

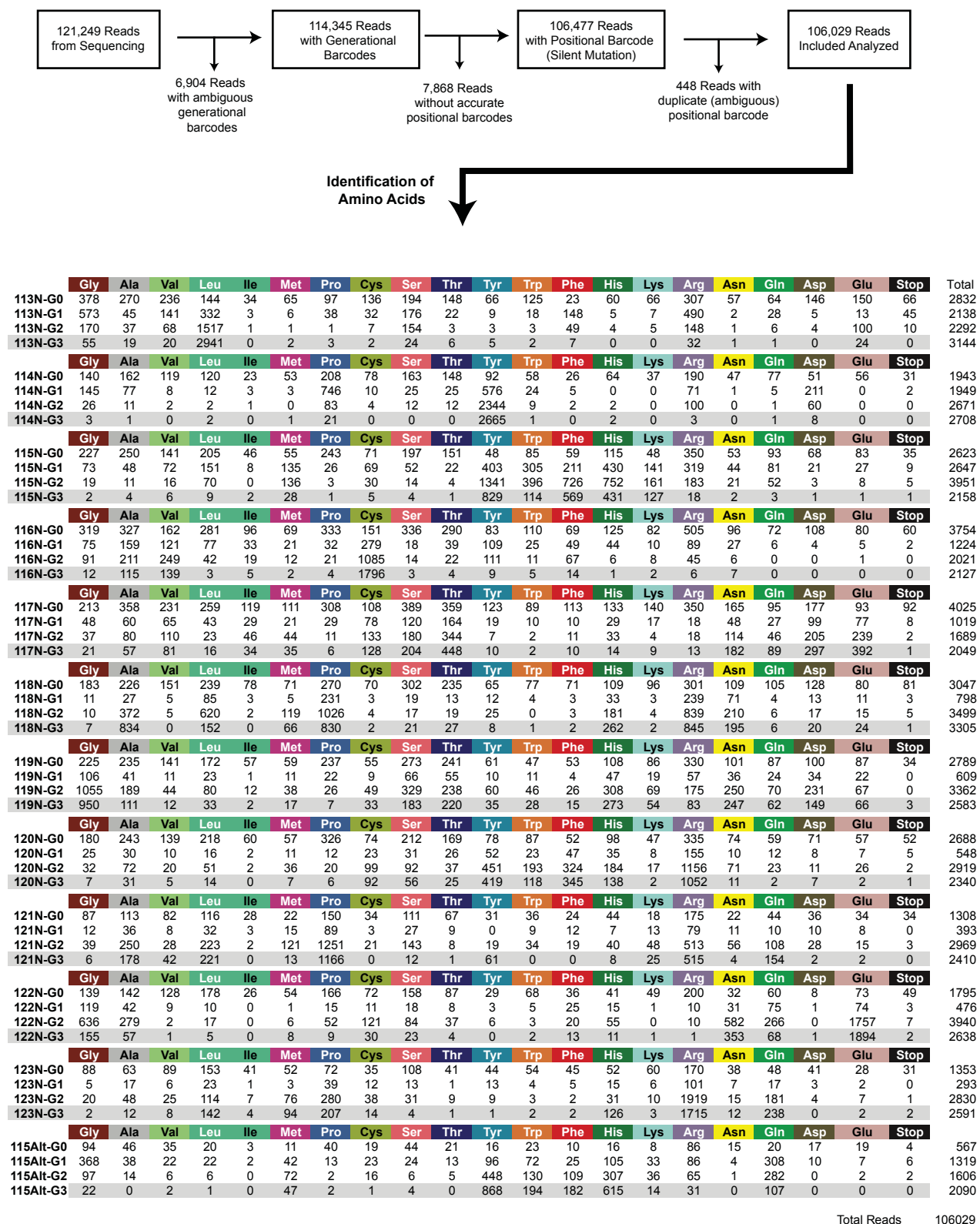
High-throughput mutagenesis reveals functional determinants for DNA targeting by Activation-Induced Deaminase

Kiran S. Gajula, Peter J. Huwe, Charlie Y. Mo, Daniel J. Crawford, James T. Stivers, Ravi Radhakrishnan and Rahul M. Kohli

SUPPLEMENTARY INFORMATION

Supplementary Figures S1-S6, Supplementary Tables S1-S2, Supplementary Video S1

SUPPLEMENTARY FIGURES



Total Reads 106029

Figure S1. Deconvolution of Sequencing Data from Sat-Sel-Seq. Shown is the total number of reads from 454 sequencing, filtered by the presence of barcodes for the generation number, by internal barcodes encoding for position number and eliminating any variants with redundant barcodes. The remaining reads were binned by the two barcodes and the identity of the codon at the mutagenized position was cataloged. The total number of reads for each condition are shown, with the relative frequencies of each amino acid reported in Figure S4.

