

1 **A multi-task convolutional deep learning method for HLA allelic imputation**
2 **and its application to trans-ethnic MHC fine-mapping of type 1 diabetes.**

3

4 Tatsuhiko Naito^{1,2}, Ken Suzuki¹, Jun Hirata^{1,3}, Yoichiro Kamatani⁴, Koichi Matsuda⁵, Tatsushi
5 Toda², Yukinori Okada^{1,6,7*}.

6

7 1) Department of Statistical Genetics, Osaka University Graduate School of Medicine,
8 565-0871, Suita, Japan.

9 2) Department of Neurology, Graduate School of Medicine, The University of Tokyo, 113-8655,
10 Tokyo, Japan.

11 3) Pharmaceutical Discovery Research Laboratories, Teijin Pharma Limited, 191-8512, Hino,
12 Japan

13 4) Laboratory of Complex Trait Genomics, Department of Computational Biology and Medical
14 Sciences, Graduate School of Frontier Sciences, The University of Tokyo, 108-8639, Tokyo,
15 Japan

16 5) Laboratory of Clinical Genome Sequencing, Department of Computational Biology and
17 Medical Sciences, Graduate School of Frontier Sciences, The University of Tokyo,
18 108-8639, Tokyo, Japan.

19 6) Laboratory of Statistical Immunology, Immunology Frontier Research Center (WPI-IFReC),
20 Osaka University, 565-0871, Suita, Japan.

21 7) Integrated Frontier Research for Medical Science Division, Institute for Open and
22 Transdisciplinary Research Initiatives, Osaka University, 565-0871, Suita, Japan.

23

24 * Corresponding author:

25 Yukinori Okada, MD, PhD

26 Address: Department of Statistical Genetics, Osaka University Graduate School of Medicine,

27 2-2 Yamadaoka, Suita, Osaka 565-0871, Japan.

28 Tel: +81-6-6879-3971

29 E-mail: yokada@sg.med.osaka-u.ac.jp

30

31 **Abstract**

32 Conventional HLA imputation methods drop their performance for infrequent alleles, which is
33 one of the factors that reduce the reliability of trans-ethnic MHC fine-mapping due to inter-ethnic
34 heterogeneity in allele frequency spectra. We developed DEEP*HLA, a deep learning method
35 for imputing HLA genotypes. Through validation using the Japanese and European HLA
36 reference panels ($n = 1,118$ and $5,122$), DEEP*HLA achieved the highest accuracies with
37 significant superiority for low-frequency and rare alleles. DEEP*HLA was less dependent on
38 distance-dependent linkage disequilibrium decay of the target alleles and might capture the
39 complicated region-wide information. We applied DEEP*HLA to type 1 diabetes GWAS data
40 from BioBank Japan ($n = 62,387$) and UK Biobank ($n = 354,459$), and successfully disentangled
41 independently associated class I and II HLA variants with shared risk among diverse
42 populations (the top signal at amino acid position 71 of HLA-DR β 1; $P = 7.5 \times 10^{-120}$). Our study
43 illustrates a value of deep learning in genotype imputation and trans-ethnic MHC fine-mapping.

44

45 **Introduction**

46 Genetic variants of the major histocompatibility complex (MHC) region at 6p21.3 contribute to
47 the genetics of a wide range of human complex traits.¹ Among the genes densely present in the
48 MHC region, human leukocyte antigen (HLA) genes are considered to explain most of the
49 genetic risk of MHC.¹ Strategies for direct typing of HLA alleles, including sequence specific
50 oligonucleotide (SSO) hybridization, Sanger sequencing, and next-generation sequencing
51 (NGS), cannot be easily scaled up for large cohorts since they are labor-intensive,
52 time-consuming, expensive, and limited in terms of allele resolution and HLA gene coverage.^{2,3}
53 As a result, in many cases, the genotypes of HLA allele are indirectly imputed from single
54 nucleotide variant (SNV)-level data using population-specific HLA reference panels.³⁻⁶ Although
55 a high-throughput alternative is HLA type inference from whole-genome sequencing data,^{7,8}
56 HLA imputation is still widely performed for existing single nucleotide polymorphism (SNP)
57 genotyping data.

58 The MHC region harbors unusually complex sequence variations and haplotypes that
59 are specific to individual ancestral populations; thus, the distribution and frequency of the HLA
60 alleles are highly variable across different ethnic groups.^{1,9} This results in heterogeneity in
61 reported HLA risk alleles of human complex diseases across diverse populations.¹⁰ For
62 instance, in type 1 diabetes (T1D), the strong association between non-Asp57 in HLA-DQβ1
63 and T1D risk has been found in European populations^{11,12} but not in the Japanese populations,
64 where the T1D susceptible HLA-DQβ1 alleles carry Asp57.¹³ Although the elucidation of risk
65 alleles across ethnicities would contribute to further understanding of the genetic architecture of
66 the MHC region associated with the pathologies of complex diseases, limited trans-ethnic MHC
67 fine-mappings have been reported to date.¹⁴ One method for conducting trans-ethnic

68 fine-mapping in the comprehensive MHC region is to newly construct a large HLA reference
69 panel that captures the complexities of the MHC region across different populations.¹⁵ Another
70 method is to integrate data of different populations that are imputed with population-specific
71 reference panels. The latter approach appears straightforward but requires an HLA imputation
72 method accurate enough for infrequent alleles to allow robust evaluation of HLA variants that
73 show highly heterogeneous in allele frequency across ethnicities.

74 Starting with a simple inference using tag SNPs,^{16,17} various methods have been
75 developed for HLA allelic imputation. Leslie et al. first reported a probabilistic approach to
76 classical HLA allelic imputation.¹⁸ HLA*IMP uses Li & Stephens haplotype model with SNP
77 data from European populations.^{19,20} A subsequently developed software program,
78 HLA*IMP:02, uses SNP data from multiple populations and can address genotypic
79 heterogeneity.²¹ The current version of HLA*IMP:02 does not provide a function for users to
80 generate an imputation model using their own reference data locally. SNP2HLA is another
81 standard software, which uses the imputation software package Beagle to impute both HLA
82 alleles and the amino acid polymorphisms for those classical alleles.²² HLA Genotype
83 Imputation with Attribute Bagging (HIBAG)²³ is also promising software, which employs multiple
84 expectation-maximization-based classifiers to estimate the likelihood of HLA alleles. Whereas
85 SNP2HLA explicitly uses reference haplotype data, for which public access is often limited,
86 HIBAG does not require these data once the trained models are generated. These methods
87 have achieved high imputation accuracy;²⁴ however, they are less accurate for rare alleles as
88 shown later. The complex linkage disequilibrium (LD) structures specific for the MHC region
89 requires a more sophisticated pattern recognition algorithm beyond simple stochastic inference.

90 After boasting of its extremely high accuracy in image recognition, deep learning has
91 been attracting attention in various fields. It can learn a representation of input data and extract
92 relevant features of high complexity through deep neural networks. Many successful
93 applications in the field of genomics have been reported.²⁵ A typical application of deep learning
94 for genomics is the prediction of the effects of non-coding and coding variants, where models
95 encode the inputs of flanking nucleotide sequence data.²⁶⁻²⁹ Another application is non-linear
96 unsupervised learning of high-dimensional quantitative data from transcriptome.^{30,31} However,
97 successful representation learning for SNV-data in the field of population genetics is limited.³²
98 Here, we developed DEEP*HLA, a multi-task convolutional deep learning method to accurately
99 impute genotypes of HLA genes from SNV-level data. Through the application to the two HLA
100 reference panels of different populations, DEEP*HLA achieved higher imputation accuracy than
101 conventional methods. Notably, DEEP*HLA was advantageous especially for imputing
102 low-frequency and rare alleles. Furthermore, DEEP*HLA showed by far the fastest total
103 processing time, which suggests its applicability to biobank-scale data. We applied the trained
104 models of DEEP*HLA to the large-scale T1D genome-wide association study (GWAS) data
105 from BioBank Japan (BBJ) and UK Biobank (UKB) and conducted trans-ethnic fine-mapping in
106 the MHC region.

107

108 **Results**

109

110 **Overview of the study**

111 An overview of our study is presented in **Fig. 1**. Our method, DEEP*HLA, is convolutional
112 neural networks that learn from an HLA referenced panel and impute genotypes of HLA genes
113 from pre-phased SNV data. Its framework uses a multi-task learning that can learn and impute
114 alleles of several HLA genes which belong to the same group simultaneously (see Methods).
115 Multi-task learning is presumed to have two advantages in this situation. First, the genotypes of
116 some flanking HLA genes, which often show strong LD for each other, are correlated, and the
117 shared features of individual tasks are likely to be informative. Second, the processing time is
118 reduced by grouping tasks, especially in our latest reference panel, which comprises more than
119 30 HLA genes. We employed the two different HLA imputation reference panels for robust
120 benchmarking: (i) our Japanese reference panel ($n = 1,118$)³ and (ii) the Type 1 Diabetes
121 Genetics Consortium (T1DGC) reference panel ($n = 5,122$).³³ We compared its performance
122 with that of other HLA imputation methods by 10-fold cross-validation and an independent HLA
123 dataset ($n = 908$).⁶ Further, we tested its imputation accuracy for multi-ethnic individuals using
124 data from the Phase III 1000 Genomes Project (1KGv3). In the latter part, we performed MHC
125 fine-mapping of the Japanese cohort from BBJ and British cohort from UKB by applying trained
126 models specific for individual populations. We integrated the imputed GWAS genotypes and
127 performed trans-ethnic HLA association analysis.

128

129 **DEEP*HLA achieved high imputation accuracy especially for low-frequency and rare**
130 **alleles**

131 First, we applied DEEP*HLA to the Japanese reference panel, a high-resolution allele catalog
132 of NGS-based HLA typing data of the 33 classical and non-classical HLA genes along with
133 high-density SNP data of the MHC region by genotyping with the Illumina HumanCoreExome
134 BeadChip for 1,118 individuals of Japanese ancestry.³ We compared the imputation accuracy
135 of DEEP*HLA in terms of sensitivity, positive predictive value (PPV), and r^2 of allelic dosage,
136 and concordance rate of best-guess genotypes (see **Methods**) with those of SNP2HLA and
137 HIBAG in 10-fold cross-validation. DEEP*HLA achieved total sensitivity of 0.987, PPV of 0.986,
138 r^2 of 0.984, and concordance rate of 0.988 in 4-digit allelic resolution. The differences in total
139 accuracy were modest among the methods; however, DEEP*HLA was more advantageous for
140 rare alleles (For alleles with a frequency < 1%, sensitivity = 0.690; PPV = 0.799; r^2 = 0.911; and
141 concordance rate = 0.691 in DEEP*HLA, compared to sensitivity = 0.628, 0.635; PPV = 0.624,
142 0.505; r^2 = 0.862, 0.792; and concordance rate = 0.621, 0.675 in SNP2HLA and HIBAG,
143 respectively; **Fig. 2a**). Further, we applied the model trained with our Japanese reference panel
144 to a dataset of 908 Japanese individuals to investigate whether DEEP*HLA could impute well
145 when applied to independent samples. The dataset comprised 4-digit alleles of 8 classical HLA
146 genes based on the SSO method and SNP data genotyped using multiple genotyping arrays.⁶
147 DEEP*HLA achieved the highest total accuracy, with a sensitivity of 0.973, PPV of 0.972, r^2 of
148 0.986, and concordance rate of 0.973. Again, DEEP*HLA was more advantageous for
149 low-frequency and rare alleles (**Fig. 2a**). For alleles with a frequency < 1%, sensitivity = 0.690;
150 PPV = 0.799; r^2 = 0.911; and concordance rate = 0.691 in DEEP*HLA, compared to sensitivity

151 = 0.628, 0.635; PPV = 0.624, 0.505; r^2 = 0.862, 0.792; and concordance rate = 0.621, 0.675 in
152 SNP2HLA and HIBAG, respectively.

153 We also applied DEEP*HLA to the Type 1 Diabetes Genetics Consortium (T1DGC)
154 reference panel of 5,122 unrelated individuals of European ancestries.³³ It comprises 2- and
155 4-digit alleles of the 8 classical HLA genes based on the SSO method, with SNP data of the
156 MHC region genotyped with the Illumina ImmunoChip. DEEP*HLA achieved a sensitivity of
157 0.979, PPV of 0.976, r^2 of 0.981, and concordance rate of 0.979 in 4-digit resolution, and these
158 values were superior to those of SNP2HLA and HIBAG. DEEP*HLA was more advantageous
159 especially in PPV and r^2 , for low-frequency and rare alleles (**Fig. 2b**). For alleles with a
160 frequency < 1%, sensitivity = 0.830; PPV = 0.863; r^2 = 0.908; and concordance rate = 0.832 in
161 DEEP*HLA, compared to sensitivity = 0.793, 0.745; PPV = 0.640, 0.753; r^2 = 0.745, 0.886; and
162 concordance rate = 0.804, 0.769 in SNP2HLA and HIBAG, respectively.

