1 _____

# Quantification of inter-sample differences in T cell receptor sequences

**Ryo Yokota** [1,2*], **Yuki Kaminaga** [2], **and Tetsuya J. Kobayashi** [1,2,3]

[1] *Institute of Industrial Science, the University of Tokyo, Tokyo, Japan*

[2] *Department of Electrical Engineering and Information Systems, Graduate School of Engineering, the University of Tokyo, Tokyo, Japan*

[3] *PRESTO, Japan Science and Technology Agency (JST), Saitama, Japan*

Correspondence*:

Ryo Yokota

Institute of Industrial Science, the University of Tokyo, 4-6-1 Meguro-ku, Komaba, Tokyo 153-8505, Japan, yokota@sat.t.u-tokyo.ac.jp

2 **ABSTRACT**

3    Inter-sample comparisons of the T cell receptor (TCR) repertoire are crucial for gaining a better

4  understanding into the immunological states determined by different collections of T cells from

5  different donor sites, cell types, and genetic and pathological backgrounds. As a theoretical

6  approach for the quantitative comparison, previous studies utilized the Poisson abundance models

7  and the conventional methods in ecology, which focus on the abundance distribution of observed

8  TCR sequences. However, these methods ignore the details of the measured sequences and are

9  consequently unable to identify sub-repertoires that might have the contributions to the observed

10  inter-sample differences. In this paper, we propose a new comparative approach based on TCR

11  sequence information, which can estimate the low-dimensional structure by projecting the pairwise

12  sequence dissimilarities in high-dimensional sequence space. The inter-sample differences are

13  then quantified according to information-theoretic measures among the distributions of data

14  estimated in the embedded space. Using an actual dataset of TCR sequences in transgenic

15  mice that have strong restrictions on somatic recombination, we demonstrate that our proposed

16  method can accurately identify the inter-sample hierarchical structure, which is consistent with

17  that estimated by previous methods based on abundance or count information. Moreover, we

18  identified the key sequences that contribute to the pairwise sample differences. Such identification

19  of the sequences contributing to variation in immune cell repertoires may provide substantial

20  insight for the development of new immunotherapies and vaccines.

21  **Keywords: T cell, TCR repertoire, inter-sample comparison, pairwise sequence alignment, sequence dissimilarity, manifold learning,**

22  **Jensen-Shannon Divergence**

## 1   INTRODUCTION

23  The development of high-throughput sequencing with next-generation sequencers has provided new

24  opportunities to quantify T cell receptor (TCR) repertoires and to compare their differences among different

25  cell types, organisms, and pathological samples. Such information is indispensable for quantitatively

26  understanding the immunological state of organisms that is shaped by the collection of immune cells.

27  Moreover, the detailed information of TCR repertoires, especially that of inter-sample differences, is

28  anticipated to significantly promote the development of immunotherapies and vaccines (Hou et al., 2016).

29  To this end, several theoretical methods have been proposed to quantify sample differences by focusing on

30  the count (abundance) distribution of unique TCR sequences in a repertoire (Greiff et al., 2015; Laydon

31  et al., 2015; Hou et al., 2016). Poisson abundance (PA) models are among the recently developed methods

32 based on the hierarchical Bayesian inference algorithm, which estimates the parameters of the models

33 from experimental data and defines the inter-sample difference according to the deviation of the estimated

34 parameters. This method overcomes the substantial sampling fluctuations derived from the huge diversity in

35 TCR repertoires, and provides a stable result related to the inter-sample distances on the basis of statistical

36 interpretations. For example, Rempala et al. (2011) used a bivariate Poisson log-normal (BPLN) model to

37 classify eight different samples of the following sample conditions: donor sites, types of T cells, and the

38 genetic backgrounds of different mouse lines. Guindani et al. (2014) used a Poisson Dirichlet process to

39 classify the types of T cells (i.e., conventional and regulatory T cells). Besides the above examples, other

40 variations of PA models have been proposed for sample classification based on the measurement of TCR

41 diversity (Sepúlveda et al., 2010; Greene et al., 2013).

42    Although these PA models can successfully quantify the inter-sample distances, they are also associated

43 with a major drawback in that some of the sequence information for each sample is lost since these models

44 focus only on the count distribution. This loss of information has hampered the ability to determine the

45 characteristic sequences of each sample, which is a requisite for further investigations of the source of the

46 difference by, for example, evaluations of the interaction with microbial peptides (Aas-Hanssen et al., 2015)

47 and the simulation of TCR crystal structures (Klausen et al., 2015).

48    As an alternative method, we can extract and count the overlapping sequences between two samples or

49 among many samples. However, the possible sequence space of the TCR repertoire is at least as large as

50 $10^{15}$ (Davis and Bjorkman, 1988), and therefore the measured sequences can only sparsely cover the entire

51 space. This sparsity substantially reduces the chance to observe the same sequence in two samples. Thus, by

52 focusing on the overlaps, it is only possible to detect the public sequences that appear very frequently

53 among the samples. Moreover, even if no overlapping among the sequences is detected, it is not possible to

54 judge whether this occurs because the two repertoires cover quite different subspaces of the sequence space

55 or because the repertoires cover the same subspace but show no overlapping by chance simply owing to the

56 sparsity of the coverage. This difference can be determined by exploiting the information of inter-sequence

57 differences in the repertoires.

58     To address this problem, we developed a new method for the quantitative comparison of TCR repertoires,

59 by focusing on the sequence information in all samples, and estimating the low-dimensional structure

60 (manifold) by projecting the high-dimensional inter-sequence relations, calculated from pairwise sequence

61 alignments, onto a low-dimensional space. The methods for manifold estimation have been successfully

62 applied in previous studies of virus evolution (Ito et al., 2011) and relationships of 16S rRNA gene sequences

63 in bacterial genomes (Hughes et al., 2012) to extract the evolutionary pathways and interconnections of

64 bacteria. Although manifold estimation has also been employed for evaluating the TCR repertoire, this was

65 mainly used only for visualization purposes (Duez et al., 2016). However, the low-dimensional embedding

66 of the original repertoire contains the information how the repertoires from different samples cover the

67 possible sequence space. Therefore, by employing such information, it may be possible to detect a subset of

68 sequences in the repertoire that has a major contribution to the inter-sample difference.

69     To quantitatively compare the embedded sequences, we estimated a probability density function of

70 the sequence distribution in the low-dimensional space. This density estimation compensates for the

71 sampling bias due to unseen sequences from the sparsity of the measured sequences. Finally, we quantified

72 the inter-sample differences between the estimated density functions of the individual samples using

73 the Jensen-Shannon divergence (JSD). This information-theoretic measure characterizes the difference

74 between two distributions by the probability of observing either one by chance with random sampling

75 from the others. Thus, this measure can effectively and quantitatively capture information on the existence

76 of non-overlapping sequences between two repertoires. By extracting the sequences that show a major

77 contribution to the information theoretic measures, the sequences most responsible for the inter-sample

78 differences can be determined, which cannot be identified with previous approaches.

79 The paper is structured as follows. We first describe the experimental data adopted to test our method

80 and the step-by-step data analysis procedure, including (i) quantification of sequence dissimilarity with

81 the pairwise sequence alignment algorithm, (ii) evaluations of four different manifold learning methods

82 for projecting the sequence distribution in low-dimensional space, (iii) adoption of the kernel density

83 estimation algorithm (KDE) to quantify the sequence distribution, and (iv) quantification of the inter-sample

84 differences and identification of the contributing sequences according to the JSD values of the distributions.

85 To validate the applicability of our method, we apply a true dataset of TCR repertoires and demonstrate that

86 similar inter-sampling clustering can be obtained by both our method and previous methods despite their use

87 of different modalities (sequence and count, respectively) of a repertoire. We further evaluate the statistical

88 significance of our results using a bootstrap algorithm to confirm the derived sample difference. Overall, we

89 aim to demonstrate the advantages of our method to previous methods by capturing more complete and

90 quantitative information on TCR repertoires. This method is expected to be of value for understanding

91 variation of the immunological states to facilitate development of immunotherapies and vaccines.

