



National Cancer Institute  
at the National Institutes of Health



Icahn  
School of  
Medicine at  
Mount  
Sinai

# DreamAI: Algorithm for the Imputation of proteomics data

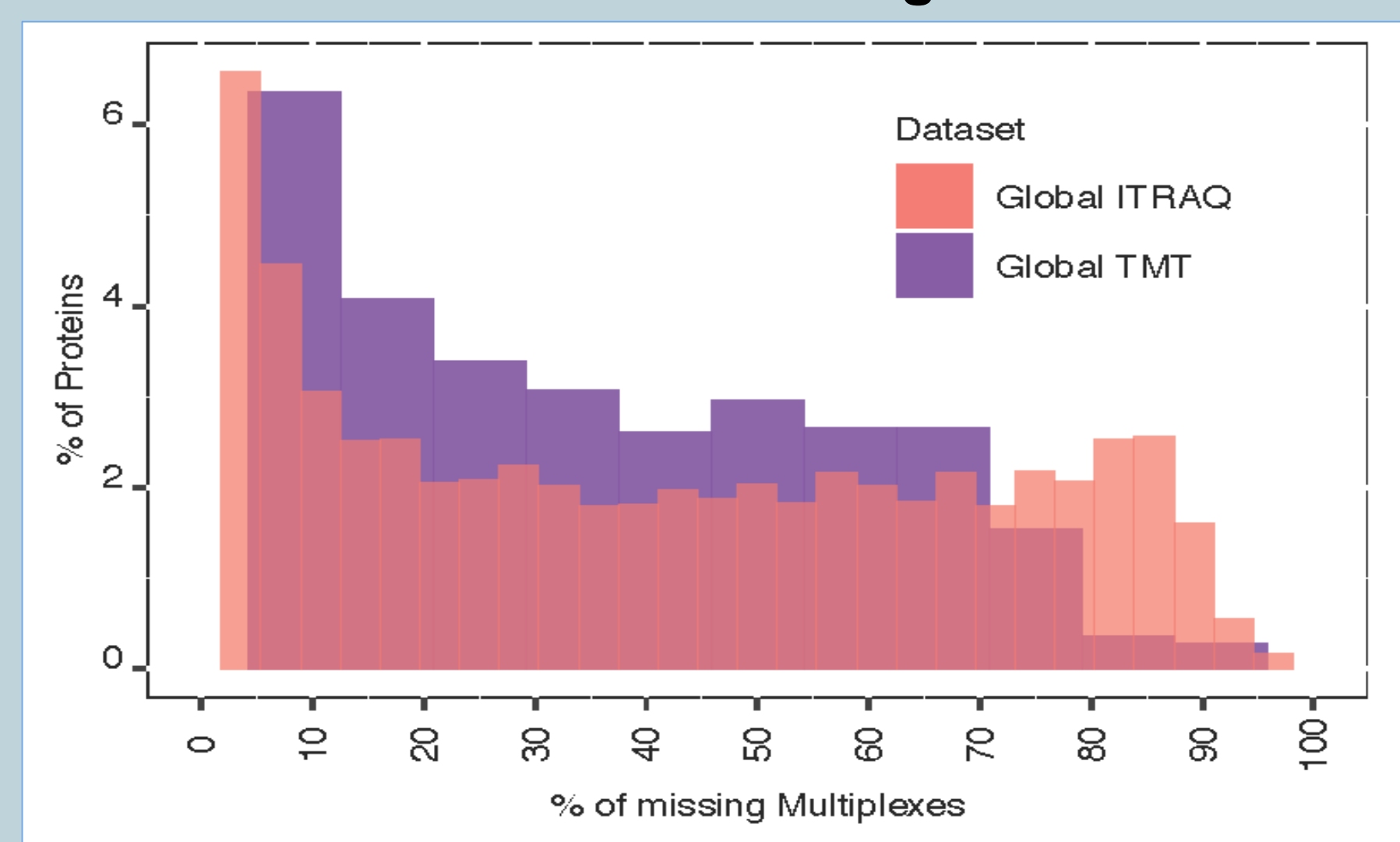
Weiping Ma<sup>1\*</sup>, Sunkyu Kim<sup>2\*</sup>, Shrabanti Chowdhury<sup>1</sup>, Zhi Li<sup>3</sup>, Mi Yang<sup>4</sup>, Seungyeul Yoo<sup>1</sup>, Francesca Petralia<sup>1</sup>, Jeremy Jacobsen<sup>5</sup>, Jingyi Jessica Li<sup>6</sup>, Xinzhou Ge<sup>6</sup>, Kexin Li<sup>7</sup>, Thomas Yu<sup>8</sup>, Nathan Edwards<sup>9</sup>, Samuel Payne<sup>10</sup>, Paul C. Boutros<sup>11</sup>, Henry Rodriguez<sup>12</sup>, Gustavo Stolovitzky<sup>13</sup>, Jun Zhu<sup>1</sup>, Jaewoo Kang<sup>2</sup>, David Fenyo<sup>3</sup>, Julio Saez-Rodriguez<sup>5</sup>, Pei Wang<sup>1#</sup>

