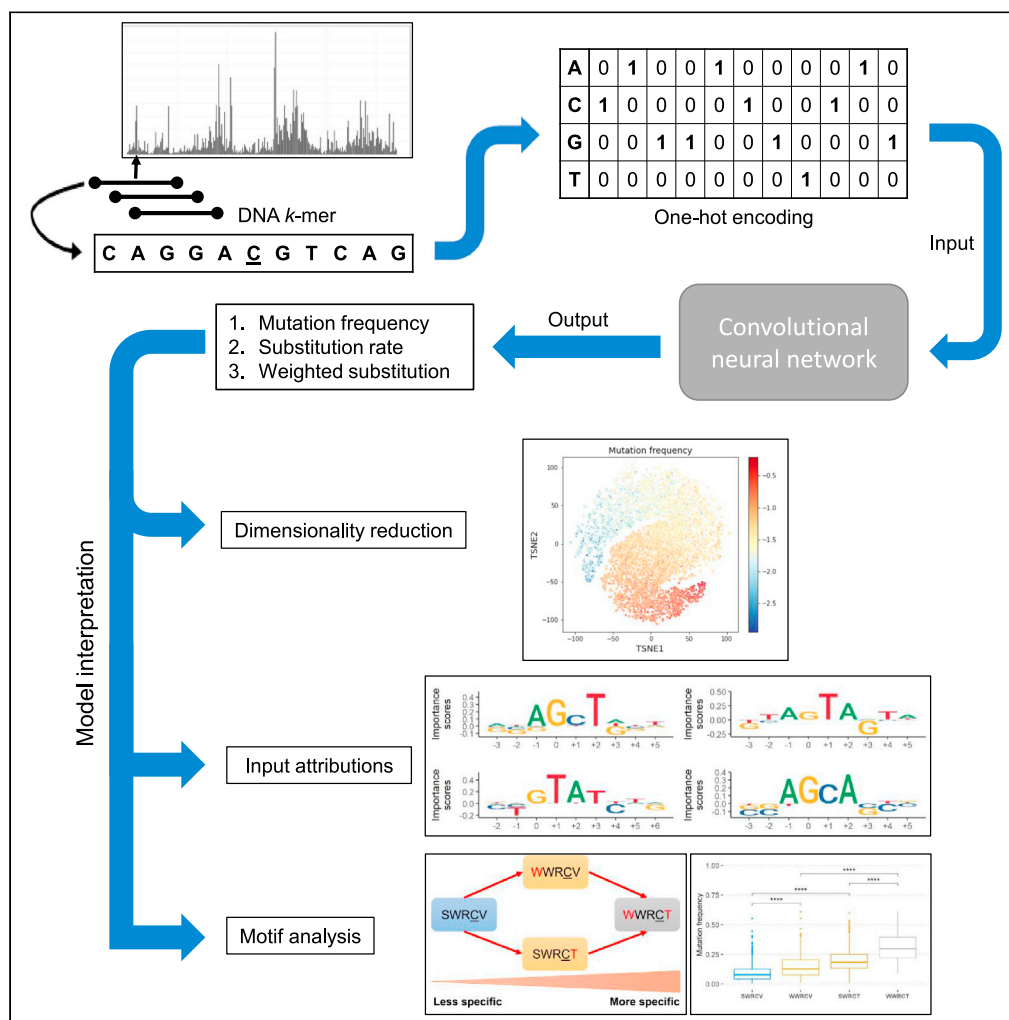


Article

Deep learning model of somatic hypermutation reveals importance of sequence context beyond hotspot targeting



Catherine Tang,
Artem
Krantsevich,
Thomas
MacCarthy

thomas.maccarthy@
stonybrook.edu

Highlights
Deep learning on somatic hypermutation (DeepSHM) improves mutation predictions

DeepSHM shows importance of wider DNA context (up to +/- 7 nucleotides)

Model interpretation reveals features contributing to SHM, such as G content

Results suggest WWRCT as a refined AID hotspot motif

Tang et al., iScience 25, 103668
January 21, 2022 © 2021
<https://doi.org/10.1016/j.isci.2021.103668>



Article

Deep learning model of somatic hypermutation reveals importance of sequence context beyond hotspot targeting

Catherine Tang,^{1,3,4} Artem Krantsevich,^{1,4} and Thomas MacCarthy^{1,2,5,*}

SUMMARY

B cells undergo somatic hypermutation (SHM) of the Immunoglobulin (Ig) variable region to generate high-affinity antibodies. SHM relies on the activity of activation-induced deaminase (AID), which mutates C>U preferentially targeting WRC (W=A/T, R=A/G) hotspots. Downstream mutations at WA Polymerase η hotspots contribute further mutations. Computational models of SHM can describe the probability of mutations essential for vaccine responses. Previous studies using short subsequences (*k*-mers) failed to explain divergent mutability for the same *k*-mer. We developed the DeepSHM (Deep learning on SHM) model using *k*-mers of size 5–21, improving accuracy over previous models. Interpretation of DeepSHM identified an extended WWRCT motif with particularly high mutability. Increased mutability was further associated with lower surrounding G content. Our model also discovered a conserved AGYCTGGGG (Y=C/T) motif within FW1 of IGHV3 family genes with unusually high T>G substitution rates. Thus, a wider sequence context increases predictive power and identifies features that drive mutational targeting.

INTRODUCTION

Upon encountering antigen, germinal center (GC) B cells undergo several programmed mutational events in secondary lymphoid organs to mount an effective humoral immune response. Somatic hypermutation (SHM) takes place in the GC dark zone whereby mostly point mutations are introduced into the Immunoglobulin (Ig) variable (V) region. Selection for mutations leading to higher binding B cell receptors to cognate antigen occurs in the GC light zone, thus, producing a diverse repertoire of high-affinity antibodies (Methot and Di Noia, 2017; Pilzecker and Jacobs, 2019; Rajewsky, 1996). The mutagenic enzyme, activation-induced deaminase (AID), initiates SHM (Muramatsu et al., 2000) by converting cytosine (C) to uracil (U) in single-stranded DNA (ssDNA), resulting in a U:G (guanine) mismatch (Bransteitter et al., 2003). AID displays preferential targeting at WRC/GYW "hotspot" motifs (where W=A/T, R=A/G, Y=C/T, and the underlined base indicates the mutated base in the top and bottom strand, respectively), whereas SYC/GRS "cold-spots" (S=C/G) are significantly less targeted (Pham et al., 2003; Rogozin and Diaz, 2004; Rogozin and Kolchanov, 1992; Yu et al., 2004). If left unrecognized, U mismatches will act as a template T and be replicated over (Pilzecker and Jacobs, 2019). The resulting C>T transition mutation is commonly referred to as the DNA "footprint" of AID (Liu et al., 2008). Downstream DNA repair further contributes to antibody diversity that is mediated by low-fidelity polymerases. During non-canonical base-excision repair (ncBER), the U:G mismatch is recognized and excised by uracil-DNA glycosylase (UNG), resulting in an abasic site (Rada et al., 2004). Repair of these abasic sites by REV1 can cause both transition and transversion mutations at C:G base-pairs (Jansen et al., 2006). In the case of non-canonical mismatch repair (ncMMR), the U:G mismatch is recognized by the MSH2/MSH6 heterodimer. Next, EXO1 exonuclease is recruited to create a patch of ssDNA, which then allows error-prone polymerases, particularly Polymerase eta (Pol η), to resynthesize. Pol η is known to create mutations at neighboring adenine (A) and thymine (T) sites of the initial AID-induced lesion, most notably at WA/TW hotspot motifs (Matsuda et al., 2001; Mayorov et al., 2005).

Several computational models have been developed for the SHM process and intrinsic biases exhibited by key proteins such as AID and Pol η . These models have mainly utilized *k*-mer subsequences, where *k* is a specified integer length, ranging between 3-mers to 7-mers (Cui et al., 2016; Elhanati et al., 2015; Shapiro et al., 1999, 2002; Yaari et al., 2013). Two of these models (Cui et al., 2016; Yaari et al., 2013) are widely used

¹Department of Applied Mathematics and Statistics, Stony Brook University, Stony Brook, NY 11794, USA

²Laufer Center for Physical and Quantitative Biology, Stony Brook University, Stony Brook, NY 11794, USA

³Present address: Exploratory Science Center, Merck & Co., Inc., Cambridge, MA, USA

⁴These authors contributed equally

⁵Lead contact

*Correspondence: thomas.maccarthy@stonybrook.edu

<https://doi.org/10.1016/j.isci.2021.103668>



and have leveraged 5-mer motifs to capture the dependency of the local surrounding sequence for the middle nucleotide to mutate, while simultaneously bypassing any influence of selection. The first of these targeting models ("S5F") evaluates all possible 5-mers and synonymous (silent) mutations derived from functionally rearranged, or productive, VDJ coding sequences (Yaari et al., 2013). The second model ("RS5NF") similarly assesses 5-mers but uses both synonymous and non-synonymous (replacement) mutations from non-productively (non-functional) rearranged sequences (Cui et al., 2016). Such models have been used to simulate B cell repertoire lineages by constructing a set of hypothetical sequences that have been mutated in a sequential manner as governed by, for example, the underlying S5F substitution scores (Krantsevich et al., 2021; Sheng et al., 2017). Although *k*-mer approaches are generally able to capture some key local intrinsic biases of SHM, such as hotspot targeting, there is evidence that shorter *k*-mers are insufficient to properly characterize differential SHM targeting. For example, a recent study extended a local sequence (5-mer) context model and improved accuracy by including parameters describing the position within the IGHV gene (Spisak et al., 2020). Another study compared the mutability of identical 5-mer (middle position ± 2 nt) motifs at different positions within an IGHV gene (Zhou and Kleinstejn, 2020), and found that the mutation frequency of these motif-allele pairs (MAPs) positively correlates with the overall mutability of a wider neighborhood of motifs, suggesting that an extended *k*-mer may better capture SHM.

Current efforts towards the development of HIV and influenza vaccines aim to recreate the mutations leading to known broadly neutralizing antibodies (bnAbs) (Haynes et al., 2012). A major obstacle to this approach is that, particularly in the case of HIV, bnAbs have acquired large numbers of mutations. A previous study used a methodology (ARMADiLLO) that combined the S5F model with simulations and found that many key bnAbs mutations had a low probability of occurring naturally (Wiehe et al., 2018). Simulation approaches using mutability models such as S5F are needed due to non-independence between sites, most obviously when mutations create or destroy hotspots at adjacent sites (Krantsevich et al., 2021). Improved mutability models can therefore be used to better evaluate the probability of generating key mutations in the context of vaccine development.