## 2  MATERIAL & METHODS

### 2.1  Sequence data

93 In this study, we used a public dataset of TCR repertoires in mice published in Rempala et al. (2011). This

94 dataset includes information of eight different TCR populations, which are classified according to donor

95 sites, types of CD4+ T-cells, and the genetic backgrounds of mice. The CD4+ T-cells were collected and

96 isolated either from the thymus or peripheral lymph nodes, which are labeled as "1" and "2", respectively. In

97  addition, these cells were categorized into either naive T-cells (TN) or regulatory T-cells (TR) in accordance

98  with the presence or absence of Foxp3 expression. The two genetic backgrounds of mice are labeled as

99  "wild type" and "Ep" in the original paper. Both groups showed strong restriction on rearrangement of

100  V(D)J genes (i.e., the two $\alpha$-chain rearrangements between J$\alpha$2.6 and J$\alpha$2 with a fixed V$\alpha$2.9 segment

101  and fixed $\beta$-chain V$\beta$14D$\beta$2J$\beta$2.6). The main difference between these groups is that the Ep mice were

102  backcrossed with mice that express transgenic class 2 major histocompatibility complex molecules bound to

103  a single "Ep" peptide (Pacholczyk et al., 2006). Thus, Ep mice are expected to show a more restricted

104  TCR repertoire than wild-type mice. To evaluate the diversity of TCR repertoire, the complementarity

105  determining region 3 (CDR3) of TCR$\alpha$ chains were sequenced and amplified. Further details on this dataset

106  are described in Section 4 of Rempala et al. (2011).

107  ## 2.2 Data analysis procedure

108  To date, no sufficient and effective method for comparing TCR repertoires among different samples

109  has been established, which is mainly due to the enormous complexity and diversity of TCR sequences

110  (Hou et al., 2016). In ecology, the diversity in a population is conventionally measured with metrics of

111  species abundance' between pairs of samples. However, these methods generally rely on observational

112  data of species abundance counts, and can therefore be vulnerable to sampling bias. One widely used

113  typical measure to quantify biological diversity is through a dissimilarity metric such as the Bray-Curtis

114  index (Bray and Curtis, 1957; Tang et al., 2016; Silverman et al., 2016). In the context of TCR repertories,

115  the abundance counts of observed sequences' can be considered to be analogous to species abundance'

116  in an ecological context. However, because of the sparsity of observed sequences, application of these

117  dissimilarity metrics to a dataset of TCR repertoires may not always work well. PA models have attracted

118  substantial attention as methods to overcome these issues, since these models can compensate for the

119    sampling bias by estimating the statistical parameters directly from the measures of the abundance counts of

120    observed sequences (Robinson and Smyth, 2007; Sepúlveda et al., 2010), and are thus expected to be more

121    robust to sampling bias. However, both methods ignore detailed information of the amino acid sequences of

122    TCRs.

123    Therefore, in this study, rather than focusing on observation counts, we instead focus on the sequence

124    similarity among repertoires. Our method consists of four steps: (i) calculate a dissimilarity matrix of

125    observed TCR sequences in all samples using the Smith-Waterman (SW) algorithm with a scoring matrix;

126    (ii) embed the data in a low-dimensional Euclidian space by dimensionality reduction methods while

127    preserving the inter-sequence relations quantified by the dissimilarity matrix; (iii) estimate the sequence

128    distributions in the low-dimensional space by the KDE algorithm; and (iv) quantify the sample differences

129    by calculating the JSD value between the probabilistic density functions of different samples, and cluster

130    the samples accordingly. Each of the above steps is described in detail in the following subsections.

131    2.2.1    Quantification of sequence dissimilarity

132    The first step of our method is the quantification of similarity for each pair of TCR sequences in all

133    samples. The SW algorithm remains the most popular pairwise local sequence alignment algorithm in

134    bioinformatics for quantifying the similarity of amino acid sequences (Smith and Waterman, 1981). In

135    recent years, improved versions of the SW algorithm have been proposed to resolve the problems related to

136    the increase in computational costs along with the rapidly increasing size of datasets that are now possible

137    from next-generation sequencing. Here, we used one of these modifications, the striped SW algorithm

138    (Farrar, 2007), which uses a single-instruction-multiple-data (SIMD) system that allows for multiple units to

139    simultaneously execute the same operation. The algorithm was implemented with Parasail, an open-source

140    software for sequence alignment (Daily, 2016).

141      The SW algorithm requires amino acid substitution matrices, which determine the cost of the replacement

142   of a single amino acid residue by another (Henikoff and Henikoff, 1992). Although the SW algorithm has

143   already been applied to TCR sequences as a mapping tool for CDR3 sequences (Shugay et al., 2014), no

144   study has yet established the best choice of substitution matrices for comparison of TCR data. Therefore, to

145   clarify the effect of the type of substitution matrix employed and determine the optimal choice for our

146   method, we tested 10 different matrices: five different point-accepted mutation matrices (PAM; 30, 100,

147   120, 160, and 250) (States et al., 1991), and five different blocks substitution matrices (BLOSUM; 45, 50,

148   62, 80, and 100) (Henikoff and Henikoff, 1992). The gap opening and extension penalties were set to 10

149   and 1, respectively (Farrar, 2007).

150      Since the substitution matrices give nonzero values for replacements between the same amino acid

151   residues, the total score of the alignment between two identical sequences depends on their sequence lengths.

152   Thus, the diagonal elements of a pairwise distance matrix will have different values even when they are

153   calculated from the alignments of two identical sequences. In other words, both the sequence similarity and

154   the sequence length determine the values of the pairwise distance matrix. To adjust for this sequence-length

155   effect, we converted the pairwise distance matrix into a dissimilarity matrix using the following equation:

$$S_{i,j} = 1 - \frac{2D_{i,j}}{D_{i,i} + D_{j,j}} \, , \tag{1}$$

156   where $D_{i,j}$ and $S_{i,j}$ are a pairwise distance matrix and dissimilarity matrix between the two sequences $i$ and

157   $j$, respectively. At this step, we calculated the pairwise distances between all pairs of unique sequences

158   observed in all samples with the striped SW algorithm. We then transformed the pairwise distance matrix

159   into the dissimilarity matrix using Eq. 1.

160   2.2.2   Dimensionality reduction with manifold learning methods

161   To visualize the structure of the high-dimensional dissimilarity matrix in a low-dimensional space, we

162   applied dimensionality reduction (manifold learning) techniques to the dissimilarity matrix described

163   above that was constructed with BLOSUM62. Here, we compared the results calculated with four different

164   methods: multidimensional scaling (MDS) (Borg and Groenen, 2005), ISOMAP (Tenenbaum et al., 2000),

165   spectral embedding (SE) (Belkin and Niyogi, 2001), and t-distributed stochastic neighbor embedding

166   (t-SNE) (Van Der Maaten and Hinton, 2008). All of these methods transform the dissimilarity matrix $S$

167   with dimensionality $N$ into a new dataset $Y$ with a lower dimensionality $d$ in such a way as to preserve

168   the structure of the dissimilarity matrix by minimizing cost functions. The major difference among these

169   methods is the cost function, which is determined according to the relative distances between all pairs of

170   sequences. MDS with a SMACOF algorithm minimizes the sum of squared errors in the relative distances

171   of all sequence pairs before and after projections (Borg and Groenen, 2005). This cost function of MDS

172   tends to preferentially retain the distances between more distant data points over those between more

173   adjacent points (Van Der Maaten et al., 2009). ISOMAP also minimizes the sum of squared errors, but rather

174   than using the relative distances, it uses the geodesic distances, which are the distances along the shortest

175   paths between two nodes on the neighborhood graph, calculated with a k-nearest neighbor algorithm

176   (Tenenbaum et al., 2000). In the present study, we calculated the geodesic distances with the Warshall-Floyd

177   algorithm (Floyd, 1962). ISOMAP retains a neighborhood structure of data points lying on a curved

178   manifold (e.g., the Swiss roll dataset (Tenenbaum et al., 2000)), which is collapsed in MDS. SE, also known

179   as Laplacian eigenmaps, minimizes the cost function based on the neighborhood graph, which ensures that

180   local neighborhood relations in a high-dimensional space are preserved in a embedded low-dimensional

181   space (Belkin and Niyogi, 2001, 2003). We regarded the adjacency matrix based on the k-nearest neighbor

182   algorithm as the weighted graph matrix to construct the Laplacian graph of SE. Finally, t-SNE converts the

183  relative distances to joint probabilities, and minimizes the Kullback-Leibler divergence between the joint

184  probabilities of the high-dimensional space and those of a embedded low-dimensional space (Van Der

185  Maaten and Hinton, 2008). For calculation of the joint probabilities, t-SNE uses different kernels for the

186  high- and low-dimensional spaces: a Gaussian kernel and a Student t-distribution, respectively. Since the

187  Student t-distribution results in heavier tails than the Gaussian kernel, the t-SNE method emphasizes the

188  local distances between data points in the low-dimensional space.

