1 **Deep learning model of somatic hypermutation reveals importance of sequence context**

2 **beyond targeting of AID and Polη hotspots**

3

4

5 **Authors and Affiliations**

6 Catherine Tang[†,1], Artem Krantsevich[†,1], and Thomas MacCarthy[1,2,*]

7

8 [†] These authors contributed equally to this work

9 [1]Department of Applied Mathematics and Statistics, Stony Brook University, Stony Brook, NY,

10 11794, USA

11 [2]Laufer Center for Physical and Quantitative Biology, Stony Brook University, Stony Brook,

12 NY, 11794, USA

13 [*]Correspondence: thomas.maccarthy@stonybrook.edu

14

15

16 **Summary**

17

18    B-cells undergo somatic hypermutation (SHM) of the Immunoglobulin (Ig) variable region

19    to generate high-affinity antibodies. SHM relies on the activity of activation-induced deaminase

20    (AID), which mutates C>U preferentially targeting WR<u>C</u> (W=A/T, R=A/G) hotspots. Downstream

21    mutations at W<u>A</u> Polymerase η hotspots contribute further mutations. Computational models of

22    SHM can describe the probability of mutations essential for vaccine responses. Previous studies

23    using short subsequences (*k*-mers) failed to explain divergent mutability for the same *k*-mer. We

24    developed the DeepSHM (Deep learning on SHM) model using *k*-mers of size 5-21, improving

25    accuracy over previous models. Interpretation of DeepSHM identified an extended DWR<u>C</u>T

26    (D=A/G/T) motif with particularly high mutability. Increased mutability was further associated

27    with lower surrounding G content. Our model also discovered a conserved AGYC<u>T</u>GGGGG

28    (Y=C/T) motif within FW1 of IGHV3 family genes with unusually high T>G substitution rates.

29    Thus, a wider sequence context increases predictive power and identifies novel features that drive

30    mutational targeting.

31

32

## Introduction

34

35    Upon encountering antigen, germinal center (GC) B cells undergo several programmed

36    mutational events in secondary lymphoid organs to mount an effective humoral immune response.

37    Somatic hypermutation (SHM) takes place in the GC dark zone whereby mostly point mutations

38    are introduced into the Immunoglobulin (Ig) variable (V) region. Selection for mutations leading

39    to higher binding B cell receptors to cognate antigen occurs in the GC light zone, thus, producing

40    a diverse repertoire of high-affinity antibodies (Methot and Di Noia, 2017; Pilzecker and Jacobs,

41    2019; Rajewsky, 1996). The mutagenic enzyme, activation-induced deaminase (AID), initiates

42    SHM (Muramatsu et al., 2000) by converting cytosine (C) to uracil (U) in single-stranded DNA

43    (ssDNA), resulting in a U:G (guanine) mismatch (Bransteitter et al., 2003). AID displays

44    preferential targeting at WRC/GYW "hotspot" motifs (where W=A/T, R=A/G, Y=C/T, and the

45    underlined base indicates the mutated base in the top and bottom strand, respectively), whereas

46    SYC/GRS "coldspots" (S=C/G) are significantly less targeted (Pham et al., 2003; Rogozin and

47    Diaz, 2004; Rogozin and Kolchanov, 1992; Yu et al., 2004). If left unrecognized, U mismatches

48    will act as a template T and be replicated over (Pilzecker and Jacobs, 2019). The resulting C>T

49    transition mutation is commonly referred to as the DNA "footprint" of AID (Liu et al., 2008).

50    Downstream DNA repair further contributes to antibody diversity that is mediated by low-fidelity

51    polymerases. During non-canonical base-excision repair (ncBER), the U:G mismatch is

52    recognized and excised by uracil-DNA glycosylase (UNG), resulting in an abasic site (Rada et al.,

53    2004). Repair of these abasic sites by REV1 can cause both transition and transversion mutations

54    at C:G base-pairs (Jansen et al., 2006). In the case of non-canonical mismatch repair (ncMMR),

55    the U:G mismatch is recognized by the MSH2/MSH6 heterodimer. Next, EXO1 exonuclease is

56    recruited to create a patch of ssDNA, which then allows error-prone polymerases, particularly

57    Polymerase eta (Polη), to resynthesize. Polη is known to create mutations at neighboring adenine

58    (A) and thymine (T) sites of the initial AID-induced lesion, most notably at WA/TW hotspot motifs

59    (Matsuda et al., 2001; Mayorov et al., 2005).

60        Several computational models have been developed for the SHM process and intrinsic

61    biases exhibited by key proteins such as AID and Polη. These models have mainly utilized *k*-mer

62    subsequences, where $k$ is a specified integer length, ranging between 3- to 7-mers (Cui et al., 2016;

63    Elhanati et al., 2015; Shapiro et al., 1999; Shapiro et al., 2002; Yaari et al., 2013). Two of these

64    models (Cui et al., 2016; Yaari et al., 2013) are widely used and have leveraged 5-mer motifs to

65    capture the dependency of the local surrounding sequence for the middle nucleotide to mutate,

66    while simultaneously bypassing any influence of selection. The first of these targeting models

67    ("S5F") evaluates all possible 5-mers and synonymous (silent) mutations derived from functionally

68    rearranged, or productive, VDJ coding sequences (Yaari et al., 2013). The second model

69    ("RS5NF") similarly assesses 5-mers but uses both synonymous and non-synonymous

70    (replacement) mutations from non-productively (non-functional) rearranged sequences (Cui et al.,

71    2016). Such models have been used to simulate B cell repertoire lineages by constructing a set of

72    hypothetical sequences that have been mutated in a sequential manner as governed by, for example,

73    the underlying S5F substitution scores (Krantsevich et al., 2021; Sheng et al., 2017). Although $k$-

74    mer approaches are generally able to capture some key local intrinsic biases of SHM, such as

75    hotspot targeting, there is evidence that shorter $k$-mers are insufficient to properly characterize

76    differential SHM targeting. For example, a recent study extended a local sequence (5-mer) context

77    model and improved accuracy by including parameters describing the position within the IGHV

78    gene (Spisak et al., 2020). Another study compared the mutability of identical 5-mer (middle

79    position +/-2nt) motifs at different positions within an IGHV gene (Zhou and Kleinstein, 2020),

80    and found that the mutation frequency of these motif-allele pairs (MAPs) positively correlates with

81    the overall mutability of a wider neighborhood of motifs, suggesting that an extended $k$-mer may

82    better capture SHM.

83        Earlier studies have shown that using deep learning is effective in different genomic

84    applications; for example, convolutional neural networks (CNNs) in extracting conserved

4

85    sequence motifs among target sequences (Alipanahi et al., 2015; Kelley et al., 2016; Zhou and

86    Troyanskaya, 2015). In this study, we adopted a deep learning approach using a 2-D CNN to

87    analyze extended *k*-mer lengths to better understand the underlying SHM process. We demonstrate

88    that our model, DeepSHM (Deep learning on SHM), can more accurately represent the SHM

89    process by evaluating longer *k*-mers of up to 21 nts. Additionally, DeepSHM using 15-mers as

90    inputs was able to recapitulate AID WR<u>C</u>/<u>G</u>YW hotspot motifs and identify an extended DWR<u>C</u>T

91    motif. Neural network predictions are notoriously difficult to explain (the "black box" problem),

92    but many new methods are available to interpret results (Koo and Ploenzke, 2020). We used one

93    such method to identify a negative association between increased mutability at a site and its

94    surrounding G content. On the other hand, lower mutation frequency was correlated with increased

95    substitution rates of certain substitution types, particularly for G>T and C>A mutations.

96    Furthermore, many highly conserved sites within G-rich sub-regions belonging to several IGHV3

97    genes display an extremely high bias towards creating G mutations, some of which may participate

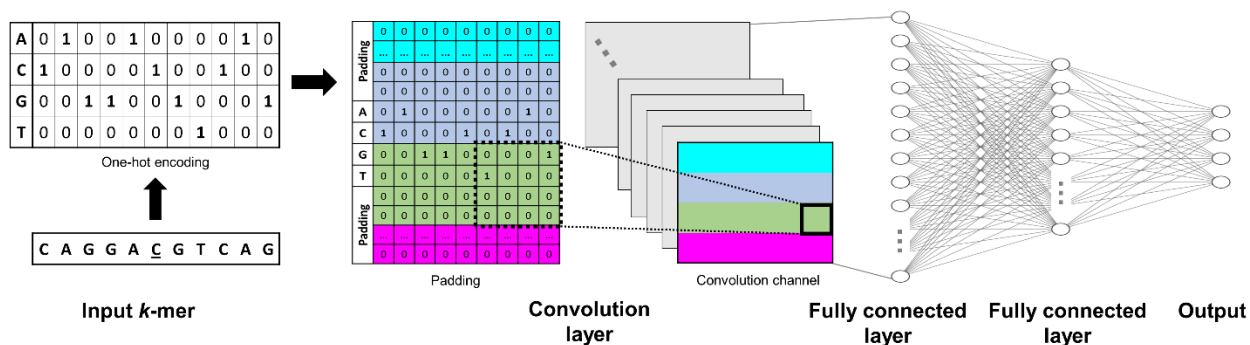98    in the formation of G-quadruplex (G4) structures.

99

100

101    **Results**

102

103    **Deep learning can more accurately represent SHM mutabilities and substitution biases**

104

5

105      The objective of our analysis was to use supervised deep learning to build an accurate

106      convolutional neural network (CNN) for SHM and, as much as possible, identify novel features

107      contributing to mutability. We chose CNNs because we still expected mutation frequency to

108      depend on recurring motifs that might occur at any position in the sequence (most obviously, AID

109      hotspots), a task CNNs are well suited to. The workflow of our network consists of an input layer

110      that processes a $k$-mer subsequence represented in its one-hot encoding format (i.e. a $4 \times k$ matrix

111      of zeros and ones), followed by a convolution layer and two fully connected layers as the hidden

112      layers, and finally the output layer of size $4 \times 1$ or $1 \times 1$, depending on the task that is being predicted

113      (**Figure 1**, see Methods). Several hyperparameters, including dropout rate and learning rate, were

114      fine-tuned with our model as well (**Supplementary Table 1**). We defined a model that would

115      separate mutations on each strand (which are predominantly at C and A on the top strand and at G

116      and T on the bottom strand) at the input level. To achieve this, we identified a simple solution

117      using padding that assigns a row in each channel of the convolution layer output separately to each

118      strand (**Figure 1**). CNNs are also often used together with attribution methods such as Integrated

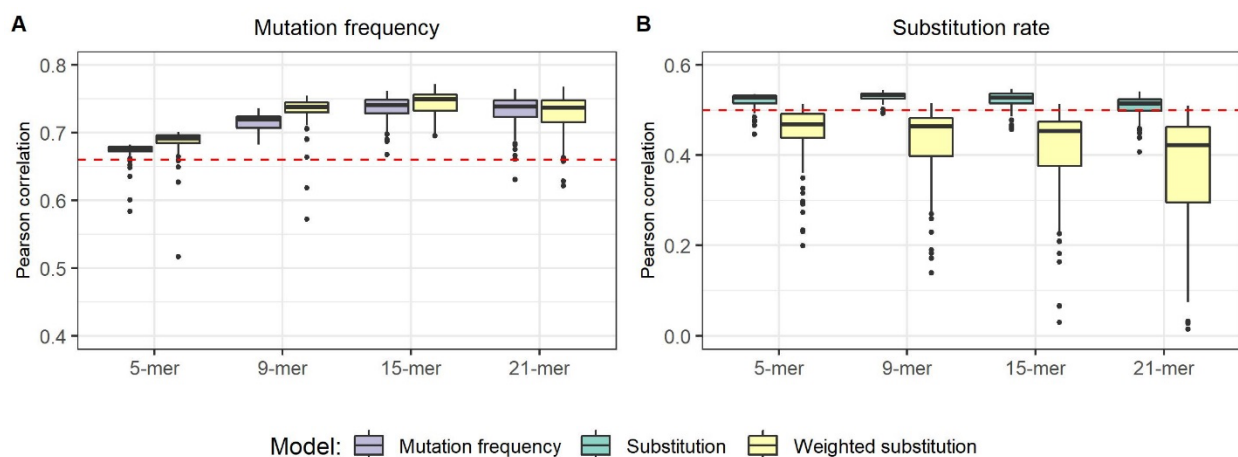119      Gradients, to help with interpretation of the results.