Earlier studies have shown that using deep learning is effective in different genomic applications; for example, convolutional neural networks (CNNs) in extracting conserved sequence motifs among target sequences (Alipanahi et al., 2015; Kelley et al., 2016; Zhou and Troyanskaya, 2015). In this study, we adopted a deep learning approach using a 2-D CNN to analyze extended *k*-mer lengths to better understand the underlying SHM process. We demonstrate that our model, DeepSHM (Deep learning on SHM), can more accurately represent the SHM process by evaluating longer *k*-mers of up to 21 nts. In addition, DeepSHM using 15-mers as inputs was able to recapitulate AID WRC/GYW hotspot motifs and identify an extended WWRCT motif. Neural network predictions are notoriously difficult to explain (the "black box" problem), but many new methods are available to interpret results (Koo and Ploenzke, 2020). We used one such method to identify a negative association between increased mutability at a site and its surrounding G content. On the other hand, lower mutation frequency was correlated with increased substitution rates of certain substitution types, particularly for G>T and C>A mutations. Furthermore, many highly conserved sites within G-rich subregions belonging to several IGHV3 genes display an extremely high bias towards creating G mutations, some of which may participate in the formation of G-quadruplex (G4) structures.

RESULTS

Deep learning can more accurately represent SHM mutability and substitution biases

The objective of our analysis was to use supervised deep learning to build an accurate convolutional neural network (CNN) for SHM and, as much as possible, identify features contributing to mutability. We chose CNNs because we still expected mutation frequency to depend on recurring motifs that might occur at any position in the sequence (most obviously, AID hotspots), a task CNNs are well suited to. The workflow of our network consists of an input layer that processes a *k*-mer subsequence represented in its one-hot encoding format (i.e., a $4 \times k$ matrix of zeros and ones), followed by a convolution layer and two fully connected layers as the hidden layers, and finally the output layer of size 4×1 or 1×1 , depending on the task that is being predicted (Figure 1, see STAR Methods). Several hyperparameters, including dropout rate and learning rate, were fine-tuned with our model as well (Table S1A). We defined a model that would separate mutations on each strand (which are predominantly at C and A on the top strand and at G and T on the bottom strand) at the input level. To achieve this, we identified a simple solution using padding that assigns a row in each channel of the convolution layer output separately to each strand (Figure 1). CNNs are also often used together with attribution methods such as Integrated Gradients, to help with interpretation of the results.

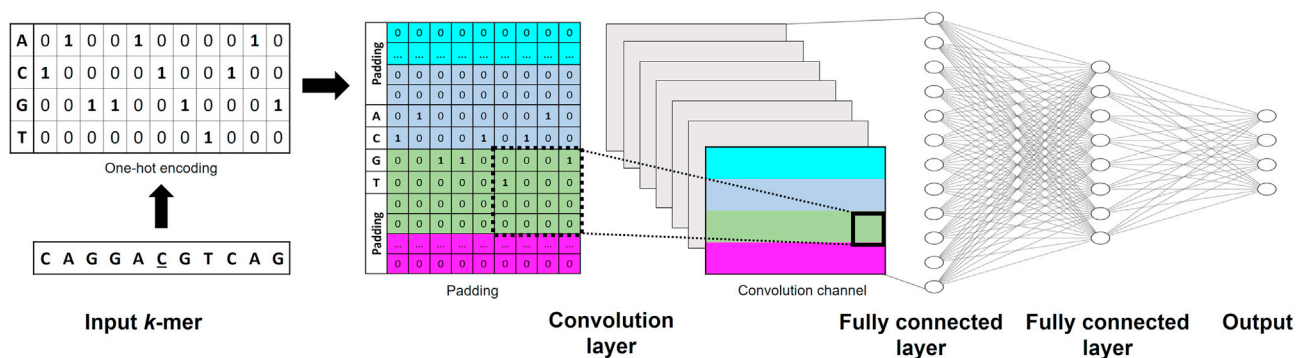


Figure 1. DeepSHM model architecture

Each model had an input layer, one convolution layer, two fully connected layers, and an output layer. The input layer was a $4 \times k$ dimensional one-hot encoded matrix (k is the length of the subsequence). The dimension of the output layer was dependent on the task: substitution (4×1), mutation frequency (1×1), or weighted substitution (4×1). For the convolutional layer, ‘same’ padding was used to allow the model to process top and bottom strand mutations separately. With ‘same’ padding, the output of each convolutional channel has the same shape as the input ($4 \times k$) with the following properties: the first and the fourth rows are populated with zeros only (there was no real input, only padding; cyan and magenta rows); the input used for the second (light blue) row contained two rows of padding and two data rows corresponding to A or C nucleotides only; and similarly, the input used for the third (green) row also contained two rows of padding and two rows of data corresponding to G or T nucleotides. Because AID and Pol η target C and A sites, respectively, this approach was taken with the expectation of helping the model distinguish top and bottom strands.

As a starting point, we trained two CNN models, which we collectively refer to as DeepSHM (Deep learning on SHM), to separately predict mutation frequency and substitution rates, calculated from previously published B cell repertoire data containing non-productively rearranged and clonally independent VDJ coding sequences (Tang et al., 2020), for varying k -mer lengths (see STAR Methods). We trained both models independently using different combinations of k -mer lengths and hyperparameters as listed in Table S1A. We found that for predicting mutation frequency, 15-mers were moderately better than 9-mers (purple boxplots in Figure 2A, Mann-Whitney U test: $P < 1.1 \times 10^{-11}$) and that further extending the motif length to 21 did not improve accuracy because both produced an overall maximum correlation (across hyperparameters) of ~ 0.81 (Figure 2A, Table 1). Thus, using k -mers of length 15 or longer outperformed shorter lengths, specifically 5-mers and 9-mers (Table 1), suggesting that an extended DNA motif can better model the SHM process. However, using longer k -mers did not substantially improve the model that predicts SHM substitution bias alone, achieving an average correlation of 0.56 for 15-mers (green boxplots in Figure 2B, Table 1), but which is similar for different lengths. For the interpretability analysis below, we chose to use the best 15-mer models to keep the k -mer length consistent for comparisons across all models. To check if the performance of the models leading to the best results was consistent, we also trained 30 different iterations of each model, keeping the hyperparameters fixed but using different random seeds. We found the standard deviation across correlations was very small, at 0.002 for the mutation frequency model and 0.001 for the substitution rate model, showing the strong consistency of our results.

We next sought to compare DeepSHM against the widely used S5F model that is based on 5-mer motifs (Yaari et al., 2013). To ensure a fair comparison, we generated an S5F targeting model using the same data set that was used to train DeepSHM, as well as the same cross-validation scheme (see STAR Methods). Using the same test set splits as above, we found that there was an average correlation of 0.69 between the predicted S5F model mutability and empirical mutation frequency, and an average correlation of 0.52 for predicted S5F substitution scores and empirical substitution rates (red dashed lines in Figure 2, Table 1). The substitution model slightly outperformed S5F for all k -mer values we analyzed. The mutation frequency model achieved a modest improvement over S5F using 5-mers as an input, and this difference became more evident for 9-mers, 15-mers, and 21-mers (Figure 2A, Table 1). To shed light on performance variation not due to changing hyperparameters, we used 30 iterations (using different random seeds) of the best 15-mer models for both mutation frequency and substitution models discussed above and found these iterations to have significantly greater accuracy than S5F both individually and in aggregate ($P < 1.8 \times 10^{-6}$ for each model, Wilcoxon signed-rank test). Overall, these results show that our deep learning approach successfully extracts meaningful information from the wider sequence context to improve predictions.

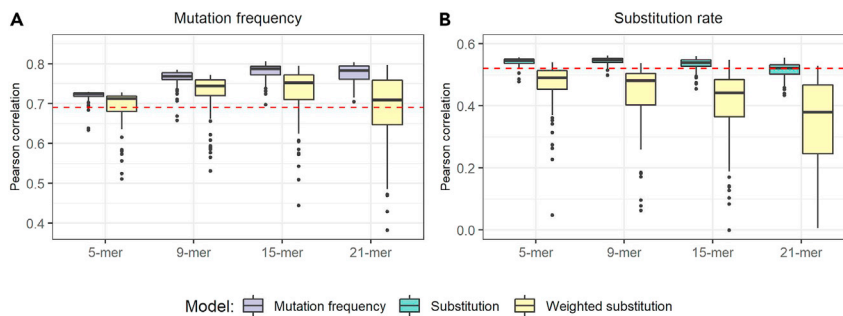


Figure 2. Performance of DeepSHM

(A and B) Boxplots describing the distribution (across random hyperparameters) of Pearson correlations between DeepSHM predictions and empirical data (y-axis) are shown for different input *k*-mers (x-axis) for (A) mutation frequencies, and (B) substitution rates, for all three models (mutation frequency, substitution, and weighted substitution). Red dashed lines signify correlations of predicted S5F values, which uses 5-mers.

To identify associations between mutation frequency and specific substitutions, we further constructed a DeepSHM model to predict the "weighted substitution" of a *k*-mer, i.e., the product of the percentage of each observable substitution type (e.g. G>N) and the mutation frequency of the *k*-mer (see [STAR Methods](#)). Note that this weighted substitution metric is a vector representing the four ordered DNA bases, with a "0" placed at the position that matches the middle nucleotide of the *k*-mer. Because weighted substitution constitutes aspects of both the observable mutation frequency and substitution rate of the middle nucleotide of a given *k*-mer, we were able to evaluate DeepSHM on each metric separately. Although this model made poorer substitution rate predictions on average (varying hyperparameters) than S5F ([Table 1](#)), the best model performed similarly to S5F for substitution rates while performing almost as well as any model in predicting mutation frequency. Cross-validation in this instance produced a range of correlations between 0.51 and 0.57 for predicted substitution rates – a level similar to that of S5F ([Figure 2B](#), [Table 1](#)). On the other hand, the DeepSHM model for weighted substitution values was almost as good at predicting mutation frequency for 15-mers (correlation: 0.79) as the standalone model that was tasked to learn mutation frequency only, as well as outperforming S5F ([Figure 2A](#), [Table 1](#)).