189  In the studies of sequence alignments for sequences with different lengths, it is impossible to know

190  the precise coordinates and the dimension of the sequence space. Thus, we cannot directly use principal

191  components analysis, which is the most widely used dimensionality reduction technique (Bishop, 2007)

192  but requires vector data with fixed dimensionality. The common advantage of the above four methods is

193  that if the distances between all pairs of data points are known, then there is no need to know the specific

194  coordinates of the sequence space (Van Der Maaten and Hinton, 2008; Van Der Maaten et al., 2009).

195  We implemented t-SNE, MDS, and SE with the Scikit-learn manifold learning library with Python

196  (Pedregosa et al., 2012). ISOMAP was implemented with our custom-written code in Python, because the

197  ISOMAP function of the Scikit-learn toolbox does not support the dissimilarity matrix as an argument. The

198  detailed parameters of all methods are described in Table S1 in the Supplementary Information.

199  2.2.3   Estimation of the probability density function with KDE

200  To compare the data points scattered in the embedded low-dimensional space among different samples,

201  the embedded discrete data can be interpolated with a probability density function (PDF). Here, we

202  estimated the PDF with the KDE algorithm (Jones et al., 1996; Heidenreich et al., 2013; Arlot and

203  Celisse, 2010). The exponential function was used as the kernel of the KDE (Christopher et al., 1997).

204  The bandwidth parameter of the exponential kernel function was optimized by maximum-likelihood

205 estimation with a cross-validation algorithm (Arlot and Celisse, 2010). To reduce the computational

206 cost of this calculation, we utilized the Kd-tree algorithm, which is an N-body algorithm that divides

207 all of the data into N clusters based on their relative Euclidean distances (Gray and Moore, 2001). KDE

208 was implemented with the parameter optimization toolbox in Scikit-learn (Pedregosa et al., 2012). For

209 application of the KDE, we discretized the embedded space with 400 bins along each axis with the following

210 range: $[\min x_i - (\max x_i - \min x_i)/10, \max x_i + (\max x_i - \min x_i)/10]$, where $x_i$ indicates the position

211 of a data point (i.e., a sequence) in the embedded space, and $i$ indicates each axis of that space.

212 2.2.4    Quantification of sample differences with JSD

213     The final step of our method involves quantification of the inter-sample differences by calculating the JSD

214 values between all pairs of the estimated PDFs (Elhanati et al., 2014). The JSD is defined as:

$$
\begin{aligned}
D_{JS}[P(\boldsymbol{x})||Q(\boldsymbol{x})] &= \int D_{JS}^{\text{local}} d\boldsymbol{x} \\
&= \int \frac{1}{2}\left\{ P(\boldsymbol{x})\log\frac{P(\boldsymbol{x})}{M(\boldsymbol{x})} + Q(\boldsymbol{x})\log\frac{Q(\boldsymbol{x})}{M(\boldsymbol{x})} \right\} d\boldsymbol{x} \\
&= \frac{1}{2}D_{KL}[P(\boldsymbol{x})||M(\boldsymbol{x})] + \frac{1}{2}D_{KL}[Q(\boldsymbol{x})||M(\boldsymbol{x})] \,,
\end{aligned}
\tag{2}
$$

215 where $P(\boldsymbol{x})$ and $Q(\boldsymbol{x})$ are the estimated PDFs and $D_{KL}$ is the Kullback-Leibler divergence; $M(\boldsymbol{x})$ is

216 $\frac{P(\boldsymbol{x})+Q(\boldsymbol{x})}{2}$ and $D_{JS}^{\text{local}}$ is the "local JSD" , whose integration with respect to $\boldsymbol{x}$ gives the JSD. Thus, the

217 "pairwise" JSDs provide a sample-difference matrix that quantifies the combinatorial differences between

218 all pairs of the samples. To categorize all samples, we utilized hierarchical cluster analysis, which converts

219 the $N \times N$-dimension sample-distance matrix into a dendrogram. Specifically, we used an agglomerative

220 hierarchical cluster technique; each sample is initially treated as a singleton cluster, and pairs of clusters are

221 repetitively merged according to a criterion until only a single cluster remains (Maimon and Rokach, 2005).

222 We here used Wards criterion (Ward, 1963) for agglomerative clustering, which was implemented using

223 the linkage function of Matlab's Statistics and Machine Learning Toolbox (The MathWorks Inc., Natick,

224 MA, USA). To compare our clustering result of the observed sequences with those obtained using other

225 count-based methods, we also quantified the inter-sample difference with BPLN and Bray-Curtis methods.

226 BPLN was applied according to the methods described in the original paper by Rempala et al. (2011).

227     To evaluate the goodness of fit of the clustering results, Rempala and colleagues Rempala et al.

228 (2011) calculated the cophenetic correlation coefficient (CCC), which quantifies the distortion due to the

229 transformation from the distance matrix to the cophenetic matrix, from which the dendrogram was derived.

230 However, the CCC does not always accurately reflect the goodness of fit of the results. Indeed, Wards

231 method tends to produce lower CCC values than other methods such as average and centering methods

232 even though it has been previously reported as the best agglomerative method (Hands and Everitt, 1987;

233 Saracli et al., 2013). Therefore, instead of the CCC, we verified the fit of the model based on the statistical

234 significance of the distance between the nodes of the dendrogram, because the significance of the estimated

235 value of JSD is unclear. Specifically, we used bootstrap methods to evaluate significance, resampled data

236 points from the naive PDF according to the number of observed read counts, and then re-estimated the PDF

237 from the resampled data points. We then calculated the JSDs between the naive and re-estimated PDFs.

238 We repeated this process 100 times to obtain a histogram of the calculated JSDs. The 2.5th and 97.5th

239 percentiles of the histogram of the JSDs between the naive and each re-estimated PDF represent both ends

240 of the 95% confidence interval, where values outside of the interval indicate a significance level of over 5%.

241     To identify the sequences with the greatest contributions to the inter-sample distances, we selected square

242 bins for the top 1% of the local JSDs. We next defined the sequences in these bins as those contributing to

243 the observed pairwise sample difference. Furthermore, to investigate the characteristics of the contributing

244 sequences, we calculated the relative frequencies of the amino acid residues in all of the contributing

245   sequences. The graphics of the relative frequencies were obtained using WebLog 3 software (Crooks et al.,

246   2004).