120

121    **Figure 1. DeepSHM model architecture.** Each model had an input layer, one convolution layer, two fully connected

122    layers, and an output layer. The input layer was a $4 \times k$ dimensional one-hot encoded matrix ($k$ is length of

123     subsequence). The dimension of the output layer was dependent on the task: substitution (4×1), mutation frequency

124     (1×1), or weighted substitution (4×1). For the convolutional layer, 'same' padding was used to allow the model to

125     process top and bottom strand mutations separately. With 'same' padding, the output of each convolutional channel

126     has the same shape as the input (4×$k$) with the following properties: the first and the fourth rows are populated with

127     zeros only (there was no real input, only padding; cyan and magenta rows); the input used for the second (light blue)

128     row contained two rows of padding and two data rows corresponding to A or C nucleotides only; and similarly, the

129     input used for the third (green) row also contained two rows of padding and two rows of data corresponding to G or

130     T nucleotides. Since AID and Polη target C and A sites respectively, this approach was taken with the expectation of

131     helping the model distinguish top and bottom strands.

132        As a starting point, we trained two CNN models, which we collectively refer to as

133     DeepSHM (Deep learning on SHM), to separately predict mutation frequency and substitution

134     rates, calculated from previously published B cell repertoire data containing non-productively

135     rearranged and clonally independent VDJ coding sequences (Tang et al., 2020), for varying $k$-mer

136     lengths (see Methods). We trained both models independently using different combinations of $k$-

137     mer lengths and hyperparameters as listed in **Supplementary Table 1**. We found that for

138     predicting mutation frequency, 15-mers were moderately better than 9-mers (purple boxplots in

139     **Figure 2A**, Mann-Whitney U test: P<2.2×10$^{-16}$) and that further extending the motif length to 21

140     did not improve accuracy since both produced an overall maximum correlation (across

141     hyperparameters) of 0.76 (**Figure 2A**, **Table 1)**. Thus, using $k$-mers of length 15 or longer

142     outperformed shorter lengths, specifically 5-mers and 9-mers (**Table 1**), suggesting that an

143     extended DNA motif can better model the SHM process. However, using longer $k$-mers did not

144     substantially improve the model that predicts SHM substitution bias alone, achieving an average

145     correlation of 0.55 for 15-mers (green boxplots in **Figure 2B**, **Table 1**), but which is similar for

146     different lengths. For the interpretability analysis below, we chose to use the best 15-mer models

147     to keep the $k$-mer length consistent for comparisons across all models. In order to check if the

148    performance of the models leading to the best results was consistent, we also trained 30 different

149    iterations of each model, keeping the hyperparameters fixed but using different random seeds. We

150    found the standard deviation across correlations was very small, at 0.002 for the mutation

151    frequency model and 0.001 for the substitution rate model, showing the strong consistency of our

152    results.



153

154    **Figure 2. Performance of DeepSHM.** Boxplots describing the distribution (across random hyperparameters) of

155    Pearson correlations between DeepSHM predictions and empirical data (y-axis) are shown for different input *k*-mers

156    (x-axis) for **(A)** mutation frequencies, and **(B)** substitution rates, for all three models (mutation frequency, substitution,

157    and weighted substitution). Red dashed lines signify correlations of predicted S5F values, which uses 5-mers.

158    We next sought to compare DeepSHM against the widely used S5F model that is based on

159    5-mer motifs (Yaari et al., 2013). To ensure a fair comparison, we generated an S5F targeting

160    model using the same data set that was used to train DeepSHM, as well as the same cross-validation

161    scheme (see Methods). Using the same test set splits as above, we found that there was an average

162    correlation of 0.66 between the predicted S5F model mutabilities and empirical mutation

163    frequency, and an average correlation of 0.50 for predicted S5F substitution scores and empirical

164    substitution rates (red dashed lines in **Figure 2**, **Table 1**). The substitution model slightly (but

165 statistically significantly) outperformed S5F for all *k*-mer values we analyzed. The mutation

166 frequency model achieved a modest improvement over S5F using 5-mers as an input, and this

167 difference became more evident for 9-, 15-, and 21-mers (**Figure 2A**, **Table 1**). We also similarly

168 computed 30 iterations (using different random seeds) of the best 15-mer models for both mutation

169 frequency and substitution models, and found these iterations to have significantly greater

170 accuracy than S5F both individually and in aggregate ($P<1.8\times10^{-6}$ for each model, Wilcoxon

171 signed-rank test). Overall, these results show that our deep learning approach successfully extracts

172 meaningful information from the wider sequence context to improve predictions.

| Model | Test set | S5F (5-mer) | DeepSHM (5-mer) | DeepSHM (9-mer) | DeepSHM (15-mer) | DeepSHM (21-mer) |
|---|---|---|---|---|---|---|
| Substitution rate | IGHV1 | 0.52 | 0.57 | 0.57 | 0.58 | 0.57 |
| | IGHV3 | 0.49 | 0.52 | 0.52 | 0.53 | 0.52 |
| | IGHV4 | 0.48 | 0.53 | 0.53 | 0.54 | 0.52 |
| | IGHV2, 5, 6, 7 | 0.52 | 0.52 | 0.54 | 0.54 | 0.53 |
| | **Avg correlation** | **0.50** | **0.54** | **0.54** | **0.55** | **0.54** |
| | Best - S5F | NA | 0.04 | 0.04 | 0.05 | 0.04 |
| | Mean - S5F | NA | 0.02 | 0.03 | 0.02 | 0.01 |
| | P-value | NA | 1.18E-13 | 1.28E-17 | 2.04E-13 | 2.25E-04 |
| Mutation frequency | IGHV1 | 0.69 | 0.72 | 0.78 | 0.8 | 0.81 |
| | IGHV3 | 0.65 | 0.69 | 0.72 | 0.74 | 0.73 |
| | IGHV4 | 0.64 | 0.69 | 0.73 | 0.78 | 0.78 |
| | IGHV2, 5, 6, 7 | 0.66 | 0.64 | 0.68 | 0.73 | 0.74 |
| | **Avg correlation** | **0.66** | **0.68** | **0.73** | **0.76** | **0.76** |

9

| | | | | | | |
|---|---|---|---|---|---|---|
| | Best - S5F | NA | 0.02 | 0.08 | 0.1 | 0.1 |
| | Mean - S5F | NA | 0.01 | 0.06 | 0.08 | 0.07 |
| | P-value | NA | 6.41E-14 | 1.28E-17 | 1.28E-17 | 1.28E-17 |
| | IGHV1 | 0.52 | 0.55 | 0.55 | 0.53 | 0.53 |
| | IGHV3 | 0.49 | 0.5 | 0.5 | 0.5 | 0.49 |
| | IGHV4 | 0.48 | 0.49 | 0.51 | 0.48 | 0.5 |
| Weighted substitution (substitution rate) | IGHV2, 5, 6, 7 | 0.52 | 0.49 | 0.52 | 0.49 | 0.51 |
| | **Avg correlation** | **0.50** | **0.51** | **0.52** | **0.50** | **0.51** |
| | Best - S5F | NA | 0.04 | 0.09 | 0.11 | 0.11 |
| | Mean - S5F | NA | 0.03 | 0.07 | 0.08 | 0.07 |
| | P-value | NA | 7.51E-16 | 9.47E-17 | 1.28E-17 | 2.33E-17 |
| | IGHV1 | 0.69 | 0.77 | 0.81 | 0.84 | 0.84 |
| | IGHV3 | 0.65 | 0.69 | 0.72 | 0.73 | 0.71 |
| | IGHV4 | 0.64 | 0.68 | 0.7 | 0.74 | 0.73 |
| Weighted substitution (mutation frequency) | IGHV2, 5, 6, 7 | 0.66 | 0.66 | 0.74 | 0.77 | 0.77 |
| | **Avg correlation** | **0.66** | **0.70** | **0.74** | **0.77** | **0.76** |
| | Best - S5F | NA | 0.01 | 0.02 | 0.01 | 0.01 |
| | Mean - S5F | NA | -0.05 | -0.07 | -0.09 | -0.14 |
| | P-value | NA | 2.06E-15 | 1.53E-16 | 3.19E-17 | 1.28E-17 |

173 **Table 1. Cross-validation of various input *k*-mer sequences.** The correlations of repeatedly trained models using

174 different random seeds (but the same hyperparameters) for neural network training had small standard deviations, in

175 all cases below 0.01. P-values are from a Wilcoxon signed-rank test comparing the training results for each model

176 with the corresponding S5F model accuracy. P-values were corrected (Benjamini-Hochberg) for multiple

177 comparisons.

178    To identify associations between mutation frequency and specific substitutions, we further

179    constructed a DeepSHM model to predict the "weighted substitution" of a *k*-mer, i.e., the product

180    of the percentage of each observable substitution type (e.g G>N) and the mutation frequency of

181    the *k*-mer (see Methods). Note that this weighted substitution metric is a vector representing the

182    four ordered DNA bases, with a "0" placed at the position that matches the middle nucleotide of

183    the *k*-mer. Since weighted substitution constitutes aspects of both the observable mutation

184    frequency and substitution rate of the middle nucleotide of a given *k*-mer, we were able to evaluate

185    DeepSHM on each metric separately. Although this model made poorer substitution rate

186    predictions on average (varying hyperparameters) than S5F (**Table 1**), the best model performed

187    similarly to S5F for substitution rates while, surprisingly, performing slightly better than any

188    model in predicting mutation frequency. Cross-validation in this instance produced a range of

189    average correlations between 0.50-0.52 for predicted substitution rates – a level similar to that of

190    S5F (**Figure 2B**, **Table 1**). On the other hand, DeepSHM of weighted substitution values was

191    marginally better at predicting mutation frequency for 15-mers (correlation: 0.77) than the

192    previous standalone model that was tasked to learn mutation frequency only as well as being better

193    than S5F. (**Figure 2A**, **Table 1**). Since the weighted substitution model was able to perform at a

194    level slightly better to the standalone mutation frequency model for longer *k*-mers and substantially

195    better for shorter (5-mer, 9-mer), this suggested a possible association between the projected

196    substitution bias of a site and overall mutability and furthermore, that interpretability methods

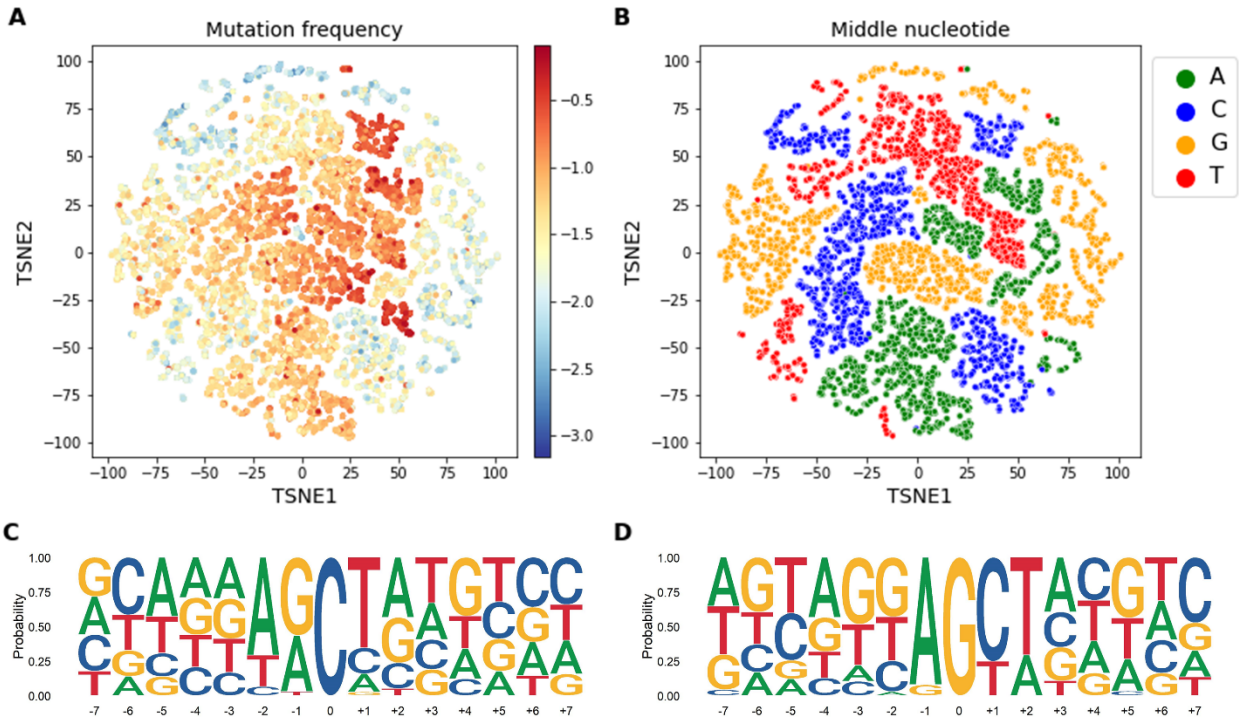197    might uncover these (see below).

