

1 **SUNi mutagenesis: scalable and uniform nicking for efficient generation of
2 variant libraries**

3

4 Taylor L. Mighell^{1*}, Ignasi Toledano¹, Ben Lehner^{1, 2, 3, 4*}

5

6 ¹ Center for Genomic Regulation (CRG), The Barcelona Institute of Science and
7 Technology, Barcelona, Spain.

8 ² Universitat Pompeu Fabra (UPF), Barcelona, Spain.

9 ³ Institut^o Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain.

10 ⁴ Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK

11 *Correspondence: taylor.mighell@crg.eu, bl11@sanger.ac.uk

12

13

14 **Abstract**

15 Multiplexed assays of variant effects (MAVEs) have made possible the
16 functional assessment of all possible mutations to genes and regulatory sequences. A
17 core pillar of the approach is generation of variant libraries, but current methods are
18 either difficult to scale or not uniform enough to enable MAVEs at the scale of gene
19 families or beyond. We present an improved method called Scalable and Uniform
20 Nicking (SUNi) mutagenesis that combines massive scalability with high uniformity to
21 enable cost-effective MAVEs of gene families and eventually genomes.

22

23 **Background**

24 Massive mutagenesis followed by functional assays, commonly known as
25 MAVEs or deep mutational scanning, is a powerful strategy for understanding the
26 effects of genetic variation[1–3], dissecting and engineering proteins[4–6], and directed
27 evolution[7]. Modern approaches for generating mutagenesis libraries generally fall into
28 two categories. First, synthesis of oligonucleotides containing programmed mutations
29 followed by subcloning, known as cassette or tile mutagenesis[2,8,9]. Second,
30 synthesis of polymerase chain reaction (PCR) primers containing programmed
31 mutations which bind template DNA and extend to form a mutated strand, followed by
32 various means of degrading and resynthesizing the opposite strand to form mutated
33 double-stranded DNA[10–12].

34 While cassette mutagenesis yields highly uniform libraries, current DNA
35 synthesis technologies can only generate oligonucleotides up to length ~300, with
36 synthesis quality decaying rapidly with increased length[13]. Since many genes exceed
37 this length, it is necessary to generate sub-libraries, which require complex
38 experimental designs that severely limit scalability. The key advantage of primer-based
39 mutagenesis is that it does not have this limitation; in theory, any number of genes of
40 any length can be mutated in a single pot. However, primer-based methods suffer from
41 differences in mutagenesis efficiency between positions, resulting in libraries with
42 highly nonuniform representation of variants[10–12]. Additionally, primer-based
43 methods can generate substantial amounts of wild-type carryover, requiring the use of
44 larger experimental volumes, increased sequencing, and sequencing errors artificially
45 inflating counts for variants[14]. These drawbacks are problematic because they
46 reduce data quality and increase the cost of every step of a MAVE experiment, thereby
47 limiting scalability.

48 The Atlas of Variant Effects (AVE) Alliance has the goal to quantify the impact
49 of variation in most human genes and regulatory elements using diverse selection
50 assays[15]. With current rates of progress this endeavor is likely to take decades to

51 achieve[16]. Here we detail a protocol that we term Scalable and Uniform Nicking
52 (SUNi) mutagenesis that represents a two-fold improvement over the existing state of
53 the art method[12] for large variant library construction. SUNi mutagenesis yields highly
54 uniform variant libraries with massive potential scalability.

55

56 **Results and Discussion**

57 Nicking mutagenesis generates mutated plasmid in four steps: degradation of
58 one DNA strand; annealing and extension of a mutagenic primer; degradation of the
59 opposite strand; and resynthesis of the opposite strand, incorporating the mutation[12]
60 (Supplementary Fig. 1). Previous data indicated that longer homology arms could
61 improve mutagenesis efficiency[17], and that the melting temperature (T_m) of the
62 mutagenic primer was correlated with mutagenesis efficiency[18]. We reasoned that
63 since binding of both homology arms to the template is required for efficient
64 mutagenesis, performance could be improved by optimizing the T_m of both arms of the
65 primer independently. Therefore, we designed a pool of primers (referred to as opt1)
66 where, for each position, the left and right homology arm had the length between 20-40
67 nucleotides that had the predicted T_m closest to 61°. These primers were designed to
68 target two 40-codon regions of the μ opioid receptor (MOR) which were chosen
69 because of very high or low GC content (MOR2 =65.8% GC, MOR6 =40.8% GC) and
70 so were expected to provide the greatest challenge for the new design. Advances in
71 DNA synthesis have made oligonucleotide pools an affordable, and therefore scalable,
72 option for synthesizing large numbers of sequences. A previous version of nicking
73 mutagenesis synthesized primers as microarray-based oligonucleotide pools, but the
74 quality of these libraries was substantially lower than the original method[18], possibly
75 due to the femtomole-scale yield of microarray synthesis. To maintain the scalability
76 advantage of oligo pools while still maximizing library quality, we synthesized our
77 primers as IDT oPools, which have picomole-scale yield.