247   All analyses were performed using custom-made codes written in Python, Matlab, and R.

## 3   RESULTS

### 3.1   Evaluation of sequence dissimilarity for pairwise sequence alignment

249   Using the pairwise sequence alignment and Eq. 1, which excludes the influence of the sequence lengths

250   from the alignment results, we calculated the dissimilarity matrix of all pairwise sequences in the dataset of

251   Rempala et al. (2011). The upper panels in Figs. 1(A) and 2(B) show the dissimilarity matrices obtained

252   with the 10 different substitution matrices, five of PAM and five of BLOSUM. As shown in these panels,

253   the components of the dissimilarity matrices are clearly separated into two distinct clusters, which are

254   considered to reflect the $\alpha$-chain rearrangements between J$\alpha$2.6 and J$\alpha$2 under the usage of the other

255   fixed VJ genes. Moreover, the low-numbered PAMs and high-numbered BLOSUMs showed more gradual

256   differences among the matrix elements than the others. This tendency was even more evident when viewing

257   their embedded spaces for the separation of clusters. The lower panels in Fig. 1(A) and (B) show the t-SNE

258   projection maps of the corresponding dissimilarity matrices in the upper panels. In this case, low-numbered

259   PAMs and high-numbered BLOSUMs tended to have merged clusters. This may be attributed to the

260   specific characteristics of these two substitution matrices, which have a higher variability in the scores for

261   replacements between a pair of amino acids. Based on these results, we used the BLOSUM62 dissimilarity

262   matrix for subsequent analyses for two main reasons. First, both the too high-numbered PAMs and too

263   low-numbered BLOSUMs seemed to lose the intra-cluster structures by trying to compress the clusters into

264   regions that were too small, while both the too low-numbered PAMs and too high-numbered BLOSUMs

265 diminished any inter-cluster differences, resulting in indistinguishable clusters. Second, BLOSUM62 has

266 been the most widely used matrix in analyses of TCRs and antigen peptides to date (Oyarzún et al., 2013;

267 Schwaiger et al., 2014; Hoffmann et al., 2015; Aas-Hanssen et al., 2015).

268 ## 3.2 Dimensionality reduction of the dissimilarity matrix

269 To evaluate the applicability of dimensionality reduction methods, we reduced the dimensionality of the

270 dissimilarity matrix into a two-dimensional space using four different dimensionality reduction methods

271 (t-SNE, MDS, ISOMAP, and SE). In Fig. 2, each point in each panel corresponds to a unique sequence

272 of TCRs, and the spatial distances between pairs of points reflect the dissimilarity of the sequences

273 corresponding to the points. Panels (i) and (ii) of Fig. 2(A–D) show the projection results of the unique

274 sequences obtained from all samples, and the subset of points (sequences) in (i) that appeared in the sample

275 denoted in the inset letters of the panel, respectively. The sample differences could be clearly reflected

276 according to the scattering pattern of the points. Moreover, the points derived from the t-SNE and MDS

277 methods spread more widely over the two-dimensional space than the others, whereas the points were more

278 locally consolidated with the ISOMAP method, and especially with SE. This result suggests that t-SNE and

279 MDS may be more appropriate than other reduction methods for larger datasets, because highly dense

280 regions can cause difficulty in comparing the probabilistic distributions between samples. Furthermore,

281 the two clear clusters in the dissimilarity matrix (the upper panel of BLOSUM62 in Fig. 1(B)) were well

282 reflected in the two clusters for the MDS and ISOMAP methods (Fig. 2(B,C)), but were not represented

283 clearly in the clusters of t-SNE. This result suggests that t-SNE emphasizes slight differences within clusters

284 rather than large differences between the clusters of the dissimilarity matrix. Since it is unclear whether this

285 visualization property of t-SNE works efficiently for comparisons between samples, we quantified and

286 compared the distributions of data points at the next step, and examined the method that would be most

287 appropriate for this purpose.

## 3.3   Hierarchical clustering of the pairwise-sample-difference matrix

289   We applied the KDE algorithms to the spatial distributions of data points to estimate their probability

290 density functions (color gradient in Fig. 2), with which JSD is calculated to quantify pairwise-sample

291 differences. The matrices of the pairwise-sample differences are shown in Fig. 3(A), and the resulting

292 dendrograms in Fig. 3(B) indicate the hierarchical clustering results with the agglomerative method. The

293 clustering results can be categorized into two groups: ISOMAP and the others (MDS, t-SNE, and SE). The

294 dendrograms of t-SNE, MDS, and SE showed good correspondence with our intuitive notions about the

295 hierarchical structure of the experimental conditions, which were ranked in order of donor sites, types of T

296 cells, and genetic background with clear biological significance (Rempala et al., 2011). By contrast, the

297 dendrogram of ISOMAP showed a mismatch in the hierarchical order between the T-cell types and donor

298 sites of Ep mice.

299   To verify the results of hierarchical clustering obtained by our method, they were compared with those

300 obtained with previous observation count-based methods, the BPLN and Bray-Curtis method. As shown

301 in Fig. 4, the sample differences and dendrograms estimated from the BPLN and Bray-Curtis methods

302 were very similar to those obtained using our approach with MDS, t-SNE, and SE. Importantly, these

303 similar results were obtained with different data modalities: sequence similarity and observation counts.

304 Therefore, this consistency suggests that there is common information between sequence similarity and

305 observation counts with respect to quantifying the differences among samples. We should note that these two

306 modalities can be combined simply by assigning the number of observed sequence counts as a weighting

307 factor for each data point (i.e., a unique sequence) in the embedded space. Indeed, the counts-weighted

308  PDFs using KDE (Fig. S1) showed no obvious change in the hierarchical clustering structure of the

309  pairwise-sample differences. Taking these results together, the MDS or t-SNE appears to be the better choice

310  as a dimension-reduction method for evaluation of differences in TCR repertoires among samples, given

311  that these methods show wide spatial distributions of the data points, and also show the most consistent

312  dendrogram structures with those of previous count-based methods.

### 3.4  Significance test for inter-sample differences

314    To verify the statistical significance of the calculated JSDs between all sample pairs, we calculated the

315  JSDs between the naive and the re-estimated PDFs using a non-parametric bootstrap algorithm. Figure 5

316  shows the histogram of the JSD values between the naive PDF of Fig. 2 (C, ii) and the re-estimated PDFs.

317  In the figure, arrows indicate the naive JSD values between the sample designated on the top of the panel

318  and the other samples. If the values indicated by the arrows are bigger than the light red region in the panel,

319  the pairwise naive PDFs deriving the naive JSD are significantly different each other. Except for EpTN2, the

320  arrows that indicate the JSD values between the pairs in proximity to the terminal nodes of the dendrogram

321  in Fig. 3 (A) were in inside of the light red regions, which means that the JSD values were not significantly

322  bigger than the JSD values of the histogram. This result suggests that the PDFs of these pairs are so similar

323  that they cannot be statistically distinguished from each other. By contrast, the arrows that indicate the naive

324  JSD values between pairs in the upper parts of the clusters, above the terminal nodes, were in outside of the

325  red light regions, which means that the JSD values were significantly different from each other. This result

326  indicates that the types of T cells and the genetic background can be discriminated with sufficient statistical

327  significance.

### 3.5 Spatial distribution of local JSD values

328

329 The main advantage of our method compared to count-based methods is the ability to identify the major

330 sequences contributing to inter-sample differences. To identify the sequences with the greatest contributions

331 to large local JSD values, we plotted the spatial distribution of the local JSDs between the WtTN1 and

332 EpTN1 sequences. As shown in Fig. 6, two regions were identified that were associated with top 1%

333 significance values. Table 1 lists the identified sequences in these regions with larger local JSDs than the

334 others. In regions 1 and 2, there were no sequences for WtTN1, whereas EpTN1 had multiple sequences

335 in these regions, suggesting that these apparent Ep-specific sequences may contribute to the observed

336 abnormality in the antigen presentation of Ep mice.

337 This type of sequence identification can provide further knowledge about the characteristics of sequence

338 alignments. Figure 7 shows the occupation probability (relative frequency) of amino acids at each position

339 of the sequences obtained from all of the sequences contributing to pairwise differences, and a consensus

340 sequence was determined from amino acid positions 6th to 11th. We note that these contributing sequences

341 and their characteristics cannot be easily identified simply by examining the overlapping sequences in two

342 samples, because there was almost no overlap between EpTN1 and WtTN1 sequences (0.352%, 1/284), and

343 because these account for only 6.34% (18/284) of the total unique sequences in the two samples.

## 4 DISCUSSION

344 We quantified the difference in TCR repertoires among different samples based on amino acid sequence

345 dissimilarity. Through a quantitative comparison of the sequence distributions in the dimension-reduction

346 spaces of the dissimilarity matrix, we estimated the inter-sample hierarchical structure, which was almost

347 identical to that estimated with previous count-based methods that did not incorporate detailed sequence

348 information. Furthermore, we identified the sequences that most strongly contribute to the pairwise sample

349 difference using the local JSD distribution.