Interpretation of the DeepSHM network reveals extended hotspot motifs

A complication often associated with deep neural networks is model interpretability (the "black box" problem). One way we interrogated the predictions made by DeepSHM, and what it has learned about the SHM process, was to analyze the output of the penultimate layer of each 15-mer based model. In particular, analyzing the output, or "encodings", of this layer can be viewed as an alternative, and more informative, way of representing the input 15-mer. To visualize the multi-dimensional encodings of the individual 15-mers, we used t-SNE, a dimensionality reduction technique, to project each onto a 2-dimensional embedding (see [STAR Methods](#)). At this point in order to make full use of the data, we merged all of the 15-mer data into one training set, and then trained three new individual models (one for each output type) using the hyperparameters which previously led to the best cross-validation results. The analyses we present below are derived from the DeepSHM models that were trained using this merged data set.

We began by identifying features learned by DeepSHM that predicted weighted substitutions. Because weighted substitution is a measure of both mutation frequency and substitution bias, the embedding should capture both metrics simultaneously. Each point in the resulting t-SNE embedding in [Figure 3A](#) represents a single 15-mer and is colored according to its corresponding mutation frequency. We found that the *k*-mers separated into a relatively small number of clusters and that within the four largest clusters there was a gradient of mutability. When we considered the middle nucleotide of each 15-mer, we observed that these clusters also shared the same middle nucleotide ([Figure 3B](#)), suggesting that the network identified as a key feature the "0" value in the weighted substitution output vector that is associated with the middle nucleotide.

Next, we applied clustering on the embedding as a way to generate smaller clusters for further analysis ([Figure S1](#), see [STAR Methods](#)). We subsequently created a sequence logo plot representing each cluster as shown in [Table S2](#). As expected, clusters with the highest mutation frequencies had inner subsequences

Table 1. Cross-validation of various input *k*-mer sequences

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	0.56	0.56
	IGHV3	0.51	0.54	0.55	0.55	0.54
	IGHV4	0.54	0.57	0.57	0.58	0.57
	IGHV2, 5, 6, 7	0.53	0.55	0.56	0.56	0.55
	Avg correlation	0.52	0.56	0.56	0.56	0.55
	Best - S5F	NA	0.04	0.04	0.04	0.03
	Mean - S5F	NA	0.02	0.03	0.01	-0.01
	p-value	NA	3.52E-15	3.46E-17	3.78E-9	0.52
Mutation frequency	IGHV1	0.69	0.74	0.79	0.82	0.82
	IGHV3	0.68	0.74	0.79	0.8	0.79
	IGHV4	0.69	0.74	0.79	0.84	0.84
	IGHV2, 5, 6, 7	0.70	0.69	0.76	0.78	0.77
	Avg correlation	0.69	0.73	0.78	0.81	0.80
	Best - S5F	NA	0.04	0.09	0.12	0.11
	Mean - S5F	NA	0.03	0.07	0.09	0.09
	p-value	NA	2.76E-15	2.31E-17	2.31E-17	2.31E-17
Weighted substitution (substitution rate)	IGHV1	0.52	0.55	0.53	0.55	0.53
	IGHV3	0.51	0.52	0.53	0.52	0.51
	IGHV4	0.54	0.55	0.54	0.57	0.54
	IGHV2, 5, 6, 7	0.53	0.53	0.54	0.55	0.53
	Avg correlation	0.52	0.54	0.54	0.55	0.53
	Best - S5F	NA	0.02	0.02	0.03	0.01
	Mean - S5F	NA	-0.06	-0.08	-0.11	-0.17
	p-value	NA	4.19E-14	2.91E-16	6.82E-17	2.31E-17
Weighted substitution (mutation frequency)	IGHV1	0.69	0.74	0.78	0.80	0.81
	IGHV3	0.68	0.73	0.78	0.80	0.78
	IGHV4	0.69	0.74	0.78	0.80	0.82
	IGHV2, 5, 6, 7	0.70	0.70	0.75	0.77	0.78
	Avg correlation	0.69	0.73	0.77	0.79	0.80
	Best - S5F	NA	0.04	0.08	0.10	0.11
	Mean - S5F	NA	0	0.04	0.04	0
	p-value	NA	0.001	2.33E-8	1.49E-7	0.25

The correlations of repeatedly trained models using different random seeds (but the same hyperparameters) for neural network training had small standard deviations, in all cases below 0.01. p-values are from a Wilcoxon signed-rank test comparing the training results for each model with the corresponding S5F model accuracy. p-values were corrected (Benjamini-Hochberg) for multiple comparisons.

containing AID (C cluster 8, G cluster 12) and Pol η (A cluster 3, T cluster 20) hotspots. For AID, these are WRC (Figure 3C) and GYW motifs (Figure 3D). Within the two most highly targeted AID hotspot clusters there is a substantial presence of both WGC/GCW and WAC/GTW contexts, rather than only the well-known WGCW overlapping hotspot motif (Tang et al., 2020; Wei et al., 2015). Furthermore, even when we include WAC/GTW, there is a preference for a T base at the 3' end of the WRC hotspot, and conversely, a bias for an A base at the 5' end of the GYW hotspot (Figures 3C and 3D). This motif is consistent with a genome-wide study of AID mutations in mice that reported observing high mutability at AACT and AGCT motifs in both strands (Álvarez-Prado Á et al., 2018). When we assessed the mutability of all possible WRCN (N = A/C/G/T) motifs separately, we observed WRCT to be the most highly mutated in each case (Figure S2). Previous studies identified WRCY/RGYW (Y=C/T, R = A/G) and later WRCH/DGYW (H = A/C/T, D = A/G/T) to be a better predictor of mutability at C:G bases (Rogozin and Diaz, 2004; Rogozin and Kolchanov, 1992). However, we discovered some inconsistencies with these definitions, as AGCC was found to

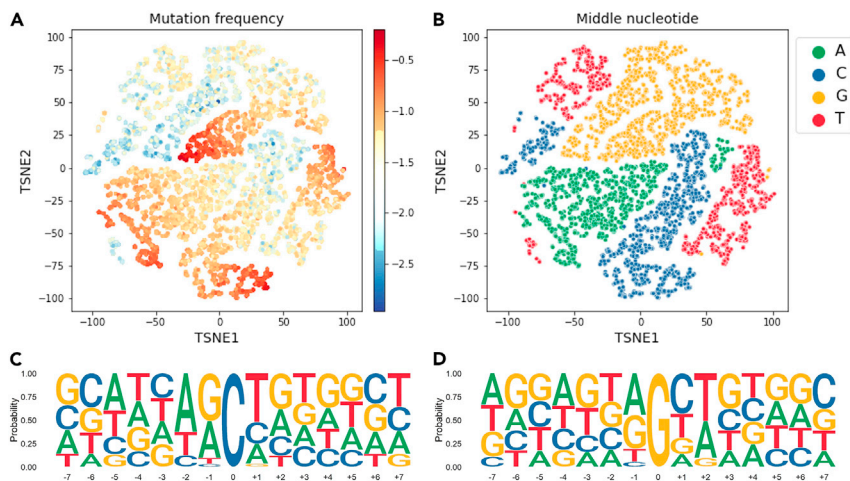


Figure 3. Neural network encodings analysis: weighted substitution model

(A and B) 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 STAR Methods) and is colored according to its corresponding (A) mutation frequency (\log_{10}) and (B) middle nucleotide. (C and D) Consensus sequences derived from the highest mutated cluster identified using PAM clustering on the embedding of 15-mers containing either (C) a middle C nucleotide or (D) a G nucleotide (clusters 8 and 12 in Table S2).

be the least mutated of the AGCN motifs and WRCG was not always the least mutated, on both strands (Figure S2). Overall, these early hotspot definitions may have been too broad, and WRCT/AGYW is a more consistent predictor of AID targeting. Lastly, we also noted that among the least mutated k -mer clusters, many were G-rich in their surrounding context (for example, C cluster 9, A cluster 5, T cluster 17), and particularly for G (G cluster 15) (Table S2).

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 STAR 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 because 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 STAR Methods). We began by analyzing the importance scores derived from Integrated Gradients (Sundararajan et al., 2017) of 15-mers whose middle nucleotide conformed to a WRC/gyw AID hotspot motif. As expected, the positively contributing sites in the set of ensuing motifs aligned with the hotspot motifs (Figures 4A and 4B). Moreover, TF-MoDISco again revealed a preference for having a T base at the +1 position of the WRC (WRCT , Figure 4A) and an A base at the -1 position of the gyw (AGYW , Figure 4B). In addition to WRC/AGYW being a well-represented motif identified by TF-MoDISco, as measured by having positive contributions to mutability (above horizontal axis on Figure 4), we also noticed many neighboring C and G bases contained negative contributions (below horizontal axis on Figure 4), most evidently at the -3 position of the WRC hotspot, and to an extent at the +3 position of the gyw hotspot (Figure 4). Here, the negative contribution at the -3 position signifies that having an S ($\text{S}=\text{C/G}$) at that position reduces mutational targeting to the middle C.

