

1 A rationally designed mimotope library for profiling of the human IgM repertoire

2

3

4 Anastas Pashov^{1,*}, Velizar Shivarov², Maya Hadzhieva¹, Victor Kostov^{1,3}, Dilyan Ferdinandov³, Karen-
5 Marie Heinz⁵, Shina Pashova^{1,4}, Milena Todorova¹, Tchavdar Vassilev¹, Thomas Kieber-Emmons⁶,
6 Leonardo A. Meza-Zepeda⁵, Eivind Hovig⁵

7

8 ¹Department of Immunology, Stephan Angeloff Institute of Microbiology, BAS, Sofia, Bulgaria

9 ²Laboratory of Clinical Immunology and Department of Clinical Hematology, Sofamed University
10 Hospital, Sofia, Bulgaria; ORCID: 0000-0001-5362-7999

11 ³Neurosurgery Clinic, St. Ivan Rilsky Hospital, MU, Sofia, Bulgaria

12 ⁴Department of Molecular Immunology, Institute of Biology and Immunology of Reproduction, BAS,
13 Sofia, Bulgaria

14 ⁵Department of Tumor Immunology, Oslo University Hospital, Oslo, Norway

15 ⁶Winthrop P. Rockefeller Cancer Research Center, UAMS, Little Rock, AR, USA

16

17

18 *Corresponding author.

19 **Mailing Address:**

20 Institute of Microbiology, BAS,

21 Acad. G Bonchev St, block 26

22 Sofia 1113, Bulgaria

23 **E-mail:** a_pashov@microbio.bas.bg (AP)

24 **Phone:** +359 897 944628

25

26

27

28

29

1 **Abstract**

2 Specific antibody reactivities are routinely used as biomarkers but the use of antibody repertoire profiles
3 is still awaiting recognition. Here we suggest to expedite the adoption of this class of system level
4 biomarkers by rationally designing a peptide array as an efficient probe for an appropriately chosen
5 repertoire compartment. Most IgM antibodies are characterized by few somatic mutations, polyspecificity
6 and physiological autoreactivity with housekeeping function. Previously, probing this repertoire with a set
7 of immunodominant self-proteins provided only coarse information on repertoire profiles. In contrast,
8 here we describe the rational selection of a peptide mimotope set, appropriately sized as a potential
9 diagnostic, that also represents optimally the diversity of the human public IgM reactivities. a 7-mer
10 random peptide phage display library was panned on pooled human IgM. Next generation sequencing of
11 the selected phage yielded a non-exhaustive set of 224087 mimotopes which clustered in 790 sequence
12 clusters. A set of 594 mimotopes, representative of the most significant clusters, was used to demonstrate
13 that this approach samples symmetrically the space of IgM reactivities. When probed with diverse
14 patients' sera in an oriented peptide array, this set produced a higher and more dynamic signal as
15 compared to 1) random peptides, 2) random peptides purged of mimotope-like sequences and 3)
16 mimotopes from a small subset of clusters. In this respect, the representative library is an optimized probe
17 of the human IgM diversity. Proof of principle predictors for randomly selected diagnoses based on the
18 optimized library demonstrated that it contains more than 10^{70} different profiles with the capacity to
19 correlate with diverse pathologies producing well separable data. Thus, an optimized small library of IgM
20 mimotopes is found to address very efficiently the dynamic diversity of the human IgM repertoire,
21 providing informationally dense, structurally interpretable, IgM reactivity profiles.

22

23

1 **Introduction**

2 The repertoire of human IgM contains a considerable proportion of moderately autoreactive antibodies
3 characterized by low intrinsic affinity/ low specificity, functioning as a first line of defense [1], as
4 scavengers of senescent cells and debris [2-6], and even in tumor surveillance [7]. It is becoming
5 increasingly clear that the human antibody repertoire has an organization similar to that of its murine
6 counterpart [8-12]. About one fourth of the murine splenic B lymphocytes that respond to
7 lipopolysaccharide have B cell receptors which are moderately autoreactive. Affected very little by
8 somatic mutations and follicular evolution, the physiological self-reactivities largely overlap with
9 germline-encoded polyspecific antibodies [13-15]. Eighty percent of murine serum IgM falls in this
10 category and is referred to as natural antibodies (nAbs) [16, 17]. Apart from the polyspecific splenic B
11 cells, the source of nAbs in mice seems to be mostly a population of B1-related IgM⁺ plasma cells residing
12 in a unique IL-5 dependent bone marrow niche [18]. By interacting with structures of self and carrying
13 housekeeping tasks, this part of the antibody repertoire reacts swiftly to and reflects changes in the
14 internal environment. Therefore, it can be a better source of biomarkers if probed as a natural biosensor
15 than other compartments of the antibody repertoire like IgG and IgE.

16 Our working hypothesis is that an essential part of the human IgM repertoire involved in homeostasis can
17 be probed by a set of mimotopes, the size of which can be tailored to the diagnostic goals by optimization.
18 The existing approaches for immunosignature [19, 20] or immunomic [21] analysis of the immunoglobulin
19 repertoires focus mostly on IgG and have used arrays of either 10² proteins or 10⁴-10⁵ random peptides.
20 The IgM repertoire has been previously probed by protein arrays [22] containing a biologically determined
21 representative set of autoantigens which is a structurally coarse approach. We set out to explore the
22 feasibility of a method that, similar to the self-protein “homunculus” arrays, targets a small set of
23 rationally selected probes, but also preserves the structural interpretability of peptides in a format
24 applicable for routine diagnostics.

1 **Materials and methods**

2 **Deep panning**

3 Human IgM was isolated from a sample of IgM enriched IVIg - IgM-Konzentrat (Biotest AG, Dreieich,
4 Germany, generously provided by Prof. Srinivasa Kaveri), while human monoclonal IgM paraprotein was
5 isolated from an IgM myeloma patient's serum selected from the biobank at the Center of Excellence for
6 Translational Research in Hematology at the National Hematology Hospital, Sofia (with the kind
7 cooperation of Dr. Lidiya Gurcheva). In both cases, IgM was purified using affinity chromatography with
8 polyclonal anti- μ antibody coupled to agarose (A9935, SIGMA-ALDRICH, USA). A 7-mer random peptide
9 library (E8100S, Ph.D. -7, New England Biolabs, USA) was panned overnight at 4°C on pooled human IgM
10 adsorbed on polystyrene plates at a concentration of 0.1 mg/ml, washed, eluted with glycine buffer at pH
11 2.7 and immediately brought to pH7. The eluate was transferred to a plate coated with monoclonal IgM
12 and incubated according to the same protocol, but this time the phage solution was collected after
13 adsorption and amplified once, according to Matochko et al. [23]. Briefly, the phage DNA was extracted
14 and the peptide-coding fragment amplified by PCR. The amplicons were subjected to deep sequencing
15 using the Next Seq platform (Illumina, USA), performed at the Sequencing Core Facility of Oslo University
16 Hospital.

17 **Mimotope selection**

18 The reads from the deep sequencing experiment were processed using the script provided by Matochko
19 et al. [23]. Reads found in 3-10 copies per million (CPM) were selected and subjected to clustering using
20 the GibbsCluster-2.0 method [24]. The number of clusters was optimized using the Kulback-Leibler
21 distance (KLD) from the background model of random sequences [24]. Position weighted matrices (PWM)
22 were defined for each cluster using pseudo counts as follows:

1

$$PWM_{k,j} = \log_2 \left(\frac{\sum I(X_{i,j} = k) + b_k \sqrt{N}}{(N + \sqrt{N})b_k} \right)$$

2 where $i=1..N$ are the rows of the alignment, $j=1:7$ are the columns of the alignment, $I(aa=k)$ is the indicator
3 function used to count the occurrences of amino acid k in column j and b_k is the background frequency of
4 amino acid k in the phage display library [25].

5 Using the PWMs, the median of the log odds (LO) scores of the peptides in each cluster was calculated.
6 Next, the probability of the occurrence of peptides with LO greater than the respective median in a set of
7 random peptides 10-fold larger than the analyzed library was determined empirically. Using this estimate,
8 the probability of the chance occurrence of as many peptides with scores higher than the median score
9 as observed in each cluster was calculated using the binomial distribution. The probabilities, thus found
10 for each cluster, were used as a criterion for their relevance to the mimotope library being generated.

11 The clustering of the mimotopes was visualized using the t-sne (t-Distributed Stochastic Neighbor
12 Embedding) algorithm [26] in its implementation using the faster Barnes-Hut algorithm (Rtsne package)
13 with theta parameter default value of 0.5 and a maximum of 1500 iterations [27].

14 The amino acid residues were described using the 5-dimensional scale of amino acid residue properties
15 published by Hellberg et al. [28]. The amino acid residue properties quantification used was the same as
16 for the pepStat binding normalization [29]. The 5 scales (z1-z5) were extracted as the latent variables
17 describing the major factors underlying the variability of amino acids in the space of 26 physicochemical
18 parameters. Thus, each peptide was represented by a vector of 35 scores corresponding to its seven
19 positions. For a comparison, the same number of random 7-mer peptides generated using the amino acid
20 background frequencies of Ph.D.-7 were clustered using the same parameters. The clusters in the t-sne
21 plot were labeled using k-mean clustering.

22 **Patients' sera**

1 Sera from randomly selected patients with glioblastoma multiforme (GBM), low grade glioma (G), brain
2 metastases of breast (MB) or lung (ML) cancers, as well as non-tumor bearing patients (C) (herniated disc
3 surgery, trauma, etc.) of the Neurosurgery Clinic of St. Ivan Rilski University Hospital, Sofia acquired
4 according to the rules of the ethics committee of the Medical University in Sofia, after its approval and
5 obtaining informed consent, were analyzed on the sets of peptides defined in microarray format. The sera
6 were aliquoted and stored at -20°C. Before staining the sera were thawed, incubated for 30 min at 37°C
7 for dissolution of IgM complexes, diluted 1:100 with PBS, pH 7.4, 0.05% Tween 20 with 0.1% BSA, further
8 incubated for 30 min at 37°C and filtered through 0.22µm filters before application on the chips.

9 **Peptide microarray**

10 The customized microarray chips were produced by PEPperPRINT™ (Heidelberg, Germany) by synthesis
11 in situ as 7-mer peptides attached to the surface through their C-terminus and a common spacer GGGs.
12 The layout was in a format of a single field of up to 5500 or five fields of up to 600 peptides in randomly
13 positioned duplicates. The chips were blocked for 60 minutes using PBS, pH 7.4, 0.05% Tween 20 with 1%
14 BSA on a rocker, washed 3x1 min with PBS, pH 7.4, 0.05% Tween 20 and incubated with sera in dilutions
15 equivalent to 0.01 mg/ml IgM (approx. 1:100 serum dilution) on a rocker overnight at 4°C. After 3x1
16 minute washing the chips were incubated with secondary antibodies at RT, washed, rinsed with distilled
17 water and dried by spinning in a vertical position in empty 50 ml test tubes at 100 x g for 2 minutes.

18 **Microarray data treatment**

19 The microarray images were acquired using a GenePix 4000 Microarray Scanner (Molecular Devices, USA).
20 The densitometry was done using the GenePix® Pro v6.0 software. All further analysis was performed
21 using publicly available packages of the R statistical environment for Windows (v3.4.1) (Bioconductor –
22 Biostrings, limma, pepStat, sva, e1071, Rtsne, clvalid, entropy, RankProd, multcomp) as well as in house
23 developed R scripts (<https://github.com/ansts/IgMimoPap1> and <https://github.com/ansts/IgMimoPap2>

1).

2 *Local normalization*

3 The spot intensities of a non-treated chip were subtracted from the stained chip spot intensities.
4 Treatment with secondary antibody only did not result in any binding. Next, the local background had to
5 be inferred due to the dense spot layout. To that end an approximated background was smoothed by
6 support vector regression and the result passed to the backgroundCorrect [30] function for local
7 normalization using the normexp (mle) method. The initial background approximation is described in the
8 supplement file 1.

9 *Global normalization*

10 The peptides with missing values (flagged “bad”) in some of the patients were removed. Log transformed
11 locally normalized microarray data were next normalized for amino acid composition dependent binding
12 using the ZpexQuad method of pepStat package [31]. This step is considered indispensable because of
13 the strong effect of amino acid composition on binding due mostly to electrostatic interactions. The amino
14 acid residue properties were quantified using the 5 dimensional descriptor z1-z5 of Sandberg et al. (1998)
15 [29]. This was followed by global normalization using normalizeCyclicLoess (method affy) from the
16 package limma [30] and subjected to batch effect compensation using the ComBat [32] function from
17 package sva whenever necessary (for different chip batches and for different channels). The data was
18 acquired in 2 different batches using the green (batch G) and the red channel (batch R). Eleven cases were
19 part of an earlier experiment and that subset of the data could be used also in this assay (batch P). For
20 batch effect compensation the groups were balanced by selecting a subset of the patients with relatively
21 even representation of each tested diagnosis. The criterion for inclusion of the 5 GBM patients from batch
22 “R” was the minimal difference of the mean and coefficient of variation from the mean and CV of the
23 whole group of GBM in that batch. Finally, each peptide binding intensity was represented by the mean
24 of its duplicates.

1 Because of a lack of a suitable negative control, the baseline binding was determined from the data. The
2 clean data used for the testing of the diagnostic potential was subjected to dimensionality reduction of
3 the peptide reactivities using t-sne [26], which clearly outlined a group of peptides with uniformly low
4 reactivities, (Supplement fig.1) that was considered background binding. The mean of the background
5 binding intensity was subtracted from the data before testing for significant reactivities.