163 We assessed the superiority of DEEP*HLA using a down-sampling approach
164 (**Supplementary Note 1a**). DEEP*HLA trained with down-sampled data also outperformed
165 other methods especially for rare allele, although there were differences between metrics
166 (**Supplementary Fig. 1**). In the cross-validation of our Japanese reference panel, DEEP*HLA
167 with sampling rates of 70%–80% and 60%–70% was almost equivalent to HIBAG and
168 SNP2HLA, respectively. In the Japanese independent samples, DEEP*HLA with a sampling
169 rate of even 70% and 60% outperformed HIBAG and SNP2HLA, respectively. In the
170 cross-validation of the T1DGC panel, DEEP*HLA with a sampling rates of 70%–80% was
171 almost equivalent to HIBAG and SNP2HLA, respectively. Notably, DEEP*HLA with a sampling
172 rate of even 50% outperformed other methods in most cases in terms of PPV.

173 Finally, we investigated differences in accuracy among different HLA genes (**Fig. 3**).
174 Whereas the accuracies for *HLA-B* and *HLA-DRB1* were lower than those for other loci
175 especially in terms of total accuracy, those in DEEP*HLA were relatively high. As a result,
176 DEEP*HLA had the highest means and lowest variances of accuracies among HLA genes in
177 most cases. Only for rare alleles in the Japanese independent samples, the variances of
178 sensitivity and concordance rate were higher than those for SNP2HLA, in which the accuracy
179 metrics of SNP2HLA were lower than those of DEEP*HLA for almost all loci.

180 In summary, although the improvement in total accuracy of DEEP*HLA might be
181 modest, DEEP*HLA was advantageous in imputing infrequent alleles especially in terms of the
182 dosage accuracy. PPV was significantly decreased in SNP2HLA, probably because the sum of
183 the allele dosages of each HLA gene in an individual can exceed the expected value (i.e. = 2.0)
184 since it imputes each allele separately as a binary allele. The improvement in dosage accuracy
185 is meaningful considering that allelic dosages are typically used for association analysis.³
186 Furthermore, its small inter-locus variation in imputation accuracy should also be advantageous
187 in MHC fine-mapping because the accuracy difference among HLA genes would result in
188 imbalanced filtering, leading to a biased result.

189
190 **DEEP*HLA achieved higher accuracy when applied to 1000 Genomes Project data using**
191 **a mixed reference panel**

192 To conduct further validation in independent samples and evaluate the effect of ethnicity
193 differences between a reference panel and target populations, we tested imputation accuracy in
194 1KGv3 cohort. First, we conducted HLA imputation using our Japanese panel and ($n = 1,118$)
195 and a mixed panel which was experimentally conducted using the Japanese and the European

196 panels ($n = 6,240$). When we used the Japanese panel, DEEP*HLA achieved the highest
197 accuracies in all the metrics in the 1KGV3 JPT cohort (sensitivity = 0.974, PPV = 0.950, $r^2 =$
198 0.995, and concordance rate = 0.975 in total alleles; **Supplementary Fig. 2a**). All the methods
199 achieved high accuracies for rare alleles, in which DEEP*HLA was still superior (sensitivity =
200 0.862, PPV = 0.865, $r^2 = 0.999$, and concordance rate = 0.862 for alleles with a frequency of <
201 1%). On the other hand, in other populations including EAS (excluding JPT), no methods were
202 found to be accurate enough for practical use. This is probably attributed to the distinct
203 haplotype structures and allele frequency spectra specific for Japanese ancestries even within
204 East Asian populations.⁶ In addition, DEEP*HLA did not always perform better than other
205 methods. Presumably, its high learning capacity of deep learning might backfire and caused
206 overfitting to the population-specific reference panel. We thus recommend empirical validation
207 of accuracy when applying DEEP*HLA to individuals mismatched with a reference panel
208 population.

209 When we used a mixed panel, despite a slight decline in accuracy in JPT (sensitivity =
210 0.965, PPV = 0.940, $r^2 = 0.996$, and concordance rate = 0.964 for total alleles), DEEP*HLA
211 achieved high accuracies in EUR populations (sensitivity = 0.964, PPV = 0.918, $r^2 = 0.983$, and
212 concordance rate = 0.963 for total alleles). DEEP*HLA also achieved the highest accuracies in
213 both JPT and EUR populations for total and rare alleles although the difference was relatively
214 modest (**Supplementary Fig. 2b**). Thanks to a significant increase in the sample sizes of the
215 reference panel, the accuracies in other populations were also improved. Notably, DEEP*HLA
216 achieved the highest accuracies in the different populations, especially for rare alleles. Although
217 the mixed panel used here is an experimental version that comprises genotypes from different

218 typing procedures, the present results would suggest the applicability of our method to a
219 multi-ethnic reference panel.

220 **DEEP*HLA can define HLA amino acid polymorphisms consistently with classical alleles**

221 DEEP*HLA separately imputes classical alleles of each HLA gene, as a multi-class
222 classification in the field of machine learning. Thus, it has an advantage that the sum of imputed
223 allele dosages of each HLA gene is definitely set as an ideal value of 2.0. This enables us to
224 define a dosage of amino acid polymorphisms from the imputed 4-digit allele dosages
225 consistently with classical alleles. We compared this method of imputing amino acid
226 polymorphisms with SNP2HLA, which imputes each allele as binary alleles. Although
227 DEEP*HLA was equivalent to SNP2HLA in imputing amino acid polymorphisms in total alleles
228 (sensitivity = 0.996, PPV = 0.996, $r^2 = 0.951$, and concordance rate = 0.996 in the Japanese
229 panel; sensitivity = 0.997, PPV = 0.995, $r^2 = 0.982$, and concordance rate = 0.997 in T1DGC
230 panel), it achieved more accurate imputation for rare alleles (sensitivity = 0.487, PPV = 0.811, r^2
231 = 0.665, and concordance rate = 0.487 in the Japanese panel; sensitivity = 0.775, PPV = 0.864,
232 $r^2 = 0.826$, and concordance rate = 0.775 in T1DGC panel for alleles with a frequency of < 1%;
233 **Fig. 2c, d**). The improvement in performance in terms of PPV was remarkable.

234 We admit that this method is only applicable to the reference panel where 4-digit
235 alleles are accurately determined. Therefore, our method could not eliminate the ambiguity in
236 the genotyping that derived from incompleteness of the original reference panel.

237 238 **High performance of DEEP*HLA in computational costs**

239 We benchmarked the computational costs of DEEP*HLA against those of SNP2HLA and
240 HIBAG using a subset of the GWAS dataset from BBJ containing $n = 1,000, 2,000, 5,000,$

241 10,000, 20,000, 50,000, and 100,000 samples (2,000 SNPs was consistent with the reference
242 panel). A model-training process with reference data is required for DEEP*HLA and HIBAG but
243 not for SNP2HLA. In addition, DEEP*HLA took an input of pre-phased GWAS data. Thus, we
244 compared the total processing time including pre-phasing of GWAS data, model training, and
245 imputation of DEEP*HLA, with the time of model training and imputation of HIBAG, and the
246 running time of SNP2HLA. As shown in **Fig. 2e**, DEEP*HLA imputation had by far the fastest
247 total processing time as the sample size increased. On comparing pure imputation times, it was
248 faster than HIBAG (**Supplementary Table 1**). Furthermore, with a state-of-the-art GPU, the
249 training time of DEEP*HLA was shortened from 153 min to 36 min. As for memory cost, all
250 methods exhibited maximum memory usage scaling roughly linearly with sample size (**Fig. 2e**
251 and **Supplementary Table 1**). HIBAG was the most memory-efficient across all sample sizes.
252 Whereas SNP2HLA could not run within our machine's 100 GB memory for sample sizes of
253 >20,000, DEEP*HLA was able to perform imputation even for biobank-scale sample sizes of
254 100,000.

255

256 **Characteristics of the alleles for which DEEP*HLA was advantageous to impute**

257 We focused on the characteristics of the HLA alleles of which accuracy was improved by
258 DEEP*HLA compared with SNP2HLA, which is a gold-standard software. SNP2HLA runs
259 Beagle intrinsically, which performs imputation based on a hidden Markov model of a localized
260 haplotype-cluster. We hypothesized that this kind of methods shows better performance for
261 imputing alleles for which LDs with the surrounding SNVs are stronger in close positions and
262 get weaker as the distance from the target HLA allele increases (we termed this feature as
263 distance-dependent LD decay). Conversely, it might show limited performance for imputing

264 alleles with sparse LD structures throughout the MHC region. We defined the area under the
265 curve (AUC) representing distance-dependent LD decay to verify this hypothesis. AUC values
266 increase when LDs with the surrounding SNVs get stronger as they get closer to the target HLA
267 allele (**Fig. 4b**). We evaluated the degree by which the accuracies of DEEP*HLA and SNP2HLA
268 were affected by the AUCs and allele frequency using multivariate linear regression analysis.
269 When calculating AUCs, we tested two different window sizes of AUCs: bilateral 1,000 SNPs
270 from a target HLA allele and input size of DEEP*HLA. As expected, all accuracy metrics of
271 SNP2HLA were positively correlated with the AUCs. Although the accuracy metrics of
272 DEEP*HLA were also correlated with AUC, the correlations were weaker than those in
273 SNP2HLA for all the metrics in both reference panels (**Fig. 4a** and **Supplementary Table 2**). In
274 addition, we assessed the correlation between a simple metric of the maximum value of LD
275 coefficients within 100 SNPs from a target allele and the accuracy of each method to examine
276 our assumption more robustly with another index. Similarly, the correlations in DEEP*HLA were
277 weaker than those in SNP2HLA (**Supplementary Table 2**).

278 Next, we used SmoothGrad to investigate our assumption that DEEP*HLA performs
279 better imputation by recognizing distant SNVs as well as close SNVs of strong LD. SmoothGrad
280 is a method for generating sensitivity maps of deep learning models.³⁴ It is a simple
281 approach based on the concept of adding noise to the input data and taking the mean of the
282 resulting sensitivity maps for each sampled data. A trained DEEP*HLA model reacted to the
283 noises of not only the surrounding SNVs with strong LD but also the distant SNVs as displayed
284 in example HLA alleles (**Fig. 4c**). Interestingly, SNVs that reacted strongly were not always
285 those of even moderate LD, but also spread across the entire the input region. While the validity
286 of SmoothGrad for a deep learning model of genomic data is presently under investigation, one

287 probable explanation is that predicting an allele using our method also means predicting the
288 absence of other alleles of the target HLA gene. Thus, any SNV positions in LD with any of the
289 other HLA alleles could be informative. Another explanation is that DEEP*HLA might recognize
290 complex combinations of multiple distinct SNVs within the region rather than the simple LD
291 correlations between HLA alleles and -SNVs.

292

293 **Empirical evaluation of imputation uncertainty**

294 A common issue in deep learning models is quantification of the reliability of their predictions.
295 One potential solution is uncertainty inferred from the concept of Bayesian deep learning.³⁵ We
296 experimentally evaluated imputation uncertainty by DEEP*HLA using Monte Carlo (MC)
297 dropout, which could be applied following the general implementation of neural networks with
298 dropout units.^{36,37} In MC dropout, uncertainty is presented as entropy of sampling variation with
299 keeping dropout turned on. This uncertainty index corresponds not to each binary allele of a
300 HLA gene, but to the prediction of genotype of each HLA gene of an individual. Thus, we
301 evaluated whether the uncertainty could guess the correctness of best-guess genotypes of the
302 target HLA genes. We compared this with a dosage-based discrimination, in which we assumed
303 that a best-guess imputation of higher genotype dosage (probability) is more likely to be correct.
304 The entropy-based uncertainty identified incorrectly-imputed genotypes with an areas under the
305 curve of the receiver operating characteristic curve (ROC-AUC) of 0.851 in the Japanese panel
306 and of 0.883 in the T1DGC reference panel in 4-digit alleles, which were superior to
307 dosage-based discrimination (ROC-AUC = 0.722 and 0.754 in the Japanese T1DGC panels,
308 respectively; **Supplementary Fig. 3**). Estimation of prediction uncertainty of a deep learning

309 model is still under development;³⁷ however, our results might illustrate its potential applicability
310 to the establishment of a reliability score for genotype imputation by deep neural networks.

311

312 **Trans-ethnic MHC fine-mapping of T1D**

313 We applied the DEEP*HLA models trained with our Japanese panel and the T1DGC panel to
314 HLA imputation of T1D GWAS data from BBJ (831 cases and 61,556 controls) and UKB (732
315 cases and 353,727 controls), respectively. T1D is a highly heritable autoimmune disease that
316 results from T cell-mediated destruction of insulin-producing pancreatic β cells.³⁸ We
317 performed imputation for GWAS data of the cohorts separately and then combined them to
318 perform trans-ethnic MHC fine-mapping (1,563 cases and 415,283 controls). We filtered
319 imputed alleles in which r^2 accuracy in 10-fold cross-validation was lower than 0.7 in the current
320 application.