78 The sequential degradation of each DNA strand of a plasmid is accomplished
79 with the nicking activity (cleavage of only one strand of double stranded DNA) of
80 engineered variants of the BbvCI restriction enzyme. We found that some plasmids
81 containing only one BbvCI site are inefficiently digested in the first nicking step,
82 potentially leading to wild-type carryover. Adding a second BbvCI site to the plasmid
83 improved digestion efficiency (Supp. Fig 2). Therefore, we engineered a plasmid
84 bearing MOR to contain two BbvCI sites, and followed the published nicking protocol
85 with minor modifications (Supplementary Protocol 1). Sequencing of the mutagenesis
86 libraries revealed similar proportions of programmed mutations (63.8 and 58.8% for
87 opt1 versus 65.3 and 64.2% for standard nicking) and slightly increased wild-type
88 percent (26.9 and 33.2% for opt1 versus 23.8 and 23.3% for standard nicking) but with
89 improved uniformity (log difference (LogDiff) between 90th and 10th percentile of
90 mutants of 0.83 and 0.92 for opt1 libraries versus 1.18 and 0.94 for standard nicking
91 libraries, Fig. 1a,b). While overall uniformity was improved, there was still substantial
92 positional bias (Fig. 1b), which we next sought to understand. However, we found no
93 relationship between mutagenesis frequency (median frequency of all programmed
94 mutations per position) and predicted T_m of left or right mutagenesis primer homology
95 arm, or for minimum, maximum, sum, or difference between left and right T_m. We also
96 found no contribution of predicted free energy of secondary structure formation of
97 primers (Supplementary Table 4).

98 Surprisingly, we did find a significant contribution of GC content of the five 5'
99 terminal bases of the primer. The strongest signal comes when considering GC content
100 of the three 5' terminal bases (Spearman $\rho=0.56$, $p=6.8\times 10^{-8}$ Fig. 1c,d). A GC-rich 3'
101 terminus of a primer (also known as “GC clamp”) is widely thought to improve priming
102 efficiency, but here we find no contribution of 3' GC clamp (Supplementary Table 4).
103 We divided primers based on the 5' terminus sequence and found that primers with
104 SSS, SWS, or SSW sequence (from 5' to 3', where S =G or C and W =A or T) have the
105 highest median mutagenesis efficiency (Fig. 1e). Conceptually, the importance of a 5'

106 GC clamp makes sense because the extension step of the mutagenesis PCR is long
107 and at a relatively high temperature (7 minutes at 72°), and if the mutagenic primer
108 terminus is dissociated from the template when the polymerase completes the
109 mutagenic strand, it may polymerize extra bases and make ligation of the mutagenic
110 strand impossible.

111 We designed a new set of nicking primers (referred to as SUNi), targeting the
112 same regions, and taking advantage of the 5' GC clamp discovery. Briefly, for each
113 position we sought to find a primer that had optimal predicted T_m and also a strong 5'
114 GC clamp (full description in Methods). Further, we reasoned that one contribution to
115 wild-type carryover is NNK primers in which the wild-type codon is encoded by NNK.
116 Since K encodes G and T, for any codon that ends in these bases the wild-type
117 sequence will be present in the NNK pool, and this fully complementary wild-type
118 primer would be expected to outcompete mutation-bearing primers. To minimize this,
119 we used NNK to mutagenize codons that end in A, C, or G, and NNS (S= G or C) if the
120 wild-type codon ended in T. Sequencing of MOR2 and MOR6 SUNi mutagenesis
121 libraries demonstrated increased percentage of programmed mutants (77.5 and 68.9%,
122 respectively) decreased percentage of wild-type (13.9 and 23.7%, respectively), and
123 improved uniformity (LogDiff = 0.65, 0.92, respectively, Fig. 2a).