6 A general linear model, followed by Tukey contrast, was used to compare the expression of reactivities in
7 the different tested libraries. The total correlation between the different peptide reactivity profiles of
8 each library with 10 patients' serum IgM was used as a measure of the redundancy of the library. The
9 profiles were viewed as points in a 10-mer space with each dimension corresponding to a patient. Total
10 correlation was calculated as KLD from the theoretical maximal entropy joint distribution of the profile
11 data, binned appropriately. To keep the number of bins comparable to the number of peptides, the
12 intensity values for each patient were discretized just to "high" and "low" relative to the median for the
13 studied library, effectively centering each library's data and, thus, removing the effect of the median
14 intensity of each library. To further limit the number of bins, randomly selected bins (the same for all
15 libraries) were aggregated in groups of 8, yielding 128 bins for the calculation of the baseline and the
16 actual joint distributions. To equalize the size of the libraries, random samples of 380 peptides were used
17 from each library, except for pepneglo which was excluded from this analysis due to its small size. The
18 mean KLD values were determined in a bootstrap procedure with the described sampling repeated with
19 replacement 100 times for the bin aggregation and 30 times for the library resampling producing 3000
20 samples. KLD for each library was thus calculated using the entropy::KL function.

21 The same approach could not be used for the transposed matrix, due to the disproportionately high
22 number of bins necessary to describe the distribution in high/low values of 500-1000 different peptides.
23 Instead, the mean correlation coefficient between the patient profiles across the peptides was used. To

1 compare the mean correlation values between libraries by GLM, the correlation values were converted
2 to z-scores.

3 Another test of the optimal sampling of the mimotope space was the mean nearest neighbor distance
4 (NND) per library which was normalized (nNND) relative to the theoretical mean distance between the
5 points in the 10-dimensional data cloud for the different libraries:

$$6 \quad nNND_L = \frac{NND_L}{\left(\frac{V}{N_L}\right)^{\frac{1}{k}}}$$

7 where $L=1..8$, N_L is the number of peptides in each library, k is the dimensionality (10 since the peptides
8 are compared on the basis to their reactivity to 10 patients' sera) and V is the volume of the data cloud
9 approximated as a k -dimensional ellipsoid [33]:

$$10 \quad V = \frac{2\pi^{\frac{k}{2}}}{k\Gamma\left(\frac{k}{2}\right)} \prod_{i=1}^k (2 * \sigma_i),$$

11 where σ_i is the standard deviation of the data along the i^{th} dimension (the values of the i^{th} serum) and Γ is
12 the gamma function. The logarithms of nNND were compared by general linear models (GLM).

13 Together with median intensity per library as a measure of the intensity of the signal, total correlation,
14 mean correlation between patients' sera reactivities and mean nearest neighbor distance were combined
15 to measure the optimization of the libraries. To make all four criteria positively correlated with the desired
16 qualities of the libraries, the sign of the correlation measures was changed. The rank product method (p
17 value for consistent "high expression" across the four tests) was used to test the optimization of the
18 libraries.

19

20 **Feature selection algorithm**

21 The composition of the library underwent a small change - 75 of the lowest scoring sequences from
22 pepnegrnd library were added to the library as a negative control. Because of the size limit of the final

1 library, this addition was done by replacement of 75 of the selected positive peptides. For the design of
2 a practical algorithm for extracting information about a particular diagnosis, a feature selection approach
3 was used based on recursive feature elimination. It used the quality of clustering of the predetermine
4 diagnostic groups when mapped to only the selected subset of features as a criterion. The cluster
5 separation of the cases of interest was measured using a combined clustering criterion:

$$6 \quad Crit = \frac{100 * Dunn(d, c) * BH\gamma(d, c)}{Conn(d, c)}$$

7 Where $Dunn(d, c)$ is the Dunn's clustering criterion [34, 35] which is based on single inter-cluster and intra-
8 cluster distances (extreme case), $BH\gamma(d, c)$ is the Baker-Hubert Gamma index [36] which gives the overall
9 agreement of distances and cluster assignment based on all the data, $Conn$ is the connectivity validation
10 measure for a given clustering partitioning based on 10 nearest neighbors [34], which emphasizes the
11 agreement of distances and partitions in the vicinity of each case; d is the matrix of distances between
12 the cases, and c is the diagnosis code for the cases. Supplement Figure 8 shows the typical traces of the
13 clustering criterion as a function of the number of features thus obtained with the best set of features
14 corresponding to the maximum of the clustering criterion.

15 Let F be the set of features used - the union of features of significant expression in any diagnosis ($n=380$).
16 The recursive elimination features selection from F was performed multiple times in a leave one out
17 setting, both as a direct validation as well as a bootstrap scheme used to improve further the
18 generalization of the proof of principle predictor.

19 The following algorithm was used for predicting GBM as contrasted to the rest of the cases. Let N be the
20 number of cases. Each of the N bootstrap samples contained the binding data for $N - 1$ cases. The
21 recursive elimination feature selection applied to each of those samples ultimately produced an
22 intersecting family of N feature sets. To ensure better generalization, only the features common for at
23 least 2 of those bootstrap generated feature sets were analyzed. The common features were pooled by

1 the number of sets they were common for (or “commonality”) in groups with commonality greater than
2 a threshold, i.e. – the group labeled 1 included all features, group n included features found in at least n
3 bootstrap sets and group 28 contained the features found in all bootstrap sets. These groups were used
4 to find the commonality level providing the best performing feature set.

5 **Results**

6 **Selection of 7-mer mimotopes**

7 We chose to pan a commercially available 7-mer random peptide phage display library of diversity 10^9 .
8 Thus, the size of the mimotopes would be in the range of the shorter linear B cell epitopes in the IEDB
9 database (<http://www.iedb.org/>). At the same time, the complete diversity of sequences of that length
10 could be interrogated. As a repertoire template we used an experimental preparation of human
11 immunoglobulins for intravenous use highly enriched in IgM, pool representative of the repertoire of at
12 least 10 000 healthy donors. The phage eluted from the IgM repertoire were adsorbed on a monoclonal
13 IgM to filter out non-specific binding of phage to the constant regions, and thereby focus binding only on
14 the mimotopes (Fig. 1). The peptide inserts were amplified and deep sequenced using the approach
15 described by Matochko et al. (2012) [23]. Two separate experiments starting with 20% of the original
16 phage library were performed (experiments A and B), while in a third one (C), a preamplified 20% sample
17 of the original phage library was used. The yield was 688 860 (experiment A), 518 533 (experiment B) and
18 131 475 (experiment C) unique reads. The Phred quality score cut off used was 32 with a probability for
19 an erroneous base call of 6×10^{-4} . The unique reads from experiments A and B were pooled so that only
20 one copy of a sequence existed in the final set ($n=1\ 114\ 845$). These were further filtered based on the
21 number of occurrences of each sequence (similar to the concept of sequencing depth). The criterion for
22 retaining a sequence was an occurrence in 3-10 CPM, yielding 224 087 unique sequences using the
23 following rationale. Considering that 79% of the single base changes lead to a change in the encoded
24 amino acid [37], the used Phred score led to about $1.1 \cdot 10^4$ sequences with one erroneous amino acid

1 residue if there is only one wrong base call per sequence and much fewer with multiple errors per
2 sequence. Since the criteria for inclusion is at least 3 CPM, the probability for the error to yield one
3 sequence already present in 2 CPM (so that it is included in the selected set, there are 182 424 sequences
4 at 2 CPM) is 7% and less than 0.3% for 2 and more incorrect sequences. This error rate could lead also to
5 a loss of about 80 potentially selected sequences when it affects a sequence presented at 3 CPM (n=89
6 167). So, this Phred score can slightly restrict the selected sequences (false negatives), but produces
7 practically no false positives. The probabilities for a change at the higher threshold (10 CPM) are negligible.
8 The limit of 10 CPM on the high copy numbers was applied after comparison of the distributions of the
9 number of clones by CPM between the original and the preamplified library (Fig. 2). It showed that 10
10 CPM was the threshold discriminating the original and preamplified libraries, with diversity in the latter
11 skewed towards highly proliferating clones. This fact was interpreted in view of the observation that the
12 affinity selection seemed to favor low CPM clones. Therefore, the high CPM clones were excluded to avoid
13 a possible contamination with non-selected clones having an advantage when they are highly
14 proliferating. This restriction led to the exclusion of 9.96% of the reads.

15 To optimally probe the mimotope sequence space, the library was clustered based on the sequences.
16 Thus, the set of 224 087 unique 7-mer mimotope sequences (from experiment A and B) were subjected
17 to clustering using the GibbsCluster-2.0 method [24] originally applied for inferring the specificity of
18 multiple peptide ligands tested on multiple MHC receptors. The algorithm was used in its stand-alone
19 version, which allowed an unlimited number of clusters and the use of large data sets. The number of
20 clusters was optimized using the Kullback-Leibler distance from the background model of random
21 sequences [24] in the range of 100 to 2500 clusters. This criterion indicated optimal clustering in 790
22 clusters (Fig. 3). PWMs were calculated from each cluster and used further to find the most relevant
23 peptide sequence as a representative of each cluster (Supplement file 2).

24 **Structural properties of the mimotope clusters**

1 The frequencies of the amino acids residues in the mimotopes selected from the phage library showed a
2 skewing in favor of the amino acid residues G,W,A,R,T,H,M,P,Q and against C,F,N,Y,I,L and S (Fig. 4A). The
3 evidence of selection stood out in the overall PWM, showing higher divergence from the random
4 distribution of the frequencies in the N-terminal part (Fig. 4B). L,F and C were not only underrepresented,
5 but they also were particularly reduced in frequency in the first position. The most pronounced deficit
6 was for C being even further selected against with respect to its already reduced representation in the
7 phage display library (accounted for by the background probabilities used to calculate the log odds).
8 Interestingly, proline was avoided at the first position, but was highly selected for in the second positions,
9 and in the C-terminus. G,W,A,T,M and E were also highly represented in the N-terminus.

10 As expected, the residues found selected for or against in particular N-terminal positions were also
11 involved in shaping the majority of the cluster features. The log odds (LO) forming the position weighted
12 matrices for the different clusters were further combined in a linearized form in an overall matrix of 140
13 positions/amino acids by 790 clusters. This matrix was used to determine the amino acid residue in each
14 position defining the clusters. Biclustering of the overall LO matrix (Supplement Figure 2) again showed
15 clearly that P, W and H in all positions stand out as the most informative residues that determine to a
16 great extent the formation of most of the 790 prototype mimotope clusters. Cysteine and several avoided
17 residues belong also to this cluster of informative residues due to their extremely low LO values. A small
18 number of mimotope clusters are defined by the presence of C, which is practically lacking from the rest
19 of the clusters.

20 **Generation of different mimotope libraries**

21 The maximal optimized library at this stage would consist of these 790 representative sequences as they
22 evenly (symmetrically) sample the mimotope sequence space, as ensured by the GibbsCluster algorithm.
23 The final clusters vary with respect to the probability of random occurrence of groups of the same size
24 and PWM in a random collection of peptides. The respective probability (based on binomial distribution)

1 was used to rank the profiles by significance. Only the top 594 clusters (with $p < 1e-4$) were retained, and
2 only the mimotope with the top LO score from each cluster was kept as a mimotope prototype for the
3 profile. This library was labeled peppos. Next, each of 2.3×10^6 random 7-mer peptide sequences was
4 tested against each of the PWM of mimotope clusters defined. For each random peptide, only the score
5 of the top scoring cluster was retained, and the peptides were ranked in the ascending order of these
6 scores. The lowest ranking peptides represented random sequences that were the least related to any of
7 the clusters in the selected library and the first 753 were, thus, used as a negative control (library
8 pepnegrnd).

9 Other libraries were also generated based on: 1) low probability for relevance to the cluster based on KLD
10 scores of the peptides in each cluster (pepneg and pepneglo); 2) 2 groups of 5 highest scoring clusters
11 (pep5 and pepother5); 3) random 7-mer peptides predicted to belong to 5 highest scoring clusters by
12 calculating the respective scores relative to the PWM of the clusters (pep5pred) and 4) random 7-mer
13 peptides (peprnd) (see Table 1 for description of all libraries). The number of peptides per library was
14 constrained by the size of the chip.

15 **Comparison between libraries**

16 Sera from patients with glioblastoma multiforme (GBM), brain metastases of breast cancers (MB) as well
17 as non-tumor bearing neurosurgery patients (C) (herniated disc surgery, trauma, etc.) were analyzed on
18 the sets of peptides synthesized in an oriented (C-terminus attached) planar microarray format. In the
19 first round of experiments, the 8 different libraries defined were compared based on the IgM reactivity in
20 the sera from 10 patients (Supplement Fig. 3 and 4). The data on the mean serum IgM reactivity of each
21 peptide across the different sera were grouped by library and the libraries were compared for their overall
22 reactivity using GLM (Fig. 5A). The proposed optimized library (peppos) had significantly higher ($p < 0.001$)
23 average reactivity than pepneg, peprnd or pepnegrnd. Interestingly, the library theoretically purged of

1 relevant reactivities (pepngrnd) had indeed the lowest reactivity, significantly lower than both the
2 weakly clustering peptides (pepneg) and the random sequences (peprnd) (Suppl. Table 1).

3 Next, the capacity of the different libraries to sample symmetrically the space of mimotope reactivities
4 was tested by comparing the total correlation of the IgM reactivity profiles of the peptide mapped on the
5 10 different sera (Fig. 5B). The total correlation is a KLD based multidimensional generalization of mutual
6 information. High total correlation would signify redundancy in the library with many peptides sharing
7 similar reactivity profiles. On the contrary, the minimal total correlation would be found in a library where
8 all reactivity profiles are as dissimilar as possible. The library peppos had the lowest total correlation,
9 while pepngrnd, and especially pep5pred had the top total correlation, indicating redundancy, while the
10 pools of 5 clusters – pep5 and pepother5 – had relatively low correlation. All differences, except between
11 the top two libraries (the library purged of sequences similar to the clusters and the library of predicted
12 sequences similar to 5 highly significant clusters), were significant (Supplement, Table 2).

13 Another way to test the optimal probing of the mimotope space in the different libraries is to compare
14 the mean nearest neighbor distance (MNND) of the scaled and centered data of mimotope staining
15 intensity mapped again to the 10 patients' sera IgM. Theoretically, the projections of the mimotope
16 reactivities of an optimal library in the space of multiple individual sera/repertoires should be widely and
17 symmetrically spaced. The MNND of peppos ranked second only to pepneglo (Fig. 5D) and was
18 significantly higher than MNND of all the other libraries (Supplement Table 4.).