We next sought to determine whether adding the 5' W or 3' T context of the canonical WRC hotspot is more influential in terms of increasing its susceptibility to AID mutagenesis. To address this, we increased the hotspot specificity step by step, starting from SWRCV ($\text{V}=\text{A/C/G}$), and assessed the impact a single change in the motif at either S or V site, causing a WWRCV or SWRCT intermediate motif to form, respectively, has on mutability (Figure 5A). We found that both WWRCV and SWRCT intermediate hotspots were shown to mutate significantly more than SWRCV (Figure 5B). We also discovered that the mutability of the WWRCT hotspot, which contains the extended hotspot in both 5' and 3' directions, was significantly higher than both intermediate hotspots (Figure 5B). Performing a pairwise comparison between the mutation frequency of all 16 individual (W/S) WRC (T/V) contexts further confirmed that those containing both a 5' W and 3' T were significantly more mutated than the remaining hotspot motifs, with WAGCT being the

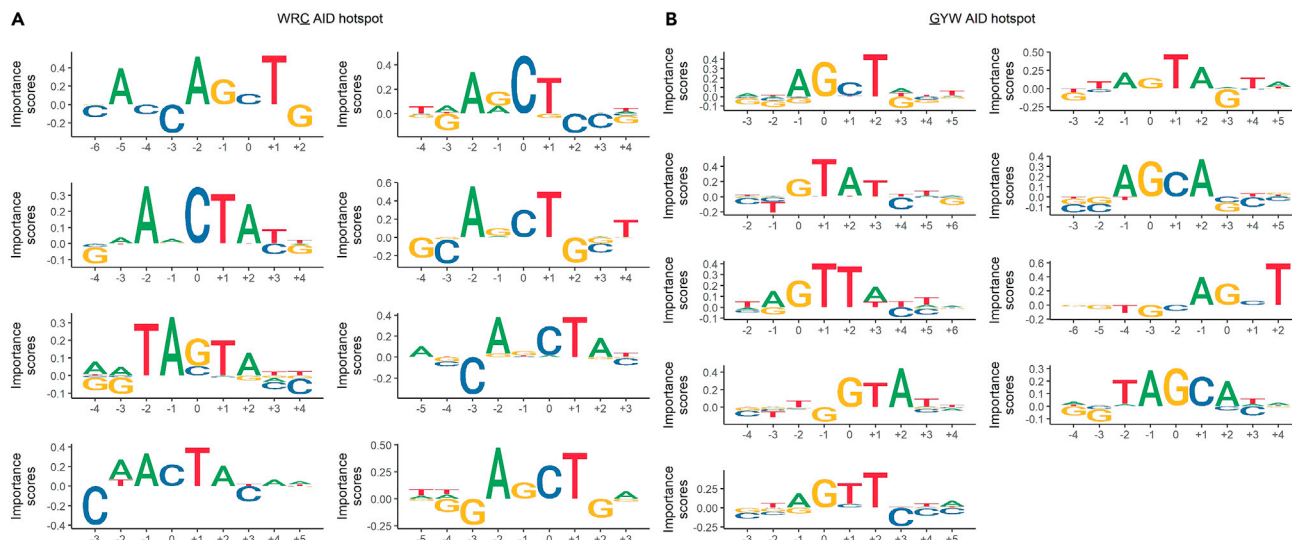


Figure 4. Recurrent motifs identified by TF-MoDISco

(A and B) TF-MoDISco results using the Integrated Gradients as base-level importance scores of 15-mers whose middle nucleotide conformed to a (A) WRQ or (B) GYW AID hotspot motif.

most mutated (Figure 5C, Table S3). In addition, the next four successively mutated hotspots followed an SWRCT context, overall suggesting the 3' T to be more impactful to AID recognition than the 5' W.

The motif TA emerged when we applied TF-MoDISco to the 15-mers conforming to either a WA (Figure S3A) or TW Pol η hotspot (Figure S3B). In addition to our model identifying the TA hotspot motif as important, it also suggested TAT, TAH (where H=A/C/T, because 3' G makes a negative contribution) on the top strand, and its bottom strand equivalent D \overline{TA} (D=A/G/T, because 5' C makes a negative contribution). However, when we compared the mutation frequencies of TAH and TAG, there was no significant difference, although we did find that AAH > AAG, DTA > CTA and DTT > CTT (Figure S4). Now considering the case of TAT, we also found more generally that WAT/ATW hotspots mutate significantly more than their WAV/BTW (B=C/G/T) counterparts (Figure 6). Thus, although TA hotspots consistently have higher mutability than AA, the presence of a 3' T individually increases the mutability of each of these Pol η hotspots.

Highly targeted sites display a lower surrounding GC content

We next applied the same t-SNE methodology to the DeepSHM model that predicted only mutation frequency. We found that the organization of the subsequent embedding followed a direction of descending mutation frequency, with the highest mutating 15-mers located at the bottom-right portion of the plot (Figure 7A) and low-mutating 15-mers at the upper-left, which was enriched with FW1 15-mers (Figure 7B). In addition, we examined the possible influence of the local surrounding sequence by calculating the individual base content of the four DNA bases in each 15-mer. However, the inner 5-mer, which contains the dominant context, was excluded when computing all base counts. When we colored the t-SNE embedding according to the GC content of each 15-mer, we observed that GC content increases along the same direction as decreasing mutability seen previously (Figure 7C). Quantifying this observation more formally, we indeed found a significant negative correlation between the GC content and the mutation frequency of the 15-mers ($R = -0.32$, $P < 2.2 \times 10^{-16}$; Figure 7D). On the other hand, when we considered each individual base count independently, we observed that the count of G nucleotides specifically shows a stronger negative correlation ($R = -0.22$) than the C nucleotide count ($R = -0.074$) alone (Figure S5A), although both correlations are highly significant ($P < 2.2 \times 10^{-16}$). This result is consistent with the cluster analysis above (Table S2) where we observed several clusters with G-rich k -mers and low mutation frequencies. If we further separate the mutation frequencies into categories defined by the middle nucleotide, we find that G content has a consistent negative correlation regardless (column G of Figure S5B). More generally, A and T richness (columns A and T of Figure S5B) mostly shows a consistent positive correlation, whereas C and G richness shows a consistent negative (or nonsignificant) correlation. In summary, it appears that low-mutating sites

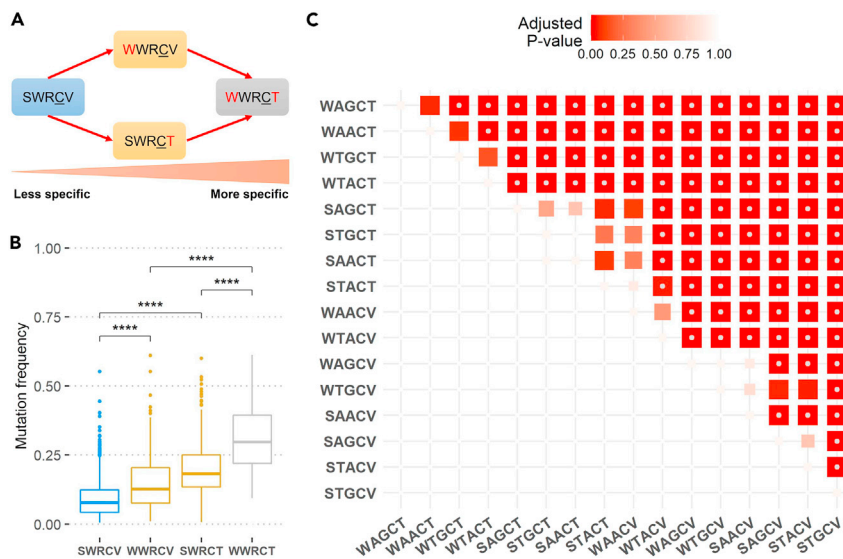


Figure 5. Mutability of extended AID hotspots

(A) Schematic showing an increase of AID hotspot specificity (left to right).

(B) Boxplots displaying the mutability of different (W/S)WRC(T/V) hotspot contexts, where W=A/T, S=C/G, R=A/G, V=A/C/G. Asterisks indicate significance ($p \leq 0.0001$) 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 (S/W)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 .

generally have a high local GC (and particularly G) content, and conversely, that highly targeted sites display an elevated local AT (particularly A) content.

Conserved FW1 sites in IGHV3 genes display a high T > G transversion bias

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 Pol η intervention during noncanonical base-excision repair (BER) and noncanonical 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 a C, G, or T middle nucleotides (Figure 8B). Because the model was tasked with learning the distributed substitution rates of each 15-mer, we next sought to evaluate the rate of each individual substitution type (e.g. C>T) within the embedding. In certain clusters, a similar gradient of high to low substitution rates could also be seen as we observed for mutation frequency (Figures 8C–8F). For instance, we noticed the rate of G>T substitutions was highest at the outer boundaries of the cluster (top-right cluster in Figure 8F), which was also associated with low mutation frequencies in the same cluster (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, IGHV4-34), so as to include sites with low mutation frequencies at high coverage, and calculated the correlation between mutation frequency and rate of substitution for each substitution type. As an example, for IGHV3-23 we found a significant negative correlation for G>T mutations ($R = -0.29$, $p = 0.0057$; Figure S6). Alternatively, we observed a significant positive correlation between mutation frequency and A>T mutations ($R = 0.35$, $p = 0.0076$; Figure S6) and G>A ($R = 0.22$, $p = 0.039$). Similar patterns to IGHV3-23 were also observed for IGHV1-18*01 and IGHV4-34*01, although there was some variation in which correlations were significant (Figure S6). The fact that the trend for C>T and G>A transitions is positive and that for C>A and G>T transversions is

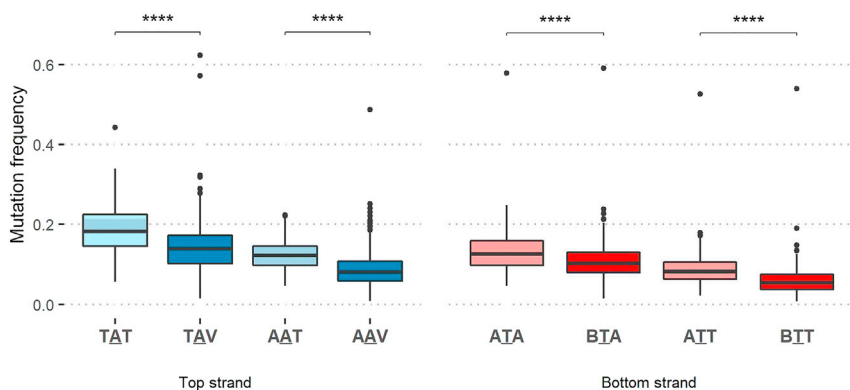


Figure 6. Mutability of extended Pol η hotspot motifs

Boxplots comparing the mutation frequency of various top strand WAT against WAV (blue), and bottom strand ATW against BTW (red) motifs. Asterisks indicate significance ($p \leq 0.0001$) of a one-sided Mann-Whitney U test comparing the greater mutation frequency of the boxplot on the left against the one on the right.

negative (though not always significant) for all three genes, is consistent with replication bypass, which predominantly causes C>T transitions, being favored over BER at sites with high mutation frequency.