321 Association analysis of the imputed HLA variants with T1D identified the most
322 significant association at the HLA-DR β 1 amino acid position 71 ($P_{\text{omnibus}} = P = 7.5 \times 10^{-120}$; **Fig.**
323 **5a and Supplementary Table 3**), one of the T1D risk-associated amino acid polymorphisms in
324 the European population.¹² As for T1D, the largest HLA gene associations were reported for a
325 combination of variants in the *HLA-DRB1*, *-DQA1*, and *-DQB1*;^{12,39} thus, we further investigated
326 independently associated variants within these tightly linked HLA genes before searching for
327 other risk-associated loci. When conditioning on HLA-DR β 1 amino acid position 71, we
328 observed the most significant independent association in HLA-DQ β 1 amino acid position 185
329 ($P_{\text{omnibus}} = 3.1 \times 10^{-69}$). Through stepwise forward conditional analysis in the class II HLA region,
330 we found significant independent associations for Tyr30 in HLA-DQ β 1 ($P_{\text{binary}} = 6.7 \times 10^{-20}$),

331 HLA-DR β 1 amino acid position 74 ($P_{\text{omnibus}} = 1.2 \times 10^{-11}$), and Arg70 in HLA-DQ β 1 ($P_{\text{omnibus}} =$
332 3.3×10^{-9} ; **Supplementary Fig. 4** and **Supplementary Table 4**).

333 These results were different from those of a previous study of a large T1D cohort of
334 European ancestries, which reported three amino acid polymorphisms, i.e., HLA-DQ β 1 position
335 57, HLA-DR β 1 position 13, and HLA-DR β 1 position 71, as the top-associated amino acid
336 polymorphisms in the *HLA-DRB1*, *-DQA1*, and *-DQB1* region. We then constructed multivariate
337 regression models for individual populations that incorporated our T1D risk-associated HLA
338 amino acid polymorphisms and classical alleles of *HLA-DRB1* and *HLA-DQB1*, and compared
339 the effects of these variants. The odds ratios of the risk-associated variants reported previously
340 did not show any positive correlation between different populations (Pearson's $r = -0.59$, $P =$
341 0.058 ; **Supplementary Fig. 5** and **Supplementary Table 5**). On the other hand, we identified a
342 set of variants with significant positive correlation by trans-ethnic fine-mapping of the integrated
343 cohort data (Pearson's $r = 0.76$, $P = 6.8 \times 10^{-3}$; **Supplementary Fig. 5**).

344 We further investigated whether the T1D risk was associated with other HLA genes
345 independently of *HLA-DRB1*, *-DQA1*, and *-DQB1*. When conditioning on *HLA-DRB1*, *-DQA1*,
346 and *-DQB1*, we identified a significant independent association at HLA-A amino acid position 62
347 ($P_{\text{omnibus}} = 5.9 \times 10^{-13}$; **Fig. 5b** and **Supplementary Table 3**). After conditioning on HLA-A
348 amino acid position 62, we did not observe any additional independent association in HLA-A
349 alleles. When we conditioned on *HLA-DRB1*, *-DQA1*, *-DQB1*, and *-A*, we identified a significant
350 independent association at HLA-B*54:01 ($P_{\text{binary}} = 1.3 \times 10^{-9}$; **Fig. 5c** and **Supplementary**
351 **Table 8**), and its unique amino acid polymorphisms (Gly45 and Val52 at HLA-B). When
352 conditioned on *HLA-DRB1*, *-DQA1*, *-DQB1*, *-A*, and *-B*, no variants in the MHC region satisfied
353 the genome-wide significance threshold ($P > 5.0 \times 10^{-8}$; **Fig. 5d** and **Supplementary Table 3**).

354 Multivariate regression analysis of the identified risk variants explained 10.3% and 27.6% of the
355 phenotypic variance in T1D under assumption of disease prevalence of 0.014%⁴⁰ and 0.4%.⁴¹
356 for the Japanese and British cohorts, respectively. Their odds ratios on T1D risk were also
357 correlated between different populations (Pearson's $r = 0.71$, $P = 4.4 \times 10^{-3}$; **Fig. 6** and **Table**
358 **1**).

359 To evaluate the advantage of the trans-ethnic fine-mapping, we performed
360 fine-mapping for each cohort separately and compared the results with those of the trans-ethnic
361 analysis. The most significant associations were observed in the *HLA-DRB1* and *-DQB1* in bot
362 h cohorts (**Supplementary Fig. 6** and **Supplementary Fig. 7**). The top signals were at the
363 HLA-DQ β 1 amino acid position 185 ($P = 8.3 \times 10^{-47}$) for the BBJ cohort and HLA-DR β 1 amino
364 acid position 71 ($P = 4.1 \times 10^{-107}$) for the UKB cohort, both of which were consistent with the
365 risk-associated variants identified through the trans-ethnic fine-mapping. On the other hand, the
366 risk-associated variants pointed in subsequent conditional analyses within this region were not
367 identical. Generally, parsimonious fine-mapping using a single population was challenging due
368 to multiple candidate variants with similar degrees of LD (and thus associations) to the top
369 signal in each iteration of the stepwise conditional analysis (**Supplementary Fig. 8** and
370 **Supplementary Fig. 9**). As a result of the trans-ethnic analysis, we successfully identified finer
371 sets of the more variants, which exhibited clearer significance by interrogating the different LD
372 patterns between the populations. When conditioning on *HLA-DRB1*, *-DQA1*, and *-DQB1*, we
373 identified significant independent associations in *HLA-B* for the BBJ cohort with the top at
374 HLA-B*54:01 ($P = 4.1 \times 10^{-10}$), and *HLA-A* for the UKB cohort with the top at HLA-A amino acid
375 position 62 ($P = 1.4 \times 10^{-8}$), respectively (**Supplementary Fig. 6** and **Supplementary Fig. 7**).
376 Both variants were identical to those originally identified in the trans-ethnic analysis. This

377 observation indicates that the trans-ethnic analysis could discover more associated loci than
378 single population-based analyses. Whereas HLA-B*54:01 was too rare and not assessed in
379 Europeans, it is notable that the T1D risk of HLA-A amino acid position 62 was shared with East
380 Asians. These observations should illustrate the value of the trans-ethnic MHC fine-mapping.

381

382 **Discussion**

383 In this study, we demonstrated that DEEP*HLA, a multi-task convolutional deep learning
384 method for HLA imputation, outperformed conventional HLA imputation methods in various
385 aspects. DEEP*HLA was more advantageous when the target HLA variants, including classical
386 alleles and amino acid polymorphisms, were low-frequent or rare. Our study demonstrated that
387 the performance of a conventional method was reduced for alleles that did not exhibit
388 distance-dependent LD decay with the target HLA allele. DEEP*HLA was less dependent on
389 this point, and might comprehensively capture the relationships among multiple distinct variants
390 regardless of LD. Taking advantage of the significant improvement of imputation accuracy in
391 rare alleles, we conducted trans-ethnic MHC fine-mapping of T1D. This approach could be
392 performed as well using the conventional HLA imputation methods. However, the results
393 obtained using DEEP*HLA should be more reliable because there were several risk-associated
394 alleles which were rare only in one population.

395 To date, technical application of deep neural networks to population genetics data has
396 been limited. In a previous attempt for genotype imputation, a sparse convolutional denoising
397 autoencoder was only compared with reference-free methods.³² There might be two possible
398 explanations for the success of our DEEP*HLA. First unlike genotype imputation by denoising
399 autoencoders, which assumes various positions of missing genotypes in a reference panel to
400 impute, the prediction targets were fixed to the HLA allele genotypes as a classification problem.
401 Second, convolutional neural networks, which leverage a convolutional kernel that is capable of
402 learning various local patterns, might be better suited for learning the complex LD structures in
403 the MHC region.

404 We filtered alleles with poor imputation quality based on the results of cross-validation
405 in the current application; however, an indicator of reliability could be further utilized. We
406 demonstrated that the prediction uncertainty inferred from a Bayesian deep learning method
407 had potential capability of identifying incorrectly-imputed alleles in a per-gene level. Our future
408 work should establish a method to quantify per-allele imputation uncertainty that can be
409 practically used as a filtering threshold for subsequent analyses.

410 As for the genetic features of the MHC region associated with T1D, the highest risk is
411 conferred by DR3-DQA1*05-DQB1*02 and DR4-DQA1*03-DQB1*03:02 haplotypes in
412 Europeans,^{39,42} and by DR9-DQA1*03-DQB1*03:03 and DR4-DQA1*03-DQB1*04:01
413 haplotypes in Japanese.⁴³ In a previous study for a large European cohort, Hu et al.
414 demonstrated that the three amino acid polymorphisms of DR β 1 and HLA-DQ β 1 explained the
415 majority of the risk in the *HLA-DRB1*, *-DQA1*, and *-DQB1* region with the top signal at
416 non-Asp57 in HLA-DQ β 1.¹² Conversely, the risk haplotypes in Japanese population carry
417 Asp57 of HLA-DQ β 1.⁴³ We obtained several additional insights in the present study. We initially
418 conducted a trans-ethnic MHC fine-mapping of T1D, and successfully disentangled a set of 5
419 risk-associated amino acid polymorphisms of position 71 and 74 in HLA-DR β 1, and 30, 70, and
420 185 in HLA-DQ β 1. Four of these positions compose the peptide-binding grooves, suggesting
421 their functional contributions to antigen-presentation ability (**Supplementary Fig. 10**). While the
422 association of HLA-DR β 1 amino acid position 71 was replicated in concordance direction with
423 Europeans, the effects in the Japanese population were not preserved in the final model.
424 Whereas the association of amino acid position 74 in HLA-DR β 1 has been reported in
425 Han-Chinese and certain European populations,^{44,45} the European study did not report its
426 independent association due to the rareness of its characterized classical

427 allele, HLA-DRB1*04:03. We successfully identified its independent association in trans-ethnic
428 cohorts with a similar effect size between the diverse populations. Although amino acid position
429 185 in HLA-DQ β 1 does not compose the peptide-binding groove, the variation of Ile/Thr is
430 suggested to alter DQ-DM anchoring by interacting with its neighboring residues, leading to the
431 susceptibility to other autoimmune diseases.^{46,47} Variant Ile185 is tagged with HLA-DQA1*03,
432 which composes the risk haplotypes in Japanese and European population respectively. A
433 correspondence table of the amino acid polymorphisms and 4-digit classical HLA alleles is
434 shown in **Supplementary Table 6**. As a result, the catalogue of the T1D risk-associated
435 variants in this region identified by our trans-ethnic approach was different from that in the
436 European study.¹² We admit the possibility that the smaller sample size in our study and
437 different definitions of the phenotypes (between studies, and between cohorts in our study)
438 might contribute to this disparity. Particularly, we note the potential distinctiveness of Japanese
439 T1D phenotypes.⁴⁸ However, considering that our observed variants shared the effects on the
440 T1D risk between different populations, we might gain a novel insight into the issue of
441 inter-ethnic heterogeneity of T1D risk alleles in the MHC region. As for class I HLA genes, the
442 independent association of amino acid position 62 in HLA-A was consistent with the previous
443 European study.¹² We found that it had similar effects on the T1D risk also in the Japanese
444 population. HLA-B*54:01 has traditionally been suggested as a potential risk allele in Japanese
445 by a candidate HLA gene approach,¹³ of which an independent association via the MHC
446 region-wide fine-mapping was first proven here.

447 While an advantage of trans-ethnic fine-mapping is the elucidation of truly
448 risk-associated signals by adjusting confound by LD of each population,⁴⁹ there are several
449 potential limitations to note. First, we need to consider population-specific LD structures and

450 allele frequency spectra, which are important especially in the MHC region. Strong
451 population-specificity may preclude removal of the effects of LD for the current purpose of
452 trans-ethnic fine-mapping when few populations are available. Conversely, some HLA alleles
453 exist only in a certain population, and fine-mapping in a single population could also be of
454 importance. Second, modeling heterogeneity in effects among diverse populations could
455 enhance the power of discovery of causal variants in trans-ethnic analysis.⁵⁰ Since the purpose
456 of the current trans-ethnic fine-mapping is to identify trans-ethnically risk-associated variants
457 rather than to discover variants with a strong effect only in one population, we did not explicitly
458 model heterogeneity. However, in an analysis using more cohorts from different populations,
459 modeling heterogeneity might be more suitable because a bias by single population would be
460 reduced.

461 Therefore, multi-ethnic MHC fine-mapping that integrates further diverse ancestry
462 should be warranted for robust prioritization of risk-associated HLA variants as a next step.¹⁵
463 Given their high learning capacity of deep neural networks, our method will be helpful not only
464 when integrating the imputation results from multiple references, but also when using a more
465 comprehensive multi-ethnic reference. We expect that highly accurate imputation realized by
466 learning of complex LDs in the MHC region using neural networks will enable us to further
467 elucidate the involvement of common genetic features in the MHC region that affect human
468 complex traits across ethnicities.

469

470

471 **Acknowledgements**

472 We would like to thank all the participants involvement in this study. We thank the members of
473 Biobank Japan and RIKEN Center for Integrative Medical Sciences for their supports on this
474 study.