<sup>1</sup>Icahn School of Medicine at Mount Sinai (USA), <sup>2</sup>Department of Computer Science and Engineering, Korea University (South Korea), <sup>3</sup>New York University (USA), <sup>4</sup>Faculty of Biosciences, Heidelberg University (Germany), <sup>5</sup>University of Colorado (USA), <sup>6</sup>Department of Statistics, University of California Los Angeles (USA), <sup>7</sup>Department of Mathematics, Tsinghua University (China), <sup>8</sup>Sage Bionetworks (USA), <sup>9</sup>Georgetown University (USA), <sup>10</sup>Pacific Northwest National Laboratory (USA), <sup>11</sup>University of California, Los Angeles (USA), <sup>12</sup>National Cancer Institute (USA), <sup>13</sup>IBM Research & Mount Sinai (USA), <sup>#</sup>Corresponding author

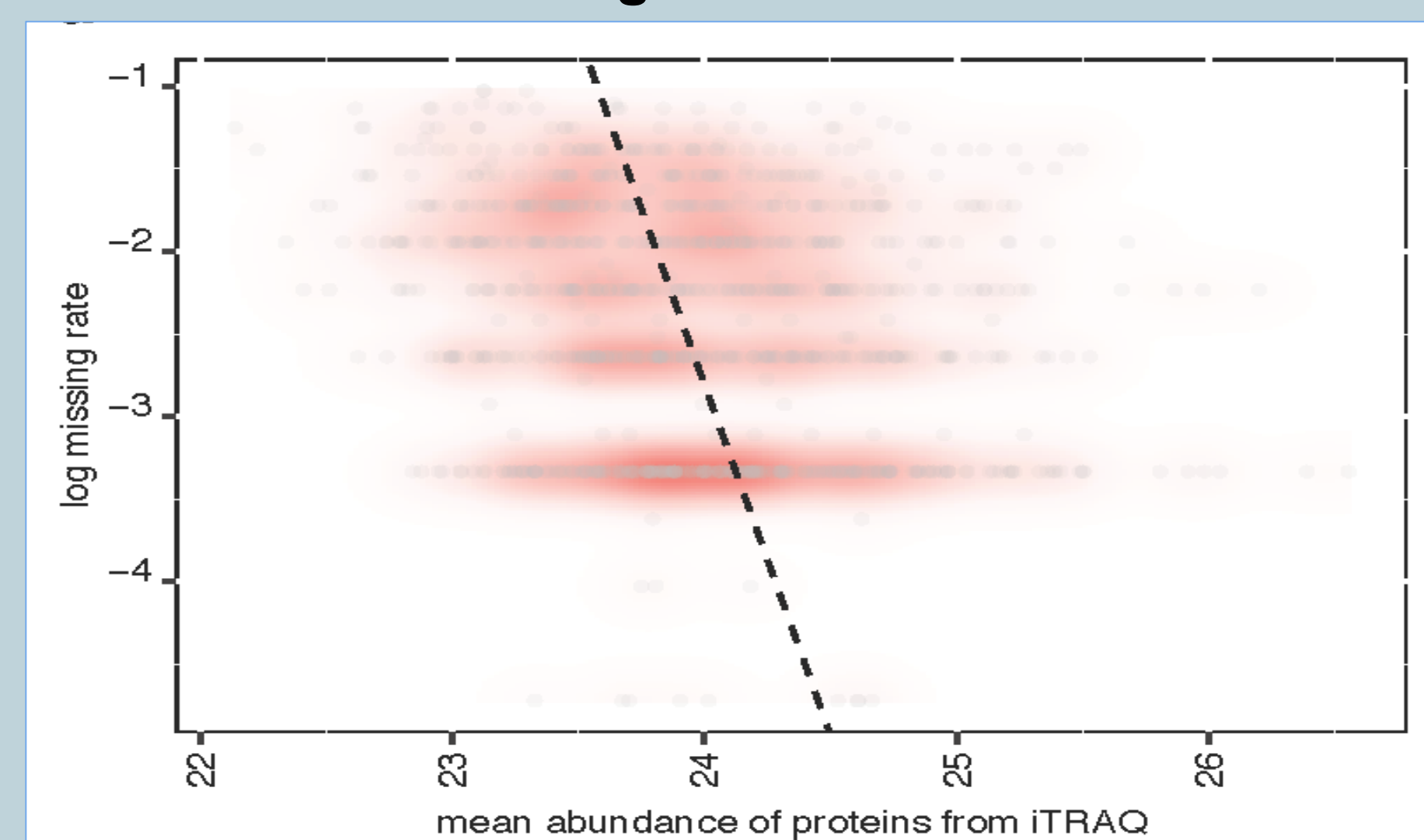
## Introduction

Deep proteomics profiling using labelled LC-MS/MS experiments has been proven to be powerful to study complex diseases. However, due to the dynamic nature of the discovery mass spectrometry, the generated data contain a substantial fraction of missing values. This poses great challenges for data analyses, as many tools, especially those for high dimensional data, cannot deal with missing values directly. To address this problem, the NCI-CPTAC Proteogenomics DREAM Challenge was carried out to develop effective imputation algorithms for labelled LC-MS/MS proteomics data through crowd learning. The final resulting algorithm, **DreamAI**, is based on an ensemble of six different imputation methods. This new tool will nicely enhance data analysis capabilities in proteomics research.