In the t-SNE analysis of the substitution model, we also discovered two small clusters of 15-mers containing a C and T as their middle nucleotide (left side of Figure 8A) that did not group with their respective larger clusters, suggesting that these particular sites might have distinct substitution patterns. Generating the consensus sequence of the outlier T k-mers revealed a partially conserved AGYCTGGGGG sequence (Figure 9A). When we examined these subsequences more closely, we discovered that they were located only in IGHV3 family genes at either position 21 or position 45 according to the IMGT unique numbering system (Lefranc, 2001). The motif was also surprisingly common. At position 21 it appeared in 26 different alleles (across 16 genes) and was fully conserved in all alleles. Coincidentally, the motif also appeared in 25 different alleles (across 15 genes) at position 45, although it sometimes differed slightly at the +4 position (Figure 9A). These two sets of alleles only partially overlap, such that 15 alleles had the motif at both positions 21 and 45. Thus, this specific motif in FW1 of the IGHV3 family genes appears to be highly conserved evolutionarily, suggesting a possible functional role. The rates of substitution at these sites were also found to be highly biased towards creating T>G mutations, with an average T>G rate of about 0.74 at position 21, and an even greater rate of 0.88 at position 45 (background rate: 0.27 ± 0.22) (Figure 9B, Table 2). A previous study using Sanger sequencing data that was limited to IGHV3-23 and the pseudogene IGHV3-h had noted similarly high T>G substitution rates at positions 21 (for IGHV3-h) and 45 (for IGHV3-23) (Ohm-Laursen and Barington, 2007). Although the T subjected to mutation at both positions did not conform to a bottom strand TW Pol η hotspot, these genes at position 45 displayed a relatively high average mutation frequency of 0.17 ± 0.08 (Table 2), which is somewhat unusual given that mutations are generally more biased towards the CDRs than FW regions (Cohen et al., 2011; Shapiro et al., 2002), and that we reported above that many sites within FW1 tended to display low mutability (Figures 7A and 7B).

While examining the C outlier cluster (left of Figure 8A), we found the consensus sequence to be more conserved compared to the outlier with a middle T (Figure S7A). On the other hand, we noticed some overlap between both outlier clusters because, in some cases, the C corresponded to position 44 that preceded the middle T of the other outlier cluster (Figure S7A, Table 2). We further found this site to have a similar elevated C>G substitution rate (mean rate of 0.87 compared to background mean of 0.35, $P < 2.2 \times 10^{-16}$) (Figure S7B, Table 2), suggesting the model distinguished sites with a general preference to create G mutations.

Given that the sites with strikingly high T>G and C>G substitution rates we identified here have adjoining G-rich subregions (Figures 9A and S7A), we evaluated the possible influence these mutations might have on the formation of G-quadruplex (G4) structures. In a recent study, we assessed the potential for DNA G4 structures to form in the Ig V region, using a pre-trained deep learning model that computes the G4 potential of a linear DNA sequence (Tang and MacCarthy, 2021). There we found that the IGHV3 family had the highest propensity to form stable G4s in the top strand. We now sought to assess the overall mutational effect on G4 assembly of the IGHV3 sites that are biased towards G mutations. Following the

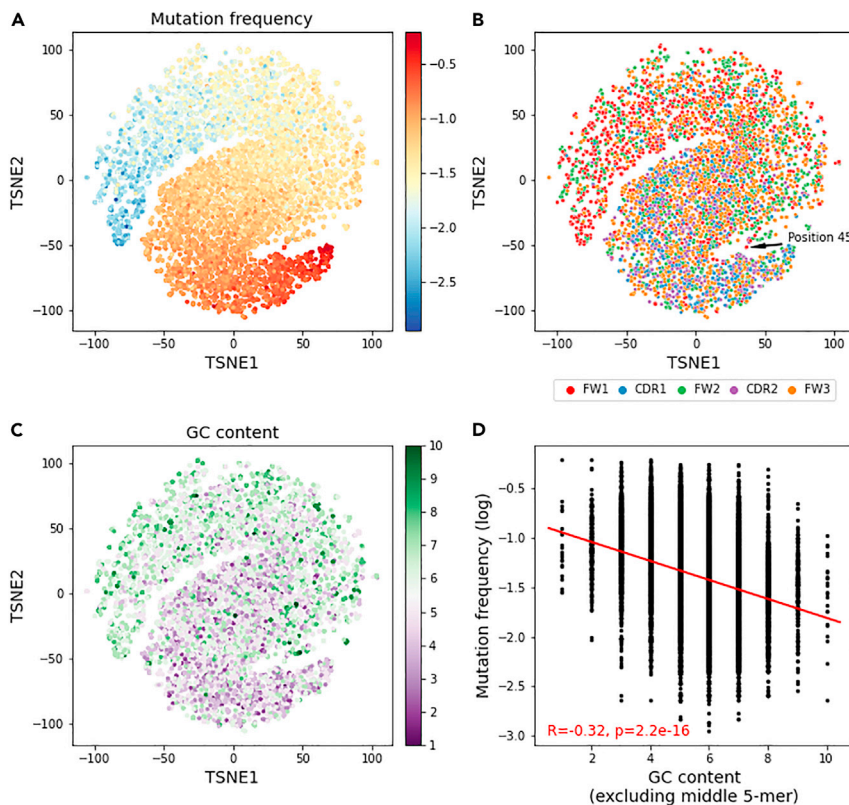


Figure 7. Neural network encodings analysis: mutation frequency model

(A and B) 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 STAR Methods) and is colored according to its corresponding (A) mutation frequency (\log_{10}), and (B) Ig V subregion 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.

methodology of our previous study, we calculated the difference between the predicted G4 potential of the germline with that of the sequence with a single mutation at either position 21 or 45. Here, we found that a T>G mutation at position 21 elevated average G4 potentials to a very high value of 0.83 ± 0.12 compared to a germline value (already relatively high) of 0.55 ± 0.21 , whereas the same mutation occurring at position 45 displayed a far smaller average increase of 0.04 ± 0.03 (Figure 9C, Table 2). As for position 44, there seemed to be little effect of C>G mutations on G4 potential (Table 2). Interestingly, we made another observation regarding the instances where an A nucleotide disrupts the run of G nucleotides at IMGT position 49 (Figure 9A) which was that these sites also displayed a high A>G substitution rate (0.77 ± 0.18 ; Table 2). This hypothetical mutation also caused a moderate, though substantial, increase in G4 potential (0.17 ± 0.02 , Table 2). These findings reveal that particular recurring mutations in this subregion may promote G4 formation, and that the bias toward generating new G sites suggests specific DNA repair enzymes may be recruited to these subregions within FW1.

DISCUSSION

In this study, we leveraged deep learning to gain insights into SHM, a key process in antibody affinity maturation. We trained multiple deep learning models using a convolutional neural network (CNN) framework to analyze DNA *k*-mer subsequences of various lengths, ranging from 5 to 21 nts, derived from human IGHV germline sequences. Using a high-quality data set containing non-productive B cell repertoire data, the model was tasked to learn two focal aspects of SHM: the frequency of mutation at a given site, and the spectrum of mutations that can arise at this site (substitution). Understanding the propensity of a site to mutate and the underlying substitution biases that ensue can lead to a better understanding of how AID

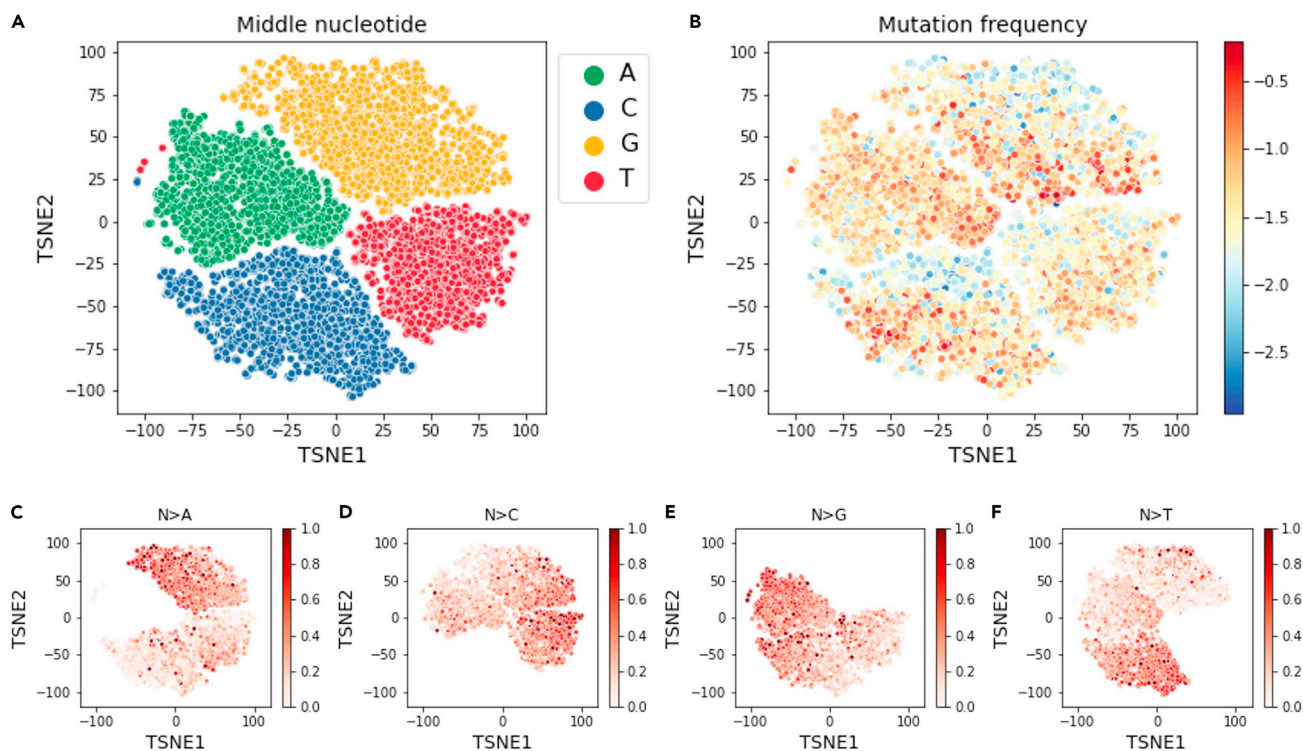


Figure 8. Neural network encodings analysis: substitution model

(A and B) 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 to learn the associated substitution rates (see [STAR Methods](#)) and is colored according to its corresponding (A) middle nucleotide and (B) mutation frequency (\log_{10}).

(C–F) The t-SNE embedding is colored by the rate of substitution for the middle 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.

is recruited to and targets the Ig V region, as well as the associated downstream DNA repair mechanisms that follow AID deamination.