124 We wanted to compare methods using a more comprehensive metric, so we
125 calculated screening efficiency = $\frac{\% \text{ programmed}}{10^{\text{LogDiff}}}$, a term which incorporates the fraction of
126 programmed sequences in the library and the uniformity of those sequences, which are
127 both important to determine the efficiency of screening the library. Screening efficiency
128 for both libraries increases from opt1 to SUNi designs, and on average SUNi is twice
129 as efficient as the standard nicking protocol (0.128 versus 0.058, respectively, Fig.
130 2b,c). We also compared a mutagenesis library made by cassette mutagenesis
131 (b2AR2, 250 nucleotide oligonucleotides introducing mutations at 70 positions). We
132 find that in the best case (MOR2), SUNi screening efficiency approaches that of

133 cassette mutagenesis (0.173 versus 0.200, respectively, Fig. 2b,c), while requiring
134 substantially less hands-on time and allowing mutagenesis of much larger and many
135 different targets in a single reaction pool. Cassette mutagenesis yields highly uniform
136 libraries, but the percent of programmed mutants is low (Fig. 2c) due to errors in DNA
137 synthesis.

138 We chose to mutagenize regions with high and low GC content, assuming
139 these would be difficult templates for mutagenesis. However, we didn't anticipate the
140 crucial importance of the 5' GC clamp. The data suggests that the mutagenesis
141 efficiency of SUNi is likely related to GC content, indicating that MOR6 is likely difficult
142 while MOR2 is likely an amenable template. We expect SUNi mutagenesis efficiency
143 for regions with intermediate GC content to be intermediate between the examples
144 shown here.

145 SUNi mutagenesis has the potential to be massively scaled, as there is no
146 theoretical limit to the length of mutated region or the number of mutated regions in a
147 single reaction. The efficiency of screening a SUNi library is twice that of the standard
148 nicking protocol, meaning that at all steps (library generation, screening, and
149 sequencing), the reagents required, and therefore cost, will be halved. We expect SUNi
150 mutagenesis coupled with a panel of selection assays[16] will allow the rapid and cost-
151 effective generation of variant effect atlases for entire gene families. The bright future
152 of MAVEs is reliant on scalable methods for generating high quality variant libraries,
153 and SUNi mutagenesis represents an important step in that direction.

154

155 **Conclusions**

156 More efficient libraries empower more scalable experiments that will be
157 necessary for generating atlases of variant effect at the gene-family or genome scale.
158 In this report, we outline design and experimental improvements that improve the
159 screening efficiency of nicking mutagenesis two-fold.

160

161 **Materials & Methods**

162 **Opt1 primer design**

163 Primers were designed to introduce all single amino acid mutations and stop codon
164 (via “NNK” codon mutagenesis) for 80 codons in the μ opioid receptor (MOR). To pick
165 a guide for each position, for each homology arm, we found the candidate between 20
166 and 40 nucleotides with T_m closest to 61° (calculated with biopython[19] using the
167 Bio.SeqUtils.MeltingTemp.Tm_NN function). The two pools of opt1 primers were
168 ordered as IDT oPools. Sequences reported in Supplementary Table 2.

169 **SUNi primer design**

170 Like opt1, we designed primers to introduce single amino acid mutations at 80
171 positions of MOR. For each position, we found the right homology arm in the same way
172 as for Library 1, i.e. the arm between 20 and 40 nucleotides that had predicted T_m
173 closest to 61°. For the left homology arm, we enumerated all arms that had predicted
174 T_m between 59° and 66°. If one or more of these arms had all three 5' terminal
175 nucleotides as S (degenerate codon notation; S=G or C, W=A or T), the shortest of
176 these was chosen. If there were no SSS 5' termini, then we looked for arms with SSW
177 or SWS termini, and if there were one or more, we chose the shortest arm. If there
178 were no suitable homology arms with SSW or SWS termini, we then found the arm
179 closest to 64° irrespective of 5' terminus. Since we would then predict this primer to be
180 suboptimal, we encoded it twice in the oPool. In this library we used NNK as the
181 degenerate mutagenic codon if the WT codon ended in A, C, or G, but we used NNS if
182 the WT codon ended in T. The two pools of SUNi primers were ordered as IDT oPools.
183 Sequences reported in Supplementary Table 2.

184 **b2AR2 mutagenesis**

185 Oligonucleotides were designed to introduce all possible single amino acid changes,
186 and many double amino acid changes, for a total of 4005 variants. These were
187 synthesized by Twist Bioscience as 250 nucleotide oligos. PCR with primers
188 dialout_tile2_[F/R] (primers used in this study reported in Supplementary Table 1) was

189 done to amplify these mutagenic oligos. PCR with primers designed to amplify the rest
190 of the vector besides the region to be mutagenized (b2AR_satmut_tile2_[F/R]) was
191 performed to prepare the vector, and then Gibson assembly was used to introduce the
192 mutagenic oligos.