SL1	SL5
<pre> Y A G R P R L Y F C E D R K A E P TTA TAT TWT TGC GAA GMT SGG ARA SCC GAG CSC </pre>	<pre> Y A G R R L Y F C E D R K R E P TTA TAT TWT TGC GAA GMT SGG ARA CGG GAG CSC </pre>
SL2	SL6
<pre> Y A G R P R L Y F C T D R K A E P TTA TAT TWT TGC ACC GMT SGG ARA SCC GAG CSC </pre>	<pre> Y A G R R L Y F C T D R K R E P TTA TAT TWT TGC ACC GMT SGG ARA CGG GAG CSC </pre>
SL3	SL7
<pre> Y R G R P R L Y F C E P R K A E P TTA TAT TWT TGC GAA CSC SGG ARA SCC GAG CSC </pre>	<pre> Y R G R R L Y F C E P R K R E P TTA TAT TWT TGC GAA CSC SGG ARA CGG GAG CSC </pre>
SL4	SL8
<pre> Y R G R P R L Y F C T P R K A E P TTA TAT TWT TGC ACC CSC SGG ARA SCC GAG CSC </pre>	<pre> Y R G R R L Y F C T P R K R E P TTA TAT TWT TGC ACC CSC SGG ARA CGG GAG CSC </pre>

Figure S2. Sub-libraries for covariation selection. Using the oligonucleotides shown, eight different sub-libraries were created. These when pooled in the ratio of 2:2:2:2:1:1:1:1 generate a starting library that contain equal amount of each of the 384-library family members.