350     Despite the fact that our method relies on sequence similarity and previous methods are based on

351 observation frequency, completely different features of the TCR repertoire, almost the same sample-

352 clustering structure was obtained. This suggests that there is a relationship between the observation counts

353 and sequence similarities, which was further confirmed by the lack of change in the structure of the

354 hierarchical clustering result when estimating the PDFs by the KDE algorithm with the weights based on

355 the number of observed sequence counts. Further studies to understand this relation in greater depth would

356 allow to cross-check the results of sample classification by investigating the consistency of two methods.

357 Moreover, clearly distinguishing the overlapping and non-overlapping information between counts and

358 sequences may allow for more detailed classification of samples.

359     Although our method and the counts-based methods provide similar classification result, there are two

360 unique merits of our method. The first is the robustness against errors derived from polymerase chain

361 reaction (PCR) amplification bias attributed to the variability in reproducibility for individual sequences

362 (Greiff et al., 2015). Previous studies have shown that the PCR efficiency is affected by sequence profiles

363 such as the length and GC content (Aird et al., 2011; Kivioja et al., 2011). Indeed, high-throughput

364 sequencing with DNA barcoding has confirmed that the PCR amplification efficiency of TCR sequences is

365 highly variable due to the differences in profiles of individual cDNA molecules (Carlson et al., 2013; Shugay

366 et al., 2014; Best et al., 2015). This fact suggests that the PCR process for TCR sequence amplification

367 induces errors in the numbers of observed sequence counts, which may eventually lead to errors in the

368 results of counts-based methods such as PA models. Alternatively, our method does not depend on the

369 sequence counts, allowing for reliability against errors due to PCR bias. The second key merit of our

370 method is the ability to identify the sequences with the greatest contributions to pairwise sample differences.

371  This sequence identification allows for targeted analyses along with the results of other studies such as the

372  simulation modeling for determining the crystal structures of the TCRs encoded by these sequences (Klausen

373  et al., 2015) or establishing alignments between CDR3 sequences and microbial genomes (Aas-Hanssen

374  et al., 2015). Such a closed-loop experimental design may help to achieve a breakthrough in the development

375  of vaccines or immunotherapies (Hou et al., 2016). Thus, the present results and advantages demonstrate the

376  potential applicability of adopting a sequence-based method in repertoire analysis, which can compensate

377  for the drawbacks in conventional count-based methods.

378  Nevertheless, there are several issues and problems that should be mentioned that are worthy of further

379  investigation in developing and improvement of sequence-based approaches for comparison of TCR

380  repertoires.

381  One issue concerns the treatment of gap penalties. When we evaluated the differences of the score matrices

382  shown in Fig. 1, we fixed the gap opening and extension penalties to 10 and 1, respectively. Although the

383  effects of the penalties have not been adequately investigated in previous studies (Wrabl and Grishin, 2004),

384  the gap opening penalty was found to affect estimations of the hierarchical data structure (data not shown).

385  The CDR3 region of TCR is a much shorter sequence than peptide sequences, and also shows frequent

386  deletions and insertions from somatic recombination events. Considering these characteristics, further

387  investigations about the effects of gap penalties are needed.

388  Another aspect worthy of further consideration is the empirical estimation of the cost functions used

389  in the dimensionality reduction methods and in the combination of projection methods and comparison

390  of the embedded results. As demonstrated in Fig. 2, the scattering patterns of the sequence data in the

391  low-dimensional space depend on the cost functions of the method adopted. Using MDS and ISOMAP,

392  we obtained two clear clusters reflecting two regions in the dissimilarity matrix. This is because the cost

393  function of MDS preferentially preserves the distances between the distant points rather than those between

394  nearby points (Van Der Maaten et al., 2009). By contrast, t-SNE emphasizes the local structures of nearby

395  points over global points by using the Student t-distribution as the kernel of the embedded space. These

396  cost functions were empirically determined for visualization purposes in the original papers, without

397  consideration of the subsequent quantitative inter-sample comparison of the embedded results. Although

398  our results suggest that the empirical combination of dimensionality reduction methods and comparison of

399  the embedded results by JSDs may work well, both the projection method and comparison method in the

400  embedded space should be consistently designed so as to best reflect the inter-sample difference in the

401  original sequence space. This method might be developed by choosing an information-theoretic measure

402  for the cost function of projection that can preserve the relevant information of repertoires in the original

403  sequence space. Because the underlying high-dimensional structures of the repertoire are difficult to capture

404  intuitively, methods based on firm theoretical rationality and biological significance are indispensable for

405  further exploitation of repertoire information.


## DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

406  The authors declare that the research was conducted in the absence of any commercial or financial

407  relationships that could be construed as a potential conflict of interest.


## AUTHOR CONTRIBUTIONS

408  R.Y., Y.K., and T.J.K: study conception and design; R.Y. and Y.K.: performed the research; R.Y. and Y.K:

409  data analysis; and R.Y. and T.J.K.: wrote the paper.

## FUNDING

## ACKNOWLEDGMENTS

## REFERENCES

417    Aas-Hanssen, K., Thompson, K. M., Bogen, B., and Munthe, L. A. (2015). Systemic Lupus Erythematosus:

418    Molecular Mimicry between Anti-dsDNA CDR3 Idiotype, Microbial and Self PeptidesAs Antigens for

419    Th Cells. *Frontiers in Immunology* 6, 382. doi:10.3389/fimmu.2015.00382

420    Aird, D., Ross, M. G., Chen, W.-S., Danielsson, M., Fennell, T., Russ, C., et al. (2011). Analyzing

421    and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome biology* 12, R18.

422    doi:10.1186/gb-2011-12-2-r18

423    Arlot, S. and Celisse, A. (2010). A survey of cross-validation procedures for model selection. *Statistics*

424    *Surveys* 4, 40–79. doi:10.1214/09-SS054

425    Belkin, M. and Niyogi, P. (2001). Laplacian eigenmaps and spectral techniques for embedding and

426    clustering. In *Proceedings of the 14th International Conference on Neural Information Processing*

427    *Systems: Natural and Synthetic* (Cambridge, MA, USA: MIT Press), NIPS'01, 585–591

428   Belkin, M. and Niyogi, P. (2003).   Laplacian Eigenmaps for Dimensionality Reduction and Data

429       Representation. *Neural Computation* 15, 1373–1396. doi:10.1162/089976603321780317

430   Best, K., Oakes, T., Heather, J. M., Shawe-Taylor, J., and Chain, B. (2015). Computational analysis of

431       stochastic heterogeneity in PCR amplification efficiency revealed by single molecule barcoding. *Scientific*

432       *Reports* 5, 14629. doi:10.1038/srep14629

433   Bishop, C. (2007). *Pattern Recognition and Machine Learning (Information Science and Statistics), 1st edn.*

434       *2006. corr. 2nd printing edn* (Springer, New York)

435   Borg, I. and Groenen, P. (2005). *Modern Multidimensional Scaling - Theory and Applications* (Springer

436       Science & Business Media). doi:10.1007/0-387-28981-X

437   Bray, J. R. and Curtis, J. T. (1957). An ordination of the upland forest communities of southern wisconsin.

438       *Ecological monographs* 27, 325–349

439   Carlson, C. S., Emerson, R. O., Sherwood, A. M., Desmarais, C., Chung, M.-w., Parsons, J. M., et al.

440       (2013). Using synthetic templates to design an unbiased multiplex PCR assay. *Nature communications* 4,

441       2680. doi:10.1038/ncomms3680

442   Christopher, A., Andrew, M., and Stefan, S. (1997). Locally weighted learning. *Artif Intell Rev* 11, 11–73

443   Crooks, G. E., Hon, G., Chandonia, J.-m., and Brenner, S. E. (2004).  WebLogo : A Sequence Logo

444       Generator. *Genome research* 14, 1188–1190. doi:10.1101/gr.849004.1

445   Daily, J. (2016). Parasail: SIMD C library for global, semi-global, and local pairwise sequence alignments.

446       *BMC Bioinformatics* 17, 81. doi:10.1186/s12859-016-0930-z

447   Davis, M. M. and Bjorkman, P. J. (1988). T-cell antigen receptor genes and T-cell recognition. *Nature* 334,

448       395–402. doi:10.1038/334395a0

449   Duez, M., Giraud, M., Herbert, R., Rocher, T., Salson, M., and Thonier, F. (2016). Vidjil: A web platform

450       for analysis of high-Throughput repertoire sequencing. *PLoS ONE* 11, 1–12. doi:10.1371/journal.pone.