19 Alternatively, the correlation of the serum profiles based on the different mimotopes (transposing the
20 matrices of the previous tests) can be viewed as a criterion for the capacity of the libraries to extract as
21 much information as possible from the IgM repertoire. Due to the extreme multidimensionality of this
22 feature matrix and the very small sample size (10 patients), the total correlation approach could not be
23 applied. Instead, the mean correlation between patient profile pairs was used to compare the libraries

1 after z transformation to allow comparison by linear models (Fig. 5C). Again, when tested against the
2 peppos library, the serum IgM reactivity profiles exhibited the lowest correlation - significantly lower
3 compared to the reactivities with the other libraries except for pepnegrnd and pepneglo (Supplement.
4 Table 3.)

5 Finally, all four criteria were summarized using a rank product test, which proved that reactivity with
6 peppos stands out from all the other tested libraries as the optimal among them for probing the IgM
7 repertoire (Table 2).

8 **Visualization of the Mimotope Space**

9 T distributed stochastic neighbor embedding (t-sne) based on the Barnes-Hut algorithm was used to
10 visualize the structure of the mimotope sequence space as represented by the general mimotope library
11 produced by deep panning (n=224 087). The sequences were represented by converting each amino acid
12 residue to a 5-dimensional vector of physical property scores as described in Materials and methods. Thus,
13 each 7-mer peptide is represented by a 35-dimensional vector. The correlation dimension showed that
14 the intrinsic dimension of the sequence space in this mapping was 11.25 for random peptides and 11.05
15 for the phage selected mimotope library. Principle component analysis indicated that the first 14 principle
16 components account for approximately 75% of its variance (Supplemental Figure 5). The t-sne mapping
17 was done after reducing the 35 dimensions to 14 by PCA (the initial.dims parameter of the Rtsne function).
18 Figure 6 shows the maps for the mimotope library and for an equal number of random 7-mer peptides
19 constructed using the residue background frequencies of the phage display library. Mapping also the five
20 most significant clusters defined by GibbsCluster (the pep5 library) on the t-sne map shows that the two
21 clustering approaches yield different results with only limited correlation between them. Nevertheless,
22 the mapping of the optimized library (peppos) showed that even on this map it covers quite symmetrically
23 the mimotope sequence space (Figure 7). To compare in more detail the mimotope space to the overall

1 random peptide space, equally sized random samples of 50 000 mimotopes and random peptides
2 together with the pepnegrnd library were merged. After removing the duplicate sequences, the mixture
3 contained 49964 mimotopes, 49950 random peptides and 684 peptides from library pepnegrnd. This set
4 of peptides was plotted again using the same approach and parameter values as for the whole library
5 (Figure 8). While most of the clusters were uniformly represented by the mimotope library, some
6 contained fewer mimotopes, although the sequence space was still sampled uniformly within these
7 clusters. The random peptide sequences selected for their minimal similarity to the mimotope clusters
8 defined in the previous section (library pepnegrnd) mapped to the areas of the low density of mimotopes.
9 The t-sne representation of the mixture of peptides was further clustered using k-means clustering in 350
10 clusters as a means of binning neighboring points. The number of random peptides and mimotopes were
11 counted in each cluster (Supplement Figure 6). This plot allowed to identify the parts of the peptide space
12 which were underrepresented by the mimotopes. According to this classification of the peptides,
13 approximately 42 % of the random peptides, 14% of the mimotopes and 85% of the pepnegrnd library
14 were in the underrepresented areas (Chi square, $p < 0.0001$). The underrepresented areas had very similar
15 sequence profiles to the normally represented areas, except for less abundant charged residues (Suppl.
16 Data file “t-sne cluster profiles”).

17 **Diagnostic potential of a rationally designed mimotope library**

18 A universal mimotope library sampling optimally the public IgM reactivities would have multiple
19 applications both in the theoretical research of antibody repertoires, as well as in the design of theranostic
20 tools. Having support for the hypothesis that the mimotope library pepnos, sampling major structural
21 clusters, is optimal when compared to a set of 8 other libraries, we next studied its diagnostic potential
22 using a larger set of patient sera ($n=34$). Due to the small data set, the main goal was a “proof of principle”
23 test demonstrating the capacity of the assay to provide feature subsets suitable for building predictors.
24 The distribution of patients by diagnosis is shown in Table 3. After cleaning, local, global normalization

1 and balancing the groups which warranted the use of ComBat [38] for the following batch compensation,
2 the reactivity data represented 28 patients' serum IgM binding to 586 peptides. The comparison between
3 the staining intensities of the mimotopes (or features) in the patients' diagnostic groups yielded
4 overlapping sets of reactivities significantly expressed in each diagnosis - 290 features for GBM, 263 for
5 ML, and 204 for C. Overall, 380 features showed significant reactivity in at least one of the diagnostic
6 groups. The "negative" peptides (library pepnegrnd) represented 49/206 non-significant and 24/380
7 significant reactivities (χ^2 , $p < 0.0001$). The finding of individuals with IgM reactive for some of them when
8 testing a larger group is not a surprise. That is why the background reactivity was considered more reliably
9 determined by the data analysis, rather than on the mean level of the pepnegrnd library.

10 A projection of the cases on the positive reactivities by multidimensional scaling (MDS) showed no
11 separation (Suppl. Fig.7). The feature space is highly multidimensional. The peptide library is not targeted
12 to any particular pathology, but represents a universal tool for IgM repertoire studies. Therefore, a feature
13 selection step is needed for each potential predictor of a given state. A recursive elimination algorithm
14 was applied, removing at each step the feature, the removal of which improved maximally the clustering
15 of the diagnostic groups (see Materials and methods section for details). Furthermore, this feature
16 selection strategy was included in a leave one out validation (LOOV) setting, used additionally as a
17 bootstrap approach for a further improvement of the feature selection.

18 Using this approach, we tested the capacity of the selected feature sets to separate dichotomously each
19 of the diagnoses from the rest. As expected, LOOV based on a very small training set did not show
20 significant prediction. In an attempt to improve the generalization of the prediction, next we used the
21 data from LOOV as a bootstrap scheme testing the predictor's efficiency as a function of the feature set
22 's "commonality" (see Materials and Methods). The predictor was constructed after dimensionality
23 reduction by MDS to two-dimensions and using a radial basis function (Gaussian) kernel-based support
24 vector machine. The performance of the model was measured using the Matthew's correlation coefficient

1 (MCC). As can be seen in Fig. 9, the most efficient predictors were composed on the basis of features with
2 intermediate “commonality”, i.e. – found in more than 50% of the bootstrap derived sets. Figure 10 shows
3 the distribution of these features after applying this algorithm to each of the three diagnoses separately.
4 Among the set of features selected to differentiate the GBM cases, we used further those of 50%
5 commonality (n=55) to predict the diagnoses of all the cases of batch “R”. By including cases omitted
6 before the batch compensation, this batch served both as a source of a validation set, as well as a control
7 for the non-confounding effect of the batch compensation. As can be seen from Fig. 11, the feature
8 selection strategy yielded a predictor classifying correctly the validation cases using an appropriately
9 tuned support vector machine model (the color shade of the plot stands for the probabilities of each point
10 to belong to the GBM or non-GBM set). The Supplemental Figure 9 shows a schematic representation of
11 the used feature selection algorithm.

1 **Discussion**

2 The introduction of high-throughput omics screening methods has expanded the knowledge base and the
3 potential for diagnosis. All these approaches have led to the identification of biomarkers as profiles
4 extracted from a particular dynamic diversity – proteome, genome, glycome, secretome, etc. A less
5 explored source of biomarker profiles is the IgM antibody repertoire. IgM antibodies appear early in the
6 course of an infection. However, they fall relatively fast, even after restimulation, providing a dynamic
7 signal. Consequently, IgM antibodies have gained interest as biomarkers of physiological or pathological
8 processes [39-43]. IgM antibodies are still underused as immunodiagnostics, although their interactions
9 with sets of antigens have been studied in a range of platforms [39, 42-45]. These studies indicate a
10 significant, but untapped, potential in systems level screening using arrays of appropriately designed
11 targets for the discovery of (natural) IgM signatures, enabling diagnosis and prognosis.

12 The use of the antibody repertoire as a source of biomarkers has been defined and approached in multiple
13 ways. First came the technically minimalistic, but conceptually loaded, semiquantitative immunoblotting,
14 developed 20 years ago. This technique served as no less than a paradigm setter for systems immunology
15 [46-51]. The further development produced methods that have been referred to as functional
16 immunomics [21] in terms of protein reactivities, as immunosignaturing [19] in terms of random peptide
17 libraries, or described as a deep panning technique [52] and in terms of IgOme of mimotopes selected
18 from random phage display libraries. Here we describe the design of the first mimotope library for the
19 analysis of the human IgM repertoire of reactivities recurrent in most individuals [12, 53, 54].

20 The deep panning approach relies on next generation sequencing (NGS), and thus requires balancing
21 between sequence fidelity and diversity. Within a number of reads limited by the technology (i.e. – on the
22 order of 10^7), high sequencing depths are impossible when the target diversity is itself on the order of 10^5 -
23 10^6 . Hence our choice of 3 CPM as a lower limit. In addition, highly represented sequences are to be
24 avoided, in order to reduce the effects of the collapse of diversity during phage display amplification.

1 Choosing reads of less than 10 CPM , we obtain a set of 224 087 mimotopes which does not contain any
2 of the parasitic sequences reported by Matochko et al. [55]. The use of monodispersed droplets-
3 compartmentalized phage amplification [56] can overcome the latter problem. But even with diversity
4 affected by discarding sequences of one and two copies per million on the one hand, and overgrowth of
5 phage clones on the other, our strategy still manages to find a general representation of the mimotope
6 sequence space by identifying clusters of mimotopes. This relatively small set of structural classes is
7 hypothesized to be related to the modular organization of the repertoire defined previously [57].
8 In general structural terms, the reduced sequence entropy of the selected mimotopes, especially of their
9 N-terminus, indicates selection and comes as a proof of the efficiency of the phage selection step. The
10 abundance of W and P is not surprising, since these amino acids are often found in epitopes and
11 carbohydrate mimotopes [58]. The central role of prolines in the natural antibody mimotopes has also
12 been observed previously [59]. Since Tchernychev et al. also used a phage display library, it is possible
13 that the high proline content is related to the structural basis of recognition by antibodies, because
14 prolines are associated with turns and flanking structures. Proline abundance may also be related to a
15 necessity to control the entropic component of the binding, which is important in the case of polyspecific
16 antibodies and peptide ligands.
17 The mimotope library of diversity 10^5 we derived by deep panning reflects the recurrent (also referred to
18 as public) IgM specificities found in the human population. It can be used as is in large arrays when
19 applicable. Using this library, a generalization of the mimotope sequence space based on Gibbs Clustering
20 was attempted, yielding the best clustering in 790 sequence clusters. The latter are viewed as broad
21 structural contexts, each including multiple polyspecificities. This classification was used to produce a
22 smaller and more applicable library for clinical use, of a subset of approx. 600 mimotopes (peppos) by
23 picking representative sequences from the most significant clusters. Thus, this library was designed to
24 optimally represent the mimotopes' main public reactivity patterns found in the phage selection

1 experiment. The proposed optimized library could be used as a tool for the study of the IgM public
2 repertoire, as a source of mimotopes for immunotherapeutics design [60-63], but it may serve also as a
3 multipurpose diagnostic tool.

4 To study the properties of the optimized small library peppos, we compared it to 7 other libraries (Table
5 1.). Far from exhaustive, this set of 7-mer libraries represented the major classes of alternatives to the
6 optimized library: - random peptides, restricted sets of mimotopes, different groups underrepresented
7 sequences and predicted sequences. When the 8 libraries were compared across patients, the reactivities
8 to the optimized library peppos had the lowest total correlation, while pepnegrnd, and especially the
9 predicted mimotopes pep5pred, had the top total correlation indicating redundancy of the reactivity
10 profiles. Thus, peppos proved the most informative, while the predicted pep5pred had the highest
11 redundancy. The latter library represents similar peptides defined by a rule (PWM) derived on the basis
12 of another library - pep5. Sequence profile captures the central structural theme, but mislead about the
13 existing diversity in the cluster. The other measures used to test peppos were two more correlation
14 measures that capture different aspects of this redundancy, as well as the mean intensity to ensure that
15 the library provides a sufficiently strong signal.

16 As a diagnostic tool, the optimized small library has some key properties that distinguish it from other
17 omic sets. Since it is designed to represent practically ubiquitous public specificities, the diagnostic profiles
18 would be quantitatively, rather than qualitatively, defined by combinations of reactivity to the different
19 mimotopes. Indeed, the sets of features (mimotope reactivities), significantly expressed in the different
20 diagnoses, were overlapping considerably. No single reactivity was correlating strongly with a whole
21 diagnostic group, but subsets of reactivities collectively could separate the diagnoses. Thus, feature
22 selection becomes essential for the design of predictors based on polyspecificities. Using the proposed
23 algorithms, the typical feature set tuned for a dichotomous separation of diagnoses contained between
24 28 and 111 sequences (median=66). The improvement of generalization by keeping only features

1 recurring in the bootstrap feature selection cycles helped reduce the overfitting of the models. The
2 optimal feature set for GBM diagnosis contained 55 mimotopes. Thus, if the library provides in the order
3 of 500 significant reactivities, the theoretical capacity of this approach is $>10^{70}$ different subsets. These
4 are further increased by at least 6 orders of magnitude when quantitative different profiles are
5 considered. Thus, the information provided by a typical IgM binding assay with the library is probably
6 enough to describe any physiological or pathological state of clinical relevance reflected in the IgM
7 repertoire. Of course, this is just an estimate of the resolution of the method. The number of naturally
8 occurring profiles and their correlation with clinically relevant states will determine the actual capacity.
9 Training with appropriately sized teaching sets for clinically relevant predictors will help clarify this.
10 The novelty of our approach is based on the combination of several previously existing concepts:
11 First, the antibody repertoire is compartmentalized in various ways (by isotype, B-cell subsets, nature or
12 origin of the antigens, relatedness to immunopathology, etc.). Early studies argued that the physiologically
13 autoreactive natural antibodies comprise such a consistent, organized compartment [48, 51, 64-67]. The
14 consistency of the natural antibody self-reactivity among individuals was considered evidence for the
15 existence of a relatively small set of preferred self-antigens. Such “public reactivities” are most probably
16 related to the germline repertoire of antibodies generated by evolutionarily encoded paratope features
17 and negative/positive selection [22, 68]. These antibodies were targeted using protein microarrays, the
18 utility of which has been previously demonstrated [21, 22, 43, 57]. Recently, the existence of structurally
19 distinct public V-regions has been analyzed using repertoire sequencing [12], noting that they are often
20 found in natural antibodies. If the repertoire should be read as a source of information providing
21 consistent patterns that can be mapped to physiological and pathological states, the public natural
22 autoreactivity seems to be a suitable but underused compartment. Respectively, IgM seems more suitable
23 than IgG or IgE (although IgA could be also very useful, especially for probing the microbiome influence).