475

476 **Conflicts of interests**

477 The authors declare no conflicts of interests.

478

479 **Data availability**

480 The Japanese HLA data have been deposited at the National Bioscience Database Center
481 (NBDC) Human Database (research ID: hum0114). Independent HLA genotype data of
482 Japanese population is available in the Japanese Genotype-phenotype archive (JGA;
483 accession ID: JGAS00000000018). T1DGC HLA reference panel can be download at a NIDDK
484 central repository with a request (<https://repository.niddk.nih.gov/studies/t1dgc-special/>).
485 GWAS data of the BBJ are available at the NBDC Human Database (research ID: hum0014).
486 The analysis of UKB GWAS data was conducted via the application number 47821
487 (<https://www.ukbiobank.ac.uk/>).

488

489 **Code availability**

490 Python scripts for training a model and performing imputation with our method are in
491 DEEP*HLA GitHub repository (<https://github.com/tatsuhikonaito/DEEP-HLA>).

492

493 **Methods**

494 **The architecture of DEEP*HLA**

495 DEEP*HLA is a multi-task convolutional neural network comprising a shared part of two
496 convolutional layers and a fully-connected layer, and individual fully-connected layers that
497 output allelic dosages of individual HLA genes to simultaneously impute HLA genes of the
498 same group (**Fig. 1a**). The grouping was based on the LD structure³ and physical distance in
499 the current application: (1) {*HLA-F*, *HLA-V*, *HLA-G*, *HLA-H*, *HLA-K*, *HLA-A*, *HLA-J*, *HLA-L*, and
500 *HLA-E*}, (2) {*HLA-C*, *HLA-B*, *MICA*, and *MICB*}, (3) {*HLA-DRA*, *HLA-DRB9*, *HLA-DRB5*,
501 *HLA-DRB4*, *HLA-DRB3*, *HLA-DRB8*, *HLA-DRB7*, *HLA-DRB6*, *HLA-DRB2*, *HLA-DRB1*,
502 *HLA-DQA1*, *HLA-DOB*, and *HLA-DQB1*}, and (4) {*TAP2*, *TAP1*, *HLA-DMB*, *HLA-DMA*,
503 *HLA-DOA*, *HLA-DPA1*, and *HLA-DPB1*}. Genes not typed or with only single alleles in individual
504 reference panels were excluded from the group. Comparisons with single-task neural networks
505 or multi-task neural networks with random groupings are shown in **Supplementary Note 1b**
506 and **Supplementary Fig. 11**.

507 DEEP*HLA takes the input of each haplotype SNV genotypes from pre-phased data,
508 and outputs the genotype dosages of individual alleles for each HLA gene. For each group,
509 SNVs within its window are encoded to one-hot vectors based on whether each genotype is
510 consistent with a reference or alternative allele. The window sizes on each side were fixed to
511 500 kb for fair comparisons in the current investigation; using different window sizes might
512 slightly change the accuracy for some loci (**Supplementary Note 1c** and **Supplementary Fig.**
513 **12**). Two convolutional layers with max-pooling layers and a fully-connected layer follow the
514 input layer as a shared part. The fully-connected layer at the end of the shared part is followed
515 by individual fully-connected layers which have nodes consistent with the number of alleles of

516 each HLA gene. Softmax activation was added before the last output to return an imputation
517 dosage that ranges from 0.0 to 1.0 for each allele of one haplotype. Thus, an individual layer
518 outputs the individual allelic dosages of the HLA gene of which the sum equals 1 for one
519 haplotype. Dropout was used on the convolutional and fully-connected layers,⁵¹ and batch
520 normalization was added to the convolutional layers.⁵²

521 During training, 5% of the data set were used for sub-validation to determine the point
522 for early-stopping training. In 10-fold cross-validation, we separated sub-validation for
523 early-stopping from a training fold to conduct valid benchmarking (**Supplementary Fig. 13**). A
524 categorical cross entropy loss function for each HLA gene was minimized using the Adam
525 optimizing algorithm.⁵³ For a multi-task learning to find a Pareto optimal solution of all tasks, we
526 used the multiple-gradient descent algorithm – upper bound (MGDA-UB), where the loss
527 function of each task was scaled based on its optimization algorithms.⁵⁴ To taking advantage of
528 the hierarchical nature of HLA alleles (i.e. 2-digit, 4-digit, and 6-digit), we implemented
529 hierarchical fine-tuning, in which parameters of the model of upper hierarchical structures were
530 transferred to those of the lower one.⁵⁵ We transferred the parameters of shared networks of
531 2-digit alleles to 4-digit alleles, and of 4-digit alleles to 6-digit alleles successively during training.
532 Although some HLA alleles in our reference panel were not determined in 4-digit or 6-digit
533 resolution, we set their upper resolution instead to maintain equivalent hierarchical levels with
534 other HLA genes. Hyperparameters, including the number of filters and kernel sizes of
535 convolutional layers, and fully-connected layer size, were tuned using Optuna.⁵⁶ The
536 hyperparameters for each reference panel were determined using a randomly sampled dataset
537 before cross-validation. Our deep learning architectures were implemented using Pytorch 1.4.1
538 (see URLs), a Python neural network library.

539

540 Empirical evaluation of HLA imputation accuracy

541 We used the accuracy metrics of sensitivity, PPV, and r^2 for imputed allelic dosage, and
542 concordance rate for best-guess genotypes to evaluate the imputation accuracy in various
543 aspects.

544 In the paper of SNP2HLA, per-locus accuracy was defined as a sum of the dosage of
545 each true allele across all individuals divided by the total number of observations.³³ This
546 definition of accuracy counts positives that are correctly identified as such and it corresponds to
547 sensitivity in a cross-tabulation table when decomposed to individual alleles (**Supplementary**
548 **Note 2 and Supplementary Fig. 14**). Thus, we termed this as sensitivity (Se) to contrast with
549 the PPV defined later

$$Se(L) = \frac{\sum_{i=1}^n (D_i(A1_{i,L}) + D_i(A2_{i,L}))}{2n}$$

550 where n denotes the number of individuals, D_i represents the imputed dosage of an allele in
551 individual i , and alleles $A1_{i,L}$ and $A2_{i,L}$ represent the true HLA alleles for individual i at locus L .
552 The calculations were based on the condition that the imputed alleles are arranged to optimize
553 for consistency with the truth alleles $A1_{i,L}$ and $A2_{i,L}$.

554 To evaluate the imputation performance in individual HLA alleles, we decomposed the
555 $Se(L)$ to evaluate the imputation performance of each allele as.

$$Se(A) = \frac{\sum_{j=1}^m D_j(A)}{m}$$

556 This metric cannot evaluate the effect of false positives; thus, we defined PPV in the same
557 manner as

$$PPV(A) = \frac{\sum_{j=1}^m D_j(A)}{\sum_{j=1}^m D_j(A) + \sum_{k=1}^{2n-m} D_k(A)}$$

558 where m denotes the number of true observations of allele A in the total sample, and D_j
559 represents imputed dosage of allele A in individual haplotype j that has allele A . D_k represents
560 imputed dosage of allele A in individual haplotype k that has an allele other than allele A . This
561 definition is also based on a cross-tabulation table (**Supplementary Fig. 14a**).

562 In addition, we calculated r^2 based on Pearson's product moment correlation
563 coefficient between imputed and typed dosages for each allele.²²

564 Further, to evaluate the accuracy of best-guess genotypes, we calculated the
565 concordance rate (CR) of best-guess genotypes and true genotypes for each allele as

$$CR(L) = \frac{\sum_{i=1}^n (B_i(A1_{i,L}) + B_i(A2_{i,L}))}{2n}$$

566 where B_i represents the best-guess genotype of an allele in individual i . By definition, it was the
567 same as the sensitivity, in which dosages were changed to best-guess genotypes. Thus, we
568 decomposed it to $CR(A)$ for accuracy for each allele in the same way. We did not evaluate PPV
569 for best-guess genotype due to redundancy.

570 When determining accuracy metrics for each locus or a certain range of allele
571 frequencies, we calculated the weighted-mean of individual allele-level accuracies based on
572 individual allele frequencies. For r^2 , we applied Fisher's Z transformation to individual values,
573 and back-transformed them after averaging to reduce bias.⁵⁷

574

575 **Estimation of HLA imputation uncertainty of DEEP*HLA using MC dropout method**

576 In order to estimate prediction uncertainty, we adopted the entropy of sampling variation of MC
577 dropout method.³⁶ In MC dropout, dropouts are kept during prediction to perform multiple model

578 calls. Different units are dropped across different model calls; thus, it can be considered as
579 Bayesian sampling with treating the parameters of a CNN model as random variables of
580 Bernoulli distribution. The uncertainty of a best-guess genotype inferred from the entropy of
581 sampling variation is determined as

$$H = -\left(\frac{t}{T} \log \frac{t}{T} + \frac{T-t}{T} \log \frac{T-t}{T}\right)$$

582 where T is the number of variational samplings and t is the number of times in which obtained
583 genotype was identical to the best-guess genotype. We set $T = 200$ in the current investigation.

584 585 **AUC metric representing distance-dependent LD decay**

586 To evaluate whether the the strength of LD between an HLA allele and its surrounding SNVs
587 weakens as the the distance between them increases, we calculated the AUC of the cumulative
588 curve of r^2 from the HLA allele (AUC for distance-dependent LD decay). When the LD of
589 flanking SNVs of an HLA allele has such a characteristic, r^2 of LD from the HLA allele tends to
590 decrease. In other words, the bilateral cumulative curve of r^2 from the HLA allele is more likely
591 to be convex upward; then, the AUC tends to be higher. We determined the AUC by
592 normalizing the maximum values of r^2 sum and window sizes to 1. We evaluated the
593 association of the AUC with allele-level accuracy metrics of each imputation method by linear
594 regression models adjusted for an allele frequency. The window size of the AUC should be set
595 to an input range for each imputation method. However, SNP2HLA does not have a clear input
596 range. Thus, we tested two different window sizes as bilateral 1,000 SNPs from a target HLA
597 allele and the input range of DEEP*HLA. We investigated the correlation between the
598 imputation accuracy and the AUC of two different window sizes, respectively.

599

600 **Regional sensitivity maps of DEEP*HLA**

601 We applied SmoothGrad to estimate which SNVs were important for DEEP*HLA imputation of
602 each HLA allele.³⁴ For each haplotype, we generated 200 samples which were added Gaussian
603 noise to encoded SNV data and input them into a trained model. Sensitivity values for individual
604 SNV positions were obtained by averaging the absolute values of gradients caused by the
605 difference from the true label. When we obtained the sensitivity of an HLA allele, we averaged
606 the maps of all haplotypes that have the target HLA allele.

607

608 **HLA imputation software and parameter settings**

609 We tested the latest version of the software available in Jun 2020 for comparison with our
610 method. SNP2HLA (v1.0.3) first arranges the strand in its own algorithm; however, we removed
611 this step data during cross-validation because the strands must be the same between training
612 and test data. The other settings of SNP2HLA were set to the default values. For HIBAG
613 (1.22.0.) the number of classifiers was set to 25, which is sufficient to achieve good
614 performance⁵⁸ for testing the Japanese data. For the T1DGC panel, the training time was
615 extremely long with 25 classifiers; thus, we set 2 classifiers after we confirmed that the
616 imputation accuracy was almost unchanged in the first set of cross-validation. The flanking
617 regions on each side were set to 500 kb. The current version of HLA*IMP:02 did not support a
618 function to generate an imputation model using own reference data in a publicly available form;
619 thus, we did not evaluate its performance in this study for fair comparison.

620

621 **Measurement of computational costs**

622 We measured the computational costs of imputation of a subset of BioBank Japan (BBJ)
623 Project data set ($n = 1,000, 2,000, 5,000, 10,000, 20,000, 50,000,$ and $100,000$ samples) using
624 our Japanese reference panel (2,000 SNVs were consistent). All our runtime analyses were
625 performed on a dedicated server running CentOS 7.2.1511, with 48 CPU cores (Intel® Xeon®
626 E5-2687W v4 @ 3.00 GHz) and 256 GB of RAM without GPU. Additionally, we measured the
627 training time of DEEP*HLA with GPU using a machine with Ubuntu 16.04.6 LTS with 20 CPU
628 cores (Intel® Core™ i9-9900X @ 3.50 GHz), 2 GPUs (NVIDIA® GeForce® RTX 2080 Ti),
629 and 128 GB of RAM. DEEP*HLA requires pre-phased GWAS data and the models trained with
630 reference data; thus, we measured the process of not only imputation, but also pre-phasing of
631 GWAS data (conducted by Eagle) and training the models with a reference panel. Similarly,
632 HIBAG requires the time for training a model, which was also measured. In SNP2HLA, the
633 maximum of available memory was set to 100 GB. The processing time and maximum memory
634 usage were measured using GNU Time software when running from a command line interface.