198

199    **Interpretation of the DeepSHM network reveals extended hotspot motifs**

200

201       A complication often associated with deep neural networks is model interpretability (the

202     "black box" problem). One way we interrogated the predictions made by DeepSHM, and what it

203     has learned about the SHM process, was to analyze the output of the penultimate layer of each 15-

204     mer based model. In particular, analyzing the output, or "encodings", of this layer can be viewed

205     as an alternative, and more informative, way of representing the input 15-mer. To visualize the

206     multi-dimensional encodings of the individual 15-mers, we used t-SNE, a dimensionality reduction

207     technique, to project each onto a 2-dimensional embedding (see Methods). At this point in order

208     to make full use of the data, we merged all of the 15-mer data into one training set, and then trained

209     three new individual models (one for each output type) using the hyperparameters which

210     previously led to the best cross-validation results. The analyses we present below are derived from

211     the DeepSHM models that were trained using this merged data set.

212       We began by identifying features learned by DeepSHM that predicted weighted

213     substitutions. Since weighted substitution is a measure of both mutation frequency and substitution

214     bias, the embedding should capture both metrics simultaneously. Each point in the resulting t-SNE

215     embedding in **Figure 3A** represents a single 15-mer and is colored according to its corresponding

216     mutation frequency. We identified several clusters of 15-mers that are mostly grouped by similar

217     mutation rates, including those expressing high mutability. Clusters with mid to high mutation

218     frequencies are similarly within close proximity but displayed no obvious groupings other than

219     being mostly located towards the center. When we considered the middle nucleotide of each 15-

220     mer, we observed that these clusters also shared the same middle nucleotide (**Figure 3B**),

221     suggesting that the network identified as a key feature the "0" value in the weighted substitution

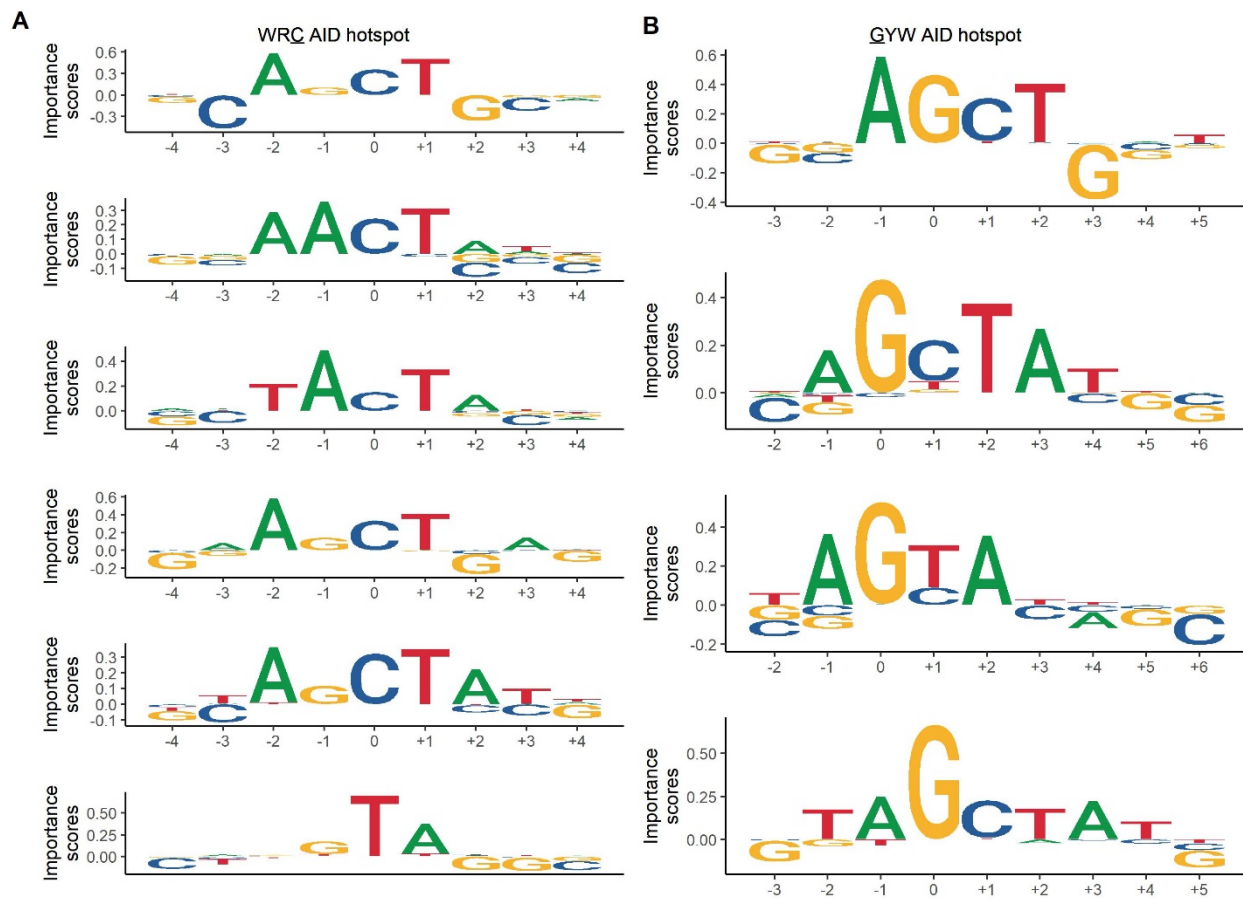222     output vector that is associated with the middle nucleotide.

12

**Figure 3. Neural network encodings analysis: weighted substitution model.** Each point in the t-SNE embedding represents a single 15-mer processed through the truncated model (to extract the output of the penultimate layer) originally trained to learn the associated weighted substitutions (see Methods) and is colored according to its corresponding **(A)** mutation frequency ($\log_{10}$), and **(B)** middle nucleotide. Consensus sequences derived from the highest mutated cluster identified using k-means clustering on the embedding of 15-mers containing either **(C)** a middle C nucleotide or **(D)** a G nucleotide (clusters 10 and 16 in Supplementary Table 2).

Next, we applied k-means clustering on the embedding as a way to isolate cluster boundaries (**Supplementary Figure 1**, see Methods). We subsequently created a sequence logo plots representing each cluster shown in **Supplementary Table 2**. As expected, clusters with the highest mutation frequencies had inner subsequences containing AID (C cluster 10, G cluster 16) and Polη (A clusters 1 and 2, T cluster 20) hotspots. For AID, these are WR<u>C</u> (**Figure 3C**) and <u>G</u>YW motifs (**Figure 3D**). Within the two most highly targeted AID hotspot clusters there is a substantial presence of both WG<u>C</u>/<u>G</u>CW and WA<u>C</u>/<u>G</u>TW contexts, rather than only the well-

13

237    known WGCW overlapping hotspot motif (Tang et al., 2020; Wei et al., 2015). Furthermore, even

238    when we include WAC/GTW, there is a preference for a T base at the 3' end of the WRC hotspot,

239    and conversely, an even stronger bias for an A base at the 5' end of the GYW hotspot (**Figure 3C,**

240    **D**). This motif is consistent with a genome-wide study of AID mutations in mice that reported

241    observing high mutability at AACT and AGCT motifs in both strands (Álvarez-Prado Á et al.,

242    2018). When we assessed the mutability of all possible WRCN (N=A/C/G/T) motifs separately,

243    we observed WRCT to be the most highly mutated in each case (**Supplementary Figure 2**).

244    Previous studies identified WRCY/RGYW (Y=C/T, R=A/G) and later WRCH/DGYW

245    (H=A/C//T, D=A/G/T) to be a better predictor of mutability at C:G bases (Rogozin and Diaz, 2004;

246    Rogozin and Kolchanov, 1992). However, we discovered some inconsistencies with these

247    definitions, as AGCC was found to be the least mutated of the AGCN motifs and WRCG was not

248    always the least mutated, on both strands. Overall, these early hotspot definitions may have been

249    too broad, and WRCT/AGYW is a more consistent predictor of AID targeting. Lastly, we also

250    noted that among the least mutated *k*-mer clusters, many were G-rich in their surrounding context

251    (for example, C clusters 8 and 9, A cluster 4, T cluster 21), and particularly for G (G clusters 12

252    and 14) (**Supplementary Table 2**).

**Figure 4. Recurrent motifs identified by TF-MoDISco.** TF-MoDISco results using the Integrated Gradients as base-level importance scores of 15-mers whose middle nucleotide conformed to a **(A)** WRC or **(B)** GYW AID hotspot motif.

As a complementary way to find sequence motifs associated with mutability, we used TF-MoDISco (Transcription Factor Motif DIScovery), a program for identifying recurring motif patterns in genomic data (see Methods) (Shrikumar et al., 2018). We applied TF-MoDISco to the standalone model that predicts only mutation frequency because we reasoned that sequence features related to mutability would be more easily identifiable since the model is only required to learn a single task. TF-MoDISco uses importance scores, which can be derived from many machine learning methods, to produce a set of unique motifs learned by the model (see Methods).
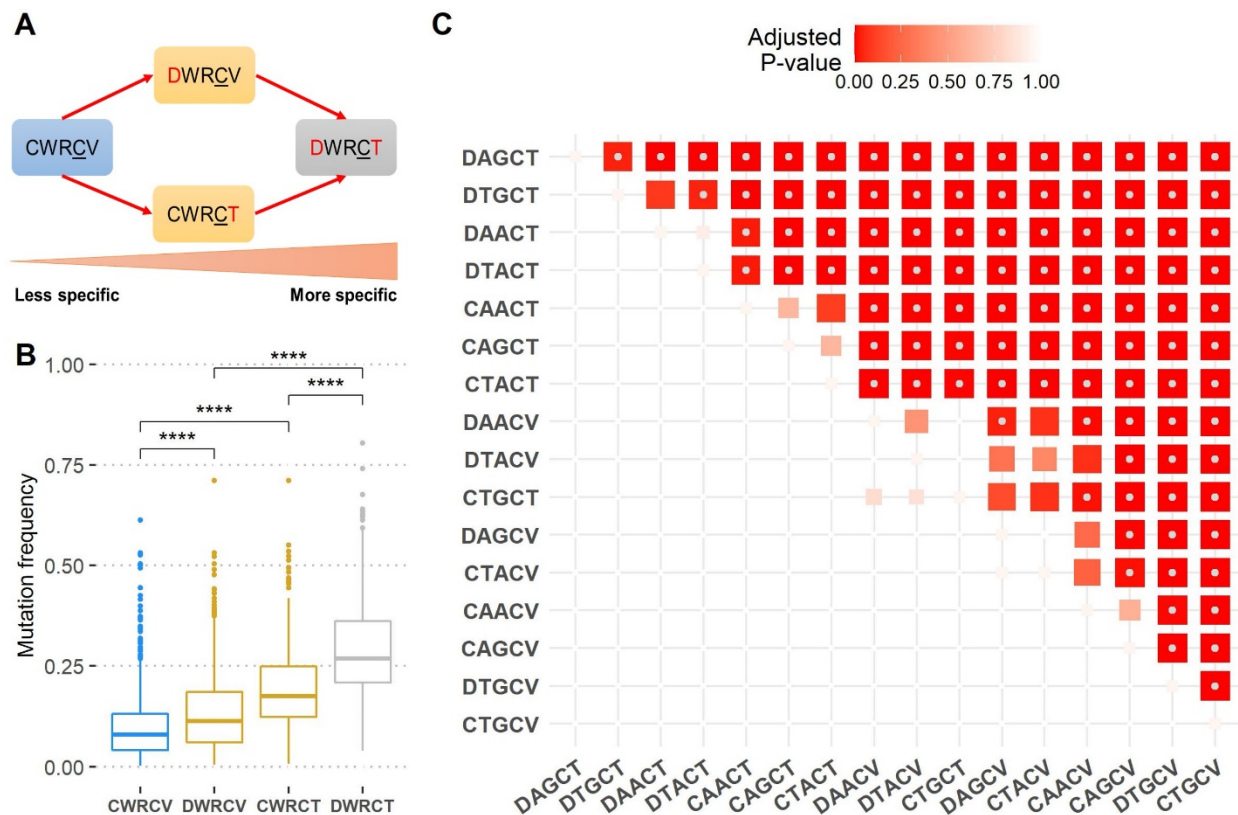
15

264    We began by analyzing the importance scores derived from Integrated Gradients (Sundararajan et

265    al., 2017) of 15-mers whose middle nucleotide conformed a WRC/GYW AID hotspot motif. As

266    expected, the positively contributing sites in the set of ensuing motifs aligned with the hotspot

267    motifs (**Figure 4A, B**). In addition, TF-MoDISco again revealed a preference for having a T base

268    at the +1 position of the WRC (WRCT, **Figure 4A**) and an A base at the -1 position of the GYW

269    (AGYW, **Figure 4B**).