1 Therefore, the template used for this study was a pooled human IgM preparation, which emphasizes
2 common specificities and dilutes out unique or rare ones.

3 Second, germline variable regions are characterized by polyspecificity or cross-reactivity with protein and
4 non-protein antigens [14]. This raises the question: Do nominal antigens actually exist for polyspecific
5 germline antibodies, and what is the structure of this convoluted (because it is polyspecific) repertoire of
6 reactivities? It seems that going for epitopes could be a way to approach the repertoire convolution. Yet,
7 the actual epitopes will be mostly conformational and hard to study. In similar tasks, mimotopes are often
8 used [69-72], although M.H. Van Regenmortel argues that mimotopes are of little use to structural
9 prediction of the B-cell epitope [71]. Therefore, their utility might be rather in the structural study of the
10 repertoire as a whole.

11 Third, the usage of peptide arrays for the analysis of the antibody repertoire is increasingly popular [52,
12 73-76]. It involves the use of random peptide arrays for extracting repertoire immunosignatures on the
13 one hand and deep panning of phage display libraries to analyze antibody responses on the other. Since
14 an antibody can often cross-react with a linear epitope that is part of the nominal conformational epitope
15 [71], the 7-residue library offers suitable short mimotopes as compared to typical B-cell epitopes.
16 Furthermore, from the Immunoepitope Database (<http://www.iedb.org>) collection of linear B cell
17 epitopes, 4821 of 45829 entries are less than 8 residues long.

18 The study of the IgM repertoire might be expected to give information about interactions that occur
19 mostly in the blood and the tissues with fenestrated vessels since, unlike IgG, IgM cannot easily cross the
20 normal vascular wall. Yet, IgM tissue deposits are a common finding in diverse inflammatory conditions
21 [77-79] and especially in the disorganized vasculature of the tumors, where they are a key element of the
22 innate immune surveillance mechanism [7, 80, 81]. Changes in the IgM repertoire further reflect B cell
23 function affected by antigenic, danger and inflammatory signals, but also by anatomical changes leading

1 to vascular permeability or disruption. Thus, IgM repertoire monitoring has the potential to provide
2 clinically relevant information about most of the pathologies involving inflammation and vascular
3 remodeling, as well as all types of cancer. The library also provides a rich source of mimotopes that can
4 be screened for different theranostic tasks focused on particular targets. On an omics scale, the optimized
5 mimotope library proposed here probes efficiently the relevant repertoire of public IgM reactivities
6 matching its dynamic diversity with potentially over 10^{70} distinct profiles. The major task ahead is
7 designing studies aimed at efficiently extracting specific diagnostic profiles and building appropriate
8 predictors, e.g. – for classifying immunotherapy responders, predicting the risk of malignancy in chronic
9 inflammation, etc.

10 **Competing interests**

11 The authors declare no competing interests.

12 **ACKNOWLEDGEMENTS** This work was performed with the support of EEA/Norway Grant BG09/D03-103
13 and the Bulgarian Fund for Scientific Research Grant D01-11/2016. The authors wish to thank Prof. Radha
14 Nagarajan, Prof. Ivanka Tsakovska and Prof. Soren Hairabedian for critically reading the manuscript and
15 a number of useful comments.

1 **Tables**

2 **Table 1. Libraries of 7-mer peptides studied.**

3

Library	Description	N
peppos	The sequence with the highest score for the respective position weighted matrix from each significant cluster (significant clusters are those for which the number of sequences with more than median PWM score is greater than the expected number of occurrences of such score in random peptides - $p < 0.0001$ by Binomial test) .	594
pep5	A group of 5 of the 288 clusters with best binomial $p < 1e-16$: clusters # 2,6,9,10,11.	600
pepothor5	A group of 5 of the 288 clusters with best binomial $p < 1e-16$: clusters # 115,61,55,53, 258.	1193
pep5pred	A hundred and fifty random peptides with log odds scores greater than the median score of the respective cluster for each of 5 clusters (# 2,6,9,10,11).	750
pepneg	The lowest scoring sequence (using KLD) from each significant cluster. These peptides are least certain to belong to the clusters.	594
pepneglo	Among the set of the lowest scoring peptides (pepneg) using GibbsCluster's own "Corrected" score - those with score < 5 ([24]).	82
pepnegrnd	From a set of 2×10^6 random 7-mer peptides, get the max score for each peptide when tested against all cluster PWM and take the 753 lowest scoring.	753
peprnd	800 random peptides.	800
	Total	5366

* The random sets are constructed with underlying frequencies in phage display library Ph.D -7 .

4

5

6

1 **Table 2.** Rank product test of four criteria for optimal mimotope library:

Library	Rank Products	p Value
pep5	4.864599	0.78057
pep5pred	5.825901	0.924705
pepneg	5.957892	0.937399
pepneglo	2.114743	0.054359
pepnegrnd	6.0548	0.945747
pepoth5	3.22371	0.318709
peppos	1.189207	0.001071
peprnd	4.864599	0.78057

2

3

4 **Table 3.** Patients tested using the optimized library.

5

Diagnosis	Abbr.	Batches			Total
		G	P	R	
Non-tumor bearing (control)	C	1	3	4	8
Glioblastoma Multiforme	GBM	2	4	9(5)*	15 (11)*
Lung Cancer (Brain Metastasis)	ML	2	4	3	9
Breast Cancer (Brain Metastasis)	MB	0	0	2(0)*	2 (0)*
Total					34

6 To balance the group sizes between batches, only 5/9 GBM samples from batch “R” were used and the
 7 breast cancer cases were omitted before batch compensation using the ComBat function. All cases in
 8 batch “R” were used in the validation step.

9

10 **Table 4.** Distribution of peptides among normally represented and underrepresented parts of the 7-mer
 11 peptide space.

	Random Peptides	Mimotopes	pepnegrnd
Normally Represented Region	35 166	43 708	100
Underrepresented Region	14 784	6 256	584

12

13

1 **Legends**

2 **Figure 1.** Schematic representation of the deep panning experiment.

3 **Figure 2.** Results from NGS. Plot of reads distribution by CPM (counts per million). Black – sample 1 –
4 selected from original library, Red – sample 2 – selected from a pre-amplified original library. Here the
5 selection caused deviation from the power law (the straight lines) with an emphasis on the lower CPM.
6 The pre-amplification enriches in better fit phages and leads to a considerable shrinkage of the highly
7 diverse compartment of low CPM where most of the targeted diversity lies.

8 **Figure 3.** Results from GibbsCluster of the mimotopes. Different predefined number of clusters were
9 screened for the quality of clustering measured by Kullback-Leibler Divergence – KLD. The inset shows
10 amplified scale around the peak KLD values.

11 **Figure 4.** Overall distribution of the amino acid residues in the selected 224 087 mimotopes as log odds
12 (LO) relative to the background probabilities. Total abundance (A) and color coded LO by position in the
13 7-mer (B). There is a clear skewing of the representation of the residues due to the selection by the IgM
14 repertoire which is indicated by the concentration of the non-random distribution in the free N-terminus
15 as compared to the C-terminus which is attached to the PIII protein of the phage.

16 **Figure 5.** Statistics testing the libraries' capacity to probe the mimotope reactivity space. A) Mean
17 reactivity of each peptide across patients grouped by library. The optimized library peppos has the
18 highest reactivity. For library content see Table 1. B) Total correlation of the peptide profiles grouped by
19 library across 10 patients. The optimized library peppos provides the least redundant information. C)
20 Mean correlation of patient profiles across the peptides in each library compared after z-transformation.
21 The optimized library peppos provides the most diverse characteristics of the patients, which indicates a
22 high potential for discrimination of different states but increases the requirements for the size of the
23 teaching sets to extract models of good generalization. D) Mean nearest neighbor distance of the
24 peptide profiles across 10 patients in each library compared after z-transformation. Again, the optimized
25 library peppos appears to sample the mimotope reactivity space evenly. The width of the bars is
26 proportional to the size of the sample.

27 **Figure 6.** Visualization of the 7-mer peptide space of the phage display selected mimotopes (A) and
28 equal number random peptides with the same background residue frequencies (B) constructed using t-
29 sne based on the Barnes-Hut algorithm. The peptides in five clusters (included in library pep5) are color
30 coded in (A). The peptide sequences were encoded using a 5-dimensional score reflecting basic
31 biophysical properties of the amino acids. Each sequence is represented by a 35 dimensional vector but
32 the t-sne mapping is performed after initial reduction to 14 dimensions by PCA. There is only moderate
33 correlation between the clustering visualized by t-sne and the GibbsCluster classification. The image is of
34 high resolution and can be zoomed for better detail inspection.

35 **Figure 7.** Visualization of the 7-mer peptide mimotopes sequence space with the optimized library
36 peppos marked in red (see fig.8 for details). Although, individual GibbsCluster defined clusters do not
37 coincide with those shown by t-sne, the mapping of the optimized library apparently probes quite
38 uniformly the mimotope sequence space.

39 **Figure 8.** Visualization of the 7-mer peptide mimotope sequence space representing a mixture of
40 random sample of 50 000 phage display selected mimotopes (red) and 50 000 random peptides(gray)
41 plus the pepnegrnd library (blue). A part of the peptide space is represented by mimotopes at a lower

1 density and the peptides unrelated to the defined 790 mimotope clusters map mostly to this area (blue
2 points). A high definition version of this figure is included in the supplemental information.

3 **Figure 9.** Matthew's correlation coefficient as a measure of the prediction quality for SVM models
4 constructed using GBM predicting feature sets of different minimal commonality. Minimal commonality
5 of n means that the features in the set are found in n or more of the bootstrap sets. The validation set
6 consists of the cases in batch "R" that were omitted from the batch compensated united sets. The
7 model predicts these cases as belonging to the same class as the rest of the respective cases in batch
8 "R". Since the values in batch "R" were not subject to batch compensation the validation also serves as a
9 control against confounding introduced by the ComBat function.

10 **Figure 10.** Venn diagram of the overlapping sets of features (peptide reactivities) of 50% commonality
11 calculated separately for each diagnosis. These are the optimal sets for the prediction of each diagnosis
12 based on the proof of principle training set. The figures correspond to the number of features in the
13 corresponding intersection.

14 **Figure 11.** Multidimensional scaling plot of cases in batch "R" based on the feature set of minimal
15 commonality of 14 (50%). See figure 9 for details. The encircled points correspond to the validation set.

16 References

- 17 1. Baumgarth N, Herman OC, Jager GC, Brown LE, Herzenberg LA, Chen J. B-1 and B-2 cell-derived
18 immunoglobulin M antibodies are nonredundant components of the protective response to influenza
19 virus infection. *J Exp Med*. 2000;192(2):271-80.
- 20 2. Ochsenbein AF, Fehr T, Lutz C, Suter M, Brombacher F, Hengartner H, et al. Control of early viral
21 and bacterial distribution and disease by natural antibodies. *Science*. 1999;286(5447):2156-9.
- 22 3. Vollmers HP, Brandlein S. Tumors: too sweet to remember? *Mol Cancer*. 2007;6:78. PubMed
23 PMID: 18053197.
- 24 4. Matter MS, Ochsenbein AF. Natural antibodies target virus-antibody complexes to organized
25 lymphoid tissue. *Autoimmun Rev*. 2008;7(6):480-6. PubMed PMID: 18558366.
- 26 5. Avrameas S, Guilbert B, Dighiero G. Natural antibodies against tubulin, actin myoglobin,
27 thyroglobulin, fetuin, albumin and transferrin are present in normal human sera, and monoclonal
28 immunoglobulins from multiple myeloma and Waldenstrom's macroglobulinemia may express similar
29 antibody specificities. *Ann Immunol (Paris)*. 1981;132C(2):231-6.
- 30 6. Panda S, Zhang J, Tan NS, Ho B, Ding JL. Natural IgG antibodies provide innate protection against
31 ficolin-opsonized bacteria. *EMBO J*. 2013;32(22):2905-19. doi: 10.1038/emboj.2013.199. PubMed PMID:
32 24002211; PubMed Central PMCID: PMC3831310.
- 33 7. Vollmers HP, Brandlein S. Natural antibodies and cancer. *N Biotechnol*. 2009;25(5):294-8. Epub
34 2009/05/16. doi: S1871-6784(09)00060-0 [pii] 10.1016/j.nbt.2009.03.016. PubMed PMID: 19442595.
- 35 8. Prieto JMB, Felipe MJB. Development, phenotype, and function of non-conventional B cells.
36 *Comparative immunology, microbiology and infectious diseases*. 2017;54:38-44. Epub 2017/09/17. doi:
37 10.1016/j.cimid.2017.08.002. PubMed PMID: 28916000.
- 38 9. Lobo PI. Role of Natural Autoantibodies and Natural IgM Anti-Leucocyte Autoantibodies in Health
39 and Disease. *Front Immunol*. 2016;7:198. doi: 10.3389/fimmu.2016.00198. PubMed PMID: 27375614;
40 PubMed Central PMCID: PMC4893492.
- 41 10. Rothstein TL, Griffin DO, Holodick NE, Quach TD, Kaku H. Human B-1 cells take the stage. *Annals*
42 *of the New York Academy of Sciences*. 2013;1285:97-114. doi: 10.1111/nyas.12137. PubMed PMID:
43 PMC4429725.