193 **Sequencing library preparation**

194 Two stage PCR was performed to amplify each mutated region and append indexed
195 Illumina sequencing adapters. Q5 High Fidelity polymerase (New England Biolabs)
196 was used for all PCRs. For MOR2 and MOR6 regions, primers
197 MOR_nicking_T[2/6]_seq_[F/R] were used in stage1 PCR to amplify the target and
198 append partial Illumina sequencing adapters, with 50 ng of purified plasmid as
199 template. Cycling protocol was 98° for 30s, followed by 17 cycles of [98° for 20s, 55°
200 for 30s, 72° for 30s]. Products were column purified and 0.2% of PCR1 was used as
201 input for PCR2 with primers indexed_i[5/7] and cycled with 98° for 30s, followed by 5
202 cycles of [98° for 15s, 64° for 30s, 72° for 30s]. Products were column purified and
203 sequenced on Illumina Nextseq 500 or Nextseq 2000 instruments. For b2AR2, 10 ng of
204 purified plasmid was used as input to PCR using primers b2AR_Tile2_PCR1_5N_[F/R]
205 and cycling with 98° for 30s, followed by 12 cycles of [98° for 15s, 66° for 30s, 72° for
206 30s]. Products were column cleaned and 0.2% of PCR1 was used as input for PCR2
207 with primers indexed_i[5/7] and cycled with 98° for 30s, followed by 10 cycles of [98°
208 for 15s, 64° for 30s, 72° for 30s]. Products were column cleaned and sequenced on
209 Illumina MiSeq instrument.

210 **Sequencing data processing**

211 We obtained raw fastq data from the original nicking paper[12] from the Short Read
212 Archive with accession numbers SRR4105481 and SRR4105482. All fastq data were
213 processed identically: first, read pairs were merged and filtered for reads which
214 contained <0.5 expected errors using vsearch[20]. Then, cutadapt[21] was used to trim
215 adapters and only those reads with matching adapters were retained. Variant counts
216 were enumerated by comparing sequencing reads to expected sequences based on

217 mutagenesis strategy (i.e. NNN, NNK, or NNS) and counting only perfect matches.

218 Read processing data in Supplementary Table 3.

219 **Availability of data and materials**

220 Code to generate SUNi mutagenesis primers is available at <https://github.com/lehner->

221 [lab/SUNi_mutagenesis](#). Raw sequencing data produced for this study can be found at

222 the Sequence Read Archive with accession number PRJNA939024.

223 **Figure 1. Optimization and analysis of nicking mutagenesis primer design**

224 **a**, Per position mutation frequency presented as fraction of all sequencing reads for
225 standard nicking. Dashed lines indicate 90th and 10th percentile of all mutation
226 frequencies.

227 **b**, Per position mutation frequency presented as fraction of all sequencing reads for
228 opt1 nicking. Dashed lines indicate 90th and 10th percentile of all mutation frequencies.

229 **c**, Spearman correlation between GC content of the 5' terminus and mutagenesis
230 efficiency, when considering between one and five terminal bases.

231 **d**, Mutagenesis frequency of positions with different GC content in the 5' terminal three
232 bases. Spearman $\rho=0.56$, $p=6.8\times 10^{-8}$.

233 **e**, Mutagenesis frequency of positions with different SW sequences (S=G or C, W=A or
234 T) in the 5' terminal three bases.

235

236 **Figure 2. Performance and comparison of SUNi mutagenesis with other methods**

237 **a**, Per position mutation frequency presented as fraction of all sequencing reads for
238 SUNi mutagenesis. Dashed lines indicate 90th and 10th percentile of all mutation
239 frequencies.

240 **b**, Screening efficiency of different mutagenesis methods.

241 **c**, Screening efficiency of different mutagenesis methods, as a function of uniformity
242 and percent programmed. Colors the same as in **b**.

