# Profiling performance of high-resolution HLA imputation with complementary metrics and visualizations using an open-source validation framework

# Alyssa N. Paynter<sup>1</sup>, Malek Kamoun<sup>2</sup>, Nicholas K. Brown<sup>2</sup>, Ryan J. Urbanowicz<sup>3</sup>, Keith McCullough<sup>4</sup>, Gerald Morris<sup>5</sup>, Martin Maiers<sup>6</sup>, Massimo Mangiola<sup>7</sup>, Brendan

Keating<sup>8</sup>, Bonnie Lonze<sup>8</sup>, Michal Mankowski<sup>8</sup>, Loren Gragert<sup>1</sup>

 Division of Biomedical Informatics and Genomics, Deming Department of Medicine, Tulane University School of Medicine, New Orleans, LA, United States. 2. Department of Pathology and Laboratory Medicine, University of Pennsylvania, Philadelphia, PA, United States. 3. Department of Computational Biomedicine, Cedars Sinai Medical Center, Los Angeles, CA, United States. 4. Arbor Research Collaborative for Health, Ann Arbor, MI, United States. 5. Department of Pathology, UC San Diego Health, San Diego, CA, United States. 6. NMDP, Minneapolis, MN, United States. 7. Department of Pathology, NYU Grossman School of Medicine, NYU Langone Health, New York, NY, United States. 8. Department of Surgery, NYU Grossman School of Medicine, NYU Langone Health, New York, NY, United States. 8. Department of Surgery, NYU Grossman School of Medicine, NYU Langone Health, New York, NY, United States.

INTRODUCTION	RESULTS	DISCUSSION
<ul> <li>Research studies evaluating high-resolution (HR) HLA imputation in transplantation have varied substantially in their approaches to measuring prediction performance.</li> <li>We aimed to develop a package for computing comprehensive imputation performance metrics and visualizations for predictions of HLA alleles, amino</li> </ul>	Calibration Plot and Prediction Probability Distribution in 100 pairs for DQ unique eplet MM Brier: 0.1, Bin Avg MSE: 0.004, Bin Avg City-Block Dist: 0.0317 Mean Predicted Probability for Quantile 0.0 0.2 0.4 0.6 0.8 1.0 1.0 Ideal Calibration True Fraction vs Probability Average for Quantile	Input Required for Framework         DQ Antigen-level Typing         DQ2 + DQ6         Possible DQ Genotypes and Probabilities

### acids, and molecular mismatch categories for donorrecipient pairs with ambiguous typing.

# METHODS

- We developed a Python framework for imputation validation where model performance is evaluated using built-in functions from the *scikit-learn* package, a best-in-class machine learning framework.
- We tested the framework on deceased kidney donors (N=217) typed at antigen level from the national registry, and then had HR typing preformed.

		Description
	Brier Score	Measures the distance of individual- level prediction probabilities and correctness for the dataset from an ideal predictor with 100% correct predictions each with 100% probability. (lower scores are better)
olot	Bin Average Mean	Used for fitting regression lines, are computed based on squared distance from the perfect diagonal and



**Figure 1.** Calibration curves compare prediction probability averages for quantiles with the true fraction. Red points close to diagonal indicate excellent calibration for DQ unique eplet mismatch count predictions.

DQB1\*02:01 + DQB1\*06:01, 0.7 2. DQB1\*02:01 + DQB1\*06:02, 0.2 3. DQB1\*02:02 + DQB1\*06:02, 0.05



- Imputation distribution and truth table are input
- Mismatched metrics derived from high resolution predictions on transplant pairs.
- When HLA imputation probabilities are shown to be well calibrated, uncertainty in HLA allele assignments can be incorporated into statistical models for HLA association studies.

# **Different Levels of Analysis**

<b>Mutli-loci Unphased</b>	Single Locus Unphased
Genotype:	Genotype:
DRB1 + DQA1 + DQB1	DRB1 + DRB1
<b>Eplet Mismatch:</b>	Amino Acid Mismatch:
Unique eplet MMs	Unique AAMMs
Counts of eplet MM	Counts of AAMMs

alibration  $\bigcirc$ the Within

from the perfect diagonal and penalize poor calibration more heavily than city-block distance. (lower Error (MSE) scores are better)

> Measures the distance between fraction of correct predictions and the quantile-averaged prediction probabilities, with perfect calibration falling on the dashed-line diagonal. (lower scores are better)

## Histogram OŤ imputation probabilities

Squared

Average

Distance

**City Block** 

Bin

Shows the full prediction probability distribution.

**Table for** the Quantiles Imputation probabilities turned into quantiles (true fraction vs probability average) and the table gives summary statistics of each one.

Receiver Operating Characteristic (ROC)-Area

Plots the true positive rate against the false positive rate for the most probable imputation prediction and quantify performance using





Classifier (AUC = 0.78)

0.8

1.0

--- Chance level (AUC = 0.5)

0.6

False Positive Rate

Figure 2. ROC curves that are towards the upper left indicate

better performance, however ROC curves are more useful



Risk categories [2]

• The framework is available on GitHub at https:// github.com/lgragert/imputation-validation/. Python PyPi package forthcoming.

# CONCLUSION

- Calibration plots provide advantages over using ROC metrics and simple accuracy measures, especially in datasets with class imbalance.
- This framework could support high quality HLA imputation studies across immunogenetics and transplant settings.
- Planned features include typing resolution score visualization at amino acid level and support for more mismatch metrics to more clearly separate the concepts of uncertainty and inaccuracy which are often conflated [1].

#### **References:**

1.Engen, R. M., Jedraszko, A. M., Conciatori, M. A. & Tambur, A. R. Substituting imputation of HLA



with a more even distribution of prediction probabilities. Most predictions were near to 100%, causing class imbalance.

0.4

antigens for high-resolution HLA typing: Evaluation of a multiethnic population and implications for clinical decision making in transplantation. American Journal of Transplantation21, 344-352 (2021). 2.Dasariraju, S. et al. HLA amino acid Mismatch-Based risk stratification of kidney allograft failure using a novel Machine learning algorithm. Journal of Biomedical Informatics142, 104374 (2023).



C-AUC



0.0

0.2



Created with BioRender Poster Builder