- 1 11. Weller S, Braun MC, Tan BK, Rosenwald A, Cordier C, Conley ME, et al. Human blood IgM
2 "memory" B cells are circulating splenic marginal zone B cells harboring a pre-diversified immunoglobulin
3 repertoire. *Blood*. 2004. PubMed PMID: 15191950.
- 4 12. Greiff V, Weber CR, Palme J, Bodenhofer U, Miho E, Menzel U, et al. Learning the High-
5 Dimensional Immunogenomic Features That Predict Public and Private Antibody Repertoires. *The Journal*
6 *of Immunology*. 2017;199(8):2985-97. doi: 10.4049/jimmunol.1700594.
- 7 13. Van Regenmortel MH. Specificity, polyspecificity, and heterospecificity of antibody-antigen
8 recognition. *J Mol Recognit*. 2014;27(11):627-39. doi: 10.1002/jmr.2394. PubMed PMID: 25277087.
- 9 14. Willis JR, Briney BS, DeLuca SL, Crowe JE, Jr., Meiler J. Human Germline Antibody Gene Segments
10 Encode Polyspecific Antibodies. *PLoS Comput Biol*. 2013;9(4):e1003045. doi:
11 10.1371/journal.pcbi.1003045.
- 12 15. Cohen IR, Young DB. Autoimmunity, microbial immunity and the immunological homunculus.
13 *Immunol Today*. 1991;12(4):105-10.
- 14 16. Avrameas S, Guilbert B, Mahana W, Matsiota P, Ternynck T. Recognition of self and non-self
15 constituents by polyspecific autoreceptors. *Int Rev Immunol*. 1988;3(1-2):1-15.
- 16 17. Avrameas S. Natural autoantibodies: from 'horror autotoxicus' to 'gnothi seauton'. *Immunol*
17 *Today*. 1991;12(5):154-9.
- 18 18. Reynolds AE, Kuraoka M, Kelsoe G. Natural IgM is produced by CD5- plasma cells that occupy a
19 distinct survival niche in bone marrow. *J Immunol*. 2015;194(1):231-42. Epub 2014/11/28. doi:
20 10.4049/jimmunol.1401203. PubMed PMID: 25429072; PubMed Central PMCID: PMC4272881.
- 21 19. Hughes AK, Cichacz Z, Scheck A, Coons SW, Johnston SA, Stafford P. Immunosignaturing Can
22 Detect Products from Molecular Markers in Brain Cancer. *PLoS ONE*. 2012;7(7):e40201. doi:
23 10.1371/journal.pone.0040201.
- 24 20. Stafford P, Halperin R, Legutki JB, Magee DM, Galgiani J, Johnston SA. Physical Characterization of
25 the "Immunosignaturing Effect". *Molecular & Cellular Proteomics*. 2012;11(4). doi:
26 10.1074/mcp.M111.011593.
- 27 21. Quintana FJ, Hagedorn PH, Elizur G, Merbl Y, Domany E, Cohen IR. Functional immunomics:
28 microarray analysis of IgG autoantibody repertoires predicts the future response of mice to induced
29 diabetes. *Proc Natl Acad Sci U S A*. 2004;101 Suppl 2:14615-21. PubMed PMID: 15308778.
- 30 22. Merbl Y, Zucker-Toledano M, Quintana FJ, Cohen IR. Newborn humans manifest autoantibodies
31 to defined self molecules detected by antigen microarray informatics. *J Clin Invest*. 2007;117(3):712-8.
32 PubMed PMID: 17332892.
- 33 23. Matochko WL, Chu K, Jin B, Lee SW, Whitesides GM, Derda R. Deep sequencing analysis of phage
34 libraries using Illumina platform. *Methods*. 2012;58(1):47-55. doi:
35 <http://dx.doi.org/10.1016/j.ymeth.2012.07.006>.
- 36 24. Andreatta M, Lund O, Nielsen M. Simultaneous alignment and clustering of peptide data using a
37 Gibbs sampling approach. *Bioinformatics*. 2013;29(1):8-14. doi: 10.1093/bioinformatics/bts621.
- 38 25. de Hoon M, Vitkup D. Comparative Systems Biology of the Sporulation Initiation Network in
39 Prokaryotes. In: Eskin E, Ideker T, Raphael B, Workman C, editors. *Systems Biology and Regulatory*
40 *Genomics: Joint Annual RECOMB 2005 Satellite Workshops on Systems Biology and on Regulatory*
41 *Genomics*, San Diego, CA, USA; December 2-4, 2005, Revised Selected Papers. Berlin, Heidelberg: Springer
42 Berlin Heidelberg; 2006. p. 62-9.
- 43 26. van der Maaten LJP, Hinton G. Visualizing Data using t-SNE. *Journ Machine Learning Res*.
44 2008;9:27.
- 45 27. van der Maaten L. Accelerating t-SNE using Tree-Based Algorithms. *J Machine Learn Res*.
46 2014;15:3221-45.

- 1 28. Hellberg S, Sjoestroem M, Skagerberg B, Wold S. Peptide quantitative structure-activity
2 relationships, a multivariate approach. *Journal of Medicinal Chemistry*. 1987;30(7):1126-35. doi:
3 10.1021/jm00390a003.
- 4 29. Sandberg M, Eriksson L, Jonsson J, Sjostrom M, Wold S. New chemical descriptors relevant for the
5 design of biologically active peptides. A multivariate characterization of 87 amino acids. *J Med Chem*.
6 1998;41(14):2481-91. Epub 1998/07/03. doi: 10.1021/jm9700575. PubMed PMID: 9651153.
- 7 30. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers differential expression
8 analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res*. 2015;43(7):e47. doi:
9 10.1093/nar/gkv007. PubMed PMID: 25605792; PubMed Central PMCID: PMC4402510.
- 10 31. Imholte G, Sauteraud R, Gottardo R. Analyzing Peptide Microarray Data with the R pepStat
11 Package. *Methods Mol Biol*. 2016;1352:127-42. Epub 2015/10/23. doi: 10.1007/978-1-4939-3037-1_10.
12 PubMed PMID: 26490472.
- 13 32. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
14 empirical Bayes methods. *Biostatistics (Oxford, England)*. 2007;8(1):118-27. Epub 2006/04/25. doi:
15 10.1093/biostatistics/kxj037. PubMed PMID: 16632515.
- 16 33. Wilson JA. Volume of n-dimensional ellipsoid. *Scientia Acta Xaveriana*. 2010;1(1):101-6.
- 17 34. Handl J, Knowles J, Kell DB. Computational cluster validation in post-genomic data analysis.
18 *Bioinformatics*. 2005;21(15):3201-12. Epub 2005/05/26. doi: 10.1093/bioinformatics/bti517. PubMed
19 PMID: 15914541.
- 20 35. Dunn JC. Well-Separated Clusters and Optimal Fuzzy Partitions. *Journal of Cybernetics*.
21 1974;4(1):95-104. doi: 10.1080/01969727408546059.
- 22 36. Baker FB, Hubert LJ. Measuring the Power of Hierarchical Cluster Analysis. *Journal of the American*
23 *Statistical Association*. 1975;70(349):31-8. doi: 10.2307/2285371.
- 24 37. Abdullah T, Faiza M, Pant P, Rayyan Akhtar M, Pant P. An Analysis of Single Nucleotide Substitution
25 in Genetic Codons - Probabilities and Outcomes. *Bioinformation*. 2016;12(3):98-104. doi:
26 10.6026/97320630012098. PubMed PMID: PMC5267951.
- 27 38. Nygaard V, Rodland EA, Hovig E. Methods that remove batch effects while retaining group
28 differences may lead to exaggerated confidence in downstream analyses. *Biostatistics (Oxford, England)*.
29 2016;17(1):29-39. Epub 2015/08/15. doi: 10.1093/biostatistics/kxv027. PubMed PMID: 26272994;
30 PubMed Central PMCID: PMC4679072.
- 31 39. Silverman GJ, Srikrishnan R, Germar K, Goodyear CS, Andrews KA, Ginzler EM, et al. Genetic
32 imprinting of autoantibody repertoires in systemic lupus erythematosus patients. *Clin Exp Immunol*.
33 2008;153(1):102-16. Epub 2008/05/31. doi: 10.1111/j.1365-2249.2008.03680.x. PubMed PMID:
34 18510544; PubMed Central PMCID: PMC432104.
- 35 40. Sharron B-Z, Dror YK, Gittit D, Asaf M, Yifat M, Francisco JQ, et al. Individual and meta-immune
36 networks. *Physical Biology*. 2013;10(2):025003.
- 37 41. Mao J, Ladd J, Gad E, Rastetter L, Johnson MM, Marzbani E, et al. Mining the pre-diagnostic
38 antibody repertoire of TgMMTV-neu mice to identify autoantibodies useful for the early detection of
39 human breast cancer. *J Transl Med*. 2014;12:121. Epub 2014/06/03. doi: 10.1186/1479-5876-12-121.
40 PubMed PMID: 24886063; PubMed Central PMCID: PMC4022541.
- 41 42. Butvilovskaya VI, Popletaeva SB, Chechetkin VR, Zubtsova ZI, Tsybulskaya MV, Samokhina LO, et
42 al. Multiplex determination of serological signatures in the sera of colorectal cancer patients using
43 hydrogel biochips. *Cancer Med*. 2016. doi: 10.1002/cam4.692. PubMed PMID: 26992329.
- 44 43. Merbl Y, Itzchak R, Vider-Shalit T, Louzoun Y, Quintana FJ, Vadai E, et al. A Systems Immunology
45 Approach to the Host-Tumor Interaction: Large-Scale Patterns of Natural Autoantibodies Distinguish
46 Healthy and Tumor-Bearing Mice. *PLoS ONE*. 2009;4(6):e6053. doi: 10.1371/journal.pone.0006053.
- 47 44. Stafford P, Wrapp D, Johnston S. General assessment of humoral activity in healthy humans.
48 *Molecular & Cellular Proteomics*. 2016. doi: 10.1074/mcp.M115.054601.

- 1 45. Campbell CT, Gulley JL, Oyelaran O, Hodge JW, Schlom J, Gildersleeve JC. Serum Antibodies to
2 Blood Group A Predict Survival on PROSTVAC-VF. *Clinical Cancer Research*. 2013;19(5):1290-9. doi:
3 10.1158/1078-0432.ccr-12-2478.
- 4 46. Haury M, Grandien A, Sundblad A, Coutinho A, Nobrega A. Global analysis of antibody repertoires.
5 1. An immunoblot method for the quantitative screening of a large number of reactivities. *Scand J*
6 *Immunol*. 1994;39(1):79-87. PubMed PMID: 8290896.
- 7 47. Stahl D, Yeshurun M, Gorin NC, Sibrowski W, Kaveri SV, Kazatchkine MD. Reconstitution of Self-
8 Reactive Antibody Repertoires of Autologous Plasma IgM in Patients with Non-Hodgkin's Lymphoma
9 Following Myeloablative Therapy. *Clinical Immunology*. 2001;98(1):31-8. doi:
10 <http://dx.doi.org/10.1006/clim.2000.4949>.
- 11 48. Mouthon L, Haury M, Lacroix-Desmazes S, Barreau C, Coutinho A, Kazatchkine MD. Analysis of the
12 normal human IgG antibody repertoire. Evidence that IgG autoantibodies of healthy adults recognize a
13 limited and conserved set of protein antigens in homologous tissues. *J Immunol*. 1995;154(11):5769-78.
14 PubMed PMID: 7751627.
- 15 49. Lacroix-Desmazes S, Mouthon L, Coutinho A, Kazatchkine MD. Analysis of the natural human IgG
16 antibody repertoire: life-long stability of reactivities towards self antigens contrasts with age-dependent
17 diversification of reactivities against bacterial antigens. *Eur J Immunol*. 1995;25(9):2598-604. PubMed
18 PMID: 7589132.
- 19 50. Nobrega A, Haury M, Grandien A, Malanchere E, Sundblad A, Coutinho A. Global analysis of
20 antibody repertoires. II. Evidence for specificity, self-selection and the immunological "homunculus" of
21 antibodies in normal serum. *Eur J Immunol*. 1993;23(11):2851-9.
- 22 51. Mouthon L, Nobrega A, Nicolas N, Kaveri SV, Barreau C, Coutinho A, et al. Invariance and
23 restriction toward a limited set of self-antigens characterize neonatal IgM antibody repertoires and prevail
24 in autoreactive repertoires of healthy adults. *Proc Natl Acad Sci U S A*. 1995;92(9):3839-43.
- 25 52. Ryvkin A, Ashkenazy H, Smelyanski L, Kaplan G, Penn O, Weiss-Ottolenghi Y, et al. Deep Panning:
26 steps towards probing the IgOme. *PLoS One*. 2012;7(8):e41469. doi: 10.1371/journal.pone.0041469.
27 PubMed PMID: 22870226; PubMed Central PMCID: PMC3409857.
- 28 53. Truck J, Ramasamy MN, Galson JD, Rance R, Parkhill J, Lunter G, et al. Identification of antigen-
29 specific B cell receptor sequences using public repertoire analysis. *J Immunol*. 2015;194(1):252-61. Epub
30 2014/11/14. doi: 10.4049/jimmunol.1401405. PubMed PMID: 25392534; PubMed Central PMCID:
31 PMC4272858.
- 32 54. Gu H, Tarlinton D, Muller W, Rajewsky K, Forster I. Most peripheral B cells in mice are ligand
33 selected. *J Exp Med*. 1991;173(6):1357-71.
- 34 55. Matochko WL, Cory Li S, Tang SK, Derda R. Prospective identification of parasitic sequences in
35 phage display screens. *Nucleic Acids Res*. 2014;42(3):1784-98. Epub 2013/11/13. doi:
36 10.1093/nar/gkt1104. PubMed PMID: 24217917; PubMed Central PMCID: PMC3919620.
- 37 56. Derda R, Tang SK, Li SC, Ng S, Matochko W, Jafari MR. Diversity of phage-displayed libraries of
38 peptides during panning and amplification. *Molecules*. 2011;16(2):1776-803. Epub 2011/02/23. doi:
39 10.3390/molecules16021776. PubMed PMID: 21339712.
- 40 57. Madi A, Hecht I, Bransburg-Zabary S, Merbl Y, Pick A, Zucker-Toledano M, et al. Organization of
41 the autoantibody repertoire in healthy newborns and adults revealed by system level informatics of
42 antigen microarray data. *Proceedings of the National Academy of Sciences*. 2009;106(34):14484-9. doi:
43 10.1073/pnas.0901528106.
- 44 58. Luo P, Agadjanyan M, Qiu J, Westerink MA, Steplewski Z, Kieber-Emmons T. Antigenic and
45 immunological mimicry of peptide mimotopes of Lewis carbohydrate antigens. *Mol Immunol*.
46 1998;35(13):865-79. PubMed PMID: 9839555.
- 47 59. Tchernychev B, Cabilly S, Wilchek M. The epitopes for natural polyreactive antibodies are rich in
48 proline. *Proc Natl Acad Sci U S A*. 1997;94(12):6335-9.