243

References

- 244 1. Mighell TL, Evans-Dutson S, O'Roak BJ. A Saturation Mutagenesis Approach to
245 Understanding PTEN Lipid Phosphatase Activity and Genotype-Phenotype Relationships. *Am J*
246 *Hum Genet.* 2018;102:943–55.
- 247 2. Weile J, Sun S, Cote AG, Knapp J, Verby M, Mellor JC, et al. A framework for exhaustively
248 mapping functional missense variants. *Mol Syst Biol.* 2017;13:957.
- 249 3. Amorosi CJ, Chiasson MA, McDonald MG, Wong LH, Sitko KA, Boyle G, et al. Massively
250 parallel characterization of CYP2C9 variant enzyme activity and abundance. *Am J Hum Genet.*
251 2021;108:1735–51.
- 252 4. Faure AJ, Domingo J, Schmiedel JM, Hidalgo-Carcedo C, Diss G, Lehner B. Mapping the
253 energetic and allosteric landscapes of protein binding domains. *Nature.* 2022;604:175–83.
- 254 5. Penn WD, McKee AG, Kuntz CP, Woods H, Nash V, Gruenhagen TC, et al. Probing biophysical
255 sequence constraints within the transmembrane domains of rhodopsin by deep mutational
256 scanning. *Sci Adv.* American Association for the Advancement of Science; 2020;6:eaay7505.
- 257 6. Whitehead TA, Chevalier A, Song Y, Dreyfus C, Fleishman SJ, De Mattos C, et al. Optimization
258 of affinity, specificity and function of designed influenza inhibitors using deep sequencing. *Nat*
259 *Biotechnol.* 2012;30:543–8.
- 260 7. Acevedo-Rocha CG, Ferla M, Reetz MT. Directed Evolution of Proteins Based on Mutational
261 Scanning. *Methods Mol Biol Clifton NJ.* 2018;1685:87–128.
- 262 8. Melnikov A, Rogov P, Wang L, Gnrke A, Mikkelsen TS. Comprehensive mutational scanning
263 of a kinase in vivo reveals substrate-dependent fitness landscapes. *Nucleic Acids Res.*
264 2014;42:e112.
- 265 9. Macdonald CB, Nedrud D, Grimes PR, Trinidad D, Fraser JS, Coyote-Maestas W. DIMPLE:
266 deep insertion, deletion, and missense mutation libraries for exploring protein variation in
267 evolution, disease, and biology. *Genome Biol.* 2023;24:36.
- 268 10. Kitzman JO, Starita LM, Lo RS, Fields S, Shendure J. Massively parallel single-amino-acid
269 mutagenesis. *Nat Methods.* Nature Publishing Group; 2015;12:203–6.
- 270 11. Firnberg E, Ostermeier M. PFunkel: Efficient, Expansive, User-Defined Mutagenesis. *PLOS*
271 *ONE.* Public Library of Science; 2012;7:e52031.
- 272 12. Wrenbeck EE, Klesmith JR, Stapleton JA, Adeniran A, Tyo KEJ, Whitehead TA. Plasmid-based
273 one-pot saturation mutagenesis. *Nat Methods.* Nature Publishing Group; 2016;13:928–30.
- 274 13. Lietard J, Leger A, Erlich Y, Sadowski N, Timp W, Somoza MM. Chemical and photochemical
275 error rates in light-directed synthesis of complex DNA libraries. *Nucleic Acids Res.*
276 2021;49:6687–701.
- 277 14. Faure AJ, Schmiedel JM, Baeza-Centurion P, Lehner B. DiMSum: an error model and
278 pipeline for analyzing deep mutational scanning data and diagnosing common experimental
279 pathologies. *Genome Biol.* 2020;21:207.

280 15. AVE Alliance Founding Members. The Atlas of Variant Effects (AVE) Alliance: understanding
281 genetic variation at nucleotide resolution [Internet]. Zenodo; 2021 Mar. Available from:
282 <https://zenodo.org/record/7508716>

283 16. Tabet D, Parikh V, Mali P, Roth FP, Claussnitzer M. Scalable Functional Assays for the
284 Interpretation of Human Genetic Variation. *Annu Rev Genet.* 2022;56:441–65.

285 17. Kirby MB, Medina-Cucurella AV, Baumer ZT, Whitehead TA. Optimization of multi-site
286 nicking mutagenesis for generation of large, user-defined combinatorial libraries. *Protein Eng Des Sel.* 2021;34:gzab017.

288 18. Medina-Cucurella AV, Steiner PJ, Faber MS, Beltrán J, Borelli AN, Kirby MB, et al. User-
289 defined single pot mutagenesis using unamplified oligo pools. *Protein Eng Des Sel.*
290 2019;32:41–5.

291 19. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available
292 Python tools for computational molecular biology and bioinformatics. *Bioinformatics.*
293 2009;25:1422–3.

294 20. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for
295 metagenomics. *PeerJ. PeerJ Inc.*; 2016;4:e2584.

296 21. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
297 *EMBnet.journal.* 2011;17:10–2.