270        In addition to WRCT/AGYW being a well-represented motif identified by TF-MoDISco,

271    as measured by having positive contributions to mutability (above horizontal axis on **Figure 4A**),

272    we also noticed many neighboring C and G bases contained negative contributions (below

273    horizontal axis on **Figure 4A**), most evidently at the C located at the -3 position of the WRC

274    hotspot, and the G located at the +3 position of the GYW hotspot (**Figure 4B**). Here, the negative

275    contribution at the -3 position signifies that having a C at that position reduces mutational targeting

276    to the middle C. By the same token, a mutation that changes the -3 position from C will increase

277    the likelihood of the middle C subsequently being targeted. This observation supports our recently

278    published study where we observed a strong positive "mutual association" – a correlation metric

279    describing the impact of mutating one site and its effect at another site – between CC (or GG) pairs

280    distanced by two nucleotides (Krantsevich et al., 2021). In that study we were able to explain most

281    of such correlations in terms of overlapping AID and/or Polη hotspots, with the CNNC/GNNG

282    motif being one the exceptions which we suggested might be explained by AID processivity (Pham

283    et al., 2003; Storb et al., 2009). However, the TF-MoDISco analysis suggests a different

284    explanation in which the absence of a C in the -3 position might be a part of an extended AID

285    hotspot, defining CWRC as being similar to a sequential overlap motif, which we previously

286    defined (Krantsevich et al., 2021) as a motif in which an initial mutation creates a new hotspot that

16

287    previously did not exist. Here, although the WRC hotspot did previously exist, a mutation in the

288    first C would create a DWRC (D=A/G/T) motif, potentially with higher mutability.

289    We next sought to determine whether adding the 5' D or 3' T context of the canonical WRC

290    hotspot is more influential in terms of increasing its susceptibility to AID mutagenesis. To address

291    this, we increased the hotspot specificity step by step, starting from CWRCV (V=A/C/G) and

292    assessed the impact a single change in the motif at either the first C or V site, causing a DWRCV

293    or CWRCT intermediate hotspot to form respectively, has on mutability (**Figure 5A**). We found

294    that both DWRCV and CWRCT intermediate hotspots were shown to mutate significantly more

295    than CWRCV (**Figure 5B**). We also discovered that the mutability of the DWRCT hotspot, which

296    contains the extended hotspot in both 5' and 3' directions, was significantly higher than both

297    intermediate hotspots (**Figure 5B**). Performing a pairwise comparison between the mutation

298    frequency of all 16 individual (D/C)WRC(T/V) contexts further confirmed that those containing

299    both a 5' D and 3' T were significantly more mutated than the remaining hotspot motifs, with

300    DAGCT being the most mutated (**Figure 5C, Supplementary Table 3**). Additionally, the next

301    three successively mutated hotspots followed a CWRCT context, overall suggesting the 3' T to be

302    more impactful to AID recognition than the 5' D, but the addition of both substantially increases
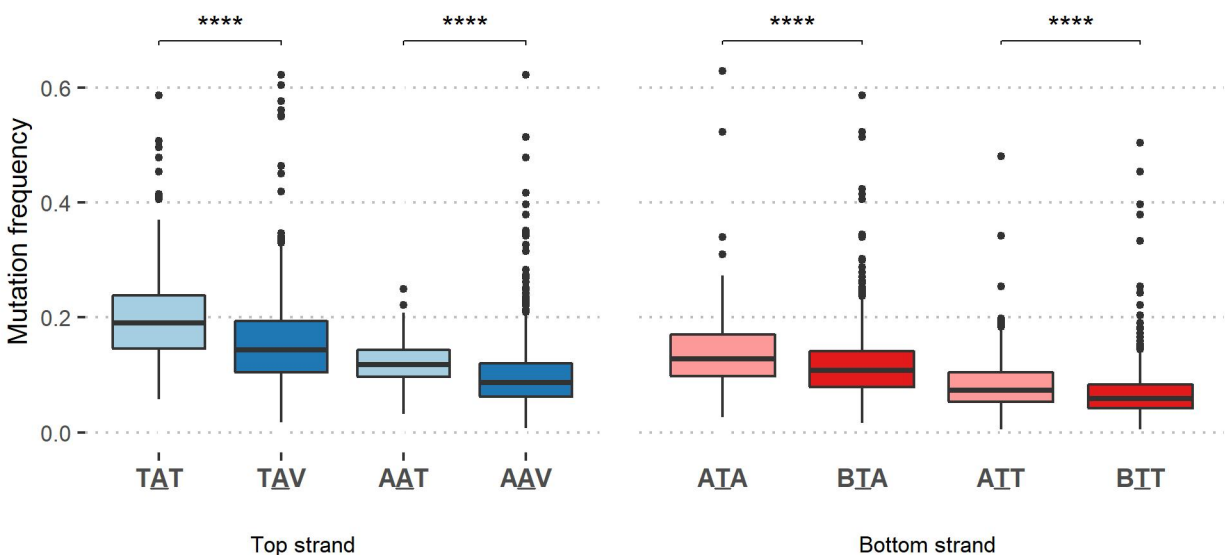
303    targeting in human V regions.

304

**Figure 5. Mutability of extended AID hotspots. (A)** Schematic showing an increase of AID hotspot specificity (left to right). **(B)** Boxplots displaying the mutability of different (C/D)WRC(T/V) hotspot contexts, where D=A/G/T, V=A/C/G. Asterisks indicate significance (p ≤ 0.0001) of a one-sided Mann-Whitney U test comparing the greater mutation frequency of the boxplot on the right against the one on the left. **(C)** Pairwise comparison of mutability for all 16 (C/D)WRC(T/V) hotspot contexts. Boxes represent the p-value - adjusted for multiple comparisons (Benjamini-Hochberg correction) - of a one-sided Mann-Whitney U test comparing the greater mutation frequency of the hotspot indicated by the row to the left, against the hotspot shown in the column below. Rows and columns are ordered by mean mutation frequency (high to low). The color and size of each box is scaled according to the adjusted p-value. Gray dots inside boxes indicate p-values ≤ 0.05.

314         In addition, another secondary motif that unexpectedly emerged from the TF-MoDISco

315     analysis of WR<u>C</u>/<u>G</u>YW 15-mers did not contain a positively contributing C nucleotide; rather it

316     conformed to a <u>TA</u> Polη hotspot (**Figure 4A**, bottom). Having a <u>TA</u> hotspot appear while

317    specifically analyzing only 15-mers containing WR<u>C</u> hotspots reveals the importance of attracting

318    Polη to these areas.  This finding is consistent with our previous analysis highlighting the

319    importance of co-localization of A<u>GC</u>T overlapping AID hotspots and Polη hotspots within the

320    CDRs (Tang et al., 2020; Wei et al., 2015).

321        The <u>TA</u> motif also emerged when we applied TF-MoDISco to all 15-mers conforming to

322    either a W<u>A</u> (**Supplementary Figure 3A**) or <u>T</u>W Polη hotspot (**Supplementary Figure 3B**). In

323    addition to our model identifying both the <u>TA</u> and A<u>A</u> hotspot motifs as important, it also identified

324    a T<u>A</u>T/A<u>T</u>A motif as a special case for both strands. Further analysis showed that W<u>A</u>T/A<u>T</u>W

325    hotspots mutate significantly more than their W<u>A</u>V/B<u>T</u>W counterparts (**Figure 6**). Thus, while T<u>A</u>

326    hotspots consistently have higher mutability than A<u>A</u>, the presence of a 3' T individually increases

327    the mutability of each of these Polη hotspots.



328

**Figure 6. Mutability of extended Polη hotspot motifs.** Boxplots comparing the mutation frequency of various top

strand W<u>A</u>T against W<u>A</u>V (blue), and bottom strand A<u>T</u>W against B<u>T</u>W (red) motifs. Asterisks indicate significance

19

331     (p ≤ 0.0001) of a one-sided Mann-Whitney U test comparing the greater mutation frequency of the boxplot on the left
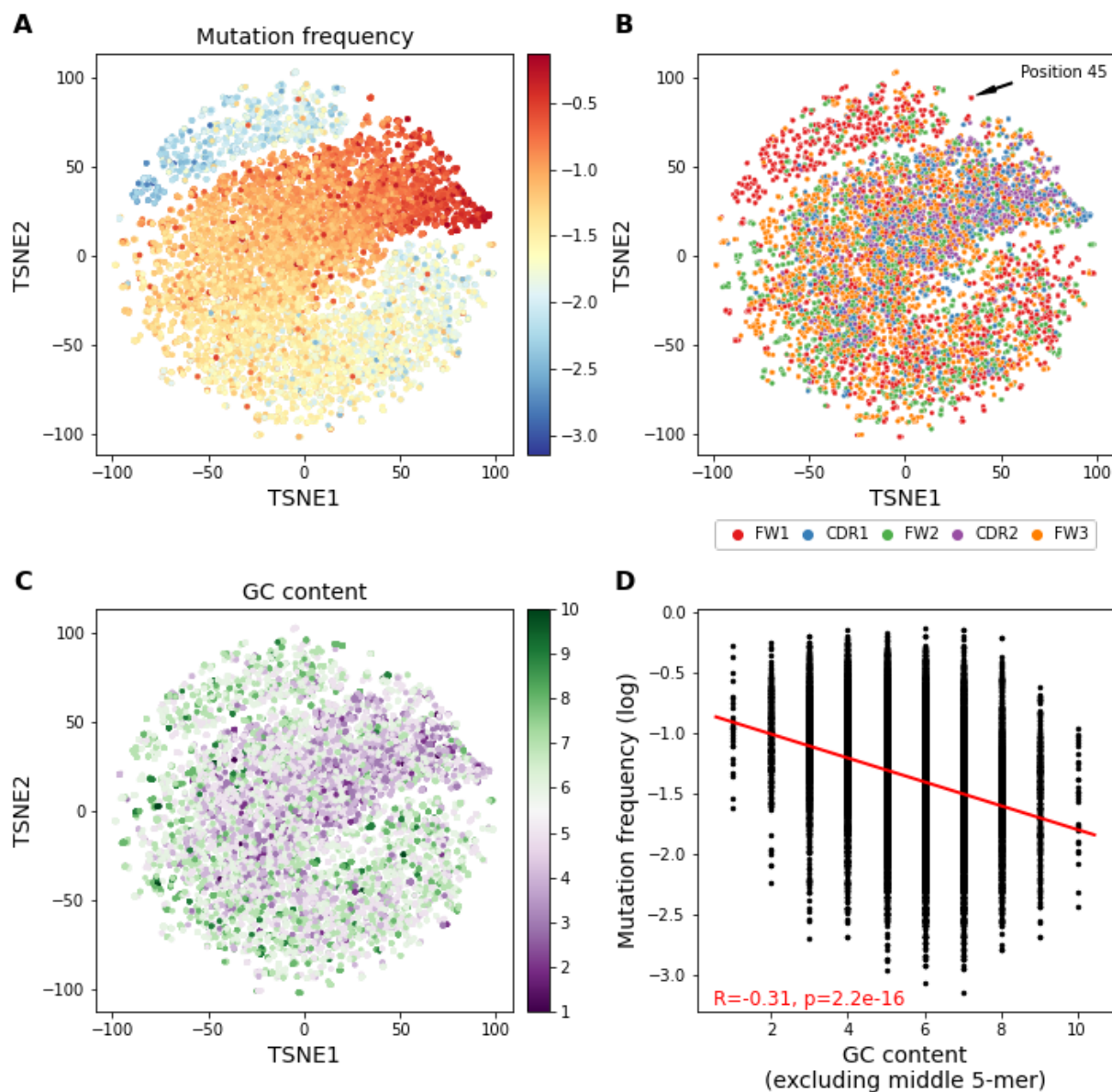
332     against the one on the right.

333

334     **Highly targeted sites display a lower surrounding GC content**

335

336     We next applied the same t-SNE methodology to the DeepSHM model that predicted only

337     mutation frequency. We found that the organization of the subsequent embedding followed a

338     direction of descending mutation frequency, with the highest mutating 15-mers located at the mid-

339     to upper-right portion of the plot (**Figure 7A**). A cluster of low-mutating 15-mers was also isolated

340     in the upper-left (**Figure 7A**), which was enriched with ~76% of FW1 15-mers (**Figure 7B**).