### Substantial missing rate

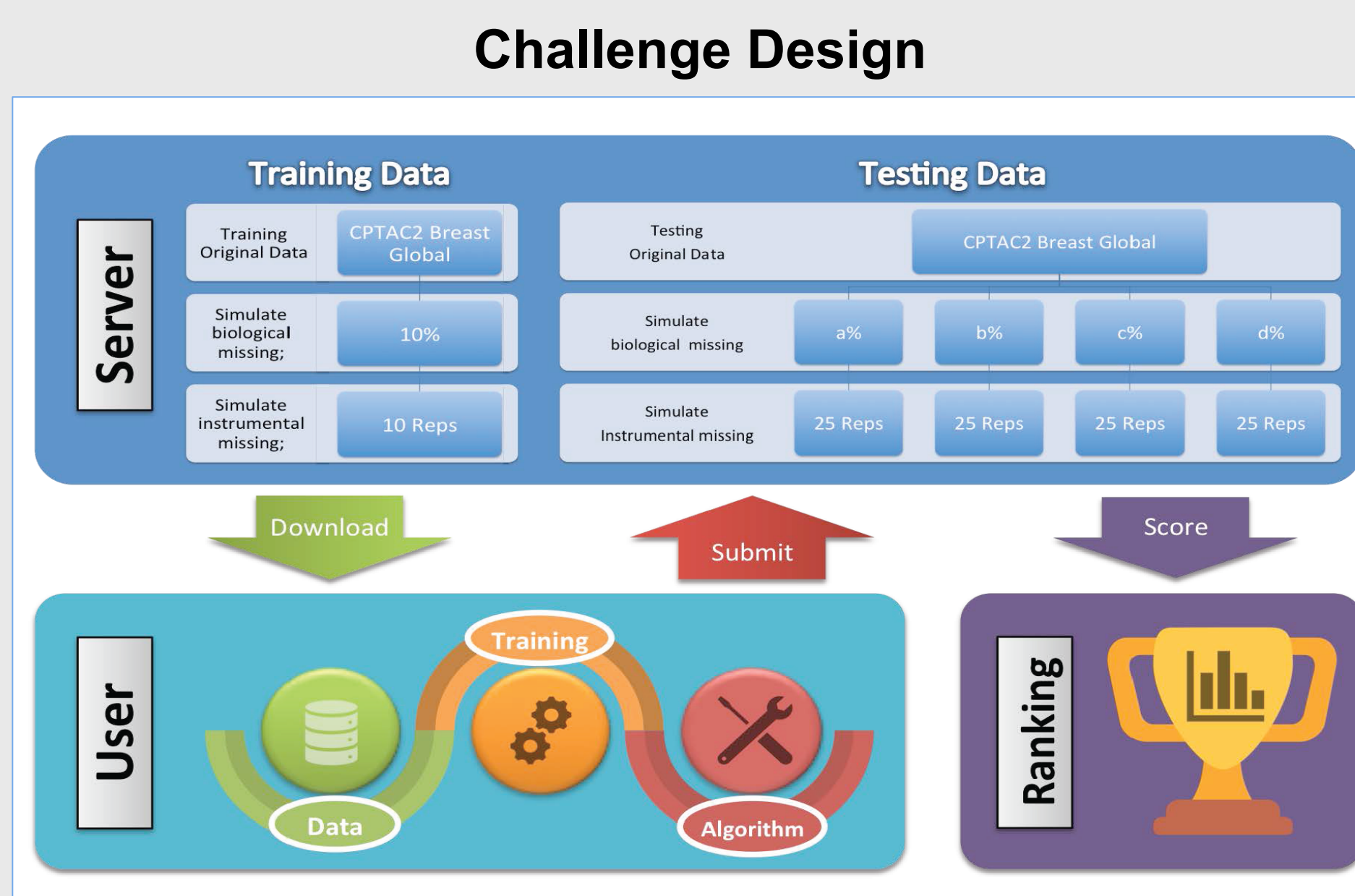


### Missing not at random

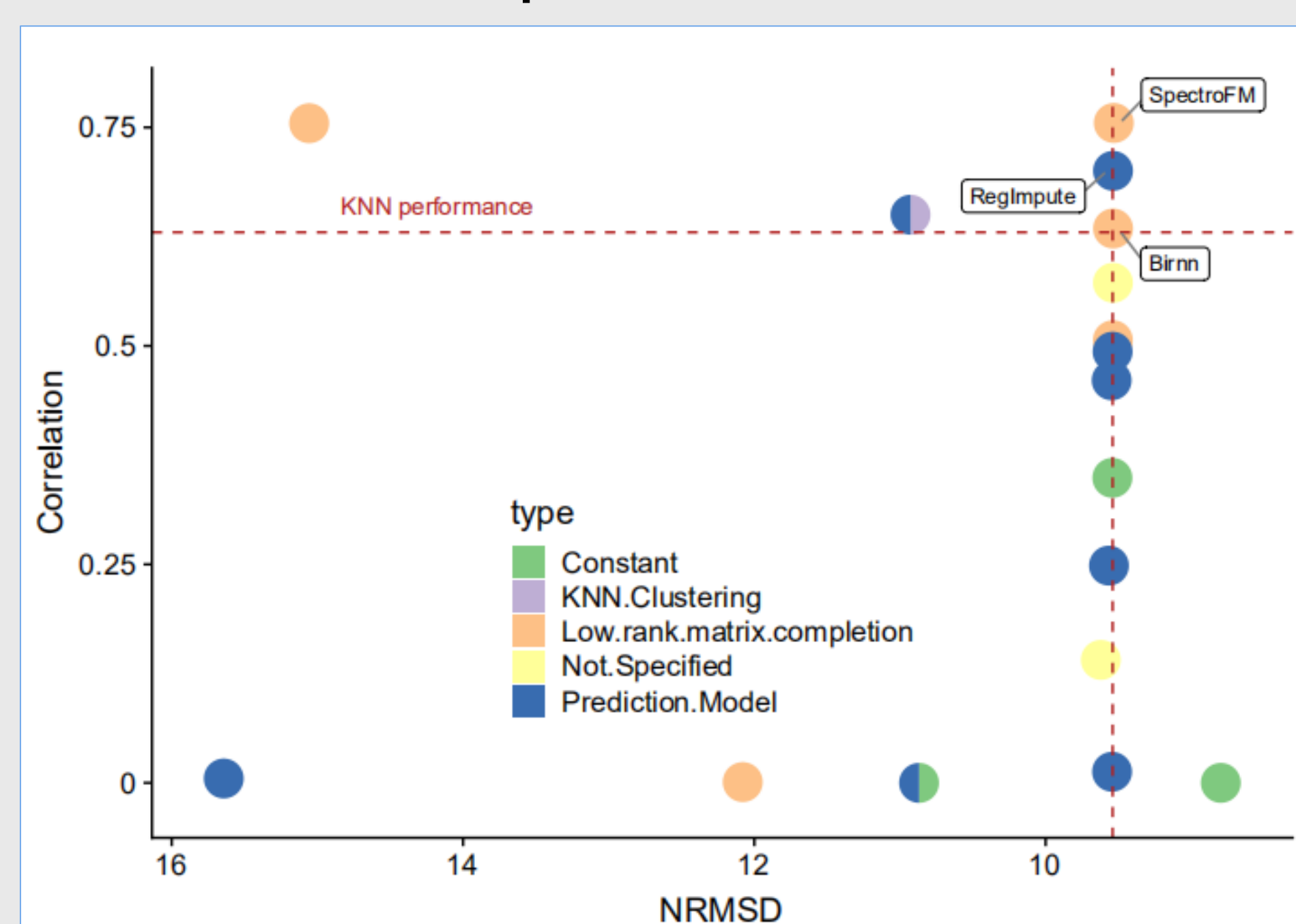


## NCI-CPTAC Dream Challenge: imputation of proteomics data

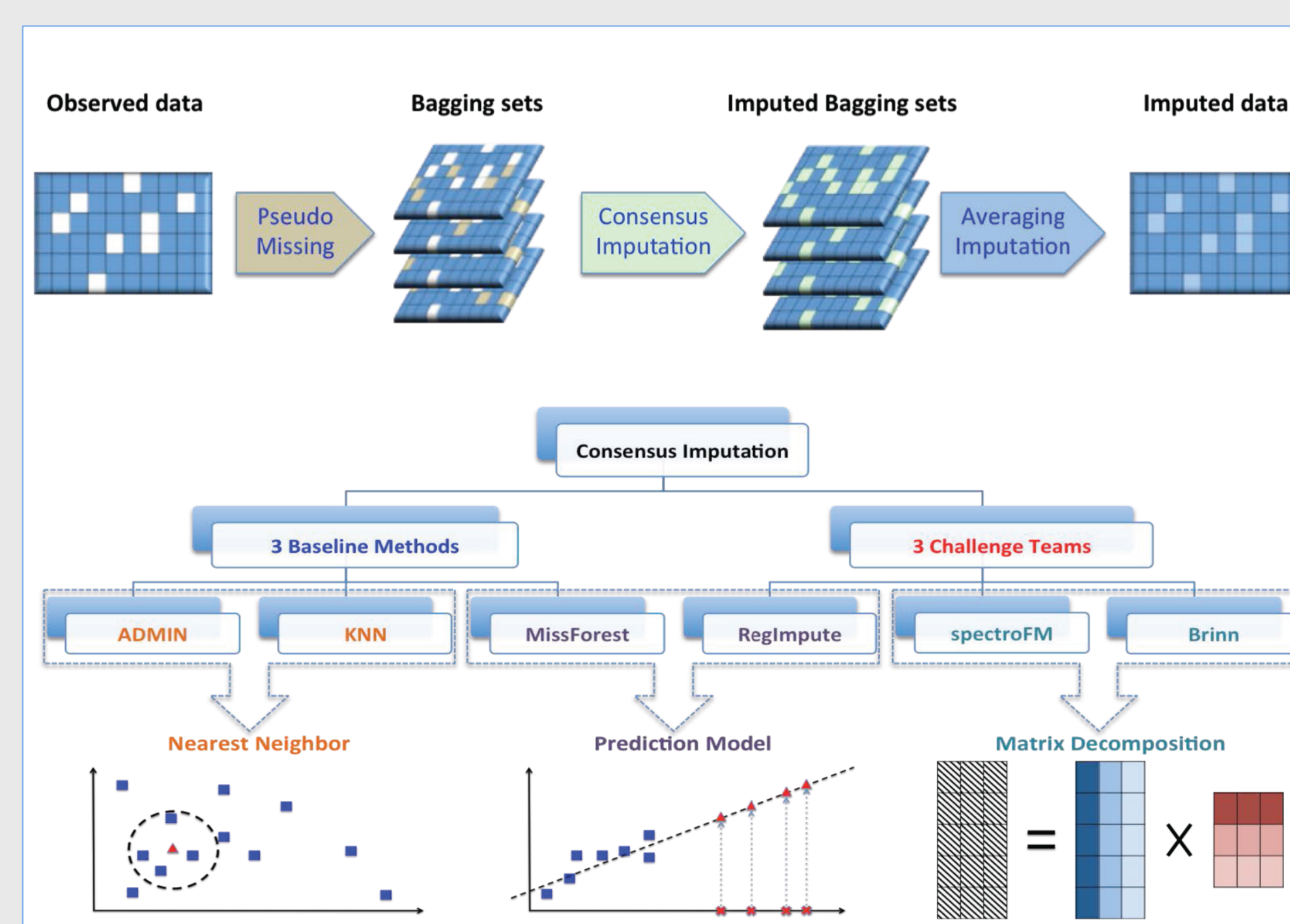
The Challenge included a competition phase and a collaborative phase. In the competition phase, participants were invited to submit imputation algorithms trained on labelled LC-MS/MS proteomics data sets, and the performances of these algorithms were evaluated on a collection of test datasets generated from the CPTAC BRCA data. In the collaborative phase, together with the three winning teams from the competition phase, we further enhanced and integrated different imputation techniques and developed the final Aggregation based Imputation algorithm --- DreamAI