451    0166126

452    Farrar, M. (2007). Striped Smith-Waterman speeds database searches six times over other SIMD

453        implementations. *Bioinformatics* 23, 156–161. doi:10.1093/bioinformatics/btl582

454    Floyd, R. W. (1962). Algorithm 97: Shortest path. *Commun. ACM* 5, 345–. doi:10.1145/367766.368168

455    Gray, A. and Moore, A. (2001). N-body problems in statistical learning. *In Advances in Neural Information

456        Processing Systems* 4, pages 521–527

457    Greene, J., Birtwistle, M. R., Ignatowicz, L., and Rempala, G. A. (2013). Bayesian multivariate Poisson

458        abundance models for T-cell receptor data. *Journal of Theoretical Biology* 326, 1–10. doi:10.1016/j.jtbi.

459        2013.02.009

460    Greiff, V., Miho, E., Menzel, U., and Reddy, S. T. (2015). Bioinformatic and Statistical Analysis of

461        Adaptive Immune Repertoires. *Trends in Immunology* 36, 738–749. doi:10.1016/j.it.2015.09.006

462    Guindani, M., Sepúlveda, N., Paulino, C. D., and Mueller, P. (2014). A Bayesian semi-parametric approach

463        for the Differential Analysis of Sequence Counts Data. *Journal of the Royal Statistical Society. Series C,

464        Applied Statistics* 63, 385–404

465    Heidenreich, N. B., Schindler, A., and Sperlich, S. (2013). Bandwidth selection for kernel density

466        estimation: A review of fully automatic selectors. *AStA Advances in Statistical Analysis* 97, 403–433.

467        doi:10.1007/s10182-013-0216-y

468    Henikoff, S. and Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. *Proceedings

469        of the National Academy of Sciences of the United States of America* 89, 10915–10919. doi:10.1073/pnas.

470        89.22.10915

471    Hoffmann, T., Krackhardt, A. M., and Antes, I. (2015). Quantitative Analysis of the Association Angle

472        between T-cell Receptor V$\alpha$/V$\beta$ Domains Reveals Important Features for Epitope Recognition. *PLoS

473        Computational Biology* 11, 1–28. doi:10.1371/journal.pcbi.1004244

474 Hou, D., Chen, C., Seely, E. J., Chen, S., and Song, Y. (2016). High-throughput sequencing-based immune

475     repertoire study during infectious disease. *Frontiers in Immunology* 7, 1–11. doi:10.3389/fimmu.2016.

476     00336

477 Hughes, A., Ruan, Y., Ekanayake, S., Bae, S.-H., Dong, Q., Rho, M., et al. (2012). Interpolative

478     multidimensional scaling techniques for the identification of clusters in very large sequence sets. *BMC*

479     *Bioinformatics* 13 Suppl 2, S9. doi:10.1186/1471-2105-13-S2-S9

480 Ito, K., Igarashi, M., Miyazaki, Y., Murakami, T., Iida, S., Kida, H., et al. (2011). Gnarled-trunk evolutionary

481     model of influenza a virus hemagglutinin. *PLoS ONE* 6. doi:10.1371/journal.pone.0025953

482 Jones, M. C., Marron, J. S., and Sheather, S. J. (1996). A brief survey of bandwidth selection for density

483     estimation. *Journal of the American Statistical Association* 91, 401–407

484 Kivioja, T., Vähärautio, A., Karlsson, K., Bonke, M., Enge, M., Linnarsson, S., et al. (2011). Counting

485     absolute numbers of molecules using unique molecular identifiers. *Nature Methods* 9, 72–74. doi:10.

486     1038/nmeth.1778

487 Klausen, M. S., Anderson, M. V., Jespersen, M. C., Nielsen, M., and Marcatili, P. (2015). LYRA, a

488     webserver for lymphocyte receptor structural modeling. *Nucleic Acids Research* 43, W349–W355.

489     doi:10.1093/nar/gkv535

490 Laydon, D. J., Bangham, C. R. M., and Asquith, B. (2015). Estimating T-cell repertoire diversity: limitations

491     of classical estimators and a new approach. *Phil Trans R Soc B* 370, 20140291–. doi:10.1098/rstb.2014.

492     0291

493 Maimon, O. and Rokach, L. (2005). *Data Mining and Knowledge Discovery Handbook* (Secaucus, NJ,

494     USA: Springer-Verlag New York, Inc.)

495 Oyarzún, P., Ellis, J. J., Bodén, M., and Kobe, B. (2013). PREDIVAC: CD4+ T-cell epitope prediction

496     for vaccine design that covers 95% of HLA class II DR protein diversity. *BMC bioinformatics* 14, 52.

497   doi:10.1186/1471-2105-14-52

498   Pacholczyk, R., Ignatowicz, H., Kraj, P., and Ignatowicz, L. (2006). Origin and T Cell Receptor Diversity

499   of Foxp3+CD4+CD25+ T Cells. *Immunity* 25, 249–259. doi:10.1016/j.immuni.2006.05.016

500   Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2012). Scikit-learn:

501   Machine Learning in Python. *Journal of Machine Learning Research* 12, 2825–2830. doi:10.1007/

502   s13398-014-0173-7.2

503   Rempala, G. A., Seweryn, M., and Ignatowicz, L. (2011). Model for comparative analysis of antigen

504   receptor repertoires. *Journal of Theoretical Biology* 269, 1–15. doi:10.1016/j.jtbi.2010.10.001

505   Robinson, M. D. and Smyth, G. K. (2007). Gene expression Moderated statistical tests for assessing

506   differences in tag abundance. *Bioinformatics* 23, 2881–2887. doi:10.1093/bioinformatics/btm453

507   Schwaiger, J., Aberle, J. H., Stiasny, K., Knapp, B., Schreiner, W., Fae, I., et al. (2014). Specificities of

508   human CD4+ T cell responses to an inactivated flavivirus vaccine and infection: correlation with structure

509   and epitope prediction. *Journal of virology* 88, 7828–42. doi:10.1128/JVI.00196-14

510   Sepúlveda, N., Daniel, C., and Carneiro, J. (2010). Estimation of T-cell repertoire diversity and clonal

511   size distribution by Poisson abundance models. *Journal of Immunological Methods* 353, 124–137.

512   doi:10.1016/j.jim.2009.11.009

513   Shugay, M., Britanova, O. V., Merzlyak, E. M., Turchaninova, M. a., Mamedov, I. Z., Tuganbaev,

514   T. R., et al. (2014). Towards error-free profiling of immune repertoires. *Nature methods* 11, 653–5.

515   doi:10.1038/nmeth.2960

516   Silverman, J. D., Washburne, A., Mukherjee, S., and David, L. A. (2016). A phylogenetic transform

517   enhances analysis of compositional microbiota data. *bioRxiv* , 072413doi:10.1101/072413

518   Smith, T. F. and Waterman, M. S. (1981). Identification of common molecular subsequences. *Journal of*

519   *molecular biology* 147, 195–197

520 States, D. J., Gish, W., and Altschul, S. F. (1991). Improved sensitivity of nucleic acid database searches

521     using application-specific scoring matrices. *Methods* 3, 66–70

522 Tang, Z. Z., Chen, G., and Alekseyenko, A. V. (2016). PERMANOVA-S: Association test for microbial

523     community composition that accommodates confounders and multiple distances. *Bioinformatics* 32,

524     2618–2625. doi:10.1093/bioinformatics/btw311

525 Tenenbaum, J. B., de Silva, V., and Langford, J. C. (2000). A global geometric framework for nonlinear

526     dimensionality reduction. *Science (New York, N.Y.)* 290, 2319–23. doi:10.1126/science.290.5500.2319