635

636 **HLA imputation reference data**

637 We used two HLA reference panels in cross-validation and HLA imputation for biobank GWAS
638 data. The panels were distributed as a phased condition; thus, they were used as input for
639 training a DEEP*HLA model as they were. When they were used as a validation set, we
640 removed the target alleles (i.e. HLA alleles and amino acid alleles) to leave only phased SNP
641 data. We discussed stricter cross-validation including the process of haplotype pre-phasing in

642 **Supplementary Note 1d.**

643 (i) Our Japanese reference panel and a validation dataset

644 Our Japanese reference panel contained NGS-based 6-digit resolution HLA typing data of 33
645 classical and non-classical HLA genes, of which 9 were classical HLA genes (*HLA-A*, *HLA-B*,
646 and *HLA-C* for class I; *HLA-DRA*, *HLA-DRB1*, *HLA-DQA1*, *HLA-DQB1*, *HLA-DPA1*, and
647 *HLA-DPB1* for class II) and 24 were non-classical HLA genes (*HLA-E*, *HLA-F*, *HLA-G*, *HLA-H*,
648 *HLA-J*, *HLA-K*, *HLA-L*, *HLA-V*, *HLA-DRB2*, *HLA-DRB3*, *HLA-DRB4*, *HLA-DRB5*, *HLA-DRB6*,
649 *HLA-DRB7*, *HLA-DRB8*, *HLA-DRB9*, *HLA-DOA*, *HLA-DOB*, *HLA-DMA*, *HLA-DMB*, *MICA*, *MICB*,
650 *TAP1*, and *TAP2*), along with high-density SNP data in the MHC region by genotyped using the
651 Illumina HumanCoreExome BeadChip (v1.1; Illumina) of 1,120 unrelated individuals of
652 Japanese ancestry.³ It was phased using Beagle imputation software. We excluded 2
653 individuals' data of which sides of some HLA alleles were inconsistent among different
654 resolutions.

655 We used 908 individuals of Japanese ancestry with 4-digit resolution alleles of
656 classical HLA genes (*HLA-A*, *HLA-B*, *HLA-C*, *HLA-DRB1*, *HLA-DQA1*, *HLA-DQB1*, *HLA-DPA1*)
657 based on SSO method to benchmark the imputation performance when the Japanese panel
658 was applied to an independent dataset. The dataset was used as an HLA reference panel in our
659 previous study.⁶ It contains high-density SNP data genotyped using four SNP genotyping
660 arrays (the Illumina HumanOmniExpress BeadChip, the Illumina HumanExome BeadChip, the
661 Illumina ImmunoChip, and the Illumina HumanHap550v3 Genotyping BeadChip). It was
662 distributed in a phased condition with Beagle format. Samples with missing genotype data for a
663 locus were excluded in the accuracy evaluation of the locus. This study was approved by the
664 ethics committee of Osaka University Graduate School of Medicine with written informed
665 consent obtained from all participants.

666

667 (ii) The Type 1 Diabetes Genetics Consortium (T1DGC) reference panel.
668 The T1DGC panel contains 5,868 SNPs (genotyped using Illumina ImmunoChip) and 4-digit
669 resolution HLA typing data of classical HLA genes (*HLA-A*, *HLA-B*, and *HLA-C* for class I;
670 *HLA-DPA1*, *HLA-DPB1*, *HLA-DQA1*, *HLA-DQB1*, and *HLA-DRB1* for class II) based on SSO
671 method of 5,225 unrelated individuals of European ancestry.²² It was distributed in a phased
672 condition with Beagle format. We excluded 103 individuals' data of which sides of some HLA
673 alleles were inconsistent among different resolutions.

674

675 **HLA imputation in 1000 Genomes Project data**

676 We used Phase III 1000 Genomes Project (1KGv3) cohort as independent data to evaluate
677 imputation accuracy. It comprises 2,554 individuals of 5 different super populations (AFR, AMR,
678 EAS, EUR, and SAS). We obtained NGS-based 4-digit resolution HLA typing data for classical
679 HLA genes (*HLA-A*, *HLA-B*, and *HLA-C* for class I; *HLA-DRB1* and *HLA-DQB1* for class II). The
680 process of HLA typing has been described elsewhere.⁵⁹ We evaluated imputation accuracy for
681 individual populations based on their allele frequencies. Samples containing ambiguous alleles
682 for a locus were excluded in the accuracy evaluation of that locus.

683 We experimentally constructed a mixed panel by merging the Japanese and T1DGC
684 panels to assess imputation accuracy in diverse populations of 1KGv3. Considering the
685 disparity in allele frequency of SNPs between two populations, we removed all palindromic
686 SNVs to align the strands correctly when merging reference panels. We used 1,445 SNPs for
687 imputation which were consistent with 1KGv3 genotype data. We used the same 1,445 SNPs
688 for imputation to compare the accuracies in the same condition when we evaluated imputation
689 accuracy using the Japanese panel.

690

691 **T1D GWAS data in the Japanese population**

692 The BioBank Japan (BBJ) is a multi-institutional hospital-based registry that comprises DNA,
693 serum, and clinical information of approximately 200,000 individuals of Japanese ancestry
694 recorded from 2003 to 2007.^{60,61} We used GWAS data from 831 cases who had record of T1D
695 diagnosis and 61,556 controls of Japanese genetic ancestry enrolled in the BBJ Project. The
696 controls were same as those enrolled in our previous study that investigated the association of
697 the MHC region with comprehensive phenotypes, and the number of T1D cases was
698 increased.³ The process of patient registration, the GWAS data, and the QC process are
699 described elsewhere.^{60–62}

700

701 **T1D GWAS data in the British population**

702 The UK Biobank (UKB) comprises health-related information approximately 500,000 individuals
703 aged between 40–69 recruited from across the United Kingdom from 2006 to 2010.⁶³ We used
704 GWAS data of 732 T1D patients and 353,727 controls of British genetic ancestry enrolled in
705 UKB. We selected T1D patients as individuals who were diagnosed with insulin-dependent
706 diabetes mellitus in hospital records, and eliminated individuals with non-insulin-independent
707 diabetes mellitus in hospital records and type 2 diabetes in self-reported diagnosis. The controls
708 were individuals with no record of any autoimmune diseases in hospital records or in
709 self-reported diagnosis. We included only individuals of British ancestry according to
710 self-identification and criteria based on principal component (PC).⁶⁴ We excluded individuals of
711 ambiguous sex (sex chromosome aneuploidy and inconsistency between self-reported and
712 genetic sex), and outliers of heterozygosity or call rate of high quality markers.

713

714 **Imputation of the HLA variants of GWAS data of T1D cases and controls**

715 In this study, we defined the HLA variants as SNVs in the MHC region, classical 2-digit and
716 4-digit biallelic HLA alleles, biallelic HLA amino acid polymorphisms corresponding to the
717 respective residues, and multiallelic HLA amino acid polymorphisms for each amino acid
718 position. We applied DEEP*HLA to the GWAS data to determine classical 2-digit and 4-digit
719 biallelic HLA alleles. The dosages of biallelic HLA amino acid polymorphisms corresponding to
720 the respective residues and multiallelic HLA amino acid polymorphisms of each amino acid
721 position were determined from the imputed 4-digit classical allele dosages. We applied
722 post-imputation filtering as the biallelic alleles in which r^2 accuracy in 10-fold cross-validation
723 was lower than 0.7. The SNVs in the MHC region were imputed using minimac3 (version 2.0.1)
724 after pre-phased with Eagle (version 2.3). We applied stringent post-imputation QC filtering of
725 the variants (minor allele frequency $\geq 0.5\%$ and imputation score $R_{sq} \geq 0.7$). For trans-ethnic
726 fine-mapping, we integrated results of the imputation of individual cohorts by including the HLA
727 genes, amino acid position, and SNVs that were typed in both reference panels. Regarding the
728 HLA alleles and amino acid polymorphisms, those present in one population were regarded as
729 absent in the other population. Considering the disparity in allele frequency of SNVs among
730 different populations, we removed all palindromic SNVs to correctly align the strands.

731

732 **Association testing of the HLA variants**

733 We assumed additive effects of the allele dosages on the log-odds scales for susceptibility to
734 T1D; and evaluated associations of the HLA variants with the risk of T1D using a logistic
735 regression model. To robustly account for potential population stratification, we included the top

736 10 PCs obtained from the GWAS genotype data of each cohort (not including the MHC region)
737 as covariates in the regression model. We also included ascertainment centre and genotyping
738 chip for UKB as covariates. For trans-ethnic analysis, PC terms for each other population were
739 set to 0, and a categorical variable indicating a population was added as a covariate. We also
740 included the sex of individuals as a covariate.

741 To evaluate independent risk among the HLA variants and genes, we conducted a
742 forward-type stepwise conditional regression analysis that additionally included the associated
743 variant genotypes as covariates. When conditioning on HLA gene(s), we included all the 4-digit
744 alleles as covariates to robustly condition the associations attributable to the HLA genes, as
745 previously described.^{3,14} When conditioning on the specific HLA amino acid position(s), we
746 included the multiallelic variants of the amino acid residues. We applied a forward stepwise
747 conditional analysis for the HLA variants and then HLA genes, based on a genome-wide
748 association significance threshold ($P = 5.0 \times 10^{-8}$). A previous study reported that the T1D risk
749 was strongly associated with a combination of variants in the region of *HLA-DRB1*, *-DQA1*, and
750 *-DQB1*, where the variants have strong LD to each other.¹² In such a situation, conditioning on
751 all the 4-digit alleles of a single HLA gene might inadvertently blind the association of alleles of
752 other HLA genes; therefore, we conditioned on a set of individual HLA variants rather than an
753 each HLA gene when analyzing this region.

754 We tested a multivariate full regression model by including the risk-associated HLA
755 variants in *HLA-DRB1*, *HLA-DQB1*, *HLA-A*, and *HLA-B*, which were identified through the
756 stepwise regression analysis. We excluded the most frequent residue in the British cohort from
757 each amino acid position as the reference allele when we included amino acid polymorphisms
758 in the model. Phenotypic variance explained by the identified risk-associated HLA variants was

759 estimated on the basis of a liability threshold model assuming a population-specific prevalence
760 of T1D and using the effect sizes obtained from the multivariate regression model.

761

762 **URLs**

763 DEEP*HLA, <https://github.com/tatsuhikonaito/DEEP-HLA>

764 Pytorch, <http://pytorch.org/>

765 SNP2HLA, <http://software.broadinstitute.org/mpg/snp2hla/>

766 HIBAG, <https://www.bioconductor.org/packages/release/bioc/html/HIBAG.html>

767 Eagle, <https://data.broadinstitute.org/alkesgroup/Eagle/>

768 Minimac3, <https://genome.sph.umich.edu/wiki/Minimac3>

769 BioBank Japan, <https://biobankjp.org/english/index.html>

770 UK Biobank, <https://www.ukbiobank.ac.uk/>

771 UCSF Chimera, <https://www.cgl.ucsf.edu/chimera/>

772

773 **References**

- 774 1. Dendrou, C. A., Petersen, J., Rossjohn, J. & Fugger, L. HLA variation and disease. *Nat.*
775 *Rev. Immunol.* **18**, 325–339 (2018).
- 776 2. Erlich, H. HLA DNA typing: Past, present, and future. *Tissue Antigens* **80**, 1–11 (2012).
- 777 3. Hirata, J. *et al.* Genetic and phenotypic landscape of the major histocompatibility complex
778 region in the Japanese population. *Nat. Genet.* **51**, 470–480 (2019).
- 779 4. Pereyra, F. *et al.* The major genetic determinants of HIV-1 control affect HLA class I
780 peptide presentation. *Science (80-.)*. **330**, 1551–1557 (2010).
- 781 5. Raychaudhuri, S. *et al.* Five amino acids in three HLA proteins explain most of the
782 association between MHC and seropositive rheumatoid arthritis. *Nat. Genet.* **44**, 291–296
783 (2012).
- 784 6. Okada, Y. *et al.* Construction of a population-specific HLA imputation reference panel and
785 its application to Graves' disease risk in Japanese. *Nat. Genet.* **47**, 798–802 (2015).
- 786 7. Lee, H. & Kingsford, C. Kourami: Graph-guided assembly for novel human leukocyte
787 antigen allele discovery. *Genome Biol.* **19**, 1–16 (2018).
- 788 8. Dilthey, A. T. *et al.* HLA*LA - HLA typing from linearly projected graph alignments.
789 *Bioinformatics* **35**, 4394–4396 (2019).
- 790 9. Gourraud, P. A. *et al.* HLA diversity in the 1000 genomes dataset. *PLoS One* **9**, (2014).
- 791 10. Okada, Y. *et al.* Risk for ACPA-positive rheumatoid arthritis is driven by shared HLA
792 amino acid polymorphisms in Asian and European populations. *Hum. Mol. Genet.* **23**,
793 6916–6926 (2014).
- 794 11. Todd JA, Bell JI & McDevitt HO. HLA-DQBeta gene contributes to susceptibility and
795 resistance to insulin-dependent diabetes mellitus. *Nature* **329**, 599–604 (1987).