We began by developing three models, collectively referred to as DeepSHM, to predict separate tasks for a given k -mer: observable mutation frequency; distributed substitution rates; and a combination of both measures (weighted substitution). We found that predicting substitution rates did not substantially depend on the k -mer size, whereas 15-mers were optimal for predicting mutation frequencies (Figure 2, Table 1). In addition, DeepSHM predicted both substitution rates and mutation frequencies more accurately than the widely used S5F targeting model for all k -mer sizes we evaluated ($k = 5, 9, 15$ and 21) (Table 1). Even though we were able to outperform S5F in representing substitution biases, the correlation between our predictions and empirical data was moderate (~ 0.56), suggesting that the processes underlying SHM substitution biases may be more fundamentally random than mutational site targeting alone. Error-prone DNA repair processes downstream of AID are highly complex. For example, although Pol η is biased towards making $WA \rightarrow WG$ mutations (Zhang et al., 2014) and plays a dominant role in generating mutations at A:T sites, many A:T mutations still occur in its absence (Saribasak et al., 2009) that are mediated by other polymerases (Maul et al., 2016). Similarly complex, BER is biased towards transversions but can also repair faithfully, with a further dependence on hotspot mutability (Pérez-Durán et al., 2012). Thus, downstream repair processes may simply be too complex, or genuinely random, to be captured well by a model that depends on sequence context alone.

To uncover some of the hidden features learned by DeepSHM, we analyzed the output, or encodings, obtained from the penultimate layer of the network predicting weighted substitution using input 15-mers, and performed t-SNE, a method of dimensionality reduction, to visualize the encodings in two dimensions. The subsequent embedding formed clusters of 15-mers that were distinguished by mutation frequency and middle nucleotide (Figures 3A and 3B). Individual clusters containing a C or G middle nucleotide that

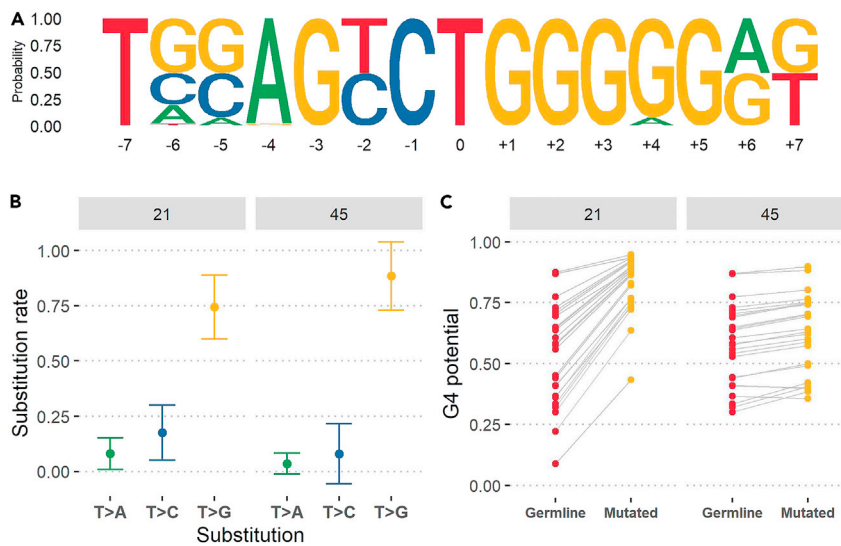


Figure 9. Evaluation of the T outlier cluster in the DeepSHM substitution model

(A) Sequence logo representation of the 15-mers appearing in the T outlier cluster in Figure 8A (left-hand side, red dots). (B) Substitution rates of T>A, T>C, and T>G for 15-mers corresponding to 26 and 25 IGHV3 alleles at IMGT positions 21 and 45, respectively. 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.

were associated with high mutability, assumed to be relevant to AID hotspots, revealed a strong preference for a T base at the +1 position of the top strand AID WRC ($W=A/T$, $R=A/G$) hotspot, including for WAC motifs that are not part of a $WGCW$ motif, and similarly, an A base at the -1 position of the bottom strand GYW ($Y=C/T$) context (Figures 3C and 3D). As an alternative way to identify sequence features, we applied TF-MoDISco (see STAR Methods) to reveal recurrent genomic patterns using importance scores extracted from the model predicting mutation frequency for each 15-mer. This approach supported the importance of the T base at the +1 position of WRC (Figure 4A) and the A base at the -1 position of the bottom strand GYW hotspot (Figure 4B). An early study by Rogozin and Diaz reported the $WRCH/DGYW$ ($H=A/C/T$, $D=A/G/T$) to be a good predictor of mutability at C:G bases (Rogozin and Diaz, 2004); however, we found $WRCT$ to be a more consistent definition (Figure 5). The authors of the SSF model also supported the $WRCH$ definition because they found their model can capture the higher mutability rate seen at certain $WRCA$ motifs (Yaari et al., 2013), presumably at the $AGCA$ overlapping hotspot. Previous hotspot definitions have largely failed to describe targeting beyond the -2 position of the WRC motif. We further identified having a C or G at the -3 position of WRC , or at the +3 position of GYW , as a strong negative contribution, i.e., as a reduced effect on targeting. Thus, our results suggest the typical AID hotspot definition might be extended to $WWRCT$. Comparing the mutation frequencies of the individual $WWRCT$ hotspot motifs showed the 3' T to be more important for AID recognition than the 5' W alone, however, together they have a synergistic effect that makes mutability between 2.2-fold (for TAC) and 4.3-fold (for TGC) higher (Figure 5C, Table S3).

We next applied the same t-SNE methodology on the two developed standalone models that separately predicted either the mutation frequency or substitution rates of the 15-mer middle nucleotides. The t-SNE embedding on the independent DeepSHM model predicting only mutation frequency revealed a significant negative correlation between the mutability of a site and the surrounding GC content of the 15-mer (Figure 7D). This finding alternatively suggests that highly mutated sites may have evolved to have a richer local AT content. This *in vivo* result is consistent with earlier *in vitro* results that considered AID targeting artificial substrates (Abdouni et al., 2018).

On the other hand, the t-SNE embedding stemming from the standalone substitution model hinted at plausible associations, both positive and negative, between mutation frequency and certain transition and transversion mutations (Figures 8B–8F). We next analyzed multiple genes representing different

Table 2. Summary statistics on outlier C and T clusters

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)
21	T	26	0.74 ± 0.15	0.05 ± 0.02	0.55 ± 0.21	0.83 ± 0.12	0.28 ± 0.11
44	C	21	0.87 ± 0.11	0.03 ± 0.01	0.58 ± 0.16	0.62 ± 0.16	0.04 ± 0.03
45	T	25	0.88 ± 0.08	0.17 ± 0.08	0.58 ± 0.16	0.62 ± 0.16	0.04 ± 0.03
49	A	4	0.77 ± 0.18	0.03 ± 0.04	0.35 ± 0.06	0.52 ± 0.05	0.17 ± 0.02

IGHV families containing the largest amounts of mutation data in order to avoid any potential sites with few observable mutations, such as coldspots. We observed a negative correlation between mutability and substitution rates specifically for C>A and G>T transversion mutations (Figures 8 and S6) and, on the other hand, positive correlations for C>T and G>A transitions (Figure S6). The trend for increased transition mutations at highly mutating AID hotspots mediated by UNG2 had previously been observed in experiments using 3T3 (mouse fibroblast) cells (Pérez-Durán et al., 2012), although the particular bias against C:G>A:T transversions was not apparent. Previous work has also shown that UNG2 is cell-cycle regulated, possibly mediated by FAM72A (Feng et al., 2021), and active primarily during G1 (Sharbeen et al., 2012). Although AID is also primarily active during G1, it may sometimes persist for slightly longer than UNG2 and thus highly targeted sites may avoid BER especially when the mutations occur just before the cell enters S phase, which would lead to fixation of C>T transitions via replication bypass. Alternative polymerases may also be preferentially recruited to some sites. For example, in DT40 (chicken) B cell lines, the POLD3 subunit of Polymerase delta (Polδ) has been proposed as a specific mechanism for both C>A and G>T mutations (Hirota et al., 2015; Pilzecker and Jacobs, 2019).

In addition, we investigated two outlier clusters from the substitution model embedding that contained 15-mers having a C and T middle nucleotide that did not group with their respective larger clusters (Figure 8A). A closer analysis revealed that the T outlier contained a highly conserved AGYCTGGGGG consensus sequence that was derived from two independent sites located in FW1 from multiple IGHV3 alleles (Figure 9A, Table 2). Both outlier clusters also displayed significantly elevated T>G (Figure 9B, Table 2) and C>G substitution rates (Figure S7B, Table 2), respectively. In our recent study on G-quadruplexes (G4s) in IGHV genes, we observed the IGHV3 family to form G4s more favorably on the top strand, as measured by their predicted G4 potential using a pre-trained CNN model (Tang and MacCarthy, 2021). Given the strong preference for creating G mutations in these FW1 subregions, we evaluated the impact of these mutations on G4 potential. In some cases, the resulting G mutation led to a strong increase in G4 potential, particularly at position 21 (Figure 9C, Table 2), whereas for other sites, the effect was mostly negligible (Table 2). Notably, however, a high A>G substitution rate was also observed at the +4 positions (Figure 9A), which were also associated with increases in G4 potential (Table 2). These biased A>G mutations may further be related to previous work that found that a repeated mutation that occurs in one IGHV allele often matches the sequence variant of a different allele (Saini and Hershberg, 2015). Alternatively, these mutations may be related to R-loop initiation, which forms in G-rich non-template DNA, possibly forming in FW1 of these IGHV3 genes. Studies have found that reducing G-density in mammalian Ig switch regions compromises class-switch recombination efficiency and R-loops from forming (Roy et al., 2008; Zhang et al., 2014). The high rate of T>G and C>G transversions also suggests that particular repair enzymes may be recruited to these subregions during SHM.

Limitations of the study

In principle, a wider range of *k*-mers, as well as a greater variety of neural network architectures, might have been considered for this study. However, because the tuning of each model takes a substantial amount of computational resources and time, we considered a reduced number of models. In addition, we limited this study to consider data only for humans, the species for which we had high quality (UMI barcoded) data in high abundance, although the approach could be extended to other species such as mice in future work.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE

- **RESOURCE AVAILABILITY**
 - Lead contact
 - Materials availability
 - Data and code availability
- **METHODS DETAILS**
 - Generating *k*-mer data
 - Calculating mutation frequencies, substitution rates, and weighted substitutions of *k*-mers
 - CNN architecture and model optimization
 - Inferring an S5F targeting model
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
 - Statistical tests
 - Neural network encodings analysis
 - Cluster identification
 - Identifying recurring genomic patterns using TF-ModISco

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2021.103668>.

