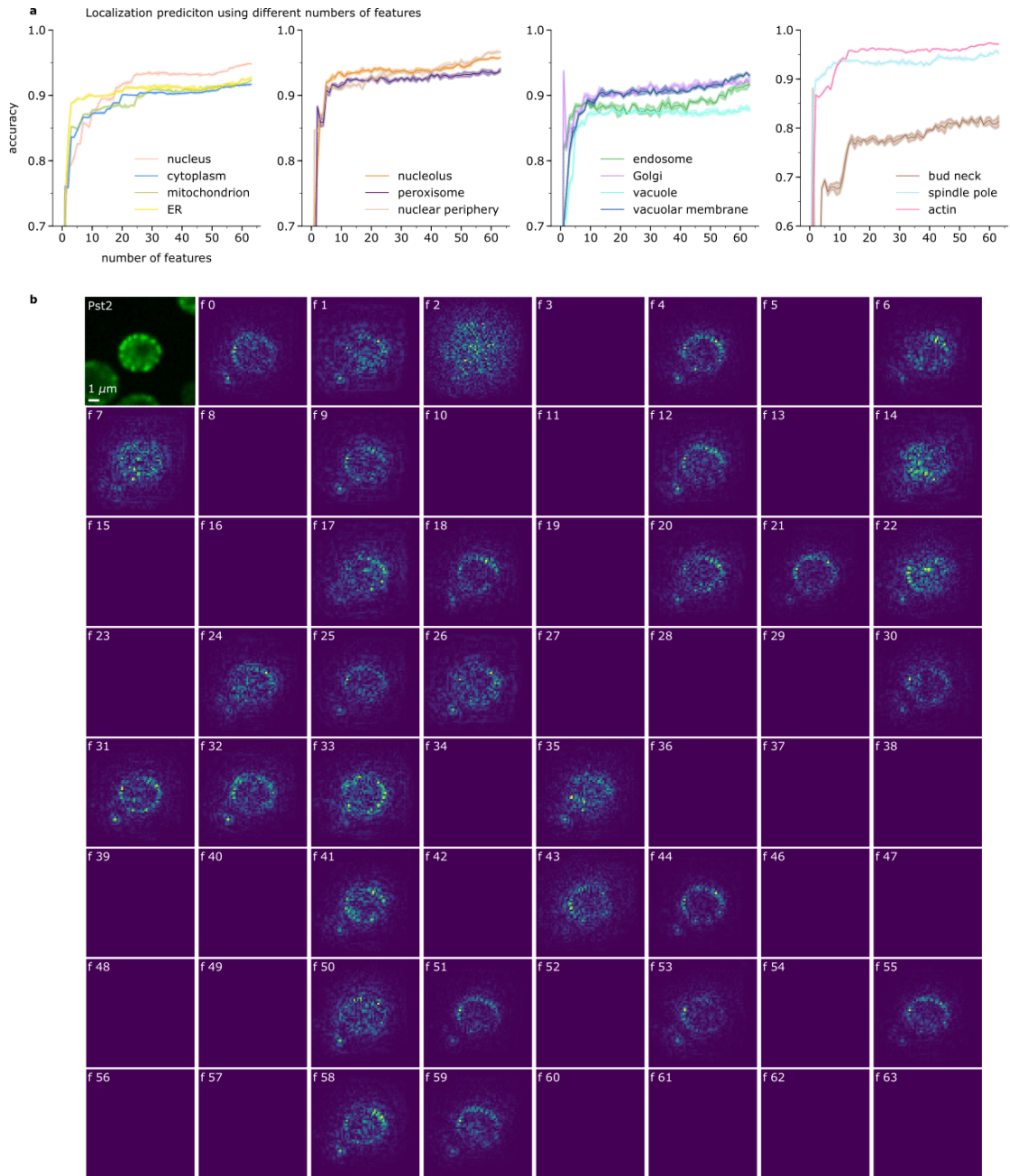
- Examples of GFP-tagged proteins from nine clusters corresponding to major cellular components.
- Plots showing adjusted mutual information across different correlation thresholds for cellular component, bioprocess and protein complex standards. Vertical line indicates a point of a dendrogram cut determined by saturation of a score.
- Bar graphs showing the top Gene Ontology Biological Process terms and corresponding functional enrichments for 21 clusters obtained from clustering by GO Biological Process Slim cutoff.

- D.** Bar graph showing entropy across 21 bioprocess clusters calculated from localization category frequencies (from Huh et al. (2003) standard).
- E.** Distribution of fold enrichments on EMBL protein complex standard for 205 clusters derived from Protein Complex cutoff.
- F.** Silhouette scores for different thresholds across single-localizing aFPs from the same localization category.



**Appendix Figure S2. Cell cycle classifier training settings. Related to Fig 4.**

- Architecture of the convolutional neural network used for cell cycle classification.
- Train and test set performance (accuracy and loss) of the cell cycle classifier across three independent network runs.



**Appendix Figure S3. Extended interpretability studies. Related to Fig 6.**

- A. Plot of accuracy of subcellular localization prediction as a function of dimensionality of feature profiles.
- B. Examples of feature-wise gradient maps obtained with SmoothGrad for all features of Pst2-GFP protein crop.