- 796 12. Hu, X. *et al.* Additive and interaction effects at three amino acid positions in HLA-DQ and
797 HLA-DR molecules drive type 1 diabetes risk. *Nat. Genet.* **47**, 898–905 (2015).
- 798 13. Kawabata, Y. *et al.* Differential association of HLA with three subtypes of type 1 diabetes:
799 Fulminant, slowly progressive and acute-onset. *Diabetologia* **52**, 2513–2521 (2009).
- 800 14. Okada, Y. *et al.* Contribution of a Non-classical HLA Gene, HLA-DOA, to the Risk of
801 Rheumatoid Arthritis. *Am. J. Hum. Genet.* **99**, 366–374 (2016).
- 802 15. Luo, Y. *et al.* A high-resolution HLA reference panel capturing global population diversity
803 enables multi-ethnic fine-mapping in HIV host response. *medRxiv Prepr.* (2020)
804 doi:<https://doi.org/10.1101/2020.07.16.20155606>.
- 805 16. De Bakker, P. I. W. *et al.* A high-resolution HLA and SNP haplotype map for disease
806 association studies in the extended human MHC. *Nat. Genet.* **38**, 1166–1172 (2006).
- 807 17. Monsuur, A. J. *et al.* Effective detection of human leukocyte antigen risk alleles in celiac
808 disease using tag single nucleotide polymorphisms. *PLoS One* **3**, 1–6 (2008).
- 809 18. Leslie, S., Donnelly, P. & McVean, G. A Statistical Method for Predicting Classical HLA
810 Alleles from SNP Data. *Am. J. Hum. Genet.* **82**, 48–56 (2008).
- 811 19. Li, Na (Department of Biostatistics, University of Washington, Seattle, W. 98195) &
812 Stephens, Matthew (Department of Statistics, University of Washington, Seattle, W.
813 98195). Modeling Linkage Disequilibrium and Identifying Recombination Hotspots Using
814 Single-Nucleotide Polymorphism Data. *Genetics* **165**, 2213–2233 (2003).
- 815 20. Dilthey, A. T., Moutsianas, L., Leslie, S. & McVean, G. HLA*IMP-an integrated framework
816 for imputing classical HLA alleles from SNP genotypes. *Bioinformatics* **27**, 968–972
817 (2011).
- 818 21. Dilthey, A. *et al.* Multi-Population Classical HLA Type Imputation. *PLoS Comput. Biol.* **9**,

- 819 e1002877 (2013).
- 820 22. Jia, X. *et al.* Imputing Amino Acid Polymorphisms in Human Leukocyte Antigens. *PLoS*
821 *One* **8**, (2013).
- 822 23. Levin, A. M. *et al.* Performance of HLA allele prediction methods in African Americans for
823 class II genes HLA-DRB1, -DQB1, and -DPB1. *BMC Genet.* **15**, 1–11 (2014).
- 824 24. Karnes, J. H. *et al.* Comparison of HLA allelic imputation programs. *PLoS One* **12**, 1–12
825 (2017).
- 826 25. Eraslan, G., Avsec, Ž., Gagneur, J. & Theis, F. J. Deep learning: new computational
827 modelling techniques for genomics. *Nat. Rev. Genet.* **20**, 389–403 (2019).
- 828 26. Alipanahi, B., Delong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence
829 specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol* **33**,
830 831–838 (2015).
- 831 27. Zhou, J. *et al.* Deep learning sequence-based ab initio prediction of variant effects on
832 expression and disease risk [supplementary]. *Nat. Genet.* **50**, 1171–1179 (2018).
- 833 28. Sundaram, L. *et al.* Predicting the clinical impact of human mutation with deep neural
834 networks. *Nat. Genet.* **50**, 1161–1170 (2018).
- 835 29. Naito, T. Predicting the impact of single nucleotide variants on splicing via
836 sequence-based deep neural networks and genomic features. *Hum. Mutat.* **40**,
837 1261–1269 (2019).
- 838 30. Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deep generative models of genetic
839 variation capture the effects of mutations. *Nat. Methods* **15**, 816–822 (2018).
- 840 31. Dwivedi, S. K., Tjörnberg, A., Tegnér, J. & Gustafsson, M. Deriving disease modules from
841 the compressed transcriptional space embedded in a deep autoencoder. *Nat. Commun.*

- 842 **11**, (2020).
- 843 32. Chen, J. & Shi, X. Sparse convolutional denoising autoencoders for genotype imputation.
844 *Genes (Basel)*. **10**, 1–16 (2019).
- 845 33. Han, B. *et al.* Fine mapping seronegative and seropositive rheumatoid arthritis to shared
846 and distinct HLA alleles by adjusting for the effects of heterogeneity. *Am. J. Hum. Genet.*
847 **94**, 522–532 (2014).
- 848 34. Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. SmoothGrad: removing
849 noise by adding noise. (2017).
- 850 35. Kendall, A. & Gal, Y. What uncertainties do we need in Bayesian deep learning for
851 computer vision? *Adv. Neural Inf. Process. Syst.* **2017-Decem**, 5575–5585 (2017).
- 852 36. Gal, Y. & Ghahramani, Z. Bayesian Convolutional Neural Networks with Bernoulli
853 Approximate Variational Inference. 1–17 (2015).
- 854 37. Gal, Y. & Ghahramani, Z. Dropout as a Bayesian approximation: Representing model
855 uncertainty in deep learning. *33rd Int. Conf. Mach. Learn. ICML 2016* **3**, 1651–1660
856 (2016).
- 857 38. Atkinson, M. A., Eisenbarth, G. S. & Michels, A. W. Type 1 diabetes. *Lancet* **383**, 69–82
858 (2014).
- 859 39. Erlich, H. *et al.* HLA DR-DQ haplotypes and genotypes and type 1 diabetes risk analysis
860 of the type 1 diabetes genetics consortium families. *Diabetes* **57**, 1084–1092 (2008).
- 861 40. Onda, Y. *et al.* Incidence and prevalence of childhood-onset Type 1 diabetes in Japan:
862 the T1D study. *Diabet. Med.* **34**, 909–915 (2017).
- 863 41. Sivertsen, B., Petrie, K. J., Wilhelmsen-Langeland, A. & Hysing, M. Mental health in
864 adolescents with Type 1 diabetes: Results from a large population-based study. *BMC*

- 865 *Endocr. Disord.* **14**, 1–8 (2014).
- 866 42. Thomson, G. *et al.* Relative predispositional effects of HLA class II DRB1-DQB1
867 haplotypes and genotypes on type 1 diabetes: a meta-analysis. *Tissue Antigens* **70**,
868 110–127 (2007).
- 869 43. Miyadera, H. & Tokunaga, K. Associations of human leukocyte antigens with autoimmune
870 diseases: Challenges in identifying the mechanism. *J. Hum. Genet.* **60**, 697–702 (2015).
- 871 44. Cucca, F. A correlation between the relative predisposition of MHC class II alleles to type
872 1 diabetes and the structure of their proteins. *Hum. Mol. Genet.* **10**, 2025–2037 (2001).
- 873 45. Zhu, M. *et al.* Identification of novel T1D risk loci and their association with age and islet
874 function at diagnosis in autoantibody-positive T1D individuals: Based on a two-stage
875 genome-wide association study. *Diabetes Care* **42**, 1414–1421 (2019).
- 876 46. Wang, H. yu *et al.* Risk HLA class II alleles and amino acid residues in
877 myeloperoxidase–ANCA-associated vasculitis. *Kidney Int.* **96**, 1010–1019 (2019).
- 878 47. Kachooei-mohaghegh-yaghoobi, L., Rezaei-rad, F. & Zamani, M. The impact of the HLA
879 DQB1 gene and amino acids on the development of narcolepsy. *Int. J. Neurosci.* **0**, 1–8
880 (2020).
- 881 48. Kawasaki, E. & Eguchi, K. Is type 1 diabetes in the Japanese population the same as
882 among Caucasians? *Ann. N. Y. Acad. Sci.* **1037**, 96–103 (2004).
- 883 49. Li, Y. R. & Keating, B. J. Trans-ethnic genome-wide association studies: Advantages and
884 challenges of mapping in diverse populations. *Genome Med.* **6**, 1–14 (2014).
- 885 50. Lee, C. H., Eskin, E. & Han, B. Increasing the power of meta-analysis of genome-wide
886 association studies to detect heterogeneous effects. *Bioinformatics* **33**, i379–i388 (2017).
- 887 51. Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout□:

- 888 A Simple Way to Prevent Neural Networks from Overfitting. *J. Mach. Learn. Res.* **15**,
889 1929–1958 (2014).
- 890 52. Ioffe, S. & Szegedy, C. Batch normalization: Accelerating deep network training by
891 reducing internal covariate shift. *Proc. ICML* 448–456 (2015).
- 892 53. Kingma, D. & Ba, J. Adam: A method for stochastic optimization. *Int. Conf. Learn.*
893 *Represent.* (2015).
- 894 54. Sener, O. & Koltun, V. Multi-task learning as multi-objective optimization. *Adv. Neural Inf.*
895 *Process. Syst.* **2018-Decem**, 527–538 (2018).
- 896 55. Shimura, K., Li, J. & Fukumoto, F. HFT-CNN: Learning Hierarchical Category Structure
897 for Multi-label Short Text Categorization. 811–816 (2019) doi:10.18653/v1/d18-1093.
- 898 56. Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A Next-generation
899 Hyperparameter Optimization Framework. *Proc. ACM SIGKDD Int. Conf. Knowl. Discov.*
900 *Data Min.* 2623–2631 (2019) doi:10.1145/3292500.3330701.
- 901 57. Silver, N. C. & Dunlap, W. P. Averaging correlation coefficients: Should Fisher’s z
902 transformation be used? *J. Appl. Psychol.* **72**, 146–148 (1987).
- 903 58. Zheng, X. *et al.* HIBAG - HLA genotype imputation with attribute bagging.
904 *Pharmacogenomics J.* **14**, 192–200 (2014).
- 905 59. Abi-Rached, L. *et al.* Immune diversity sheds light on missing variation in worldwide
906 genetic diversity panels. *PLoS One* **13**, e0206512 (2018).
- 907 60. Nagai, A. *et al.* Overview of the BioBank Japan Project: Study design and profile. *J.*
908 *Epidemiol.* **27**, S2–S8 (2017).
- 909 61. Hirata, M. *et al.* Cross-sectional analysis of BioBank Japan clinical data: A large cohort of
910 200,000 patients with 47 common diseases. *J. Epidemiol.* **27**, S9–S21 (2017).

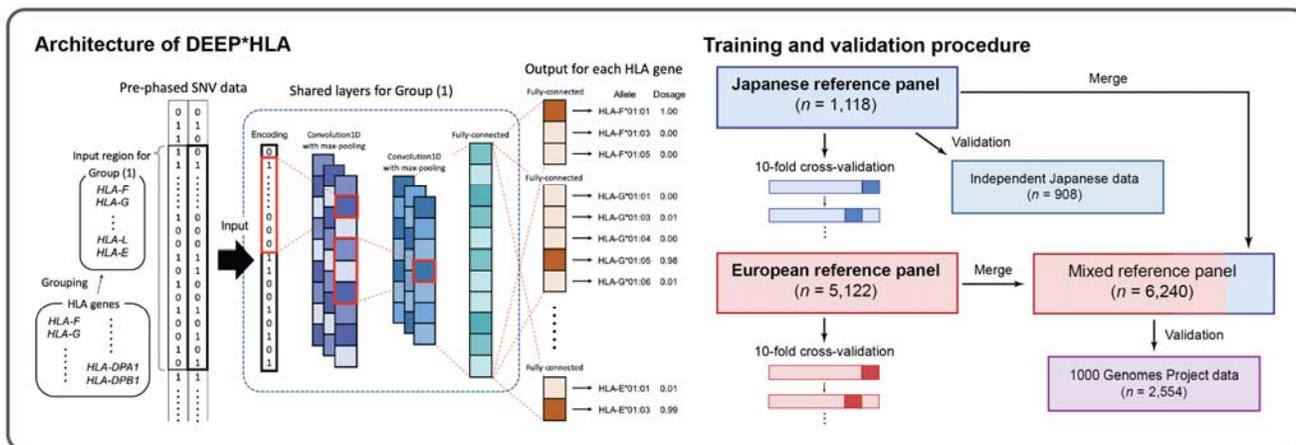
- 911 62. Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese population links cell
912 types to complex human diseases. *Nat. Genet.* **50**, 390–400 (2018).
- 913 63. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a
914 Wide Range of Complex Diseases of Middle and Old Age. *PLoS Med.* **12**, 1–10 (2015).
- 915 64. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
916 *Nature* **562**, 203–209 (2018).
- 917

918 **Figure Legends**

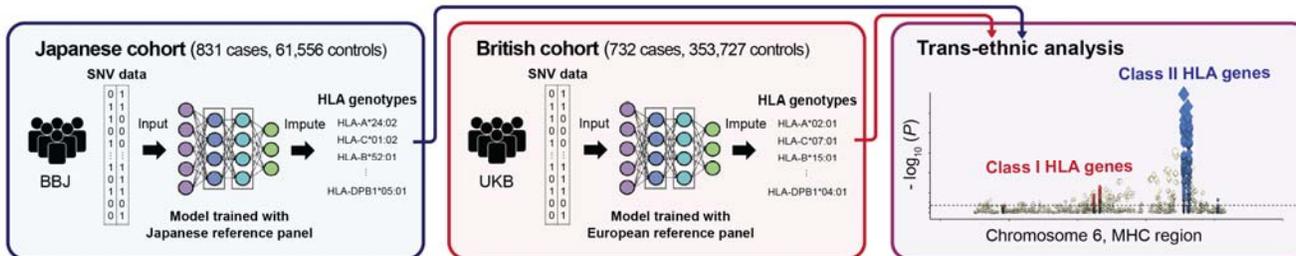
919

920 **Figure 1. An overview of the study**

a Constructing models with HLA references panels



b MHC fine-mapping in T1D GWAS data of biobanks



921

922 (a) DEEP*HLA is a deep learning architecture that takes an input of pre-phased genotypes of