- 1 60. Scott JK. Discovering peptide ligands using epitope libraries. [Review]. Trends in Biochemical
2 Sciences. 1992;17(7):241-5.
- 3 61. Westerink MA, Giardina PC, Apicella MA, Kieber-Emmons T. Peptide mimicry of the
4 meningococcal group C capsular polysaccharide. Proc Natl Acad Sci U S A. 1995;92(9):4021-5.
- 5 62. Kieber-Emmons T. Peptide mimotopes of carbohydrate antigens. Immunol Res. 1998;17(1-2):95-
6 108. PubMed PMID: 9479572.
- 7 63. Pashov A, Canziani G, Monzavi-Karbassi B, Kaveri SV, Macleod S, Saha R, et al. Antigenic properties
8 of peptide mimotopes of HIV-1-associated carbohydrate antigens. J Biol Chem. 2005;280(32):28959-65.
9 PubMed PMID: 15955803.
- 10 64. Cohen IR. The cognitive paradigm and the immunological homunculus. Immunol Today.
11 1992;13(12):490-4. doi: 10.1016/0167-5699(92)90024-2. PubMed PMID: 1463581.
- 12 65. Cohen IR. Autoantibody repertoires, natural biomarkers, and system controllers. Trends Immunol.
13 2013;34(12):620-5. doi: 10.1016/j.it.2013.05.003. PubMed PMID: 23768955.
- 14 66. Lacroix-Desmazes S, Mouthon L, Pashov A, Barreau C, Kaveri SV, Kazatchkine MD. Analysis of
15 antibody reactivities toward self antigens of IgM of patients with Waldenstrom's macroglobulinemia. Int
16 Immunol. 1997;9(8):1175-83.
- 17 67. Mouthon L, Lacroix-Desmazes S, Nobrega A, Barreau C, Coutinho A, Kazatchkine MD. The self-
18 reactive antibody repertoire of normal human serum IgM is acquired in early childhood and remains stable
19 throughout life. Scand J Immunol. 1996;44:243-51.
- 20 68. Hardy RR, Hayakawa K. Positive and negative selection of natural autoreactive B cells. Adv Exp
21 Med Biol. 2012;750:227-38. doi: 10.1007/978-1-4614-3461-0_17. PubMed PMID: 22903678.
- 22 69. Putterman C, Deocharan B, Diamond B. Molecular analysis of the autoantibody response in
23 peptide-induced autoimmunity. J Immunol. 2000;164(5):2542-9.
- 24 70. Pashov A, Monzavi-Karbassi B, Kieber-Emmons T. Immune surveillance and immunotherapy:
25 Lessons from carbohydrate mimotopes. Vaccine. 2009;27(25-26):3405-15. PubMed PMID: 19200843.
- 26 71. Van Regenmortel MH. What is a B-cell epitope? Methods Mol Biol. 2009;524:3-20. doi:
27 10.1007/978-1-59745-450-6_1. PubMed PMID: 19377933.
- 28 72. Huang J, He B, Zhou P. Mimotope-based prediction of B-cell epitopes. Methods Mol Biol.
29 2014;1184:237-43. doi: 10.1007/978-1-4939-1115-8_13. PubMed PMID: 25048128.
- 30 73. Weber LK, Palermo A, Kugler J, Armant O, Isse A, Rentschler S, et al. Single amino acid
31 fingerprinting of the human antibody repertoire with high density peptide arrays. J Immunol Methods.
32 2017;443:45-54. Epub 2017/02/09. doi: 10.1016/j.jim.2017.01.012. PubMed PMID: 28167275.
- 33 74. Weiss-Ottolenghi Y, Gershoni JM. Profiling the IgOme: meeting the challenge. FEBS Lett.
34 2014;588(2):318-25. Epub 2013/11/19. doi: 10.1016/j.febslet.2013.11.005. PubMed PMID: 24239539.
- 35 75. Navalkar KA, Johnston SA, Stafford P. Peptide based diagnostics: Are random-sequence peptides
36 more useful than tiling proteome sequences? J Immunol Methods. 2014. Epub 2014/12/17. doi:
37 10.1016/j.jim.2014.12.002. PubMed PMID: 25497701.
- 38 76. Legutki JB, Zhao ZG, Greving M, Woodbury N, Johnston SA, Stafford P. Scalable high-density
39 peptide arrays for comprehensive health monitoring. Nature communications. 2014;5:4785. doi:
40 10.1038/ncomms5785. PubMed PMID: 25183057.
- 41 77. Kulthanan K, Pinkaew S, Suthipinittharm P. Diagnostic value of IgM deposition at the dermo-
42 epidermal junction. Int J Dermatol. 1998;37(3):201-5. Epub 1998/04/29. PubMed PMID: 9556108.
- 43 78. Borrelli M, Maglio M, Agnese M, Paparo F, Gentile S, Colicchio B, et al. High density of
44 intraepithelial $\gamma\delta$ lymphocytes and deposits of immunoglobulin (Ig)M anti-tissue transglutaminase
45 antibodies in the jejunum of coeliac patients with IgA deficiency. Clinical and experimental immunology.
46 2010;160(2):199-206. doi: 10.1111/j.1365-2249.2009.04077.x. PubMed PMID: PMC2857942.

- 1 79. Chan RK, Ding G, Verna N, Ibrahim S, Oakes S, Austen WG, Jr., et al. IgM binding to injured tissue
2 precedes complement activation during skeletal muscle ischemia-reperfusion. *J Surg Res.* 2004;122(1):29-
3 35. Epub 2004/11/04. doi: 10.1016/j.jss.2004.07.005. PubMed PMID: 15522311.
- 4 80. Hensel F, Hermann R, Schubert C, Abe N, Schmidt K, Franke A, et al. Characterization of
5 glycosylphosphatidylinositol-linked molecule CD55/decay-accelerating factor as the receptor for antibody
6 SC-1-induced apoptosis. *Cancer Res.* 1999;59(20):5299-306. PubMed PMID: 10537313.
- 7 81. Vollmers HP, Brandlein S. The "early birds": natural IgM antibodies and immune surveillance.
8 *Histol Histopathol.* 2005;20(3):927-37. PubMed PMID: 15944943.

9

10 **Author Contribution:**

11 A.P. conceptualized the project, analyzed the results performing all the in silico work, supervised
12 experiments except for the sequencing as well as the overall project execution and prepared the
13 manuscript;

14 M. Hadzhieva ran the phage display experiments;

15 V.K. and M. T. ran the microarray experiments up to data processing, catalogued and maintained the
16 seroteque;

17 V.S., M. Heinz and L.A.M.Z. carried out the DNA isolation, PCR and sequencing;

18 E.H. supervised the sequencing task, participated in conceptualizing the project and the preparation of
19 the manuscript;

20 S. P. and M.T. performed the data processing of microarray scans;

21 T.V. and T.K.E participated in conceptualizing the project, analysis of the results and the preparation of
22 the manuscript;

23 D.F. was responsible for the patient selection, informed consent, ethics committee protocol preparation,
24 blood collection and serum preparation.