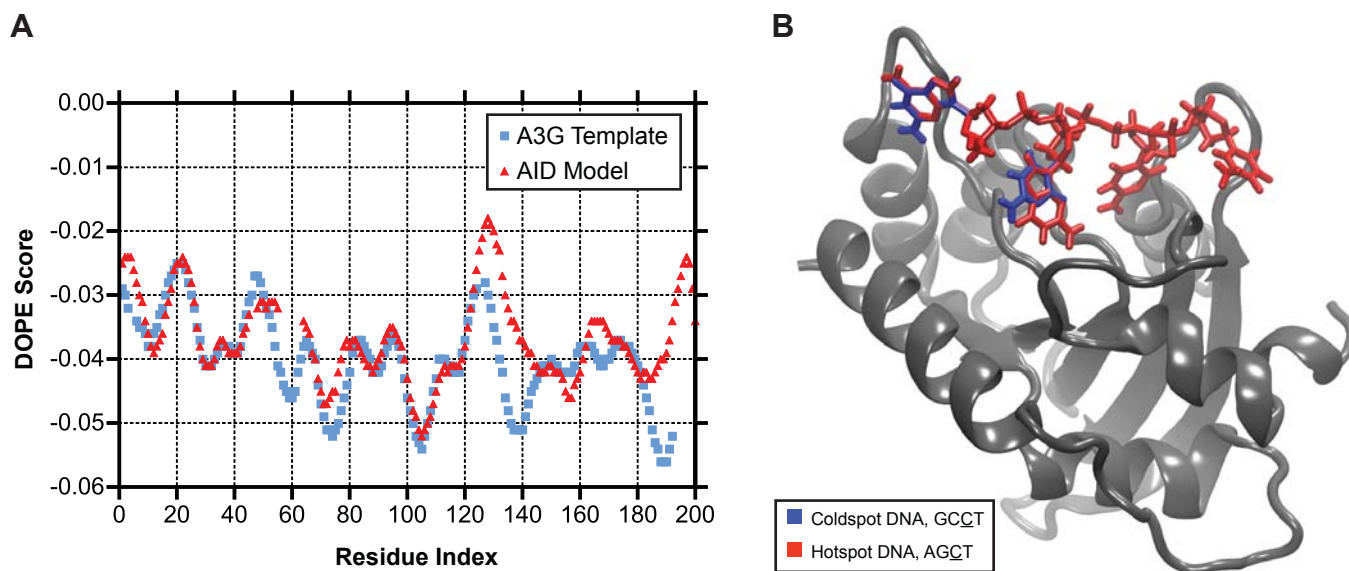


Figure S3. AID Structural Model (**A**) The Discrete Optimized Protein Energy (DOPE) profiles for the template (A3G) and selected target (AID homology) structures are shown. The selected AID homology structure in red showed good fit to the template structure in blue. The selected model was further refined with extensive MD simulations. (**B**) Model of AID bound to DNA. Shown is the homology model of AID(1-181), based on the structure of A3G (PDB 3IQS), with bound ssDNA containing either a hotspot (AGCT) or coldspot (GCCT) sequence, colored red and blue respectively. The alternations in nucleobase composition were done after pre-equilibration of the model, making the starting point for MD simulations identical (RMSD of 0 Å between AID in two structures).