ACKNOWLEDGMENTS

We would like to extend our gratitude to Sergio Roa Gómez for his comments and suggestions. This work was supported by NIH grant A1132507.

AUTHOR CONTRIBUTIONS

C.T., A.K., and T.M. conceived the idea, analyzed the results, and wrote the manuscript. C.T. and A.K. developed the model and performed computational analysis. All authors contributed to the article and approved the submitted version.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: August 4, 2021

Revised: November 8, 2021

Accepted: December 16, 2021

Published: January 21, 2022

REFERENCES

- Abdouni, H.S., King, J.J., Ghorbani, A., Fifield, H., Berghuis, L., and Larjani, M. (2018). DNA/RNA hybrid substrates modulate the catalytic activity of purified AID. *Mol. Immunol.* *93*, 94–106.
- Alipanahi, B., DeLong, A., Weirauch, M.T., and Frey, B.J. (2015). Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat. Biotechnol.* *33*, 831–838.
- Álvarez-Prado Á, F., Pérez-Durán, P., Pérez-García, A., Benguria, A., Torroja, C., de Yébenes, V.G., and Ramiro, A.R. (2018). A broad atlas of somatic hypermutation allows prediction of activation-induced deaminase targets. *J. Exp. Med.* *215*, 761–771.
- Bransteitter, R., Pham, P., Scharff, M.D., and Goodman, M.F. (2003). Activation-induced cytidine deaminase deaminates deoxycytidine on single-stranded DNA but requires the action of RNase. *Proc. Natl. Acad. Sci. U S A* *100*, 4102–4107.
- Cohen, R.M., Kleinstein, S.H., and Louzoun, Y. (2011). Somatic hypermutation targeting is influenced by location within the immunoglobulin V region. *Mol. Immunol.* *48*, 1477–1483.
- Cui, A., Di Niro, R., Vander Heiden, J.A., Briggs, A.W., Adams, K., Gilbert, T., O'Connor, K.C., Vigneault, F., Shlomchik, M.J., and Kleinstein, S.H. (2016). A model of somatic hypermutation targeting in mice based on high-throughput Ig sequencing data. *J. Immunol.* *197*, 3566–3574.
- Elhanati, Y., Sethna, Z., Marcou, Q., Callan, C.G., Jr., Mora, T., and Walczak, A.M. (2015). Inferring processes underlying B-cell repertoire diversity. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* *370*, 20140243.
- Feng, Y., Li, C., Stewart, J., Barbulescu, P., Desivo, N.S., Álvarez-Quilón, A., Pezo, R.C., Perera, M.L.W., Chan, K., Tong, A.H.Y., et al. (2021). FAM72A antagonizes UNG2 to promote mutagenic uracil repair during antibody maturation. *Nature* *600*, 324–328. <https://doi.org/10.1038/s41586-021-04144-4>.
- Haynes, B.F., Kelsoe, G., Harrison, S.C., and Kepler, T.B. (2012). B-cell-lineage immunogen design in vaccine development with HIV-1 as a case study. *Nat. Biotechnol.* *30*, 423–433.
- Hirota, K., Yoshikiyo, K., Guilbaud, G., Tsurimoto, T., Murai, J., Tsuda, M., Phillips, L.G., Narita, T., Nishihara, K., Kobayashi, K., et al. (2015). The POLD3 subunit of DNA polymerase δ can promote translesion synthesis independently of DNA polymerase ζ . *Nucleic Acids Res.* *43*, 1671–1683.
- Jansen, J.G., Langerak, P., Tsaalbi-Shtylik, A., van den Berk, P., Jacobs, H., and de Wind, N. (2006). Strand-biased defect in C/G transversions in hypermutating immunoglobulin genes in Rev1-deficient mice. *J. Exp. Med.* *203*, 319–323.
- Kelley, D.R., Snoek, J., and Rinn, J.L. (2016). Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* *26*, 990–999.
- Koo, P.K., and Ploenzke, M. (2020). Deep learning for inferring transcription factor binding sites. *Curr. Opin. Syst. Biol.* *19*, 16–23.

- Krantsevich, A., Tang, C., and MacCarthy, T. (2021). Correlations in somatic hypermutation between sites in IGHV genes can be explained by interactions between AID and/or pol η hotspots. *Front. Immunol.* *11*, 3751.
- Lefranc, M.P. (2001). IMGT, the international ImMunoGeneTics database. *Nucleic Acids Res.* *29*, 207–209.
- Liu, M., Duke, J.L., Richter, D.J., Vinuesa, C.G., Goodnow, C.C., Kleinstein, S.H., and Schatz, D.G. (2008). Two levels of protection for the B cell genome during somatic hypermutation. *Nature* *451*, 841–845.
- Matsuda, T., Bebenek, K., Masutani, C., Rogozin, I.B., Hanaoka, F., and Kunkel, T.A. (2001). Error rate and specificity of human and murine DNA polymerase ϵ . *J. Mol. Biol.* *312*, 335–346.
- Maul, R.W., MacCarthy, T., Frank, E.G., Donigan, K.A., McLenigan, M.P., Yang, W., Saribasak, H., Huston, D.E., Lange, S.S., Woodgate, R., et al. (2016). DNA polymerase ϵ functions in the generation of tandem mutations during somatic hypermutation of antibody genes. *J. Exp. Med.* *213*, 1675–1683.
- Mayorov, V.I., Rogozin, I.B., Adkison, L.R., and Gearhart, P.J. (2005). DNA polymerase ϵ contributes to strand bias of mutations of A versus T in immunoglobulin genes. *J. Immunol.* *174*, 7781–7786.
- Methot, S.P., and Di Noia, J.M. (2017). Molecular mechanisms of somatic hypermutation and class switch recombination. *Adv. Immunol.* *133*, 37–87.
- Muramatsu, M., Kinoshita, K., Fagarasan, S., Yamada, S., Shinkai, Y., and Honjo, T. (2000). Class switch recombination and hypermutation require activation-induced cytidine deaminase (AID), a potential RNA editing enzyme. *Cell* *102*, 553–563.
- Ohm-Laursen, L., and Barington, T. (2007). Analysis of 6912 unselected somatic hypermutations in human VDJ rearrangements reveals lack of strand specificity and correlation between phase II substitution rates and distance to the nearest 3' activation-induced cytidine deaminase target. *J. Immunol.* *178*, 4322–4334.
- Pérez-Durán, P., Belver, L., de Yébenes, V.G., Delgado, P., Pisano, D.G., and Ramiro, A.R. (2012). UNG shapes the specificity of AID-induced somatic hypermutation. *J. Exp. Med.* *209*, 1379–1389.
- Pham, P., Bransteitter, R., Petruska, J., and Goodman, M.F. (2003). Processive AID-catalysed cytosine deamination on single-stranded DNA simulates somatic hypermutation. *Nature* *424*, 103–107.
- Pilzecker, B., and Jacobs, H. (2019). Mutating for good: DNA damage responses during somatic hypermutation. *Front Immunol.* *10*, 438.
- Rada, C., Di Noia, J.M., and Neuberger, M.S. (2004). Mismatch recognition and uracil excision provide complementary paths to both Ig switching and the A/T-focused phase of somatic mutation. *Mol. Cell* *16*, 163–171.
- Rajewsky, K. (1996). Clonal selection and learning in the antibody system. *Nature* *381*, 751–758.
- Rogozin, I.B., and Diaz, M. (2004). Cutting edge: DGYW/WRCH is a better predictor of mutability at G:C bases in Ig hypermutation than the widely accepted RGYW/WRCY motif and probably reflects a two-step activation-induced cytidine deaminase-triggered process. *J. Immunol.* *172*, 3382–3384.
- Rogozin, I.B., and Kolchanov, N.A. (1992). Somatic hypermutagenesis in immunoglobulin genes. II. Influence of neighbouring base sequences on mutagenesis. *Biochim. Biophys. Acta* *1171*, 11–18.
- Roy, D., Yu, K., and Lieber, M.R. (2008). Mechanism of R-loop formation at immunoglobulin class switch sequences. *Mol. Cell Biol* *28*, 50–60.
- Saini, J., and Hershberg, U. (2015). B cell variable genes have evolved their codon usage to focus the targeted patterns of somatic mutation on the complementarity determining regions. *Mol. Immunol.* *65*, 157–167.
- Saribasak, H., Rajagopal, D., Maul, R.W., and Gearhart, P.J. (2009). Hijacked DNA repair proteins and unchained DNA polymerases. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* *364*, 605–611.
- Shapiro, G.S., Aviszus, K., Ikle, D., and Wysocki, L.J. (1999). Predicting regional mutability in antibody V genes based solely on di- and trinucleotide sequence composition. *J. Immunol.* *163*, 259–268.
- Shapiro, G.S., Aviszus, K., Murphy, J., and Wysocki, L.J. (2002). Evolution of Ig DNA sequence to target specific base positions within codons for somatic hypermutation. *J. Immunol.* *168*, 2302–2306.
- Sharbeen, G., Yee, C.W., Smith, A.L., and Jolly, C.J. (2012). Ectopic restriction of DNA repair reveals that UNG2 excises AID-induced uracils predominantly or exclusively during G1 phase. *J. Exp. Med.* *209*, 965–974.
- Sheng, Z., Schramm, C.A., Kong, R., Program, N.C.S., Mullikin, J.C., Mascola, J.R., Kwong, P.D., and Shapiro, L. (2017). Gene-specific substitution profiles describe the types and frequencies of amino acid changes during antibody somatic hypermutation. *Front Immunol.* *8*, 537.
- Shrikumar, A., Tian, K., Shcherbina, A., Avsec, Z., Banerjee, A., Sharmin, M., Nair, S., and Kundaje, A. (2018). Technical note on transcription factor motif discovery from importance scores (TF-MoDISco). *bioRxiv*. <https://arxiv.org/abs/1811.00416>.
- Spisak, N., Walczak, A.M., and Mora, T. (2020). Learning the heterogeneous hypermutation landscape of immunoglobulins from high-throughput repertoire data. *Nucleic Acids Res.* *48*, 10702–10712.
- Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic Attribution for Deep Networks, *arXiv:1703.01365*.
- Tang, C., Bagnara, D., Chiorazzi, N., Scharff, M.D., and MacCarthy, T. (2020). AID overlapping and poleta hotspots are key features of evolutionary variation within the human antibody heavy chain (IGHV) genes. *Front Immunol.* *11*, 788.
- Tang, C., and MacCarthy, T. (2021). Characterization of DNA G-quadruplex structures in human immunoglobulin heavy variable (IGHV) genes. *Front. Immunol.* *12*, 671944.
- Wei, L., Chahwan, R., Wang, S., Wang, X., Pham, P.T., Goodman, M.F., Bergman, A., Scharff, M.D., and MacCarthy, T. (2015). Overlapping hotspots in CDRs are critical sites for V region diversification. *Proc. Natl. Acad. Sci. U S A* *112*, E728–E737.
- Wiehe, K., Bradley, T., Meyerhoff, R.R., Hart, C., Williams, W.B., Easterhoff, D., Faison, W.J., Kepler, T.B., Saunders, K.O., Alam, S.M., et al. (2018). Functional relevance of improbable antibody mutations for HIV broadly neutralizing antibody development. *Cell Host Microbe* *23*, 759–765.e6.
- Yaari, G., Vander Heiden, J.A., Uduman, M., Gadala-Maria, D., Gupta, N., Stern, J.N., O'Connor, K.C., Hafler, D.A., Laserson, U., Vigneault, F., et al. (2013). Models of somatic hypermutation targeting and substitution based on synonymous mutations from high-throughput immunoglobulin sequencing data. *Front Immunol.* *4*, 358.
- Yu, K., Huang, F.T., and Lieber, M.R. (2004). DNA substrate length and surrounding sequence affect the activation-induced deaminase activity at cytidine. *J. Biol. Chem.* *279*, 6496–6500.
- Zhang, Z.Z., Pannunzio, N.R., Hsieh, C.L., Yu, K., and Lieber, M.R. (2014). The role of G-density in switch region repeats for immunoglobulin class switch recombination. *Nucleic Acids Res.* *42*, 13186–13193.
- Zhou, J., and Troyanskaya, O.G. (2015). Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods* *12*, 931–934.
- Zhou, J.Q., and Kleinstein, S.H. (2020). Position-dependent differential targeting of somatic hypermutation. *J. Immunol.* *ij2000496*.