527 Van Der Maaten, L., Postma, E., and Van den Herik, J. (2009). Dimensionality reduction: a comparative

528     review. *Journal of Machine Learning Research* 10, 66–71

529 Van Der Maaten, L. J. P. and Hinton, G. E. (2008). Visualizing high-dimensional data using t-sne. *Journal*

530     *of Machine Learning Research* 9, 2579–2605. doi:10.1007/s10479-011-0841-3

531 Ward, J. H. (1963). Hierarchical grouping to optimize an objective function. *Journal of the American*

532     *Statistical Association* 58, 236–244. doi:10.1080/01621459.1963.10500845

533 Wrabl, J. O. and Grishin, N. V. (2004). Gaps in Structurally Similar Proteins : Towards Improvement of

534     Multiple Sequence Alignment. *Proteins: Structure, Function, and Bioinformatics* 54, 71–87

## FIGURES AND TABLES

**Figure 1.** Dissimilarity matrices and their embedded distributions with 10 different score matrices; (A) PAM, (B) BLOSUM.
The upper and lower panels show the dissimilarity matrices and projection maps in two-dimensional space, respectively. All of the rows and columns in each dissimilarity matrix were sorted according to the sum of their elements. The colors of points in the lower panels of (A) and (B) correspond to the clusters in PAM250 and BLOSUM45 that were discriminated by k-means algorithms (k = 7).

**Figure 2.** Dimensional reduction with four different dimensionality reduction methods: (A) t-SNE, (B) MDS, (C) ISOMAP, and (D) SE.

Panel (i) includes the points of the total unique sequences observed in all samples. Panel (ii) includes only the portions of sequences that were observed in each sample. "Ep" and "Wt" denote two different genetic backgrounds of mice. "TN" and "TR" denote naive and regulatory T-cells. "1" and "2" denote the thymus and peripheral lymph nodes, respectively. As an instance, EpTN1 denotes the naive T-cells that were collected from the thymus in the "Ep" mice.

(A) Matrices of pairwise-sample differences



(B) dendrogram



**Figure 3.** JSD matrices and their clustering results with four different methods: (i) tSNE, (ii) MSD, (iii) ISOMAP, and (iv) SE. (A) Matrices of pairwise-sample differences and (B) the dendrogram constructed from the matrices.
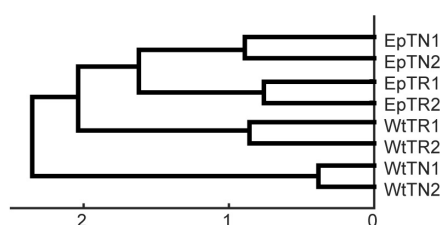
(A) Matrices of pairwise-sample differences

(i) BPLN

(ii) Bray - Curtis

(B) dendrogram
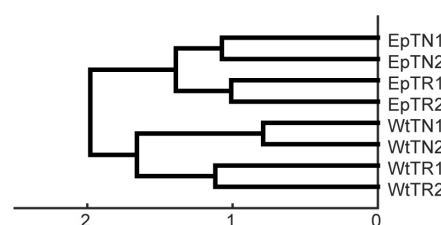
(i) BPLN

(ii) Bray - Curtis



**Figure 4.** Sample-distance matrices constructed with two methods: (i) BPLN, (ii) Bray-Curtis. (A) Matrices of pairwise-sample differences and (B) the dendrogram constructed from the matrices.

**Table 1.** Sequences contributing to the JSD between EpTN1 and WtTN1.

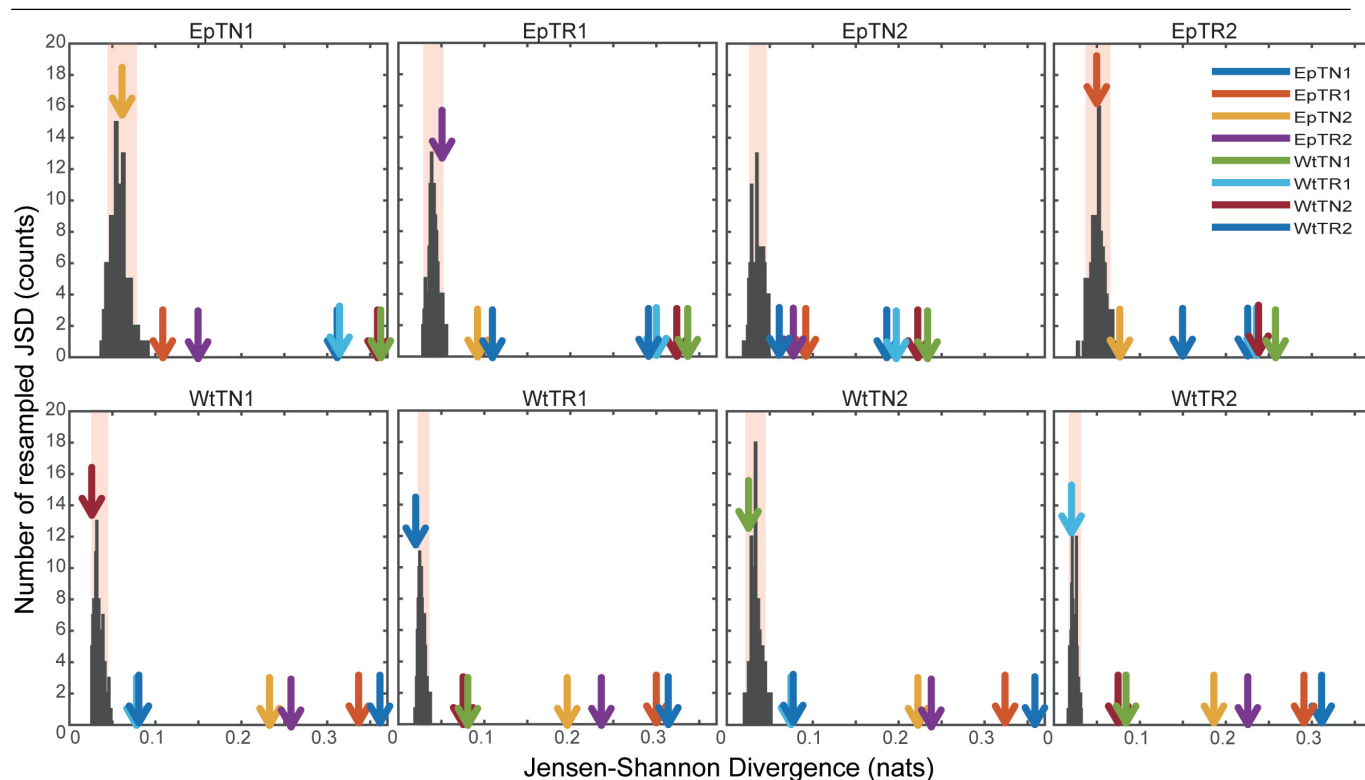| Regions | EpTN1 | WtTN1 |
|---------|-------|-------|
| 1 | CAASAYQLIWG<br>CAASCYQLIWG<br>CAASRYQLIWG<br>CAASTYQLIWG | |
| 2 | CAAGNYQLIWG<br>CAAHNYQLIWG<br>CAANNYQLIWG<br>CAARNYQLIWG<br>CAASNYQLIWG<br>CAATNYQLIWG<br>CADLNYQLIWG<br>CADSNYQLIWG<br>CAGSNYQLIWG<br>CASHNYQLIWG<br>CASSNYQLIWG<br>CATSNYQLIWG<br>CAVSNYQLIWG<br>CVGSNYQLIWG | |