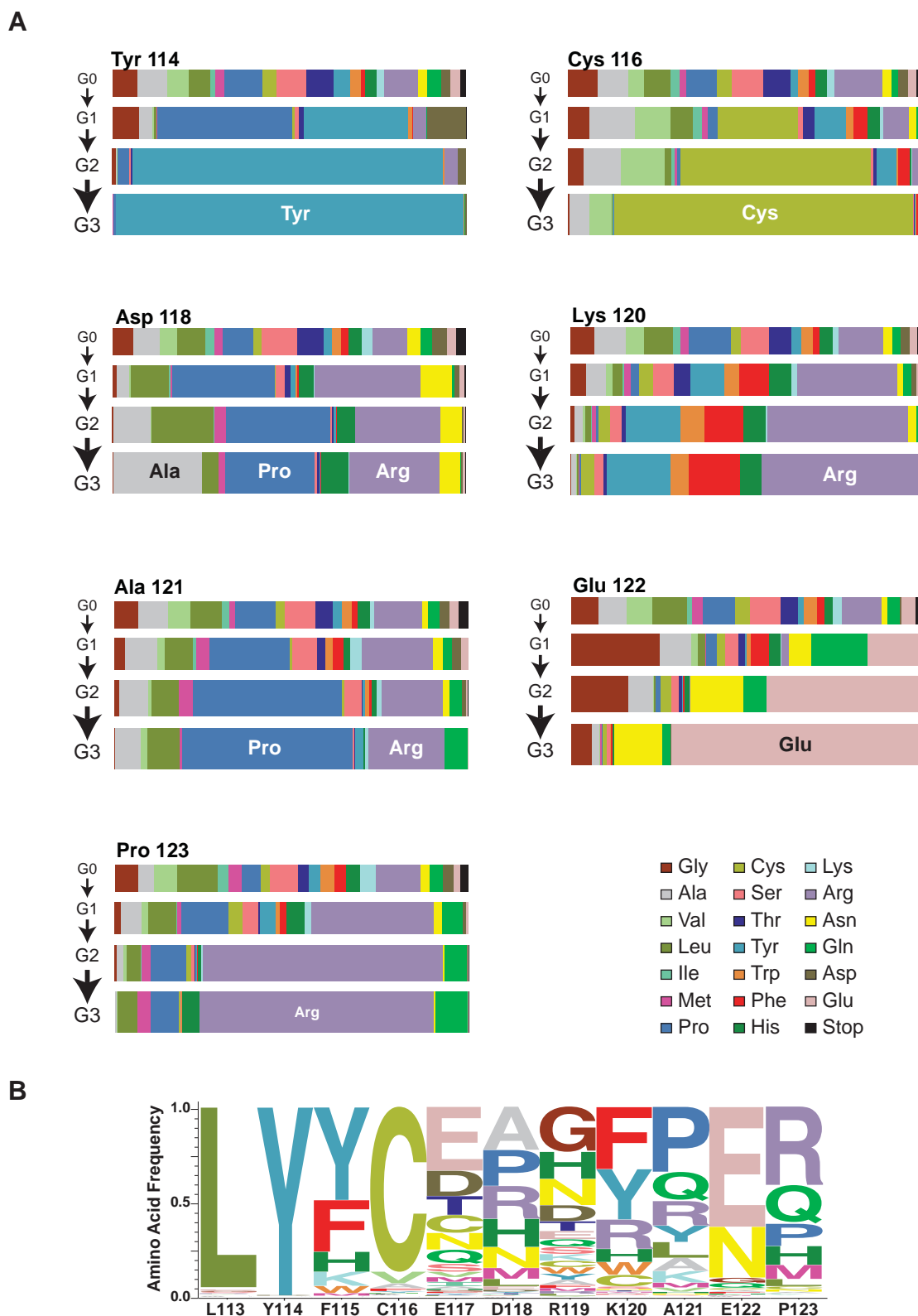
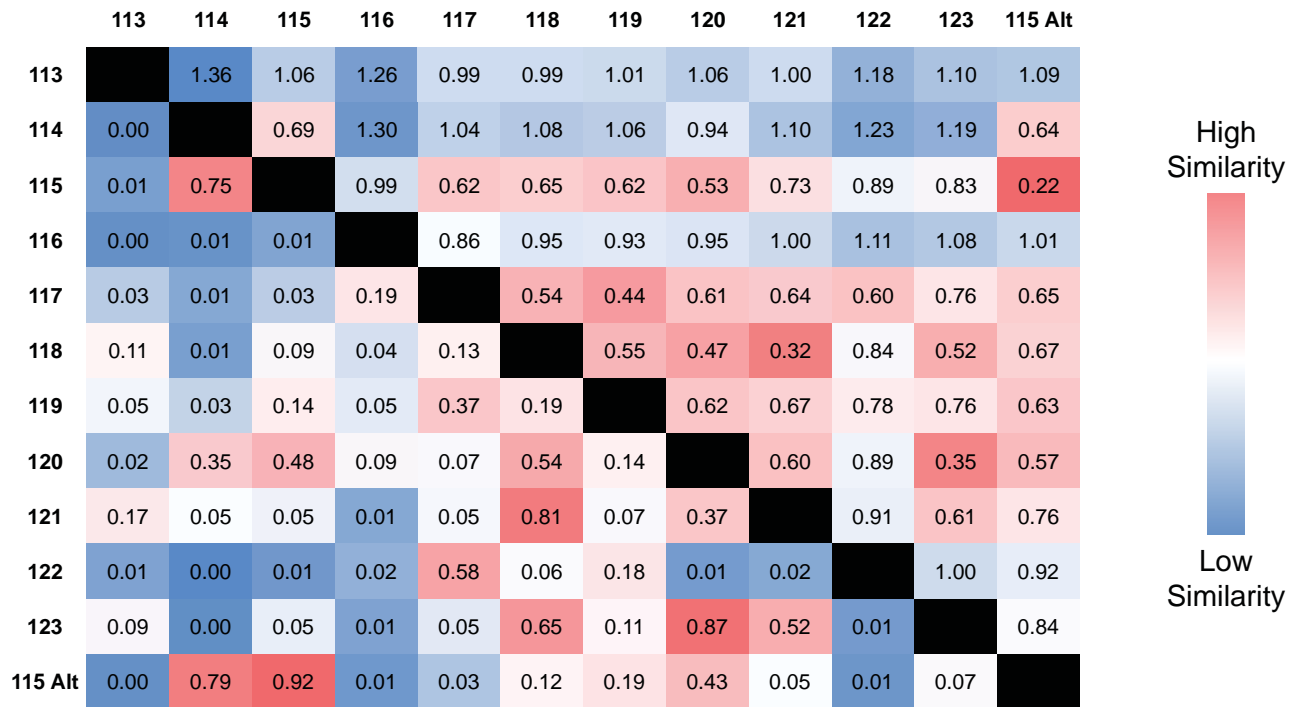


Figure S4. Functional determinants across the entire Leu113-Pro123 loop revealed by Sat-Sel-Seq. **(A)** Horizontal bar plots depicting the prevalence of amino acids from the four generational libraries, G0 to G3, for the positions not depicted in Figure 3. **(B)** Shown is a normalized logo plot for residues 113-123 summarizing the frequency of each amino acid at each position in the final selected library, adjusted for the prevalence of each amino acid in the G0 library. The corrected counts for each amino acid X in G3 were made by $G3_{Cor}^X = G3^X (G3^X/G3^{total}) / (G0^X/G0^{total})$.