923 SNVs and outputs the genotype dosages of HLA genes. To train a model and benchmark its

924 performance, we used Japanese and European HLA reference panels respectively. We

925 evaluated its accuracies in cross-validation with other methods. For the Japanese panel, we

926 also evaluated its accuracy by applying the trained model to an independent Japanese HLA

927 dataset. Further, we experimentally generated a mixed panel and validated its accuracy using

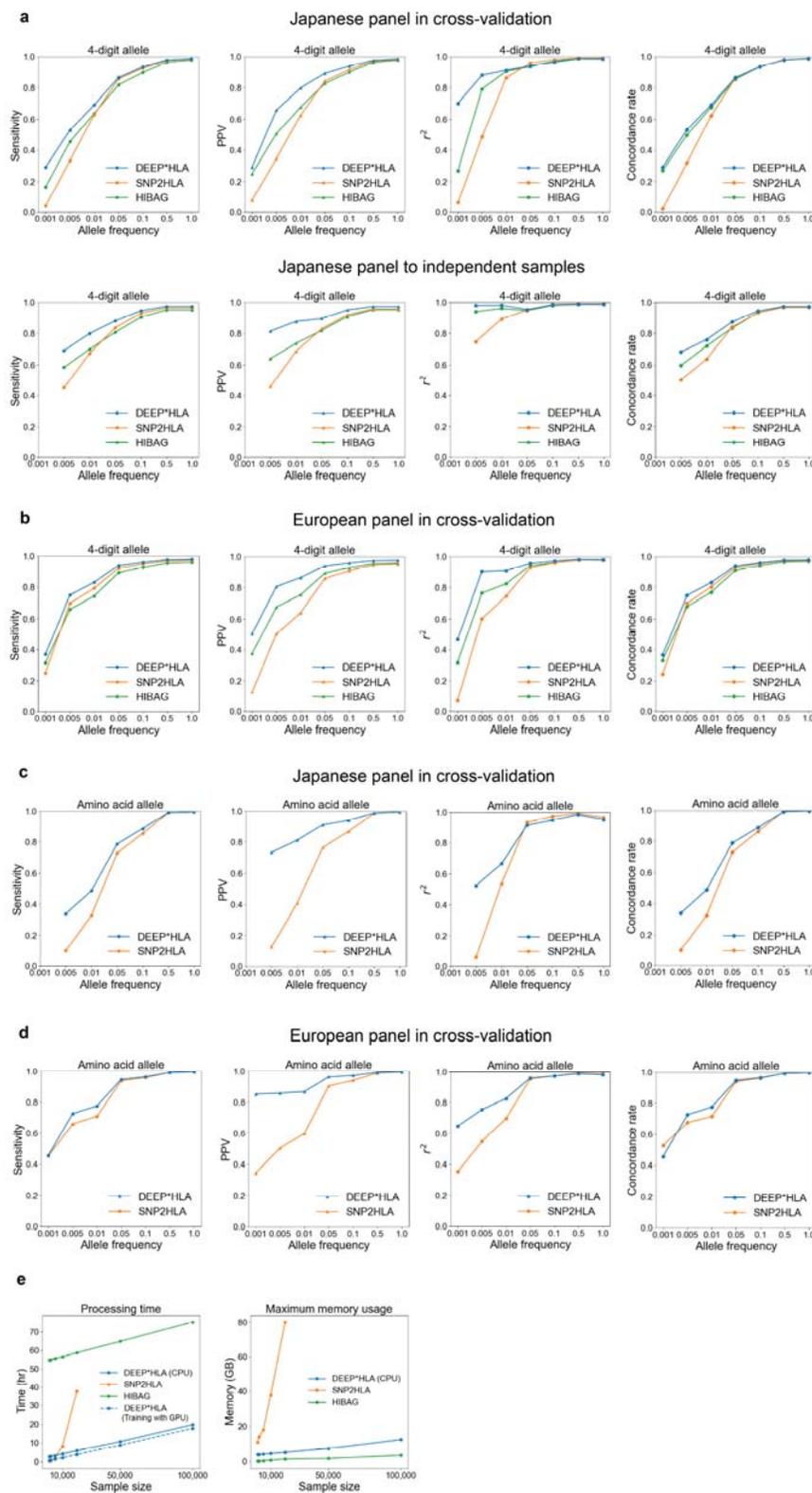
928 1KGV3 data. (b) We conducted trans-ethnic MHC fine-mapping in T1D GWAS data. We

929 performed HLA imputation for the Japanese cohort from BBJ and the British cohort from UKB

930 using models specific for individual populations. We integrated the individual results of imputed
931 genotypes and performed trans-ethnic association analysis.

932

933 **Figure 2. Performance evaluations of HLA imputation methods**

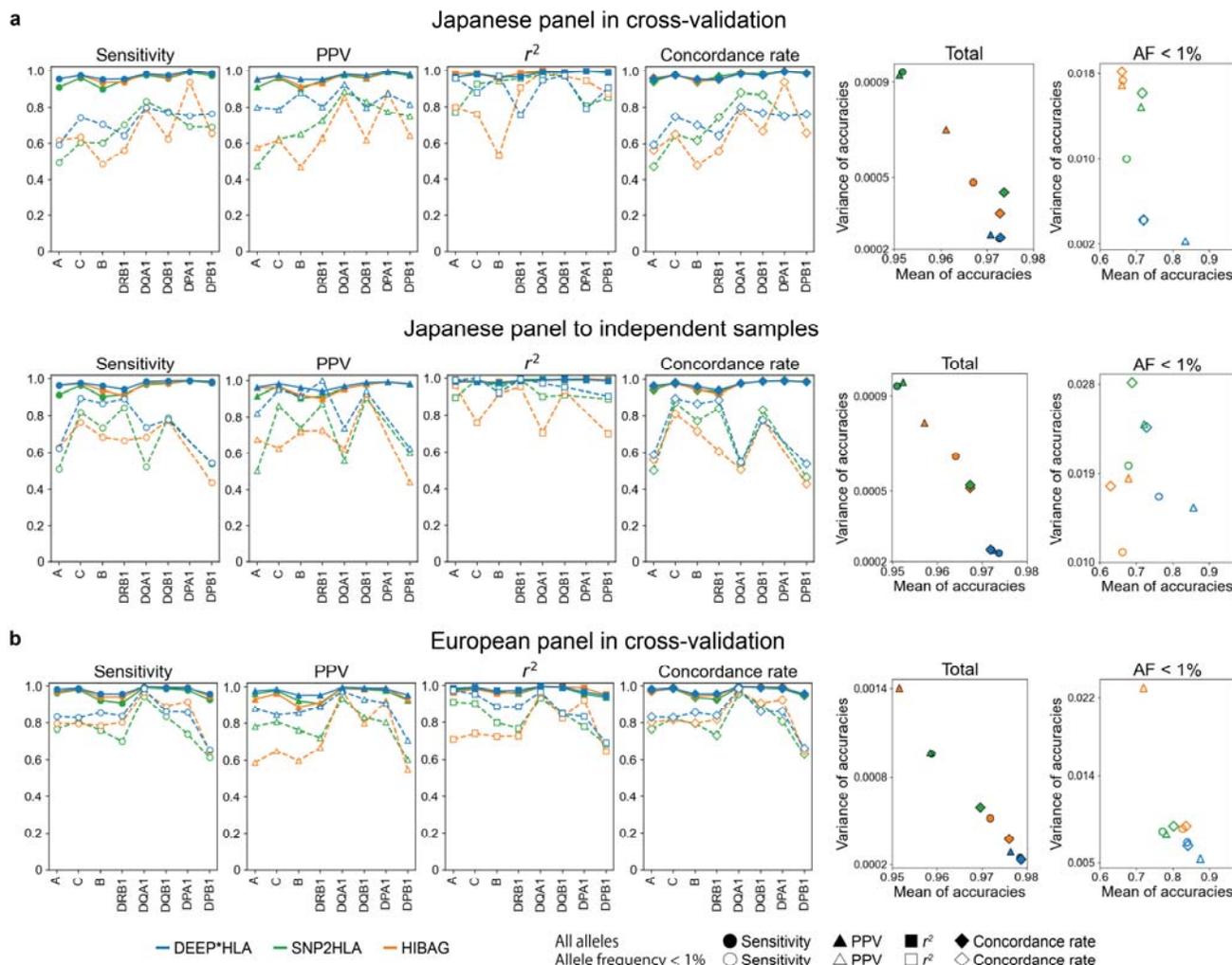


934

935 **(a-d)** Sensitivity, PPV, and r^2 of allelic dosage and concordance rate of best-guess genotypes
936 for the 4-digit alleles **(a, b)** and amino acid polymorphisms **(c, d)** evaluated in our Japanese
937 reference panel **(a, c)** and the T1DGC reference panel **(b, d)**. For each metric, mean values of
938 alleles with a frequency less than a value on the horizontal axis are shown on the vertical axis.
939 DEEP*HLA was advantages especially for rare alleles. **(e)** Processing time (left) and maximum
940 memory usage (right) evaluated on imputing BBJ samples using the Japanese panel.
941 DEEP*HLA imputed by far the fastest in total processing time as the sample size increased.
942 The dashed blue line in the processing time represents a case when DEEP*HLA used GPU
943 only in training a model. All methods exhibited maximum memory usage scaling roughly linearly
944 with sample size. SNP2HLA did not work within 100 GB in our machine for the sample sizes
945 greater than 20,000.

946

947 **Figure 3. Comparison of imputation accuracy between different HLA genes**



948

949 Each panel represents accuracy in 8 classical HLA genes evaluated in the Japanese panel in

950 cross-validation (a, upper), the Japanese panel to the independent data (a, lower), and the

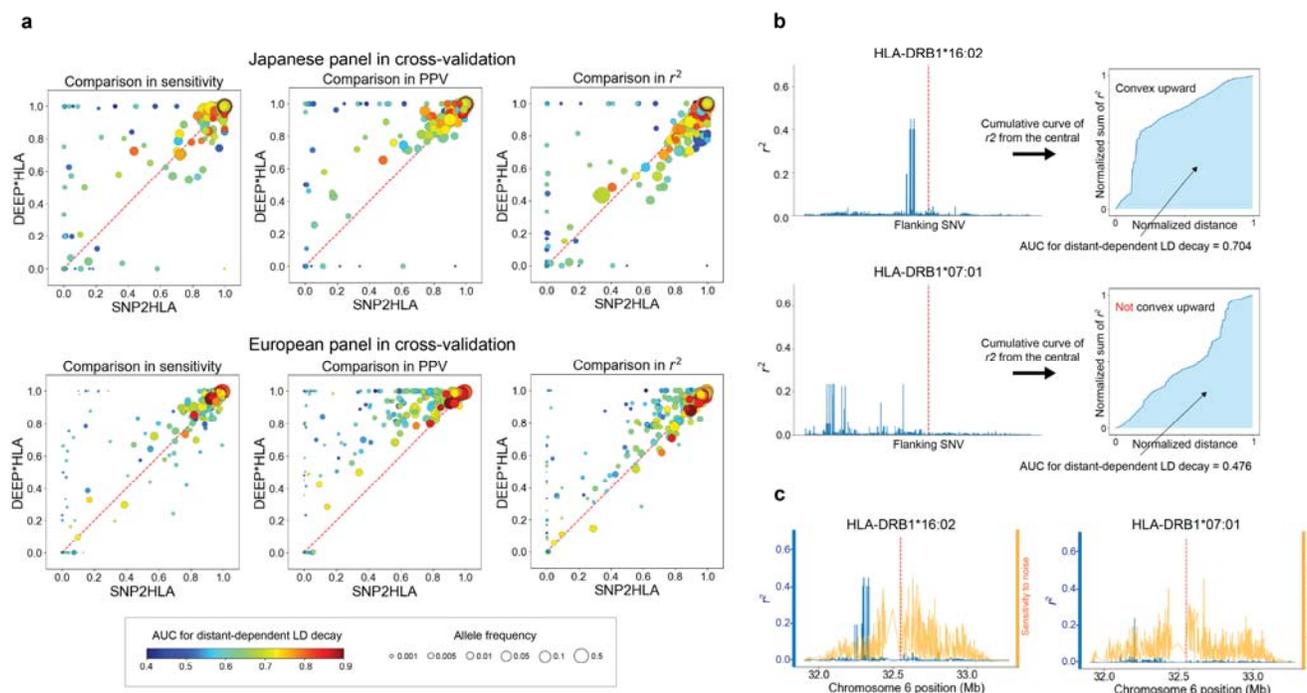
951 European panel in cross-validation (b). Solid and dashed lines correspond to the accuracy of all

952 allele and allele with frequency < 1%, respectively. The right two scatter plots represent the

953 relation between the mean and variance of each metric among different HLA genes for

954 individual methods. R^2 metric is not shown because it is not an additive statistic.