### Participants Performance



### Structure of Dream AI



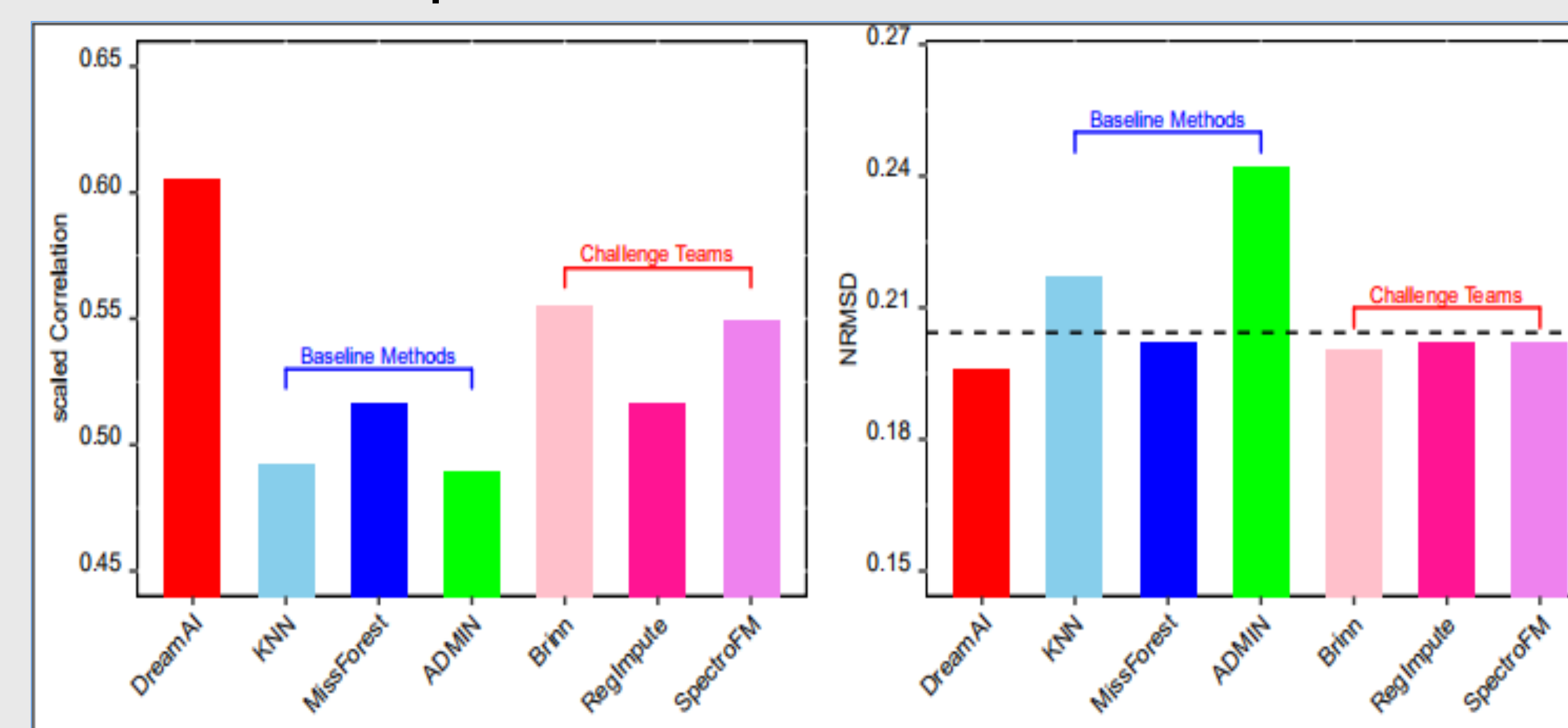
## Performance of DreamAI

DreamAI has better performance than the 6 individual algorithms on imputation of CPTAC2 OV proteomics data set (performed on samples produced by PNNL, and evaluated with the same replicated samples produced by JHU). Evaluation has been assessed by NRMSE and Correlation.