Euclidean Distance



Cosine Similarity

$F_{n,x}$ = Frequency of amino acid n in G3 Library at Position x

Euclidean Distance between Position x and Position y results

$$\sqrt{\sum_{n = \text{all amino acids}} (F_{n,x} - F_{n,y})^2}$$

Cosine Similarity between Position x and Position y results

$$\frac{\sum_{n = \text{all amino acids}} (F_{n,x} * F_{n,y})}{\sqrt{\sum_{n = \text{all amino acids}} (F_{n,x})^2} * \sqrt{\sum_{n = \text{all amino acids}} (F_{n,y})^2}}$$

Figure S5. Similarity measures between positions. For each starting library, the G3 library was converted into a vector encoding the frequency ($F_{n,x}$) of each of the twenty amino acids (n) at that position (x). The distance between the library x and library y vectors was determined by two different metrics. Cosine similarity ranges between 0 (different) to 1 (identical). For the Euclidean distance, identical sequences have a metric of 0, while a larger metric means data are more dissimilar. The upper right of the table represents the Euclidean distance metric while the lower left of the table represents the cosine similarity, and entries are shaded from most identical (red) to most dissimilar (blue). The reproducibility of selection in the primary and alternative position 115 libraries is evident in the results.

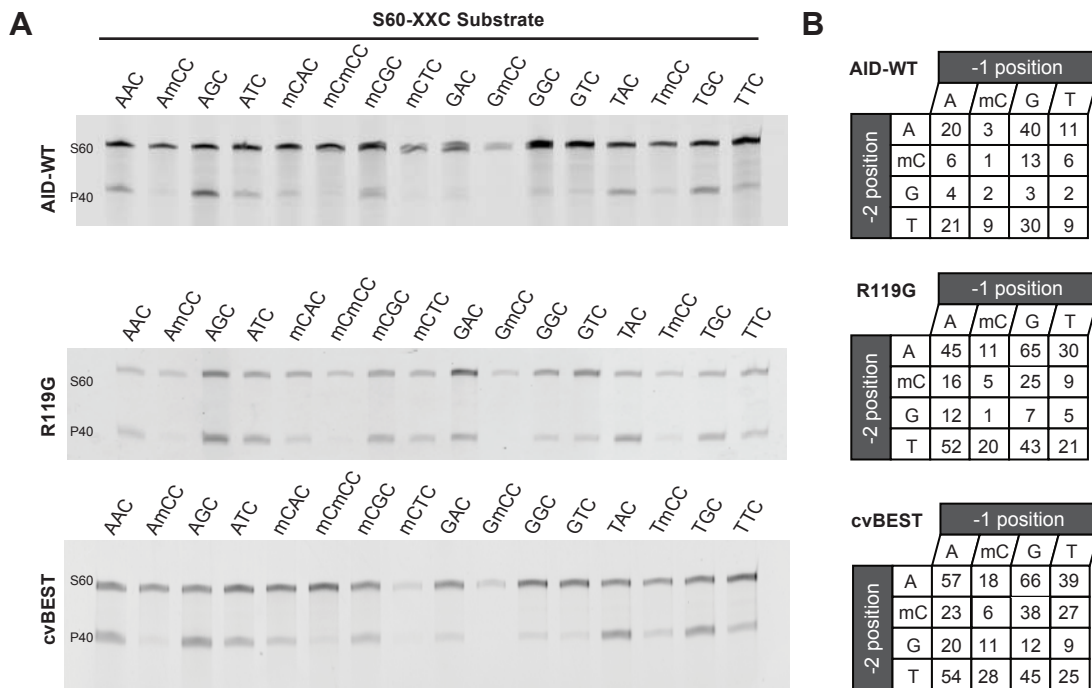


Figure S6. Sequence Preference Profiles (**A**) To determine the sequence preference profiles, AID-WT, R119G and cvBEST were assayed against a panel of sixteen 3'-fluorophore labeled ssDNA 60 bp substrates containing a single cytosine in the context of varying -1 and -2 position residues (S60-XXC). The residue X was varied between A, 5-methylcytosine (mC), G and T. The reactions were run with 150 nM DNA for 3 hrs with enzyme concentrations selected such that substrate turnover was not complete under any condition (AID-WT 200 nM, R119G 80 nM, cvBEST 80 nM). The S60 substrates were fragmented to generate a 40 bp product by treatment with UDG and base. Shown is a representative gel from one assay. (**B**) Product formation from two replicates was averaged to give the percent product formed under standard conditions with each substrate, with associated error of <25% for each measurement. Sequence preference profiles (Figure 5B) were generated by a two-step calculation. First, the product formation was averaged for all sequences that contain the same nucleotide at either the -1 or -2 position. For example, P60-XAC (the average product formation for the S60-XAC substrates) was calculated by averaging product formation with S60-AAC, -mCAC, -GAC and -TAC. Next, the probability of deamination for each nucleotide at the -1 or -2 position was calculated relative to the other nucleotide variants. For example, the percent preference for A at -1 was calculated as $(P60-XAC / (P60-XAC + P60-XmCC + P60-XGC + P60-XTC)) \times 100$.