**Figure 5.** Significance tests of JSD values using bootstraps.
Each colored arrow indicates the naive JSD values between the resampled sample and another. The light red region indicates the 95% confidence interval.
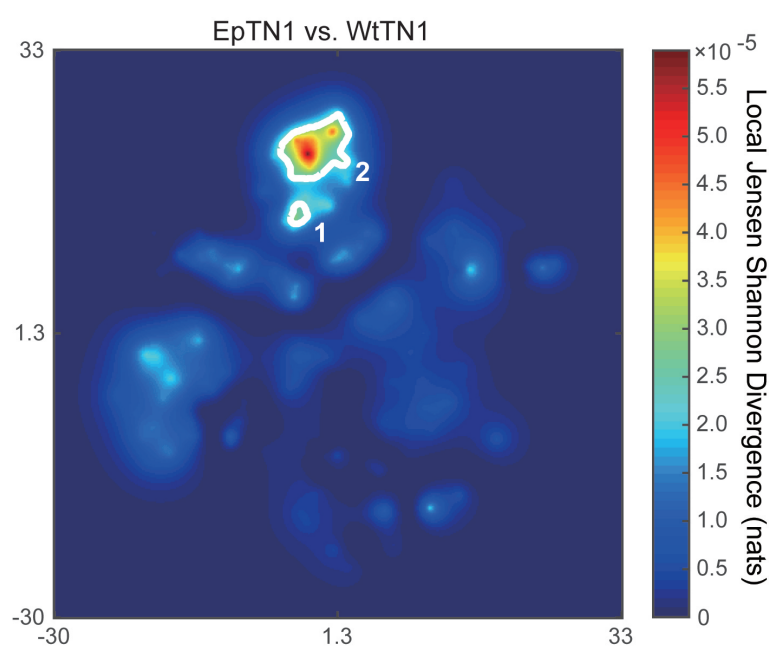


**Figure 6.** Spatial distribution of local JSD values between EpTN1 and WtTN1.
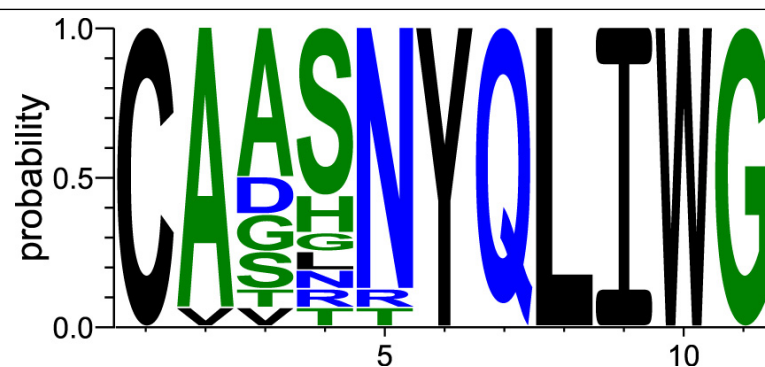The white line shows the contours of the regions with significant local JSDs.

**Figure 7.** Relative frequencies of observed amino acids at each position in the contributing sequences.
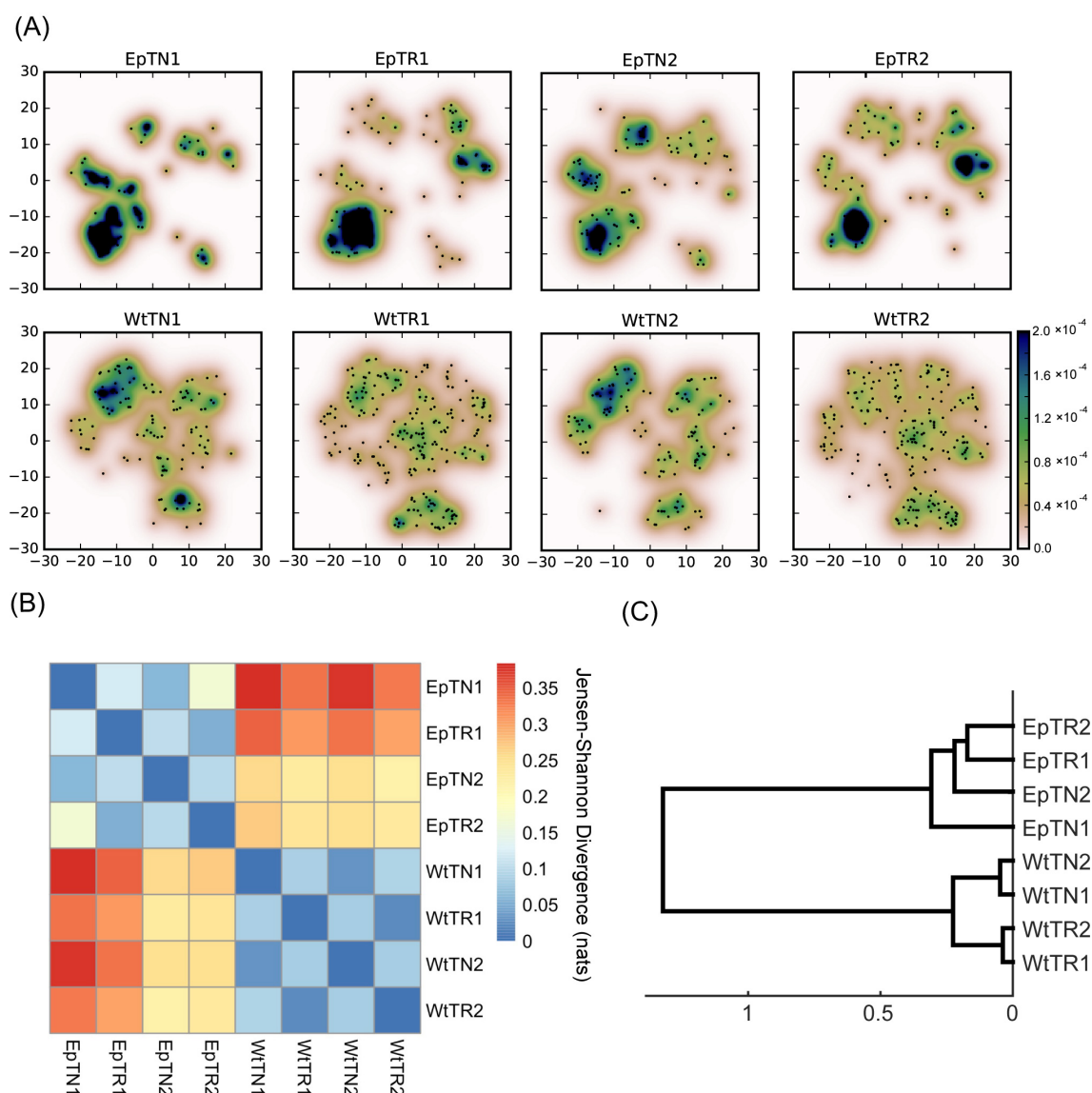
## SUPPLEMENTARY INFORMATION



**Figure S1.** Count-weighted PDF estimation with tSNE-embeding data.
Conventions comply with Fig.2 and 3. (A) PDFs weighted by the number of sequence counts. (B) Sample distance matrix estimated with the weighted PDFs. (C) the dendrogram constructed from the matrix in (B).

**Table S1.** Detail parameters of manifold learning methods.

Except for ISOMAP, we used the functions in the class of sklearn.manifold in the scikit-learn toolbox. The parameters of each function are described in the Table S1. For ISOMAP, first, we calculated the geodesic distances by using k-nearest neighbor algorithm and Floyd-Warshall method. Then, we applied the geodesic distances to MDS with the following parameters.

| tSNE | parameters | value |
|---|---|---|
| | n_components | 2 |
| | random_state | 0 |
| | metric | precomputed' |
| | the other parameters | default |
| MDS | parameters | value |
| | n_components | 2 |
| | maximum iteration | 400 |
| | relative tolerance w.r.t stress to declare converge | 1.00E−05 |
| | dissilarity | precomputed' |
| | the other parameters | default |
| SE | parameters | value |
| | n_components | 2 |
| | n_neighbors | 50 |
| | affinity | precomputed' |
| | the other parameters | default |

| ISOMAP | parameters | value |
|---|---|---|
| | n_components | 2 |
| | maximum iteration | 1000 |
| | relative tolerance w.r.t stress to declare converge | 1.00E−05 |
| | dissilarity | precomputed' |
| | n_neighbors | 50 |
| | the other parameters | default |