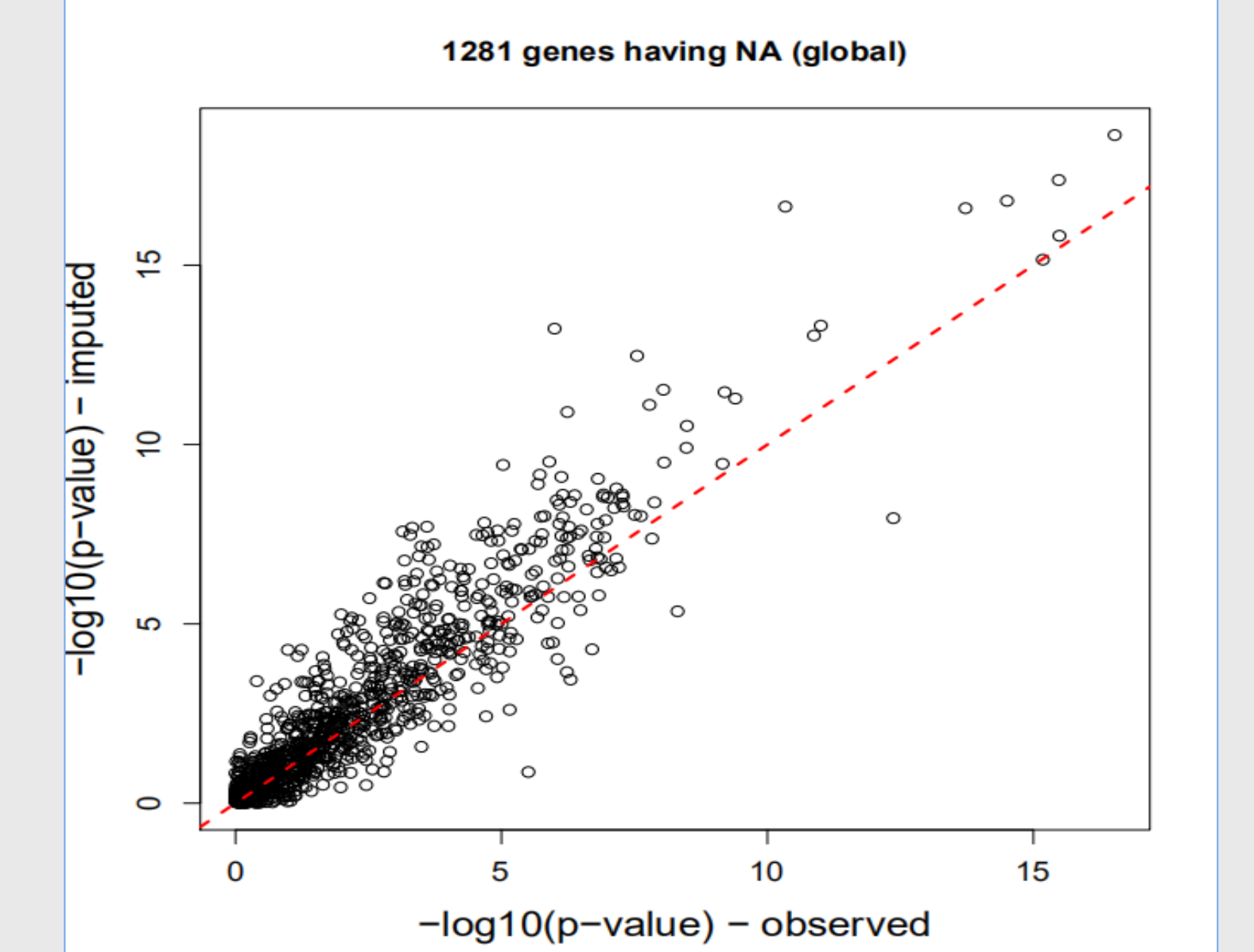
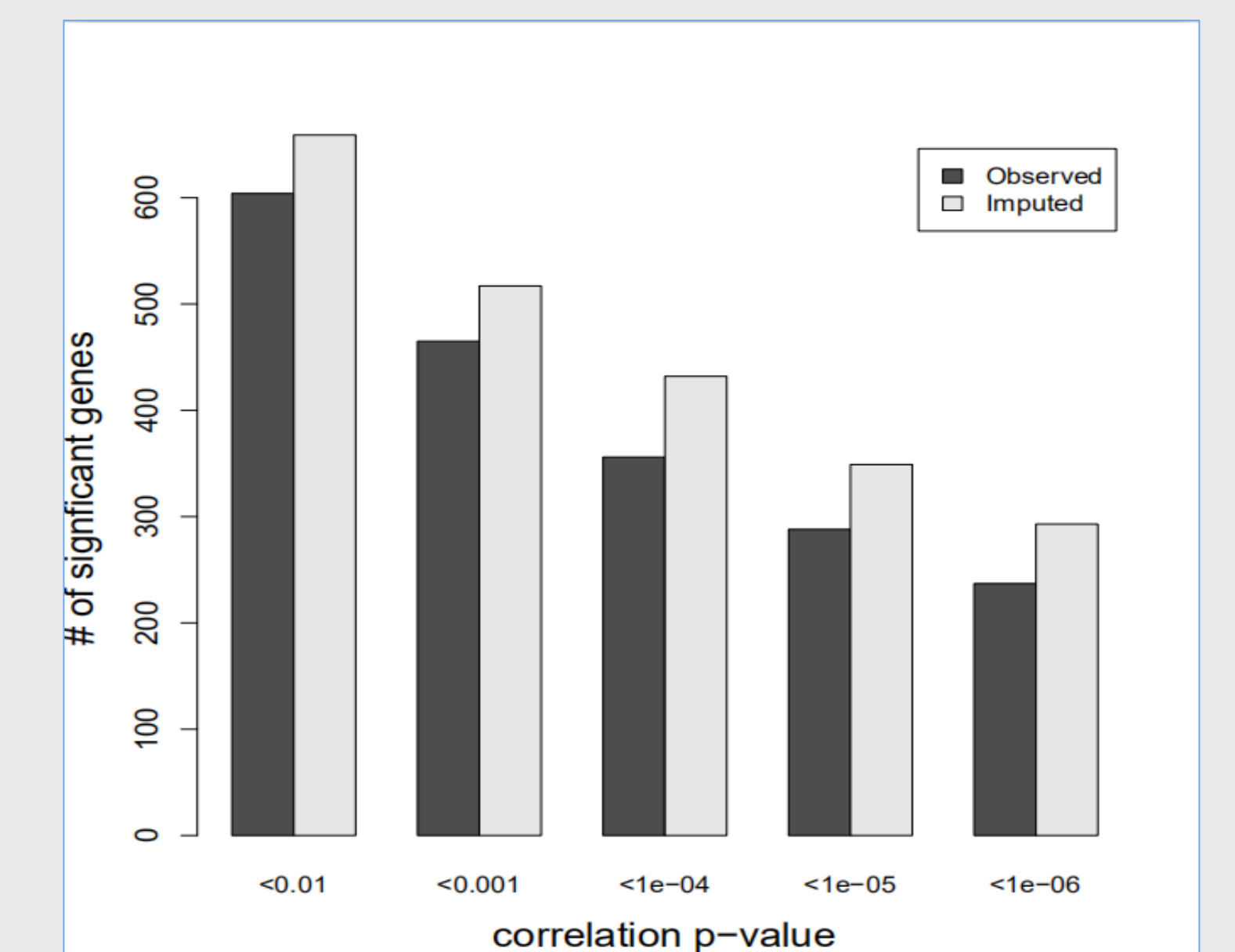
$$NRMSE = \frac{\sqrt{\sum_{i=1}^{n_{missing}} (y_i - x_i)^2 / n_{missing}}}{y_{max} - y_{min}} \quad r = \frac{1}{n_{missing} - 1} \sum_{i=1}^{n_{missing}} \frac{(x_i - \bar{X})(y_i - \bar{Y})}{S_x S_y}$$

DreamAI Imputation increased protein-rna concordance at gene level in CPTAC3 CCRCC proteomics data set. Concordance was estimated by the pairwise correlation between proteomics and gene expression