SUPPLEMENTARY TABLES

Table S1 . Hydrogen Bonding Interactions Between AID and 5'-AGCT-3'.

Residue	AID-WT			R119G			cvBEST		
113	-			GUA-Side-N2	LEU113-Main-O	15.07%	GUA-Side-N2	LEU113-Main-O	30.95%
				GUA-Side-N1	LEU113-Main-O	10.10%	GUA-Side-N1	LEU113-Main-O	53.35%
114	-			-			GUA-Side-N2	TYR114-Main-O	3.03%
117	GUA-Side-N2	GLU117-Side-OE1	3.93%	GUA-Side-N2	GLU117-Side-OE2	5.82%	-		
	GUA-Side-N2	GLU117-Side-OE2	3.68%	GUA-Side-N2	GLU117-Side-OE1	4.83%			
	GUA-Side-N1	GLU117-Side-OE2	3.45%	GUA-Side-N1	GLU117-Side-OE2	3.70%			
	GUA-Side-N1	GLU117-Side-OE1	2.40%	GUA-Side-N1	GLU117-Side-OE1	3.18%			
119	ARG119-Side-NH2	GUA-Side-N7	16.65%	-			-		
	ARG119-Side-NH1	GUA-Side-O6	9.48%						
	ARG119-Side-NE	GUA-Side-O6	6.48%						
	ARG119-Side-NH2	ADE-Side-O5'	2.92%						
	ARG119-Side-NH1	GUA-Side-N7	2.85%						
	ARG119-Side-NH2	GUA-Side-O2P	2.27%						

Hydrogen bond occupancy analysis was performed using HBonds Plugin (Version 1.2) in VMD. The hydrogen bond occupancy reflects the percentage of a simulation that a particular hydrogen bond exists with given cutoff criteria. Moderate and strong hydrogen bonds were included by defining a bond cutoff length of 3.2 Å (between centers of heavy atoms) and cutoff angle of >150 degrees. Listed are the hydrogen bonding interactions that occurred at > 2% frequency during the simulation.

Table S2 . Solvent accessible surface area for side chain residues

Residue	Side Chain Size, Å ²	AID-WT, Å ² (% Accessible)	Y114F, Å ² (% Accessible)	R119G, Å ² (% Accessible)	cvBEST, Å ² (% Accessible)
113	163.7	36.9 (22.5%)	-	14.2 (8.7%)	35.0 (21.4%)
114	190.6 (F) 209.6 (Y)	116.5 (55.6%)	73.2 (38.4%)	-	-
115	196.7	5.6 (2.9%)	-	-	-

Solvent accessible surface area (SASA) for Residues 113, 114 and 115 were computed in VMD using a probe with radius of 1.4 Å. Reported values represent the average solvent exposed surface area of a given residue side chain over the NVT trajectory (in Å²) or are scaled (for % Accessible) relative to the average total surface area of the residue (solvent exposed area plus buried surface area).

SUPPLEMENTARY VIDEO

Video S1. Molecular Dynamics Simulations of AID with 5'-AGCT-3' DNA. The AID model is shown in gray with highlighted residues L113 (yellow), Y114 (green), F115 (purple) and R119 (blue). The DNA is shaded from red to white in the 5' to 3' direction. The model was equilibrated for 40 ns with constrained DNA. Shown are the subsequent 120 ns of the simulation between AID-WT and unconstrained DNA.