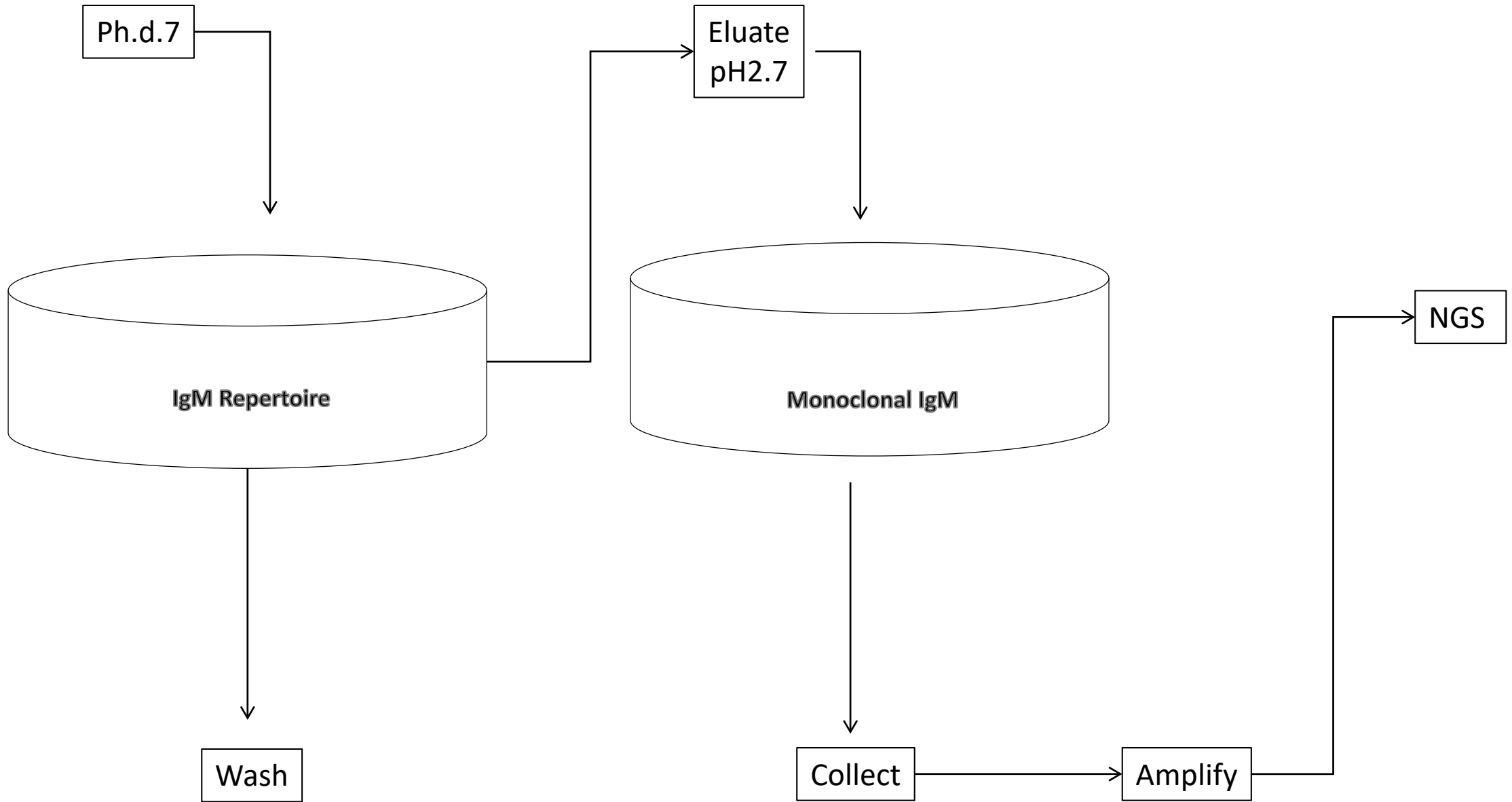


Fig.1

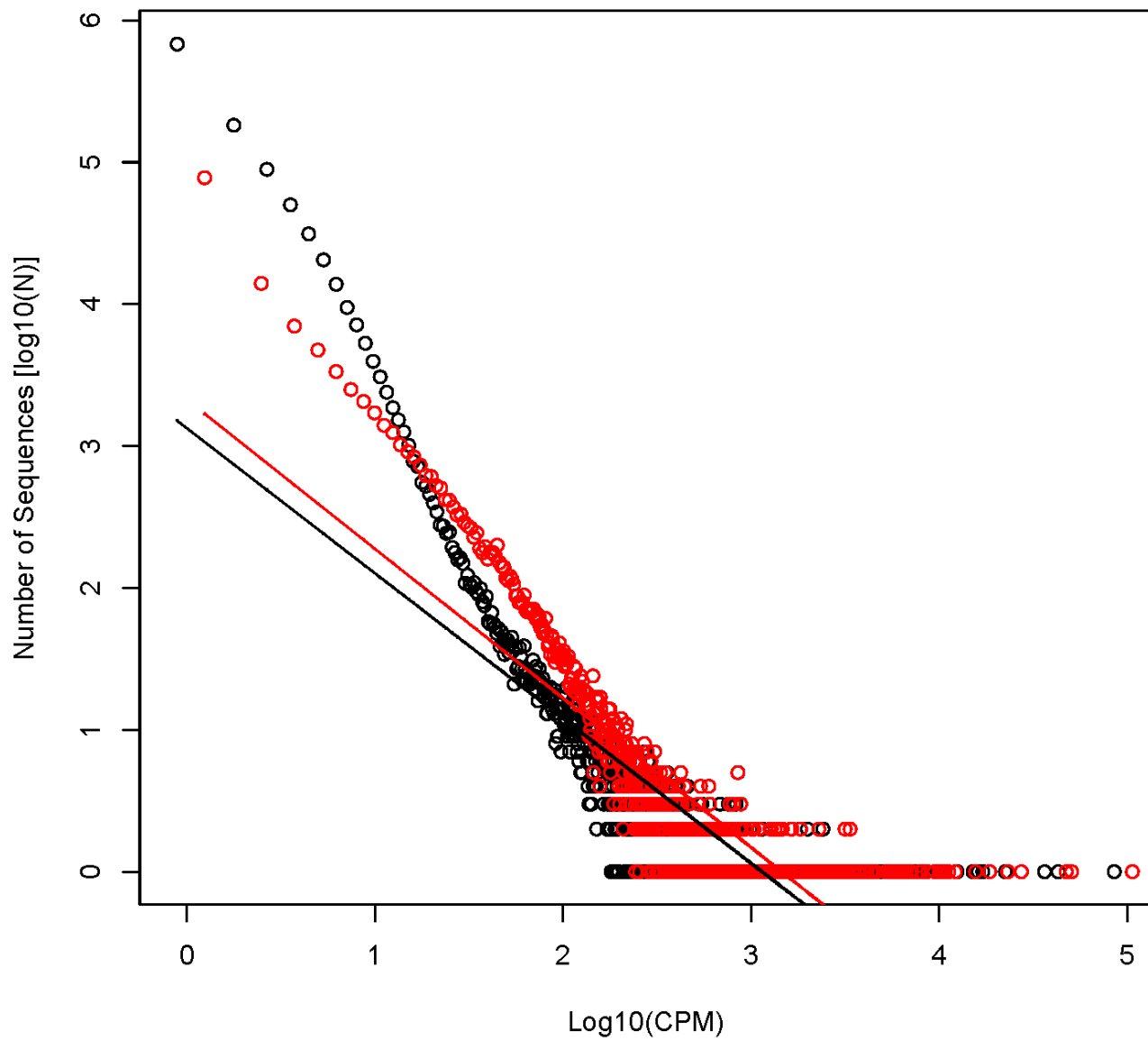


Fig. 2

Fig. 3

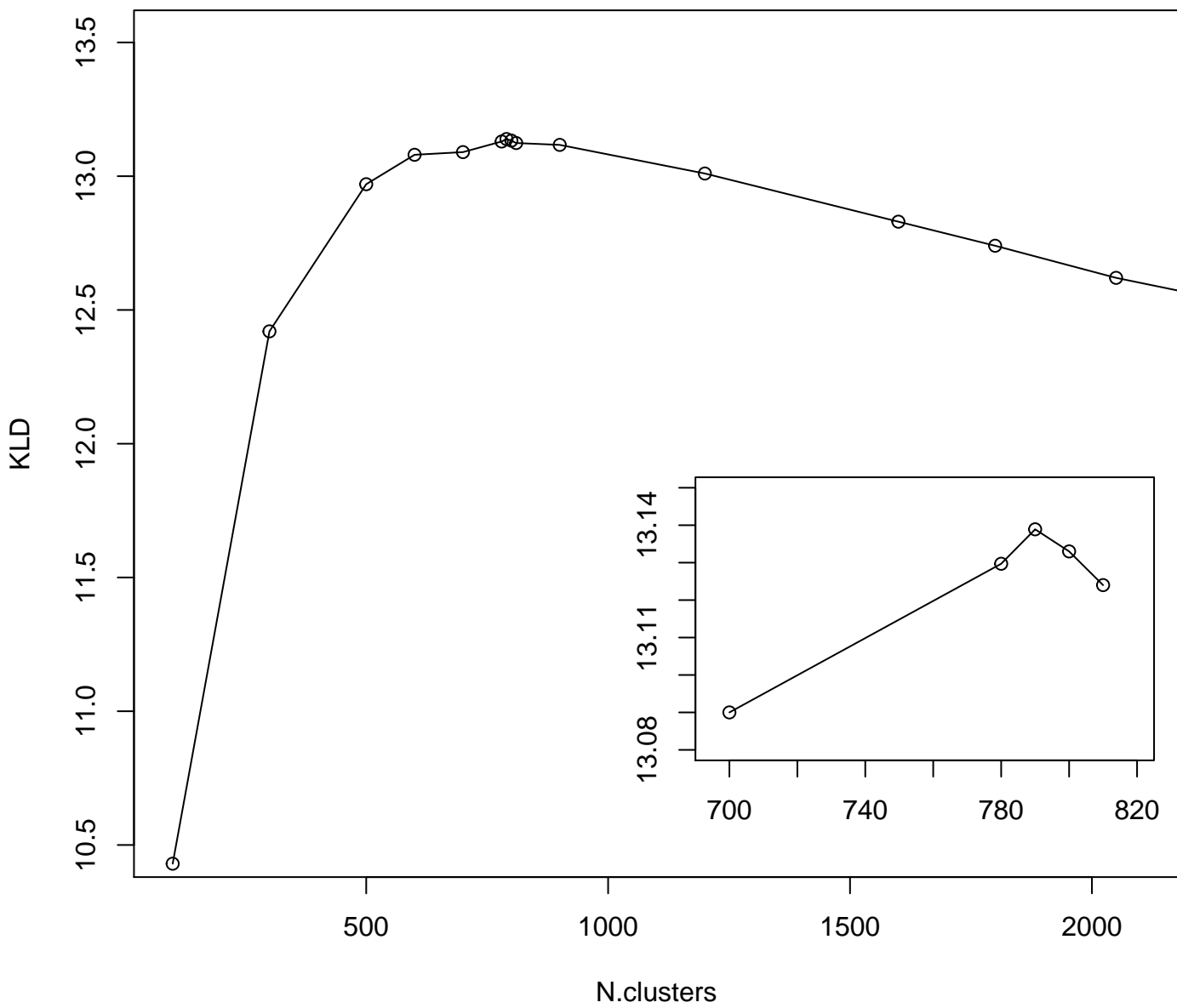
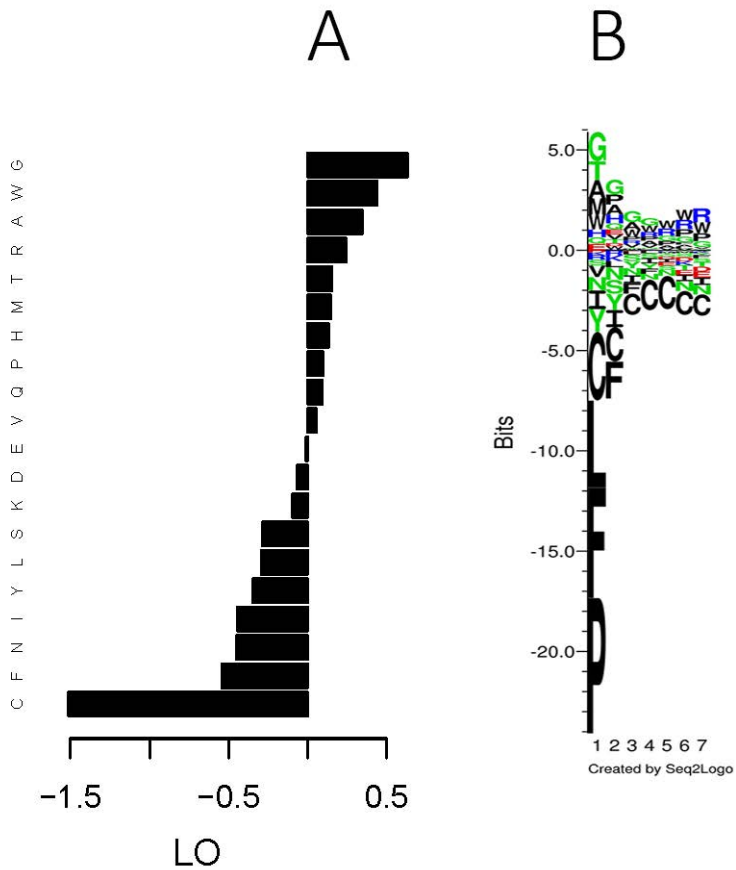
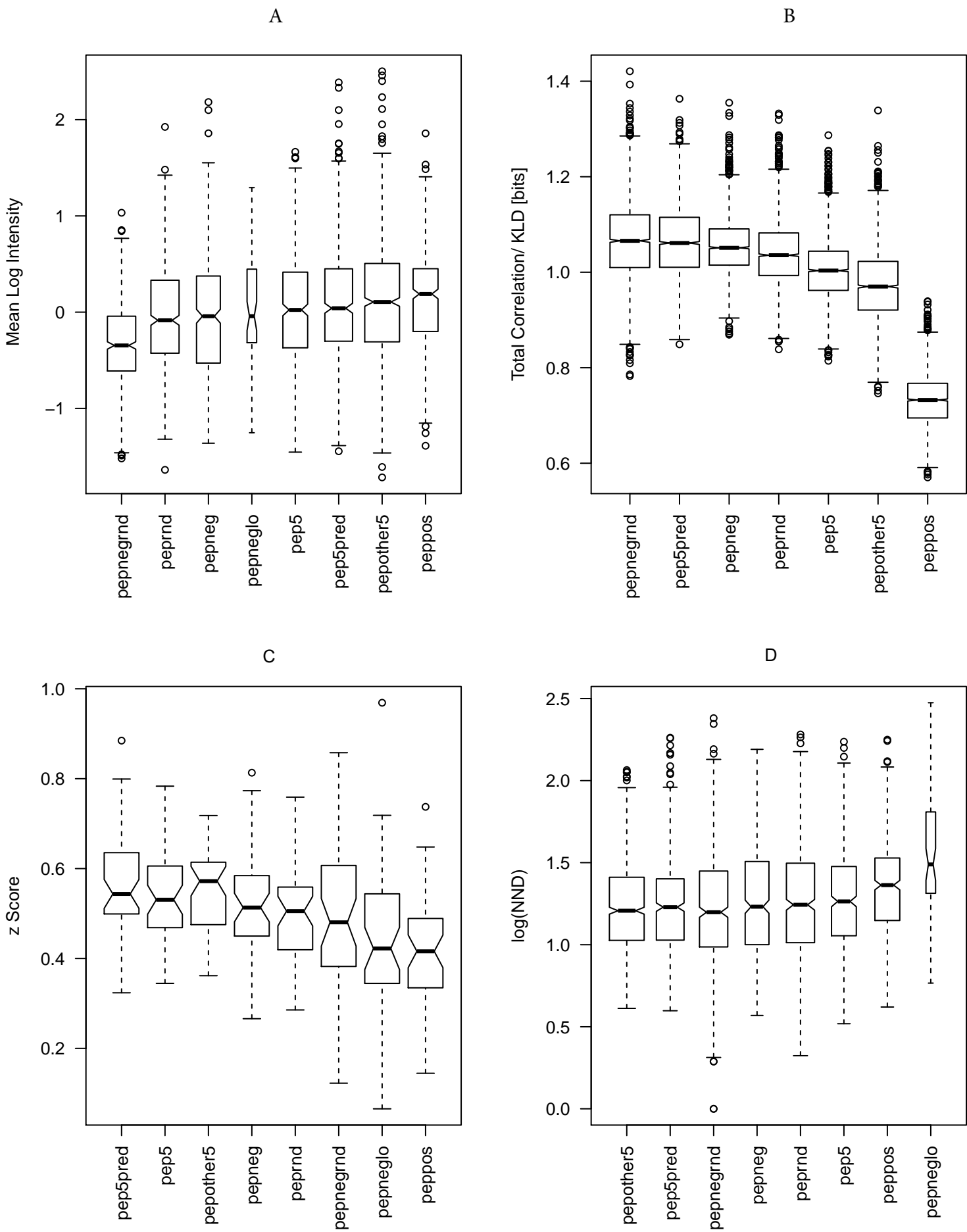