STAR★METHODS

KEY RESOURCES TABLE

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: PRJNA381394, PRJNA591804)
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/

RESOURCE AVAILABILITY

Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by the lead contact, Thomas MacCarthy (thomas.maccarthy@stonybrook.edu).

Materials availability

This study did not generate new unique reagents.

Data and code availability

- Data used for this research was published previously by Tang et al. (2020).
- A custom Python package developed for this project is available at <https://gitlab.com/maccarthyslab/deepshm>.
- Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

METHODS DETAILS

Generating *k*-mer data

Germline IGHV reference sequences were downloaded from the international ImMunoGeneTics information system (IMGT) website (Lefranc, 2001). The leader portion of each reference sequence was also extracted if available. To generate the *k*-mers of a given germline sequence, $\pm[k/2]$ nt sequences were extracted from the start of the V exon, where *k* is the length of the subsequence, and $\lfloor k/2 \rfloor$ represents the greatest integer less than or equal to $k/2$. This process was continued, moving 1 nt at a time, until the end of the exon was reached. Next, all *k*-mers were converted to their respective one-hot encodings. A one-hot encoding is a transformation of a DNA sequence using a 2-D matrix containing only zeros and ones, where each row represents one of the four ordered DNA bases and each column is an individual site in the sequence. For each column, a "1" is filled in the row that matches the nucleotide of that site and a "0" in the remaining unmatched rows (Figure 1).

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

We used a high-quality dataset previously published by us (Tang et al., 2020). This dataset contains only non-productive sequences that contain a frameshift in CDR3 to minimize selection effects (Ohm-Laursen and Barington, 2007) and which were derived from marginal, memory, and plasma B cell types. To avoid any mutational artifacts due to B cell clonal relatedness, we sampled one sequence per clonal group for each IGHV allele. We calculated the mutation frequencies of every *k*-mer in a germline sequence as the number of observed mutations at each site (corresponding to a single *k*-mer), divided by the total number of sequences for that germline IGHV allele. The substitution rate of each *k*-mer was computed as the

number of times the middle nucleotide mutated from the germline nucleotide to the other three DNA bases, divided by the total number of overall mutations. Note that a zero was recorded if the mutated base was the same as germline. Lastly, the weighted substitution of a k -mer was calculated as the observed mutation frequency multiplied by the substitution rate vector. The data used for this study can be accessed via the NCBI SRA (BioProject: PRJNA381394, PRJNA591804).

CNN architecture and model optimization

We implemented a convolutional neural network (CNN) to analyze the k -mer input data. Three separate architectures were used to predict different SHM outcomes: mutation frequency, substitution rate, and weighted substitution (see above). Although the hyperparameters that were ultimately selected varied from model-to-model, all CNNs followed the same general architecture, which consisted of one convolution layer, followed by two fully connected layers (Figure 1). Additional parameters, such as dropout and batch normalization, were optimized by generating 100 separate models with randomly selected hyperparameters for each k -mer and corresponding model architecture we generated. The range of values for all parameters and hyperparameters that were tested for each architecture and output type are specified in Table S1A, as well as the values used for the best 15-mer models.

Next, we utilized 4-fold cross-validation to evaluate the performance of the model on unseen (test) data. In total, there are seven IGHV families (IGHV1-7), where each IGHV family consists of genes that share a high percentage of sequence similarity (Lefranc, 2001). The k -mers derived from the three largest IGHV families, IGHV1, IGHV3, and IGHV4, formed three separate groups, and the k -mers belonging to the remaining 4 smaller IGHV families constituted the final group in order to create a data set comparable in size with the other groups (Table S1B). Thus, we separated the data by their respective IGHV family to reduce the chances of model overfitting, since it is likely that k -mers from the same IGHV family will be similar even if they come from different genes and, therefore, bias the results if they appear in both training and test sets. In every cross-validation fold, three of the data groups were used as training set, and the fourth used as test set. We also evaluated the model performance, for each fold, by calculating the Pearson correlation(s) between the predicted mutation frequency and/or substitution rate of the test set k -mers and the equivalent output type of the empirical data. The average correlation across the 4 validation folds was reported for the model, as in Figure 2.

As an additional step, we wrote a custom, universal Python script (available at <https://gitlab.com/maccarthyslab/deepshm>) to automatically generate the CNN architecture, parameters, and hyperparameters of each model, regardless of the output specified, to ensure that all models were constructed in a consistent manner. All CNNs were generated using the built-in Keras API in Tensorflow 2.4.1 and trained on GPU processors using three Nvidia GeForce RTX 2080 graphics cards.

Inferring an S5F targeting model

In order to ensure a fair comparison between S5F values and our deep learning predictions, we used the SHazaM R package (Yaari et al., 2013) to create an S5F targeting model, which provides analogous 5-mer mutability and substitution scores based on the same data set we used to train our CNN models with. We specified the S5F targeting model to count both silent and replacement mutations ("rs" parameter) since the mutation data we used was derived from non-functionally rearranged VDJ coding sequences (i.e. in the absence of selection) and with each sequence being clonally independent (see above). Multiple mutations were handled specifying the "independent" parameter, which treats each mutation independently. Default values were used for all other parameters.

QUANTIFICATION AND STATISTICAL ANALYSIS

Statistical tests

Since training of Deep Learning models entails some stochasticity, we used the Wilcoxon signed-rank test to compare the overall performance of DeepSHM against S5F, which we used as a benchmark. Thus, for each value of k (5, 9, 15 or 21) and each type of model (mutation frequency, substitution, or weighted substitution) the performance of 100 models with randomly generated hyperparameters was compared against the corresponding S5F accuracy. In order to compare the performance of the models using different k -mers, we used a one-sided Mann-Whitney U test.

We also used a one-sided Mann-Whitney U test to compare the specificity of different potential refined AID and Pol η hotspot motifs. Note that in [Figures 5](#) and [6](#), asterisks indicate significance ($p \leq 0.0001$) of a one-sided Mann-Whitney U test.

Neural network encodings analysis

The output (encoding) of the penultimate layer of the CNN model was used as a way to explain the SHM patterns learned by the model. To generate the encodings from this layer, we removed the last layer of the CNN while keeping the remaining layers intact. Next, we processed the k -mers through the truncated model to retrieve the ensuing output values. We then applied t-distributed stochastic neighbor embedding (t-SNE) in Python on these multi-dimensional encodings to visually represent the resulting embedding in two dimensions.

Cluster identification

We implemented partitioning around medoids (PAM) clustering to identify clusters within the t-SNE embedding of the weighted substitution model ([Figure S1](#)). The goal of PAM clustering is to minimize the sum of dissimilarities between the objects in a cluster and the center of the same cluster (medoid). Pairwise dissimilarities between k -mers were computed using Gower's distance as a way to group k -mers with similar mutation frequencies. We separated all k -mers sharing the same middle nucleotide and then applied PAM clustering independently on each group to facilitate the clustering process. All clustering assignments were performed using the *cluster* package in R. For each middle nucleotide, we specified the algorithm to identify 5 distinct clusters.

Identifying recurring genomic patterns using TF-MoDISco

We applied TF-MoDISco ([Shrikumar et al., 2018](#)), a machine learning interpretability method, to identify recurring motifs our model detected in the 15-mer data. From the data, we isolated four groups of 15-mers based on the middle nucleotide (C, G, A, or T) of the 15-mer, with the additional condition that the middle nucleotide conformed to WRC or GYW AID hotspots, or WA or TW Pol η hotspots, respectively. TF-MoDISco requires importance scores to be used as input, which can be generated by utilizing one of several attribution methods. Here we generated the importance scores for each group by applying Integrated Gradients ([Sundararajan et al., 2017](#)) to the most accurate 15-mer mutation frequency model. Using the resulting importance scores, we then ran TF-MoDISco for all groups separately, still subject to the hotspot constraint, and requiring each of the identified patterns to be associated with at least 10 input subsequences (or "sequelets").