341     Additionally, we examined the possible influence of the local surrounding sequence by calculating

342     the individual base content of the four DNA bases in each 15-mer. However, the inner 5-mer,

343     which contains the dominant context, was excluded when computing all base counts. When we

344     colored the t-SNE embedding according to the GC content of each 15-mer, we observed that GC

345     content increases along the same direction as decreasing mutability seen previously (**Figure 7C**).

346     Quantifying this observation more formally, we indeed found a significant negative correlation

347     between the GC content and the mutation frequency of the 15-mers (R=-0.31, P<$2.2\times10^{-16}$; **Figure**

348     **7D**). On the other hand, when we considered each individual base count independently, we

349     observed that the count of G nucleotides specifically shows a stronger negative correlation (R=-

350     0.19) than the C nucleotide count (R=-0.084) alone (**Supplementary Figure 4A**), although both

351     correlations are highly significant (P<$2.2\times10^{-16}$). This result is consistent with the cluster analysis

352     above (**Supplementary Table 2**) where we observed several clusters with G-rich *k*-mers and low

353    mutation frequencies. If we further separate the mutation frequencies into categories defined by

354    the middle nucleotide, we find that G content has a consistent negative correlation regardless

355    (column G of **Supplementary Figure 4B**). More generally, A and T richness (columns A and T

356    of **Supplementary Figure 4B**) shows a consistent positive or sometimes non-significant

357    correlation, whereas C and G richness shows a consistent negative (or non-significant) correlation.

358    In summary, it appears that low-mutating sites generally have a high local GC (and particularly G)

359    content, and conversely, that highly targeted sites display an elevated local AT (particularly A)

360    content.

**Figure 7. Neural network encodings analysis: mutation frequency model.** Each point in the t-SNE embedding represents a single 15-mer processed through the truncated model (to extract the output of the penultimate layer) trained on mutation frequencies (see Methods) and is colored according to its corresponding **(A)** mutation frequency ($\log_{10}$), and **(B)** by Ig V sub-region location as defined by IMGT. **(C)** The t-SNE embedding is colored according to the GC content of each 15-mer. The calculated GC content excludes the middle 5-mer context of the 15-mer to remove any confounding AID hotspot or coldspot bias. **(D)** Computed Pearson correlation between mutation frequency and GC content, again excluding the middle 5-mer.
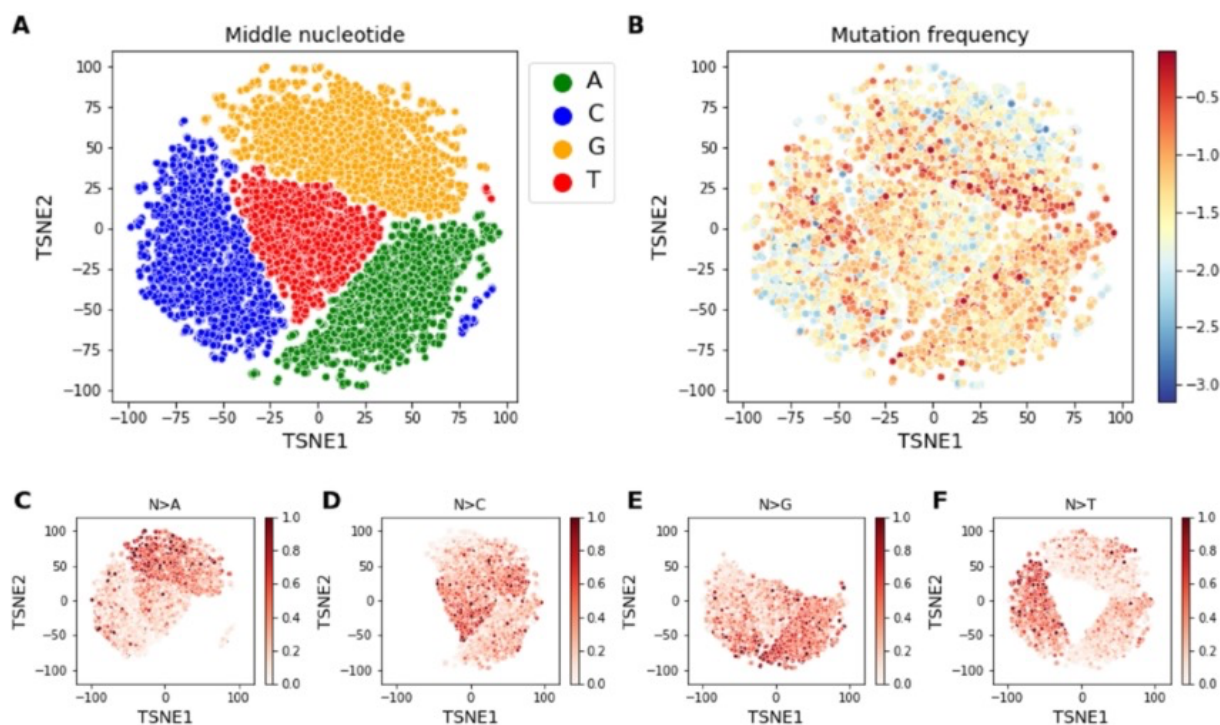
369

**Conserved FW1 sites surrounded by clusters of AID coldspots in IGHV3 genes display a**

**high T>G transversion bias**

372

We now analyzed the standalone model predicting only substitution rates to gain possible insight into additional substitution biases exhibited by AID or downstream error-prone DNA damage response pathways, for example, as a result of REV1 or Polη intervention during non-canonical base-excision repair (BER) and non-canonical mismatch repair (MMR), respectively. The resulting t-SNE embedding from this model identified four main clusters, as well as two much smaller satellite clusters, with each cluster containing 15-mers that share a common middle nucleotide (**Figure 8A**). A distinction between 15-mers with high and low mutation frequencies could also be observed based on their location on opposite ends of the cluster, especially for clusters containing either a C or G middle nucleotide, with high-mutating 15-mers typically located on the side closest to the center (**Figure 8B**). Since the model was tasked with learning the distributed substitution rates of each 15-mer, we next sought to evaluate the embedding by the rate of each individual substitution type (e.g. C>T). In certain clusters, a similar gradient of high to low substitution rates could also be seen as we observed for mutation frequency (**Figure 8C-F**). For instance, we noticed the rate of G>T substitutions increasing from the side nearest to the origin towards the outer boundaries of the cluster (top-right cluster in **Figure 8F**), which was associated with a shift towards decreasing mutation frequencies in the same cluster while proceeding in the same direction (**Figure 8B**). To evaluate this trend more closely, we analyzed three human IGHV genes from different families for which we had the most data (IGHV1-18, IGHV3-23, IGHV4-

391  34), so as to include sites with low mutation frequencies at high coverage, and calculated the

392  correlation between mutation frequency and rate of substitution for each substitution type. As an

393  example, for IGHV3-23 we found the most significant negative correlations to be at C>A

394  mutations (R=-0.33, p=0.0058), and the reverse, G>T (R=-0.24, p=0.022; **Supplementary Figure**

395  **5**). Alternatively, we observed a significant positive correlation between mutation frequency and

396  C>T transition mutations (R=0.29, p=0.018; **Supplementary Figure 5**). Similar patterns were also

397  observed for IGHV1-18*01 and IGHV4-34*01 (**Supplementary Figure 5**). These results are

398  consistent with replication bypass (predominantly causing C>T) being favored over BER at sites

399  with high mutation frequency.
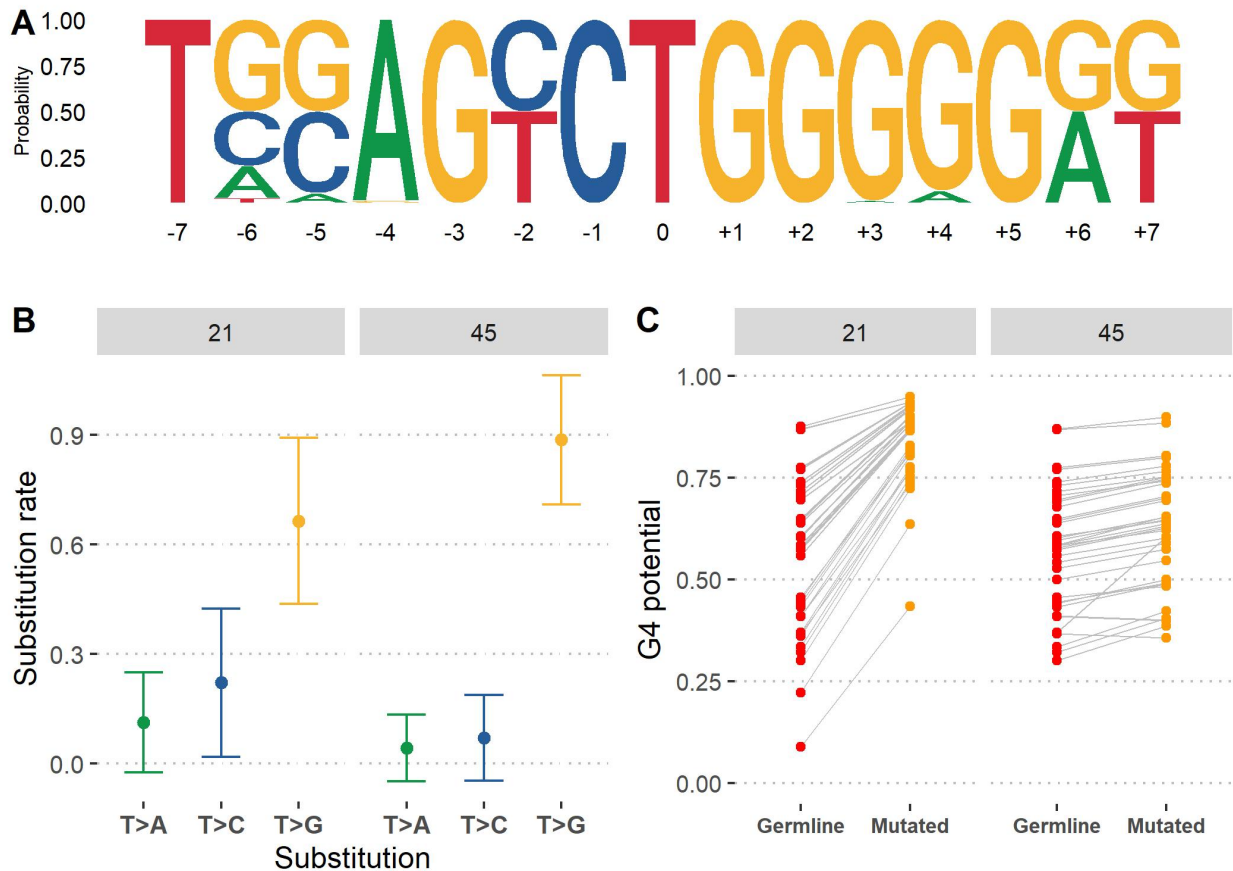


400

401  **Figure 8. Neural network encodings analysis: substitution model.** Each point in the t-SNE embedding represents

402  a single 15-mer processed through the truncated model (to extract the output of the penultimate layer) trained to learn

403  the associated substitution rates (see Methods) and is colored according to its corresponding **(A)** middle nucleotide,

24

404    and **(B)** mutation frequency ($\log_{10}$). **(C-F)** The t-SNE embedding is colored by the rate of substitution for the middle

405    nucleotide of every 15-mer to mutate to A (N>A); to C (N>C); to G (N>G); and to T (N>T), respectively.

406        In the t-SNE analysis of the substitution model, we also discovered two small clusters of

407    15-mers containing a C and T as their middle nucleotide (**Figure 8A**) that did not group with their

408    respective larger clusters, suggesting that these particular sites might have distinct substitution

409    patterns. Generating the consensus sequence of the outlier T cluster revealed a partially conserved

410    AGYC<u>T</u>GGGGG sequence (**Figure 9A**). When we examined these subsequences more closely,

411    we discovered that they were located only in IGHV3 family genes at either position 21 or position

412    45 according to the IMGT unique numbering system (Lefranc, 2001) (**Supplementary Figure 6**).