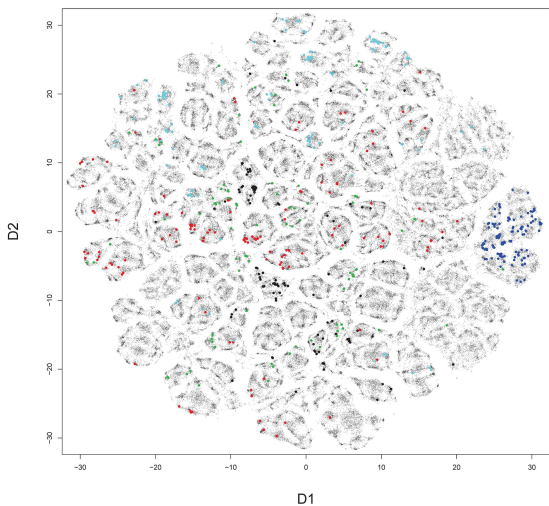
Fig. 4





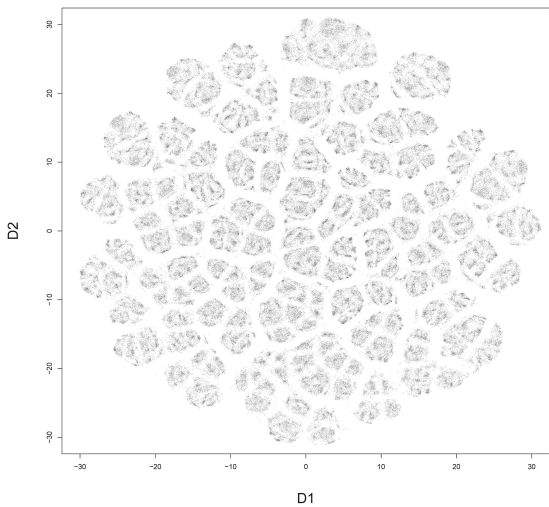
A

Selected Mimotopes with Five Highly Significant Clusters



B

Random Peptides



Selected Mimotopes with the Optimized Library

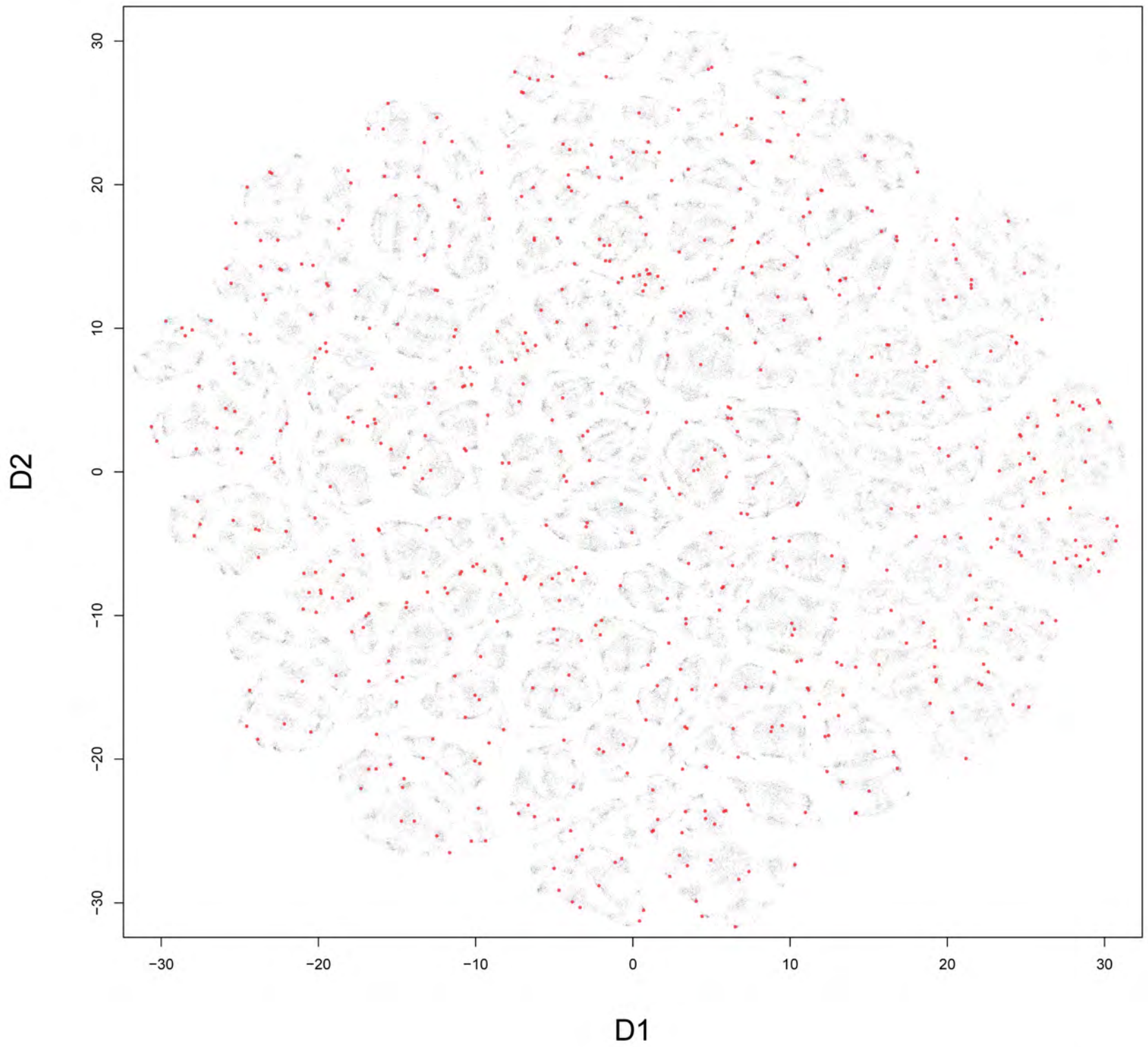


Fig. 7

Mixture of Selected Mimotopes and Random Peptides

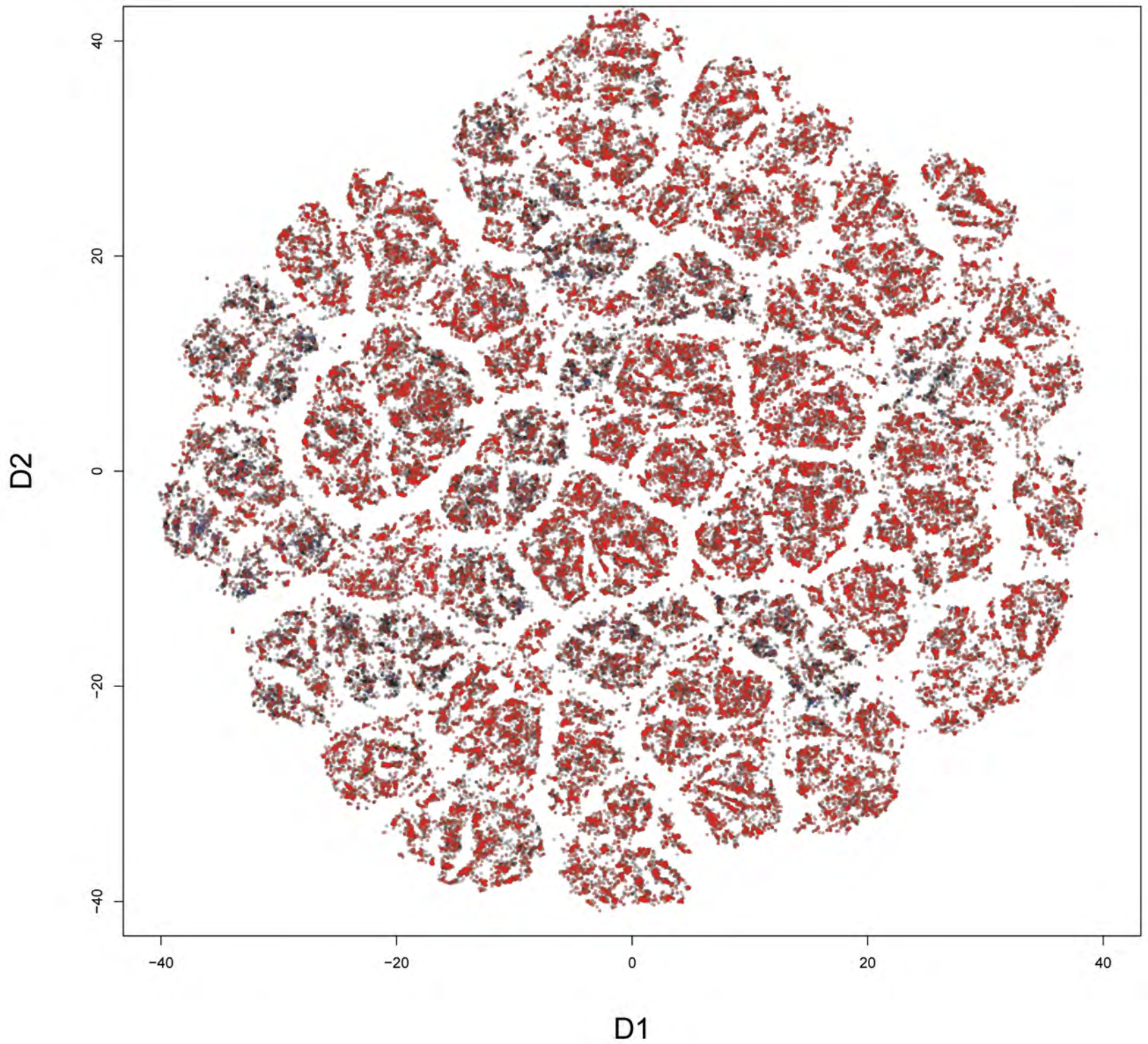
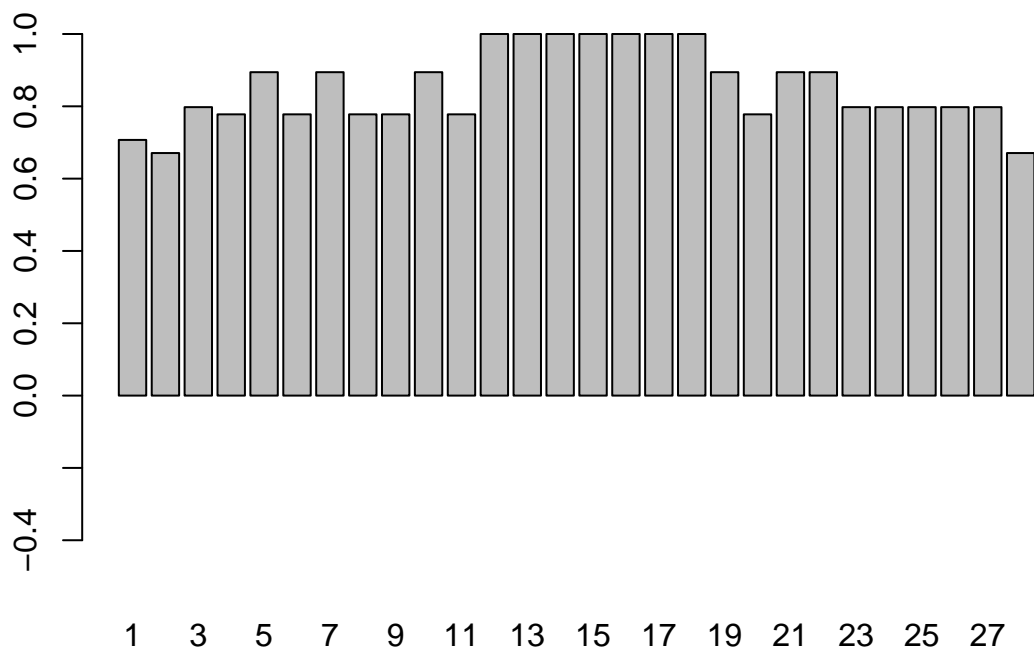
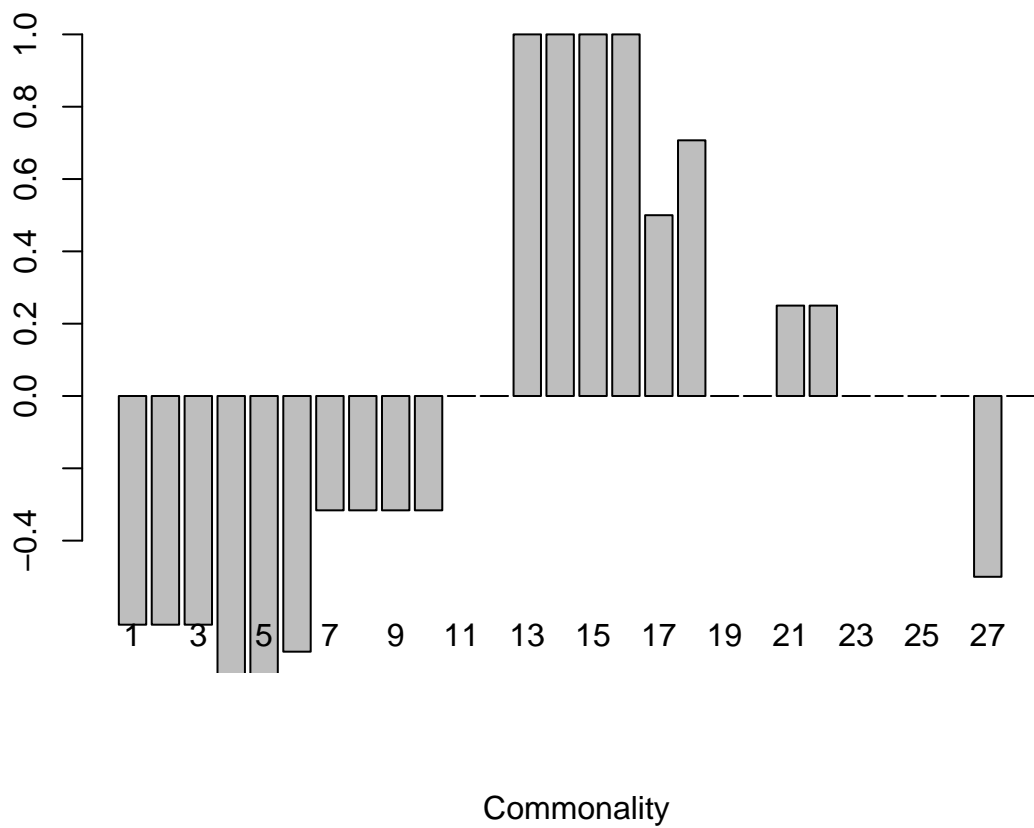


Fig. 8

Training



Validation



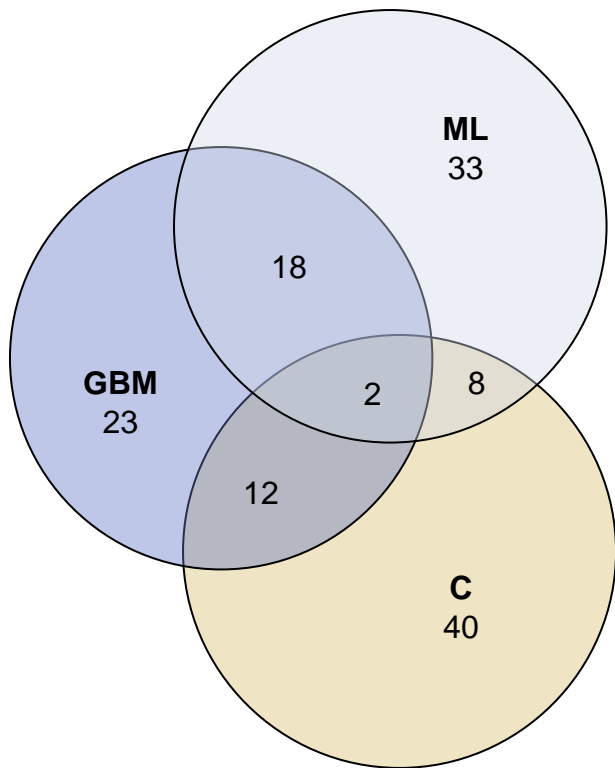


Fig. 10

Fig. 11