### DreamAI out performed all single algorithm components in proteomics data of CPTAC2 OV cohort



### DreamAI Imputation increases gene level Protein-RNA concordance in CPTAC3 CCRCC cohort

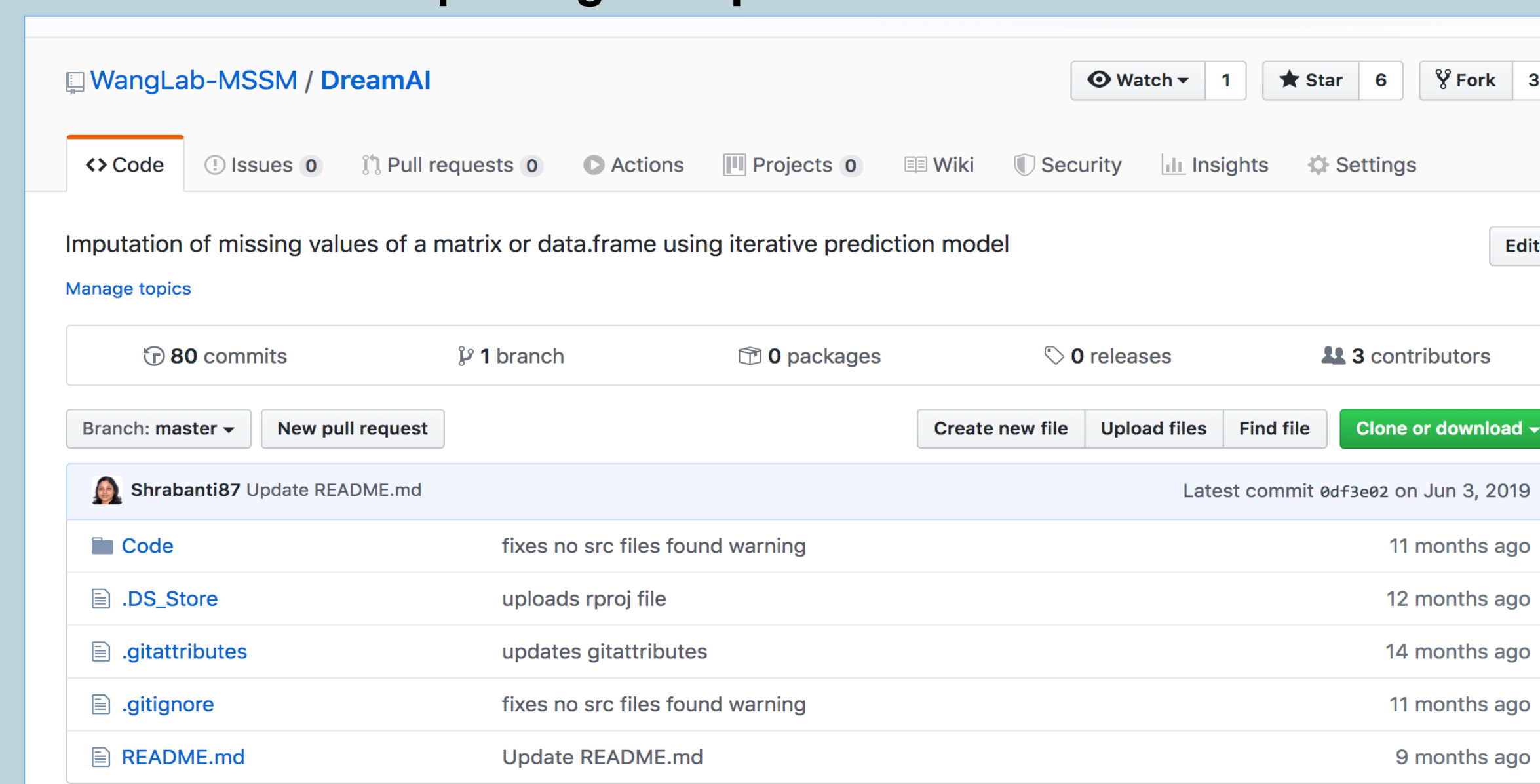


## Acknowledgement and References

R Package of the algorithm is publicly available through Github: <https://github.com/WangLab-MSSM/DreamAI>

Paper draft is also available on biorxiv: <https://www.biorxiv.org/content/10.1101/2020.07.21.214205v1>

### The DreamAI R package is open source and available on Github



- Mertins, Philipp, et al. "Proteogenomics connects somatic mutations to signalling in breast cancer." *Nature* 534, no. 7605 (2016): 55.
- Zhang, Hui, et al. "Integrated proteogenomic characterization of human high-grade serous ovarian cancer." *Cell* 166, no. 3 (2016): 755-765.
- CPTAC (NCI/NIH). "CPTAC Ovarian Cancer Confirmatory Study." Distributed by NCI Proteomic Data Commons. <https://cptac-data-portal.georgetown.edu/cptac/s/S038>
- CPTAC (NCI/NIH). "CPTAC Breast Cancer Confirmatory Study." Distributed by NCI Proteomic Data Commons. <https://cptac-data-portal.georgetown.edu/cptac/s/S039>
- Clark, David J. et al. "Integrated proteogenomic characterization of clear cell renal cell carcinoma." *Cell* 179, no. 4 (2019): 964-983.
- Hastie, Trevor, Robert Tibshirani, Gavin Sherlock, Michael Eisen, Patrick Brown, and David Botstein. "Imputing missing data for gene expression arrays." (1999).
- Stekhoven, Daniel J., and Peter Bühlmann. "MissForest—non-parametric missing value imputation for mixed-type data." *Bioinformatics* 28, no. 1 (2011): 112-118.
- Wang, Minghui, et al. "The Mount Sinai cohort of large-scale genomic, transcriptomic and proteomic data in Alzheimer's disease." *Scientific data* 5 (2018): 180185.