413    The motif was also surprisingly common. At position 21 it appeared in 37 different alleles (across

414    19 genes) and was fully conserved in all alleles. Coincidentally, the motif also appeared in 37

415    different alleles (across 18 genes) at position 45, although it differed slightly at the +3 and +4

416    positions (**Figure 9A**). These two sets of alleles only partially overlap, such that 15 alleles had the

417    motif at both positions 21 and 45. Thus, this specific motif in FW1 of the IGHV3 family genes

418    appears to be highly conserved evolutionarily, suggesting a possible functional role. The rates of

419    substitution at these sites were also found to be highly biased towards creating T>G mutations,

420    with an average T>G rate of about 0.66 at position 21, and an even greater rate of 0.89 at position

421    45 (background rate: $0.28 \pm 0.23$) (**Figure 9B, Table 2**). A previous study using Sanger sequencing

422    data that was limited to IGHV3-23 and the pseudogene IGHV3-h had noted similarly high T>G

423    substitution rates at positions 21 (for IGHV3-h) and 45 (for IGHV3-23) (Ohm-Laursen and

424    Barington, 2007). Although the T subjected to mutation at both positions did not conform to a

425    bottom strand <u>T</u>W Polη hotspot, these genes at position 45 displayed a relatively high average

426    mutation frequency of $0.17 \pm 0.08$ (**Table 2**), which is somewhat unusual given that mutations are

427 generally more biased towards the CDRs than FW regions (Cohen et al., 2011; Shapiro et al.,

428 2002), and that we reported above that many sites within FW1 tended to display low mutability

429 (**Figure 7A, B**).



430

**Figure 9. Evaluation of the T outlier cluster in the DeepSHM substitution model. (A)** Sequence logo
representation of the 15-mers appearing in the T outlier cluster in **Figure 8A** (right-hand side, red dots). **(B)**
Substitution rates of T>A, T>C, and T>G for 15-mers corresponding to 37 IGHV3 alleles separately at IMGT positions
21 and 45. Bars represent ±1 standard deviation. **(C)** G-quadruplex (G4) formation potential for the same IGHV3
alleles in **(B)**. G4 potentials (y-axis) are computed using the germline IGHV sequence ("Germline") and the mutated
sequence ("Mutated") containing a single simulated T>G mutation at either IMGT positions 21 or 45.

437 While examining the C outlier cluster (**Figure 8A**), we found the consensus sequence to

438 be more diverse compared to the outlier with a middle T (**Supplementary Figure 7A**). The

26

439    sequence variation seen here was partly due to the fact that the 15-mers that constituted this cluster

440    belonged to many other IGHV families besides IGHV3 and across different sub-regions of the IgV

441    (**Supplementary Figure 6**). On the other hand, we noticed some overlap between both outlier

442    clusters since, in some cases, the C corresponded to positions 20 and 44 that preceded the middle

443    T of the other outlier cluster (**Supplementary Figure 7A, Table 2**). We further found these sites

444    to have a similar elevated C>G substitution rate (mean rate of 0.62 compared to background mean

445    of 0.33, $P<2.2\times10^{-16}$) (**Supplementary Figure 7B, Table 2**), suggesting the model distinguished

446    sites with a general preference to create G mutations.

447        Given that the sites with strikingly high T>G and C>G substitution rates we identified here

448    are in adjoining G-rich sub-regions (**Figure 9A, Supplementary Figure 7A**), we evaluated the

449    possible influence these mutations might have on the formation of G-quadruplex (G4) structures.

450    In a recent study, we assessed the potential for DNA G4 structures to form in the IgV region, using

451    a pre-trained deep learning model that computes the G4 potential of a linear DNA sequence (Tang

452    and MacCarthy, 2021). There we found that the IGHV3 family had the highest propensity to form

453    stable G4s in the top strand. We now sought to assess the overall mutational effect on G4 assembly

454    of the IGHV3 sites that are biased towards G. Following the methodology of our previous study,

455    we calculated the difference between the predicted G4 potential of the germline with that of the

456    sequence with a single mutation at either position 21 or 45. Here, we found that a T>G mutation

457    at position 21 elevated average G4 potentials to a very high value of $0.84 \pm 0.10$ compared to a

458    germline value (already relatively high) of $0.54 \pm 0.19$ , whereas the same mutation occurring at

459    position 45 displayed a far smaller average increase of $0.05 \pm 0.04$ (**Figure 9C, Table 2**). As for

460    the remaining cases, there seemed to be little effect of C>G mutations on G4 potential (**Table 2**).

461   Interestingly, we made another observation regarding the instances where an A nucleotide disrupts

462   the run of G

| IMGT position | 15-mer middle nucleotide | n | Avg. substitution rate to G | Avg. mutation frequency | Avg. germline G4 potential | Avg. mutated G4 potential | Avg. difference in G4 potential (mutated - germline) |
|---|---|---|---|---|---|---|---|
| 1 | C | 1 | 0.67 | 0.02 | 0.01 | 0.01 | 0.00 |
| 20 | C | 31 | 0.50 ± 0.32 | 0.01 ± 0.01 | 0.54 ± 0.20 | 0.64 ± 0.19 | 0.09 ± 0.02 |
| 21 | T | 37 | 0.66 ± 0.23 | 0.05 ± 0.03 | 0.54 ± 0.19 | 0.84 ± 0.10 | 0.29 ± 0.10 |
| 34 | C | 19 | 0.60 ± 0.37 | 0.03 ± 0.05 | 0.07 ± 0.08 | 0.07 ± 0.08 | 0.01 ± 0.00 |
| 44 | C | 38 | 0.77 ± 0.24 | 0.02 ± 0.01 | 0.56 ± 0.17 | 0.64 ± 0.18 | 0.07 ± 0.04 |
| 45 | T | 37 | 0.89 ± 0.18 | 0.17 ± 0.08 | 0.58 ± 0.15 | 0.63 ± 0.14 | 0.05 ± 0.04 |
| 48 | A | 1 | 0.79 | 0.11 | 0.37 | 0.56 | 0.19 |
| 49 | A | 5 | 0.79 ± 0.12 | 0.07 ± 0.03 | 0.37 ± 0.07 | 0.53 ± 0.05 | 0.16 ± 0.02 |
| 61 | C | 45 | 0.67 ± 0.14 | 0.04 ± 0.02 | 0.54 ± 0.19 | 0.54 ± 0.19 | 0.01 ± 0.01 |
| 101 | C | 1 | 0.52 | 0.11 | 0.01 | 0.01 | 0.00 |
| 167 | C | 1 | 0.46 | 0.52 | 0.37 | 0.41 | 0.04 |
| 173 | C | 3 | 0.64 ± 0.29 | 0.09 ± 0.09 | 0.04 ± 0.02 | 0.04 ± 0.02 | 0.00 ± 0.00 |
| 180 | C | 1 | 0.20 | 0.03 | 0.01 | 0.01 | 0.00 |
| 214 | C | 4 | 0.51 ± 0.26 | 0.13 ± 0.06 | 0.08 ± 0.06 | 0.10 ± 0.07 | 0.01 ± 0.01 |
| 249 | C | 15 | 0.59 ± 0.24 | 0.06 ± 0.02 | 0.08 ± 0.09 | 0.08 ± 0.09 | 0.00 ± 0.01 |
| 268 | C | 44 | 0.55 ± 0.19 | 0.06 ± 0.03 | 0.55 ± 0.18 | 0.53 ± 0.18 | -0.01 ± 0.01 |

463   **Table 2. Summary statistics on outlier C and T clusters from Figure 8A.**

464         nucleotides at the +3 or +4 positions (**Figure 9A**) which was that these sites also displayed

465   high A>G substitution rates (0.79 ± 0.11; **Table 2,** positions 48 and 49). This hypothetical mutation

466    also caused a moderate, though substantial, increase in G4 potential ($0.17 \pm 0.02$, **Table 2**). These

467    findings reveal that particular recurring mutations in this sub-region may promote G4 formation,

468    and that the bias towards generating new G sites suggests specific DNA repair enzymes may be

469    recruited to these sub-regions within FW1.

470

471

472    **Discussion**

473

474        In this study, we leveraged deep learning to gain novel insights into SHM, a key process

475    in antibody affinity maturation. We trained multiple deep learning models using a convolutional

476    neural network (CNN) framework to analyze DNA $k$-mer subsequences of various lengths, ranging

477    from 5 to 21 nts, derived from human IGHV germline sequences. Using a high-quality data set

478    containing non-productive B cell repertoire data, the model was tasked to learn two focal aspects

479    of SHM: the frequency of mutation at a given site, and the spectrum of mutations that can arise at

480    this site (substitution). Understanding the propensity of a site to mutate and the underlying

481    substitution biases that ensue can lead to a better understanding of how AID is recruited to and

482    targets the Ig V region, as well as the associated downstream DNA repair mechanisms that follow

483    AID deamination.

484        We began by developing three models, collectively referred to as DeepSHM, to predict

485    separate tasks for a given $k$-mer: observable mutation frequency; distributed substitution rates; and

486    a combination of both measures (weighted substitution). We found that predicting substitution

487    rates did not substantially depend on the *k*-mer size, while 15-mers were optimal for predicting

488    mutation frequencies (**Figure 2, Table 1**). Additionally, DeepSHM predicted both substitution

489    rates and mutation frequencies more accurately than the widely used S5F targeting model for all

490    *k*-mer sizes we evaluated (*k* = 5, 9, 15 and 21) (**Table 1**). Even though we were able to outperform

491    S5F in representing substitution biases, the correlation between our predictions and empirical data

492    was moderate (~0.55), suggesting that the processes underlying SHM substitution biases may be

493    more fundamentally random than mutational site targeting alone. Error-prone DNA repair

494    processes downstream of AID are highly complex. For example, while Polη is biased towards

495    making W<u>A</u>>W<u>G</u> mutations (Zhang et al., 2014) and plays a dominant role in generating mutations

496    at A:T sites, many A:T mutations still occur in its absence (Saribasak et al., 2009) that are mediated

497    by other polymerases (Maul et al., 2016). Similarly complex, BER is biased towards transversions

498    but can also repair faithfully, with a further dependence on hotspot mutability (Pérez-Durán et al.,

499    2012). Thus, downstream repair processes may simply be too complex, or genuinely random, to

500    be captured well by a model that depends on sequence context alone.

501    In order to uncover some of the hidden features learned by DeepSHM, we analyzed the

502    output, or encodings, obtained from the penultimate layer of the network predicting weighted

503    substitution using input 15-mers, and performed t-SNE, a method of dimensionality reduction, to

504    visualize the encodings in two dimensions. The subsequent embedding formed clusters of 15-mers

505    that were distinguished by mutation frequency and middle nucleotide (**Figure 3A, B**). Individual

506    clusters containing a C or G middle nucleotide that were associated with high mutability, assumed

507    to be relevant to AID hotspots, revealed a strong preference for a T base at the +1 position of the

508    top strand AID WR<u>C</u> (W=A/T, R=A/G) hotspot, including for WA<u>C</u> motifs that are not part of a

509    WG<u>C</u>W motif, and similarly, an A base at the -1 position of the bottom strand <u>G</u>YW (Y=C/T)

510    context (**Figure 3C, D**). As an alternative way to identify sequence features, we applied TF-

511    MoDISco (see Methods) to reveal recurrent genomic patterns using importance scores extracted

512    from the model for each 15-mer. This approach confirmed the importance of the T base at the +1

513    position of WRC (**Figure 4A**) and the A base at the -1 position of the bottom strand GYW hotspot

514    (**Figure 4B**). An early study by Rogozin and Diaz reported the WRCH/DGYW (H=A/C/T,

515    D=A/G/T) to be a good predictor of mutability at C:G bases (Rogozin and Diaz, 2004), but we

516    found WRCT to be a more consistent definition. The authors of the S5F model also supported the

517    WRCH definition since they found their model can capture the higher mutability rate seen at

518    certain WRCA motifs (Yaari et al., 2013), presumably at the AGCA overlapping hotspot.

519    However, previous hotspot definitions have largely failed to describe targeting beyond the -2

520    position of the WRC motif. We further identified having a C at the -3 position of WRC or a G at

521    the +3 position of GYW as a strong negative contribution, i.e., as a reduced effect on targeting.

522    Thus, our results suggest the typical AID hotspot definition might be extended to DWRCT

523    (D=A/G/T). Comparing the mutation frequencies of the individual DWRCT hotspot motifs

524    showed the 3' T to be more important for AID recognition than the 5' D alone, however, together

525    they have a synergistic effect that makes mutability between 1.8-fold (for TAC) and 4.7-fold (for

526    TGC) higher (**Figure 5C, Supplementary Table 3**).

527         We next applied the same t-SNE methodology on the two developed standalone models

528    that separately predicted either the mutation frequency or substitution rates of the 15-mer middle

529    nucleotides. The t-SNE embedding on the independent DeepSHM model predicting only mutation

530    frequency revealed a significant negative correlation between the mutability of a site and the

531    surrounding GC content of the 15-mer (**Figure 7D**). This finding alternatively suggests that highly

532    mutated sites may have evolved to have a richer local AT content. This *in vivo* result is consistent

31

533    with earlier *in vitro* results that considered AID targeting on artificial substrates (Abdouni et al.,

534    2018).