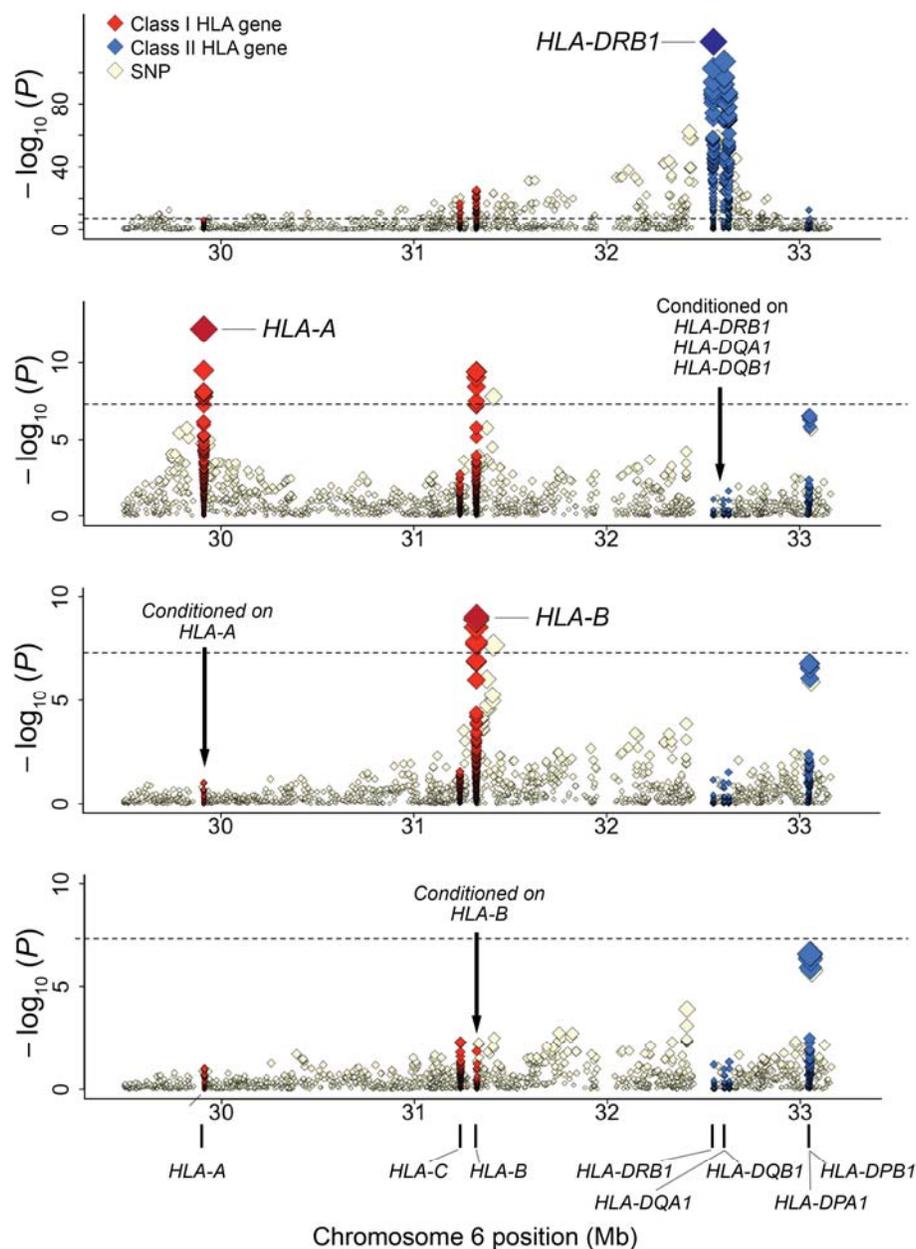
955 **Figure 4. Comparison between DEEP*HLA and SNP2HLA displayed with allele**
 956 **frequencies and AUC for distance-dependent LD decay**



957
 958 **(a)** Comparisons of imputation accuracy between DEEP*HLA and SNP2HLA for 4-digit allele
 959 imputation for cross-validation with the Japanese panel (upper) and T1DGC panels (lower).
 960 Each dot corresponds to one allele, displayed with allele frequencies (size) and AUC for
 961 distance-dependent LD decay (color). The AUC was calculated based on bilateral 1,000 SNPs.
 962 Comparisons in concordance rate are not shown because they were almost the same as those
 963 in sensitivity. The performance of SNP2HLA was limited when imputing the alleles with
 964 low-frequency and low AUC; DEEP*HLA was relatively accurate even for the less frequent
 965 alleles regardless of AUC. **(b)** Example illustrations of AUC for distance-dependent LD decay.
 966 The left figures illustrate r^2 of LD between an HLA allele (red dash line in the central) and
 967 flanking SNVs. HLA-DRB1*16:02 has strong LD in close positions and weaker LD in the distant
 968 positions. The cumulative curve of r^2 of bilateral SNVs becomes convex upward; and the AUC

969 increases. In contrast, HLA-DRB1*07:01 has moderate LD in distant or sparse positions, the
970 curve does not become convex upward, and the AUC becomes smaller. (c) Comparison
971 between r^2 (blue line) and sensitivity maps of DEEP*HLA (orange line) for example alleles (red
972 dash line in the center). The sensitivities are normalized for visibility. In both examples,
973 DEEP*HLA reacted to noise across an extensive area regardless of LD.
974

975 **Figure 5. Trans-ethnic association plots of HLA variants with T1D in the MHC region.**



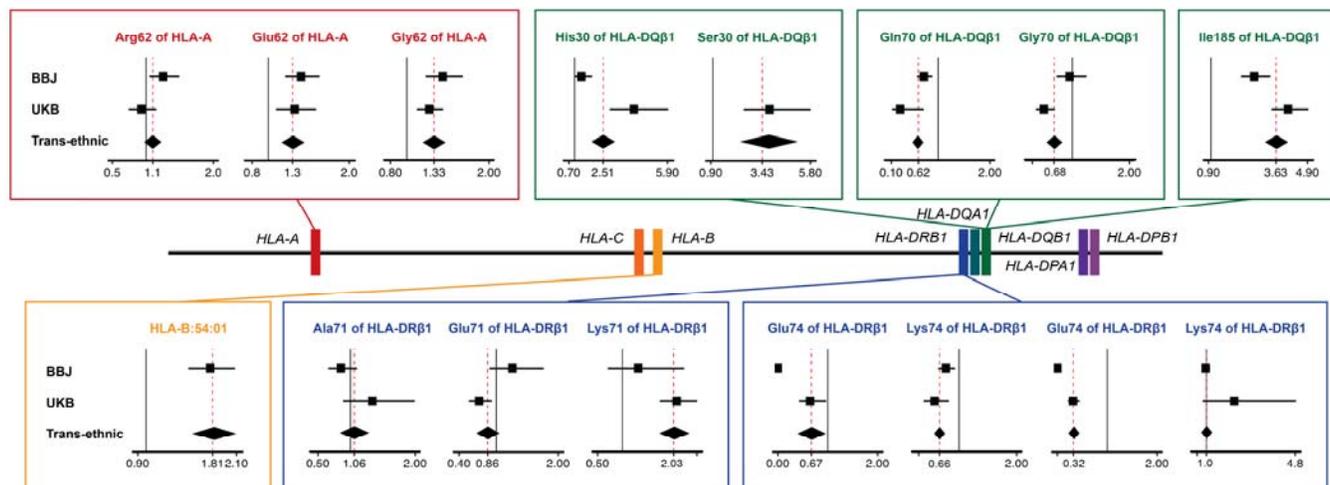
976

977 Diamonds represent $-\log_{10}(P)$ values for the tested HLA variants, including SNPs, classical
978 alleles, and amino acid polymorphisms of the HLA genes. Dashed black horizontal lines
979 represent the genome-wide significance threshold of $P = 5.0 \times 10^{-8}$. The physical positions of
980 the HLA genes on chromosome 6 are shown at the bottom. (a-e) Each panel shows the

981 association plot in the process of stepwise conditional regression analysis: nominal results. **(a)**
982 Results conditioned on *HLA-DRB1*, *HLA-DQA1*, and *HLA-DRB1*. **(b)** Results conditioned on
983 *HLA-DRB1*, *HLA-DQA1*, *HLA-DRB1*, and *HLA-A*. **(c)** Results conditioned on *HLA-DRB1*,
984 *HLA-DQA1*, *HLA-DRB1*, *HLA-A*, and *HLA-B*. **(d)** Our study identified the independent
985 contribution of multiple HLA class I and class II genes to the T1D risk in a trans-ethnic cohort, in
986 which the impacts of class II HLA genes were more evident. Detailed association results are
987 shown in **Supplementary Table 3**.

988

989 **Figure 6. HLA variants associated with the T1D risk identified through trans-ethnic**
 990 **fine-mapping.**



991
 992 Forest plots for individual risk-associated alleles are displayed along with a location map of
 993 classical HLA genes. Each forest plot shows the estimated odds ratio (OR) and 95% confidence
 994 interval from cohort-specific logistic model for BBJ and UKB, and the trans-ethnic logistic model.
 995 Red dashed lines indicate OR in trans-ethnic cohorts. Black solid lines represent OR = 1.
 996 Colored square boxes represent amino acid polymorphisms of the same position or a classical
 997 allele.

998

999 **Tables 1. Associations of the HLA variants with the T1D risk identified through**
 1000 **trans-ethnic fine-mapping study.**

HLA variant	Frequency (BBJ)		Frequency (UKB)		OR (95% CI)		P†	
	Case n = 831	Control n = 61,556	Case n = 732	Control n = 353,727	BBJ	UKB	BBJ	UKB
HLA-DRβ1 amino acid position 71								
Alanine	0.10	0.18	0.04	0.15	0.85 (0.66-1.10)	1.34 (0.89-1.99)	0.23	0.16
Arginine	0.82	0.73	0.33	0.45	(reference)			
Glutamic acid	0.073	0.074	0.083	0.12	1.26 (0.89-1.77)	0.72 (0.56-0.93)	0.019	0.0013
Lysine	0.0096	0.011	0.54	0.28	1.31 (0.71-2.24)	2.11 (1.77-2.53)	0.035	1.9 × 10 ⁻¹⁶
HLA-DQβ1 amino acid position 185								
Isoleucine	0.39	0.57	0.68	0.83	2.74 (2.21-3.40)	4.12 (3.49-4.99)	3.5 × 10 ⁻²⁰	7.0 × 10 ⁻⁵⁵
Threonine	0.61	0.43	0.32	0.17	(reference)			
HLA-DQβ1 amino acid position 30								
Histidine	0.16	0.19	0.18	0.23	1.36 (0.97-1.93)	4.16 (2.86-5.96)	0.0078	3.0 × 10 ⁻¹⁴
Serine	0.0042	0.0038	0.34	0.25	inf	3.82 (2.53-5.87)	0.079	3.8 × 10 ⁻¹⁰
Tyrosine	0.83	0.80	0.48	0.52	(reference)			
HLA-DRβ1 amino acid position 74								
Alanine	0.56	0.59	0.59	0.65	(reference)			
Arginine	0.0018	0.00088	0.28	0.15	0 (0-0.045)	0.64 (0.42-0.96)	0.08	0.0036
Glutamic acid	0.32	0.27	0.021	0.036	0.77 (0.64-0.93)	0.57 (0.38-0.82)	0.00065	0.0004
Glutamine	0.0024	0.0030	0.0795	0.15	0 (0-0.0029)	0.31 (0.21-0.44)	0.079	4.5 × 10 ⁻¹⁰
Leucine	0.12	0.14	0.023	0.023	0.97 (0.81-1.16)	2.20 (0.85-4.84)	0.074	0.0077
HLA-DQβ1 amino acid position 70								
Arginine	0.60	0.62	0.79	0.63	(reference)			
Glutamic acid	0.26	0.17	0.020	0.020	0.73 (0.59-0.9)	0.27 (0.11-0.71)	0.00020	0.0052
Glycine	0.14	0.20	0.19	0.35	0.95 (0.72-1.25)	0.50 (0.36-0.70)	0.073	3.1 × 10 ⁻⁵
HLA-A amino acid position 62								
Arginine	0.19	0.20	0.06	0.09	1.25 (1.05-1.49)	0.93 (0.74-1.16)	0.0012	0.53
Glutamic acid	0.39	0.37	0.09	0.09	1.40 (1.21-1.63)	1.33 (1.10-1.60)	9.2 × 10 ⁻⁶	0.0025
Glutamine	0.15	0.19	0.46	0.49	(reference)			
Glycine	0.26	0.24	0.33	0.29	1.44 (1.23-1.68)	1.27 (1.12-1.44)	6.6 × 10 ⁻⁶	1.6 × 10 ⁻⁴
Leucine	0	0	0.055	0.044	-	2.01 (1.57-2.55)	1.5 × 10 ⁻¹²	1.9 × 10 ⁻⁸
HLA-B*54:01	0.14	0.073	0	0	1.78 (1.51-2.08)	-	-	-

HLA, human leucocyte antigen; OR, odds ratio; 95% CI, 95% confidence interval.

†Obtained from the multivariate regression model that included all the variants listed here.