535            On the other hand, the t-SNE embedding stemming from the standalone substitution model

536    hinted at plausible associations, both positive and negative, between mutation frequency and

537    certain transition and transversion mutations (**Figure 8B-F**). We next analyzed multiple genes

538    representing different IGHV families containing the largest amounts of mutation data in order to

539    avoid any potential sites with few observable mutations, such as coldspots. We observed a negative

540    correlation between mutability and substitution rates specifically for C>A and G>T transversion

541    mutations (**Figure 8, Supplementary Figure 6**) and, on the other hand, positive correlations for

542    C>T and G>A transitions (**Supplementary Figure 6**). The trend for increased transition mutations

543    at highly mutating AID hotspots mediated by UNG2 had previously been observed in experiments

544    using 3T3 (mouse fibroblast) cells (Pérez-Durán et al., 2012), although the particular bias against

545    C:G>A:T transversions was not apparent. Previous work has also shown that UNG2 is cell-cycle

546    regulated, possibly mediated by FAM72A (Feng et al., 2020), and active primarily during G1

547    (Sharbeen et al., 2012). Although AID is also primarily active during G1, it may sometimes persist

548    for slightly longer than UNG2 and thus highly targeted sites may avoid BER especially when the

549    mutations occur just before the cell enters S phase, which would lead to fixation of C>T transitions

550    via replication bypass. Alternative polymerases may also be preferentially recruited to some sites.

551    For example, in DT40 (chicken) B-cell lines, the POLD3 subunit of Polymerase delta (Polδ) has

552    been proposed as a specific mechanism for both C>A and G>T mutations (Hirota et al., 2015;

553    Pilzecker and Jacobs, 2019).

554            Additionally, we investigated two outlier clusters from the substitution model embedding

555    that contained 15-mers having a C and T middle nucleotide that did not group with their respective

556     larger clusters (**Figure 8A**). A closer analysis revealed that the T outlier contained a highly

557     conserved AGYC<u>T</u>GGGGG consensus sequence that was derived from two independent sites

558     located in FW1 from multiple IGHV3 alleles (**Figure 9A, Table 2**). Both outlier clusters also

559     displayed significantly elevated T>G (**Figure 9B, Table 2**) and C>G substitution rates

560     (**Supplementary Figure 7, Table 2**) respectively. In our recent study on G-quadruplexes (G4s)

561     in IGHV genes, we observed the IGHV3 family to form G4s more favorably on the top strand, as

562     measured by their predicted G4 potential using a pre-trained CNN model (Tang and MacCarthy,

563     2021). Given the strong preference for creating G mutations in these FW1 sub-regions, we

564     evaluated the impact of these mutations on G4 potential. In some cases, the resulting G mutation

565     led to a strong increase in G4 potential, particularly at position 21 (**Figure 9C, Table 2**), whereas

566     for other sites, the effect was mostly negligible (**Table 2**). Notably however, a high A>G

567     substitution rate was also observed at the +3 or +4 positions (**Figure 9A**), which were also

568     associated with increase in G4 potential (**Table 2**). These biased A>G mutations may further be

569     related to previous work that found that a repeated mutation that occurs in one IGHV allele often

570     matches the sequence variant of a different allele (Saini and Hershberg, 2015). Alternatively, these

571     mutations may be related to R-loop initiation, which forms in G-rich non-template DNA, possibly

572     forming in FW1 of these IGHV3 genes. Studies have found that reducing G-density in mammalian

573     Ig Switch regions compromises class-switch recombination efficiency and R-loops from forming

574     (Roy et al., 2008; Zhang et al., 2014). The high rate of T>G and C>G transversions also suggests

575     that particular repair enzymes may be recruited to these sub-regions during SHM.


576


577     **Limitations of the study**

578

579 In principle, a wider range of k-mers, as well as a greater variety of neural network architectures,

580 might have been considered for this study. However, since the tuning of each model takes a

581 substantial amount of computational resources and time, we considered a reduced number of

582 models. Additionally, we limited this study to consider data only for human, the species for

583 which we had high quality (UMI barcoded) data in high abundance, although the approach could

584 be extended to other species such as mouse in future work.

585

586 **Acknowledgements**

587 We would like to extend our gratitude to Sergio Roa Gómez for his comments and suggestions.

588

589 **Author contributions**

590 C.T., A.K., and T.M. conceived the idea, analyzed the results, and wrote the manuscript. C.T. and

591 A.K. developed the model and performed computational analysis. All authors contributed to the

592 article and approved the submitted version.

593

594 **Declaration of interests**

595 The authors declare no competing interests.

**References**

596

597

598 Abdouni, H.S., King, J.J., Ghorbani, A., Fifield, H., Berghuis, L., and Larijani, M. (2018).
599 DNA/RNA hybrid substrates modulate the catalytic activity of purified AID. Molecular
600 Immunology *93*, 94-106.

601 Alipanahi, B., Delong, A., Weirauch, M.T., and Frey, B.J. (2015). Predicting the sequence
602 specificities of DNA- and RNA-binding proteins by deep learning. Nat Biotechnol *33*, 831-838.

603 Álvarez-Prado Á, F., Pérez-Durán, P., Pérez-García, A., Benguria, A., Torroja, C., de Yébenes,
604 V.G., and Ramiro, A.R. (2018). A broad atlas of somatic hypermutation allows prediction of
605 activation-induced deaminase targets. J Exp Med *215*, 761-771.

606 Bransteitter, R., Pham, P., Scharff, M.D., and Goodman, M.F. (2003). Activation-induced cytidine
607 deaminase deaminates deoxycytidine on single-stranded DNA but requires the action of RNase.
608 Proc Natl Acad Sci U S A *100*, 4102-4107.

609 Cohen, R.M., Kleinstein, S.H., and Louzoun, Y. (2011). Somatic hypermutation targeting is
610 influenced by location within the immunoglobulin V region. Mol Immunol *48*, 1477-1483.

611 Cui, A., Di Niro, R., Vander Heiden, J.A., Briggs, A.W., Adams, K., Gilbert, T., O'Connor, K.C.,
612 Vigneault, F., Shlomchik, M.J., and Kleinstein, S.H. (2016). A Model of Somatic Hypermutation
613 Targeting in Mice Based on High-Throughput Ig Sequencing Data. J Immunol *197*, 3566-3574.

614 Elhanati, Y., Sethna, Z., Marcou, Q., Callan, C.G., Jr., Mora, T., and Walczak, A.M. (2015).
615 Inferring processes underlying B-cell repertoire diversity. Philos Trans R Soc Lond B Biol Sci
616 *370*.

617 Feng, Y., Li, C., Stewart, J., Barbulescu, P., Desivo, N.S., Álvarez-Quilón, A., Pezo, R.C., Perera,
618 M.L.W., Chan, K., Tong, A.H.Y.*, et al.* (2020). FAM72A antagonizes UNG2 to promote
619 mutagenic uracil repair during antibody maturation. bioRxiv, 2020.2012.2023.423975.

620 Hirota, K., Yoshikiyo, K., Guilbaud, G., Tsurimoto, T., Murai, J., Tsuda, M., Phillips, L.G., Narita,
621 T., Nishihara, K., Kobayashi, K.*, et al.* (2015). The POLD3 subunit of DNA polymerase δ can
622 promote translesion synthesis independently of DNA polymerase ζ. Nucleic acids research *43*,
623 1671-1683.

624 Jansen, J.G., Langerak, P., Tsaalbi-Shtylik, A., van den Berk, P., Jacobs, H., and de Wind, N.
625 (2006). Strand-biased defect in C/G transversions in hypermutating immunoglobulin genes in
626 Rev1-deficient mice. J Exp Med *203*, 319-323.

627 Kelley, D.R., Snoek, J., and Rinn, J.L. (2016). Basset: learning the regulatory code of the
628 accessible genome with deep convolutional neural networks. Genome Res *26*, 990-999.

629 Koo, P.K., and Ploenzke, M. (2020). Deep learning for inferring transcription factor binding sites.
630 Curr Opin Syst Biol *19*, 16-23.

631  Krantsevich, A., Tang, C., and MacCarthy, T. (2021). Correlations in Somatic Hypermutation
632  Between Sites in IGHV Genes Can Be Explained by Interactions Between AID and/or Polη
633  Hotspots. Frontiers in Immunology *11*.

634  Lefranc, M.P. (2001). IMGT, the international ImMunoGeneTics database. Nucleic Acids Res *29*,
635  207-209.

636  Liu, M., Duke, J.L., Richter, D.J., Vinuesa, C.G., Goodnow, C.C., Kleinstein, S.H., and Schatz,
637  D.G. (2008). Two levels of protection for the B cell genome during somatic hypermutation. Nature
638  *451*, 841-845.

639  Matsuda, T., Bebenek, K., Masutani, C., Rogozin, I.B., Hanaoka, F., and Kunkel, T.A. (2001).
640  Error rate and specificity of human and murine DNA polymerase eta. J Mol Biol *312*, 335-346.

641  Maul, R.W., MacCarthy, T., Frank, E.G., Donigan, K.A., McLenigan, M.P., Yang, W., Saribasak,
642  H., Huston, D.E., Lange, S.S., Woodgate, R.*, et al.* (2016). DNA polymerase iota functions in the
643  generation of tandem mutations during somatic hypermutation of antibody genes. J Exp Med *213*,
644  1675-1683.

645  Mayorov, V.I., Rogozin, I.B., Adkison, L.R., and Gearhart, P.J. (2005). DNA polymerase eta
646  contributes to strand bias of mutations of A versus T in immunoglobulin genes. J Immunol *174*,
647  7781-7786.

648  Methot, S.P., and Di Noia, J.M. (2017). Molecular Mechanisms of Somatic Hypermutation and
649  Class Switch Recombination. Adv Immunol *133*, 37-87.

650  Muramatsu, M., Kinoshita, K., Fagarasan, S., Yamada, S., Shinkai, Y., and Honjo, T. (2000). Class
651  switch recombination and hypermutation require activation-induced cytidine deaminase (AID), a
652  potential RNA editing enzyme. Cell *102*, 553-563.

653  Ohm-Laursen, L., and Barington, T. (2007). Analysis of 6912 unselected somatic hypermutations
654  in human VDJ rearrangements reveals lack of strand specificity and correlation between phase II
655  substitution rates and distance to the nearest 3' activation-induced cytidine deaminase target. J
656  Immunol *178*, 4322-4334.

657  Pérez-Durán, P., Belver, L., de Yébenes, V.G., Delgado, P., Pisano, D.G., and Ramiro, A.R.
658  (2012). UNG shapes the specificity of AID-induced somatic hypermutation. J Exp Med *209*, 1379-
659  1389.

660  Pham, P., Bransteitter, R., Petruska, J., and Goodman, M.F. (2003). Processive AID-catalysed
661  cytosine deamination on single-stranded DNA simulates somatic hypermutation. Nature *424*, 103-
662  107.

663  Pilzecker, B., and Jacobs, H. (2019). Mutating for Good: DNA Damage Responses During Somatic
664  Hypermutation. Front Immunol *10*, 438.

665 Rada, C., Di Noia, J.M., and Neuberger, M.S. (2004). Mismatch recognition and uracil excision
666 provide complementary paths to both Ig switching and the A/T-focused phase of somatic mutation.
667 Mol Cell *16*, 163-171.

668 Rajewsky, K. (1996). Clonal selection and learning in the antibody system. Nature *381*, 751-758.

669 Rogozin, I.B., and Diaz, M. (2004). Cutting edge: DGYW/WRCH is a better predictor of
670 mutability at G:C bases in Ig hypermutation than the widely accepted RGYW/WRCY motif and
671 probably reflects a two-step activation-induced cytidine deaminase-triggered process. J Immunol
672 *172*, 3382-3384.

673 Rogozin, I.B., and Kolchanov, N.A. (1992). Somatic hypermutagenesis in immunoglobulin genes.
674 II. Influence of neighbouring base sequences on mutagenesis. Biochim Biophys Acta *1171*, 11-18.

675 Roy, D., Yu, K., and Lieber, M.R. (2008). Mechanism of R-loop formation at immunoglobulin
676 class switch sequences. Mol Cell Biol *28*, 50-60.

677 Saini, J., and Hershberg, U. (2015). B cell variable genes have evolved their codon usage to focus
678 the targeted patterns of somatic mutation on the complementarity determining regions. Mol
679 Immunol *65*, 157-167.

680 Saribasak, H., Rajagopal, D., Maul, R.W., and Gearhart, P.J. (2009). Hijacked DNA repair proteins
681 and unchained DNA polymerases. Philos Trans R Soc Lond B Biol Sci *364*, 605-611.

682 Shapiro, G.S., Aviszus, K., Ikle, D., and Wysocki, L.J. (1999). Predicting regional mutability in
683 antibody V genes based solely on di- and trinucleotide sequence composition. J Immunol *163*,
684 259-268.

685 Shapiro, G.S., Aviszus, K., Murphy, J., and Wysocki, L.J. (2002). Evolution of Ig DNA sequence
686 to target specific base positions within codons for somatic hypermutation. J Immunol *168*, 2302-
687 2306.

688 Sharbeen, G., Yee, C.W., Smith, A.L., and Jolly, C.J. (2012). Ectopic restriction of DNA repair
689 reveals that UNG2 excises AID-induced uracils predominantly or exclusively during G1 phase. J
690 Exp Med *209*, 965-974.

691 Sheng, Z., Schramm, C.A., Kong, R., Program, N.C.S., Mullikin, J.C., Mascola, J.R., Kwong,
692 P.D., and Shapiro, L. (2017). Gene-Specific Substitution Profiles Describe the Types and
693 Frequencies of Amino Acid Changes during Antibody Somatic Hypermutation. Front Immunol *8*,
694 537.

695 Shrikumar, A., Tian, K., Shcherbina, A., Avsec, Z., Banerjee, A., Sharmin, M., Nair, S., and
696 Kundaje, A. (2018). Technical Note on Transcription Factor Motif Discovery from Importance
697 Scores (TF-MoDISco). bioRxiv.

698 Spisak, N., Walczak, A.M., and Mora, T. (2020). Learning the heterogeneous hypermutation
699 landscape of immunoglobulins from high-throughput repertoire data. Nucleic Acids Res *48*,
700 10702-10712.

701 Storb, U., Shen, H.M., and Nicolae, D. (2009). Somatic hypermutation: processivity of the cytosine
702 deaminase AID and error-free repair of the resulting uracils. Cell Cycle *8*, 3097-3101.

703 Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic Attribution for Deep Networks. ArXiv
704 *abs/1703.01365*.

705 Tang, C., Bagnara, D., Chiorazzi, N., Scharff, M.D., and MacCarthy, T. (2020). AID Overlapping
706 and Poleta Hotspots Are Key Features of Evolutionary Variation Within the Human Antibody
707 Heavy Chain (IGHV) Genes. Front Immunol *11*, 788.

708 Tang, C., and MacCarthy, T. (2021). Characterization of DNA G-Quadruplex Structures in Human
709 Immunoglobulin Heavy Variable (IGHV) Genes. Frontiers in Immunology *12*.

710 Wei, L., Chahwan, R., Wang, S., Wang, X., Pham, P.T., Goodman, M.F., Bergman, A., Scharff,
711 M.D., and MacCarthy, T. (2015). Overlapping hotspots in CDRs are critical sites for V region
712 diversification. Proc Natl Acad Sci U S A *112*, E728-737.

713 Yaari, G., Vander Heiden, J.A., Uduman, M., Gadala-Maria, D., Gupta, N., Stern, J.N., O'Connor,
714 K.C., Hafler, D.A., Laserson, U., Vigneault, F*., et al.* (2013). Models of somatic hypermutation
715 targeting and substitution based on synonymous mutations from high-throughput immunoglobulin
716 sequencing data. Front Immunol *4*, 358.

717 Yu, K., Huang, F.T., and Lieber, M.R. (2004). DNA substrate length and surrounding sequence
718 affect the activation-induced deaminase activity at cytidine. J Biol Chem *279*, 6496-6500.

719 Zhang, Z.Z., Pannunzio, N.R., Hsieh, C.L., Yu, K., and Lieber, M.R. (2014). The role of G-density
720 in switch region repeats for immunoglobulin class switch recombination. Nucleic Acids Res *42*,
721 13186-13193.

722 Zhou, J., and Troyanskaya, O.G. (2015). Predicting effects of noncoding variants with deep
723 learning-based sequence model. Nat Methods *12*, 931-934.

724 Zhou, J.Q., and Kleinstein, S.H. (2020). Position-Dependent Differential Targeting of Somatic
725 Hypermutation. The Journal of Immunology, ji2000496.
726

727

728    **STAR Methods**

729

730    **Resource availability**

731

732    **Lead contact**

733    Further information and requests for resources and reagents should be directed to and will be

734    fulfilled by the lead contact, Thomas MacCarthy (thomas.maccarthy@stonybrook.edu).

735    **Materials availability**

736    This study did not generate new unique reagents.

737    **Data and code availability**

738    Data used for this research was published previously by Tang et al, 2020. A custom Python

739    package developed for this project is available at https://gitlab.com/maccarthyslab/deepshm.

740

741    **Methods details**

742

743    **Generating *k*-mer data**

744       Germline IGHV reference sequences were downloaded from the international

745    ImMunoGeneTics information system (IMGT) website (Lefranc, 2001). The leader portion of each

746    reference sequence was also extracted if available. To generate the $k$-mers of a given germline

747    sequence, $\pm\lfloor k/2 \rfloor$ nt sequences were extracted from the start of the V exon, where $k$ is the length of

748    the subsequence, and $\lfloor k/2 \rfloor$ represents the greatest integer less than or equal to $k/2$. This process

749    was continued, moving 1 nt at a time, until the end of the exon was reached. Next, all $k$-mers were

750    converted to their respective one-hot encodings. A one-hot encoding is a transformation of a DNA

751    sequence using a 2-D matrix containing only zeros and ones, where each row represents one of the

752    four ordered DNA bases and each column is an individual site in the sequence. For each column,

753    a "1" is filled in the row that matches the nucleotide of that site and a "0" in the remaining

754    unmatched rows (**Figure 1**).

755    **Calculating mutation frequencies, substitution rates, and weighted substitutions of $k$-mers**

756       Using a high-quality data set previously published by us (Tang et al., 2020), we calculated

757    the mutation frequencies of every $k$-mer in a germline sequence as the number of observed

758    mutations at each site (corresponding to a single $k$-mer), divided by the total number of sequences

759    the germline IGHV allele contained. The substitution rate of each $k$-mer was computed as the

760    number of times the middle nucleotide mutated from the germline nucleotide to the other four

761    DNA bases, divided by the total number of overall mutations. Note that a zero was recorded in the

762    instance the mutated base was the same as the germline context. Lastly, the weighted substitution

763    of a $k$-mer was simply calculated as the observed mutation frequency multiplied by the substitution

764    rate vector.

765    **CNN architecture and model optimization**

766    We implemented a convolutional neural network (CNN) to analyze the *k*-mer input data.

767    Three separate architectures were used to predict different SHM outcomes: mutation frequency,

768    substitution rate, and weighted substitution (see above). Although the hyperparameters that were

769    ultimately selected varied from model-to-model, all CNNs followed the same general architecture,

770    which consisted of one convolution layer, followed by two fully connected layers (**Figure 1**).

771    Additional parameters, such as dropout and batch normalization, were optimized by generating

772    100 separate models with randomly selected hyperparameters for each *k*-mer and corresponding

773    model architecture we generated. The range of values for all parameters and hyperparameters that

774    were tested for each architecture and output type are specified in **Supplementary Table 1**.

775    Next, we utilized 4-fold cross-validation to evaluate the performance of the model on

776    unseen (test) data. In total, there are seven IGHV families (IGHV1-7), where each IGHV family

777    consists of genes that share a high percentage of sequence similarity (Lefranc, 2001). The *k*-mers

778    derived from the three largest IGHV families, IGHV1, IGHV3, and IGHV4, formed three separate

779    groups, and the *k*-mers belonging to the remaining 4 smaller IGHV families constituted the final

780    group in order to create a data set comparable in size with the other groups. Thus, we separated

781    the data by their respective IGHV family to reduce the chances of model overfitting, since it is

782    likely that *k*-mers from the same IGHV family will be similar even if they come from different

783    genes and, therefore, bias the results if they appear in both training and test sets. In every cross-

784    validation fold, three of the data groups were used as training set, and the fourth used as test set.

785    We also evaluated the model performance, for each fold, by calculating the Pearson correlation(s)

786    between the predicted mutation frequency and/or substitution rate of the test set *k*-mers and the

787    equivalent output type of the empirical data. The average correlation across the 4 validation folds

788    was reported for the model, as in **Figure 2**.

789    As an additional step, we wrote a custom, universal Python script (available at

790    https://gitlab.com/maccarthyslab/deepshm) to automatically generate the CNN architecture,

791    parameters, and hyperparameters of each model, regardless of the output specified, to ensure that

792    all models were constructed in a consistent manner. All CNNs were generated using the built-in

793    Keras API in Tensorflow 2.4.1 and trained on GPU processors using three Nvidia GeForce RTX

794    2080 graphics cards.

**Inferring an S5F targeting model**

796    In order to ensure a fair comparison between S5F values and our deep learning predictions,

797    we used the SHazaM R package (Yaari et al., 2013) to create an S5F targeting model, which

798    provides analogous 5-mer mutability and substitution scores based on the same data set we used

799    to train our CNN models with. We specified the S5F targeting model to count both silent and

800    replacement mutations ("rs" parameter) since the mutation data we used was derived from non-

801    functionally rearranged VDJ coding sequences (i.e. in the absence of selection) and with each

802    sequence being clonally independent (Tang et al., 2020). Multiple mutations were handled

803    specifying the "independent" parameter, which treats each mutation independently. Default values

804    were used for all other parameters.

805

**Quantification and statistical analysis**

807

**Neural network encodings analysis**

809    The output (encoding) of the penultimate layer of the CNN model was used as a way to

810    explain the SHM patterns learned by the model. To generate the encodings from this layer, we

811    removed the last layer of the CNN while keeping the remaining layers intact. Next, we processed

812    the *k*-mers through the truncated model to retrieve the ensuing output values. We then applied t-

813    distributed stochastic neighbor embedding (t-SNE) in Python on these multidimensional encodings

814    to visually represent the resulting embedding in two dimensions.

## Cluster identification

816    We implemented k-means clustering to identify clusters within the t-SNE embedding of

817    the weighted substitution model (**Supplementary Figure 1**). We separated all *k*-mers sharing the

818    same middle nucleotide and then applied k-means clustering independently on each group to

819    facilitate the clustering process. All clustering assignments were performed using the *kmeans*

820    function in R. For each middle nucleotide, we specified the algorithm to identify 5 distinct clusters

821    and subsequently inspected the clusters to ensure a proper separation between clusters of distinct

822    mutabilities occurred. In the case of G and T nucleotides, there were resulting clusters (one for

823    each nucleotide) containing hot and cold sequences (i.e one "cold" and one "hot" subcluster per

824    cluster), so we manually split each of these clusters into two distinct ("cold" and "hot") clusters to

825    reduce the disparity in mutation frequencies.

## Identifying recurring genomic patterns using TF-MoDISco

827    We applied TF-MoDISco (Shrikumar et al., 2018), a machine learning interpretability

828    method, to identify recurring motifs our model detected in the 15-mer data. From the data, we

829    isolated four groups of 15-mers based on the middle nucleotide (A, C, G, or T) of the 15-mer, with

830    the additional condition that the middle nucleotide conformed to WR<u>C</u> or <u>G</u>YW AID hotspots, or

831 W<u>A</u> or <u>T</u>W Polη hotspots, respectively. TF-MoDISco requires importance scores to be used as

832 input, which can be generated by utilizing one of several attribution methods. Here we generated

833 the importance scores for each group by applying Integrated Gradients (Sundararajan et al., 2017)

834 to the most accurate 15-mer mutation frequency model. Using the resulting importance scores, we

835 then ran TF-MoDISco for all groups separately, still subject to the hotspot constraint, and requiring

836 each of the identified patterns to be associated with at least 20 input sub-sequences (or "sequelets").

837

838 **Key resources table**

839

| REAGENT or RESOURCE | SOURCE | IDENTIFIER |
|---|---|---|
| Deposited data | | |
| Data from the memory, marginal zone, and plasma cell subsets (B10-B14, B16-21, HD001-10) | Tang et al., 2020 | NCBI SRA BioProject IDs 381394, 591804 |
| Software and algorithms | | |
| DeepSHM | This paper | https://gitlab.com/maccarthyslab/deepshm |
| TF-MoDISco | Shrikumar et al., 2018 | https://github.com/kundajelab/tfmodisco |
| SHazaM | Yaari et al., 2013 | https://shazam.readthedocs.io/en/stable/ |

840